

EDGE OF STOCHASTIC STABILITY: REVISITING THE EDGE OF STABILITY FOR SGD

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent findings by [Cohen et al. \(2021\)](#) demonstrate that when training neural networks with full-batch gradient descent with step size η , the largest eigenvalue λ_{\max} of the full-batch Hessian consistently stabilizes around $2/\eta$. These results have significant implications for convergence and generalization. This, however, is not the case of mini-batch stochastic gradient descent (SGD), limiting the broader applicability of its consequences. We show that SGD trains in a different regime we term Edge of Stochastic Stability (EOSS). In this regime, what stabilizes at $2/\eta$ is *Batch Sharpness*: the expected directional curvature of mini-batch Hessians along their corresponding stochastic gradients. As a consequence, λ_{\max} —which is generally smaller than *Batch Sharpness*—is suppressed, aligning with the long-standing empirical observation that smaller batches and larger step sizes favor flatter minima. We further discuss implications for mathematical modeling of SGD trajectories.

1 INTRODUCTION

The choice of training algorithm is a key ingredient in the deep learning recipe. Extensive evidence, e.g. [\(Keskar et al., 2016\)](#), indeed shows that performance consistently depends on the optimizer and hyperparameters. What machinery induces this optimizer-dependence is a central question of theory of deep learning.

[Cohen et al. \(2021; 2024\)](#) answered this question for Gradient Descent (GD): it optimizes neural networks in a regime of instability, they termed Edge of Stability (EOS). With a constant step size η , the highest eigenvalue of the Hessian of the full-batch loss—denoted here as λ_{\max} —grows until $2/\eta$ and hovers right above, subject to small oscillations ([Cohen et al., 2021; 2022; Jastrzębski et al., 2019; 2020; Xing et al., 2018](#)). Although, classical convex optimization theory call this step size “too large”, the loss continues to decrease. These works established a number of surprising facts: **(1)** that we require an optimization theory which works in more general scenarios than the classical $\eta < 2/\lambda_{\max}$; **(2)** that NN training happens in a special regime of instability, establishing what the source of it is; **(3)** how geometry of the local landscape around the solution found depends on the choice of hyperparameters.

Our finding: EOSS. While real-world training is almost always *mini-batch*—given the large amounts of data—existing EOS analyses **explicitly** do not apply to this case: no curvature-type quantities, such as λ_{\max} , are known to similarly affect SGD while training neural networks. We bridge this gap by establishing that:

Mini-batch SGD trains in a regime of instability akin to EOS which we term Edge of Stochastic Stability (EOSS). Precisely, *Batch Sharpness*, our notion of curvature,

$$\text{Batch Sharpness}(\theta) := \mathbb{E}_{B \sim \mathcal{P}_b} \left[\frac{\nabla L_B(\theta)^\top \mathcal{H}(L_B) \nabla L_B(\theta)}{\|\nabla L_B(\theta)\|^2} \right], \quad \begin{array}{l} \text{with } L_B \text{ being loss on the} \\ \text{batch } B \text{ sampled from } \mathcal{P}_b. \\ \mathcal{H}(\cdot) \text{ Hessian matrix.} \end{array}$$

hovers around $2/\eta$ and implicitly functions as sharpness for SGD. This implies that:
stability for SGD is stability on the mini-batch landscape

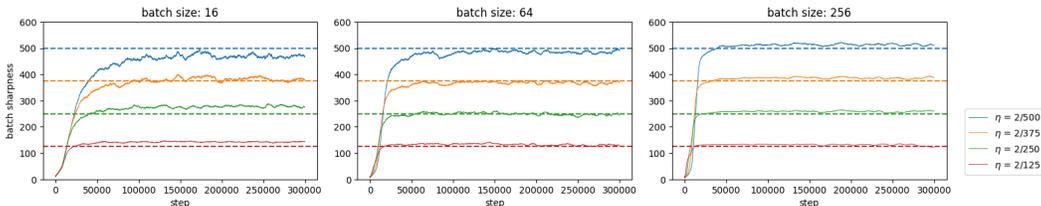


Figure 1: SGD at EOSS under different step sizes and batch sizes. MLP on an 8k subset of CIFAR-10 with step size $\eta > 0$. Batch Sharpness stabilizes at the $2/\eta$ threshold across varying batch sizes and step sizes.

Organization and Contributions. Section 2 reviews related work and outlines the key open questions we tackle. Oscillations are central to these phenomena, as a necessary step, in Section 3 we distinguish SGD oscillations between noise-driven (as in Robbins-Monro-type of stochastic optimization when the step size is kept fixed) and curvature-driven—which are the ones we are interested in. In Section 4, we introduce, properly characterize, and empirically validate the phenomenon of Edge of Stochastic Stability. In Section 5 we give a mathematical treatment of SGD stability. Finally, our results are yet another proof of the fact that the dynamics of noise-injected GD or SDEs and the dynamics of mini-batch SGD are qualitatively different and studying the firsts could be misleading for inducing properties of the second. We discuss this implication in Section 6.

Throughout the rest of this paper $B \subset \mathcal{D}$ denotes a random mini-batch of size b drawn from a fixed sampling distribution \mathcal{P}_b . For model parameters $\theta \in \mathbb{R}^d$ let $L_B(\theta) = \frac{1}{b} \sum_{(x_i, y_i) \in B} \ell(f_\theta(x_i), y_i)$, $L(\theta) = \mathbb{E}_{B \sim \mathcal{P}_b}[L_B(\theta)]$ be the mini-batch and full-batch losses, respectively, where ℓ is the loss function. Write $\mathcal{H}(L_B) = \nabla_\theta^2 L_B(\theta)$ and $\mathcal{H}(L)$ or \mathcal{H} the full-batch Hessian $\nabla_\theta^2 L(\theta)$.

2 RELATED WORK

Progressive sharpening. Early studies observed that the local shape of the loss landscape changes rapidly at the beginning of the training, by means of growth of different estimators of the curvature (Keskar et al., 2016; Jastrzębski et al., 2019; LeCun et al., 2012; Achille et al., 2017; Jastrzębski et al., 2018; Fort & Ganguli, 2019; Sagun et al., 2016; Fort & Scherlis, 2019). Subsequently, Jastrzębski et al. (2019; 2020) and Cohen et al. (2021) precisely characterized this behavior, demonstrating a steady rise in λ_{\max} along GD and SGD trajectories, typically following a brief initial decline. This phenomenon was termed *progressive sharpening* by Cohen et al. (2021).

Full-batch edge of stability. Prior research (Goodfellow et al., 2016; Li et al., 2019; Jiang et al., 2019; Lewkowycz et al., 2020) found that large initial learning rates often enhance generalization despite delaying initial loss reduction. Jastrzębski et al. (2020) attributed this effect to a phase transition, termed the break-even point, marking the end of progressive sharpening. Unlike progressive sharpening, this phenomenon is considered to result from algorithmic instability rather than inherent landscape properties. Indeed, Jastrzębski et al. (2019; 2020); Cohen et al. (2021; 2022) demonstrated that this phase transition comes at different points for different algorithms on the same landscapes. Cohen et al. (2021; 2022) later showed that it comes at the instability thresholds, in the case of full-batch optimization algorithms. Precisely, GD and full-batch Adam train in the EOS oscillatory regime (Cohen et al., 2021; 2022), where the λ_{\max} stabilizes and oscillates around a characteristic value. The name is due to the fact that, in the case of full-batch GD, the λ_{\max} hovers at $2/\eta$ which is the stability threshold for optimizing quadratics. Observations from Cohen et al. (2021; 2022) indicate that, under mean square error (MSE), the bulk of training dynamics occur within this regime, effectively determining λ_{\max} of the final solution. Lee & Jang (2023) explained why in this regime λ_{\max} often slightly exceeds $2/\eta$: this deviation arises primarily from nonlinearity of the loss gradient, which shifts the required value depending on higher-order derivatives, and the EOS being governed by the Hessian along the gradient direction, rather than λ_{\max} alone. A growing body of research analyzes the surprising mechanism underlying EOS dynamics observed during training with GD. Classically, when gradients depend linearly on parameters, divergence occurs locally if $\eta > \frac{2}{\lambda_{\max}}$, as illustrated by one-dimensional quadratic models (Cohen et al., 2021). In contrast, neural networks often converge despite violating this classical stability condition, presumably due

to the problem’s non-standard geometry. [Damian et al. \(2023\)](#) propose an explanation under some, empirically tested, assumptions of alignment of third derivatives and gradients.

Existing EOS work is limited to full-batch methods. While the empirical behavior of EOS for full-batch algorithms is relatively well-understood, neural networks are predominantly trained using mini-batch methods. As explicitly noted by [Cohen et al. \(2021\)](#), Section 6, Appendices G and H), their observations and analysis do not directly apply to mini-batch training, and EOS in SGD “does not center around the (full-batch) sharpness.” We show that the EOS phenomenon does indeed generalize to SGD, and we identify the key quantity governing this generalization (*Batch Sharpness* in Definition 3). We model stability of SGD on the neural networks landscapes: our results show that *SGD is stable if on average the step is stable on the mini-batch landscape—not on the full-batch landscape*.

What was empirically known for SGD. In the context of mini-batch algorithms, (i) [Jastrzębski et al. \(2019; 2020\)](#) noticed that for SGD the phase transition happens earlier for smaller η or smaller batch size b , but they did not quantify when. (ii) [Cohen et al. \(2021\)](#); [Gilmer et al. \(2021\)](#) established that initialization and architecture choices affect stability of SGD, without providing a definitive condition. (iii) When λ_{\max} stabilizes, that always happens at a level they could not quantify which is below the $2/\eta$ threshold ([Cohen et al., 2021](#); [Keskar et al., 2016](#)), see Figure 2, often without a proper progressive sharpening phase. This leaves the most basic questions open: *In what way the location of convergence of SGD acclimates to the choice of hyperparameters? What are the key quantities involved?* To be more specific, can we characterize the training phenomena in (i), (ii), (iii) above? What determines them? Does SGD train in an unstable regime?

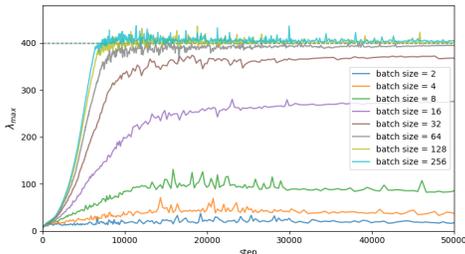


Figure 2: **SGD on CIFAR-10:** $\eta = 1/400$. The full-batch Hessian’s λ_{\max} plateaus below $2/\eta$. Smaller batch sizes lead to lower plateau values.

Previous works on SGD stability. A series of works, ([Wu et al., 2018](#); [Ma & Ying, 2021](#); [Granzio et al., 2021](#); [Wu et al., 2022](#); [Mulayoff & Michaeli, 2024](#)), studies constant-step-size SGD on quadratic losses via *linear stability* near the minima. From the perspective of Appendix B, these are valid instability criteria for a particular Lyapunov function, but two limitations remain for our purposes: (i) the resulting thresholds are expressed via d^2 -dimensional operators and are not computable for modern neural networks; and (ii) they do not address whether, and in what sense, SGD on neural networks actually *trains* in an EOS-like regime of instability along its real trajectory (limited by the aforementioned incomputability). For further discussion see Appendix B and L. A number of works ([Wu & Su, 2023](#); [Agarwala & Pennington, 2024](#)) showed that for SGD the regime of instability might be governed by the trace of loss Hessian/NTK, further discussed in Section M. Empirically, several works have documented oscillatory SGD dynamics in deep networks ([Xing et al., 2018](#); [Cohen et al., 2021](#); [Ahn et al., 2022](#); [Lee & Jang, 2023](#)). However, these works do not establish whether any of those oscillations constitute a regime of instability; in particular, they do not distinguish between noise-driven (Type-1) and curvature-driven (Type-2) oscillations, see Section 3. Our work complements these efforts by (i) placing candidate criteria such as λ_{\max} , GNI, and Batch Sharpness within a unified instability framework, and (ii) identifying Batch Sharpness as a valid, empirically saturating, and computationally tractable instability criterion for SGD on neural networks.

Flatness and Generalization. SGD-trained networks consistently generalize better than GD-trained ones, with smaller batch sizes further enhancing generalization performance ([Keskar et al., 2016](#); [LeCun et al., 2012](#); [Jastrzębski et al., 2018](#); [Goyal et al., 2017](#); [Masters & Luschi, 2018](#); [Smith et al., 2021](#); [Beneventano et al., 2024](#)). This advantage has been widely attributed to some notion of flatness of the minima ([Jiang et al., 2019](#); [Jastrzębski et al., 2021](#); [Hochreiter & Schmidhuber, 1994](#); [Neyshabur et al., 2017](#); [Wu et al., 2017](#); [Kleinberg et al., 2018](#); [Xie et al., 2020](#)). Training algorithms explicitly designed to find flat minima have indeed demonstrated strong performance across various tasks ([Izmailov et al., 2019](#); [Foret et al., 2021](#)). Our result is inherently a result about mini-batch training improving flatness. Specifically, we explain why: *Training with smaller batches constraints*

the dynamics to areas with smaller eigenvalues of the full-batch Hessian. This quantifies and characterizes prior observations that SGD tends to locate flat minima and that smaller batch sizes result in reduced Hessian sharpness (Keskar et al., 2016; Jastrzębski et al., 2021).

3 PRELIMINARIES: NOISE-DRIVEN vs CURVATURE-DRIVEN

The key defining aspect of EOS is about the solutions found by the algorithm adapting to the optimizer’s hyperparameters. In the case of full-batch algorithms, this manifests through the emergence of an oscillatory regime. Mini-batch SGD, however, always oscillates because its gradient is noisy and the step size is not annealed. The central question, therefore, is *which* oscillations signal curvature-limited dynamics (EOS-like). We define stable and unstable oscillations based on the induction of catapults.

Definition 1 (Quadratic instability and Catapults). Consider the quadratic approximation of all the data point loss landscapes $\frac{1}{2}(\theta - x_i)^\top \mathcal{H}_i(\theta - x_i)$. We say that a set of hyperparameters is unstable if the trajectory exits¹ all the compact subsets of the region in which the quadratic approximation holds up to $\mathcal{O}(\eta)$. We say the algorithm experienced a *catapult* when this event happened.

We define *Type-1* (Noise-Driven Oscillation) those that are stable under the definition above, e.g., when we increase the step size and the trajectory re-stabilizes within the neighborhood. We call *Type-2* (Curvature-induced) the oscillations which saturate stability, i.e., the ones for which a small change in the hyperparameters induces a catapult as defined in Definition 1. Interestingly, both types of oscillation involve quantities stabilizing near the critical threshold of $2/\eta$, yet they differ.

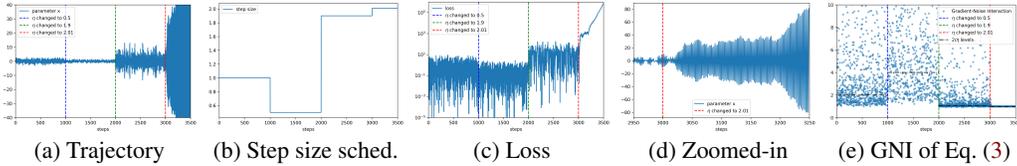


Figure 3: **Quadratics:** Dynamics of SGD on a 1-D quadratic with N datapoints, $L(x) = \frac{1}{2N} \sum_i (x - a_i)^2$, where $a_i \sim \mathcal{N}(0, 1)$. Oscillations are present for any step size. Yet, only when the step size becomes larger than $2/\lambda_{\max} = 2$ (after the red line), the oscillations become unstable (d) and the loss diverges (c). Meanwhile, GNI consistently stays at $2/\eta$.

3.1 Type-1 OSCILLATION AND GNI

SGD can wobble around a stationary point simply because gradients vary across batches and the step size is not annealed. This occurs even if the Hessian is small as, with fixed step-size, mini-batch noise has lower-bounded variance. Such noise-driven behavior is well-studied in classical stochastic approximation (Robbins & Monro, 1951; Mandt et al., 2016; Bottou et al., 2018; Mishchenko et al., 2020). We call *Type-1* oscillations any stochastic or chaotic trajectory which does not leave the region defined in Definition 1. We now introduce GNI , as a certificate of presence of oscillations:

Definition 2. We define Gradient-Noise Interaction (GNI):².

$$\text{Gradient-Noise Interaction}(\theta) := \frac{\mathbb{E}_{B \sim \mathcal{P}_b} [\nabla L_B(\theta)^\top \mathcal{H}(L) \nabla L_B(\theta)]}{\|\nabla L(\theta)\|^2} \quad (1)$$

GNI is defined by dividing the two terms in the classical descent lemma. *The SGD trajectory oscillates, no matter the reason, if and only if $GNI \approx 2/\eta$* —see Proposition 2 for a more rigorous statement and the proof—indeed:

Lemma 1. $\mathbb{E}_{B \sim \mathcal{P}_b} [L(\theta_{t+1})] \approx L(\theta_t)$ if and only if

$$-\eta \|\nabla L(\theta)\|^2 + \frac{\eta^2}{2} \mathbb{E}_{B \sim \mathcal{P}_b} [\nabla L_B(\theta)^\top \mathcal{H} \nabla L_B(\theta)] \approx 0 \iff GNI \approx \frac{2}{\eta}. \quad (2)$$

¹This means that either SGD seen as a linear dynamical system is diverging or that the re-stabilization would happen be at a level which exits the largest region in which the quadratic approximation holds and so the dynamics changes region.

²Note that both the Hessian \mathcal{H} and the gradient at the denominator are on the full-batch loss.

Notably, GNI is a quantity that is *centered* around $2/\eta$ whenever the trajectory oscillates. No matter the reason of the oscillation, see Figure 5. The regime of oscillation of SGD has been previously documented by measuring the expected total loss decrease by e.g., Cohen et al. (2021, Appendix H), Ahn et al. (2022), and Lee & Jang (2023) that tracked GNI explicitly³.

3.2 SGD ALWAYS OSCILLATES *Type-1*

Type-1 oscillations generally occur for SGD with fixed step size—even for simple quadratics⁴.

Proposition 1. Assume L_B are quadratic. Around a local minimum θ^* , fix $\eta > 0$ such that $\|(I - \eta\mathcal{H})^2 + \frac{\eta^2}{b}\mathbb{E}_B[\mathcal{H}(L_B)^2 - \mathcal{H}^2]\|_2^2 < 1$. Then the trajectory of SGD settles in a stationary distribution $\theta \sim \pi$ characterized by *Type-1* oscillations but not *Type-2* and satisfying

$$\frac{\mathbb{E}_{\theta \sim \pi}[\mathbb{E}_{B \sim \mathcal{P}_b}[\nabla L_B(\theta)^\top \mathcal{H} \nabla L_B(\theta)]]}{\mathbb{E}_{\theta \sim \pi}[\|\nabla L(\theta)\|^2]} = \frac{2}{\eta} [1 + \mathcal{O}(\eta)], \quad (3)$$

Independently of the moments of the Hessians \mathcal{H} and $\mathcal{H}(L_B)$.

See Appendix E for a formal statement and a proof. Crucially, the appearance of some quantity—GNI, as defined below—being $2/\eta$ implies the system is oscillating, not *why*. It does not mean, in principle, that the landscape or the curvature **adapted** to the hyper parameters. In this case (of *Type-1*), $2/\eta$ is about the ratio between the covariance of the gradients and the size of the full-batch gradient. Importantly, in this setting by perturbing the hyper parameters the system does not show catapults (as defined in Definition 1). When the size of oscillations increases (bigger step or smaller batch) the dynamics just increases the size of the oscillations—quickly restabilizing.

3.3 *Type-2*: CURVATURE-DRIVEN OSCILLATION

Once the *local*, or *perceived*, curvature saturates with respect to the hyperparameters, the updates become unstable in a manner analogous to the classic EoS (Cohen et al., 2021). We define *Type-2* oscillation the trajectories for which a small perturbation of the hyperparameters induces a Catapult.

Instability criteria, catapults, and EO(S)S. As we further formalize in Appendix B, we view a training algorithm with fixed hyperparameters h as a stochastic dynamical system $(\theta_t)_{t \geq 0}$ and summarize its local stability near a region U by a scalar *instability criterion* $f(\theta)$ with threshold c (Definition 4): if $f(\theta_0) > c$, the trajectory leaves every compact subset of U in finite time. In the deterministic quadratic full-batch case this recovers the classical condition $\lambda_{\max}(\nabla^2 L) \leq 2/\eta$: once $\lambda_{\max} > 2/\eta$, gradient descent is linearly unstable. We say that SGD trains at the *Edge of (Stochastic) Stability*—analogously, shows *Type-2* oscillations—when such a criterion empirically *saturates* along the trajectory, i.e. $f(\theta_t) \approx c$ (up to $O(\eta \text{poly}(\log(\eta)))$) for extended periods.

On the local quadratic approximation around θ_t , as we prove in Appendices B.2 and H, divergence has three equivalent manifestations: (i) breaking a valid instability criterion $f > c$; (ii) a *catapult* in the sense of Definition 1 (the quadratic trajectory leaves every compact subset of the trusted region U_t); and (iii) a *loss spike of sufficient size* on that quadratic model. If $f(\theta_t)$ is monotone in a destabilizing hyperparameter direction (e.g. $\eta \uparrow$ or $b \downarrow$) and is near c , Lemma 2 implies that a small perturbation of h pushes $f > c$ and forces catapults / loss spikes on the quadratic approximation. We call an oscillatory trajectory *Type-2* precisely when this happens: the dynamics oscillates (so $\text{GNI}(\theta_t) \approx 2/\eta$, cf. Lemma 1), and a curvature-based instability criterion is saturated so that small destabilizing perturbations reliably induce catapults. This is what we mean by *curvature-driven* oscillations.

Batch Sharpness as the curvature criterion. For mini-batch SGD, different Lyapunov functions lead to different scalar statistics of the Hessians (depending on various moments or cumulants of the random mini-batch Hessians), and there is no reason a priori for all of them to saturate. The mathematical treatment in Appendix B and Section 5 and the empirical in Section 4 single out *Batch Sharpness* as the relevant instability criterion:

³In their notations $\text{tr}(HS_b)/\text{tr}(S_n)$. See Appendix C for further comparison with previous work.

⁴See Figure 3 and Appendix D

Definition 3 (Batch Sharpness). We define

$$\text{Batch Sharpness}(\theta) := \mathbb{E}_{B \sim \mathcal{P}_b} \left[\frac{\nabla L_B(\theta)^\top \mathcal{H}(L_B) \nabla L_B(\theta)}{\|\nabla L_B(\theta)\|^2} \right]. \quad (4)$$

The SGD update follows $\nabla L_B(\theta)$, and *Batch Sharpness* is the expected Rayleigh quotient of $\mathcal{H}(L_B)$ along these directions. It measures the *average directional curvature of the mini-batch landscapes along the steps SGD actually takes*, in contrast to *GNI* (Definition 2), which mixes mini-batch gradients with the full-batch Hessian.

4 SGD TYPICALLY OCCURS AT THE EOSS

We characterize here the phenomenon of the Edge of Stochastic Stability. We verify the emergence of EOSS across a range of step sizes, batch sizes and architectures (Figure 6 and Appendix Q); datasets (CIFAR-10 and SVHN, Appendix R); and dataset sizes (8k and 32k subsets, Figure 7).

1. Stabilization of *Batch Sharpness*. SGD typically trained in an EOS-like regime:

*SGD tends to train in a regime we call Edge of Stochastic Stability. Precisely, after a phase of progressive sharpening, *Batch Sharpness* reaches a stability level of $2/\eta$, and hovers there.*

In particular, the level of plateau of *Batch Sharpness* is $2/\eta$ independent of the batch size (Figure 1). Importantly, *Type-1* oscillations happen throughout most of the training as highlighted by the quantity of Proposition 1, see Figure 5, but they do not impact progressive sharpening which leads to the second phase of EOSS stabilization and *Type-2* oscillations. Importantly, analogously to EOS, training continues and the loss continues to decrease while *Batch Sharpness* is constrained by the step size magnitude.

2. Stabilization of λ_{\max} and *GNI*. Crucially, stabilization of *Batch Sharpness* around $2/\eta$ happens while *GNI* has stabilized at $2/\eta$ already, and induces a corresponding stabilization of λ_{\max} . However, λ_{\max} consistently settles at a lower level, due to a batch-size-dependent gap between the two. This is also influenced by the specific optimization trajectory, Figures 6 and 7. See Section I for factors determining their gap.

3. Catapults. Unlike in EOS, in the EOSS regime what stabilizes is an *expectation* of a quantity which the algorithm sees one observation at time. Occasionally, a sequence of sampled batches exhibits anomalously high sharpness—that is too high for the stable regime—and steps overshoot, triggering a catapult (Figure 4). This causes a spike in the loss (Section B.2), after which the trajectory either diverges or re-enters a stable region. If in this region *Batch Sharpness* is strictly less than $2/\eta$, this leads to a renewed phase of progressive sharpening, eventually returning to the EOSS regime. This aligns with, and maybe explains, previous observations about catapult behaviors, e.g., (Lewkowycz et al., 2020; Zhu et al., 2024).

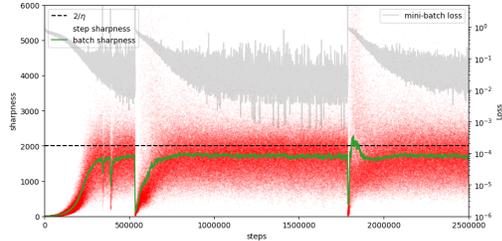


Figure 4: **Catapults at EOSS.** During EOSS, randomness in batch sampling might cause catapults, leading to renewed PS, and EOSS again. Notations follow Figure 6.

4.1 *Batch Sharpness* GOVERNS EOSS

Following Cohen et al. (2021) and the discussion in Section 3, we track how the training dynamics change when perturbing the hyperparameters mid-training. Overall, we find that *Batch Sharpness* governs EOSS behavior—mirroring how λ_{\max} operates in the full-batch EOS—while the full-batch λ_{\max} lags behind or settles inconsistently, underlining the mini-batch nature of SGD stability, see Appendix I. Increasing the step size η or decreasing the batch size b triggers a *catapult* spike in all the quantities in considerations and the training loss, before *Batch Sharpness* re-stabilizes near the

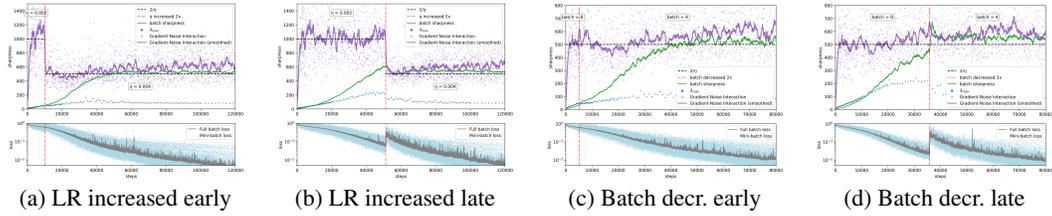


Figure 5: (1) The whole training happens with *Type-1* oscillations (see Proposition 1, $GNI \approx 2/\eta$), however, (2) GNI being $2/\eta$ does not govern *Type-2* oscillations—in particular, highlighting the difference in the two types of oscillations. (3) *Batch Sharpness* is instead an indicator of *Type-2* oscillations, as illustrated by the fact that catapults happen only when the shift in hyperparameters occurs **after** *Batch Sharpness* reaches $2/\eta$.

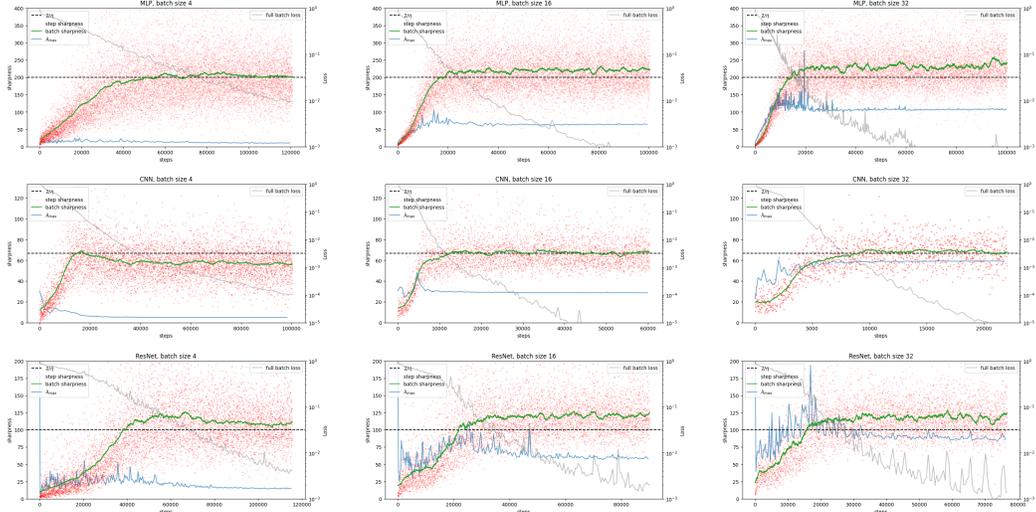


Figure 6: Comparing different sharpness measures. Red: *step sharpness*, observed sharpness on the current step’s mini batch—essentially *Batch Sharpness* without the expectation; Green: *Batch Sharpness* (Definition 3); Blue: full-batch λ_{\max} . Top row: MLP (2 hidden layers of width 512); middle: 5-layer CNN; bottom: ResNet-14; all trained on an 8k subset of CIFAR-10.

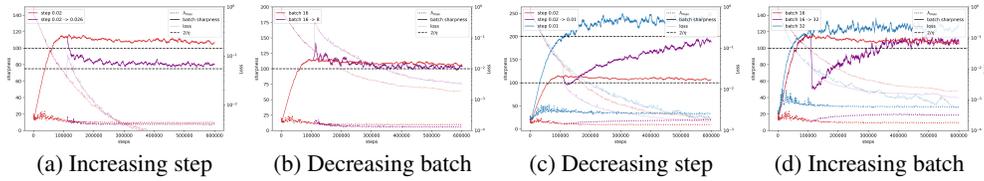


Figure 7: Effects of changing step size or batch size in EOSS. *Catapults*: (a) Increasing the step size η causes a catapult spike before *Batch Sharpness* re-settles at the new $2/\eta$. (b) Decreasing the batch size b increases *Batch Sharpness* and causes a catapult. *Restarting PS*: (c) Decreasing η prompts renewed progressive sharpening. (d) Increasing b lowers *Batch Sharpness* and re-starts progressive sharpening. The experiments are conducted on a 32k subset of CIFAR-10 to ensure sufficient complexity remains in the dataset, which is necessary for observing renewed progressive sharpening, consistent with observations by Cohen et al. (2021).

updated threshold $2/\eta$, see Figures 7a and 7b. This therefore pushes λ_{\max} lower. Conversely, reducing η raises the $2/\eta$ threshold. Analogously, increasing the batch size leaves λ_{\max} unchanged but reduces *Batch Sharpness*. These changes prompt a new phase of progressive sharpening, see Figures 7c and 7d. Notice that, *instantaneously*, the change in batch size does not change the full-batch loss landscape, but only changes the mini-batch landscapes—the fact that this causes a catapult/restarts PS is an indicator that it is indeed the mini-batch landscape (and therefore *Batch Sharpness*) that governs the stability/instability of SGD. Here, λ_{\max} also rises, but ultimately stabilizes at a lower value than if the entire training had run with the smaller step size/larger. Again, if stability was governed by λ_{\max} , this step-size adjustment would have had the same effect as starting from scratch with the new step size.

5 ON STABILITY

The previous section empirically demonstrated that mini-batch SGD generally settles into the EOSS regime, where *Batch Sharpness* hovers around $2/\eta$. Importantly, there exist many stochastic notions of stability, depending on different moments of the *random variable* $\mathcal{H}(L_B)$, see Appendix B. Some depend on quantities that can not be computed in high-dimensional experiments, others do not saturate empirically or their saturation does not induce EOSS, see Appendix M. In classical (full-batch) gradient descent, the condition $\eta < 2/\lambda_{\max}$ guarantees local stability by preventing divergence along the direction of the largest eigenvalue of a fixed Hessian. Here, is *Batch Sharpness* at $2/\eta$ *saturating the stability regime*? We answered positively empirically by showing that when you perturb hyperparameters you have explosions, see Figure 7. This proves empirically that we are at the Edge of Stability according to Definition 1. We show this mathematically in this section. Analogously, is *Batch Sharpness hovering at $2/\eta$ the cause of EOSS or a byproduct of something else happening*?. We already established empirically that the stabilization of λ_{\max} is a byproduct of it, see Figure 7 and Section I. We establish this causality proving Theorem 1 below. Precisely, Theorem 1 shows that the trajectory is unstable with respect to Definition 1 when *Batch Sharpness* is bigger than $2/\eta$. We proceed to discuss the meanings of *Batch Sharpness*, how it relates to previous criteria of instability, and why it has that form.

5.1 Type-2 IS ABOUT *Batch Sharpness*

The next theorem implies that SGD is unstable on quadratics if *Batch Sharpness* is bigger than $2/\eta$.

Theorem 1. *If *Batch Sharpness* is strictly bigger than $(2 + \epsilon)/\eta$ then $\mathbb{E}[\|\nabla L_B(\theta_{t+1})\|^2/\|\nabla L_B(\theta_t)\|^2] > (1 + \epsilon)^2$. On the second-order Taylor approximation, the norm of the mini-batch gradients increases exponentially with the SGD step and the trajectory is unstable in the sense of Definition 1.*

The proof relies on Jensen and Cauchy-Schwarz inequalities, see Appendix G. The use of these inequalities is its main limitation—we can not show the if and only if. However, Theorem 1 is the first (in)stability results that relies on a quantity we can efficiently estimate or compute in high-dimensional settings as neural networks. Note indeed that stability for quadratics is classically established by checking when $\mathbb{E}[\|\theta\|^2]$ diverges and when does not, see (Ma & Ying, 2021; Mulayoff & Michaeli, 2024) and Appendix L. *Batch Sharpness* is not directly related to these proofs and to the size of $\mathbb{E}[\|\theta\|^2]$.

5.2 CONNECTING BATCH SHARPNESS WITH EARLIER NOTIONS

In the full-batch case. In the case of deterministic GD, EOS and EOSS are equivalent. If $\lambda_{\max} \geq 2/\eta$, then *Batch Sharpness*—the Rayleigh quotient—quickly becomes bigger than $2/\eta$ as the gradients align with the top-eigenvectors. Viceversa, for full-batch, $\lambda_{\max} \geq$ *Batch Sharpness*.

Batch Sharpness governs instability. Theorem 1 allows to claim that *Batch Sharpness* generalizes λ_{\max} to the mini-batch case. Precisely, it generalizes that as it is a valid instability certificate. Indeed, we showed that when either of those cross $2/\eta$ the system becomes unstable.

Stability on the mini-batch landscape. The descent lemma on L_B shows that one SGD step on the mini-batch landscape is locally stable iff its directional curvature is below $2/\eta$:

$$\frac{\nabla L_B(\theta)^\top \mathcal{H}(L_B) \nabla L_B(\theta)}{\|\nabla L_B(\theta)\|^2} \leq \frac{2}{\eta} \iff -\eta \|\nabla L_B(\theta)\|^2 + \frac{\eta^2}{2} \nabla L_B(\theta)^\top \mathcal{H}(L_B) \nabla L_B(\theta) \leq 0. \quad (5)$$

Batch Sharpness at $2/\eta$ thus can be interpret as average such local stability.

Alignment does not come for free. λ_{\max} is *also* the averaged directional sharpness in the direction of the gradients for full-batch GD, because when it crosses $2/\eta$, the gradients align with the top eigenvectors. In the case of SGD, the gradients on different batches do not align though, as we show in Appendix K. *Batch Sharpness* is thus a natural generalization of λ_{\max} as the averaged directional sharpness. Moreover, we observe how, in some large batch size cases, *Batch Sharpness* hovers at $2/\eta$ but λ_{\max} is higher, as it does not govern stability due to mis-alignment of the gradients—e.g., Figure 44.

6 IMPLICATIONS: HOW NOISE-INJECTED GD DIFFERS FROM SGD

SGD vs. Noisy Gradient Descent. A common belief is that SGD’s regularization stems from its “noisy” gradients, which find flatter minima. Our analysis highlights how the noise in the Hessians as crucial. To test this, we compare mini-batch SGD (batch size 16) against three noisy GD variants—see details in Appendix J: (i) *Gaussian reweighting on the samples* (Wu et al., 2020) which maintains the noise structure in the Hessians; (ii) *Isotropic/Anisotropic diagonal noise* (Zhu et al., 2019); and SDE dynamics (Li et al., 2017). As shown in Figure 8, only noise which maintains the higher moments of the Hessian(s) (and thus implicitly preserves the mini-batch landscape structure) leads to an EOSS-like regime with λ_{\max} stabilizing well below $2/\eta$. Classical analyses of neural network optimization often assume noisy trajectories on a single, static, landscape. This is a further proof that the community has to be careful when modeling SGD as noise-injected GD or SDEs.

