

# DOES LLM PRE-TRAINING TYPICALLY OCCUR AT THE EDGE OF STABILITY?

Anonymous authors

Paper under double-blind review

## ABSTRACT

Quadratic approximations are a common lens for neural network optimization, but recent evidence challenges their predictive validity. In full-batch gradient descent with LR  $\eta$ , Cohen et al. (2021) observed the Edge of Stability (EoS), where the largest Hessian eigenvalue concentrates near  $2/\eta$ , in tension with classical stability conditions. In this work, we revisit the fidelity of quadratic approximation as a model of neural network training dynamics, with particular focus on its failure modes in LLM training. We first identify and decouple a distinct failure mechanism of the quadratic approximation regardless of the LR choice, which arises from persistent negative curvature during training, which we term the *Edge of Convexity* (EoC). Based on the decoupling from EoC, we then extend the definition of EoS to large-scale stochastic training with adaptive optimizers. Across different LLM pretraining with various model sizes up to 1.7B, we find: (1) EoC is always observed across LLM pretraining. (2) EoS is also prevalent but not universal; it disappears when the LR becomes sufficiently small (e.g., after decay) or when the batch size falls below a critical threshold that is linearly related to the critical batch size. Together, these findings characterize when and how quadratic approximations fail and serve as foundations for future work on understanding the training dynamics of modern neural networks.

## 1 INTRODUCTION

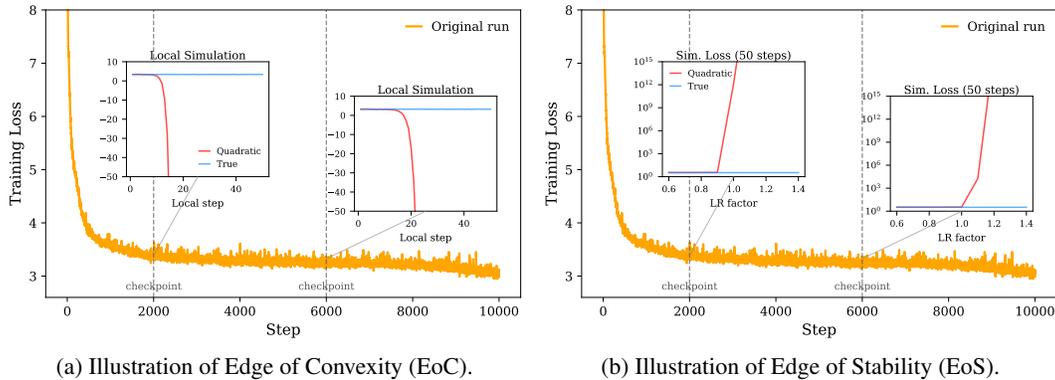


Figure 1: **Edge of Convexity (EoC) and Edge of Stability (EoS).** We train a 130M GPT-like model on FinWeb-10B and run local simulations at selected checkpoints. (a) The quadratic approximation diverges while the true dynamics remain stable. (b) With convexity compensation, a small LR increase induces a sharp loss transition, revealing the quadratic instability boundary.

In optimization theory, approximating the loss function with its *quadratic approximation*, namely the second-order Taylor expansion, is a common approach to study the local dynamics of the optimization process.

For example, based on local quadratic approximation, the classical descent lemma shows that gradient descent is able to decrease the loss when the learning rate (LR)  $\eta$  is smaller than  $2/\lambda_{\max}$ , where  $\lambda_{\max}$  is the maximal Hessian eigenvalue.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

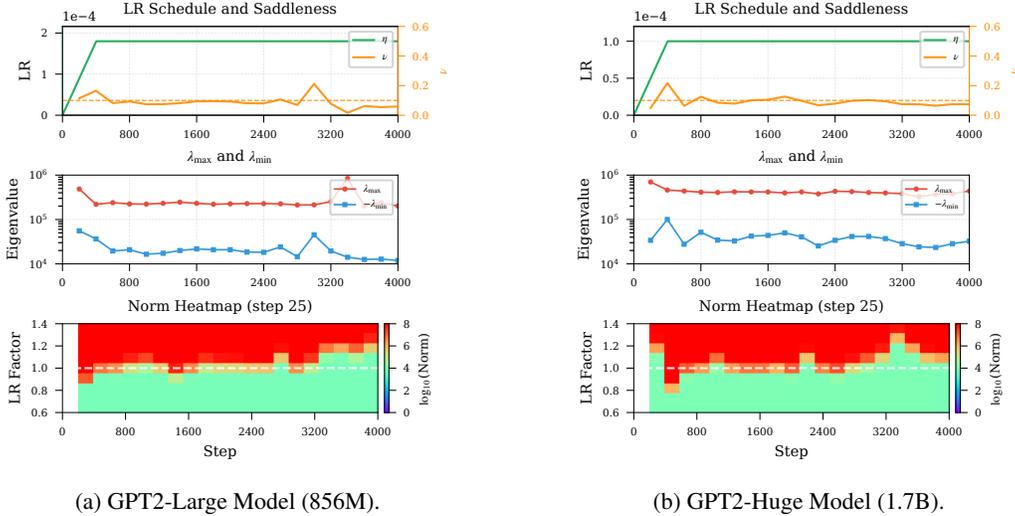


Figure 2: **EoS and EoC tests of GPT-Large and GPT-Huge models.** For each training run, we report: (i) the LR schedule; (ii) the top Hessian eigenvalue  $\lambda_{\max}$ ; (iii) the minimal Hessian eigenvalue  $\lambda_{\min}$ ; (iv) the saddleless  $\nu$ ; and (v) results of the convexity-compensated local simulation, shown as the final parameter norm across a grid of nine LR factors  $\{0.6, 0.7, \dots, 1.4\}$ .

However, in deep learning, the quadratic approximation is not always a good model for the training dynamics. Notably, Cohen et al. (2021) demonstrates that full-batch gradient descent (GD) typically occurs at *Edge of Stability* (EoS), where for a given LR  $\eta$ , the maximal Hessian eigenvalue  $\lambda_{\max}$  first progressively increases, and then oscillates near  $2/\eta$  while loss decreases non-monotonically. Follow-up work further shows an analogous EoS behavior arising for full-batch adaptive optimizers (e.g., Adam/RMSProp) when considering preconditioned curvature (Cohen et al., 2022).

This behavior fundamentally challenges existing optimization analyses for deep learning: The descent lemma arising from quadratic approximation is usually essential for these analyses to hold, but it would suggest zero decrement in loss under EoS. This has motivated a line of work to bridge this gap between theory and practice (Arora et al., 2022; Damian et al., 2022; Cohen et al., 2024).

Despite that this fundamental challenge of applying quadratic approximation to deep learning, key gaps remain in understanding *when* and *how* it fails to predict the true training dynamics in modern LLM training.

### Our Contributions.

- **Edge of Convexity (EoC) and negative curvature.** We identify and formalize Edge of Convexity (EoC) as an LR-independent failure mode of the quadratic approximation caused by persistent negative curvature, and we introduce a compensation scheme that provably decouples EoC from EoS while only slightly perturbing the EoS threshold (Figure 1a and Section 2.2).
- **Generalized EoS definition and detection.** After decoupling the impact of negative curvature, we extend the EoS from full-batch GD to stochastic training with adaptive optimizers via second-moment stability analysis of a local stochastic quadratic approximation (Figure 1b and Section 2).
- **Scalable diagnostics and large-scale characterization.** We develop an eigendecomposition-free local simulation framework to detect both EoC and EoS at scale with theoretical guarantees (Appendix F), and use it to study LLM pre-training. We observe that (1) EoC is always observed across neural network training; (2) EoS is prevalent but not universal, and it disappears when the LR becomes sufficiently small (e.g., after decay) or when the batch size falls below a critical threshold (Appendix E).

## 2 MAIN RESULTS

In this section, we derive the definition of Edge of Convexity (EoC) and Edge of Stability (EoS) and introduce practical tests for identifying them. We mainly use SGD as an example to demonstrate our EoC and EoS test methods. Our framework, however, naturally extends to widely used adaptive optimizers such as Adam and its variants with appropriate modifications. This extension is validated

by our LLM experiments in Appendix E and further discussed in Appendix H.2. Detailed derivations of the main results are provided in Appendix H, with complete proofs included thereafter.

## 2.1 STOCHASTIC STABILITY ANALYSIS

Compared to the full-batch setting, mini-batch sampling introduces additional variance that affects stability beyond what  $\lambda_{\max}$  alone can capture. To analyze this, for any fixed checkpoint  $\mathbf{w}_{t_0}$ , we consider a local quadratic approximation with stochastic Hessians:

$$Q_{\mathcal{B}_t}(\mathbf{w}) := \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + (\mathbf{w} - \mathbf{w}_{t_0})^\top \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_{t_0})^\top \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})(\mathbf{w} - \mathbf{w}_{t_0}), \quad (1)$$

where  $\mathcal{B}_t$  is a mini-batch sampled at time step  $t \geq t_0$ . Applying SGD with LR  $\eta$  to this approximation as  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}} Q_{\mathcal{B}_t}(\mathbf{w})$  yields a linear dynamical system whose stability depends on both the population Hessian and its covariance across batches. To be specific, instability can arise in two distinct ways. (1) If  $\mathbf{H}_{t_0}$  has negative curvature, the objective is unbounded below, and divergence can occur for any  $\eta$ ; (2) Even in the nonnegative-curvature subspace, mini-batch noise can make the projected second moment blow up once  $\eta$  is too large. We formalize both mechanisms next.

**Corollary 2.1** (Negative Curvature Divergence). *Let  $\mathbf{H}_{t_0} := \nabla^2 \mathcal{L}(\mathbf{w}_{t_0})$ . If  $\lambda_{\min}(\mathbf{H}_{t_0}) < 0$ , then for any  $\eta > 0$ , the quadratic model is unbounded below, and the SGD iterates on  $Q_{\mathcal{B}_t}$  diverge to  $-\infty$ .*

This corollary isolates instability caused purely by negative curvature, independent of  $\eta$ . In the remainder, we factor out these directions and focus on the locally convex subspace, where instability arises from mini-batch noise and admits a critical LR.

**Definition 2.2** (Critical LR for SGD). *Let  $\mathbf{H}_{t_0} := \nabla^2 \mathcal{L}(\mathbf{w}_{t_0})$ , we denote  $\Pi_{t_0}^+$  as the orthogonal projection onto its eigenspace of the nonnegative eigenvalues. We define*

$$\eta_*(\mathbf{w}_{t_0}) := \sup\{\eta > 0 : \sup_t \mathbb{E} [\|\Pi_{t_0}^+(\mathbf{w}_t)\|^2] < \infty\}.$$

Definition 2.2 considers the SGD trajectory projected onto the nonnegative-curvature subspace of  $\mathbf{H}_{t_0}$ , isolating the effect of EoC. In this restricted setting, the dynamics mirror full-batch EoS behavior and admit a critical LR. The following corollary formalizes this intuition.

**Corollary 2.3** (Positive Curvature Divergence). *Let  $\mathbf{H}_{t_0} := \nabla^2 \mathcal{L}(\mathbf{w}_{t_0})$  and let  $\eta^*(\mathbf{w}_{t_0})$  be defined in Definition 2.2, then the SGD trajectory is unbounded above in projected second moment  $\mathbb{E}\|\Pi_{t_0}^+(\mathbf{w}_t)\|^2$  diverges to  $+\infty$  when  $\eta > \eta^*(\mathbf{w}_{t_0})$ .*

For neural network (NN) training, once we replace the actual loss with  $Q_{\mathcal{B}_t}(\mathbf{w})$ , it is typically observed that the linear system formed by quadratic approximation diverges quickly, while the true dynamics of the NN remain stable (Figure 1a, Figure 1b).