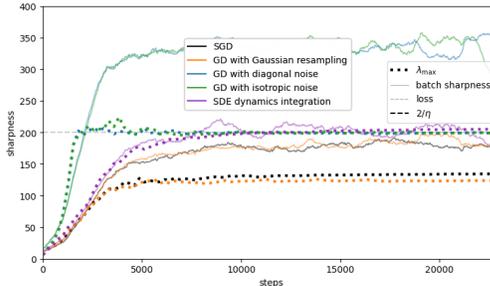


Figure 8: **SGD vs. Noisy GD vs SDE.** Only noise preserving the mini-batch structure of SGD leads to λ_{\max} plateauing below $2/\eta$ (akin to EOSS and as observed by (Keskar et al., 2016)). Noise injection fails to reproduce this behavior even with the same covariance SGD’s.

Standard SDE, noise-injected GD, or analogous approximations of SGD cannot describe the solution found by SGD or its behavior under the assumption of progressive sharpening. Indeed, they typically ignore any statistics of the Hessians except for the mean.

7 CONCLUSIONS, LIMITATIONS, AND FUTURE WORK

Conclusions. We have addressed the longstanding question of *if* and *how* mini-batch SGD enters a regime reminiscent of the “Edge of Stability” previously observed in full-batch methods. Contrary to the usual focus on the global Hessian’s top eigenvalue, we uncovered that *Batch Sharpness*—the expected directional curvature of the mini-batch landscape in the direction of its own gradient—consistently rises (progressive sharpening) and then hovers around $2/\eta$, independent of batch size. This behavior characterizes a new regime “Edge of Stochastic Stability”, which explains how mini-batch training can exhibit catapult-like surges and settle into flatter minima even when the full-batch Hessian remains below $2/\eta$. Our analysis clarifies why smaller batch sizes and larger step sizes both constrain the final curvature to a lower level, thereby linking these hyperparameters to flatter solutions and often improved generalization. Furthermore, we show that this phenomenon depends on the noise injected into the Hessians by mini-batch optimizers, highlighting important limitations of SDE-based approximations. Overall, the EOSS framework unifies several empirically observed effects—catapult phases, dependence on batch size, and progressive sharpening—under a single perspective focused on the *mini-batch* landscape and its directional curvature.

Limitations. (i) We have tested only image-classification tasks, leaving open whether similar phenomena arise in NLP, RL, or other domains. (ii) Our experiments mainly use fixed step sizes and standard architectures, so very large-scale or large-batch settings remain less explored. (iii) We have not analyzed momentum-based or adaptive methods (e.g. Adam), even though full-batch EOS has been seen there (Cohen et al., 2022).

Future Work. Beyond addressing these limitations, several directions remain: Understanding (i) *where* λ_{\max} stabilizes; (ii) how EOSS and EOS affect performances and properties of the neural network, e.g. (Lyu et al., 2023; Arora et al., 2022; Ahn et al., 2023; Zhu et al., 2023; Wang et al., 2022; Beneventano & Woodworth, 2025); (iii) consequently if it is benign effect or not; (iv) what the *other* sources of instability are there in the (pre-)training; (v) better describing the phenomenon of progressive sharpening; and (vi) understanding its causes.

REFERENCES

- 486
487
488 Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical Learning Periods in Deep Neural
489 Networks, 2017. URL <https://arxiv.org/abs/1711.08856v3>.
- 490
491 Atish Agarwala and Jeffrey Pennington. High dimensional analysis reveals conservative sharpening
492 and a stochastic edge of stability. *arXiv preprint arXiv:2404.19261*, 2024.
- 493
494 Kwangjun Ahn, Jingzhao Zhang, and Suvrit Sra. Understanding the unstable convergence of gra-
495 dient descent. In *Proceedings of the 39th International Conference on Machine Learning*, June
2022. URL <https://proceedings.mlr.press/v162/ahn22a.html>.
- 496
497 Kwangjun Ahn, Sébastien Bubeck, Sinho Chewi, Yin Tat Lee, Felipe Suarez, and Yi Zhang. Learn-
498 ing threshold neurons via the "edge of stability", October 2023. URL <http://arxiv.org/abs/2212.07469>. arXiv:2212.07469 [cs, math].
- 499
500 Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding Gradient Descent on Edge of
501 Stability in Deep Learning, October 2022. URL <http://arxiv.org/abs/2205.09745>.
502 arXiv:2205.09745 [cs].
- 503
504 Pierfrancesco Beneventano. On the Trajectories of SGD Without Replacement, December 2023.
505 URL <http://arxiv.org/abs/2312.16143>. arXiv:2312.16143.
- 506
507 Pierfrancesco Beneventano and Blake Woodworth. Gradient Descent Converges Linearly to Flatter
508 Minima than Gradient Flow in Shallow Linear Networks, January 2025. URL <http://arxiv.org/abs/2501.09137>. arXiv:2501.09137 [cs].
- 509
510 Pierfrancesco Beneventano, Andrea Pinto, and Tomaso Poggio. How Neural Networks Learn the
511 Support is an Implicit Regularization Effect of SGD. *arXiv:2406.11110 [cs, math, stat]*, June
512 2024. doi: 10.48550/arXiv.2406.11110. URL <http://arxiv.org/abs/2406.11110>.
513 arXiv:2406.11110 [cs, math, stat].
- 514
515 Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization Methods for Large-Scale Ma-
516 chine Learning, February 2018. URL <http://arxiv.org/abs/1606.04838>. arXiv:
1606.04838.
- 517
518 Jeremy M. Cohen, Simran Kaur, Yuanzhi Li, J. Zico Kolter, and Ameet Talwalkar. Gradient Descent
519 on Neural Networks Typically Occurs at the Edge of Stability. *arXiv:2103.00065 [cs, stat]*, June
520 2021. URL <http://arxiv.org/abs/2103.00065>. arXiv: 2103.00065.
- 521
522 Jeremy M. Cohen, Behrooz Ghorbani, Shankar Krishnan, Naman Agarwal, Sourabh Medapati,
523 Michal Badura, Daniel Suo, David Cardoze, Zachary Nado, George E. Dahl, and Justin Gilmer.
524 Adaptive Gradient Methods at the Edge of Stability, July 2022. URL <http://arxiv.org/abs/2207.14484>. arXiv:2207.14484 [cs].
- 525
526 Jeremy M. Cohen, Alex Damian, Ameet Talwalkar, Zico Kolter, and Jason D. Lee. Understanding
527 Optimization in Deep Learning with Central Flows, October 2024. URL <http://arxiv.org/abs/2410.24206>. arXiv:2410.24206.
- 528
529 Alex Damian, Tengyu Ma, and Jason Lee. Label Noise SGD Provably Prefers Flat Global Minimizers.
530 *arXiv:2106.06530 [cs, math, stat]*, June 2021. URL <http://arxiv.org/abs/2106.06530>.
531 arXiv: 2106.06530.
- 532
533 Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-Stabilization: The Implicit Bias of Gradient
534 Descent at the Edge of Stability, April 2023. URL <http://arxiv.org/abs/2209.15594>.
arXiv:2209.15594 [cs, math, stat].
- 535
536 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-Aware Minimization
537 for Efficiently Improving Generalization, April 2021. URL <http://arxiv.org/abs/2010.01412>. arXiv:2010.01412 [cs, stat].
- 538
539 Stanislav Fort and Surya Ganguli. Emergent properties of the local geometry of neural loss land-
scapes, 2019. URL <https://arxiv.org/abs/1910.05929v1>.

- 540 Stanislav Fort and Adam Scherlis. The Goldilocks Zone: Towards Better Understanding of Neural
541 Network Loss Landscapes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):
542 3574–3581, July 2019. ISSN 2374-3468. doi: 10.1609/aaai.v33i01.33013574. URL <https://ojs.aaai.org/index.php/AAAI/article/view/4237>. Number: 01.
- 544 Justin Gilmer, Behrooz Ghorbani, Ankush Garg, Sneha Kudugunta, Behnam Neyshabur, David Car-
545 doze, George Dahl, Zachary Nado, and Orhan Firat. A loss curvature perspective on training
546 instability in deep learning, 2021. URL <https://arxiv.org/abs/2110.04369>.
- 548 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 549 Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, An-
550 drew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet
551 in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- 552 Diego Granzio, Stefan Zohren, and Stephen Roberts. Learning rates as a function of batch size:
553 A random matrix theory approach to neural network training, 2021. URL <https://arxiv.org/abs/2006.09092>.
- 556 Jeff Z. HaoChen, Colin Wei, Jason D. Lee, and Tengyu Ma. Shape Matters: Understanding the
557 Implicit Bias of the Noise Covariance. *arXiv:2006.08680 [cs, stat]*, June 2020. URL <http://arxiv.org/abs/2006.08680>. arXiv: 2006.08680.
- 559 Sepp Hochreiter and Jürgen Schmidhuber. SIMPLIFYING NEURAL NETS BY DISCOVER-
560 ING FLAT MINIMA. In *Advances in Neural Information Processing Systems*, volume 7.
561 MIT Press, 1994. URL <https://proceedings.neurips.cc/paper/1994/hash/01882513d5fa7c329e940dda99b12147-Abstract.html>.
- 564 Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson.
565 Averaging Weights Leads to Wider Optima and Better Generalization, February 2019. URL
566 <http://arxiv.org/abs/1803.05407>. arXiv:1803.05407 [cs, stat].
- 567 Stanisław Jastrzębski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Ben-
568 gio, and Amos Storkey. Three Factors Influencing Minima in SGD. *arXiv:1711.04623 [cs, stat]*,
569 September 2018. URL <http://arxiv.org/abs/1711.04623>. arXiv: 1711.04623.
- 570 Stanisław Jastrzębski, Zachary Kenton, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos
571 Storkey. On the Relation Between the Sharpest Directions of DNN Loss and the SGD Step Length,
572 December 2019. URL <http://arxiv.org/abs/1807.05031>. arXiv:1807.05031 [stat].
- 574 Stanisław Jastrzębski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun
575 Cho, and Krzysztof Geras. The Break-Even Point on Optimization Trajectories of Deep Neural
576 Networks. *arXiv:2002.09572 [cs, stat]*, February 2020. URL <http://arxiv.org/abs/2002.09572>. arXiv: 2002.09572.
- 578 Stanisław Jastrzębski, Devansh Arpit, Oliver Astrand, Giancarlo Kerg, Huan Wang, Caiming Xiong,
579 Richard Socher, Kyunghyun Cho, and Krzysztof Geras. Catastrophic Fisher Explosion: Early
580 Phase Fisher Matrix Impacts Generalization. *arXiv:2012.14193 [cs, stat]*, June 2021. URL
581 <http://arxiv.org/abs/2012.14193>. arXiv: 2012.14193.
- 582 Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic
583 Generalization Measures and Where to Find Them. *arXiv:1912.02178 [cs, stat]*, December 2019.
584 URL <http://arxiv.org/abs/1912.02178>. arXiv: 1912.02178.
- 586 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Pe-
587 ter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv*
588 *preprint arXiv:1609.04836*, 2016.
- 589 Robert Kleinberg, Yuanzhi Li, and Yang Yuan. An Alternative View: When Does SGD Escape Local
590 Minima?, August 2018. URL <http://arxiv.org/abs/1802.06175>. arXiv:1802.06175
591 [cs].
- 592 Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In
593 *Neural networks: Tricks of the trade*, pp. 9–48. Springer, 2012.

- 594 Sungyoon Lee and Cheongjae Jang. A new characterization of the edge of stability based on a
595 sharpness measure aware of batch gradient distribution. In *International Conference on Learn-*
596 *ing Representations*, 2023. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:259298833)
597 [259298833](https://api.semanticscholar.org/CorpusID:259298833).
- 598 Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large
599 learning rate phase of deep learning: the catapult mechanism. *arXiv:2003.02218 [cs, stat]*, March
600 2020. URL <http://arxiv.org/abs/2003.02218>. arXiv: 2003.02218.
- 601
- 602 Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic
603 gradient algorithms. *arXiv:1511.06251 [cs, stat]*, June 2017. URL [http://arxiv.org/](http://arxiv.org/abs/1511.06251)
604 [abs/1511.06251](http://arxiv.org/abs/1511.06251). arXiv: 1511.06251.
- 605
- 606 Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and dynamics of stochastic
607 gradient algorithms i: Mathematical foundations, 2018. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1811.01558)
608 [1811.01558](https://arxiv.org/abs/1811.01558).
- 609
- 610 Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large
611 learning rate in training neural networks. In *Advances in Neural Information Processing Systems*,
612 pp. 11669–11680, 2019.
- 613 Zhiyuan Li, Sadhika Malladi, and Sanjeev Arora. On the Validity of Modeling SGD with Stochastic
614 Differential Equations (SDEs). *arXiv:2102.12470 [cs, stat]*, June 2021. URL [http://arxiv.](http://arxiv.org/abs/2102.12470)
615 [org/abs/2102.12470](http://arxiv.org/abs/2102.12470). arXiv: 2102.12470.
- 616
- 617 Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What Happens after SGD Reaches Zero Loss? –A
618 Mathematical Framework. *arXiv:2110.06914 [cs, stat]*, February 2022. URL [http://arxiv.](http://arxiv.org/abs/2110.06914)
619 [org/abs/2110.06914](http://arxiv.org/abs/2110.06914). arXiv: 2110.06914.
- 620
- 621 Kaifeng Lyu, Zhiyuan Li, and Sanjeev Arora. Understanding the Generalization Benefit of Normal-
622 ization Layers: Sharpness Reduction, January 2023. URL [http://arxiv.org/abs/2206.](http://arxiv.org/abs/2206.07085)
623 [07085](http://arxiv.org/abs/2206.07085). arXiv:2206.07085 [cs].
- 624
- 625 Chao Ma and Lexing Ying. On linear stability of SGD and input-smoothness of neural networks.
626 In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neu-*
627 *ral Information Processing Systems*, 2021. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=yAvCV6NwWQ)
[yAvCV6NwWQ](https://openreview.net/forum?id=yAvCV6NwWQ).
- 628
- 629 Stephan Mandt, Matthew Hoffman, and David Blei. A variational analysis of stochastic gradient
630 algorithms. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd*
631 *International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning*
632 *Research*, pp. 354–363, New York, New York, USA, 20–22 Jun 2016. PMLR. URL [https:](https://proceedings.mlr.press/v48/mandt16.html)
[//proceedings.mlr.press/v48/mandt16.html](https://proceedings.mlr.press/v48/mandt16.html).
- 633
- 634 Dominic Masters and Carlo Luschi. Revisiting Small Batch Training for Deep Neural Networks,
635 April 2018. URL <http://arxiv.org/abs/1804.07612>. arXiv:1804.07612.
- 636
- 637 Konstantin Mishchenko, Ahmed Khaled, and Peter Richtarik. Random Reshuffling: Simple Analy-
638 sis with Vast Improvements. In *Advances in Neural Information Processing Systems*, volume 33,
639 pp. 17309–17320. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper/2020/hash/c8cc6e90ccbfff44c9cee23611711cdc4-Abstract.html)
640 [cc/paper/2020/hash/c8cc6e90ccbfff44c9cee23611711cdc4-Abstract.](https://proceedings.neurips.cc/paper/2020/hash/c8cc6e90ccbfff44c9cee23611711cdc4-Abstract.html)
[html](https://proceedings.neurips.cc/paper/2020/hash/c8cc6e90ccbfff44c9cee23611711cdc4-Abstract.html).
- 641
- 642 Rotem Mulayoff and Tomer Michaeli. Exact mean square linear stability analysis for sgd, 2024.
643 URL <https://arxiv.org/abs/2306.07850>.
- 644
- 645 Behnam Neyshabur, Srinadh Bhojanapalli, David Mcallester, and Nati Srebro. Exploring Gener-
646 alization in Deep Learning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus,
647 S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*,
pp. 5947–5956. Curran Associates, Inc., 2017. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/7176-exploring-generalization-in-deep-learning.pdf)
[7176-exploring-generalization-in-deep-learning.pdf](http://papers.nips.cc/paper/7176-exploring-generalization-in-deep-learning.pdf).

- 648 Vardan Papyan. The Full Spectrum of Deepnet Hessians at Scale: Dynamics with SGD Training and
649 Sample Size, June 2019. URL <http://arxiv.org/abs/1811.07062>. arXiv:1811.07062
650 [cs, stat].
- 651 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- 652 Daniel A. Roberts. SGD Implicitly Regularizes Generalization Error. *arXiv:2104.04874 [cs, stat]*,
653 April 2021. URL <http://arxiv.org/abs/2104.04874>. arXiv: 2104.04874.
- 654 Levent Sagun, Leon Bottou, and Yann LeCun. Eigenvalues of the Hessian in Deep Learning:
655 Singularity and Beyond, November 2016. URL <https://openreview.net/forum?id=B186cP9gx>.
- 656 Samuel L. Smith, Benoit Dherin, David G. T. Barrett, and Soham De. On the Origin of Implicit
657 Regularization in Stochastic Gradient Descent. *arXiv:2101.12176 [cs, stat]*, January 2021. URL
658 <http://arxiv.org/abs/2101.12176>. arXiv: 2101.12176.
- 659 Yuqing Wang, Minshuo Chen, Tuo Zhao, and Molei Tao. Large learning rate tames homogeneity:
660 Convergence and balancing effect. In *International Conference on Learning Representations*,
661 2022. URL <https://openreview.net/forum?id=3tbDrs77LJ5>.
- 662 Jingfeng Wu, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu. On
663 the noisy gradient descent that generalizes as sgd, 2020. URL <https://arxiv.org/abs/1906.07405>.
- 664 Lei Wu and Weijie J. Su. The Implicit Regularization of Dynamical Stability in Stochastic Gradient
665 Descent, June 2023. URL <http://arxiv.org/abs/2305.17490>. arXiv:2305.17490
666 [stat].
- 667 Lei Wu, Zhanxing Zhu, and Weinan E. Towards Understanding Generalization of Deep Learning:
668 Perspective of Loss Landscapes. *arXiv:1706.10239 [cs, stat]*, November 2017. URL <http://arxiv.org/abs/1706.10239>. arXiv: 1706.10239.
- 669 Lei Wu, Chao Ma, and Weinan E. How SGD Selects the Global Minima in Over-
670 parameterized Learning: A Dynamical Stability Perspective. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL
671 https://proceedings.neurips.cc/paper_files/paper/2018/hash/6651526b6fb8f29a00507de6a49ce30f-Abstract.html.
- 672 Lei Wu, Mingze Wang, and Weijie Su. The alignment property of sgd noise and how it helps select
673 flat minima: A stability analysis, 2022.
- 674 Zeke Xie, Issei Sato, and Masashi Sugiyama. A Diffusion Theory For Deep Learning Dynamics:
675 Stochastic Gradient Descent Exponentially Favors Flat Minima. *arXiv:2002.03495 [cs, stat]*,
676 November 2020. URL <http://arxiv.org/abs/2002.03495>. arXiv: 2002.03495.
- 677 Chen Xing, Devansh Arpit, Christos Tsirigotis, and Yoshua Bengio. A Walk with SGD, May 2018.
678 URL <http://arxiv.org/abs/1802.08770>. arXiv:1802.08770 [cs, stat].
- 679 Sho Yaida. Fluctuation-dissipation relations for stochastic gradient descent. *arXiv preprint*
680 *arXiv:1810.00004*, 2018.
- 681 Hanlin Zhang, Depen Morwani, Nikhil Vyas, Jingfeng Wu, Difan Zou, Udaya Ghai, Dean Foster,
682 and Sham Kakade. How Does Critical Batch Size Scale in Pre-training?, November 2024. URL
683 <http://arxiv.org/abs/2410.21676>. arXiv:2410.21676 [cs].
- 684 Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Quadratic models
685 for understanding catapult dynamics of neural networks, May 2024. URL <http://arxiv.org/abs/2205.11787>. arXiv:2205.11787 [cs].
- 686 Xingyu Zhu, Zixuan Wang, Xiang Wang, Mo Zhou, and Rong Ge. UNDERSTANDING EDGE-
687 OF-STABILITY TRAINING DYNAMICS WITH A MINIMALIST EXAMPLE, 2023.

702 Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. The Anisotropic Noise in Stochastic
703 Gradient Descent: Its Behavior of Escaping from Sharp Minima and Regularization Effects, June
704 2019. URL <http://arxiv.org/abs/1803.00195>. arXiv:1803.00195 [cs, stat].
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756	CONTENTS	
757		
758	1 Introduction	1
759		
760	2 Related Work	2
761		
762	3 Preliminaries: Noise-Driven vs Curvature-Driven	4
763		
764	3.1 <i>Type-1</i> Oscillation and GNI	4
765		
766	3.2 SGD Always Oscillates <i>Type-1</i>	5
767		
768	3.3 <i>Type-2</i> : Curvature-Driven Oscillation	5
769		
770	4 SGD Typically Occurs at the EOSS	6
771		
772	4.1 <i>Batch Sharpness</i> Governs EOSS	6
773		
774	5 On Stability	8
775		
776	5.1 <i>Type-2</i> is about <i>Batch Sharpness</i>	8
777		
778	5.2 Connecting <i>Batch Sharpness</i> with Earlier Notions	8
779		
780	6 Implications: How Noise-Injected GD Differs from SGD	9
781		
782	7 Conclusions, Limitations, and Future Work	9
783		
784	A Acknowledgement of LLMs usage	17
785		
786	B Theory Preliminaries: A Framework for Instability	17
787		
788	B.1 Defining Instability	17
789		
790	B.2 Identifying a Regime of Instability	18
791		
792	B.3 Guiding Questions for EOSS	19
793		
794	C Comparison with Previous Empirical Work	20
795		
796	D On the Two Types of Oscillations in SGD Dynamics	22
797		
798	D.1 A Minimalistic Quadratic Example.	23
799		
800	D.2 Proof of Lemma 1	23
801		
802	E Proof of Proposition 1	24
803		
804	E.1 Setup and notation for Proposition 1	24
805		
806	E.2 Formal Version of Proposition 1	26
807		
808	E.3 Proof of Proposition 3	27
809		
	F On the Two Types of Oscillations in NNs	30
	F.1 On the Importance of <i>Type-2</i> Oscillations Compared to <i>Type-1</i>	31
	G Gradients explode above the EOSS: Proof of Theorem 1	32
	G.1 Part 1: Explosion	32

810	G.2 Part 2: SGD with Replacement	33
811	G.3 Part 2 for SGD <i>Without</i> Replacement	34
812		
813		
814	H Proof of the Equivalence of Section B.2	36
815	H.1 Setting and basic assumption	36
816	H.2 Breaking a valid instability criterion \iff catapult	36
817	H.3 Catapults \iff loss spikes of sufficient size	37
818	H.4 Putting the pieces together	38
819		
820		
821		
822	I On the fate of λ_{\max}	39
823	I.1 Empirical facts	39
824	I.2 Critical Batch Size	40
825	I.3 Why $2/\eta - C/b^\alpha$ fails.	40
826	I.4 Conclusion & Outlook: Why Path-Dependence Matters	41
827	I.5 Implications and Open Questions	42
828		
829		
830		
831	J Implications: How Noise-Injected GD Differs from SGD	44
832	J.1 Noisy GD	44
833	J.1.1 Noisy GD with Anisotropic Noise (Gaussian Resampling)	44
834	J.1.2 Noisy GD with Diagonal Noise	45
835	J.1.3 Noisy GD with Isotropic Noise	45
836	J.2 SDE	45
837		
838		
839		
840	K Alignment	47
841		
842	L Linear Stochastic Stability	50
843	L.1 Notions of Linear Stability	50
844	L.2 Empirical Behavior of Wu et al. (2018) Criterion	50
845	L.3 Implications	51
846		
847		
848	M Other Quantities of SGD Dynamics	53
849	M.1 Trace of the Loss Hessian	53
850	M.2 λ_{\max}^b : Expected Highest Eigenvalue of Mini-Batch Hessians	54
851	M.3 Modified Batch Sharpness	56
852		
853		
854		
855	N Modified Batch Sharpness is a Valid Instability Criterion	56
856		
857	O Hardware & Compute Requirements	59
858		
859	P The Hessian and the Fisher Information Matrix Overlap	59
860		
861	Q Illustration of EOSS in Variety of Settings: <i>Batch Sharpness</i>	62
862		
863	R Illustration of EOSS for the SVHN dataset	67

S Illustration of EOSS in Variety of Settings: λ_{\max}^b

70

A ACKNOWLEDGEMENT OF LLMs USAGE

We acknowledge the use of DeepSeek, Claude Code, Codex, and ChatGPT for code assistance. We used ChatGPT and Claude for text editing suggestions, proof-reading, and LaTeX editing help.

B THEORY PRELIMINARIES: A FRAMEWORK FOR INSTABILITY

In this section, we give a minimal formal framework for the way we deal with (*in*)stability in the rest of the paper. Our goal is to isolate (i) *instability criteria*—sufficient certificates for divergence of a (stochastic) dynamical system—and (ii) ways to recognize the regime where such a criterion *empirically saturates*, which we call the Edge of (Stochastic) Stability (EO(S)S). We explain why, on the quadratic approximation of the loss, the following three viewpoints are equivalent for us:

- breaking an instability criterion,
- observing a *catapult* (a large excursion out of the approximately quadratic region),
- and observing a spike of appropriate size in the loss.

This justifies why, in experiments, we both perturb hyperparameters and look for loss spikes to diagnose EOSS and what we consider catapults.

B.1 DEFINING INSTABILITY

On the notion of stability. Different notions of stability are in principle possible for training algorithms. If for every data point z_i there exist $\mathcal{H}_i \in \mathbb{R}^{d \times d}$ and $x_i \in \mathbb{R}^d$ such that the loss on z_i is $L_i(\theta) = \frac{1}{2}(\theta - x_i)^\top \mathcal{H}_i(\theta - x_i)$, we can consider the evolution of, e.g., the squared distance to a solution θ^* , or more generally any Lyapunov function $V(\theta)$. Typical stability requirements have the form

$$\mathbb{E}[V(\theta_{t+1})|\theta_t] \leq V(\theta_t), \quad \mathbb{E} \left[\frac{V(\theta_{t+1})}{V(\theta_t)} \mid \theta_t \right]^5 \leq 1, \quad \text{or} \quad \lim_{t \rightarrow \infty} \frac{1}{t} \log(V(\theta_t)) \leq 0, \quad \text{etc.} \quad (6)$$

In all these cases one can expand θ_{t+1} in a Taylor series around θ_t and obtain an equivalent inequality that can be expressed in terms of a *scalar quantity built from the Hessians*—generally termed *notion of curvature* (for instance, an eigenvalue, an operator norm, or a combination of low-order cumulants of the distribution of $\{\mathcal{H}_i\}_i$). As an example, non-expansion of the second moment of the parameters near a minimum—*linear stability*—was studied by Wu et al. (2018); Ma & Ying (2021); Mulayoff & Michaeli (2024) and boils down to a bound on an operator of the form

$$\|\mathbb{E}_{B \sim \mathcal{P}_b} [(I - \eta H(L_B))^{\otimes 2}]\| \leq 1. \quad (7)$$

We refer to Appendix L for a detailed comparison to linear stochastic stability.

Criteria for instability. We now isolate the abstract notion of stability, or better, certificate of instability, that will be used throughout the paper. Crucially, we define an instability criterion only through *sufficiency* for divergence. We do not require necessity for divergence (equivalently, sufficiency for stability), which is the type of condition typically required for convergence proofs and is, for example, the nature of the condition in Wu et al. (2018).

Definition 4 (Instability criterion). Consider a training algorithm (a discrete-time dynamical system) $(\theta_t)_{t \geq 0}$ on a parameter space $\Theta \subseteq \mathbb{R}^d$ with fixed hyperparameters h (e.g. learning rate, batch size). Let $U \subseteq \Theta$ be an open set (typically, a region where a given local approximation of the loss is trusted). Let $f : U \rightarrow \mathbb{R}$ and $c \in \mathbb{R}$. We say that f is a *valid instability criterion with threshold c* for the algorithm on U if the following holds:

$$f(\theta_0) > c \implies (\theta_t)_{t \geq 0} \text{ leaves every compact subset of } U.$$

That is, for any compact $K \subset U$ containing θ_0 , there exists a finite time T such that $\theta_T \notin K$. We say that the instability criterion f is *saturated* at θ if $f(\theta)$ is (approximately) equal to c ; in practice, up to an $O(\eta)$ tolerance.

⁵This is the one we find in our empirical work, with $V(\theta) = \|\nabla L_B(\theta)\|^2$.

In words: an instability criterion is a scalar quantity f together with a threshold c such that crossing $f > c$ is *sufficient* to force divergence from the region U we trust as a local model. For a specific f we generally want the *smallest* such c , which depends both on f and on the underlying dynamical system (algorithm, data, loss, architecture). A canonical example is full-batch EOS, where $U = \mathbb{R}^d$, $f(\theta) = \lambda_{\max}(\nabla^2 L(\theta))$ and $c = 2/\eta$: crossing $\lambda_{\max} > 2/\eta$ makes GD linearly unstable on any compact set.

The question if an optimizer acts at the *Edge of (Stochastic) Stability* thus becomes:
Are there criteria of instability, e.g., Eqs. (6, 7), that empirically saturate during training?

Quadratic and deterministic case. To connect what above with Cohen et al. (2021), in the case in which we use a full-batch method (deterministic) over a quadratic loss L , the criteria of instability obtainable from a Lyapunov function as in Eq. (6) simplify up to higher order in η as

- (i) only full-batch quantities appear in the Taylor expansion, no higher cumulants of the Hessians;
- (ii) shortly after reaching instability, the gradient aligns with the eigenvector of the highest eigenvalue λ_{\max} of the Hessian.

In particular, GD is an asymptotic stable (and Lyapunov stable) linear dynamical system if and only if $\lambda_{\max} \leq 2/\eta$. Cohen et al. (2021) empirically showed GD trains neural networks in the regime in which the inequality $\lambda_{\max} \leq 2/\eta$ is saturated, Damian et al. (2023) showed this mathematically under the assumption of progressive sharpening, making precise that saturation means that $\lambda_{\max} - 2/\eta \in [-c\eta|\log(\eta)|, +c\eta|\log(\eta)|]$ for some constant c .

B.2 IDENTIFYING A REGIME OF INSTABILITY

But how do we find if such an instability criterion is saturated? In practice, we diagnose this through the observation of *catapults* and *loss spikes* when the hyperparameters are perturbed. In this subsection we formalize this connection on the local quadratic approximation.

On the Taylor approximation of the loss, the underlying property is *divergence of the dynamics on U_t* , and the three diagnostics introduced—which we use throughout the article—are just different ways to detect it.

The three phenomena below are *equivalent manifestations of divergence on the quadratic approximation* rather than three logically independent assumptions:
 Breaking an instability criterion \iff Catapult (Def. 5) \iff Loss spike of sufficient size.

This equivalence—which we prove in Appendix H, in this precise sense, explains why, in practice, we use both catapults and loss spikes as our main empirical signatures of instability at the EO(S)S.

Catapults on the quadratic approximation. Fix a time t and a point θ_t . Assume that for each i the per-sample loss L_i is twice differentiable with β -Lipschitz Hessian in a neighborhood of θ_t , and let $\tilde{L}_i(\theta) := \frac{1}{2}(\theta - x_i)^\top \mathcal{H}_i(\theta - x_i)$ be the second-order Taylor approximation of L_i at θ_t . Let U_t denote an open neighborhood of θ_t where the quadratic approximation is accurate, e.g. where $|\mathcal{H}(L_i(\theta)) - \mathcal{H}_i| \leq C\eta$ for all $\theta \in U_t$ and all i and some constant $C > 0$ depending on β and the local geometry. Consider SGD (or GD) run on the quadratic model $\tilde{L}(\theta) := \frac{1}{N} \sum_i \tilde{L}_i(\theta)$ with the same hyperparameters as the original dynamics.

Definition 5 (Catapults on the quadratic model). We say that the algorithm *experiences a catapult at time t* if, when run on \tilde{L} from initialization θ_t , the resulting trajectory $(\theta_s)_{s \geq t}$ leaves every compact subset of U_t in finite time.

This definition is deliberately phrased in the same language as Definition 4: a catapult is precisely a divergence event for the quadratic dynamics on the region where the approximation is trusted.

Hyperparameter perturbations and saturation. Let f be a valid instability criterion as in Definition 4, with threshold c . In practice, f and c depend on hyperparameters h (e.g. $f(\theta; \eta, b)$ and $c = c(\eta, b)$). We are interested in settings where:

- 972 (a) $f(\cdot; h)$ is a valid instability criterion for the quadratic model on U_t ;
 973 (b) $f(\theta_t; h)$ is *monotone* in some destabilizing direction in h (e.g. increasing η or decreasing b
 974 increases $f(\theta_t; h) - c(h)$);
 975 (c) along the observed trajectory at hyperparameters h_0 , $f(\theta_t; h_0)$ *saturates*, i.e. $f(\theta_t; h_0) \approx$
 976 $c(h_0)$.

977 **Lemma 2** (Tight instability criterion \Rightarrow catapult under perturbation). *Assume h_0 saturates f at*
 978 *θ_t . Under assumptions (a)–(c) above, any sufficiently small destabilizing perturbation of h_0 (e.g.*
 979 *$\eta \uparrow$ or $b \downarrow$) produces hyperparameters h such that $f(\theta_t; h) > c(h)$. By validity of the criterion,*
 980 *the quadratic-model trajectory from θ_t leaves every compact subset of U_t in finite time, i.e., the*
 981 *algorithm experiences a catapult at time t in the sense of Definition 5.*

982 Thus a instability criterion that is monotone in a hyperparameter gives a concrete way to test whether
 983 we are at the EO(S)S: if it saturates along training, then a small destabilizing perturbation must
 984 trigger a catapult on the quadratic model. Conversely, if a quantity empirically saturates but small
 985 destabilizing perturbations do *not* lead to catapults, this quantity cannot be a valid instability criterion
 986 for the dynamics of interest.
 987

988 **From catapults to loss spikes.** On the quadratic model \tilde{L} and within U_t , the loss is a quadratic
 989 form in $\theta - \theta_t$. If a catapult occurs at time t , the quadratic trajectory leaves every compact subset of
 990 U_t , so $\|\theta_s - \theta_t\|$ eventually exceeds any fixed radius R for which $B_R(\theta_t) \subset U_t$. In particular, before
 991 (or as) the iterate exits such a ball, $\tilde{L}(\theta_s)$ must increase by at least a fixed factor $\mathcal{O}_\eta(1)$ compared to
 992 $\tilde{L}(\theta_t)$. We refer to such an increase as a *loss spike of sufficient size* on the quadratic model⁶.

993 Conversely, a loss spike of sufficiently large ($\mathcal{O}_\eta(1)$) relative size on the quadratic model cannot
 994 occur if the dynamics remains linearly stable and confined to a fixed compact subset of U_t : in that
 995 regime, the quadratic dynamics is a contraction in a suitable norm, and both $\|\theta_s - \theta_t\|$ and $\tilde{L}(\theta_s)$
 996 remain uniformly controlled. Thus such a spike implies that the quadratic dynamics is divergent
 997 on U_t in the sense above, and therefore that some valid instability criterion for the quadratic model
 998 must be broken along the trajectory.
 999

1000 B.3 GUIDING QUESTIONS FOR EOSS

1001 Understanding *if and how* SGD trains at the EOSS is thus inherently linked to:

- 1002
- 1003 • what scalar Hessian-based statistic (if any) saturates as a valid instability criterion (Q1);
 - 1004 • how progressive sharpening and self-stabilization act on the moments/cumulants of the
 - 1005 Hessian distribution and thereby select that statistic (Q2);
 - 1006 • and whether this statistic is computable and usable in high-dimensional practice (Q3).
- 1007

1008 The rest of the paper answers these questions, at least partially: we will show that Batch Sharpness
 1009 is a valid and tractable instability criterion for SGD on the quadratic approximation, that it empirically
 1010 saturates in the practice of neural network training. Importantly, *Batch Sharpness* is the first
 1011 proposed such quantity which does not fail any of these desiderata.
 1012

1013 **Stochastic and non-quadratic case.** The discussion above shows that for deterministic gradient
 1014 descent on a quadratic loss, all the usual stability notions collapse to the same scalar quantity

$$1015 f(\theta) = \lambda_{\max}(\nabla^2 L(\theta)), \quad f(\theta) \leq 2/\eta$$

1016 and training at the edge of stability means that this inequality saturates. In contrast, for *stochastic*
 1017 training on a *non-quadratic* loss, the situation is substantially more delicate: different Lyapunov
 1018 functions lead to different scalar quantities built from the Hessians (e.g. different moments or cumu-
 1019 lants over $\{\mathcal{H}_i\}_i$), and there is no reason a priori for them to agree or to saturate at the same values.
 1020 Note for instance, that in (7) the first two cumulants appear, while in *Batch Sharpness* the first three.
 1021 We organize the rest of the paper around three guiding questions.
 1022

1023 ⁶We are imprecise on purpose here, because the size required for the spike to be a catapult depends in
 1024 practice on the *global* geometry. As we see in Corollary 1 and Figure 3 for quadratics, when you perturb a
 1025 stable step size to a new stable step size, the loss jumps by $\mathcal{O}(\eta/(2 - \eta\lambda_{\max}))$ before restabilizing, that size is
 not a catapult spike.

(Q1) Existence and (non-)uniqueness of a criterion. In principle, there may exist infinitely many instability criteria that saturate along the SGD trajectory.⁷ This leads to the first question:

Is there a “distinguished” scalar Hessian-based quantity that both (i) is a valid instability criterion in the sense of Definition 4 and (ii) empirically saturates at the EOSS for SGD?

Even if such a quantity exists, it need not be unique in principle; part of our contribution is to show that *Batch Sharpness* is one such quantity and to argue that it is preferred over several alternatives (e.g. GNI).

(Q2) Mechanism: progressive sharpening and self-stabilization. Different candidate criteria depend on different statistics of the Hessian distribution. For example, the operator in (7) only involves the first and second cumulants of $\{\mathcal{H}_i\}_i$, whereas *Batch Sharpness* (even in the quadratic case) also depends on the third, and λ_{\max} depends only on the first. Which of these is actually driven to saturation by the dynamics depends on how *Progressive Sharpening* and *self-stabilization* act on these statistics:

- progressive sharpening pushes up certain Hessian statistics during training;
- self-stabilization constrains them by preventing blow-up.

This leads to the second question:

Which scalar Hessian-based statistic does the combination of progressive sharpening and self-stabilization actually push to a saturated instability threshold in practice?

If, for instance, progressive sharpening mostly increases the third cumulant while leaving the first two relatively unchanged, we should expect *Batch Sharpness* (which depends on the third cumulant) to be the relevant instability criterion, rather than quantities that only see the first or second cumulant.

(Q3) Computability and usefulness in high dimension. Finally, some natural instability criteria are not efficiently estimable in high dimension. For instance, conditions based directly on operators like (7) or the exact instability criterion of [Mulayoff & Michaeli \(2024\)](#) live in a d^2 -dimensional space and are infeasible to evaluate in modern neural networks. This motivates the third question:

Is there a valid instability criterion for SGD on neural networks that is both empirically saturating and computationally tractable in high dimension?

In practice, this means looking for criteria that can be estimated in roughly $\mathcal{O}(d \cdot \text{poly}(\log(d)))$ time.

C COMPARISON WITH PREVIOUS EMPIRICAL WORK

[Lee & Jang \(2023\)](#) introduce several quantities crucial for understanding neural network training dynamics. Below, we discuss the relationships among λ_{\max} , *Batch Sharpness*, and *Interaction-Aware Sharpness* (IAS, [Lee & Jang \(2023\)](#)), emphasizing that a comprehensive theory of mini-batch dynamics should explain their distinct plateau timings and interconnected behaviors. We conjecture that a complete theory of stochastic gradient descent (SGD) dynamics would elucidate these metrics’ precise interrelations and their different plateau timings.

Interaction-Aware Sharpness. [Lee & Jang \(2023\)](#) introduce *Interaction-Aware Sharpness* (IAS), denoted $\|\mathcal{H}\|_{S_b}$:

$$\|\mathcal{H}\|_{S_b} := \frac{\mathbb{E}_{B \sim \mathcal{P}_b} [\nabla L_B(\mathbf{x})^\top \mathcal{H} \nabla L_B(\mathbf{x})]}{\mathbb{E}_{B \sim \mathcal{P}_b} [\|\nabla L_B\|^2]}.$$

This quantity shares structural similarities with both *Batch Sharpness* (Definition 3) and the *Gradient-Noise Interaction* (Proposition 1), differing from the latter only in the denominator. The

⁷For instance, in Section 4 we show that $\lambda_{\max} \leq 2/\eta$ does not saturate for batch sizes smaller than a problem-dependent critical batch size. It is an open problem to characterize the level at which λ_{\max} stabilizes in that regime.

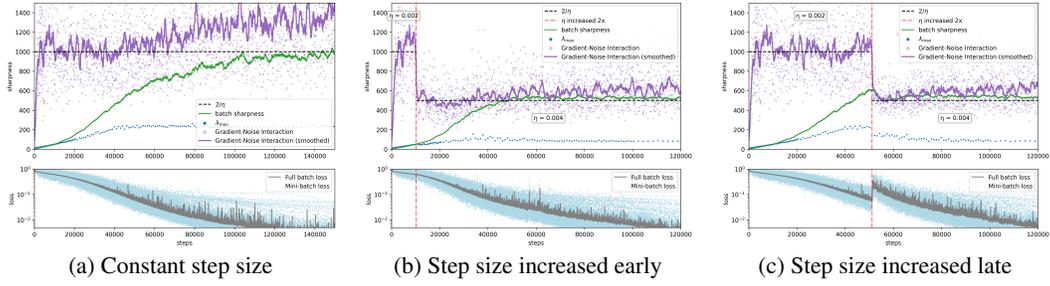


Figure 9: We demonstrate that the saturation of *GNI* does not govern a sharpness-related regime of instability typical of Type-2 oscillations - and in particular, highlighting the difference in the two types of oscillations. When we double the step size after *batch sharpness* is at least half of $2/\eta$ threshold (so that it is beyond the new $2/\eta$ level), training exhibits a catapult surge in the loss (c). But if we make the same change *before* *batch sharpness* crosses that level—despite *GNI* already saturating—no catapult occurs. (b)

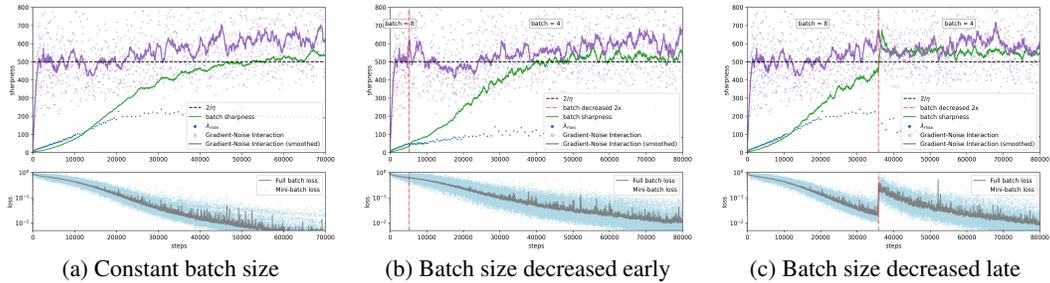


Figure 10: Similarly, reducing the batch size only triggers catapults if *batch sharpness*, not *GNI*, exceeds the threshold.

key distinction from *Batch Sharpness* lies in which Hessian is evaluated: IAS measures the directional curvature of the **full-batch** loss landscape L along mini-batch gradient directions, while *Batch Sharpness* measures the directional curvature of **mini-batch** loss landscape L_B along their corresponding gradients. This distinction is crucial, as mini-batch Hessians vary with batch selection while the full-batch Hessian remains fixed.

Notably, with full-batch GD, IAS serves as a directional alternative to the maximal Hessian eigenvalue, λ_{\max} , introduced by Cohen et al. (2021). IAS aligns closely with the $2/\eta$ threshold, unlike λ_{\max} , which often remains slightly above this threshold during EOS, especially at the beginning of it. Since IAS measures *directional* curvature, we have $\|\mathcal{H}\|_{S_n} \leq \lambda_{\max}$. Consequently, in the mini-batch setting, IAS stabilizes below $2/\eta$, consistent with empirical observations from Jastrzebski et al. (2019; 2020); Cohen et al. (2021) and our Figure 2. Notably, when $B = n$, our *Batch Sharpness* coincides with IAS rather than λ_{\max} , reinforcing the interpretation of *Batch Sharpness* as the relevant metric stabilizing at $2/\eta$ even under full-batch conditions.

Relation to Gradient-Noise Interaction. Another metric from Lee & Jang (2023) is defined as:

$$\frac{\text{tr}(HS_b)}{\text{tr}(S_n)} = \frac{\mathbb{E}_{B \sim \mathcal{P}_b} [\nabla L_B(\mathbf{x})^\top \mathcal{H} \nabla L_B(\mathbf{x})]}{\|\nabla L\|^2}$$

which coincides exactly with our definition of *GNI* (Proposition 1). As detailed in Section 3 and Appendix D, the stabilization of *GNI* around $2/\eta$ signals the presence of oscillations, at least Type-1 oscillations. Lee & Jang (2023) provide extensive empirical evidence demonstrating that neural networks spend much of their training within this oscillatory regime (see also Figures 9a and 10a). This contrasts traditional theoretical analyses (Bottou et al. (2018); Mandt et al. (2016)), which consider oscillations only near the manifold of minima.

Distinguishing oscillation types. It is crucial to note that GNI around $2/\eta$ does not inherently indicate instability. As clarified in Sections 3, 4 and Appendix E, not all oscillations are inherently unstable. Figures 5, 9b, 10b illustrate that altering hyperparameters when GNI is around $2/\eta$ typically does not trigger instability (catapult-like divergence), contrary to expectations if the system was in an EOS-like regime of instability. Instead, as shown in Figures 5, 9c, 10c, *Batch Sharpness* more reliably predicts a regime of instability. Additionally, Figure 11 highlights GNI's independence from progressive sharpening, a necessary precursor to Type-2 (curvature-driven) oscillations and EOS-like instabilities, as detailed in Appendix F.

Missing Progressive Sharpening. Extensively, both in our experiments and in the ones of Lee & Jang (2023), GNI grows to $2/\eta$ in a few initial steps (and sometimes from the very beginning if the initialization size is large) without ever being in subject to a phase of progressive sharpening unlike *Batch Sharpness* and λ_{\max} . The phase of growth of GNI is generally short and independent of the size, the behavior, and the phase in which *Batch Sharpness* and λ_{\max} are.

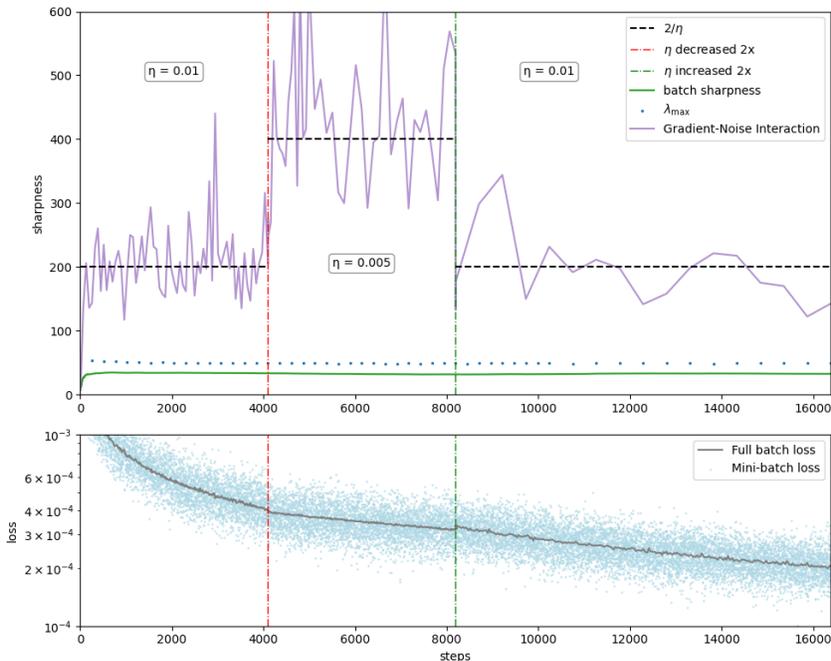


Figure 11: We construct a 32k-point "easy" CIFAR-10, where we "pull apart" all the 10 classes, so the classes become linearly separable. In this case, there is virtually no "learning" to be done, and therefore, there is barely any progressive sharpening happening (as established Cohen et al. (2021), progressive sharpening does not happen if the dataset "is not complex enough"). Yet, *GNI* still stabilizes at the initial level of $2/\eta$. More importantly, when we decrease and then increase the step size, the *GNI* measure restabilizes to the corresponding new thresholds, while λ_{\max} does not change. That means that *GNI* is independent of the curvature of the loss landscape and is unrelated to progressive sharpening, and thus Type-2 oscillations and EOS-like instability regimes.

D ON THE TWO TYPES OF OSCILLATIONS IN SGD DYNAMICS

A fundamental challenge in analyzing SGD compared to GD stems from the inherent oscillations induced by mini-batch gradient noise. This appendix, together with Appendix C (also see proofs in E and G), extends the discussion in Section 3 by formally distinguishing between two distinct types of oscillations: noise-driven (Type-1) and curvature-driven (Type-2). This distinction is crucial because Type-1 oscillations occur independently of the loss landscape's curvature and thus do not exert a regularizing effect on the sharpness of the final solution. In contrast, Type-2 oscillations are directly caused by landscape curvature and induce an implicit regularization effect by discouraging convergence towards sharp minima.

We begin with a minimalistic example to illustrate the nature of Type-1

D.1 A MINIMALISTIC QUADRATIC EXAMPLE.

Here we show mathematically what we see empirically in Figure 3 (the simplified version—only two data points). Consider a regression problem with two datapoints, 1 and -1 , and a linear model $f(x) = x$ under the quadratic loss. The (scaled) full-batch loss is given by:

$$L(x) = \frac{1}{4}(x-1)^2 + \frac{1}{4}(x+1)^2.$$

Batch-1 SGD updates with step-size $0 < \eta < 2$ result in oscillatory behavior around the optimum $x = 0$ due entirely to gradient noise, with amplitude approximately $\sqrt{\frac{\eta}{2-\eta}}$. Crucially, the Hessian in this example is small ($\frac{d^2L}{dx^2} = 1$), demonstrating that these persistent oscillations are entirely noise-driven (Type-1).

Formally, the SGD update is:

$$x_{t+1} = x_t - \eta \nabla \ell_{i_t} = (1 - \eta)x_t + \eta \xi_t$$

where ℓ_{i_t} s are the individual datapoint losses, and ξ_t s are i.i.d Rademacher random variables. Thus, we obtain the first two moments explicitly:

$$\mathbb{E}[x_t] = (1 - \eta)\mathbb{E}[x_{t-1}] = (1 - \eta)^t x_0$$

$$\mathbb{E}[x_t^2] = (1 - \eta)^2 \mathbb{E}[x_{t-1}^2] + \eta^2 = (1 - \eta)^{2t} x_0^2 + \frac{\eta^2}{1 - (1 - \eta)^2} (1 - (1 - \eta)^{2t})$$

This implies convergence in expectation for $0 < \eta < 2$, with a limiting variance given by:

$$\lim_{t \rightarrow \infty} \mathbb{E}[x_t^2] = \frac{\eta}{2 - \eta}$$

and divergence for $\eta > 2$.

A key observation is that increasing η to any value $\eta_1 < 2$ merely changes the amplitude of oscillations to $\sqrt{\frac{\eta_1}{2-\eta_1}}$ without triggering any catapult-like behavior. The only step size for which we observe Type-2 (curvature-driven) oscillations and an EOS-like⁸ instability is precisely $\eta = 2$, where the dynamics effectively become a random walk, and any larger step size leads to divergence.

Crucially, when $\eta < 2$ oscillations occur persistently on the full-batch loss, despite the individual steps on the mini-batch loss remaining stable.

The oscillation is due to the fact that the mini-batch loss landscape shifts from step to step, not to the fact that the steps are unstable.

D.2 PROOF OF LEMMA 1

We propose here the formal version of Lemma 1.

Proposition 2 (Loss increment and GNI). *Assume L is three times continuously differentiable and its Hessian is L_2 -Lipschitz in a neighborhood of θ . Then for any mini-batch B and step size $\eta > 0$ small enough we have*

$$\mathbb{E}_B [L(\theta - \eta \nabla L_B(\theta)) - L(\theta) \mid \theta] = -\eta \|\nabla L(\theta)\|^2 + \frac{\eta^2}{2} \mathbb{E}_B [\nabla L_B(\theta)^\top \mathcal{H}(\theta) \nabla L_B(\theta)] + \mathcal{O}(\eta^3), \quad (8)$$

⁸The key difference between these oscillations and genuine EOS behavior in neural networks is that, in the quadratic case, the full-batch loss does not decrease, making this scenario inherently less informative. In contrast, neural networks exhibit a surprising, albeit non-monotonic, decrease in loss within this instability regime, an effect arising from the multidimensional nature of their optimization landscape (Damian et al., 2023)

where the $\mathcal{O}(\eta^3)$ constant depends only on L_2 and an upper bound on $\|\nabla L_B(\theta)\|$. Equivalently,

$$\mathbb{E}_B[L(\theta_{t+1}) - L(\theta_t) \mid \theta_t] = -\eta\|\nabla L(\theta_t)\|^2 \left(1 - \frac{\eta}{2} \text{GNI}(\theta_t) + \mathcal{O}(\eta^2)\right), \quad (9)$$

with

$$\text{GNI}(\theta) := \frac{\mathbb{E}_B[\nabla L_B(\theta)^\top \mathcal{H}(\theta) \nabla L_B(\theta)]}{\|\nabla L(\theta)\|^2}.$$

In particular, there exists $c > 0$ such that for all sufficiently small η : if

$$|\text{GNI}(\theta_t) - 2/\eta| \geq c\eta,$$

then the sign of the expected loss increment satisfies

$$\text{sign } \mathbb{E}[L(\theta_{t+1}) - L(\theta_t) \mid \theta_t] = \text{sign}(2/\eta - \text{GNI}(\theta_t)).$$

Proof of Proposition 2. Fix θ and a mini-batch B . Write a third-order Taylor expansion of L around θ :

$$L(\theta - \eta\nabla L_B(\theta)) = L(\theta) - \eta\nabla L(\theta)^\top \nabla L_B(\theta) + \frac{\eta^2}{2} \nabla L_B(\theta)^\top \mathcal{H}(\theta) \nabla L_B(\theta) + R(\eta, \nabla L_B(\theta)),$$

where $\nabla L_B(\theta) := \nabla L_B(\theta)$ and $\mathcal{H}(\theta)$ is the full-batch Hessian. The third-order remainder satisfies the standard bound $|R(\eta, \nabla L_B(\theta))| \leq CL_2\eta^3\|\nabla L_B(\theta)\|^3$ for some numerical C , because the Hessian is L_2 -Lipschitz.

Now take the expectation over $B \sim \mathcal{P}_b$. Using $\mathbb{E}_B[\nabla L_B(\theta)] = \nabla L(\theta)$ we obtain

$$\mathbb{E}_B[L(\theta - \eta\nabla L_B(\theta)) - L(\theta) \mid \theta] = -\eta\|\nabla L(\theta)\|^2 + \frac{\eta^2}{2} \mathbb{E}_B[\nabla L_B(\theta)^\top \mathcal{H}(\theta) \nabla L_B(\theta)] + \mathcal{O}(\eta^3),$$

which is (8).

Divide the right-hand side by $-\eta\|\nabla L(\theta)\|^2$ to get

$$\mathbb{E}_B[L(\theta_{t+1}) - L(\theta_t) \mid \theta_t] = -\eta\|\nabla L(\theta_t)\|^2 \left(1 - \frac{\eta}{2} \text{GNI}(\theta_t) + \mathcal{O}(\eta^2)\right),$$

since

$$\text{GNI}(\theta_t) = \frac{\mathbb{E}_B[\nabla L_B(\theta)^\top \mathcal{H}(\theta) \nabla L_B(\theta)]}{\|\nabla L(\theta_t)\|^2}.$$

The last claim (sign agreement) follows immediately: for sufficiently small η , the $\mathcal{O}(\eta^2)$ term is dominated whenever $|\text{GNI} - 2/\eta| \geq c\eta$ for a fixed constant $c > 0$. \square

E PROOF OF PROPOSITION 1

E.1 SETUP AND NOTATION FOR PROPOSITION 1

Let

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \ell_i(\theta)$$

be three-times continuously differentiable, and let θ^* be a (possibly non-isolated) local minimiser. Denote the full-batch Hessian at θ^* by

$$\mathcal{H} := \nabla^2 L(\theta^*) \succeq 0.$$

For each sample i , define its (local) Hessian at θ^* ,

$$\mathcal{H}_i := \nabla^2 \ell_i(\theta^*),$$

so that $\mathcal{H} = \frac{1}{n} \sum_{i=1}^n \mathcal{H}_i$. We also define the (single-sample) gradient noise covariance at θ^* as

$$\Sigma_g := \mathbb{E}_i[\nabla \ell_i(\theta^*) \nabla \ell_i(\theta^*)^\top].$$

1296 We decompose the parameter space as

$$1297 \mathbb{R}^d = E_+ \oplus E_0, \quad E_+ := \text{Im}(\mathcal{H}), \quad E_0 := \ker(\mathcal{H}),$$

1299 with associated orthogonal projectors P_+ and P_0 . We will only require control of the dynamics in
1300 E_+ and assume that gradient noise in the flat subspace E_0 is not too large.

1301 For each iteration t , a mini-batch B_t of size b is drawn (with or without replacement) and the SGD
1302 update is

$$1303 \theta_{t+1} = \theta_t - \eta \nabla L_{B_t}(\theta_t), \quad \nabla L_{B_t}(\theta) := \frac{1}{b} \sum_{i \in B_t} \nabla \ell_i(\theta).$$

1306 Finally, define the Kronecker–sum operator

$$1307 \mathcal{K} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}, \quad \mathcal{K}(X) := \mathcal{H}X + X\mathcal{H}.$$

1309 On $E_+ \otimes E_+$, \mathcal{K} is positive definite and has a Moore–Penrose pseudoinverse \mathcal{K}^\dagger .

1311 We work under the following assumptions in a neighbourhood of θ^* .

1312 **(A1) Local quadratic approximation.** Each ℓ_i is twice differentiable with L_2 –Lipschitz Hessian
1313 near θ^* , and admits the Taylor expansion

$$1314 \nabla \ell_i(\theta) = \nabla \ell_i(\theta^*) + \mathcal{H}_i(\theta - \theta^*) + R_i(\theta),$$

1316 where $\|R_i(\theta)\| = \mathcal{O}(\|\theta - \theta^*\|^2)$ uniformly in i .

1317 **(A2) Compatible noise in flat directions.** The gradient noise covariance in the flat subspace is
1318 small:

$$1319 \|P_0 \Sigma_g P_0\| \lesssim \eta,$$

1320 so that the iterates do not perform an unbounded random walk along $\ker(\mathcal{H})$.

1322 **(A3) Linear stability of the linear dynamics.** The SGD on the quadratic approximation is lin-
1323 early stable on E_+ , i.e.

$$1324 \rho \mathbb{E}[(I - \eta \mathcal{H}(L_B))^{\otimes 2}]_{|E_+} < 1.$$

1325 *Remark 1* (Remarks on the assumptions).

1326 *Exact vs. Lipschitz Hessian (on (A1))* When each ℓ_i is strictly quadratic, the local linearity

$$1328 \nabla \ell_i(x) = \nabla \ell_i(x^*) + \mathcal{H}_i(x - x^*), \quad \mathcal{H}_i := \nabla^2 \ell_i(x^*),$$

1329 holds exactly, and $\mathcal{H} = \frac{1}{n} \sum_i \mathcal{H}_i$. In the general case, if $\nabla^2 \ell_i$ is L_2 –Lipschitz in a neigh-
1330 borhood of x^* , a second–order Taylor expansion gives a remainder $\mathcal{O}(\|x - x^*\|^2)$. For suffi-
1331 ciently small η , the SGD iterates typically remain in an $\mathcal{O}(\sqrt{\eta})$ –neighborhood of x^* , so these
1332 higher–order terms contribute only $\mathcal{O}(\eta^2)$ corrections in the discrete Lyapunov equation, which
1333 are dominated by the main $\mathcal{O}(\eta)$ term in Proposition 3.

1335 *Small drift in flat directions (on (A2))* The requirement $P_0 \Sigma_g P_0 = 0$ can be relaxed to
1336 $\|P_0 \Sigma_g P_0\| \leq \delta$. A standard discrete–Lyapunov analysis shows that the stationary covariance
1337 Σ_x remains finite provided $\delta = \mathcal{O}(\eta)$: roughly, if $\|P_0 \Sigma_g P_0\|$ is at most a constant multiple of
1338 η times the curvature scale on E_+ , then the null–space covariance Σ_x^{00} grows at most on the
1339 same $\mathcal{O}(\eta)$ scale as the covariance on $\text{Im}(\mathcal{H})$. If instead $P_0 \Sigma_g P_0$ is large, the dynamics ex-
1340 ecutes an (uncontrolled) random walk along E_0 , and no finite stationary covariance exists in
1341 those directions.

1342 *Linear stability and spectral gap (on (A3))* The condition

$$1343 \rho \left(\mathbb{E}[(I - \eta \mathcal{H}(L_B))^{\otimes 2}]_{|E_+} \right) < 1$$

1344 is the standard linear–stability condition for the second moment of SGD on a quadratic objective
1345 (cf. Ma & Ying (2021); Wu et al. (2018)). It ensures that the linearized error dynamics on E_+ is
1346 mean–square contractive and that the discrete Lyapunov equation

$$1347 \Sigma_x = \mathbb{E}[(I - \eta \mathcal{H}(L_B)) \Sigma_x (I - \eta \mathcal{H}(L_B))^{\top}] + \frac{\eta^2}{b} \Sigma_g$$

1348

1349

has a unique finite solution. As the spectral radius ρ approaches 1 from below, the spectral gap $1 - \rho$ controls both the mixing time and the size of the stationary covariance, with $\|\Sigma_x\| = \mathcal{O}((1 - \rho)^{-1})$. In our use of Proposition 3 and Corollary 1 we implicitly assume that the step size η is chosen so that the dynamics remains in this linearly stable regime on the time scales of interest, i.e. $\rho < 1$ (and often ρ bounded away from 1), so that the $\mathcal{O}(\eta)$ expansion and the GNI $\approx 2/\eta$ law are accurate.

E.2 FORMAL VERSION OF PROPOSITION 1

Proposition 3 (Gradient–Noise Interaction at a stable stationary regime). *Assume (A1) and (A2) above, and run mini-batch SGD with fixed batch size b and fixed step size η satisfying the linear stability condition (A3).*

Then the linearised error process $\Delta_t := \theta_t - \theta^$ admits a unique stationary covariance matrix Σ_x on E_+ , given by*

$$\Sigma_x = \frac{\eta}{b} \mathcal{K}^\dagger(\Sigma_g) + \mathcal{O}(\eta^2), \quad (10)$$

where the $\mathcal{O}(\eta^2)$ term depends only on L_2 , $\|\Sigma_g\|$ and $(\lambda_{\min}^+)^{-1}$.

Moreover, if $\theta \sim \pi$ is distributed according to this stationary law, and B is an independent fresh mini-batch of size b , then

$$\frac{\mathbb{E}_{\theta \sim \pi} \mathbb{E}_B [\nabla L_B(\theta)^\top \mathcal{H} \nabla L_B(\theta)]}{\mathbb{E}_{\theta \sim \pi} [\|\nabla L(\theta)\|^2]} = \frac{2}{\eta} (1 + \mathcal{O}(\eta)). \quad (11)$$

In particular, to leading order in η , the Gradient–Noise Interaction

$$\text{GNI}(\theta) := \frac{\mathbb{E}_B [\nabla L_B(\theta)^\top \mathcal{H} \nabla L_B(\theta)]}{\|\nabla L(\theta)\|^2}$$

is centred at $2/\eta$ under the (linearly) stable stationary distribution, and this leading behaviour is independent of the individual Hessians $\{\mathcal{H}_i\}_i$, depending on them only through Σ_g and \mathcal{H} .

Corollary 1 (Stable η –changes induce only bounded loss jumps in the Type-1 regime). *Assume (A1)–(A3) and let $\eta_0, \eta_1 > 0$ be two step sizes such that the linearized SGD dynamics on the quadratic approximation is linearly stable on E_+ for both η_0 and η_1 :*

$$\rho\left(\mathbb{E}[(I - \eta_k \mathcal{H}(L_B))^{\otimes 2}]_{|E_+}\right) < 1, \quad k \in \{0, 1\}.$$

Consider the SGD trajectory (θ_t) that is run with step size η_0 up to some time T , and with step size η_1 for all $t \geq T$.

Then:

- (i) *For each $k \in \{0, 1\}$ the linearized error process $\Delta_t^{(\eta_k)} := \theta_t - \theta^*$ admits a unique stationary covariance $\Sigma_x(\eta_k)$ on E_+ satisfying*

$$\Sigma_x(\eta_k) = \frac{\eta_k}{b} \mathcal{K}^\dagger(\Sigma_g) + \mathcal{O}(\eta_k^2),$$

as in Proposition 3.

- (ii) *Let $A_1 := \mathbb{E}[(I - \eta_1 \mathcal{H}(L_B))^{\otimes 2}]_{|E_+}$ and $\rho_1 := \rho(A_1) < 1$. There exist constants $C_1, C_2 < \infty$, depending only on Σ_g and ρ_1 (but not on T), such that for all $s \geq 0$,*

$$\mathbb{E}[\|\theta_{T+s} - \theta^*\|^2] \leq \frac{C_1}{1 - \rho_1}, \quad \mathbb{E}[L(\theta_{T+s}) - L(\theta^*)] \leq \frac{C_2}{1 - \rho_1}.$$

In particular, after switching from η_0 to η_1 the loss trajectory remains uniformly bounded and converges to the finite stationary level

$$L_\infty(\eta_1) := \mathbb{E}_{\theta \sim \pi_{\eta_1}} [L(\theta) - L(\theta^*)] = \frac{1}{2} \text{tr}(\mathcal{H} \Sigma_x(\eta_1) \mathcal{H}),$$

which itself is of order $\mathcal{O}((1 - \rho_1)^{-1})$.

Thus, in the Type-1 (noise-driven) regime, any change of step size that preserves linear stability (A3) can produce at most a finite “jump” in the expected loss, of the same order as the new stationary level $L_\infty(\eta_1)$, but cannot generate catapult-like divergence.

1404 E.3 PROOF OF PROPOSITION 3

1405 *Proof of Proposition 3.* We proceed in four steps. Throughout the proof we work on the subspace
 1406 $E_+ = \text{Im}(\mathcal{H})$; all covariances and operators are implicitly restricted to E_+ (the flat subspace E_0 is
 1407 controlled by Assumption (A2) and does not contribute to the quantities involving \mathcal{H}).
 1408

1409 **Step 1: Linearised dynamics on the quadratic approximation.**

1410 By Assumption (A1), near θ^* each per-sample loss ℓ_i admits the expansion

$$1411 \nabla \ell_i(\theta) = \nabla \ell_i(\theta^*) + \mathcal{H}_i(\theta - \theta^*) + R_i(\theta),$$

1412 where $\|R_i(\theta)\| = \mathcal{O}(\|\theta - \theta^*\|^2)$ uniformly in i . Let us denote the single-sample gradient at θ^* by

$$1413 g_i := \nabla \ell_i(\theta^*),$$

1414 so that $\Sigma_g = \mathbb{E}_i[g_i g_i^\top]$.

1415 For each iteration t , a mini-batch B_t of size b is drawn and the SGD update is

$$1416 \theta_{t+1} = \theta_t - \eta \nabla L_{B_t}(\theta_t), \quad \nabla L_{B_t}(\theta) := \frac{1}{b} \sum_{i \in B_t} \nabla \ell_i(\theta).$$

1417 Define the error vector

$$1418 \Delta_t := \theta_t - \theta^*.$$

1419 Then

$$1420 \begin{aligned} \nabla L_{B_t}(\theta_t) &= \frac{1}{b} \sum_{i \in B_t} [g_i + \mathcal{H}_i \Delta_t + R_i(\theta_t)] \\ &=: \underbrace{\xi_t}_{\text{zero mean}} + \underbrace{\mathcal{H}_{B_t}}_{\text{batch Hessian}} \Delta_t + r_t, \end{aligned}$$

1421 where

$$1422 \xi_t := \frac{1}{b} \sum_{i \in B_t} g_i, \quad \mathcal{H}_{B_t} := \frac{1}{b} \sum_{i \in B_t} \mathcal{H}_i, \quad r_t := \frac{1}{b} \sum_{i \in B_t} R_i(\theta_t).$$

1423 By construction we have

$$1424 \mathbb{E}[\xi_t] = 0, \quad \mathbb{E}[\xi_t \xi_t^\top] = \frac{1}{b} \Sigma_g,$$

1425 and (using $\nabla L(\theta^*) = 0$)

$$1426 \mathbb{E}[\mathcal{H}_{B_t}] = \mathcal{H}.$$

1427 The exact SGD recursion can therefore be written as

$$1428 \Delta_{t+1} = \Delta_t - \eta \nabla L_{B_t}(\theta_t) = (I - \eta \mathcal{H}_{B_t}) \Delta_t - \eta \xi_t - \eta r_t. \quad (12)$$

1429 For the purposes of the leading-order analysis, it is convenient to first *ignore* the nonlinear remainders r_t and consider the purely linearised dynamics on the quadratic approximation (i.e. we replace each ℓ_i by its quadratic Taylor polynomial at θ^*). On this quadratic model we have $R_i \equiv 0$, hence $r_t \equiv 0$, and (12) becomes

$$1430 \Delta_{t+1} = C_t \Delta_t - \eta \xi_t, \quad C_t := I - \eta \mathcal{H}_{B_t}. \quad (13)$$

1431 Note that in this linearised model the random matrices C_t and the noise vectors ξ_t are independent of Δ_t (they depend only on the batch B_t and the fixed Hessians $\{\mathcal{H}_i\}_i$).

1432 We will first solve the covariance structure of the linear recursion (13), and then argue that restoring the remainder r_t only introduces $\mathcal{O}(\eta^2)$ corrections.

1433 **Step 2: Discrete Lyapunov equation and existence of a stationary covariance.**

1434 Let $\Sigma_t := \mathbb{E}[\Delta_t \Delta_t^\top]$ denote the covariance of Δ_t under the linear recursion (13). Using independence of C_t and ξ_t from Δ_t , we compute

$$1435 \begin{aligned} \Sigma_{t+1} &= \mathbb{E}[\Delta_{t+1} \Delta_{t+1}^\top] \\ &= \mathbb{E}[(C_t \Delta_t - \eta \xi_t)(C_t \Delta_t - \eta \xi_t)^\top] \\ &= \mathbb{E}[C_t \Delta_t \Delta_t^\top C_t^\top] - \eta \mathbb{E}[C_t \Delta_t \xi_t^\top] - \eta \mathbb{E}[\xi_t \Delta_t^\top C_t^\top] + \eta^2 \mathbb{E}[\xi_t \xi_t^\top]. \end{aligned}$$

1458 Conditioning on Δ_t and using $\mathbb{E}[\xi_t | \Delta_t] = 0$, the two cross-terms vanish:

$$1459 \mathbb{E}[C_t \Delta_t \xi_t^\top] = \mathbb{E}[C_t \Delta_t \mathbb{E}[\xi_t^\top | \Delta_t]] = 0, \quad \mathbb{E}[\xi_t \Delta_t^\top C_t^\top] = 0.$$

1461 Thus we obtain

$$1462 \Sigma_{t+1} = \mathbb{E}[C_t \Sigma_t C_t^\top] + \eta^2 \frac{1}{b} \Sigma_g. \quad (14)$$

1464 Assuming that the linear recursion admits a stationary distribution on E_+ , we denote the stationary covariance by

$$1465 \Sigma_x := \lim_{t \rightarrow \infty} \Sigma_t$$

1467 and it must satisfy the discrete Lyapunov equation

$$1469 \Sigma_x = \mathbb{E}[C_t \Sigma_x C_t^\top] + \eta^2 \frac{1}{b} \Sigma_g. \quad (15)$$

1471 To show that such a Σ_x exists and is unique on E_+ , we vectorise (15). Recall that for any matrices A, X, B of compatible dimensions we have

$$1472 \text{vec}(AXB) = (B^\top \otimes A) \text{vec}(X),$$

1475 where \otimes denotes the Kronecker product. Applying this to $C_t \Sigma_x C_t^\top$ we get

$$1477 \text{vec}(C_t \Sigma_x C_t^\top) = (C_t \otimes C_t) \text{vec}(\Sigma_x).$$

1479 Taking expectations in (15) and using linearity of $\text{vec}(\cdot)$ we obtain

$$1480 \text{vec}(\Sigma_x) = T \text{vec}(\Sigma_x) + \eta^2 \frac{1}{b} \text{vec}(\Sigma_g), \quad T := \mathbb{E}[C_t \otimes C_t]. \quad (16)$$

1482 Rearranging gives

$$1483 (I - T) \text{vec}(\Sigma_x) = \eta^2 \frac{1}{b} \text{vec}(\Sigma_g). \quad (17)$$

1485 Assumption (A3) states that the linearised SGD on the quadratic approximation is linearly stable on E_+ , i.e.

$$1488 \rho\left(\mathbb{E}[(I - \eta \mathcal{H}(L_B))^{\otimes 2}]_{|E_+}\right) < 1.$$

1489 In our notation this means precisely that the restriction of T to $E_+ \otimes E_+$ has spectral radius strictly less than 1. Hence $I - T$ is invertible on $E_+ \otimes E_+$ and the vector equation (17) has a unique solution there, which corresponds to the unique stationary covariance Σ_x on E_+ .