So a natural question is: which regime leads to the failure of quadratic approximation in NN training, negative or negative curvature divergence? We identify: NN training sits on an *edge* where either divergence mechanism is *just at the threshold* of breaking the local quadratic approximation.

For the first negative divergence regime, we first define a quantity termed saddleness.

**Definition 2.4** (Saddleness). *Given any parameter point  $\mathbf{w}_{t_0}$  whose population Hessian  $\mathbf{H}_{t_0}$  has minimal, maximal eigenvalue denoted by  $\lambda_{\min}$  and  $\lambda_{\max}$ . We define the saddleness  $\nu$  at this point as*

$$\nu(\mathbf{w}_{t_0}) = \frac{\min\{0, -\lambda_{\min}\}}{\lambda_{\max}}.$$

Intuitively,  $\nu(\mathbf{w}_{t_0})$  measures how ‘‘saddle-like’’ the local curvature is, by comparing the most negative curvature to the strongest positive curvature. If  $\mathbf{H}_{t_0} \succeq 0$  (locally convex), then  $\nu(\mathbf{w}_{t_0}) = 0$ . A larger  $\nu$  means the negative curvature is comparable to  $\lambda_{\max}$  (strong saddle structure), while a small  $\nu$  means the negative curvature exists but is weak relative to the dominant positive directions.

One important observation is: For LLM pre-training, we observe that  $\nu \leq 0.1$  after the warm-up phase, which leads to the first edge area that LLM pre-training is in

**Edge of Convexity (EoC).** Negative eigenvalues persist during training, but  $\nu \approx 0$ , and true dynamics remains stable, while the quadratic approximation is unbounded below, and the iterates diverge to  $-\infty$  regardless of the choice of LR.

Even after removing the effect of EoC (i.e., restricting to the positive-curvature subspace in Definition 2.2), the quadratic approximation can still diverge. This divergence is caused by positive curvature when the learning rate exceeds the critical threshold  $\eta_*(\mathbf{w}_{t_0})$  (see Corollary 2.1). More subtly, neural network training often operates near this boundary: a tiny positive increase in the learning rate can already trigger divergence.

**Edge of Stability (EoS) at  $\mathbf{w}_{t_0}$ .** The current LR  $\eta_{t_0} \approx \eta_*(\mathbf{w}_{t_0})$ , a small positive perturbation of  $\eta$  leads the local quadratic approximation  $Q_{\mathcal{B}_t}(\mathbf{w})$  to  $+\infty$ .

One can notice that our definition recovers the full-batch EoS condition  $\lambda_{\max} \approx 2/\eta$  in Cohen et al. (2021). In the full-batch case,  $\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) \equiv \mathbf{H}_{t_0}$  and the update on the locally convex subspace reduces to  $\mathbf{x}_{t+1} = (\mathbf{I} - \eta \mathbf{H}_{t_0}) \mathbf{x}_t$ . Hence stability is equivalent to  $\rho(\mathbf{I} - \eta \mathbf{H}_{t_0}) < 1$ , i.e.,  $0 < \eta < 2/\lambda_{\max}(\mathbf{H}_{t_0})$ , so  $\eta_*(\mathbf{w}_{t_0}) = 2/\lambda_{\max}(\mathbf{H}_{t_0})$ .

In the following subsections, we propose practical methods to detect each regime that simulates local trajectories while avoiding explicit eigendecomposition

## 2.2 EOC TEST VIA MINIMAL EIGENVALUE ESTIMATION

At a checkpoint  $\mathbf{w}_{t_0}$ , EoC corresponds to persistent but small negative curvature, i.e.,  $\lambda_{\min}(\mathbf{H}_{t_0}) < 0$  and  $\nu \approx 0$ , under which the quadratic approximation is unbounded below and fails regardless of the LR. Our goal is to estimate  $\lambda_{\min}(\mathbf{H}_{t_0})$  using only stochastic mini-batch access, without costly eigendecomposition.

**Stochastic Power Iteration.** We propose a stochastic power iteration method with bisection to estimate the smallest eigenvalue of a large matrix from noisy observations, without costly eigendecomposition. The key idea is to probe the stability of a shifted stochastic iteration: for a shift coefficient  $\lambda$ , we run a few steps of a power iteration driven by the sampled matrices and track the growth of the iterate norm. When  $\lambda$  is smaller than the critical value, the dynamics admit an expansive mode and  $\|\mathbf{x}_t\|$  rapidly blows up; once  $\lambda$  is large enough (with a suitably small step size), the iteration becomes stable and the iterates remain bounded. This sharp “explode or stable” behavior enables a bisection search over  $\lambda$ . The full procedure is given in Algorithm 1, and its finite error bound is provided in Theorem F.1; empirical EoC results are shown in Figures 2 and 7.

**EoC Test Framework.** We use the above method to estimate the minimal eigenvalue of  $\mathbf{H}_{t_0}$  for any chosen checkpoint  $\mathbf{w}_{t_0}$  and use Lanczos (Lanczos, 1950) to estimate the maximal eigenvalue. Combining these two gives the estimated  $\hat{\nu}$ . When  $\hat{\nu} < \nu_{\text{tol}}$ , we identify that the current point is at EoC. The detailed description for this framework is provided in Algorithm 2. Actually, we observe the consistent negative eigenvalues occur during training across various training setups. See Appendix E for more empirical discussion.

## 2.3 EOS TEST VIA CONVEXITY-COMPENSATED LOCAL SIMULATION

**Compensated local loss.** We propose an efficient method to test EoS at scale when the saddleness is small. In detail, we define the compensated local loss as:

$$G_t(\mathbf{w}) := \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + (\mathbf{w} - \mathbf{w}_{t_0})^\top \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_{t_0})^\top \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})(\mathbf{w} - \mathbf{w}_{t_0}) + \frac{\lambda_P}{2} \|\mathbf{w} - \mathbf{w}_{t_0}\|^2, \quad (2)$$

where  $\lambda_P \geq 0$  is called the *compensation coefficient*. We first set  $\lambda_P \geq |\lambda_{\min}|$ . Likewise, we extend the critical LR in Definition 2.2 with the compensation coefficient  $\lambda_P$  as  $\eta_*^{\lambda_P}(\mathbf{w}_{t_0})$  (See Equation (14) for rigorous definition). Since the saddleness is observed to be small, we can prove that  $\eta_*^{\lambda_P}(\mathbf{w}_{t_0})$  is close to  $\eta_*(\mathbf{w}_{t_0})$  (Theorem F.2). Hence, the phase transition in the projected space spanned by positive eigenvectors can also be observed directly in this compensated local dynamics.

**EoS Test Framework.** We summarize the above discussion into a unified framework for the detection of EoS. We set the compensation coefficient using the estimation by Algorithm 1 proposed in Section 2.2 as  $\lambda_P = \hat{\lambda}$ . Then, we fix a single mini-batch sequence  $\{\mathcal{B}_t\}_{t=0}^{T-1}$ , for all LR factor  $\alpha \in \{\alpha_1 < \dots < \alpha_K\}$ , and simulate  $\mathbf{w}_{t+1}^{(\alpha)} = \mathbf{w}_t^{(\alpha)} - \alpha \eta \nabla G_t(\mathbf{w}_t)$ ,  $\mathbf{w}_0^{(\alpha)} = 0$ , as summarized in Algorithm 3. We claim a reference point  $\mathbf{w}_{t_0}$  is in EoS, if there exists a phase transition in the norms  $\|\mathbf{w}_T^{(\alpha)}\|$  (see the norm heat maps in Figure 2 and Figure 1b).

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

---

## REFERENCES

- Ahn, K., Zhang, J., and Sra, S. Understanding the unstable convergence of gradient descent. In *International conference on machine learning*, pp. 247–257. PMLR, 2022.
- Andreyev, A. and Beneventano, P. Edge of stochastic stability: Revisiting the edge of stability for sgd. *arXiv preprint arXiv:2412.20553*, 2024.
- Arora, S., Li, Z., and Panigrahi, A. Understanding gradient descent on the edge of stability in deep learning. In *International Conference on Machine Learning*, pp. 948–1024. PMLR, 2022.
- Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J., Sagun, L., and Zecchina, R. Entropy-sgd: Biasing gradient descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124018, 2019.
- Clement, C. B., Bierbaum, M., O’Keeffe, K. P., and Alemi, A. A. On the use of arxiv as a dataset. *arXiv preprint arXiv:1905.00075*, 2019. URL <https://arxiv.org/abs/1905.00075>.
- Cohen, J. M., Kaur, S., Li, Y., Kolter, J. Z., and Talwalkar, A. Gradient descent on neural networks typically occurs at the edge of stability. *arXiv preprint arXiv:2103.00065*, 2021.
- Cohen, J. M., Ghorbani, B., Krishnan, S., Agarwal, N., Medapati, S., Badura, M., Suo, D., Cardoze, D., Nado, Z., Dahl, G. E., et al. Adaptive gradient methods at the edge of stability. *arXiv preprint arXiv:2207.14484*, 2022.
- Cohen, J. M., Damian, A., Talwalkar, A., Kolter, J. Z., and Lee, J. D. Understanding optimization in deep learning with central flows. *arXiv preprint arXiv:2410.24206*, 2024.
- Damian, A., Nichani, E., and Lee, J. D. Self-stabilization: The implicit bias of gradient descent at the edge of stability. *arXiv preprint arXiv:2209.15594*, 2022.
- Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., and Bengio, Y. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *Advances in neural information processing systems*, 27, 2014.
- Dong, Z., Zhang, Y., Yao, J., and Sun, R. Towards quantifying the hessian structure of neural networks. *arXiv preprint arXiv:2505.02809*, 2025.
- Ghorbani, B., Krishnan, S., and Xiao, Y. An investigation into neural net optimization via hessian eigenvalue density. In *International Conference on Machine Learning*, pp. 2232–2241. PMLR, 2019.
- Gokaslan, A. and Cohen, V. Openwebtext corpus. <https://skylion007.github.io/OpenWebTextCorpus/>, 2019. URL <https://skylion007.github.io/OpenWebTextCorpus/>.
- Granziol, D., Zohren, S., and Roberts, S. Learning rates as a function of batch size: A random matrix theory approach to neural network training. *Journal of Machine Learning Research*, 23(173):1–65, 2022.
- Henry, A., Dachapally, P. R., Pawar, S. S., and Chen, Y. Query-key normalization for transformers. In Cohn, T., He, Y., and Liu, Y. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, Online, November 2020. Association for Computational Linguistics.
- Hu, S., Tu, Y., Han, X., He, C., Cui, G., Long, X., Zheng, Z., Fang, Y., Huang, Y., Zhao, W., Zhang, X., Thai, Z. L., Zhang, K., Wang, C., Yao, Y., Zhao, C., Zhou, J., Cai, J., Zhai, Z., Ding, N., Jia, C., Zeng, G., Li, D., Liu, Z., and Sun, M. Minicpm: Unveiling the potential of small language models with scalable training strategies, 2024.
- Jastrzębski, S., Kenton, Z., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. On the relation between the sharpest directions of dnn loss and the sgd step length. *arXiv preprint arXiv:1807.05031*, 2018.
- Jordan, K., Bernstein, J., Rappazzo, B., et al. modded-nanogpt: Speedrunning the nanogpt baseline, 2024. URL <https://github.com/KellerJordan/modded-nanogpt>.