1493 Step 3: Small-stepsize expansion and explicit form of Σ_x .

1494 We now compute the leading behaviour of Σ_x as a function of η for small η . Recall that

$$1496 C_t = I - \eta \mathcal{H}_{B_t}, \quad \mathcal{H}_{B_t} = \frac{1}{b} \sum_{i \in B_t} \mathcal{H}_i.$$

1498 Thus

$$1500 C_t \otimes C_t = (I - \eta \mathcal{H}_{B_t}) \otimes (I - \eta \mathcal{H}_{B_t}) = I \otimes I - \eta(\mathcal{H}_{B_t} \otimes I + I \otimes \mathcal{H}_{B_t}) + \eta^2(\mathcal{H}_{B_t} \otimes \mathcal{H}_{B_t}).$$

1502 Taking expectations and using $\mathbb{E}[\mathcal{H}_{B_t}] = \mathcal{H}$, we obtain

$$1503 T = I - \eta K + \eta^2 M, \quad (18)$$

1504 where K is the Kronecker-sum operator

$$1506 K := \mathcal{H} \otimes I + I \otimes \mathcal{H},$$

1507 and M is the operator defined by

$$1508 M := \mathbb{E}[\mathcal{H}_{B_t} \otimes \mathcal{H}_{B_t}].$$

1509 Substituting (18) into (17) gives

$$1510 (I - T) \text{vec}(\Sigma_x) = (\eta K - \eta^2 M) \text{vec}(\Sigma_x) = \eta^2 \frac{1}{b} \text{vec}(\Sigma_g).$$

On $E_+ \otimes E_+$ the operator K is positive definite (its eigenvalues are $\lambda_i + \lambda_j$ where $\lambda_i, \lambda_j > 0$ are eigenvalues of \mathcal{H}), hence invertible. Restricting to $E_+ \otimes E_+$, we can rewrite this as

$$(K - \eta M) \text{vec}(\Sigma_x) = \eta \frac{1}{b} \text{vec}(\Sigma_g). \quad (19)$$

For sufficiently small η , the operator $K - \eta M$ remains invertible and admits a Neumann-series expansion of its inverse. More precisely, on $E_+ \otimes E_+$ we have

$$(K - \eta M)^{-1} = K^{-1} + \eta K^{-1} M K^{-1} + \mathcal{O}(\eta^2),$$

where the $\mathcal{O}(\eta^2)$ term is understood in operator norm (depending on $\|K^{-1}\|$ and $\|M\|$). Applying this to (19) we obtain

$$\begin{aligned} \text{vec}(\Sigma_x) &= (K - \eta M)^{-1} \left(\eta \frac{1}{b} \text{vec}(\Sigma_g) \right) \\ &= \eta \frac{1}{b} (K^{-1} + \eta K^{-1} M K^{-1} + \mathcal{O}(\eta^2)) \text{vec}(\Sigma_g) \\ &= \eta \frac{1}{b} K^{-1} \text{vec}(\Sigma_g) + \mathcal{O}(\eta^2). \end{aligned}$$

Rewriting in matrix form and denoting by $\mathcal{K}(X) = \mathcal{H}X + X\mathcal{H}$ the corresponding Kronecker-sum operator on matrices, this says that on E_+

$$\Sigma_x = \frac{\eta}{b} \mathcal{K}^\dagger(\Sigma_g) + \mathcal{O}(\eta^2),$$

where \mathcal{K}^\dagger is the Moore–Penrose inverse of \mathcal{K} on $E_+ \otimes E_+$. This proves (10) in the statement of the proposition for the quadratic model.

The effect of the nonlinear remainders r_t can be treated as follows: by Assumption (A1), $\|r_t\| = \mathcal{O}(\|\Delta_t\|^2)$, and the stationary covariance Σ_x of the linear system is $\mathcal{O}(\eta)$, so typical $\|\Delta_t\|^2$ is $\mathcal{O}(\eta)$ and the additive perturbation $-\eta r_t$ to the dynamics has magnitude $\mathcal{O}(\eta^2)$. Its contribution to the noise covariance in the Lyapunov equation is thus $\mathcal{O}(\eta^4)$, which in turn produces a $\mathcal{O}(\eta^2)$ perturbation to Σ_x . Therefore the formula (10) remains valid up to an $\mathcal{O}(\eta^2)$ error for the true (non-quadratic) dynamics.

Step 4: Gradient–Noise Interaction ratio.

We now compute the ratio in (11). Fix $\theta = \theta^* + \Delta$ and consider a fresh mini-batch B . On the quadratic approximation we have

$$\nabla L(\theta) = \mathcal{H}\Delta,$$

and

$$\nabla L_B(\theta) = \frac{1}{b} \sum_{i \in B} (g_i + \mathcal{H}_i \Delta) = \underbrace{\xi_B}_{\text{zero mean}} + \underbrace{\mathcal{H}_B}_{\text{batch Hessian}} \Delta,$$

where $\xi_B := \frac{1}{b} \sum_{i \in B} g_i$ and $\mathcal{H}_B := \frac{1}{b} \sum_{i \in B} \mathcal{H}_i$. Conditioning on Δ and using $\mathbb{E}_B[\xi_B \mid \Delta] = 0$, we get

$$\begin{aligned} \mathbb{E}_B[\nabla L_B(\theta)^\top \mathcal{H} \nabla L_B(\theta) \mid \Delta] &= \mathbb{E}_B[(\mathcal{H}_B \Delta + \xi_B)^\top \mathcal{H} (\mathcal{H}_B \Delta + \xi_B) \mid \Delta] \\ &= \Delta^\top \mathbb{E}_B[\mathcal{H}_B^\top \mathcal{H} \mathcal{H}_B] \Delta + \mathbb{E}_B[\xi_B^\top \mathcal{H} \xi_B] \end{aligned} \quad (20)$$

(the cross term vanishes because $\mathbb{E}_B[\xi_B \mid \Delta] = 0$). Taking expectation over $\theta \sim \pi$ (the stationary law of the linear system) and recalling $\Sigma_x = \mathbb{E}_\pi[\Delta \Delta^\top]$ and $\mathbb{E}_B[\xi_B \xi_B^\top] = \frac{1}{b} \Sigma_g$, we obtain

$$\begin{aligned} N &:= \mathbb{E}_{\theta \sim \pi} \mathbb{E}_B[\nabla L_B(\theta)^\top \mathcal{H} \nabla L_B(\theta)] \\ &= \text{tr}(\mathcal{H} \Sigma_g) \frac{1}{b} + \text{tr}(\mathbb{E}_B[\mathcal{H}_B^\top \mathcal{H} \mathcal{H}_B] \Sigma_x). \end{aligned} \quad (21)$$

Similarly, the denominator is

$$D := \mathbb{E}_{\theta \sim \pi} [\|\nabla L(\theta)\|^2] = \mathbb{E}_{\theta \sim \pi} [\Delta^\top \mathcal{H}^2 \Delta] = \text{tr}(\mathcal{H}^2 \Sigma_x). \quad (22)$$

To relate N and D , we use the Lyapunov equation (15) for the quadratic model. Expanding the right-hand side of (15), we have

$$\begin{aligned}\mathbb{E}[C_t \Sigma_x C_t^\top] &= \mathbb{E}\left[(I - \eta \mathcal{H}_{B_t}) \Sigma_x (I - \eta \mathcal{H}_{B_t})^\top\right] \\ &= \Sigma_x - \eta \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x] - \eta \mathbb{E}[\Sigma_x \mathcal{H}_{B_t}] + \eta^2 \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x \mathcal{H}_{B_t}].\end{aligned}$$

Using $\mathbb{E}[\mathcal{H}_{B_t}] = \mathcal{H}$, this simplifies to

$$\mathbb{E}[C_t \Sigma_x C_t^\top] = \Sigma_x - \eta(\mathcal{H} \Sigma_x + \Sigma_x \mathcal{H}) + \eta^2 \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x \mathcal{H}_{B_t}].$$

Substituting this back into (15) and cancelling Σ_x on both sides, we obtain

$$0 = -\eta(\mathcal{H} \Sigma_x + \Sigma_x \mathcal{H}) + \eta^2 \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x \mathcal{H}_{B_t}] + \eta^2 \frac{1}{b} \Sigma_g.$$

Dividing by η yields the identity

$$\mathcal{H} \Sigma_x + \Sigma_x \mathcal{H} = \eta \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x \mathcal{H}_{B_t}] + \eta \frac{1}{b} \Sigma_g. \quad (23)$$

Now multiply both sides of (23) on the left by \mathcal{H} and take traces. Using the cyclicity of the trace and the fact that \mathcal{H} is symmetric, we get

$$\begin{aligned}\text{tr}(\mathcal{H}(\mathcal{H} \Sigma_x + \Sigma_x \mathcal{H})) &= \text{tr}(\mathcal{H}^2 \Sigma_x) + \text{tr}(\mathcal{H} \Sigma_x \mathcal{H}) \\ &= 2 \text{tr}(\mathcal{H}^2 \Sigma_x) = 2D,\end{aligned} \quad (24)$$

and

$$\text{tr}\left(\mathcal{H}(\eta \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x \mathcal{H}_{B_t}] + \eta \frac{1}{b} \Sigma_g)\right) = \eta \text{tr}\left(\mathcal{H} \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x \mathcal{H}_{B_t}]\right) + \eta \frac{1}{b} \text{tr}(\mathcal{H} \Sigma_g). \quad (25)$$

Equating (24) and (25) (they come from the two sides of (23)) gives

$$2D = \eta \text{tr}\left(\mathcal{H} \mathbb{E}[\mathcal{H}_{B_t} \Sigma_x \mathcal{H}_{B_t}]\right) + \eta \frac{1}{b} \text{tr}(\mathcal{H} \Sigma_g). \quad (26)$$

Comparing (21) and (26), we see that

$$N = \frac{1}{b} \text{tr}(\mathcal{H} \Sigma_g) + \text{tr}\left(\mathbb{E}[\mathcal{H}_{B_t}^\top \mathcal{H} \mathcal{H}_{B_t}] \Sigma_x\right) = \frac{2}{\eta} D.$$

Therefore, on the quadratic approximation we have the *exact* identity

$$\frac{N}{D} = \frac{2}{\eta}. \quad (27)$$

Restoring the higher-order terms $R_i(\theta)$ from Assumption (A1) introduces an additional third-order remainder in the Taylor expansion of the loss over one SGD step. For a Lipschitz Hessian with constant L_2 , standard Taylor estimates give a per-step remainder of order $\mathcal{O}(\eta^3 \|\nabla L_B(\theta_t)\|^3)$, which is $\mathcal{O}(\eta^3)$ in expectation under the stationary law (since the stationary covariance scales like $\mathcal{O}(\eta)$ by (10)). This adds a term of order $\mathcal{O}(\eta^3)$ to the stationarity condition $\mathbb{E}[L(\theta_{t+1}) - L(\theta_t)] = 0$, and hence contributes only a factor $\mathcal{O}(\eta)$ to the ratio N/D . Thus

$$\frac{N}{D} = \frac{2}{\eta} (1 + \mathcal{O}(\eta)),$$

which is precisely (11). This completes the proof. \square

F ON THE TWO TYPES OF OSCILLATIONS IN NNS

Differentiating Oscillations in Neural Network Optimization Our analytical treatment of SGD on one-dimensional quadratic objectives in Appendix D.1 leverages the simplicity of having a single curvature measure—the second derivative—which facilitates a precise landscape characterization and explicit stability conditions. However, extending this analysis to multidimensional quadratics

1620 already introduces significantly more intricate dynamics, necessitating advanced analytical frame-
 1621 works as developed by (Wu et al., 2018; Ma & Ying, 2021; Mulayoff & Michaeli, 2024). Transi-
 1622 tioning further to neural network optimization increases this complexity dramatically, since training
 1623 predominantly occurs away from the manifold of minima, including the EOS-like instabilities them-
 1624 selves (as evidenced by the continuous reduction in loss)—and therefore requires to go beyond linear
 1625 stability of quadratics near the manifold of minima.

1626 Given the current absence of robust theoretical tools to comprehensively analyze such dynamics,
 1627 distinguishing between curvature-driven and noise-driven oscillations necessitates empirical exper-
 1628 imentation. Specifically, we probe the dynamics by systematically varying hyperparameters (e.g.,
 1629 step size or batch size), as illustrated in Figure 3, allowing us to differentiate curvature-induced
 1630 (Type-2) oscillations from purely noise-induced (Type-1) oscillations (Figure 5).

1631
 1632 **Type-2 Oscillations Are Unique to NN Optimization** This complexity inherent in neural net-
 1633 work optimization is not merely an analytical inconvenience; rather, it is intrinsically tied to the
 1634 emergence and significance of Type-2 oscillations and EOS-style phenomena. Notably, Type-2 os-
 1635 cillations emerge naturally⁹ only in the case of neural network optimization, but not in the case
 1636 of quadratic objectives. In the one-dimensional quadratic scenario analyzed previously, curvature-
 1637 driven oscillations require the step size to precisely match the stability threshold $2/\lambda_{\max}$, or exceed
 1638 it, in which case we have divergence—in either case, it means that optimization of quadratics does
 1639 naturally enter a regime of instability. In contrast, neural network optimization uniquely exhibits
 1640 *progressive sharpening*, a third-order derivative phenomenon (Damian et al., 2023), where curva-
 1641 ture naturally increases during training. This progressive increase in curvature means that training
 1642 with a fixed step size can transition into an EOS-like regime of instability without any explicit adjust-
 1643 ment of the hyperparameters, and stay there due to self-stabilization effects (Damian et al., 2023).
 1644 Hence, Type-2 oscillations emerge naturally and robustly within neural network training dynamics
 1645 due to this intrinsic change of the loss landscape. Consequently, Type-2 oscillations and EOS-like
 1646 regimes are fundamentally driven by progressive sharpening, which does not happen in quadratics,
 1647 making it a purely neural network optimization phenomena.

1648 F.1 ON THE IMPORTANCE OF TYPE-2 OSCILLATIONS COMPARED TO TYPE-1

1649 Noise-induced (Type-1) oscillations are not unstable when introducing slight perturbations (increase
 1650 step size or decrease batch size), as showcased in Figure 9 and 10. Therefore, they do not constitute
 1651 an EoS-type phenomena, where slight perturbations do cause divergence (“complete” divergence as
 1652 long as we consider just the quadratic terms and can ignore higher terms — the fact that it doesn’t
 1653 fully diverge is exactly the higher-terms effect). Instead, after a perturbation, noise-induced oscilla-
 1654 tions quickly re-stabilize at a higher level.

1655 Crucially, a lack of such divergence means that noise-induced oscillations wouldn’t exhibit the self-
 1656 stabilization mechanism of Damian et al. (2023) characteristic of EoS (differing it from classical
 1657 convex optimization). Moreover, as shown in the quadratic example and in the proofs, noise-
 1658 induced oscillations happen for any quadratic, for a wide range of step sizes, making them in-
 1659 herently “unsurprising”, while EoS is a beyond-quadratic phenomena (and, as far as we know, a
 1660 deep-learning-specific phenomena), as it relies on both progressive sharpening and the aforemen-
 1661 tioned self-stabilization, both being an effects of higher order terms. And the reason why we care
 1662 specifically about effects of beyond-quadratic terms is specifically the adaptation of the landscape
 1663 to the hyper-parameters, which is, by definition, an effect of higher order terms. That is the reason
 1664 we specifically care about curvature-driven oscillations.

1665 Now, with all of the above, GNI, being an indicator of those noise-induced oscillations, is therefore
 1666 not an indicator of EoS-like regime. This is despite the fact that GNI in SGD comes from the same
 1667 place as λ_{\max} /Rayleigh quotient in GD — i.e. from the descent lemma; yet, it does not mean that
 1668 the two quantities serve the same role. Instead, it is the presence of the natural noise in SGD that
 1669 makes the analysis much more complex. Instead, GNI has its usefulness as a measure of the level
 1670 of noise coming from SGD. That is, noise-induced oscillations are influenced by the Hessian, but
 1671 are also strongly influenced by the ratio between the noise covariance and the norm of full batch

1672 ⁹We define an as emerging *naturally* if it arises inherently from the training dynamics, and not a result of
 1673 precisely-selected hyperparameters or initializations, reflecting a fundamental characteristic of the optimization
 process itself. Formally, it needs to happen over a range of hyperparameter choices and initializations.

1674 gradient, with the latter being the leading cause of change. In particular, GNI is decoupled from
 1675 the Hessian, and can change drastically without any change of landscape sharpness, as showcased
 1676 in our experiments. Lastly, another important consequence of EoS is that the landscape adapts to
 1677 the hyper-parameters (rather than the other way around in classical optimization). With GNI being
 1678 decoupled from the Hessian, GNI being at $2/\eta$ is not an indication of landscape adopting to the
 1679 hyper-parameters, as is the case with λ_{\max} being at $2/\eta$ during GD.

1680 G GRADIENTS EXPLODE ABOVE THE EOSS: PROOF OF THEOREM 1

1681 G.1 PART 1: EXPLOSION

1682
 1683 **Setting.** For each minibatch index i , let $L_i : \mathbb{R}^d \rightarrow \mathbb{R}$ be twice differentiable with (possibly
 1684 index-dependent) positive semidefinite Hessian H_i that is *constant in θ* (quadratic model). Write

$$1685 Y_i(\theta) := \nabla L_i(\theta) = H_i(\theta - x_i), \quad \|\cdot\| = \|\cdot\|_2.$$

1686 At iteration t , stochastic gradient descent (SGD) draws $j_t \sim \mathcal{P}_b$ independently of the past and
 1687 performs

$$1688 \theta_{t+1} = \theta_t - \eta Y_{j_t}(\theta_t), \quad \eta > 0.$$

1689 We define basic statistics evaluated at θ_t :

$$1690 r_i(\theta_t) := \frac{Y_i(\theta_t)^\top H_i Y_i(\theta_t)}{\|Y_i(\theta_t)\|^2} \in [0, \infty).$$

1691 Note that *Batch Sharpness* is

$$1692 \text{Batch Sharpness}(\theta) := \mathbb{E}_i[r_i(\theta)].$$

1693 All randomness lives on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. At each step $t \geq 0$, SGD draws $j_t \sim \mathcal{P}_b$ i.i.d.
 1694 and independent of the σ -algebra of the past

$$1695 \mathcal{F}_t := \sigma(\theta_0, j_0, \dots, j_{t-1}).$$

1696 We assume

$$1697 \Lambda := \operatorname{ess\,sup}_{i \sim \mathcal{P}_b} \|H_i\| < \infty \quad \text{and} \quad \mathbb{E}_{i \sim \mathcal{P}_b} [\|Y_i(\theta)\|^2] < \infty \quad \text{for all } \theta \text{ reached by SGD.}$$

1698 (These ensure that the Rayleigh quotients below are well defined and integrable.) We adopt the
 1699 convention that, for any i and θ with $Y_i(\theta) = 0$,

$$1700 \frac{Y_i(\theta)^\top H_i Y_i(\theta)}{\|Y_i(\theta)\|^2} := 0, \quad \frac{\|Y_i(\theta - \eta Y_i(\theta))\|^2}{\|Y_i(\theta)\|^2} := 1.$$

1701 We provide a rigorous instability statement: a multiplicative explosion controlled solely by *Batch
 1702 Sharpness*.

1703 **Proposition 4** (On-batch multiplicative explosion under *Batch Sharpness* $> 2/\eta$ for quadratics.).
 1704 For each t define the on-batch factor

$$1705 R_t := \mathbb{E}_{j_t} \left[\frac{\|Y_{j_t}(\theta_{t+1})\|^2}{\|Y_{j_t}(\theta_t)\|^2} \mid \mathcal{F}_t \right].$$

1706 Then, for all t ,

$$1707 R_t \geq \left(1 - \eta \text{Batch Sharpness}(\theta_t)\right)^2.$$

1708 Consequently, if there exist indices $t_0, \dots, t_0 + T - 1$ and a constant $\epsilon > 0$ such that

$$1709 \text{Batch Sharpness}(\theta_t) \geq \frac{2 + \epsilon}{\eta} \quad \text{for all } t_0 \leq t \leq t_0 + T - 1,$$

1710 then

$$1711 \mathbb{E} \left[\prod_{s=t_0}^{t_0+T-1} \frac{\|Y_{j_s}(\theta_{s+1})\|^2}{\|Y_{j_s}(\theta_s)\|^2} \right] \geq (1 + \epsilon)^{2T},$$

1712 i.e., the product of on-batch factors grows exponentially in expectation.

Beyond quadratic. Beyond the quadratic model: with θ -dependent Hessians, the identity $Y_i(\theta_{t+1}) = (I - \eta H_i)Y_i(\theta_t)$ incurs Taylor remainders of order $O(\eta \|Y_i\| \|\nabla H_i\|)$; standard Lipschitz–Hessian assumptions then yield perturbative versions of Proposition 4 with additional $O(\eta^3)$ terms. We do not invoke these here to keep the statements exact.

Proof of Proposition 4. Fix t and condition on $j_t = i$. By the quadratic model,

$$Y_i(\theta_{t+1}) = H_i(\theta_t - \eta Y_i(\theta_t) - x_i) = (I - \eta H_i) Y_i(\theta_t).$$

Hence

$$\frac{\|Y_i(\theta_{t+1})\|^2}{\|Y_i(\theta_t)\|^2} = 1 - 2\eta r_i(\theta_t) + \eta^2 q_i(\theta_t), \quad q_i(\theta_t) := \frac{Y_i(\theta_t)^\top H_i^2 Y_i(\theta_t)}{\|Y_i(\theta_t)\|^2}.$$

By Cauchy–Schwarz in the H_i -inner product, $q_i(\theta_t) \geq r_i(\theta_t)^2$. Therefore, pointwise in i ,

$$\frac{\|Y_i(\theta_{t+1})\|^2}{\|Y_i(\theta_t)\|^2} \geq (1 - \eta r_i(\theta_t))^2.$$

Averaging over j_t and applying Jensen’s inequality to the convex map $x \mapsto (1 - \eta x)^2$ yields

$$R_t \geq \mathbb{E}_i[(1 - \eta r_i(\theta_t))^2] \geq (1 - \eta \mathbb{E}_i[r_i(\theta_t)])^2 = (1 - \eta \text{Batch Sharpness}(\theta_t))^2.$$

If $\text{Batch Sharpness}(\theta_t) \geq (2 + \epsilon)/\eta$ for each $t \in [t_0, t_0 + T - 1]$, then $R_t \geq (1 + \epsilon)^2 > 1 + 2\epsilon$ deterministically given \mathcal{F}_t .

Now, let

$$Z_s := \frac{\|Y_{j_s}(\theta_{s+1})\|^2}{\|Y_{j_s}(\theta_s)\|^2}, \quad \mathcal{F}_s := \sigma(\theta_0, j_0, \dots, j_{s-1}), \quad U_t := \prod_{r=t_0}^{t_0+t-1} Z_r \quad (U_0 := 1).$$

From the single-step bound proved above we have, for each s ,

$$\mathbb{E}[Z_s | \mathcal{F}_s] \geq (1 - \eta \text{Batch Sharpness}(\theta_s))^2 \geq (1 + \epsilon)^2.$$

By the tower property,

$$\mathbb{E}[U_{t+1}] = \mathbb{E}[U_t Z_{t_0+t}] = \mathbb{E}[U_t \mathbb{E}[Z_{t_0+t} | \mathcal{F}_{t_0+t}]] \geq \mathbb{E}[U_t \gamma_{t_0+t}].$$

$S(\theta_s) \geq 2/\eta + \epsilon \eta$ (with $\epsilon > 0$), then $\gamma_s = (1 - \eta S(\theta_s))^2 \geq (1 + \epsilon \eta^2)^2 =: \gamma^2$, hence

$$\mathbb{E}[U_{t+1}] \geq \gamma^2 \mathbb{E}[U_t] \implies \mathbb{E}[U_T] \geq \gamma^{2T} = (1 + \epsilon)^{2T}.$$

we obtain the stated exponential lower bound. \square

Discussion after Proposition 4. The proof uses only: (i) the exact closure $Y_i(\theta_{t+1}) = (I - \eta H_i)Y_i(\theta_t)$ (quadratic model), (ii) Cauchy–Schwarz ($q_i \geq r_i^2$), and (iii) Jensen’s inequality. No cross-batch interaction is needed; hence the result holds assumption-free, and it naturally yields a multiplicative, path-wide instability witness.

G.2 PART 2: SGD WITH REPLACEMENT

In Proposition 4 we showed that if Batch Sharpness is bigger than $2/\eta$ for prolonged time a certain quantity explodes exponentially fast. What is missing is to show that if Batch Sharpness is bigger than $2/\eta$ on the quadratic setting at one step t_0 , then it will be also at the following time steps.

Lipschitz drift of batch sharpness in the quadratic model. Fix a minibatch i and write $r_i(\theta) := \frac{Y_i(\theta)^\top H_i Y_i(\theta)}{\|Y_i(\theta)\|^2}$ with $Y_i(\theta) = H_i(\theta - x_i)$ and $H_i = H_i^\top \succeq 0$. Along any segment in parameter space on which $\|Y_i(\theta)\| \geq g_{\min} > 0$ and $\|H_i\| \leq \Lambda$, the map $\theta \mapsto r_i(\theta)$ is Lipschitz with

$$\|\nabla_\theta r_i(\theta)\| \leq \frac{4\|H_i\|^2}{\|Y_i(\theta)\|} \leq \frac{4\Lambda^2}{g_{\min}},$$

hence, for any θ, θ' in that segment,

$$|r_i(\theta') - r_i(\theta)| \leq \frac{4\Lambda^2}{g_{\min}} \|\theta' - \theta\|.$$

Averaging over $i \sim \mathcal{P}_b$ yields the Lipschitz bound for the $\text{Batch Sharpness} := \mathbb{E}_i[r_i(\theta)]$:

$$|\text{Batch Sharpness}(\theta') - \text{Batch Sharpness}(\theta)| \leq L_S \|\theta' - \theta\|, \quad L_S := \frac{4\Lambda^2}{g_{\min}}. \quad (28)$$

Single-time margin implies a uniform window margin. Assume that along the SGD trajectory on $[t_0, t_0 + T]$ we have the uniform bounds $\|H_i\| \leq \Lambda$ for all i and $\|Y_i(\theta_t)\| \geq g_{\min} > 0$ for all i, t , and that the per-step move is bounded by $\|\theta_{t+1} - \theta_t\| = \eta \|Y_{j_t}(\theta_t)\| \leq \eta G_{\max}$ for some $G_{\max} > 0$ (these three constants define the quadratic region under consideration). If at time t_0

$$\text{Batch Sharpness}(\theta_{t_0}) \geq \frac{2 + \epsilon}{\eta} \quad (\epsilon > 0),$$

then for every k with $0 \leq k \leq T_* := \lfloor \frac{\epsilon}{2\eta L_S G_{\max}} \rfloor$,

$$\text{Batch Sharpness}(\theta_{t_0+k}) \geq \frac{2 + \epsilon/2}{\eta}. \quad (29)$$

Proof. By (28) and the step bound, $\text{Batch Sharpness}(\theta_{t+1}) \geq \text{Batch Sharpness}(\theta_t) - L_S \eta G_{\max}$. Iterating k times gives $\text{Batch Sharpness}(\theta_{t_0+k}) \geq \text{Batch Sharpness}(\theta_{t_0}) - k L_S \eta G_{\max} \geq \frac{2}{\eta} + \epsilon\eta - k L_S \eta G_{\max}$. If $k \leq \epsilon/(2\eta L_S G_{\max})$ this yields (29). \square

Consequence for the product (plug into the tower/induction step). On the whole window $t \in \{t_0, \dots, t_0 + T_* - 1\}$ we thus have $S(\theta_t) \geq (2 + \epsilon/2)/\eta$, hence from the one-step bound $\mathbb{E}[Z_t | \mathcal{F}_t] \geq (1 - \eta S(\theta_t))^2$ and the tower/induction argument,

$$\mathbb{E} \left[\prod_{s=t_0}^{t_0+T_*-1} Z_s \right] \geq \left(1 + \frac{\epsilon}{2}\right)^{2T_*}.$$

This upgrades a *single-time* margin at t_0 into an explicit *uniform window* of length T_* over which the exponential lower bound holds.

G.3 PART 2 FOR SGD Without REPLACEMENT

RR setting and remaining-set sharpness. Fix an epoch with a finite pool $\mathcal{I} = \{1, \dots, n\}$. In random reshuffling (RR), within an epoch we draw a uniform random permutation of \mathcal{I} and visit each index exactly once. Let $R_t \subseteq \mathcal{I}$ denote the *remaining* set at step t (those not yet visited in the current epoch) and let $m_t := |R_t|$. Define the *remaining-set sharpness* at θ_t by

$$S_{\text{rem}}(\theta_t) := \frac{1}{m_t} \sum_{i \in R_t} r_i(\theta_t), \quad r_i(\theta) := \frac{Y_i(\theta)^\top H_i Y_i(\theta)}{\|Y_i(\theta)\|^2}.$$

Let $\mathcal{F}_t^{\text{RR}} := \sigma(\theta_0, \text{the permutation prefix up to step } t-1)$; conditionally on $\mathcal{F}_t^{\text{RR}}$, j_t is uniform over R_t .

Step 1: One-step RR bound. In the quadratic model, $Y_{j_t}(\theta_{t+1}) = (I - \eta H_{j_t}) Y_{j_t}(\theta_t)$ and

$$\frac{\|Y_{j_t}(\theta_{t+1})\|^2}{\|Y_{j_t}(\theta_t)\|^2} = 1 - 2\eta r_{j_t}(\theta_t) + \eta^2 q_{j_t}(\theta_t) \geq (1 - \eta r_{j_t}(\theta_t))^2, \quad q_i(\theta) := \frac{Y_i(\theta)^\top H_i^2 Y_i(\theta)}{\|Y_i(\theta)\|^2} \geq r_i(\theta)^2,$$

by Cauchy–Schwarz. Averaging uniformly over R_t and using convexity of $x \mapsto (1 - \eta x)^2$ gives

$$\mathbb{E} \left[\frac{\|Y_{j_t}(\theta_{t+1})\|^2}{\|Y_{j_t}(\theta_t)\|^2} \middle| \mathcal{F}_t^{\text{RR}} \right] \geq \frac{1}{m_t} \sum_{i \in R_t} (1 - \eta r_i(\theta_t))^2 \geq (1 - \eta S_{\text{rem}}(\theta_t))^2. \quad (30)$$

Step 2: Persistence of a single-time margin (high probability). Assume the quadratic region bounds used in Part 1: $\|H_i\| \leq \Lambda$ for all i , $\|Y_i(\theta_t)\| \geq g_{\min} > 0$ for all i, t , and $\|\theta_{t+1} - \theta_t\| = \eta \|Y_{j_t}(\theta_t)\| \leq \eta G_{\max}$. As established there,

$$|r_i(\theta') - r_i(\theta)| \leq \frac{4\Lambda^2}{g_{\min}} \|\theta' - \theta\| \quad (31)$$

$$|\text{Batch Sharpness}(\theta') - \text{Batch Sharpness}(\theta)| \leq L_S \|\theta' - \theta\|, \quad L_S := \frac{4\Lambda^2}{g_{\min}}.$$

Suppose at the beginning of the epoch (time t_0) we have a single-time margin

$$\text{Batch Sharpness}(\theta_{t_0}) \geq \frac{2}{\eta} + \varepsilon \eta \quad (\varepsilon > 0). \quad (32)$$

Fix integers $K \in \{1, \dots, n-1\}$ and define the minimal remaining size $m_{\min} := n - K$. For any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the RR permutation,

$$\max_{0 \leq k \leq K} \left| \frac{1}{|R_{t_0+k}|} \sum_{i \in R_{t_0+k}} r_i(\theta_{t_0}) - \text{Batch Sharpness}(\theta_{t_0}) \right| \leq \Delta_K(\delta), \quad \Delta_K(\delta) := \Lambda \sqrt{\frac{2 \log(2K/\delta)}{m_{\min}}}, \quad (33)$$

(Hoeffding–Serfling inequality; we only use the simple range bound $r_i(\theta_{t_0}) \in [0, \Lambda]$). Combining (31) and (33), for each $k \in \{0, \dots, K\}$ we obtain on the same event

$$S_{\text{rem}}(\theta_{t_0+k}) \geq \underbrace{\frac{1}{|R_{t_0+k}|} \sum_{i \in R_{t_0+k}} r_i(\theta_{t_0})}_{\text{finite-pop mean at } t_0} - L_S \|\theta_{t_0+k} - \theta_{t_0}\| \geq \text{Batch Sharpness}(\theta_{t_0}) - \Delta_K(\delta) - L_S \eta G_{\max} k. \quad (34)$$

Hence, if K and δ are chosen so that

$$\Delta_K(\delta) + L_S \eta G_{\max} K \leq \frac{\varepsilon}{2} \eta, \quad (35)$$

then (32) and (34) imply the uniform window margin

$$S_{\text{rem}}(\theta_{t_0+k}) \geq \frac{2}{\eta} + \frac{\varepsilon}{2} \eta \quad \text{for all } k \in \{0, \dots, K\}, \quad \text{with probability at least } 1 - \delta. \quad (36)$$

Step 3: Exponential growth over the RR window. On the event (36), the one-step bound (30) yields

$$\mathbb{E}[Z_{t_0+k} | \mathcal{F}_{t_0+k}^{\text{RR}}] \geq (1 - \eta S_{\text{rem}}(\theta_{t_0+k}))^2 \geq (1 + (\varepsilon/2) \eta^2)^2 \quad \text{for all } k \in \{0, \dots, K-1\},$$

where $Z_t := \frac{\|Y_{j_t}(\theta_{t+1})\|^2}{\|Y_{j_t}(\theta_t)\|^2}$. By the tower property and induction,

$$\mathbb{E} \left[\prod_{s=t_0}^{t_0+K-1} Z_s \right] \geq (1 - \delta) (1 + (\varepsilon/2) \eta^2)^{2K}.$$

In words: a single-time margin (32) at the start of the epoch persists, with high probability under RR, over a whole window whose length K is explicitly controlled by the drift budget (31) and the finite-population deviation (33). On that window the product of on–batch factors explodes exponentially.

1890 H PROOF OF THE EQUIVALENCE OF SECTION B.2

1891
1892 In this appendix we make precise the informal statement in Section B.2 that, on the local quadratic
1893 approximation of the loss, the following three viewpoints are equivalent:
1894

- 1895 (i) breaking a valid instability criterion (Definition 4);
- 1896 (ii) experiencing a catapult on the quadratic model (Definition 5);
- 1897 (iii) observing a loss spike of “sufficient” size.

1898
1899 Throughout we work on the local quadratic model introduced in Section B, and make explicit the
1900 mild conditions under which the equivalence holds.

1901 H.1 SETTING AND BASIC ASSUMPTION

1902
1903 Fix a time t and a point θ_t , and consider the quadratic model \tilde{L} of the loss around θ_t as defined in
1904 Section B. We recall that \tilde{L} is of the form
1905

$$1906 \quad \tilde{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \tilde{L}_i(\theta), \quad \tilde{L}_i(\theta) := \frac{1}{2}(\theta - x_i)^\top \mathcal{H}_i(\theta - x_i)$$

1907
1908 for some matrices \mathcal{H}_i and points x_i .

1909
1910 We continue to denote by U_t an open neighborhood of θ_t on which the quadratic approximation is
1911 accurate, in the sense described informally in Section B. For the arguments below, we only need the
1912 following simple condition.

1913 **Assumption 1** (Coercive quadratic model on U_t). There exists a symmetric positive semi-definite
1914 matrix $\hat{\mathcal{H}}$ and constants $0 < \mu \leq L < \infty$ such that:
1915

- 1916 (a) $\tilde{L}(\theta) = \frac{1}{2}(\theta - \theta^*)^\top \hat{\mathcal{H}}(\theta - \theta^*)$ for some θ^* (i.e. \tilde{L} is exactly quadratic);
- 1917 (b) all non-zero eigenvalues of $\hat{\mathcal{H}}$ lie in $[\mu, L]$;
- 1918 (c) the connected component of θ_t in the sublevel set $\{\theta : \tilde{L}(\theta) \leq R\}$ is contained in U_t for
1919 all R in a neighborhood of $\tilde{L}(\theta_t)$.

1920
1921 Assumption 1 is satisfied, for instance, when the full-batch Hessian of the original loss at θ_t is
1922 positive definite in the relevant directions and the neighborhood U_t is chosen small enough. It
1923 implies in particular that:
1924

- 1925 • for any $R < \infty$, the sublevel set $\{\theta : \tilde{L}(\theta) \leq R\}$ is compact;
- 1926 • the level sets of \tilde{L} define a family of nested compacts around θ_t inside U_t .