---

270 Kingma, D. P. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.  
271

272 Lanczos, C. An iteration method for the solution of the eigenvalue problem of linear differential and  
273 integral operators. *Journal of research of the National Bureau of Standards*, 45(4):255–282, 1950.

274 LeCun, Y., Bottou, L., Orr, G. B., and Müller, K.-R. Efficient backprop. In *Neural networks: Tricks*  
275 *of the trade*, pp. 9–50. Springer, 2002.  
276

277 Li, J., Fang, A., Smyrnis, G., Ivgi, M., Jordan, M., Gadre, S. Y., et al. Datacomp-lm: In search of  
278 the next generation of training sets for language models. *arXiv preprint arXiv:2406.11794*, 2024.  
279 URL <https://arxiv.org/abs/2406.11794>.

280 Li, X., Wen, H., and Lyu, K. Adam reduces a unique form of sharpness: Theoretical insights near the  
281 minimizer manifold. *arXiv preprint arXiv:2511.02773*, 2025.  
282

283 Li, Z., Wang, Z., and Li, J. Analyzing sharpness along gd trajectory: Progressive sharpening and  
284 edge of stability. *arXiv preprint arXiv:2207.12678*, 2022.

285 Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*,  
286 2017.  
287

288 Ma, C. and Ying, L. On linear stability of sgd and input-smoothness of neural networks. *Advances in*  
289 *Neural Information Processing Systems*, 34:16805–16817, 2021.

290 Ma, C., Kunin, D., Wu, L., and Ying, L. Beyond the quadratic approximation: The multiscale  
291 structure of neural network loss landscapes. *arXiv preprint arXiv:2204.11326*, 2022.  
292

293 Malladi, S., Lyu, K., Panigrahi, A., and Arora, S. On the SDEs and Scaling Rules for Adaptive  
294 Gradient Algorithms, November 2024. URL <http://arxiv.org/abs/2205.10287>.

295 Marek, M., Lotfi, S., Somasundaram, A., Wilson, A. G., and Goldblum, M. Small batch size training  
296 for language models: When vanilla sgd works, and why gradient accumulation is wasteful. *arXiv*  
297 *preprint arXiv:2507.07101*, 2025.  
298

299 Merrill, W., Arora, S., Groeneveld, D., and Hajishirzi, H. Critical batch size revisited: A simple  
300 empirical approach to large-batch language model training. *arXiv preprint arXiv:2505.23971*,  
301 2025.  
302

303 Mulayoff, R. and Michaeli, T. Exact mean square linear stability analysis for sgd. In *The Thirty*  
304 *Seventh Annual Conference on Learning Theory*, pp. 3915–3969. PMLR, 2024.

305 Penedo, G., Kydlíček, H., Ben Allal, L., Lozhkov, A., Mitchell, M., Raffel, C., Von Werra, L., and  
306 Wolf, T. The fineweb datasets: Decanting the web for the finest text data at scale. *arXiv preprint*  
307 *arXiv:2406.17557*, 2024. URL <https://arxiv.org/abs/2406.17557>.

308 Qiu, Z., Wang, Z., Zheng, B., Huang, Z., Wen, K., Yang, S., Men, R., Yu, L., Huang, F., Huang, S.,  
309 Liu, D., Zhou, J., and Lin, J. Gated attention for large language models: Non-linearity, sparsity,  
310 and attention-sink-free. In *The Thirty-ninth Annual Conference on Neural Information Processing*  
311 *Systems*, 2025.  
312

313 Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu,  
314 P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of*  
315 *Machine Learning Research*, 21(140):1–67, 2020. URL [https://www.jmlr.org/papers/](https://www.jmlr.org/papers/v21/20-074.html)  
316 [v21/20-074.html](https://www.jmlr.org/papers/v21/20-074.html).

317 Sagun, L., Bottou, L., and LeCun, Y. Eigenvalues of the hessian in deep learning: Singularity and  
318 beyond. *arXiv preprint arXiv:1611.07476*, 2016.  
319

320 Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the hessian of  
321 over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.  
322

323 Tropp, J. A. Freedman’s inequality for matrix martingales. *Electronic Communications in Probability*,  
16:262–270, 2011.

---

324 Wikimedia Foundation. Wikipedia: Database download. <https://dumps.wikimedia.org/>,  
325 2024. URL <https://dumps.wikimedia.org/>. Accessed for English Wikipedia text  
326 dumps.

327  
328 Wu, J., Braverman, V., and Lee, J. D. Implicit bias of gradient descent for logistic regression at the  
329 edge of stability. *Advances in Neural Information Processing Systems*, 36:74229–74256, 2023.

330  
331 Wu, L., Ma, C., et al. How sgd selects the global minima in over-parameterized learning: A dynamical  
332 stability perspective. *Advances in Neural Information Processing Systems*, 31, 2018.

333  
334 Wu, L., Wang, M., and Su, W. The alignment property of sgd noise and how it helps select flat  
335 minima: A stability analysis. *Advances in Neural Information Processing Systems*, 35:4680–4693,  
2022.

336  
337 Xing, C., Arpit, D., Tsirigotis, C., and Bengio, Y. A walk with sgd. *arXiv preprint arXiv:1802.08770*,  
2018.

338  
339 Yao, Z., Gholami, A., Keutzer, K., and Mahoney, M. W. Pyhessian: Neural networks through the lens  
340 of the hessian. In *2020 IEEE international conference on big data (Big data)*, pp. 581–590. IEEE,  
341 2020.

342  
343 Zhang, Y., Chen, C., Ding, T., Li, Z., Sun, R., and Luo, Z. Why transformers need adam: A hessian  
344 perspective. *Advances in neural information processing systems*, 37:131786–131823, 2024.

345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

---

## 378 A CONCLUSION

379  
380 This paper extends the definition of EoS to the LLM training under a unified framework, together  
381 with another clearly clarified phenomenon EoC, accounting for the two failure modes of the quadratic  
382 approximation.

383 Several directions remain open for future study. (i) Our experiments mainly focus to SGD and  
384 Adamw. It is interesting to see how would EoS would behave when training with other optimizers  
385 like Muon. (ii) Fully understanding EoS and EoC also calls for solid theoretical analyses of the  
386 realistic neural network training dynamics.  
387

## 388 B ROADMAP OF APPENDIX

389 Appendix C provides a detailed discussion of related work. In Appendix D, we introduce the  
390 necessary notations and background of full-batch Edge of Stability used in the following empirical  
391 and theoretical derivations. In Appendix E, we present our detailed experimental results, together  
392 with all the figures, and Appendix G further clarifies and compensates the experimental details.  
393 Appendix F shows the theoretical guarantees for the EoC and EoS tests in practice. Appendix H and  
394 the following sections together give a self-contained derivation of our main results of EoC and EoS,  
395 with proofs and algorithm descriptions provided.  
396  
397

## 398 C RELATED WORK

399 **Edge of Stability.** Early studies connected sharpness and stability to optimizer behavior in neural  
400 network training: Wu et al. (2018) analyzed SGD minima selection through dynamical stability using  
401 sharpness and non-uniformity, and Xing et al. (2018) observed SGD “bouncing” between valley walls  
402 during training. Jastrzębski et al. (2018) then showed that SGD visits increasingly sharp regions  
403 up to a maximum set by LR and batch size, and that the step is often too large along the sharpest  
404 directions. Cohen et al. (2021) later demonstrated that full-batch gradient descent typically operates  
405 at an Edge of Stability (EoS), where the largest Hessian eigenvalue oscillates near  $2/\eta$  while loss  
406 decreases non-monotonically. Cohen et al. (2022) extended this phenomenon to adaptive methods  
407 via saturation of a preconditioned stability threshold. And theoretical explanations for EoS were  
408 discussed in simplified models or realistic assumptions (Arora et al., 2022; Li et al., 2022; Damian  
409 et al., 2022).  
410

411 **Comparison with Andreyev & Beneventano (2024).** The recent work (Andreyev & Beneventano,  
412 2024) is among the most relevant to this work, which defines Edge of Stochastic Stability for SGD  
413 via the batch sharpness, where Batch Sharpness hovers around  $2/\eta$ . Compared to Andreyev &  
414 Beneventano (2024), our work (1) In the full-batch setting, our definition for EoS can be reduced to  
415 the original definition of EoS by Cohen et al. (2021; 2022), while theirs cannot; (2) We also identify  
416 another failure mode of the quadratic approximation called Edge of Convexity (EoC); (3) We provide  
417 scalable EoS and EoC test methods with theoretical guarantees; (4) We conduct experiments on  
418 large-scale LLM pretraining with AdamW, while most of their experiments are limited to ResNets on  
419 CIFAR-10 and SVHN with vanilla SGD.  
420

421 **Training Stability of Neural Networks** Series works studied the training stability of constant-  
422 learning-rate SGD in quadratic losses via a linear system analysis of the second moment of the  
423 parameters (Wu et al., 2018; Ma & Ying, 2021; Granzio et al., 2022; Wu et al., 2022; Mulayoff &  
424 Michaeli, 2024). They obtained tight conditions in the form of  $\|\mathbb{E}[(\mathbf{I} - \eta \mathbf{H}_t)^{\otimes 2}]\| \leq 1$ , which are  
425 valid instability criteria for a particular Lyapunov function. However, in the training neural networks  
426 with (S)GD, the unstable convergence phenomenon was observed (Wu et al., 2018; Xing et al., 2018;  
427 Jastrzębski et al., 2018), formally identified and termed *Edge of Stability* (EoS) (Cohen et al., 2021).  
428 Subsequently, several studies sought to explain this unstable convergence (Ma et al., 2022; Ahn et al.,  
429 2022; Wu et al., 2023; Arora et al., 2022; Damian et al., 2022; Cohen et al., 2024). Among them,  
430 Cohen et al. (2024) gave a general characterization of the oscillatory dynamics of deep learning in the  
431 EoS regime via a time-averaged differential equation termed *central flow*. However, the effectiveness  
of the quadratic approximation was never carefully studied in the large-scale neural network training,  
hence it is not clear when and how the quadratic approximation fails or not in the realistic regime.

**Hessian Spectrum of Neural Networks.** Most studies on the Hessian spectrum of neural networks reported that the Hessian spectra of neural networks consist of a “bulk” together with a few large “outliers” in magnitude (LeCun et al., 2002; Dauphin et al., 2014; Sagun et al., 2016; 2017; Ghorbani et al., 2019; Chaudhari et al., 2019; Yao et al., 2020). Among them, Sagun et al. (2016; 2017) mentioned that there are still negative eigenvalues during training, but their magnitude is much smaller than the large outliers. More recent works reported and explained a near-block diagonal structure in the Hessian spectrum of CNNs and Transformers (Zhang et al., 2024; Dong et al., 2025). However, one thing that has long been overlooked is the characterization of the negative eigenvalues in the population Hessian during large-scale training. In this work, we present detailed empirical findings of the negative eigenvalues across different setups in large-scale neural network training.

## D PRELIMINARY

### D.1 NOTATION

Throughout,  $\|\cdot\|$  denotes the spectral norm for matrices and the  $\ell_2$  norm for vectors. For a symmetric matrix  $\mathbf{A}$ , let  $\lambda_{\max}(\mathbf{A})$  and  $\lambda_{\min}(\mathbf{A})$  be its largest and smallest eigenvalues. We write  $\text{vec}(\cdot)$  for vectorization and  $\otimes$  for the Kronecker product;  $\text{Cov}(\mathbf{M}_t)$  denotes the covariance of  $\text{vec}(\mathbf{M}_t)$ . For adaptive methods, we use momentum  $\mathbf{m}_t$ , a preconditioner  $\mathbf{P}_t$  updated by a map  $\Psi$ , and weight decay coefficient  $\lambda_{\text{WD}}$ ; the relevant curvature is the spectrum of  $\mathbf{P}_t^{-1}\nabla^2\mathcal{L}(\mathbf{w}_t)$ , equivalently that of its symmetric similarity transform  $\mathbf{P}_t^{-1/2}\nabla^2\mathcal{L}(\mathbf{w}_t)\mathbf{P}_t^{-1/2}$ . We use  $\mathbb{S}^d$  to denote the set of  $d \times d$  symmetric matrices and  $\mathbb{S}_{++}^d$  to narrow the range to symmetric positive definite matrices. For a positive semidefinite matrix  $\mathbf{S}$ , define the seminorm  $\|\mathbf{u}\|_{\mathbf{S}} := \sqrt{\mathbf{u}^\top \mathbf{S} \mathbf{u}}$ .

### D.2 PROBLEM SETUP

We consider the training dynamics of a neural network  $f(\mathbf{x}; \mathbf{w})$ , where  $\mathbf{x} \in \mathbb{R}^d$  denotes the input and  $\mathbf{w} \in \mathbb{R}^p$  denotes the trainable parameter. Given a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and a per-sample loss function  $\ell$ , the full-batch loss and mini-batch loss are defined as

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &:= \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{w}), y_i), \\ \mathcal{L}_{\mathcal{B}}(\mathbf{w}) &:= \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \ell(f(\mathbf{x}_i; \mathbf{w}), y_i),\end{aligned}$$

where  $\mathcal{B} \subseteq [n]$  denotes the mini-batch of size  $|\mathcal{B}|$ . Let  $\nabla\mathcal{L}(\mathbf{w})$  and  $\nabla^2\mathcal{L}(\mathbf{w})$  denote the full-batch gradient and Hessian, respectively.

### D.3 FULL-BATCH EDGE OF STABILITY

Consider full-batch gradient descent with LR  $\eta$ :

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla\mathcal{L}(\mathbf{w}_t).$$

In a recent work, Cohen et al. (2021) observed that during neural network training, the largest eigenvalue of the Hessian, denoted by  $\lambda_{\max}$ , typically increases over time and eventually approaches the threshold  $2/\eta$ . After reaching this value,  $\lambda_{\max}$  oscillates around it while the loss decreases non-monotonically. This regime is referred to as the *Edge of Stability* (EoS).

To understand why this behavior is surprising, consider the second-order Taylor expansion of the loss around a reference point  $\mathbf{w}_{t_0}$ :

$$\begin{aligned}Q(\mathbf{w}) &:= \mathcal{L}(\mathbf{w}_{t_0}) + (\mathbf{w} - \mathbf{w}_{t_0})^\top \nabla\mathcal{L}(\mathbf{w}_{t_0}) \\ &\quad + \frac{1}{2}(\mathbf{w} - \mathbf{w}_{t_0})^\top \nabla^2\mathcal{L}(\mathbf{w}_{t_0})(\mathbf{w} - \mathbf{w}_{t_0}).\end{aligned}$$

Applying gradient descent to this quadratic objective yields the linear dynamical system

$$(\mathbf{w}_{t+1} - \mathbf{w}_\star) = (I - \eta \nabla^2\mathcal{L}(\mathbf{w}_{t_0}))(\mathbf{w}_t - \mathbf{w}_\star),$$

where  $\mathbf{w}_\star$  satisfies  $\nabla\mathcal{L}(\mathbf{w}_{t_0}) + \nabla^2\mathcal{L}(\mathbf{w}_{t_0})\mathbf{w}_\star = 0$ .

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

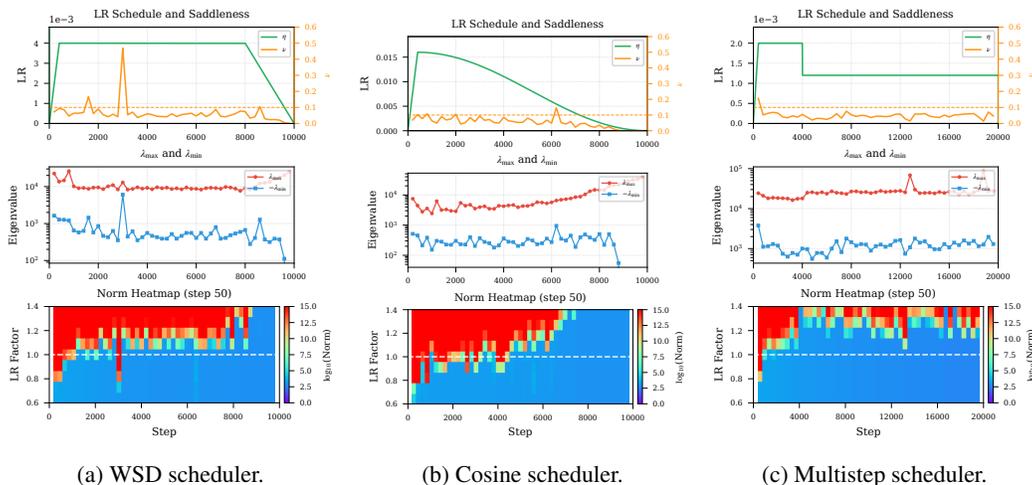
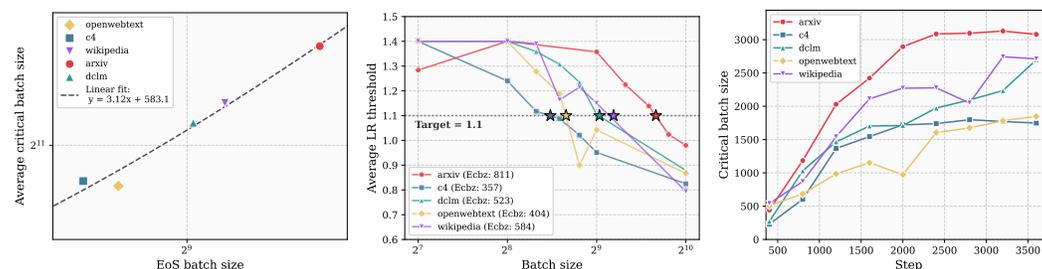


Figure 3: **EoS and EoC test under different LR schedulers.** We train a 43M GPT-like model on FineWeb-10B using AdamW with different LR schedulers. For WSD and Cosine scheduler, we tune the peak LR via a separate sweep. All runs use 400 linear warmup steps and decay the LR to zero unless otherwise specified. (A) WSD: 400 warmup / 7600 stable / 2000 decay steps, peak LR = 0.004. (B) Cosine: peak LR = 0.016. (C) Multistep: peak LR = 0.002, LR drop at step 4000 to 0.0012, no decay phase.

Classical analysis of this system shows that gradient descent is stable if and only if the Hessian is positive semidefinite and the LR is smaller than a critical threshold as  $\eta \leq 2/\lambda_{\max}$ .

From this perspective, the observation that  $\lambda_{\max}$  concentrates near  $2/\eta$  throughout training is striking. It places training dynamics precisely at the boundary of stability predicted by the quadratic model, rather than safely within it. The EoS phenomenon therefore, exposes a fundamental mismatch between the predictions of quadratic approximation and the actual behavior of neural network training.

## E EXPERIMENTS



(a) EoS batch size vs. critical batch (b) Average critical LR ratio vs. batch (c) Critical batch size vs. checkpoints.

Figure 4: **Relationship between the EoS batch-size and the critical batch size.** (a) EoS batch size versus critical batch size across settings. (b) Estimated EoS batch-size threshold as a function of  $\eta_*^{\lambda_p} / \eta$ . We use the target ratio  $\eta_*^{\lambda_p} / \eta = 1.1$ . (c) Critical batch-size estimates across checkpoints; we report the average over the last 4 checkpoints as the final estimate.

### E.1 EXPERIMENTAL SETUP.

We pretrain a family of GPT-style models with size ranging from 43M to 1.7B parameters (GPT-extra-small to GPT-Huge). Unless otherwise specified, we train on FineWeb (Penedo et al., 2024) for 5B tokens (following the Modded GPT setup (Jordan et al., 2024)) with AdamW (WD= 0.1,

---

540  $\beta_1 = 0.9$ ), global batch size 512, and the Warm–Stable–Decay (WSD) schedule (Hu et al., 2024) for  
541 10,000 steps (400 warmup / 7,600 stable / 2,000 linear decay to zero), excepted for GPT-extra-small  
542 to GPT-Huge models, for which we run a shorted horizon of 4000 steps with 400 warmup steps due to  
543 the compute and time limit. We additionally repeat key experiments on ArXiv (Clement et al., 2019),  
544 OpenWebText (Gokaslan & Cohen, 2019), C4 (Raffel et al., 2020), DCLM (Li et al., 2024), and  
545 Wikipedia (Wikimedia Foundation, 2024) to verify robustness across data sources. Full architectural  
546 specifications, HVP/simulation implementation details, and additional training/runtime settings are  
547 provided in Appendix G.

## 548

## 549 E.2 LARGE-SCALE EOS AND EOC VERIFICATION

## 550

551 **EoS and EoC are prevalent across model scales** We examine the presence of Edge of Stability  
552 (EoS) and Edge of Convexity (EoC) across model scales ranging from 43M to 1.7B parameters,  
553 as shown in Figure 7 (130M, 391M), Figure 3a (43M), and Figure 2 (856M, 1.7B). We perform a  
554 fine-grained LR search and identify an optimal LR based on training loss for 43M model. Training  
555 near this optimal regime consistently shows that training is approximately at the critical threshold,  
556 exhibiting a clear EoS signature, while persistent negative curvature and a near-zero saddleness  $\nu$   
557 (typically less than 0.1 in all our experiments) is always observed throughout training, suggesting a  
558 consistent EoC phenomenon in LLM pre-training. Additionally, for larger models with model size up  
559 to 1.7B, EoS is observed under standard LR choices. The qualitative behavior is consistent with that  
560 of smaller models, and EoC remains present throughout training. Overall, these results indicate that  
561 EoC and EoS are prevalent across model scales in practical training settings.

## 562 E.3 ABLATION STUDY

## 563

564 We perform a series of ablation studies to examine the robustness of EoC and EoS across training  
565 settings.

566 **EoS is Robust across training configurations.** Across variations in peak LR (Figure 5), batch  
567 size (Figure 6), scheduler (Figure 3), and model scale (Figures 2, 3a and 7), we consistently observe  
568 the presence of EoC. Specifically, the saddleness  $\nu$  remains below 0.1 for most of the training time,  
569 despite some unstable spikes in a few checkpoints. On the other hand, the minimal eigenvalue of the  
570 population preconditioned Hessian can attain a non-negligible magnitude.

571 **EoS emerges around and below the optimal LR.** We vary the peak LR and report the corre-  
572 sponding EoS/EoC tests in Figure 5. When the peak LR is at or below the empirically optimal value  
573 (Figures 5a and 5b), we observe a clear EoS signature. In contrast, when the peak LR is set above the  
574 optimum (Figure 5c), the quadratic approximation becomes highly unstable.