1927 H.2 BREAKING A VALID INSTABILITY CRITERION \iff CATAPULT

1928
1929 We first clarify the relationship between valid instability criteria (Definition 4) and catapults (Defi-
1930 nition 5).
1931

1932 Recall that Definition 4 is stated for a discrete-time dynamical system $(\theta_s)_{s \geq 0}$ on a parameter space
1933 Θ , with an open set $U \subseteq \Theta$, a scalar map $f : U \rightarrow \mathbb{R}$ and a threshold $c \in \mathbb{R}$. It says that f is a *valid*
1934 *instability criterion with threshold c on U* if
1935

$$1936 \quad f(\theta_0) > c \implies (\theta_s)_{s \geq 0} \text{ leaves every compact subset of } U \text{ in finite time.}$$

1937
1938 On the quadratic model at time t , Definition 5 simply re-uses this notion with $U := U_t$ and $\theta_0 := \theta_t$:
1939

1940 We say that the algorithm *experiences a catapult at time t* if, when run on \tilde{L} from
1941 initialization θ_t , the resulting trajectory $(\theta_s)_{s \geq t}$ leaves every compact subset of U_t
1942 in finite time.
1943

Thus, formally:

Lemma 3 (Breaking a valid criterion gives a catapult). *Let $f : U_t \rightarrow \mathbb{R}$ be a valid instability criterion with threshold c for the quadratic dynamics on U_t in the sense of Definition 4. Fix a time t and consider the trajectory of the quadratic model initialized at θ_t .*

If $f(\theta_t) > c$, then the quadratic trajectory experiences a catapult at time t in the sense of Definition 5.

Proof. Apply Definition 4 with $U := U_t$ and initial condition $\theta_0 := \theta_t$. Since $f(\theta_t) > c$ and f is a valid instability criterion on U_t , we have that the trajectory $(\theta_s)_{s \geq t}$ leaves every compact subset of U_t in finite time. But this is precisely the definition of a catapult at time t on the quadratic model (Definition 5). \square

Conversely, if a catapult occurs at time t , one can always *construct* a (possibly highly non-smooth) instability criterion which is broken at θ_t :

Lemma 4 (Catapult implies existence of a valid criterion). *Assume that the quadratic trajectory from θ_t experiences a catapult on U_t in the sense of Definition 5. Then there exists a map $f : U_t \rightarrow \mathbb{R}$ and a threshold $c \in \mathbb{R}$ such that:*

- (a) f is a valid instability criterion with threshold c for the quadratic dynamics on U_t ;
- (b) $f(\theta_t) > c$.

Proof. Define $f : U_t \rightarrow \mathbb{R}$ by

$$f(\theta_0) := \begin{cases} 1, & \text{if the quadratic trajectory initialized at } \theta_0 \\ & \text{leaves every compact subset of } U_t \text{ in finite time;} \\ 0, & \text{otherwise.} \end{cases}$$

Set $c := \frac{1}{2}$. Then, by construction, $f(\theta_0) > c$ if and only if the trajectory initialized at θ_0 leaves every compact subset of U_t in finite time. Hence f is a valid instability criterion with threshold c in the sense of Definition 4. Since, by assumption, the trajectory from θ_t experiences a catapult, we have $f(\theta_t) = 1 > c$, as desired. \square

Lemmas 3 and 4 make precise the first equivalence in the box of Section B.2: on the quadratic model, “breaking a (valid) instability criterion” is just another way to state the occurrence of a catapult, and conversely any catapult defines at least one (possibly non-unique) valid instability criterion which is broken.

H.3 CATAPULTS \iff LOSS SPIKES OF SUFFICIENT SIZE

We now formalize the relationship between catapults and loss spikes for the quadratic model under Assumption 1.

Definition 6 (Loss spike of size α on the quadratic model). *Fix t and θ_t , and let \tilde{L} be the quadratic model as above. For $\alpha > 1$ we say that the quadratic trajectory $(\theta_s)_{s \geq t}$ has a *loss spike of relative size at least α at time t* if there exists a time $s \geq t$ such that*

$$\tilde{L}(\theta_s) \geq \alpha \tilde{L}(\theta_t).$$

We say that it has loss spikes of *arbitrarily large* relative size if this holds for all $\alpha > 1$.

Under Assumption 1, sublevel sets of \tilde{L} are compact, and they form a convenient family of compact sets to test the definition of catapult.

Lemma 5 (Catapult \iff unbounded quadratic loss). *Under Assumption 1, the following are equivalent for the quadratic trajectory $(\theta_s)_{s \geq t}$:*

- (i) *the trajectory experiences a catapult at time t on U_t , i.e. it leaves every compact subset of U_t in finite time;*
- (ii) *the quadratic loss along the trajectory is unbounded above, i.e. $\sup_{s \geq t} \tilde{L}(\theta_s) = +\infty$.*

In particular, under Assumption 1, a catapult is equivalent to the existence of loss spikes of arbitrarily large relative size in the sense of Definition 6.

Proof. (i) \Rightarrow (ii). Assume the trajectory experiences a catapult on U_t . Suppose for contradiction that $\sup_{s \geq t} \tilde{L}(\theta_s) \leq R$ for some finite R . Then the whole trajectory is contained in the sublevel set

$$K_R := \{\theta \in U_t : \tilde{L}(\theta) \leq R\}.$$

By Assumption 1(b)–(c), K_R is compact and contained in U_t . Hence the trajectory *never* leaves K_R , which contradicts the definition of a catapult (it must leave every compact subset of U_t in finite time). Thus $\sup_{s \geq t} \tilde{L}(\theta_s) = +\infty$.

(ii) \Rightarrow (i). Conversely, assume that $\sup_{s \geq t} \tilde{L}(\theta_s) = +\infty$. Let $K \subset U_t$ be any compact set containing θ_t . Compactness and continuity of \tilde{L} imply that

$$R_K := \sup_{\theta \in K} \tilde{L}(\theta) < \infty.$$

Since $\sup_{s \geq t} \tilde{L}(\theta_s) = +\infty$, there exists a time $s_K \geq t$ with $\tilde{L}(\theta_{s_K}) > R_K$, hence $\theta_{s_K} \notin K$.

It remains to check that the trajectory eventually *stays* outside K : but this follows immediately from the fact that once $\tilde{L}(\theta_s) > R_K$, the iterate cannot re-enter K without violating the definition of R_K as the supremum of \tilde{L} on K . Thus, for every compact $K \subset U_t$ containing θ_t , there is a finite time s_K such that $\theta_s \notin K$ for all $s \geq s_K$. This is exactly the definition of a catapult on U_t . \square

Combining Lemma 5 with Definition 6 immediately yields:

Corollary 2 (Catapult \iff loss spikes of sufficient size). *Under Assumption 1, for the quadratic trajectory $(\theta_s)_{s \geq t}$ the following are equivalent:*

- (i) the trajectory experiences a catapult at time t on U_t ;
- (ii) for every $\alpha > 1$ there exists $s \geq t$ such that $\tilde{L}(\theta_s) \geq \alpha \tilde{L}(\theta_t)$, i.e. the trajectory has loss spikes of arbitrarily large relative size in the sense of Definition 6.

From the point of view of Section B.2, we can now interpret a *loss spike of sufficient size* as a spike whose relative height exceeds any upper bound that would be compatible with bounded (linearly stable) dynamics on the quadratic model. In particular, in regimes where the quadratic dynamics is linearly stable and confined to a given compact subset of U_t , the loss is uniformly bounded and only admits spikes of bounded relative size; Corollary 1 in Appendix D quantifies this explicitly in a one-dimensional example.

H.4 PUTTING THE PIECES TOGETHER

We can now summarize the equivalence more succinctly.

Theorem 2 (Equivalence of the three viewpoints on the quadratic model). *Work on the quadratic approximation \tilde{L} at time t under Assumption 1. For the quadratic trajectory $(\theta_s)_{s \geq t}$ on U_t , the following statements are equivalent:*

- (i) there exists a valid instability criterion $f : U_t \rightarrow \mathbb{R}$ with threshold c in the sense of Definition 4 such that $f(\theta_t) > c$;
- (ii) the trajectory experiences a catapult at time t on U_t in the sense of Definition 5;
- (iii) the quadratic loss along the trajectory has loss spikes of arbitrarily large relative size, i.e. for every $\alpha > 1$ there exists $s \geq t$ with $\tilde{L}(\theta_s) \geq \alpha \tilde{L}(\theta_t)$.

Proof. (i) \Rightarrow (ii) is Lemma 3. (ii) \Rightarrow (i) is Lemma 4.

(ii) \Leftrightarrow (iii) is exactly Lemma 5 together with the definition of loss spikes (Definition 6) and Corollary 2. \square

Theorem 2 formalizes the slogan in Section B.2: on the quadratic approximation, the underlying property is *divergence of the dynamics on U_t* , and the three diagnostics we use throughout the paper—breaking an instability criterion, observing a catapult, and observing a large enough loss spike—are just different ways of detecting the same phenomenon.

I ON THE FATE OF λ_{\max}

In this section we examine how λ_{\max} behaves once EOSS is reached and clarify its relationship to *Batch Sharpness*. A key aspect of the original EOS analysis is, indeed, that the controlling quantity—the largest eigenvalue of the full-batch Hessian λ_{\max} —has an immediate geometric interpretation. There exists an extensive literature about λ_{\max} size and role in neural networks, and it is a main ingredient of any proof of convergence. The EOSS picture replaces λ_{\max} with *Batch Sharpness*, a statistic whose connection to generalization and role in optimization theory is largely unexplored.

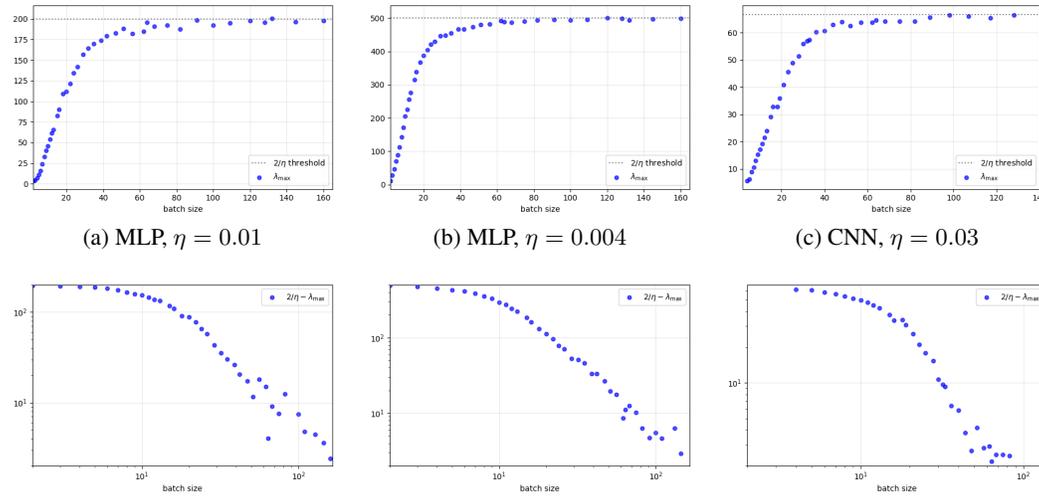


Figure 12: **Stabilisation level of λ_{\max} across step sizes and architectures.** **Top:** final-epoch λ_{\max} vs. batch size. **Bottom:** log-log plots of the gap $2/\eta - \lambda_{\max}$ for the same runs. All experiments use CIFAR-10 8k.

I.1 EMPIRICAL FACTS

Below the phenomena we extensively observe in vision classification tasks trained with MSE, ablating on batch sizes, step sizes, architectures, datasets. See Figure 6 for a good reference of what generally goes on.

- **Fact 1: Progressive Sharpening.** λ_{\max} increases at most as long as *Batch Sharpness* increases.
- **Fact 2: Phase Transition.** Once *Batch Sharpness* plateaus at $2/\eta$, λ_{\max} stops increasing. If it moves, it only decreases from this time on.
- **Fact 3: Path-dependence.** If changes to hyper parameters are made, *Batch Sharpness* changes abruptly or restart growing and λ_{\max} also changes. Stabilization of both happen as *Batch Sharpness* reaches $2/\eta$. The trajectory of λ_{\max} is not fully determined by the size of hyper parameters (see Figure 7). That is, the level of λ_{\max} is *path-dependent*: it inherits the history of progressive sharpening up to the moment EOSS is reached.
- **Fact 4: Smaller batches \Rightarrow flatter minima.** Across every setting we tested, reducing the batch size monotonically decreases the plateau level of λ_{\max} . This aligns with the long-standing empirical observation that smaller batches locate flatter minima see, e.g., Keskar et al. (2016); Jastrzębski et al. (2021)).
- **Fact 5: A critical batch size marks the SGD \rightarrow GD crossover.** Each curve in Figure 12 exhibits a bend at $b \approx b_c(\eta)$: for $b < b_c$ the plateau falls rapidly with b , while for $b > b_c$ it

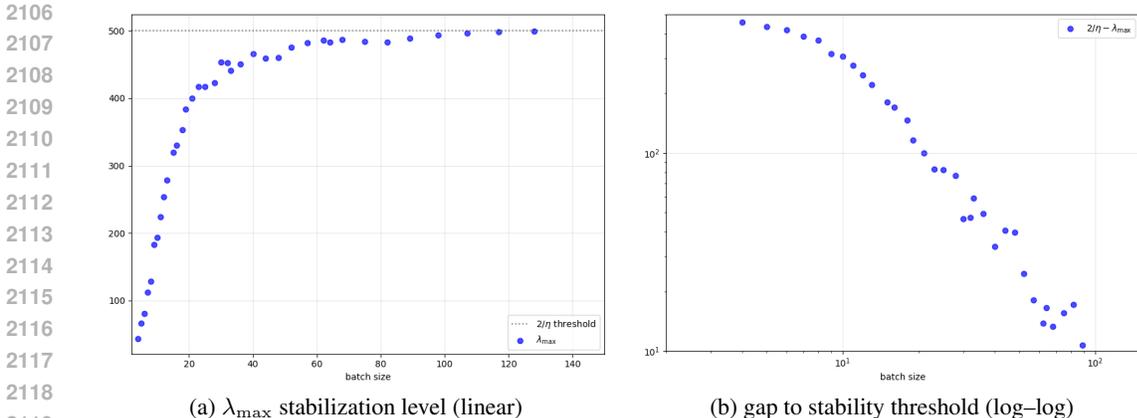


Figure 13: **Baseline MLP: stabilization of λ_{\max} as a function of batch size.** Baseline MLP (2 hidden layers, width 512) trained on an 8k-subset of CIFAR-10 with step size 0.004 until convergence. **(a)** Final λ_{\max} (linear axes). Smaller batches settle to flatter minima. For batch sizes *below* the *critical batch size* b_c the level of stabilization is significantly below the $2/\eta$ level of full-batch, indicating strong *implicit regularization*. Moreover, the curve is steep, making the the final landscape sensitive to the choice of batch size. For *larger* batches ($b > b_c$) the slope flattens and λ_{\max} plateaus close to $2/\eta$, so the dynamics resemble full-batch GD, implicit regularization is **weak**. **(b)** Log-log plot of the gap $2/\eta - \lambda_{\max}$, used to test for any power-law decay.

flattens and approaches the full-batch value. This b_c corresponds to the regime in which the mini-batch landscapes approximate *well enough* the full-batch landscape, restoring GD-like dynamics (Appendix I.2).

- **Fact 6: No universal power law.** From *static* analysis, one would expect a scaling $2/\eta - \lambda_{\max} = O(b^{-\alpha})$ for some α . The log-log plots (bottom row of Figure 12) show no robust straight-line behaviour, ruling out such law for any possible exponent $-\alpha$.

I.2 CRITICAL BATCH SIZE

We can characterize two regimes for the stabilization levels (see Figure 13):

- Small-batch regime ($b \leq b_c$):** λ_{\max} stabilizes well *below* the full-batch threshold $2/\eta$, signaling strong implicit regularization by SGD. The stabilization level rises steeply with batch size, so even modest changes in b materially affect the final curvature of the loss landscape of the solution
- Large-batch regime ($b \geq b_c$):** the growth of λ_{\max} with b becomes much slower and the curve asymptotically approaches $2/\eta$ from below, mirroring full-batch gradient descent and reflecting weak implicit regularization.

The *critical batch size* b_c is therefore the point at which the training dynamics cross over into a full-batch-like regime. Works as Zhang et al. (2024) study the following notion of *critical batch size*: "the point beyond which increasing batch size may result in computational efficiency degradation". Likewise, works focusing on generalization performance depending on the batch size (Masters & Luschi, 2018) identify a cut-off batch sizes above which test performance degrades significantly. We conjecture there may be a relation between these quantities and leave a systematic investigation to future work.

I.3 WHY $2/\eta - C/b^\alpha$ FAILS.

From linear stability analyses near the manifold of minima (Wu et al., 2018; Ma & Ying, 2021; Granzio et al., 2021; Mulayoff & Michaeli, 2024) or random matrix theory (together with the fact that we have *Batch Sharpness* stabilize at $2/\eta$) one would expect to have a law of the form $\lambda_{\max} \approx 2/\eta - O(1/b^\alpha)$. Log-log plots of the gap $2/\eta - \lambda_{\max}$ in Figure 13b shows no robust power law (for the lack of any linear dependency), invalidating this prediction (see also Figures 14-20). Importantly, this does not invalidate the findings of those theories, instead showcases the insufficiency of a *static*

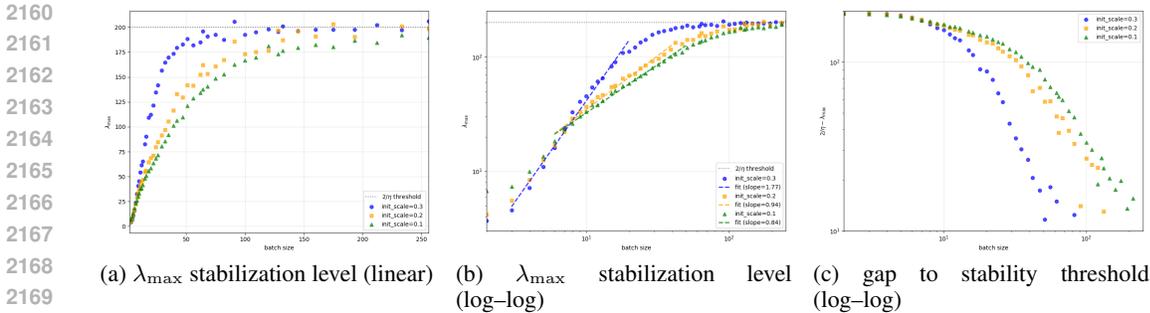


Figure 14: **Effect of weight-scale at initialization on the EOSS stabilization of λ_{\max} .** We train the same network and dataset under identical hyperparameters, varying only a global rescaling ($\times 0.1, 0.2, 0.3$ of He) of the initial weights. (a) Final-epoch λ_{\max} as a function of batch size (linear axes). Smaller batches always converge to flatter solutions, yet the absolute level—and the critical batch size at which the curve begins to approach the full-batch limit $2/\eta$ (horizontal dashed line)—shift markedly with the initialization scale. This demonstrates that the landscape geometry at convergence is already seeded by early-training choices. (b) Same data in log-log scale. The three curves exhibit distinct slopes, ruling out a single power-law exponent and confirming strong path-dependence. Linear fit is provided to the linear portion (c) Log-log plot of the gap, $2/\eta - \lambda_{\max}$. The absence of a straight line contradicts the prediction $2/\eta - \lambda_{\max} \propto b^{-\alpha}$ that follows from linear stability analyses near a minimum.

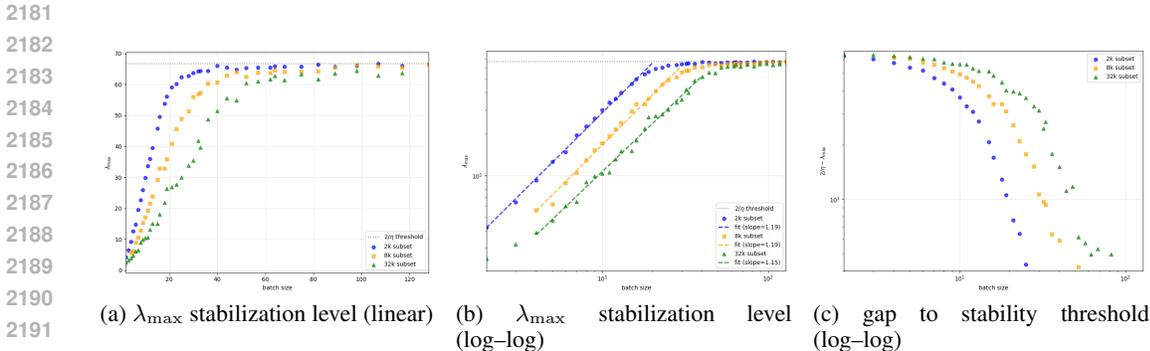


Figure 15: **Varying dataset size alters the EOSS plateau of λ_{\max} for a CNN.** We use the same setup as Fig. 14 but instead varying the number of training examples (2k, 8k, 32k). Larger datasets drive λ_{\max} to lower plateaus—i.e. flatter minima—and push the critical batch size (the knee toward the full-batch limit $2/\eta$) to higher b , as expected from b/N scaling. Plateau heights also differ from the MLP results in Fig. 14 or 13, highlighting architectural sensitivity. Panel order and axes mirror Fig. 14; see that caption for sub-plot details.

analysis. Indeed, those estimates are taken from changing the batch size *statically*, without making any training steps. In particular, linear stability analysis does accommodate virtually any law, as long as there is change in alignment between the mini-batch gradients. The fact that the static law does not apply means that there is a change to the alignment also happening. Therefore, as will be discussed further in detail, the fact that these estimates do not apply means that to give faithful description of the loss landscape at convergence one has to undertake an analysis that is path-dependent.

I.4 CONCLUSION & OUTLOOK: WHY PATH-DEPENDENCE MATTERS

With all of the above, we arrive at a **negative answer** to the question posed at the start:

There is no single, path-independent law that fixes the stabilization level of λ_{\max} from basic hyper-parameters alone.

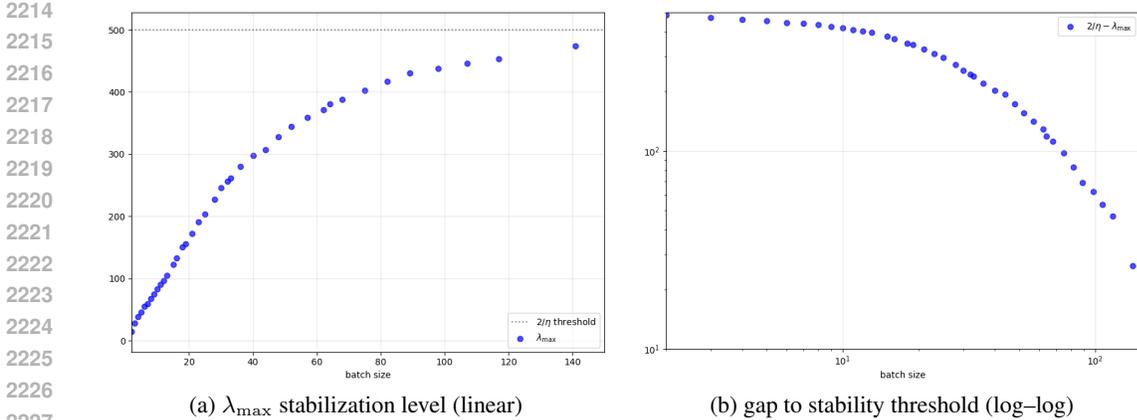


Figure 16: **Level of stabilization of λ_{\max} .** Same setup as Fig. 13 but the initial weights are rescaled by 1/3; see Fig. 14 for the broader effect of initialization. (See Fig. 13 for sub-plot explanations.)

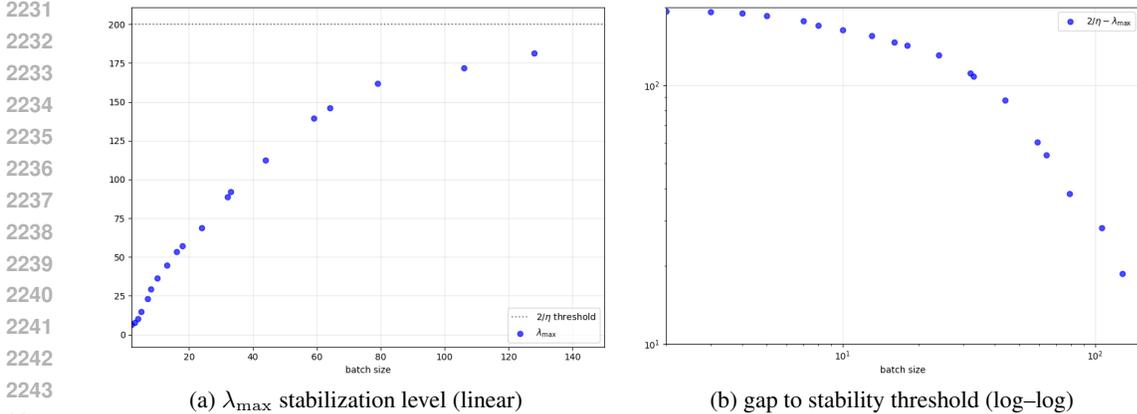


Figure 17: **Level of stabilization of λ_{\max} .** Identical to Fig. 13 except for a larger step size of 0.01. (See Fig. 13 for sub-plot explanations.)

I.5 IMPLICATIONS AND OPEN QUESTIONS

The findings above lead to the following main conclusions.

- (C1) λ_{\max} is *not* the stability limiter for mini-batch training. *Batch Sharpness* governs EOSS; λ_{\max} follows. λ_{\max} is capped from above by the value it reaches at the phase transition characterized by *Batch Sharpness* reaching $2/\eta$. This and Facts 1—3 above imply that:

The stabilization of λ_{\max} is a *by-product* of EOSS, not the quantity that governs it.

- (C2) *A theory of λ_{\max} has to account for the correct progressive sharpening.* By fixing the model and changing batch size b , the gap between the maximal eigenvalue of $\mathbb{E}[\lambda_{\max}(\mathcal{H}(L_B))]$ and $\lambda_{\max} = \lambda_{\max}(\mathcal{H}) = \lambda_{\max}\mathbb{E}[\mathcal{H}(L_B)]$ scales as $1/b$. Any theory that keeps the parameter vector fixed and only varies b , or anyways leads to a power law, misses the path-dependent descent that determines where training arrives and where λ_{\max} stabilizes. Facts 3 and 6 thus imply that analysis of λ_{\max} is insufficient *if* it does not account for (1) the precise and correct effect of progressive sharpening on the higher moments of the Hessian and (2) the correct alignment between mini-batch steps and Hessians.

Quantifying the plateau of λ_{\max} is thus still an (important) open problem. A complete account will require a dynamical theory *through* the progressive-sharpening phase and beyond. Not just properties at its endpoint as for full-batch methods.

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282

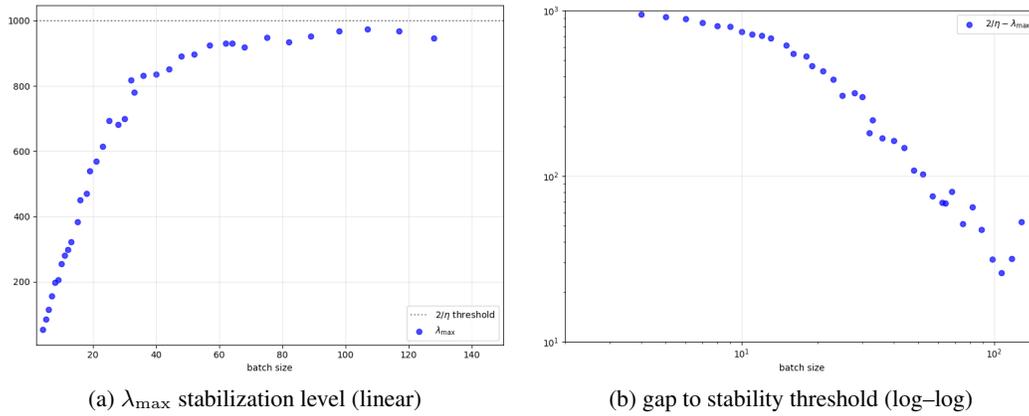


Figure 18: **Level of stabilization of λ_{\max} .** Baseline network trained on a 32k-subset of CIFAR-10 subset with step size 0.002. (See Fig. 13 for sub-plot explanations.)

2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300

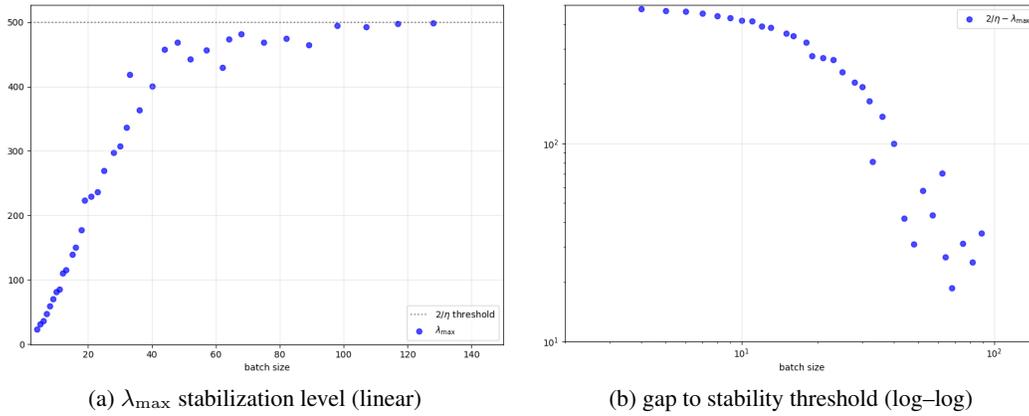


Figure 19: **Level of stabilization of λ_{\max} .** Deeper MLP (the `m1p_1`: 4 hidden layers, width 512) on the 8k-subset, step size 0.004. (See Fig. 13 for sub-plot explanations.)

2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319

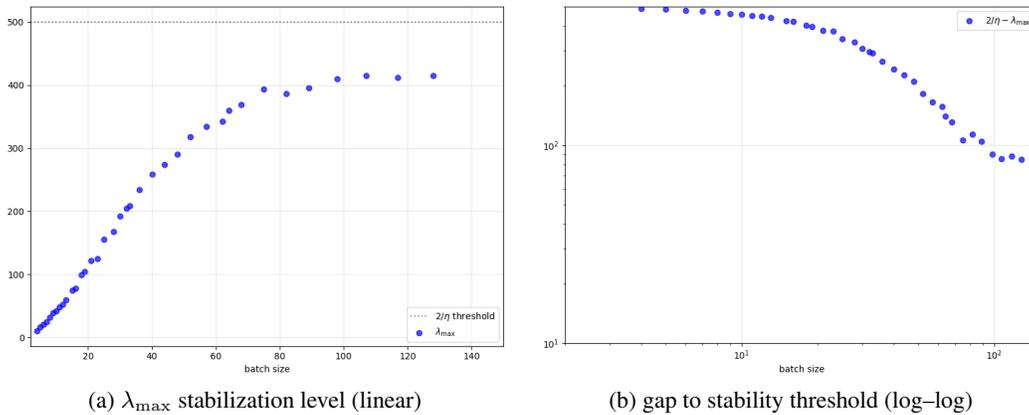


Figure 20: **Level of stabilization of λ_{\max} .** Same deeper MLP as in Fig. 19 but trained on a 32k subset. (See Fig. 13 for sub-plot explanations.)

2320
2321

J IMPLICATIONS: HOW NOISE-INJECTED GD DIFFERS FROM SGD

SGD vs. Noisy Gradient Descent. A common belief is that SGD’s regularization stems from its “noisy” gradients, which find flatter minima. However, our analysis points to the “noisy” Hessians as crucial. To test this, we compare mini-batch SGD (batch size 16) against three noisy GD variants: (see details in Appendix J)

- *Anisotropic Sampling Noise:* Gaussian reweighting on the samples (Wu et al., 2020), which is different from SGD but maintains the mini-batch structure (and injects noise in the Hessians).
- *Diagonal Noise:* Gaussian noise restricted to the diagonal part of the SGD noise covariance (Zhu et al., 2019).
- *Isotropic Noise:* Gaussian noise with isotropic covariance (Zhu et al., 2019).
- *SDE dynamics integration* (Li et al., 2017)

As shown in Figure 8, only noise which maintains the higher moments of the Hessian(s) (and thus preserves the mini-batch landscape structure) leads to an EOSS-like regime with λ_{\max} stabilizing well below $2/\eta$. More generic (e.g., diagonal or isotropic) noise fails to reproduce this behavior. These experiments suggest that stability thresholds differ fundamentally between mini-batch SGD (governed by *Batch Sharpness*) and noise-injected GD (governed by λ_{\max}). Notably, these results are consistent with the findings of Zhu et al. (2019)—although their focus is on generalization. Unsurprisingly, in the case in which the noise affects only the gradients—not the Hessians—indeed, EOSS comes for $\lambda_{\max} = 2/\eta$ as for GD (Ma & Ying, 2021; Mulayoff & Michaeli, 2024). Even in the quadratic setting, the appearance of *Type-1* oscillations and *GNI* are not affected by the structure and distribution of the Hessian on the mini-batches, see Appendix E. The stability threshold, however, is affected. It depends on the Hessian’s higher moments, see Theorem 1 or (Ma & Ying, 2021; Mulayoff & Michaeli, 2024).

Challenges for SDE Modeling. Classical analyses of neural network optimization often assume a single, static landscape: (i) **Online** perspective, modeling each step’s gradient as a noisy unbiased estimator of the expected gradient, or (ii) **Offline** perspective, treating the dataset as fixed and SGD as noisy GD on the empirical loss. In both views, it is the *full-batch* Hessian that supposedly drives curvature. Our results instead highlight that each update sees a Hessian $\mathcal{H}(L_{\mathbf{B}})$ that generally differs significantly from \mathcal{H} , leading to *Batch Sharpness* stabilizing at $2/\eta$ when λ_{\max} is smaller.

Standard SDE—or analogous—approximations of SGD cannot thus describe the location of convergence of SGD or its behavior for neural networks under the assumption of progressive sharpening. Indeed, they typically ignore any statistics of the Hessians except for the mean.

Prior works already note limitations of SDE-based approaches for SGD implicit regularization: they may be mathematically ill-posed (Yaida, 2018), fail except under restrictive conditions (Li et al., 2021), converge to qualitatively different minima (HaoChen et al., 2020), or miss higher-order effects (Damian et al., 2021; Li et al., 2022). Recent discrete analyses (Smith et al., 2021; Benevenuto, 2023; Roberts, 2021) attempt to address some of these issues. Nonetheless, our findings expose a deeper gap: when batch sizes are small, the *geometry of the mini-batch Hessian* differs markedly from that of the full-batch, altering both eigenvalues and eigenvector alignments. Conventional SDE models, which assume a static or average Hessian, cannot easily capture these rapid fluctuations.

J.1 NOISY GD

We are running a number of noisy GD implementations.

J.1.1 NOISY GD WITH ANISOTROPIC NOISE (GAUSSIAN RESAMPLING)

This version of noisy GD essentially preserves the mini-batch landscape structure by averaging the landscapes using Gaussian sampling noise. In particular, it takes a Gaussian sampling vector with the same first and second moments as the sampling vector of SGD. Now, this trivially forces the expectation of the mini-batch Hessians to be the same between SGD and Gaussian resampling (and essentially equal to the full-batch Hessians. Importantly, though, this also makes the covariance of the mini-batch Hessians to be the same between SGD and GD with Gaussian resampling noise (as per linearity of the mini-batch Hessians in the weights of the sampling vector). Together with the

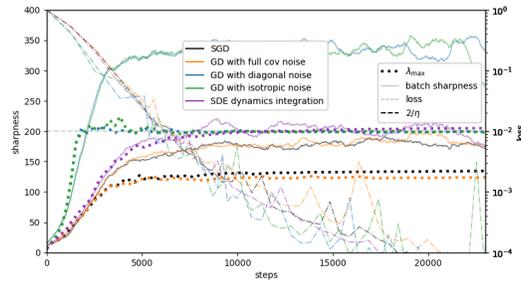


Figure 21: Version of Figure 8, with the loss curves added.

fact that GD with the Gaussian resampling behaves in the same manner as SGD from the point of view of stability, *Batch Sharpness*, and suppression of λ_{\max} —it is an indicator of the fact that it is the higher moments of the mini-batch Hessians that determine the dynamics SGD; and it is indeed the noise in the Hessians that creates the instability regime of EoSS and its consequences. As a weaker consequence, it also preserves the covariance of the noise of SGD.

For implementation details refer to [Wu et al. \(2020\)](#). In summary, we re-draw the sampling vector at each step with the corresponding covariance.