575 **Schedulers shape the temporal profile of EoS.** As shown in Figure 3, the LR schedule strongly  
576 affects when and how long the EoS signature is observed. Under WSD, EoS is consistently present  
577 throughout the stable phase, but it fades during the subsequent decay phase. Under cosine decay, EoS  
578 is concentrated in the earlier part of training and disappears noticeably earlier than in WSD. Under  
579 the multistep schedule, EoS becomes less prominent after the abrupt LR drop.

580 **EoS appears only above a batch size threshold.** As shown in Figure 6, EoS is tightly coupled to  
581 the global batch size. Specifically, with a smaller batch size ( $B = 256$ ; Figure 6a), the EoS signature  
582 is not observed. Once the batch size exceeds a threshold ( $B = 512$  and  $B = 1024$ ; Figures 6b  
583 and 6c), EoS emerges clearly and is observed reliably across checkpoints. We call this batch size  
584 threshold the EoS batch size. Figure 4c further shows that when the batch size is small (less than  
585  $2^8$ ), the EoS would not happen (The averaged LR threshold reaches the maximal factor 1.4 in grid  
586 searching, thus no grid-searched LR factor would lead to divergence), suggesting that EoS similarly  
587 disappears when the batch size is below some thresholds for other commonly used pre training  
588 datasets other than fineweb. We then term the above batch size by the EoS batch size.

589 **The EoS batch size is smaller but linearly related to the critical batch size.** We examine the  
590 relationship between the EoS batch size and the classical critical batch size ( $B_c$ ) in Figure 4. EoS is  
591 identified using the averaged critical LR ratio  $\eta_*^{\text{cp}}(\mathbf{w}_{t_0})/\eta$  over checkpoints, with a fixed threshold of  
592 1.1. Fitting this metric against batch size (Figure 4b) yields the EoS batch size. The critical batch  
593 size is estimated following Merrill et al. (2025) via local simulation and averaged over the last four  
checkpoints (Figure 4c). We find that the EoS batch size is consistently smaller than the critical

batch size, while the two vary proportionally across datasets. Implementation details are deferred to Appendix G.4.

## F THEORETICAL GUARANTEES BEHIND THE MEASUREMENTS OF EOC AND EOS

In this section, we provide theoretical justification for the definitions and tests of EoC and EoS introduced in Appendix 2. A self-contained derivation that unifies the empirical procedures from Appendix 2 with the theory developed here is presented in Appendix H.

### F.1 ERROR BOUND FOR STOCHASTIC POWER ITERATION

**Theorem F.1** (informal). *Let  $\mathbf{A} \in \mathbb{S}^d$  and  $\mathbf{A}_t = \mathbf{A} + \Delta_t$  be noisy symmetric observations with  $\mathbb{E}[\Delta_t] = \mathbf{0}$  and  $|\Delta_t| \leq \sigma$  almost surely. Given a nontrivial initialization overlap with the minimal eigenvector of  $\mathbf{A}$ , Algorithm 1 with  $\eta \in [\omega(1/T), o(1/\sqrt{T})]$ , horizon  $T$ , and  $N$  bisection rounds outputs an estimate  $\hat{\lambda}_{\min}$  such that, with high probability,*

$$|\hat{\lambda}_{\min} - \lambda_{\min}(\mathbf{A})| = \tilde{\mathcal{O}}\left(\frac{1}{\eta T} + \frac{\sigma\sqrt{\log d}}{\sqrt{T}}\right) + \mathcal{O}(2^{-N}).$$

This theorem provides a finite-sample error bound for our stochastic power iteration estimator of  $\lambda_{\min}$  in Appendix 2.2. It guarantees that Algorithm 1 returns a reliable estimate of  $|\lambda_{\min}|$ . A formal statement is deferred to Appendix H.1.2.

### F.2 EOS TEST ERROR FOR SGD

**Theorem F.2.** *Given  $\text{Cov}(\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})) \preceq \sigma^2 \mathbf{I}$  and the vanilla SGD, we have*

$$\frac{\eta_*^{\lambda_P}}{\eta_*^0} \in \left[ \frac{\lambda_{\max}(\lambda_{\max} + \lambda_P)}{(\lambda_{\max} + \lambda_P)^2 + \sigma^2}, \frac{\lambda_{\max}^2 + \sigma^2}{\lambda_{\max}(\lambda_{\max} + \lambda_P)} \right].$$

*Specifically, if  $\lambda_P = -\lambda_{\min}$ , then we have*

$$\frac{\eta_*^{\lambda_P}}{\eta_*^0} \in \left[ \frac{\kappa^2 + \kappa}{(\kappa + 1)^2 + \alpha}, \frac{\kappa^2 + \alpha}{\kappa^2 + \kappa} \right]$$

*where  $\alpha = \frac{\sigma^2}{|\lambda_{\min}|^2}$ ,  $\kappa := 1/\nu = \lambda_{\max}/|\lambda_{\min}|$ .*

This theorem quantifies how much the critical LR (hence the EoS location) changes after adding the convexity-compensation term. In particular, when we choose  $\lambda_P = -\lambda_{\min}$  to cancel negative curvature, the ratio  $\eta_*^{\lambda_P}/\eta_*^0$  is controlled by the condition number  $\kappa$  and the normalized noise level  $\alpha = \sigma^2/|\lambda_{\min}|^2$ . Therefore, our “remove EoC then test EoS” procedure does not arbitrarily shift the EoS threshold; its distortion is explicitly bounded by  $\nu$  and  $\alpha$ .

## G DETAILED EXPERIMENTS

### G.1 EXPERIMENTAL SETUPS

**Model architectures.** We use the following models in our experiments.

We further use GPT-2 vocabulary for all models except GPT-extra-small (which uses an 8k vocabulary for fast ablations), and apply QK-normalization (Henry et al., 2020) and attention gating (Qiu et al., 2025) for stability.

To accelerate the computation, we run most experiments on GPT-extra-small. Unless otherwise specified, for GPT-extra-small we use 100 iterations for stochastic power iteration and 50 iterations for convexity-compensated local simulation. For larger models (GPT-small/medium/large/huge), we use 50 iterations for stochastic power iteration and 25 iterations for convexity-compensated local simulation.

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

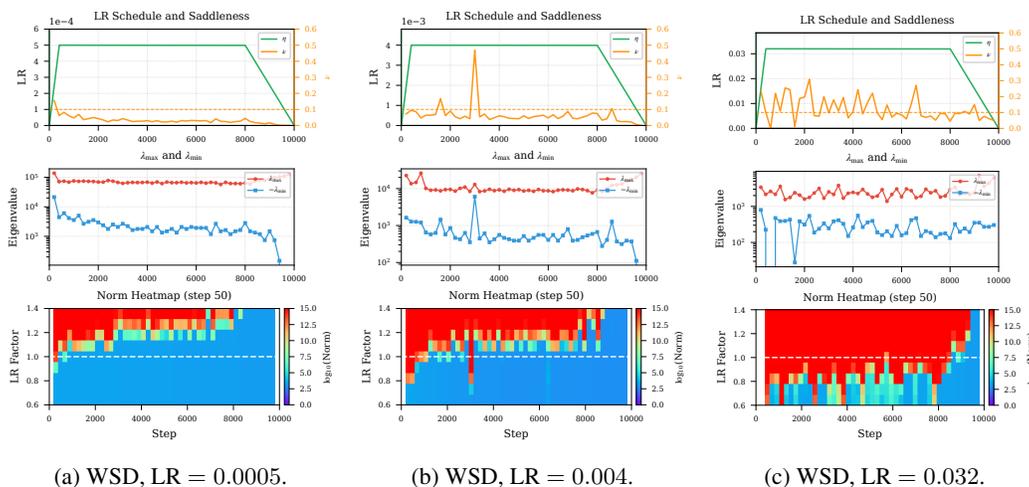


Figure 5: **EoS and EoC measurements under different LRs.** We train a 43M GPT-like model on FineWeb-10B using AdamW with the WSD scheduler and vary the LR. All runs use 400 linear warmup steps, 7600 stable steps, and 2000 linear decay steps. we identify 0.004 as the optimal peak LR via a seperate sweep.

Table 1: Model architectures used in our experiments.

Model	#Layers	Hidden Dim	#Heads	Head Dim	Vocab Size	Params
GPT-extra-small	4	768	12	64	8,192	43M
GPT-small	12	768	12	64	50,257	130M
GPT-medium	24	1,024	16	64	50,257	391M
GPT-large	36	1,280	20	64	50,257	856M
GPT-huge	48	1,600	25	64	50,257	1.7B

We log results at a fixed set of checkpoints. For GPT-extra-small, we evaluate every 200 checkpoints, yielding 50 checkpoints from step 0 to 10,000. For GPT-small and GPT-medium, we evaluate every 400 checkpoints, yielding 25 checkpoints from step 0 to 10,000. For GPT-large and GPT-huge, due to computational constraints, we evaluate every 200 steps from step 0 to 4,000.

**HVP computation.** We compute Hessian–vector products (HVPs) using PyTorch’s autograd by differentiating a gradient–vector inner product:

$$\mathbf{g}(\theta) = \nabla_{\theta} \mathcal{L}(\theta), \quad \text{HVP}(\mathbf{v}) = \nabla_{\theta} (\mathbf{g}(\theta)^{\top} \mathbf{v}) = \nabla_{\theta}^2 \mathcal{L}(\theta) \mathbf{v}. \quad (3)$$

Below we estimate its compute cost at the level of a single *linear layer*; this is representative because transformers are composed predominantly of linear operators (attention projections and MLP projections), and their runtime is largely dominated by GEMMs.

**Per-layer GEMM counting (linear layer).** Consider a linear layer

$$\mathbf{y} = \mathbf{x}\mathbf{W}, \quad (4)$$

where  $\mathbf{x} \in \mathbb{R}^{B \times d_{\text{in}}}$ ,  $\mathbf{W} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}}$ , and  $\mathbf{y} \in \mathbb{R}^{B \times d_{\text{out}}}$ . Let  $\mathbf{g} = \partial \mathcal{L} / \partial \mathbf{y}$  be the upstream gradient. The forward pass uses one GEMM:

$$\mathbf{y} = \mathbf{x}\mathbf{W} \Rightarrow 1 \text{ GEMM}. \quad (5)$$

The first-order backward pass computes gradients w.r.t.  $\mathbf{x}$  and  $\mathbf{W}$ :

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \mathbf{g}\mathbf{W}^{\top}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{x}^{\top} \mathbf{g}, \quad (6)$$

which costs two GEMMs. Hence one forward+backward through a linear layer costs  $\approx 3$  GEMMs (ignoring elementwise ops).

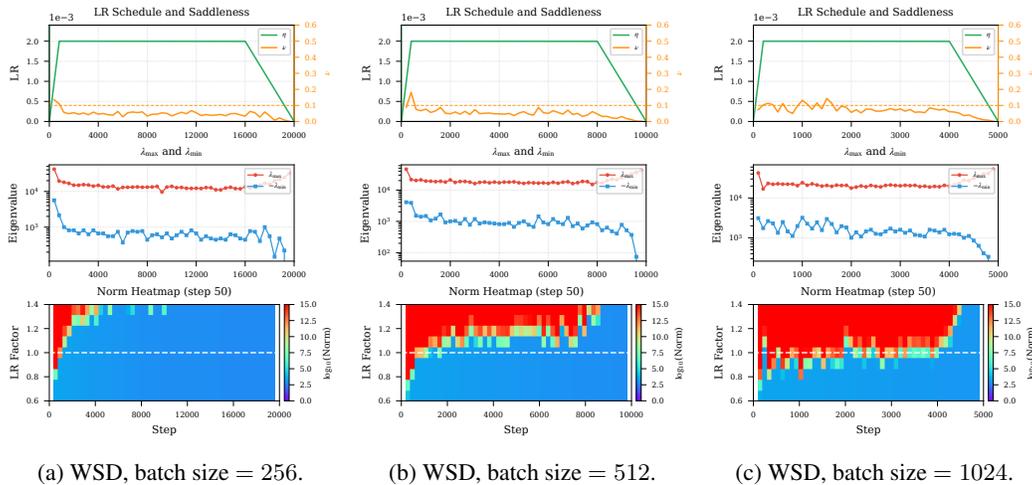


Figure 6: **EoS and EoC measurements under different batch sizes.** We train a 43M GPT-like model on FineWeb-10B using AdamW with the WSD scheduler and vary the global batch size, while fixing the LR to 0.002. (A)–(C) Batch sizes = 256, 512, and 1024, respectively.

The above autograd recipe computes an HVP by effectively backpropagating through the gradient computation (a second reverse-mode pass). For matmul-dominated layers, this incurs a constant-factor overhead comparable to an additional “forward+backward” through the same linear operators, plus extra GEMMs from differentiating the backward equations. Concretely, augmenting the forward/backward with the required second-order terms yields approximately

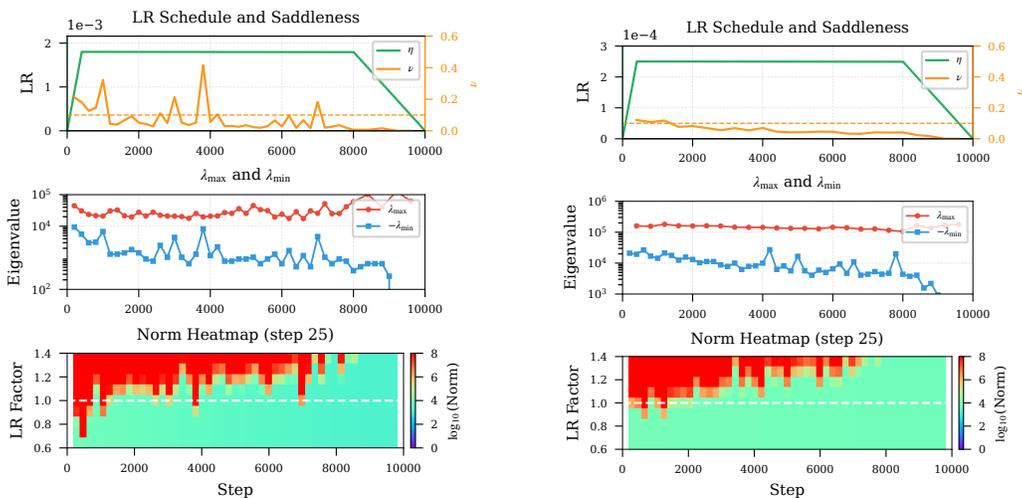
$$\underbrace{1}_{\text{forward}} + \underbrace{2}_{\text{backward}} + \underbrace{2}_{\text{differentiate backward (w.r.t. } W)} + \underbrace{4}_{\text{propagate directional adjoints}} \approx 9 \text{ GEMMs}, \quad (7)$$

i.e., about  $3 \times$  **the cost of one gradient evaluation** (9 vs. 3 GEMMs per linear layer), or equivalently  $9 \times$  **the cost of a forward pass** in the GEMM-dominated regime. Since transformer blocks are largely composed of such linear layers (QKV projections, attention output projection, and the two MLP projections), this constant-factor overhead provides a practical estimate of HVP cost for our models.

**Single simulation.** A single EoS simulation under a fixed training setting consists of the following procedure.

- **Step 1 (pretraining and checkpointing).** We pretrain the model under the specified setting (architecture, dataset, optimizer, scheduler, batch size, etc.). During training, we save checkpoints at a fixed frequency. Each checkpoint will be analyzed independently in the local-simulation stage below.
- **Step 2 (estimating the top eigenvalue).** For each checkpoint, we estimate the population top eigenvalue  $\lambda_{\max}$  of the Hessian. Concretely, we form an empirical *average Hessian* by averaging Hessian-vector products over a fixed set of mini-batches, and run the Lanczos algorithm on this averaged operator to obtain an estimate of  $\lambda_{\max}$ .
- **Step 3 (locating the convexity compensation).** For each checkpoint, we run Algorithm 1 to estimate the *optimal convexity compensation*  $\lambda_P$ . Intuitively,  $\lambda_P$  is the smallest additional regularization (in our local model) that removes the dominant negative-curvature effect at this checkpoint. We use this value as an operational measure of the strength of EoC.
- **Step 4 (LR-factor sweep and EoS heatmaps).** For each checkpoint, we run the local simulation on the modified approximate loss (with the estimated  $\lambda_P$ ) across a fixed grid of learning-rate factors  $\{0.6, 0.8, \dots, 1.3, 1.4\}$ . We record the final parameter norm produced by each run and visualize the results as EoS heatmaps. This sweep provides a coarse but stable picture of the local stability boundary.
- **Step 5 (estimating the critical LR).** For each checkpoint, we run Algorithm 5 to estimate the critical LR  $\eta_{\star}^{\lambda_P}$  under convexity compensation  $\lambda_P$ . We compare the ratio  $\eta_{\star}^{\lambda_P}/\eta$  with the theoretical quantity  $\kappa/(1+\kappa)$ . Here  $\eta$  is the original LR.

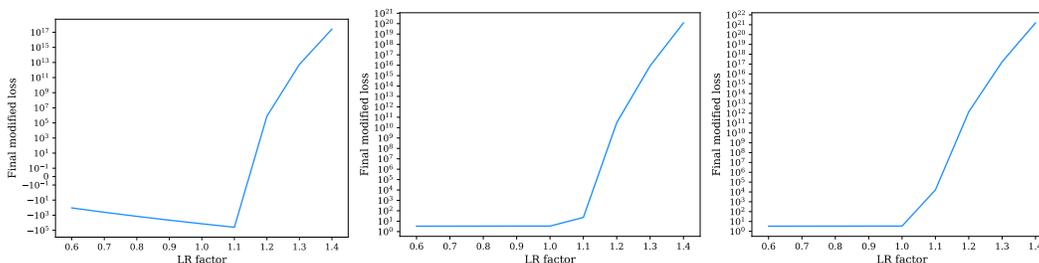
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809



(a) GPT2-Small Model.

(b) GPT2-Medium Model.

Figure 7: EoS and EoC Tests of GPT-Small and GPT-Medium models.



(a) No convexity compensation (= 0). (b) Optimal compensation ( $\approx 408$ ). (c) Compensation = 600.

Figure 8: Effects of different convexity compensations. We train a 43M GPT-like model with AdamW on FineWeb-10B and run EoS tests with different convexity compensation values and report loss on local step 25: (a) 0 (no compensation), (b)  $\approx 408$  (optimal), and (c) 600.

**Operational criteria.** In practice, we declare a checkpoint to be in an EoS phase if the estimated ratio  $\eta_*^{\lambda_p} / \eta_*$  falls within a tolerance band (we use  $[0.9, 1.1]$ ). Separately, we declare a checkpoint to exhibit EoC if the estimated convexity compensation  $\lambda_p$  is consistently bounded away from zero. We provide a full example report of this simulation pipeline in Figure 11.

## G.2 RESULTS OF GPT-SMALL AND GPT-MEDIUM

Figure 7 shows the EoS/EoC measurements for GPT-small and GPT-medium. For both model sizes, we observe a consistent EoC phenomenon across training: the estimated convexity compensation  $\lambda_p$  stays reliably bounded away from zero over a wide range of checkpoints. In addition, the EoS signature remains visible through a large portion of training. Overall, the qualitative behaviors observed in GPT-extra-small persist for GPT-small and GPT-medium.

## G.3 ROBUSTNESS OF THE CONVEXITY COMPENSATION

Here is an example to show the robustness of the convexity compensation. Even if the convexity compensation is a little above the threshold, we can still observe the sharp phase transition.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

#### G.4 DETAILED EoS BATCH-SIZE EXPERIMENT SETTING

We study five datasets (ArXiv, C4, DCLM, Wikipedia, and OpenWebText) using GPT-extra-small. We fix the optimizer (AdamW with WD = 0.1) and use the WSD scheduler throughout. For each dataset, we first estimate the critical batch size (CBZ) using the local-simulation method of Merrill et al. (2025).

To vary the batch size, we follow the scaling rules in Malladi et al. (2024); Marek et al. (2025): when increasing the batch size from  $B$  to  $kB$ , we keep all hyperparameters fixed except we set

$$\beta_2 \leftarrow \beta_2^k, \quad \eta \leftarrow \sqrt{k} \eta.$$

For each batch size  $B$ , we estimate the EoS diagnostic ratio  $\eta_*^{\lambda^p} / \eta_*^0$ , where  $\eta_*^0$  denotes the critical LR without compensation. We then interpolate this ratio as a function of  $B$  and define the *EoS batch size* (EBZ) as the smallest  $B$  for which the ratio reaches 1.1.

Figure 4a compares EBZ to CBZ. Across datasets, EBZ is consistently smaller than CBZ, while the two quantities are strongly correlated. Empirically, we observe an approximately linear relation, with  $\text{CBZ} \approx 3 \times \text{EBZ}$  in our settings.

#### G.5 ABLATION STUDY ON LR ABRUPT DROP

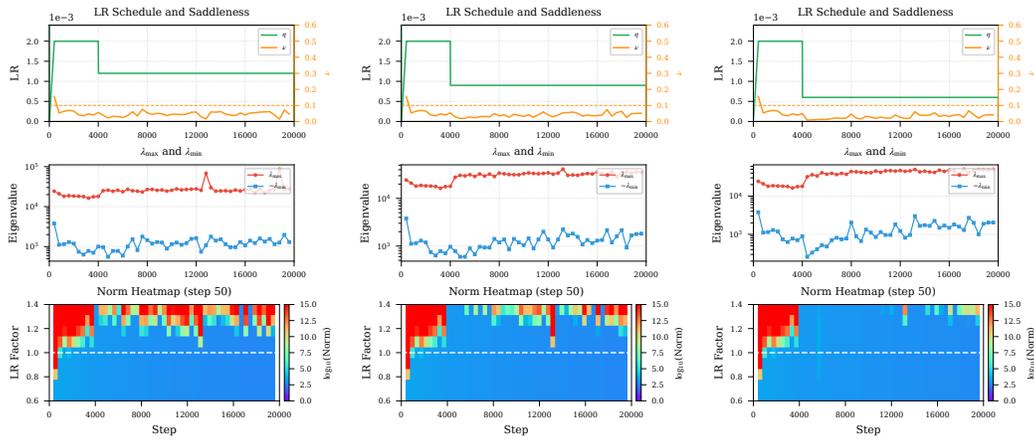
We evaluate two-step scheduler variants in Figure 9. We find that the drop ratio strongly affects whether EoS re-emerges after the drop. In particular, when the drop ratio is sufficiently small (typically below 0.5 in our experiments), the EoS signature returns quickly and remains stable. Within the range of peak LRs tested in Figure 9, this re-emergence behavior is qualitatively robust.