J.1.2 NOISY GD WITH DIAGONAL NOISE

This implementation follows [Zhu et al. \(2019\)](#) — it recreates what they refer to as "GLD diagonal". This is essentially noisy GD with the noise covariance being equal to the diagonal of the covariance of the noise produced by SGD. This preserves each parameter’s marginal variance while ignoring off-diagonal correlations. Conceptually, we are approximating SGD’s noise by $\mathcal{N}(0, \frac{1}{b} \text{diag}(\Sigma(\theta)))$ and add it to the full-batch gradient before the optimizer step. Essentially, this is one step further from a true SGD than the aforementioned GD with anisotropic noise. In particular, it does not preserve the mini-batch landscape structure. As a result, the behavior of GD with diagonal noise differs from SGD from the point of view of λ_{\max} stabilizing below $2/\eta$, and instead stabilizing at $2/\eta$. We refer the reader to ([Zhu et al., 2019](#)) for the details of implementation. In our implementation, we compute the diagonal of the covariance every 30 steps and reuse it on those 30 steps (as it is too computationally expensive to compute it at every step).

J.1.3 NOISY GD WITH ISOTROPIC NOISE

This implementation follows [Zhu et al. \(2019\)](#) — it recreates what they refer to as "GLD dynamic". This is essentially noisy GD with the noise covariance being identity (hence the "isotropic"), scaled such that the magnitude of the noise coincides with that of SGD. That is, this is isotropic gradient noise that matches the average variance of SGD noise but ignores both parameter-wise variability and correlations. Conceptually, we are approximating SGD’s noise by $\mathcal{N}(0, \frac{\sigma^2}{b} I)$ add it to the full-batch gradient before the optimizer step, where $\sigma^2 = \frac{\text{tr}(\Sigma)}{d}$ is the mean per-parameter variance from the per-sample gradient covariance Σ , b is the target batch size, and d is the number of parameters. This is one step "further" from SGD than the noisy GD with diagonal noise. Consequently, this sort of noisy GD does not preserve the regularization effect of SGD on λ_{\max} either.

J.2 SDE

We are taking the standard SDE approximation of SGD: (see e.g. [Li et al. \(2018\)](#))

$$d\theta_t = -\nabla f(\theta_t) dt + \sqrt{\eta} \Sigma^{1/2}(\theta_t) dW_t$$

where dW_t is the standard d -dimensional Wiener process, and Σ is the covariance matrix of mini-batch gradients.

To simulate its dynamics, we are using the Euler–Maruyama discretization with a step size of 0.0005, chosen to be sufficiently small compared to η (1/20th of $\eta = 0.01$ in this example). In [Figure 22](#) we are showing a number of sample paths of the SDE trajectory illustrate the similarity in the properties

2430
 2431
 2432
 2433
 2434
 2435
 2436
 2437
 2438
 2439
 2440
 2441
 2442
 2443
 2444
 2445
 2446
 2447
 2448
 2449
 2450
 2451
 2452
 2453
 2454
 2455
 2456
 2457
 2458
 2459
 2460
 2461
 2462
 2463
 2464
 2465
 2466
 2467
 2468
 2469
 2470
 2471
 2472
 2473
 2474
 2475
 2476
 2477
 2478
 2479
 2480
 2481
 2482
 2483

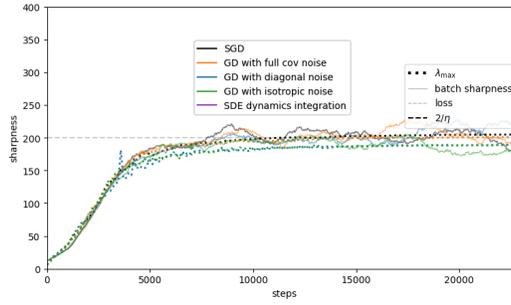


Figure 22: **SDE sample paths** Multiple realizations of SDE trajectory to showcase the similarity of the solutions found by SDE dynamics

of the solutions found by those dynamics – in particular, that λ_{\max} stabilizes around $2/\eta$, rather than below as it does for SGD dynamics. In all the experiments, batch size is 16, and η is 0.01.

2484 K ALIGNMENT

2485
2486 The stability of mini-batch SGD is governed by the geometry of the mini-batch loss landscape, rather
2487 than solely the full-batch landscape. As further discussed in Sections B, 5, L, M, the stability of SGD
2488 depends not only on the magnitude of mini-batch Hessians but also critically on their alignment
2489 (both pairwise and with the loss gradients). This appendix offers a limited characterization of the
2490 alignment structure relevant to mini-batch Hessians.

2491 We approximate both the full-batch and mini-batch Hessians by their Gauss–Newton matrices, an
2492 approximation commonly used in analyses of SGD (e.g., (Wu et al., 2018; Ma & Ying, 2021; Mu-
2493 layoff & Michaeli, 2024)), valid at convergence, and supported empirically (e.g., (Papayan, 2019)
2494 and Appendix P). Concretely,

$$2495 \quad \nabla_{\theta}^2 L_B(\theta) \approx \frac{1}{b} \sum_{i=1}^B J_i^{\top} H_{z,i} J_i$$

2496 where J_i is the Jacobian of the model output with respect to θ and $H_{z,i}$ is the loss Hessian with
2497 respect to the output evaluated at the i -th sample. We work with MSE and, for simplicity, consider
2498 a two-class setting (i.e., a single model output), which yields:

$$2499 \quad \nabla_{\theta}^2 L_B(\theta) \approx \frac{1}{b} \sum_{i=1}^b \nabla_{\theta} f_i \nabla_{\theta} f_i^{\top}$$

2500 where ∇f_i is the per-sample model gradient.

2501 Under this structure, properties of mini-batch Hessians—including their top eigenvalues and the
2502 cross-batch alignment/commutativity—are controlled (though not fully determined) by the pairwise
2503 alignment of the per-sample model gradients. We report the empirical distributions of these pair-
2504 wise alignments in Figure 23 (MLP) and Figure 24 (CNN), plotting pairwise dot products, plotting
2505 pairwise dot products, cosine similarity and individual norms. The evolution of the distributions
2506 throughout the training indicates the effects of progressive sharpening—for example, the growth of
2507 norm of model gradients corresponds to the increase of λ_{\max} and λ_{\max}^b , see Appendix M. A com-
2508 plete description of the dynamics would require a precise account of progressive sharpening and
2509 falls outside the scope of this work.

2510 The important observation for our purposes is that per-sample model gradients are only *weakly*
2511 aligned, with cosine similarities clustered around 0.1—which is still much higher than random d -
2512 dimensional vectors would have. Two immediate implications follow. First, mini-batch Hessians are
2513 generically non-commuting (as the model gradients are not orthogonal or completely collinear)—an
2514 aspect that matters for linear stochastic stability (Appendix L). L. Second, if we fix the parameters
2515 θ and vary the batch size b , the eigenspaces of the mini-batch Hessian mix gradually, which induces
2516 a gap between *Batch Sharpness* and λ_{\max}^b (see Appendix M). We leave full characterization of the
2517 mini-batch Hessians, which would depend on the higher moments of model gradients for future
2518 work. Still, these observations underscore that a comprehensive, training-time characterization of
2519 the structure of mini-batch Hessians is an important future direction of research for understanding
2520 SGD dynamics.

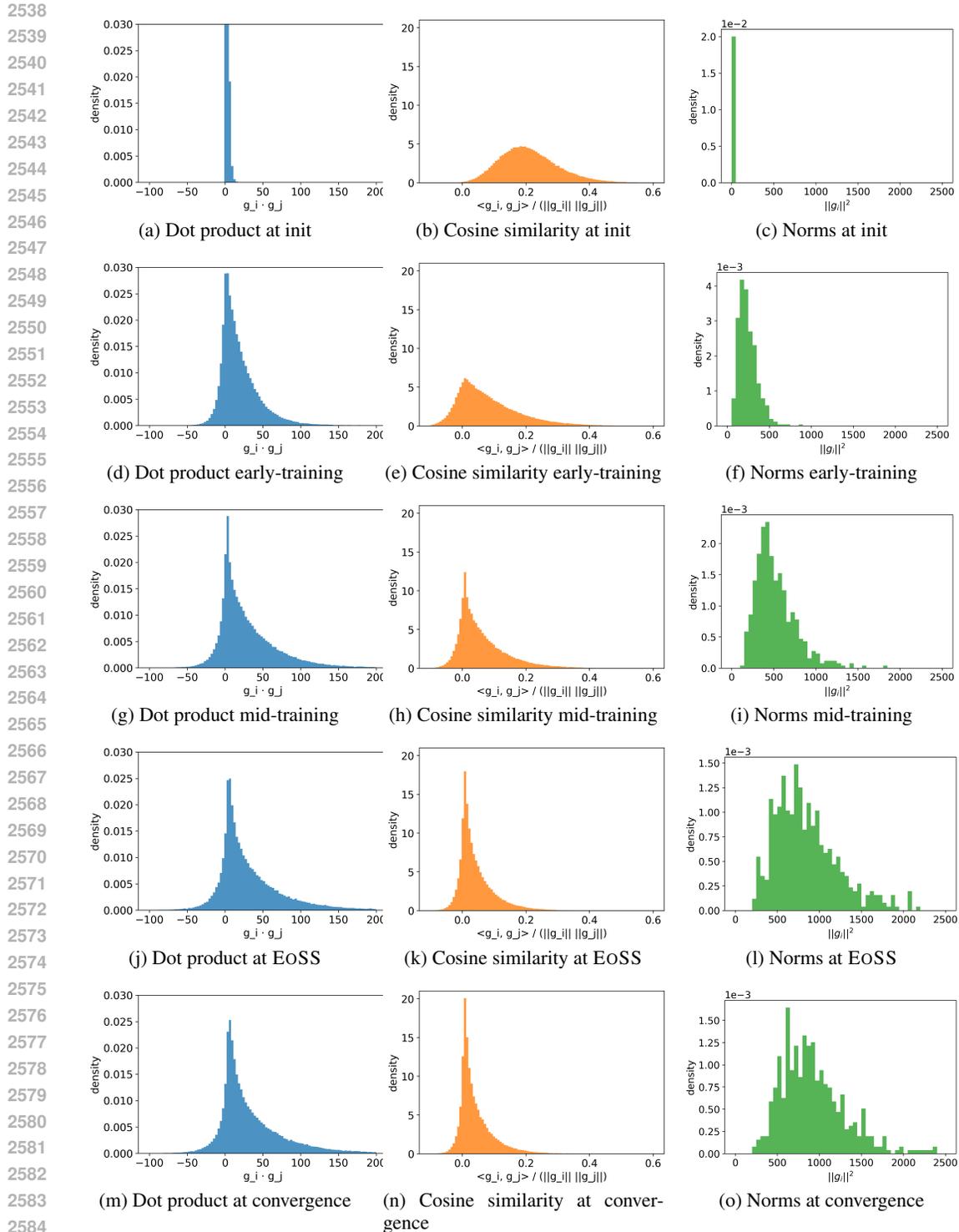


Figure 23: **Model gradients alignment.** Pairwise alignment between model gradients forming the Hessian as the training progresses. We show dot products ($i \neq j$), cosine similarities and the squared norms of the model gradients. Each row corresponds to a stage of training—from initialization, to mid-training (during progressive sharpening), to the later stages (at EOSS and convergence). Notice the gradients become weakly aligned throughout the training (with the cosine similarities clustered around 0.1), but not completely orthogonal, as it would have been with random vectors. MLP, CIFAR-8k (2 classes), $\eta = 0.02$, batch size 32

2592
 2593
 2594
 2595
 2596
 2597
 2598
 2599
 2600
 2601
 2602
 2603
 2604
 2605
 2606
 2607
 2608
 2609
 2610
 2611
 2612
 2613
 2614
 2615
 2616
 2617
 2618
 2619
 2620
 2621
 2622
 2623
 2624
 2625
 2626
 2627
 2628
 2629
 2630
 2631
 2632
 2633
 2634
 2635
 2636
 2637
 2638
 2639
 2640
 2641
 2642
 2643
 2644
 2645

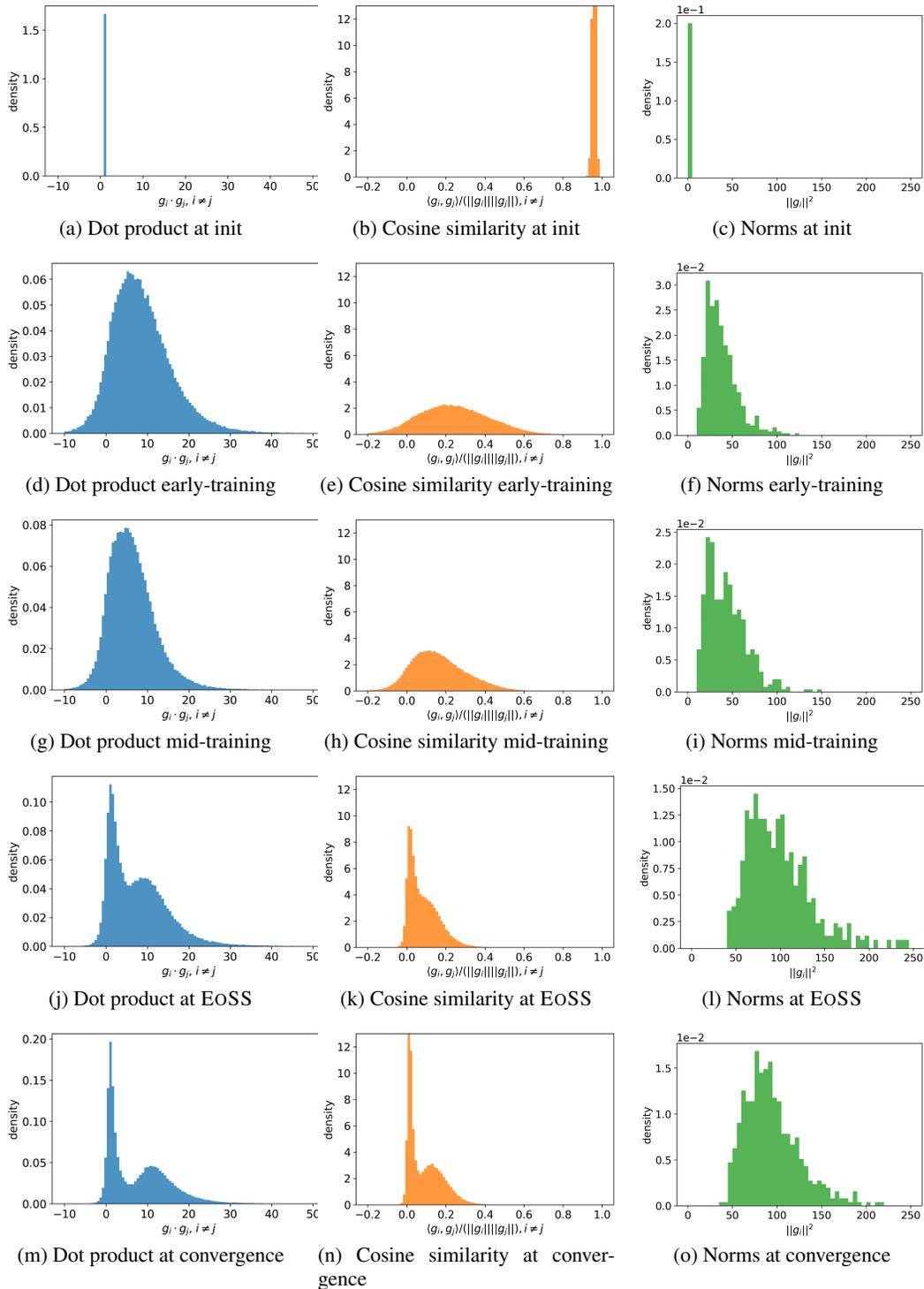


Figure 24: **Model gradients alignment.** Same setup as Figure 23, but for CNN. CIFAR-8k (2 classes), $\eta = 0.05$, batch size 8.

L LINEAR STOCHASTIC STABILITY

This appendix extends the discussion from Section B, particularly regarding works addressing SGD stability from a linear stability viewpoint, including an analysis of the behavior of the bound introduced by Wu et al. (2018).

L.1 NOTIONS OF LINEAR STABILITY

As the discussion in Section B highlights, we need two conditions to establish a regime of instability in mini-batch training: (i) a valid notion of stability as per Definition 4—an inequality whose violation leads to divergence—and (ii) empirical saturation of this stability notion during SGD training, with continued training as the condition remains at saturation.

Wu et al. (2018) were the first to analyze linear stability of SGD, establishing a *sufficient* (but, in general, necessary) condition for SGD stability. In particular, they prove that mini-batch gradient descent is stable when

$$\lambda_{\max} \left((I - \eta\mathcal{H})^2 + \frac{\eta^2}{b} \mathbb{E}[\mathcal{H}_i^2 - \mathcal{H}^2] \right) = \lambda_{\max} (\mathbb{E}_B [(I - \eta\mathcal{H}_B)^2]) \leq 1. \quad (37)$$

This criterion upper-bounds the spectral radius of the second-moment update operator. Since this condition is only sufficient (and necessary solely when $d = 1$), it does not strictly satisfy our criteria for a valid stability notion (Definition 4). Importantly, while Wu et al. (2018) explicitly note this limitation, nothing a priori excludes this bound from being *empirically* tight—after all, EOSS is fundamentally an empirical phenomenon. We show here that this criterion does not govern the EOSS—that is, that the "not necessary" part is not vacuous.

Necessary stability conditions. Conditions derived from linear stochastic stability theory that are indeed valid stability notions often suffer from computational intractability. Indeed, Ma & Ying (2021) proved that mini-batch gradient descent is 2nd order linearly stable *if and only if* the operator

$$T_k := \mathbb{E}_B \left[(I - \eta H(L_B))^{\otimes 2} \right] \quad (38)$$

is a contraction on the cone of PSD matrices. Importantly for us, this a necessary condition for stability, which would constitute a notion of stability. While this condition is necessary and would represent a valid stability notion, it is computationally infeasible in high-dimensional neural networks, as it involves operations on $d^2 \times d^2$ tensors (making even tensor-vector product unfeasible).

Mulayoff & Michaeli (2024) showed that the PSD condition is inactive, which reduces the criterion to one on a spectral norm of this operator. Moreover, they express this notion in the form of (notion of curvature) \leq (step-size dependent threshold). Although potentially useful to find the corresponding maximum stable learning rate, this reformulation did not solve the incomputability problem. (Mulayoff & Michaeli, 2024) also construct elegant lower bounds, which therefore also serve as a necessary condition for stability, and thus a valid notion of stability. However, as their empirical results show, these bounds never saturate, and thus do not effectively capture the empirical presence of an instability regime in mini-batch training.

L.2 EMPIRICAL BEHAVIOR OF WU ET AL. (2018) CRITERION

Always above 1. In neural networks negative Hessian eigenvalues are typically present, thus the quantity in Equation (37), which we term *second moment contraction*, is always bigger than 1 (Figure 25). Ideally, when "the phase transition at EOSS" happens this quantity keeps being bigger than 1, but the highest singular value becomes the biggest eigenvalue instead of the smallest. In the deterministic full-batch algorithm case this can be seen cleanly:

$$\begin{aligned} \lambda_{\max} (\mathbb{E}_B [(I - \eta\mathcal{H}_B)^2]) &= \lambda_{\max} ((I - \eta\mathcal{H})^2) = \\ &\begin{cases} 1 - \eta\lambda_{\min}(\mathcal{H}) & \sim 1 + \epsilon_1 & \text{when } \lambda_{\max} \leq 2/\eta \text{ and } \lambda_{\min} < 0 \\ |\eta\lambda_{\max} - 1| & \sim 1 + \epsilon_2 & \text{when } \lambda_{\max} \geq 2/\eta \text{ and } \lambda_{\min} < 2 - \eta\lambda_{\max} \leq 0. \end{cases} \end{aligned}$$

Thus we can think of plotting the quantity

$$\lambda_{\max} (-2\eta\mathcal{H} + \eta^2 \mathbb{E}_B [\mathcal{H}(L_B)^2]) \leq 0. \quad (39)$$

2700
 2701
 2702
 2703
 2704
 2705
 2706
 2707
 2708
 2709
 2710
 2711
 2712
 2713
 2714
 2715
 2716
 2717
 2718
 2719
 2720
 2721
 2722
 2723
 2724
 2725
 2726
 2727
 2728
 2729
 2730
 2731
 2732
 2733
 2734
 2735
 2736
 2737
 2738
 2739
 2740
 2741
 2742
 2743
 2744
 2745
 2746
 2747
 2748
 2749
 2750
 2751
 2752
 2753

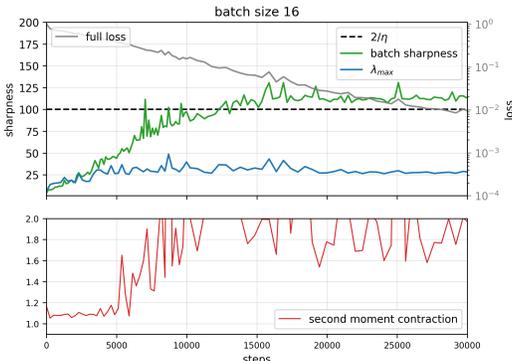


Figure 25: *Linear stochastic stability*. Tracking the condition of equation (37) (“second-moment contraction”) during training. The value stays strictly above 1.

As illustrated in Figure 26, the condition (39) is initially slightly above zero (due to negative eigenvalues). Then it undergoes the aforementioned phase transition and starts growing—although the exact moment of that transition is unclear. After that it seems to stabilize around the time that λ_{\max} stabilizes, and *Batch Sharpness* stabilizes around $2/\eta$.

Not a stability measure empirically. Empirically, the criterion (39) behaves in a way that precludes it from being a quantity that governs stability in a EOS-like fashion, as evidenced by the experiments in Figure 26:

- **The bound is not tight.** In particular, if *empirically* the *second moment contraction* was governing stability of SGD training, the condition of equation (37) would only be violated by a small margin, just like is the case with that condition in the full-batch GD—equivalently, as it is the case with λ_{\max} being just slightly above $2/\eta$ during GD training. Instead, *second moment contraction* hovers at 2 (equivalently, the quantity of Equation (39) hovering around 1).
- **No up and down oscillations.** The higher-order-term driven EOS stabilization mechanism of Damian et al. (2023) and the dynamics in Cohen et al. (2024) prescribe a notion of stability to go up and down around the stability threshold. On top of the above point of being significantly above the threshold level, *second moment contraction* does not oscillate in a way prescribed by a stabilization that’s based on higher order terms.
- **Inconsistent level of stabilization.** Finally, a notion that would govern SGD dynamics would present a consistent level of stabilization, independent of hyperparameters. This does not happen in the case of *second moment contraction*.

L.3 IMPLICATIONS

Alignment matters. The reason why the condition of Wu et al. (2018) (equation (37)) is only a sufficient one, while the condition of Ma & Ying (2021) (equation (7)) is necessary and sufficient, is that the mini-batch Hessians are not commuting/not simultaneously diagonalizable. In particular, they would have been simultaneously diagonalizable if either the model gradients forming their Gauss-Newton approximation were either all the same or all orthogonal, which is hypothetically possible as we have $N \ll d$. Now, in Appendix K it was show that neither is the case—they are misaligned, but not completely orthogonal. Now, the condition (37) not being a governing quantity of EOSS, and not being tight as an upper bound on instability condition, is evidence that this not-complete misalignment has non-trivial effects. That is, unlike in deterministic full-batch gradient descent settings, instability in SGD is dependent on the alignment between notions of curvatures. Therefore, a true notion of stability has to involve a notion of alignment, not only the magnitude of curvature—and, correspondingly, *Batch Sharpness* does consider the alignment, as opposed to, for example, λ_{\max} or λ_{\max}^b .

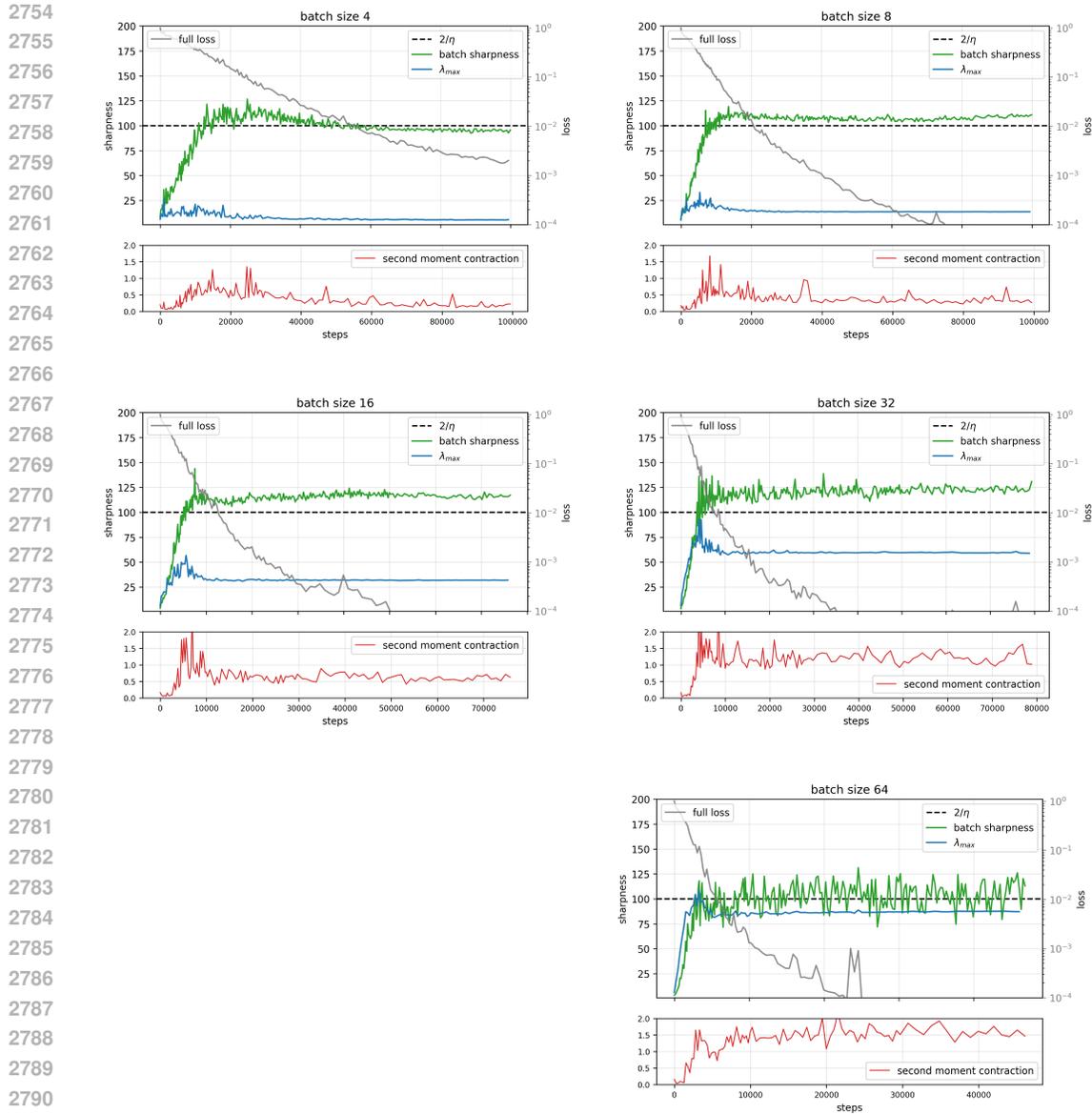


Figure 26: Tracking the condition of equation (39) (the one without the identity), *second moment contraction*, for different batch sizes. MLP, CIFAR-8k, $\eta = 0.02$. Note how the stabilization levels is significantly above the threshold, and inconsistent across batch sizes.

Importance of instability (not "stability"). We define EOSS as being at the edge of instability. This quantity of equation (37) smaller than 1 is *only* a sufficient condition for stability. The fact that breaking Eq.(37) is not enough for assessing the behavior of SGD is a further proof that what matters is the instability, not the stability. What matters is an inequality that implies divergence if broken, not that implies convergence if satisfied.

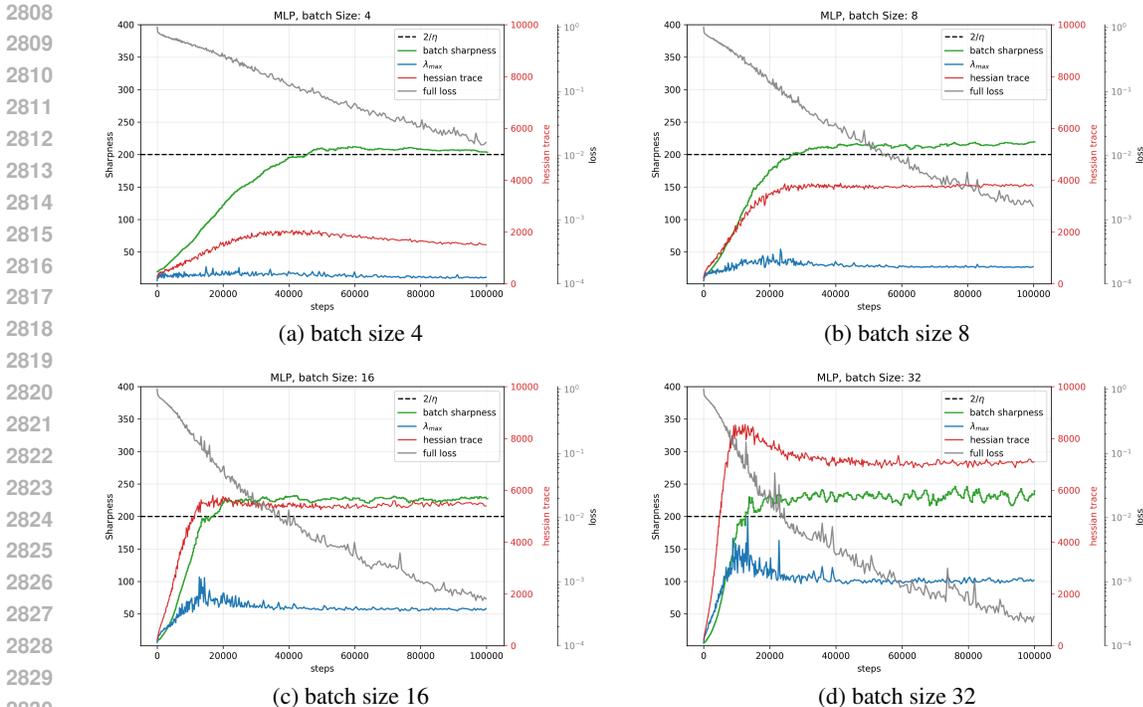


Figure 27: **Trace of the Hessian.** We plot the trace of the full-batch loss Hessian (red), together with the usual *Batch Sharpness* (green) and λ_{\max} (blue). Notice that the scale of the trace of the Hessian is much bigger than the rest of the quantities, and it follows the axis on the right (in particular, has no particular relation to $2/\eta$). The plots showcase that trace behaves in a similar manner as λ_{\max} —its level of stabilization is highly dependent on the batch size, it raises as long as *Batch Sharpness* is rising, and it is stabilizes as batch sharpness stabilizes. Here, we are doing experiments with MLP on CIFAR-10-8k and $\eta = 0.01$

M OTHER QUANTITIES OF SGD DYNAMICS

In this Appendix we explore other quantities that describe the SGD dynamics, and discuss their role from the point of view of governing stability. In particular, we are covering the following quantities: trace of full-batch loss Hessian, λ_{\max}^b (average max eigenvalues of mini-batch Hessians) and a modified version of *Batch Sharpness*.

M.1 TRACE OF THE LOSS HESSIAN

A number of works [Ma & Ying \(2021\)](#); [Wu & Su \(2023\)](#); [Agarwala & Pennington \(2024\)](#) have linked the trace of the full-batch loss Hessian to implicit regularization by SGD. We plot in Figure 27 and 28: λ_{\max} , *Batch Sharpness*, and the trace of the Hessian along the training for a variety of models and batch sizes. We observe here that trace of the Hessian behaves very similarly to the previously studied λ_{\max} . In particular, it does not have a consistent stabilization level, and depends significantly on the batch size—with smaller batch sizes leading to lower stabilization level of the trace (aka flatter solutions). Also analogous to λ_{\max} , it undergoes progressive sharpening, as long as *Batch Sharpness* is under $2/\eta$. Analogously, the stabilization of *Batch Sharpness* leads to stabilization of the trace. All of this showcases that trace of the Hessian is not the quantity that governs stability of the SGD dynamics. Yet, it might be a useful indicator of the end phase of progressive sharpening—in the potential situation when we have λ_{\max} stabilize, but other eigenvalues continue growing, as illustrated, for example, in [Cohen et al. \(2024\)](#).

It is noteworthy that, in the context of MSE loss combined with piecewise-linear activation functions (e.g., ReLU), the trace of the full-batch loss Hessian coincides with the trace of its Gauss–Newton approximation. Furthermore, under MSE loss, the trace of the Gauss–Newton matrix is equal to the

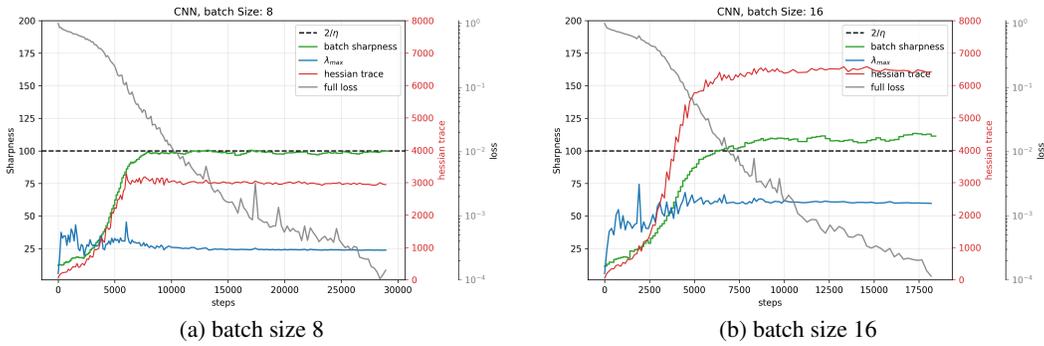


Figure 28: **Trace of the Hessian.** Similar to 27, but for CNN, and with $\eta = 0.02$

trace of the NTK. Consequently, in our setup (MLP with ReLU under MSE) evaluating the trace of the loss Hessian subsumes these cases.

M.2 λ_{\max}^b : EXPECTED HIGHEST EIGENVALUE OF MINI-BATCH HESSIANS

In the early versions of this work we have looked at another promising quantity that we term λ_{\max}^b :

$$\lambda_{\max}^b := \mathbb{E}_{B \sim \mathcal{P}_b} \left[\lambda_{\max}(\mathcal{H}(L_B)) \right].$$

In particular, the significance of this quantity lies in its characterization of the worst-case sharpness of mini-batch loss landscapes. Yet, the reason why this quantity does not govern SGD dynamics arises from the very phenomenon distinguishing SGD dynamics from full-batch gradient descent—the misalignment of the mini-batch Hessians, see K. Specifically, while individual mini-batch Hessians may exhibit considerable sharpness in their individual directions, these directions typically fail to align, preventing the emergence of a single dominant sharp direction. This scenario closely mirrors the behavior of the operator analyzed in Appendix L, illustrating why *Batch Sharpness*, which dictates the stability of SGD dynamics, relies on both the size of the mini-batch Hessians together with their alignment with the step direction.

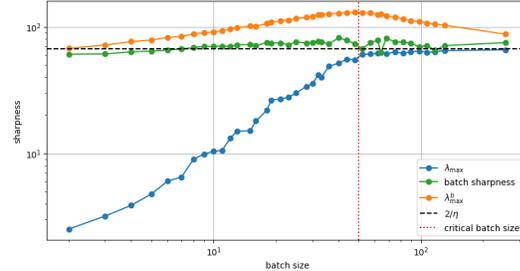


Figure 29: *Final stabilization levels of Batch Sharpness, λ_{\max} and λ_{\max}^b vs batch size.* Only the stabilization level of *Batch Sharpness* does not depend on batch size. Setting: MLP, CIFAR-10-8k, $\eta = 0.03$

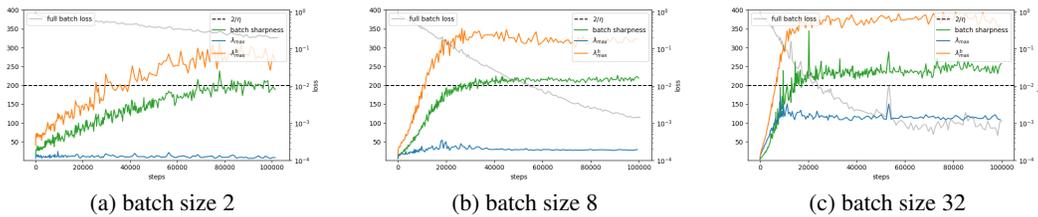


Figure 30: **Behavior of λ_{\max}^b :** λ_{\max}^b stabilizes higher than $2/\eta$, with its stabilization level dependent on the batch size. *Batch Sharpness* and λ_{\max} also shown for comparison. Setting: MLP, $\eta = 0.01$, CIFAR10-8k

Consequently, λ_{\max}^b stabilizes at a level higher than the threshold $2/\eta$ and *Batch Sharpness*. Moreover, the precise stabilization level is sensitive to the chosen batch size, as showcased in Figure 30, with the dependence of the level of stabilization on batch size shown in Figure 29. For additional

experiments with illustration of behavior of λ_{\max}^b refer to Appendix S. Now, this inconsistency of stabilization means that λ_{\max}^b does not govern the stability of SGD, and its stabilization is a by-product of EOSS and *Batch Sharpness* stabilization.