#### G.6 ABLATION STUDY ON LR DECAY FOR EoC

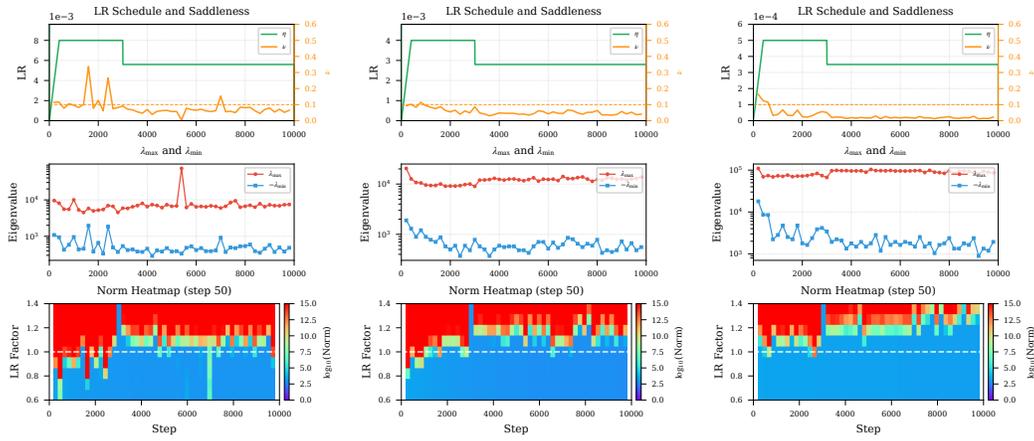
**EoC vanishes in LR decay.** We further find that EoC is influenced by LR decay, illustrated in Figure 3. As the LR becomes sufficiently small, the minimum eigenvalue gradually approaches zero, suggesting that EoC diminishes only in the small-LR regime.

#### G.7 ABLATION STUDY ON LR LINEAR DECAY

We evaluate different linear-decay configurations in Figure 10. Consistent with the two-step ablations, the *decay ratio* plays a primary role in determining whether and when EoS re-emerges during the decay phase. Moreover, the transition time depends on the decay length: longer decay schedules lead to a slower and more gradual transition, roughly scaling with the number of decay steps in our settings.



(a) Twostep, drop ratio = 0.6. (b) Twostep, drop ratio = 0.45. (c) Twostep, drop ratio = 0.3.



(d) Twostep, peak LR = 0.008. (e) Twostep, peak LR = 0.004. (f) Twostep, peak LR = 0.0005.

Figure 9: **Ablations for the Peak LR of abrupt drop.** We train a 43M GPT-like model and vary the peak LR and the drop rate of the two-step scheduler. (a)–(c) Peak LR is fixed at 0.002, with drop rates 0.6, 0.45, and 0.3; the drop step is 4000. (d)–(f) Drop rate is fixed at 0.7, with peak LR values 0.008, 0.004, and 0.0005; the drop step is 3000.

## H DETAILED DERIVATION OF EoC AND EoS

In this part, we reorganize our main empirical methods for the measurement of EoC & EoS and the corresponding theoretical analysis throughout the empirical methods.

### H.1 VANILLA SGD

Fix a training checkpoint  $w_{t_0}$ . We study the training dynamics through a stochastic quadratic surrogate built from mini-batches, which induces a linear stochastic dynamical system under SGD. This viewpoint reveals two distinct ways the quadratic model can fail: (i) *negative-curvature instability* due to  $\lambda_{\min}(\nabla^2 \mathcal{L}(w_{t_0})) < 0$  (Edge of Convexity, EoC), and (ii) *positive-curvature instability* when the LR exceeds a critical threshold (Edge of Stability, EoS). Since these two instabilities are coupled in practice, we further introduce a convexity-compensated surrogate that suppresses the negative-curvature instability while preserving the positive-curvature stability threshold, enabling a scalable local characterization of EoS.

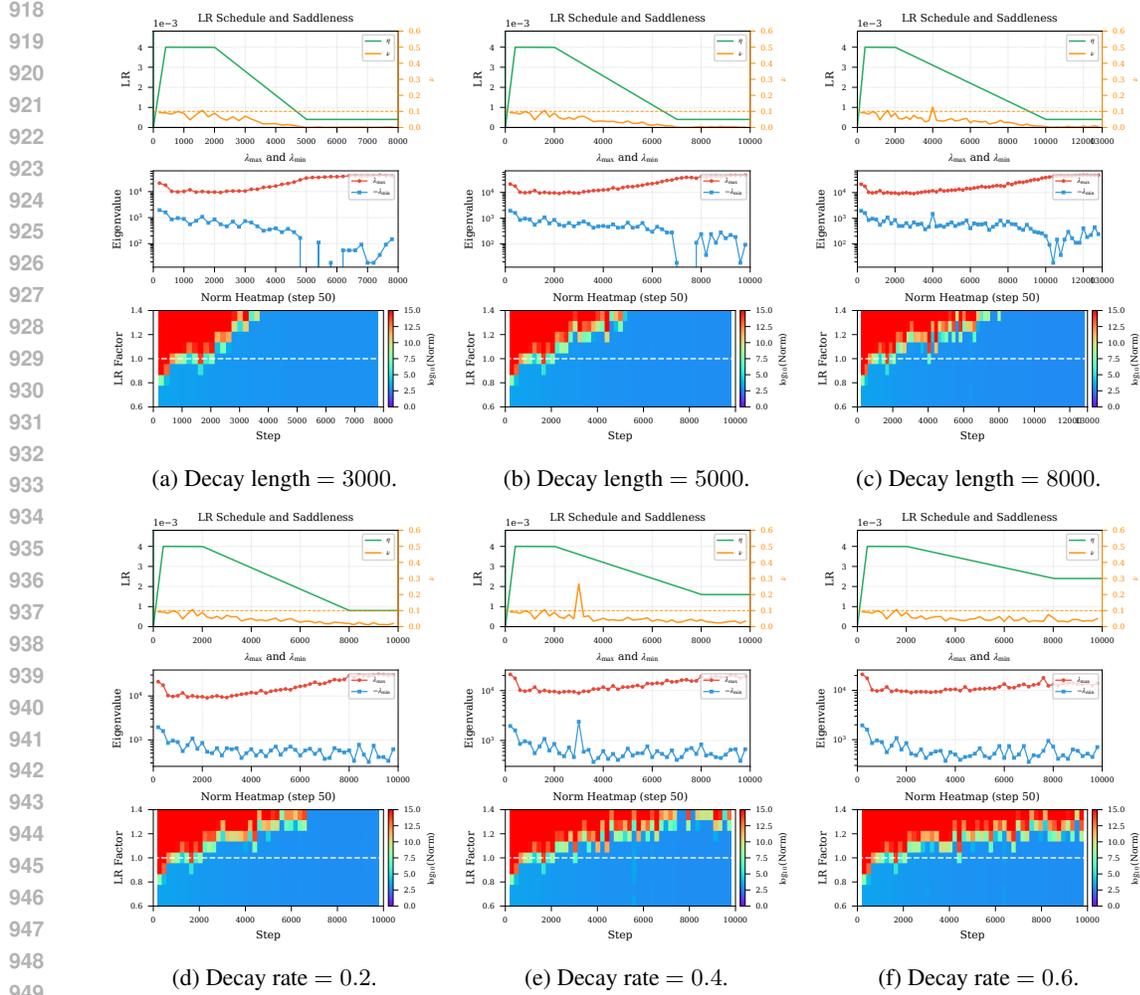


Figure 10: **Ablations for the linear decay phase.** We train a 43M GPT-like model and vary the decay length and the decay rate of the linear decay phase. (a)–(c) The decay rate is fixed at 0.1, with decay lengths of 3000, 5000, and 8000 steps. (d)–(f) The decay length is fixed at 6000 steps, with decay rates of 0.2, 0.4, and 0.6.

### H.1.1 STOCHASTIC STABILITY ANALYSIS

For stochastic optimization, mini-batch sampling introduces additional variance that affects stability beyond what  $\lambda_{\max}$  alone can capture. To analyze this, for any fixed checkpoint  $w_{t_0}$  during training, we consider a local quadratic approximation with stochastic Hessians:

$$Q_{\mathcal{B}_t}(w) := \mathcal{L}_{\mathcal{B}_t}(w_{t_0}) + (w - w_{t_0})^\top \nabla \mathcal{L}_{\mathcal{B}_t}(w_{t_0}) + \frac{1}{2}(w - w_{t_0})^\top \nabla^2 \mathcal{L}_{\mathcal{B}_t}(w_{t_0})(w - w_{t_0}), \quad (8)$$

where  $\mathcal{B}_t$  is a mini-batch sampled at time step  $t \geq t_0$ . Applying SGD with LR  $\eta$  to this approximation as

$$w_{t+1} = w_t - \eta \nabla_w Q_{\mathcal{B}_t}(w), \quad (9)$$

which yields a linear dynamical system whose stability depends on both the expected Hessian and its covariance across batches. The proof of Theorem H.1 is given in Appendix I.

**Theorem H.1** (Linear System for SGD). *Assume the batches are i.i.d., and there exists a  $w_*$  such that:*

$$\nabla \mathcal{L}(w_{t_0}) + \nabla^2 \mathcal{L}(w_{t_0})(w_* - w_{t_0}) = 0. \quad (10)$$

972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 981  
 982  
 983  
 984  
 985  
 986  
 987  
 988  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001  
 1002  
 1003  
 1004  
 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025