In particular, we establish that:

- (i) λ_{\max}^b also stabilizes.
- (ii) λ_{\max}^b stabilizes at a level ranging between $2/\eta$ and $4/\eta$. The level is lower for very small and very large batch sizes, and higher for intermediate batch sizes.
- (iii) λ_{\max}^b increases concurrently with *Batch Sharpness* and stabilizes simultaneously, indicating insights into the nature of *Batch Sharpness* growth and progressive sharpening.
- (iv) λ_{\max}^b is by construction greater than both λ_{\max} and *Batch Sharpness*. The stabilization of *Batch Sharpness* around $2/\eta$ for SGD and λ_{\max} for GD ensures that λ_{\max}^b stabilizes at or above $2/\eta$ in the EOS/EOSS regime.

Concerning (iv), the inequality $\lambda_{\max}^b \geq \text{Batch Sharpness}$ follows directly from the definition of *Batch Sharpness* as an expectation of Rayleigh quotients. Furthermore, the inequality $\lambda_{\max}^b \geq \lambda_{\max}$ results from the following reasoning. The largest singular value of the Hessian matrix derived from single data points is positive. This observation is crucial in establishing the following well-known property of matrix eigenvalues.

Lemma 6. *Let $m, b \in \mathbb{N}$ and consider m matrices $M_1, M_2, \dots, M_b \in \mathbb{R}^{m \times m}$ satisfying $\lambda_{\max} > |\lambda_{\min}|$. Then, the largest eigenvalue of their sum satisfies*

$$\lambda_{\max} \left(\sum_{i=1}^b M_i \right) \leq \sum_{i=1}^b \lambda_{\max} (M_i) \quad (40)$$

with equality only if all M_i are identical.

This lemma is a direct consequence of the convexity of the operator norm in matrices and the fact that the largest eigenvalue is positive in our setting. In our setting, it implies that with non-simultaneously-diagonalizable matrices, the maximum eigenvalue of the sum is strictly less than the sum of the maximum eigenvalues of the individual matrices. To illustrate, consider eigenvalue sequences for batch sizes that are powers of four, though the result generalizes to any $b_1 < b_2$:

$$\lambda_{\max}^1 > \lambda_{\max}^4 > \lambda_{\max}^{16} > \lambda_{\max}^{64} > \lambda_{\max}^{256} > \dots \quad (41)$$

Importantly, this ordering is the case only for "static" model – i.e. when we take a model, and without changing the weights, evaluate λ_{\max}^b as we change the batch size.

As noted in point (ii) above, this ordering does not hold in the *trained* case, as different batch sizes affect also the progressive sharpening and the nature of the mini-batch Hessians. Specifically, since *Batch Sharpness* stabilizes at $2/\eta$ at EOSS, the level of stabilization of λ_{\max}^b depends on its relation to *Batch Sharpness*. Since the two are quite similar, with the difference that *Batch Sharpness* also takes into account the alignment of the mini-batch landscapes sharpest directions with the mini-batch gradients—the gap between λ_{\max}^b and *Batch Sharpness* is governed precisely by this alignment. As illustrated in Figure 29, the level of stabilization of λ_{\max}^b is similar to that of *Batch Sharpness* for very small and very large batch sizes. For large batch sizes, this result is straightforward, as the dynamics approach full-batch GD, in which all relevant quantities equalize at EOS. Conversely, the small-batch case emerges because, for smaller batch sizes, the mini-batch Hessian (or its Gauss-Newton approximation) comprises averages of only a few per-sample *model* gradient outer products, causing mini-batch gradients to align closely with the largest eigenvalues.

This alignment diminishes as the batch size increases, leading to a widening gap between *Batch Sharpness* and λ_{\max}^b . Intriguingly, our experiments reveal that this gap only widens up to the aforementioned *critical batch size*, which also serves as a switch between SGD and GD dynamics from the point of view of λ_{\max} stabilization. Beyond the *critical batch size* the gap begins to narrow again, as depicted in Figure 29. Clarifying this phenomenon fully would be a key outcome of a comprehensive theory of progressive sharpening and SGD stability.

Another significant consequence of the stabilization of λ_{\max}^b is that it provides insights into the mechanisms underlying the progressive sharpening of *Batch Sharpness*. Specifically, the growth in

Batch Sharpness could be attributed either to a general increase in the sharpness of the mini-batch landscapes or to an increase in alignment between mini-batch Hessians and gradients. Notably, throughout the period of *Batch Sharpness* increase, both λ_{\max}^b and the trace of the loss Hessian consistently rise and stabilize simultaneously with *Batch Sharpness*. This suggests that at least portion of the increase in *Batch Sharpness* arises from the overall sharpening of the mini-batch landscapes, rather than solely from alignment of the mini-batch gradients and Hessians. Consequently, *Batch Sharpness* appears closely linked with the progressive sharpening phenomenon itself, with its eventual stabilization marking the end of progressive sharpening. This points to the fact that *Batch Sharpness* is closely connected to progressive sharpening.

M.3 MODIFIED BATCH SHARPNESS

In the earlier versions of this work we also looked at a modified definition of *Batch Sharpness*:

Definition 7 (Modified *Batch Sharpness*). We call Modified *Batch Sharpness* the quantity defined as

$$\text{Modified } \textit{Batch Sharpness}(\theta) := \frac{\mathbb{E}_{B \sim \mathcal{P}_b} \left[\nabla L_B(\theta)^\top \mathcal{H}(L_B) \nabla L_B(\theta) \right]}{\mathbb{E}_{B \sim \mathcal{P}_b} \left[\|\nabla L_B(\theta)\|^2 \right]}.$$

The difference from the definition of *Batch Sharpness* is that in this one the expectation over batches is taken inside the fraction. The intuition for this quantity comes from a notion of average stability on mini-batch landscapes. That is,

$$\frac{\mathbb{E}_{B \sim \mathcal{P}_b} \left[\nabla L_B(\theta)^\top \mathcal{H}(L_B) \nabla L_B(\theta) \right]}{\mathbb{E}_{B \sim \mathcal{P}_b} \left[\|\nabla L_B(\theta)\|^2 \right]} \leq 2/\eta \iff \mathbb{E} \left[L_B(\theta_{t+1}^B) - L_B(\theta_t) \right] \geq 0$$

where θ_{t+1}^B is the parameters that we are getting if we are stepping on the given mini-batch. This means that *Modified Batch Sharpness* $< 2/\eta$ is equivalent to "on average, the mini-batch loss does not increase when stepping the corresponding landscape". This formulation is an attempt to extend the descent lemma to the mini-batch landscapes that govern the SGD dynamics instead of the descent lemma on the full-batch landscape that govern GD. Empirically, it turns out that *Modified Batch Sharpness* also stabilizes, but its stabilization level is higher than that of the *Batch Sharpness* and therefore $2/\eta$, as illustrated in Figure 31. Moreover, its stabilization level is dependent on the batch size.

Modified *Batch Sharpness* and mini-batch gradients. Importantly, we show in Proposition 5 that Modified *Batch Sharpness* is a valid Instability Criterion and it governs the explosion of the expectation of the norm squared of the mini-batch gradients.

$$\text{Modified } \textit{Batch Sharpness}(\theta_t) > 2/\eta + c\eta \implies \mathbb{E}_{B \sim \mathcal{P}_b} [\|\nabla L_B(\theta_{t+1})\|_2^2] > \mathbb{E}_{B \sim \mathcal{P}_b} [\|\nabla L_B(\theta_t)\|_2^2].$$

N MODIFIED BATCH SHARPNESS IS A VALID INSTABILITY CRITERION

We show here that Modified *Batch Sharpness* (Definition 7) is a valid instability criterion (Definition 4).

*Importantly, while it is a valid instability criterion, it does not stabilize at $2/\eta$ in practice, thus it is not the quantity that self-stabilization tames, but its stabilization is a byproduct of *Batch Sharpness* stabilizing.*

We now compute what the update of the norm of the gradients $\mathbb{E}_i[\|Y_i\|_2^2]$ after one step is. Precisely, with the notations of Appendix G, we are computing here the value of $\mathbb{E}_t \mathbb{E}_i[\|Y_i^{t+1}\|_2^2]$ so the average over the iterations of the update to the quantity \mathcal{C} above. Precisely we here prove the following Proposition.

Proposition 5. *There exists an absolute constant $c > 0$ such that when Modified *Batch Sharpness* $> 2/\eta + c\eta$, then $\mathbb{E}\|\nabla L_B\|_2^2$ increases in size exponentially and the trajectory diverges (is quadratically unstable, see Definition 5).*

3024
 3025
 3026
 3027
 3028
 3029
 3030
 3031
 3032
 3033
 3034
 3035
 3036
 3037
 3038
 3039
 3040
 3041
 3042
 3043
 3044
 3045
 3046
 3047
 3048
 3049
 3050
 3051
 3052
 3053
 3054
 3055
 3056
 3057
 3058
 3059
 3060
 3061
 3062
 3063
 3064
 3065
 3066
 3067
 3068
 3069
 3070
 3071
 3072
 3073
 3074
 3075
 3076
 3077

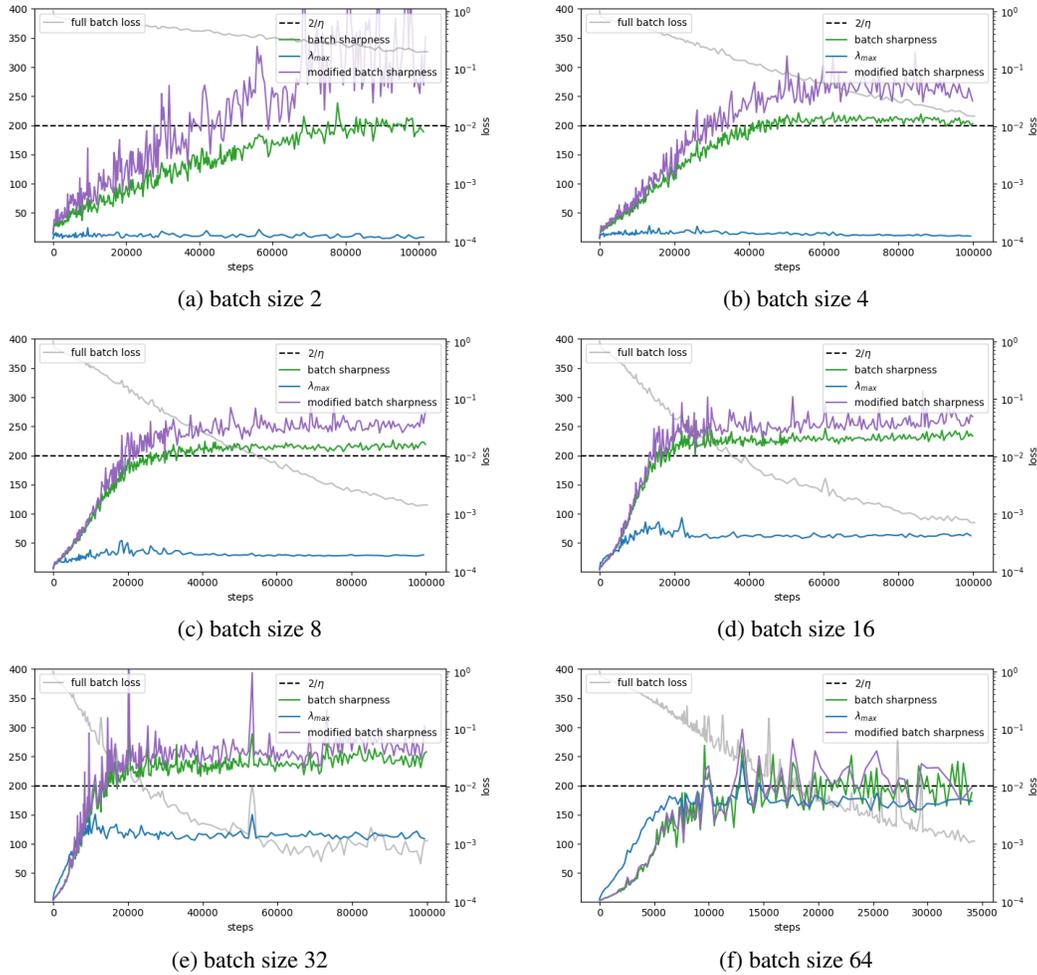


Figure 31: **Modified Batch Sharpness:** Behavior of *Modified Batch Sharpness*, definition of which is similar to *Batch Sharpness*, but with the expectation "inside" the fraction $(\mathbb{E}[\nabla L_B(\theta)^T H_B \nabla L_B(\theta)] / \mathbb{E}[\|\nabla L_B(\theta)\|^2])$. It stabilizes above $2/\eta$ with its stabilization level highly dependent on the batch size.

3078 *Proof.* In the proof we use the notations of Appendix G.

3079
3080 **Step 1: One step on the gradient's second moment** . Remind that the SGD iterate satisfy

$$3082 \theta_{t+1} = \theta_t - \eta Y_{j_t}(\theta_t), \quad i_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{D},$$

3083 and define a *fresh*, independent index j used only for the outer expectation in \mathcal{C}^{t+1} . Because $j \perp i_t$

$$3084 \text{ we may write } Y_i(\theta_{t+1}) = H_i(\theta_{t+1} - x_i) = Y_i(\theta_t) - \eta H_i Y_{j_t}(\theta_t).$$

3085 Squaring, expanding, and averaging over j gives

$$3088 \mathcal{C}^{t+1} = \mathbb{E}_i \|Y_i(\theta_t) - \eta H_i Y_{j_t}(\theta_t)\|^2$$

$$3089 = \mathcal{C}^t - 2\eta \underbrace{\mathbb{E}_{i,j_t}[Y_i(\theta_t)^\top H_i Y_{j_t}(\theta_t)]}_{\text{cross term}} + \eta^2 \underbrace{\mathbb{E}_{i,j_t}[Y_{j_t}(\theta_t)^\top H_i^2 Y_{j_t}(\theta_t)]}_{\text{variance term}}. \quad (42)$$

3093 **Step 2: Decoupling the indices.** Note that

$$3094 2\mathbb{E}_{i,j_t}[Y_i(\theta_t)^\top H_i Y_{j_t}(\theta_t)] = \mathcal{A} - \mathcal{B} - \tilde{\Delta}. \quad (43)$$

3096 This implies that we can rewrite

$$3097 \mathcal{C}^{t+1} = \mathcal{C}^t - \eta(\mathcal{A} - \mathcal{B} - \tilde{\Delta}) + \eta^2 (\text{variance term}). \quad (44)$$

3099 Next note that if we are at the EOSS, then $\mathcal{A} \approx \frac{2}{\eta}(1 + \delta)\mathcal{C}$ for some $\delta \in \mathbb{R}$. This implies that we

$$3100 \text{ can rewrite the term above as}$$

$$3101 \mathcal{C}^{t+1} \approx -(1 + 2\delta)\mathcal{C}^t + \underbrace{\eta\mathcal{B} + \eta\tilde{\Delta} + \eta^2 (\text{variance term})}_{\text{rest}}. \quad (45)$$

3104 Let us now understand the size of the rest, the trajectory diverges if and only if:

$$3105 \eta\mathcal{B} + \eta\tilde{\Delta} + \eta^2 \mathbb{E}_{i,j_t}[Y_{j_t}(\theta_t)^\top H_i^2 Y_{j_t}(\theta_t)] > 2(1 + \delta)\mathcal{C}^t. \quad (46)$$

3108 Next note that by applying Jensen inequality to the term multiplied by η^2 we obtain that

$$3109 \sqrt{\underbrace{\mathbb{E}_{i,j_t}[Y_{j_t}(\theta_t)^\top H_i^2 Y_{j_t}(\theta_t)]}_{\text{variance term}} \cdot \underbrace{\mathbb{E}_i[Y_i(\theta_t)^\top Y_i(\theta_t)]}_{\mathcal{C}}} \geq \underbrace{\mathbb{E}_{i,j_t}[Y_{j_t}(\theta_t)^\top H_i \cdot Y_i(\theta_t)]}_{\mathcal{D}}. \quad (47)$$

3113 **Step 3: Final algebra.** Plugging this above, we obtain that the trajectory diverges when

$$3114 \eta\mathcal{B} + \eta\tilde{\Delta} + \eta^2 \frac{\mathcal{D}^2}{\mathcal{C}} > 2(1 + \delta)\mathcal{C}. \quad (48)$$

3117 Again applying (43) we obtain that this is equivalent to

$$3118 \eta\mathcal{B} + \eta\tilde{\Delta} + \eta^2 \frac{(\mathcal{A} - \mathcal{B} - \tilde{\Delta})^2}{4\mathcal{C}} > 2(1 + \delta)\mathcal{C}. \quad (49)$$

3121 Since $\eta\mathcal{A} = 2(1 + \delta)\mathcal{C}$, then $\eta^2\mathcal{A}^2 = 4(1 + \delta)^2\mathcal{C}^2$ to asking

$$3122 \eta\mathcal{B} + \eta\tilde{\Delta} + \eta^2 \frac{\mathcal{B}^2 + \tilde{\Delta}^2 - 2\mathcal{A}\tilde{\Delta} - 2\mathcal{A}\mathcal{B} + 2\mathcal{B}\tilde{\Delta}}{4\mathcal{C}} > 2(1 + \delta)\mathcal{C} - \frac{4(1 + \delta)^2\mathcal{C}^2}{4\mathcal{C}}. \quad (50)$$

3125 Furthermore, equivalent to asking

$$3126 \eta\mathcal{B} + \eta\tilde{\Delta} - \frac{2(1 + \delta)}{2}\eta\tilde{\Delta} - \frac{2(1 + \delta)}{2}\eta\mathcal{B} + \eta^2 \frac{\mathcal{B}^2 + \tilde{\Delta}^2 + 2\mathcal{B}\tilde{\Delta}}{4\mathcal{C}} > (1 - \delta + \delta^2)\mathcal{C} \quad (51)$$

3129 or, even further simplified

$$3130 \eta\delta(\mathcal{B} + \tilde{\Delta}) + \eta^2 \frac{(\mathcal{B} + \tilde{\Delta})^2}{4\mathcal{C}} > (1 - \delta + \delta^2)\mathcal{C}. \quad (52)$$

We can rewrite this as

$$\eta\delta(\mathcal{A} - 2\mathcal{D}) + \eta^2 \frac{(\mathcal{A} - 2\mathcal{D})^2}{4\mathcal{C}} > (1 - \delta + \delta^2)\mathcal{C}. \quad (53)$$

By plugging, as before, $\eta\mathcal{A} = 2(1 + \delta)\mathcal{C}$ we obtain

$$2\delta(1 + \delta)\mathcal{C} - 2\eta\delta 2\mathcal{D} - 2\eta(1 + \delta)\mathcal{D} + \eta^2 \frac{\mathcal{D}^2}{\mathcal{C}} > (1 - \delta + \delta^2 - (1 + \delta)^2)\mathcal{C} \quad (54)$$

which simplifies as

$$\underbrace{2\eta(1 + 2\delta)\mathcal{D}}_{\mathcal{O}(\eta^2)} - \underbrace{\eta^2 \frac{\mathcal{D}^2}{\mathcal{C}}}_{\mathcal{O}(\eta^4)} < \delta(5 + 2\delta) \underbrace{\mathcal{C}}_{\mathcal{O}_\eta(1)}. \quad (55)$$

Thus there exists a constant $c > 0$, such that if $\delta > c\eta^2$ the trajectory diverges exponentially, if $\delta < c\eta^2$ the trajectory is stable.

□

O HARDWARE & COMPUTE REQUIREMENTS

All experiments were executed on a single NVIDIA A100 GPU (80 GB) with 256 GB of host RAM. The software stack comprises Python 3.12 and PyTorch 2.5.1 (built with the default CUDA tool-chain supplied by the wheel).

Baseline MLP (2M parameters, Section Q) Training for 100k steps on the 8 k-image CIFAR-10 subset finishes in ≈ 5 min wall-clock while computing step sharpness every 8 steps, batch sharpness every 128 steps and λ_{\max} every 256 steps. Peak device memory is 14 GB during ordinary training and ≈ 70 GB while estimating λ_{\max} on a 32k subset, comfortably fitting the 80 GB card.

Algorithmic caveats. We rely on power iteration for λ_{\max} ; while Lanczos would reduce the number of Hessian–vector products, the official PyTorch implementation remains CPU-only. To offset the extra memory incurred by double backward, we cache the first forward pass; batching λ_{\max} is left to future work.

P THE HESSIAN AND THE FISHER INFORMATION MATRIX OVERLAP

In the theoretical analysis of stability of SGD dynamics it is assumed that the loss Hessian can be well-approximated by its Gauss-Newton approximation—in particular, it is often assumed we are at the minima, where there is an equality between the two. Concretely, having C classes, we have:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(z_i(\theta), y_i), \quad z_i(\theta) = f_\theta(x_i) \in \mathbb{R}^C.$$

$$J_i := \frac{\partial z_i(\theta)}{\partial \theta} \in \mathbb{R}^{C \times d}, \quad g_{z,i} := \nabla_z \ell(z_i(\theta), y_i) \in \mathbb{R}^C, \quad H_{z,i} := \nabla_z^2 \ell(z_i(\theta), y_i) \in \mathbb{R}^{C \times C},$$

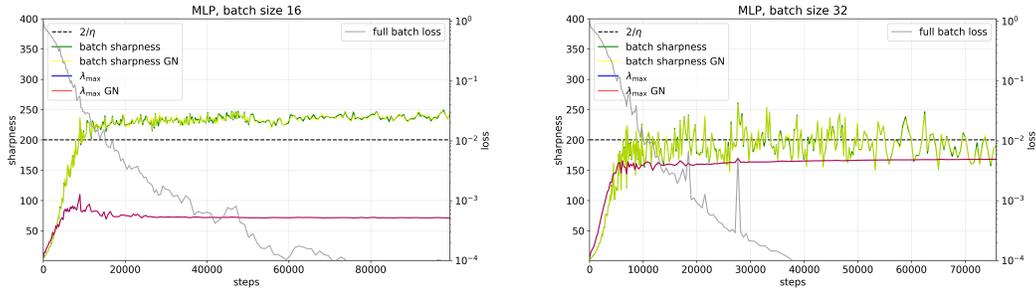
and for $j = 1, \dots, C$, let $\nabla_\theta^2 f_j(x_i) \in \mathbb{R}^{d \times d}$ denote the Hessian (w.r.t. θ) of the j -th output component. With this notation, we have:

$$\nabla_\theta L(\theta) = \frac{1}{N} \sum_{i=1}^N J_i^\top g_{z,i}.$$

and

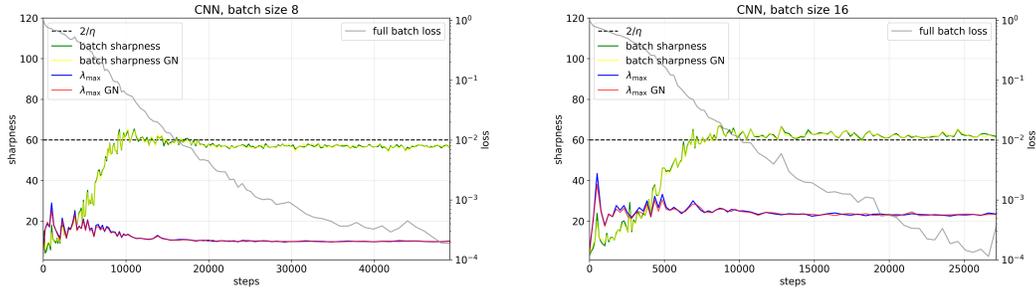
$$\begin{aligned} \nabla_\theta^2 L(\theta) &= \frac{1}{N} \sum_{i=1}^N \left(J_i^\top H_{z,i} J_i + \sum_{j=1}^C [g_{z,i}]_j \nabla_\theta^2 f_j(x_i) \right) \\ &= \underbrace{\frac{1}{N} \sum_{i=1}^N J_i^\top H_{z,i} J_i}_{\text{Gauss-Newton approx}} + \underbrace{\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C [g_{z,i}]_j \nabla_\theta^2 f_j(x_i)}_{\text{remainder / model-curvature term}}. \end{aligned}$$

3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197



3198 **Figure 32: Gauss–Newton approximation (MLP).** Comparison of *Batch Sharpness* and λ_{\max} computed with the true loss Hessian and its Gauss–Newton approximation, showing the validity of approximation. Both of the lines overlap almost perfectly

3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211



3212 **Figure 33: Gauss–Newton approximation (CNN).** Comparison of *Batch Sharpness* and λ_{\max} computed with the true loss Hessian and its Gauss–Newton approximation, showing the validity of approximation. Both of the lines overlap almost perfectly

3213
3214

3215 For MSE, this simplifies to:

3216
3217
3218
3219
3220
3221

$$\nabla_{\theta}^2 \mathcal{L}(\theta) = \underbrace{\frac{1}{N} \sum_{i=1}^N J_i^T J_i}_{\text{Gauss-Newton for MSE}} + \underbrace{\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C r_{i,j} \nabla_{\theta}^2 f_j(x_i)}_{\text{remainder / model-curvature term}}$$

3222

3222 in particular, if we are at minima, the residuals are 0, so the second terms completely disappears.

3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239

3223 Yet, the dynamics enter the EOSS regime away from minima (which is showcased by the continued decrease of the loss), where the Gauss-Newton approximation might not hold. In this Appendix we illustrate empirically that the Gauss-Newton approximation is close to the actual loss Hessian—at least from the perspective of EOSS and SGD stability. In particular, we compare *Batch Sharpness*, λ_{\max} and λ_{\max}^b when computed on the actual loss Hessian and on its Gauss-Newton approximations. Figures 32 to 34 illustrate that the computed quantities coincide throughout the *whole* training. Notice that due to the fact that the Gauss-Newton approximation and the NTK have the same spectrum, the λ_{\max} and λ_{\max}^b results also apply to the highest eigenvalues of the full-batch and mini-batch NTKs ($\frac{1}{B} J_B J_B^T$). Note that this agrees with the findings in literature, see e.g. Pappas (2019), but it was not clear whether the Gauss-Newton approximation holds before convergence. Our experiments demonstrate that it does the throughout the whole training, in particular, during EOSS.

3240
 3241
 3242
 3243
 3244
 3245
 3246
 3247
 3248
 3249
 3250
 3251
 3252
 3253
 3254
 3255
 3256
 3257
 3258
 3259
 3260
 3261
 3262
 3263
 3264
 3265
 3266
 3267
 3268
 3269
 3270
 3271
 3272
 3273
 3274
 3275
 3276
 3277
 3278
 3279
 3280
 3281
 3282
 3283
 3284
 3285
 3286
 3287
 3288
 3289
 3290
 3291
 3292
 3293

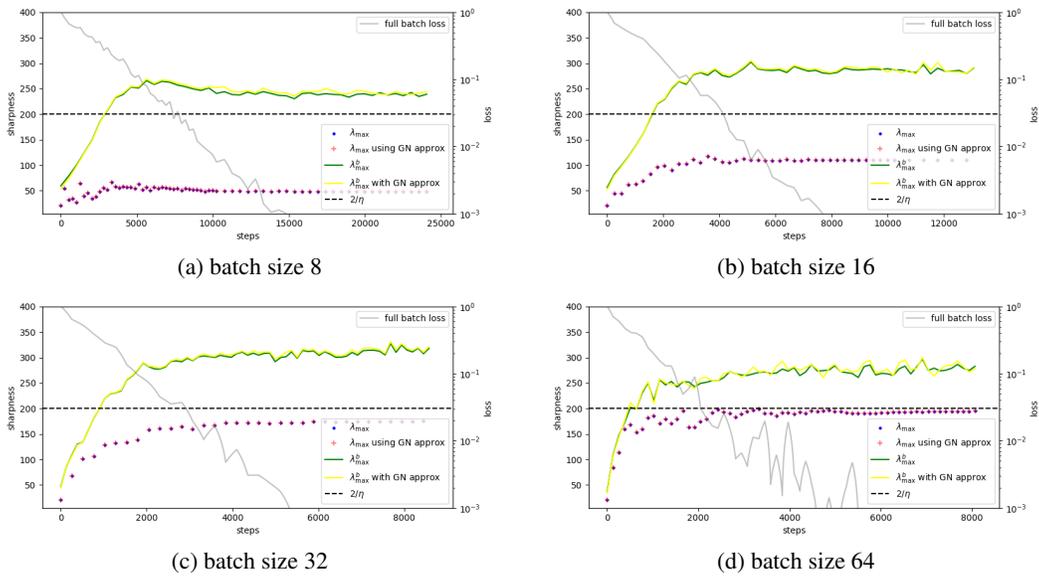


Figure 34: **Gauss-Newton approximation.** Showcasing the correctness of using the Gauss-Newton approximation to the loss Hessian by comparing the λ_{\max} and λ_{\max}^b computed with the true loss Hessian and its GN-approximation across different batch sizes.

3294 Q ILLUSTRATION OF EOSS IN VARIETY OF SETTINGS: *Batch Sharpness*

3295

3296

3297

3298

3299

3300

3301

3302

3303

3304

3305

3306

3307

3308

In this appendix, we provide further empirical evidence that EOSS arises robustly across a variety of architectures, step sizes, and batch sizes. For each experiment, we plot three quantities: λ_{\max} , *Batch Sharpness*, and *step sharpness* as a point cloud, which constitutes *Batch Sharpness* without expectation, and measured only on the *current* batch. Notice that time-averaging *step sharpness* is approximately the same as taking expectation over batches (albeit with slowly changing parameters), so it is approximately equal to *Batch Sharpness*, which takes this expectation at a point. Consistent with our main observations, we find that *Batch Sharpness* invariably stabilizes around $2/\eta$. Also refer to Section S for additional experiments illustrating EOSS

3309

3310

3311

3312

3313

3314

MLP (2-Layer) Baseline. Figure 36 illustrates EOSS for our baseline network, an MLP with two hidden layers of dimension 512, trained on an 8192-sample subset of CIFAR-10 with step size $\eta = 0.01$. As the training proceeds, *Batch Sharpness* stabilizes around $2/\eta$, whereas λ_{\max} plateaus strictly below *Batch Sharpness*.

3315

3316

3317

3318

3319

3320

3321

3322

3323

3324

3325

3326

3327

3328

3329

3330

3331

3332

3333

3334

3335

3336

3337

3338

3339

3340

3341

3342

3343

3344

3345

3346

3347

5-Layer CNN. We further confirm the EOSS regime in a five-layer CNN. As depicted in Figures 37, *Batch Sharpness* continues to plateau near the instability threshold for two distinct step sizes, while λ_{\max} once again settles at a lower level. Notably, as we vary the batch size, the gap between *Batch Sharpness* and λ_{\max} increases for smaller batches, mirroring the patterns described in Section I.

ResNet-14. Finally, we demonstrate that the EOSS regime also emerges for a deeper, residual architectures. In our case we are using RESNET-14 without BatchNor. Figure 41 highlights the same qualitative behavior, with *Batch Sharpness* stabilizing at $2/\eta$.

Overall, these experiments provide further confirmation that EOSS is a robust phenomenon across different architectures, step sizes, and batch sizes.

CNN with Full CIFAR-10. We also demonstrate in Figure 35 the emergence of EOSS when training on the full CIFAR-10 dataset. Consistent with the rest of the experiments, *Batch Sharpness* consistently stabilizes at $2/\eta$. Notably, in these experiments we also include a plot of the accuracy on the training set, to illustrate that EOSS happens away from the manifold of minima, and thus cannot be attributed solely to the structure around the manifold of minima.

3348
 3349
 3350
 3351
 3352
 3353
 3354
 3355
 3356
 3357
 3358
 3359
 3360
 3361
 3362
 3363
 3364
 3365
 3366
 3367
 3368
 3369
 3370
 3371
 3372
 3373
 3374
 3375
 3376
 3377
 3378
 3379
 3380
 3381
 3382
 3383
 3384
 3385
 3386
 3387
 3388
 3389
 3390
 3391
 3392
 3393
 3394
 3395
 3396
 3397
 3398
 3399
 3400
 3401

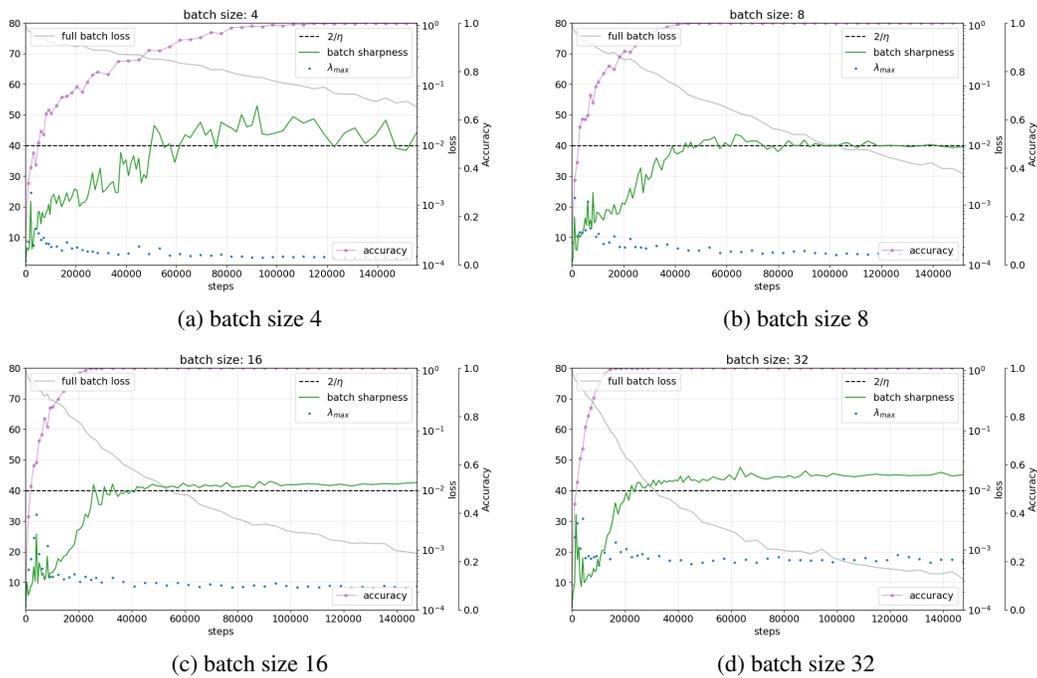


Figure 35: **CNN on Full CIFAR-10:** Same architecture as in Figure 37, but trained at a larger $\eta = 0.05$ on the full CIFAR-10 dataset, illustrating the emergence of EOSS away from manifold of minima.

3402
 3403
 3404
 3405
 3406
 3407
 3408
 3409
 3410
 3411
 3412
 3413
 3414
 3415
 3416
 3417
 3418
 3419
 3420
 3421
 3422
 3423
 3424
 3425
 3426
 3427
 3428
 3429
 3430
 3431
 3432
 3433
 3434
 3435
 3436
 3437
 3438
 3439
 3440
 3441
 3442
 3443
 3444
 3445
 3446
 3447
 3448
 3449
 3450
 3451
 3452
 3453
 3454
 3455

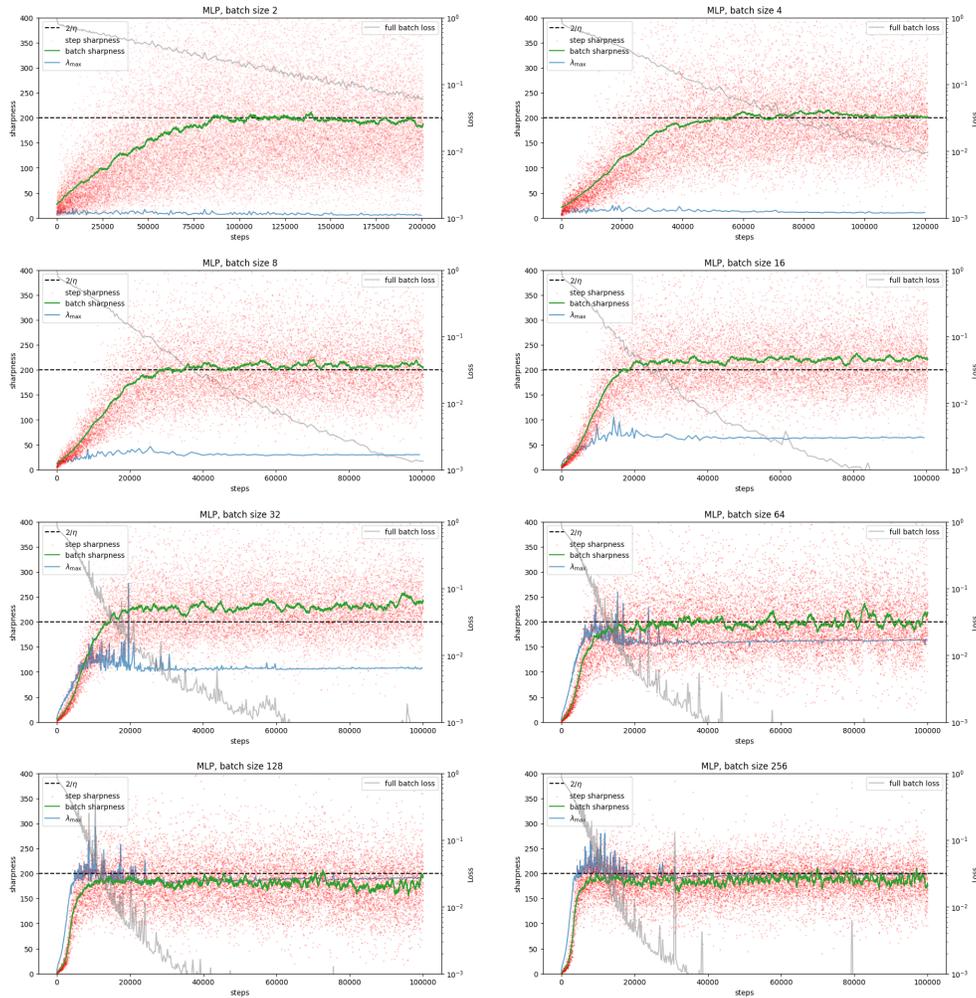


Figure 36: **MLP**: 2 hidden layers, hidden dimension 512; **step size 0.01**, 8k subset of CIFAR-10. Comparison between: step sharpness, aka batch sharpness without expectation over batches and measured on the current batch (red dots, time-smoothing would be \approx *Batch Sharpness*), the empirical *Batch Sharpness* (green line), the λ_{\max} (blue line).

3456
 3457
 3458
 3459
 3460
 3461
 3462
 3463
 3464
 3465
 3466
 3467
 3468
 3469
 3470
 3471
 3472
 3473
 3474
 3475
 3476
 3477
 3478
 3479
 3480
 3481
 3482
 3483
 3484
 3485
 3486
 3487
 3488
 3489
 3490
 3491
 3492
 3493
 3494
 3495
 3496
 3497
 3498
 3499
 3500
 3501
 3502
 3503
 3504
 3505
 3506
 3507
 3508
 3509

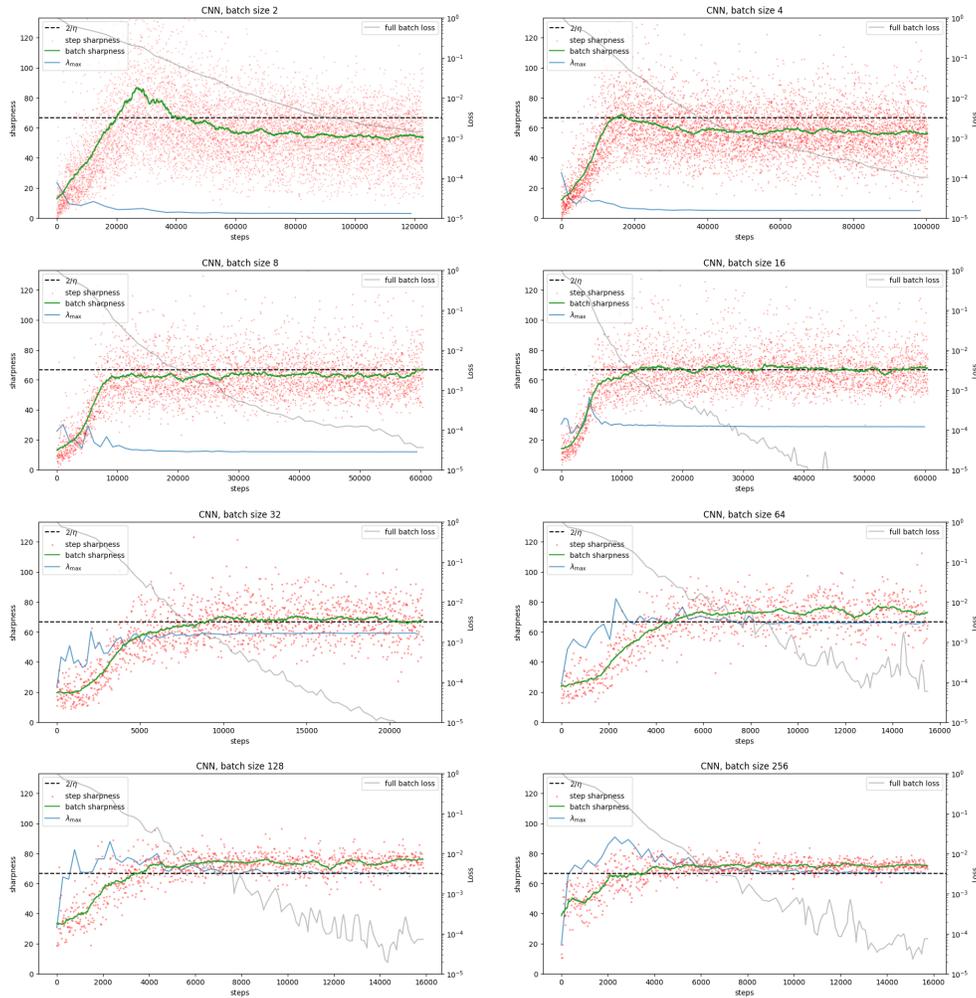


Figure 37: CNN: 5 layers (3 convolutional, 2 fully-connected), **step size 0.03**, 8k subset of CIFAR-10. Comparison between: step sharpness, aka batch sharpness without expectation over batches and measured on the current batch (red dots, time-smoothing would be \approx *Batch Sharpness*), the empirical *Batch Sharpness* (green line), the λ_{max} (blue line).

3510
 3511
 3512
 3513
 3514
 3515
 3516
 3517
 3518
 3519
 3520
 3521
 3522
 3523
 3524
 3525
 3526
 3527
 3528
 3529
 3530
 3531
 3532
 3533
 3534
 3535
 3536
 3537
 3538
 3539
 3540
 3541
 3542
 3543
 3544
 3545
 3546
 3547
 3548
 3549
 3550
 3551
 3552
 3553
 3554
 3555
 3556
 3557
 3558
 3559
 3560
 3561
 3562
 3563

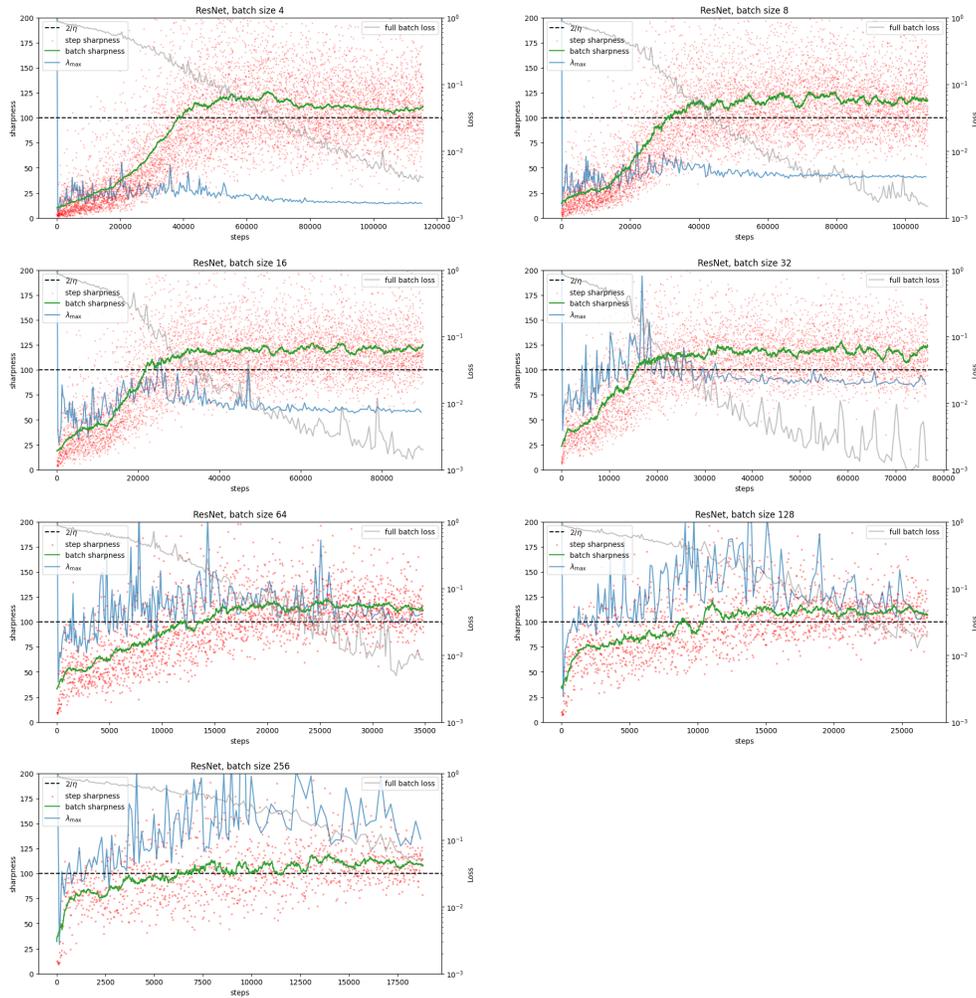


Figure 38: **ResNet-10**, step size 0.005, 8k subset of CIFAR-10. Comparison between: step sharpness, aka batch sharpness without expectation over batches and measured on the current batch (red dots, time-smoothing would be \approx *Batch Sharpness*), the empirical *Batch Sharpness* (green line), the λ_{\max} (blue line).

R ILLUSTRATION OF EOSS FOR THE SVHN DATASET

This appendix complements Appendix Q by verifying that the EOSS phenomena are not specific to CIFAR-10 but persist under a change of dataset. We repeat the experiments of Appendix Q—sweeping architectures (MLP, CNN, ResNet) and batch sizes—on an 8k subset of the SVHN dataset, and track step sharpness, Batch Sharpness, and the full-batch λ_{\max} along the training trajectory. Across all settings we again observe progressive sharpening followed by stabilization of Batch Sharpness at $2/\eta$, catapult-like spikes, and suppression of λ_{\max} below $2/\eta$, mirroring the behavior seen on CIFAR-10 and supporting the claim that EOSS is a robust feature of mini-batch SGD on standard vision benchmarks.

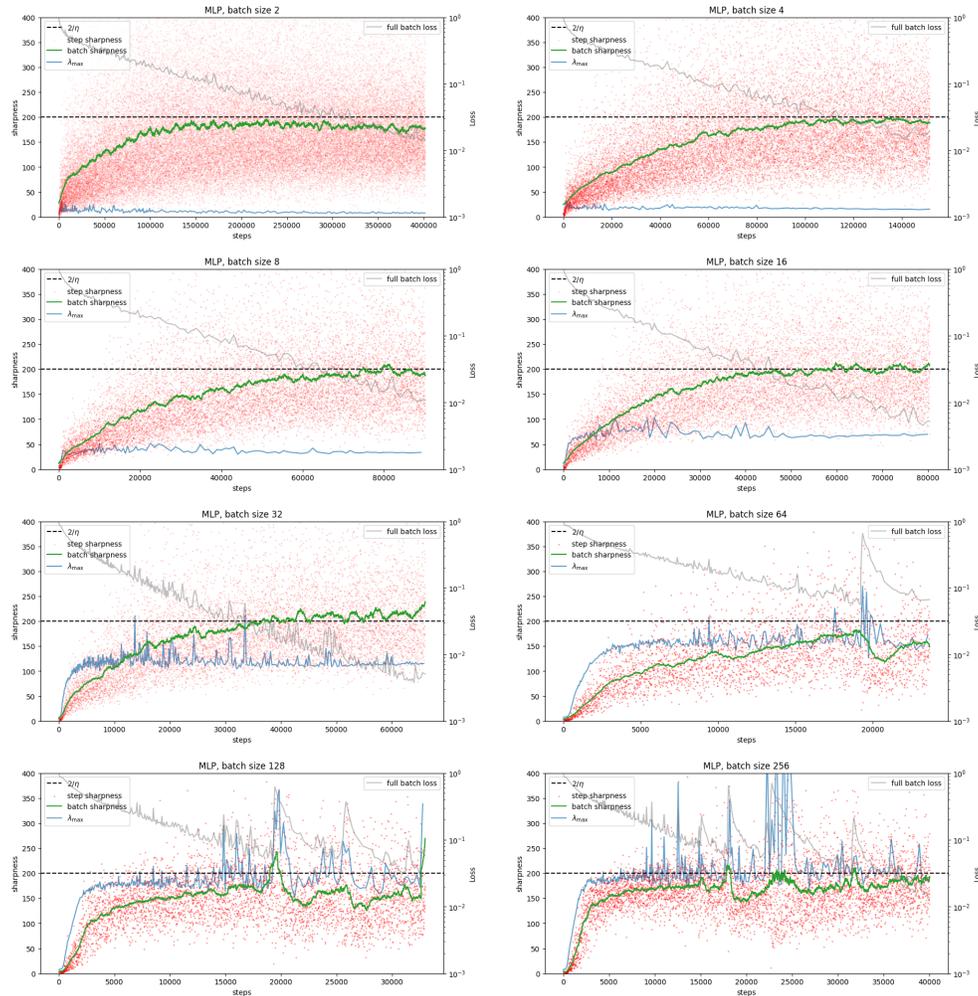


Figure 39: **MLP**: 2 hidden layers, hidden dimension 512; **step size 0.01**, 8k subset of **SVHN**. Comparison between: step sharpness, aka batch sharpness without expectation over batches and measured on the current batch (red dots, time-smoothing would be \approx Batch Sharpness), the empirical Batch Sharpness (green line), the λ_{\max} (blue line).

3618
 3619
 3620
 3621
 3622
 3623
 3624
 3625
 3626
 3627
 3628
 3629
 3630
 3631
 3632
 3633
 3634
 3635
 3636
 3637
 3638
 3639
 3640
 3641
 3642
 3643
 3644
 3645
 3646
 3647
 3648
 3649
 3650
 3651
 3652
 3653
 3654
 3655
 3656
 3657
 3658
 3659
 3660
 3661
 3662
 3663
 3664
 3665
 3666
 3667
 3668
 3669
 3670
 3671

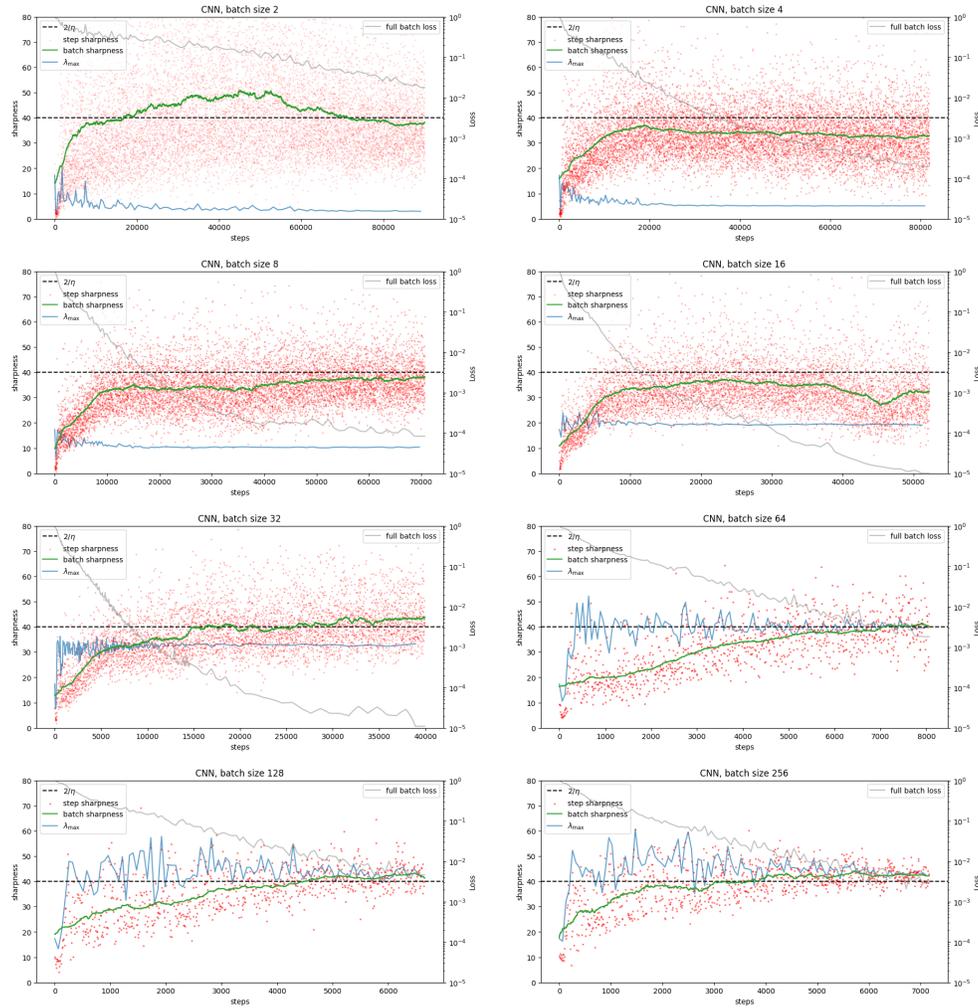


Figure 40: CNN: 5 layers (3 convolutional, 2 fully-connected), **step size 0.05**, 8k subset of SVHN. Comparison between: step sharpness, aka batch sharpness without expectation over batches and measured on the current batch (red dots, time-smoothing would be \approx Batch Sharpness), the empirical Batch Sharpness (green line), the λ_{\max} (blue line).

3672
 3673
 3674
 3675
 3676
 3677
 3678
 3679
 3680
 3681
 3682
 3683
 3684
 3685
 3686
 3687
 3688
 3689
 3690
 3691
 3692
 3693
 3694
 3695
 3696
 3697
 3698
 3699
 3700
 3701
 3702
 3703
 3704
 3705
 3706
 3707
 3708
 3709
 3710
 3711
 3712
 3713
 3714
 3715
 3716
 3717
 3718
 3719
 3720
 3721
 3722
 3723
 3724
 3725

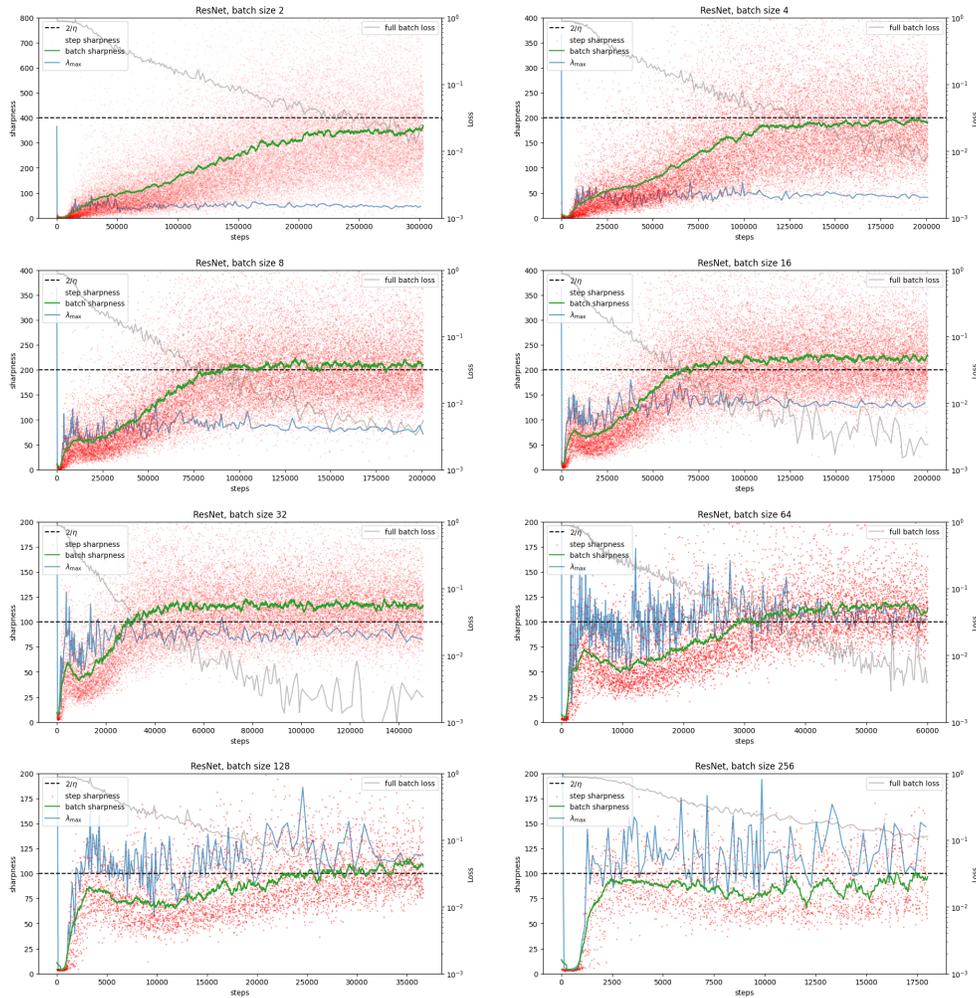


Figure 41: **ResNet-10**, step size 0.005, 8k subset of **SVHN**. Comparison between: step sharpness, aka batch sharpness without expectation over batches and measured on the current batch (red dots, time-smoothing would be \approx *Batch Sharpness*), the empirical *Batch Sharpness* (green line), the λ_{\max} (blue line).

3726 S ILLUSTRATION OF EOSS IN VARIETY OF SETTINGS: λ_{\max}^b
3727

3728 In this appendix, we provide additional empirical evidence for both emergence of EOSS and to Ap-
3729 pendix M.2 add!!!, varying across models, step sizes, and batch sizes. Consistent with our primary
3730 findings, we observe that λ_{\max}^b consistently stabilizes within the interval $(2/\eta, 2 \times 2/\eta]$, in partic-
3731 ular always higher than *Batch Sharpness* and λ_{\max} . We are conducting experiments on MLP, CNN
3732 and ResNet-20 in Figures 42, 43, 44 respectively.

3733 Note that the fact that λ_{\max}^b consistently stabilizes above $2/\eta$ implies that the supremum of Lipchitz
3734 constants of the gradient of individual mini-batch losses also stabilizes above $2/\eta$, thus clearly indi-
3735 cating that the usual assumptions in theory works on SGD about step size break. Same applies for
3736 supremum of Lip constants of gradients of per-sample losses. The former fact trivially follows from
3737 the inequality between sup and mean, and the second one from the same plus Lemma 6.

3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779

3780
 3781
 3782
 3783
 3784
 3785
 3786
 3787
 3788
 3789
 3790
 3791
 3792
 3793
 3794
 3795
 3796
 3797
 3798
 3799
 3800
 3801
 3802
 3803
 3804
 3805
 3806
 3807
 3808
 3809
 3810
 3811
 3812
 3813
 3814
 3815
 3816
 3817
 3818
 3819
 3820
 3821
 3822
 3823
 3824
 3825
 3826
 3827
 3828
 3829
 3830
 3831
 3832
 3833

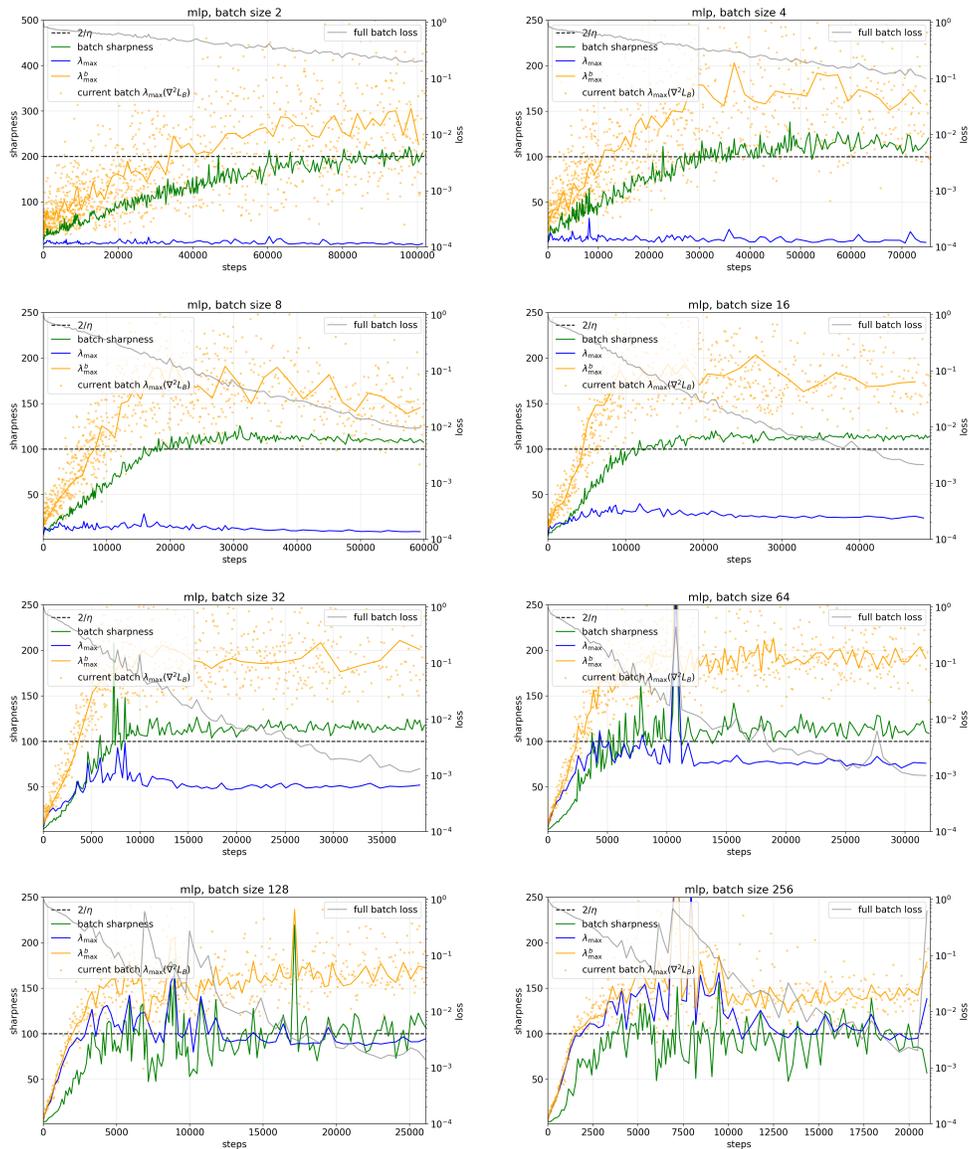


Figure 42: **Tracking λ_{\max}^b , MLP.** MLP with 2 hidden layers of width 512, step size 0.02, trained on an 8k subset of CIFAR-10. Comparison between the highest eigenvalue of the Hessian of the current mini-batch loss (orange dots, time-smoothed $\approx \lambda_{\max}^b$), *Batch Sharpness* (green line), λ_{\max} (blue line), and λ_{\max}^b (orange line). Note that *Batch Sharpness* stabilizes as $2/\eta$, while λ_{\max}^b is above it, and λ_{\max} is below for small batch sizes. Note that for batch size 2 we were using a lower step size of 0.01, as otherwise the network wasn't converging.

3834
 3835
 3836
 3837
 3838
 3839
 3840
 3841
 3842
 3843
 3844
 3845
 3846
 3847
 3848
 3849
 3850
 3851
 3852
 3853
 3854
 3855
 3856
 3857
 3858
 3859
 3860
 3861
 3862
 3863
 3864
 3865
 3866
 3867
 3868
 3869
 3870
 3871
 3872
 3873
 3874
 3875
 3876
 3877
 3878
 3879
 3880
 3881
 3882
 3883
 3884
 3885
 3886
 3887

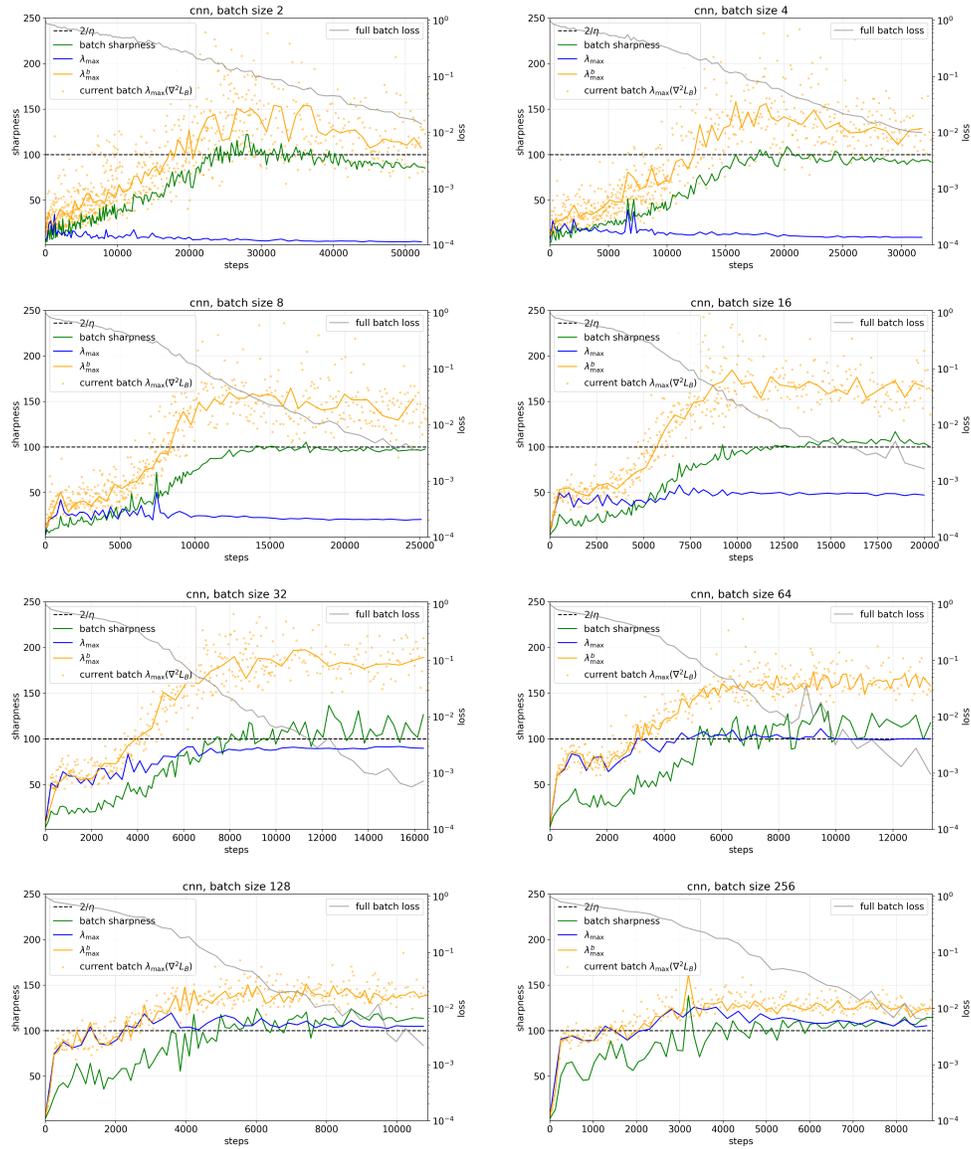


Figure 43: **Tracking λ_{\max}^b , CNN.** CNN with 5 layers (3 convolutional, 2 fully connected), step size 0.03, trained on an 8k subset of CIFAR-10. Comparison between the highest eigenvalue of the Hessian of the current mini-batch loss (orange dots, time-smoothed $\approx \lambda_{\max}^b$), *Batch Sharpness* (green line), λ_{\max} (blue line), and λ_{\max}^b (orange line). Note that *Batch Sharpness* stabilizes as $2/\eta$, while λ_{\max}^b is above it, and λ_{\max} is below for small batch sizes.

3888
 3889
 3890
 3891
 3892
 3893
 3894
 3895
 3896
 3897
 3898
 3899
 3900
 3901
 3902
 3903
 3904
 3905
 3906
 3907
 3908
 3909
 3910
 3911
 3912
 3913
 3914
 3915
 3916
 3917
 3918
 3919
 3920
 3921
 3922
 3923
 3924
 3925
 3926
 3927
 3928
 3929
 3930
 3931
 3932
 3933
 3934
 3935
 3936
 3937
 3938
 3939
 3940
 3941

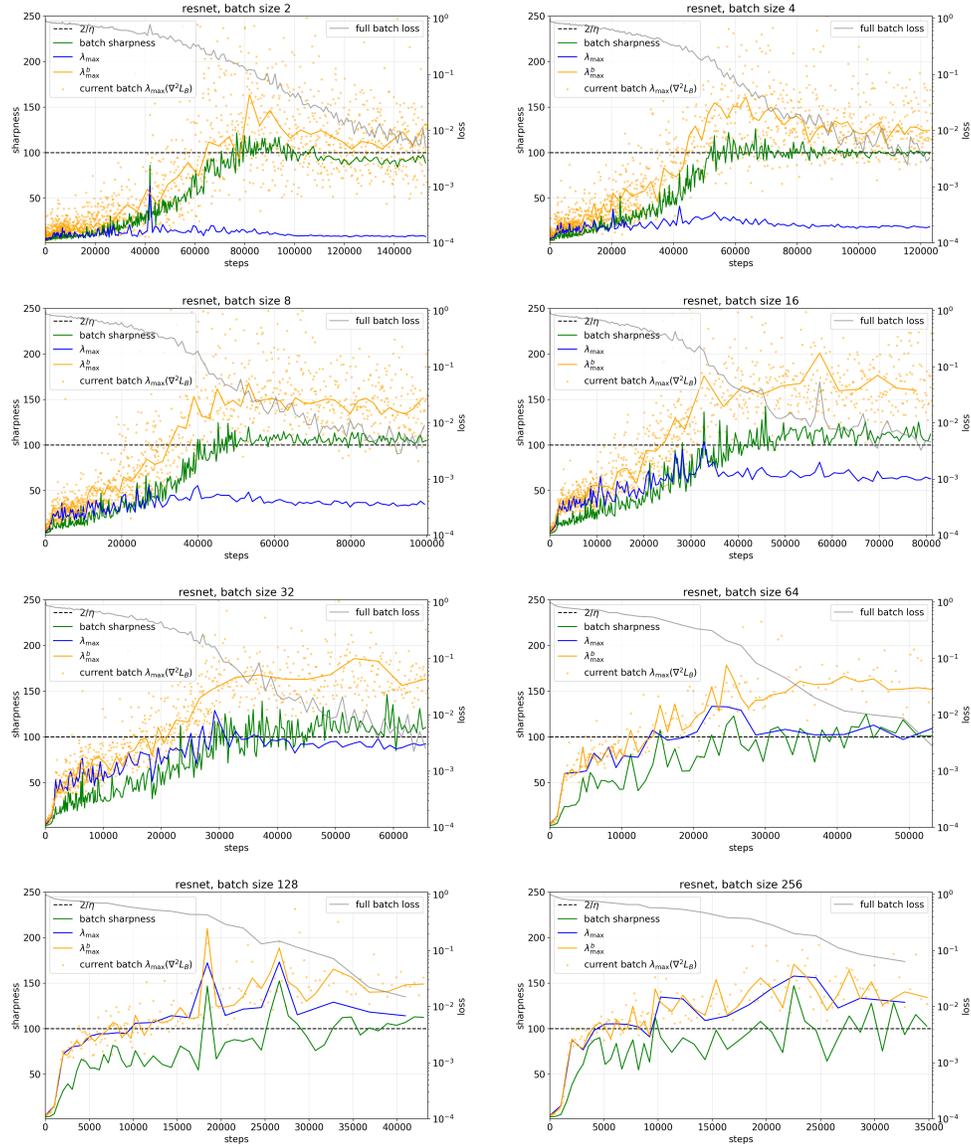


Figure 44: **Tracking λ_{\max}^b , ResNet-20.** ResNet-20 on CIFAR-10 with step size 0.02, trained on an 8k subset of the dataset. Comparison between the highest eigenvalue of the Hessian of the current mini-batch loss (orange dots, time-smoothed $\approx \lambda_{\max}^b$), *Batch Sharpness* (green line), λ_{\max} (blue line), and λ_{\max}^b (orange line). Note that *Batch Sharpness* stabilizes as $2/\eta$, while λ_{\max}^b is above it, and λ_{\max} is below for small batch sizes.