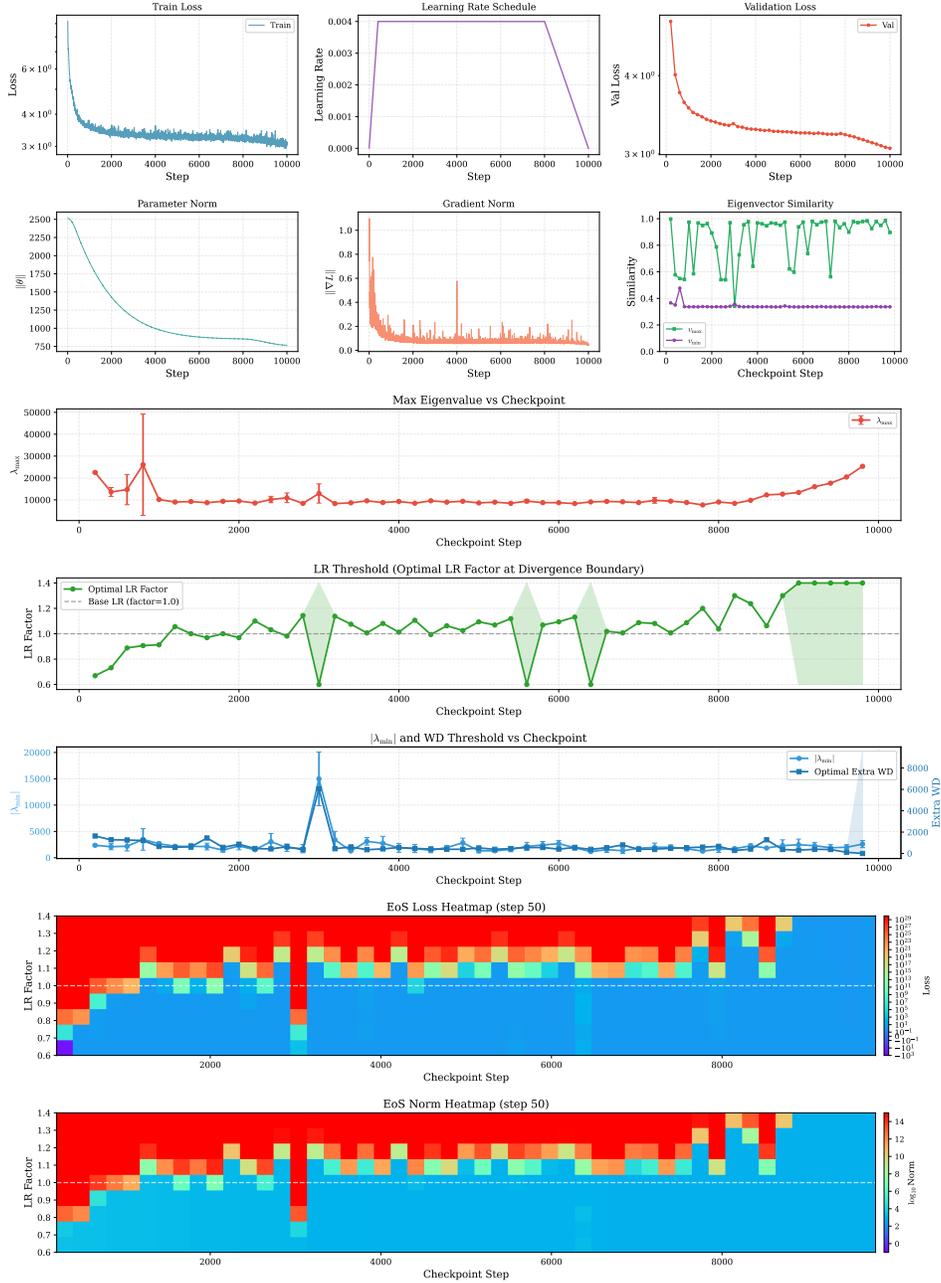


Figure 11: **Summary of a representative training run.** We show training dynamics for GPT-extra-small trained on FineWeb-10B using AdamW with LR 0.004, weight decay 0.1, and the WSD learning-rate scheduler.

If we apply SGD to  $Q_{\mathcal{B}_t}(\mathbf{w})$  as (9), the evolution of the first and second moments of the parameter trajectory  $\mathbf{w}_t$ ,

$$\begin{aligned}\mathbf{m}_t^{(1)} &:= \mathbb{E}[\mathbf{w}_t - \mathbf{w}_*], \\ \mathbf{m}_t^{(2)} &:= \text{vec}\left(\mathbb{E}[(\mathbf{w}_t - \mathbf{w}_*)(\mathbf{w}_t - \mathbf{w}_*)^\top]\right),\end{aligned}$$

is given by:

$$\begin{bmatrix} \mathbf{m}_{t+1}^{(1)} \\ \mathbf{m}_{t+1}^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbb{E}[\mathbf{M}_t], & 0 \\ \mathbf{S}, & \mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t] \end{bmatrix} \begin{bmatrix} \mathbf{m}_t^{(1)} \\ \mathbf{m}_t^{(2)} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix}. \quad (11)$$

where

$$\mathbf{M}_t := I - \eta \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}),$$

and  $\mathbf{S}$  and  $\mathbf{b}$  are constants depending on the batch distribution,  $\mathbf{w}_{t_0}$ , and  $\mathbf{w}_*$ .

The stability of this linear system is governed by  $\|\mathbb{E}[\mathbf{M}_t]\|$  and  $\|\mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t]\|$ . The first controls the mean trajectory, while the second controls the variance. The following corollaries give the stability conditions; see Appendix J for the proof.

**Corollary H.2** (Negative Curvature Divergence). *Let  $\mathbf{H}_{t_0} := \nabla^2 \mathcal{L}(\mathbf{w}_{t_0})$ . If  $\lambda_{\min}(\mathbf{H}_{t_0}) < 0$ , then for any  $\eta > 0$ , the quadratic model is unbounded below, and the SGD iterates on  $Q_{\mathcal{B}_t}$  diverge to  $-\infty$ .*

The above corollary accounts for the impact of strictly negative curvature, which would lead the quadratic objective to decrease without bound along the direction of the negative eigenvalues. After that, the surrogate iterates run away, and the quadratic loss diverges to  $-\infty$ .

As we mentioned, this failure mode of negative curvature is separated from the instability that is typically induced by  $\eta$ . To study the latter, we factor out the negative-curvature directions and focus on the subspace where the quadratic model is locally convex and define the corresponding critical LR.

**Definition H.3** (Critical LR for SGD). *Let  $\mathbf{H}_{t_0} := \nabla^2 \mathcal{L}(\mathbf{w}_{t_0})$ , we denote  $\Pi_{t_0}^+$  as the orthogonal projection onto its eigenspace of the nonnegative eigenvalues. We define*

$$\eta_*(\mathbf{w}_{t_0}) := \sup\{\eta > 0 : \sup_t \mathbb{E}\|\Pi_{t_0}^+(\mathbf{w}_t - \mathbf{w}_*)\|^2 < \infty\}.$$

Definition H.3 is intentionally projected onto the nonnegative-curvature subspace of  $\mathbf{H}_{t_0}$ . It is therefore well-defined even when  $\lambda_{\min}(\mathbf{H}_{t_0}) < 0$  and should be interpreted as the critical LR threshold for the positive-curvature instability. The following corollary shows this intuition clearly.

**Corollary H.4** (Positive Curvature Divergence). *Let  $\mathbf{H}_{t_0} := \nabla^2 \mathcal{L}(\mathbf{w}_{t_0})$  and let  $\eta^*(\mathbf{w}_{t_0})$  be defined in Definition 2.2, then the SGD trajectory is unbounded above in projected second moment  $\mathbb{E}\|\Pi_{t_0}^+(\mathbf{w}_t - \mathbf{w}_*)\|^2$  diverges to  $+\infty$  when  $\eta > \eta^*(\mathbf{w}_{t_0})$ .*

For neural network training, once we replace the actual loss with  $Q_{\mathcal{B}_t}(\mathbf{w})$ , it is typically observed that the linear system formed by quadratic approximation diverges quickly, while the true dynamics of the neural network remain stable (Figure 1a, Figure 1b).

So a natural question is: which regime is the neural network training in, negative curvature divergence or positive curvature divergence? We identify: it sits on an *edge* where either divergence mechanism is *just at the threshold* of breaking the local quadratic approximation.

For the first negative divergence regime, we first define a quantity termed saddleness.

**Definition H.5.** *Given any parameter point  $\mathbf{w}_{t_0}$  whose population Hessian  $\mathbf{H}_{t_0}$  has minimal, maximal eigenvalue denoted by  $\lambda_{\min}$  and  $\lambda_{\max}$ . We define the saddleness  $\nu$  at this point as*

$$\nu(\mathbf{w}_{t_0}) = \frac{\min\{0, -\lambda_{\min}\}}{\lambda_{\max}}.$$

Intuitively,  $\nu(\mathbf{w}_{t_0})$  measures how ‘‘saddle-like’’ the local curvature is, by comparing the magnitude of the most negative curvature to the strongest positive curvature. If  $\mathbf{H}_{t_0} \succeq 0$  (locally convex), then

$\nu(\mathbf{w}_{t_0}) = 0$ . A larger  $\nu$  means the negative curvature is comparable to  $\lambda_{\max}$  (strong saddle structure), while a small  $\nu$  means the negative curvature exists but is weak relative to the dominant positive directions.

One important observation is: For LLM pre-training and image classification, for all steps  $t$  after the warmup phase, we observe that  $\nu \leq 0.1$ , which leads to the first edge area that neural network training is in

**Edge of Convexity (EoC).** Negative eigenvalues persist during training, but  $\nu \approx 0$ , and true dynamics remains stable, while the quadratic approximation is unbounded below, and the iterates diverge to  $-\infty$  regardless of the choice of LR.

And for the positive divergence regime, we identify that the neural network training is at the following edge area

**Edge of Stability (EoS) at  $\mathbf{w}_{t_0}$ .** The current LR  $\eta_{t_0} \approx \eta_*(\mathbf{w}_{t_0})$ , a small positive perturbation of  $\eta$  leads the local quadratic approximation  $Q_{\mathcal{B}_t}(\mathbf{w})$  to  $+\infty$ .

In the full-batch (deterministic) case,  $\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) \equiv \mathbf{H}_{t_0}$  and the update on the locally convex subspace reduces to  $\mathbf{x}_{t+1} = (\mathbf{I} - \eta \mathbf{H}_{t_0}) \mathbf{x}_t$ . Hence stability is equivalent to  $\rho(\mathbf{I} - \eta \mathbf{H}_{t_0}) < 1$ , i.e.,  $0 < \eta < 2/\lambda_{\max}(\mathbf{H}_{t_0})$ , so  $\eta_*(\mathbf{w}_{t_0}) = 2/\lambda_{\max}(\mathbf{H}_{t_0})$ . Therefore our criterion  $\eta_{t_0} \approx \eta_*(\mathbf{w}_{t_0})$  recovers the classical EoS condition  $\eta \lambda_{\max} \approx 2$  in Cohen et al. (2021).

These two edge areas are geometrically distinct but coupled in the training dynamics. A naive approach is to decompose the parameter space  $\mathbb{R}^P$  into orthogonal subspaces:  $U_+$  spanned by eigenvectors with positive eigenvalues, and  $U_-$  spanned by eigenvectors with non-positive eigenvalues. One could then project the trajectory onto  $U_+$  and  $U_-$  to detect each mode. However, computing this eigendecomposition is prohibitive for high-dimensional neural networks. This leads to our first challenge: *How to separate the effect of EoS and EoC in quadratic approximation?*

Even if we can separate EoS and EoC, unlike the small-scale full-batch setting—where the EoS can be characterized via  $\eta$  and  $2/\lambda_{\max}$ —the stability of quadratic approximation for large-scale stochastic training is determined by:

$$\mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t] = \mathbb{E}[\mathbf{M}_t] \otimes \mathbb{E}[\mathbf{M}_t] + \text{Cov}(\mathbf{M}_t),$$

where  $\mathbf{M}_t := \mathbf{I} - \eta \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_t)$  and  $\text{Cov}(\mathbf{M}_t)$  is the covariance of  $\text{vec}(\mathbf{M}_t)$  induced by mini-batch sampling. The computation is prohibitive due to two reasons:

Furthermore, unlike the full-batch case where  $\eta_* = 2/\lambda_{\max}$ , we cannot give an explicit expression for  $\eta_*$  in the stochastic setting for two reasons:

**(i) Batch noise shifts the stability threshold.** The second-moment stability depends on  $\|\mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t]\|$ , which decomposes as:

$$\mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t] = \mathbb{E}[\mathbf{M}_t] \otimes \mathbb{E}[\mathbf{M}_t] + \text{Cov}(\mathbf{M}_t),$$

where  $\text{Cov}(\mathbf{M}_t)$  is the covariance of  $\text{vec}(\mathbf{M}_t)$  induced by mini-batch sampling. Without the covariance term, the spectral radius would be  $\max_i (1 - \eta \lambda_i)^2$ , yielding  $\eta_* = 2/\lambda_{\max}$ . However,  $\text{Cov}(\mathbf{M}_t)$  can shift the spectral radius significantly, causing  $\eta_*$  to deviate from  $2/\lambda_{\max}$ .

**(ii) The true  $\lambda_{\max}$  is computationally inaccessible.** Even if the threshold were  $2/\lambda_{\max}$ , computing  $\lambda_{\max}$  of the full-batch Hessian  $\nabla^2 \mathcal{L}(\mathbf{w})$  is infeasible for large-scale training, as it requires access to the entire dataset (or population). While one can estimate  $\lambda_{\max}$  from mini-batch Hessians, these estimates have high variance and may not reflect the true full-batch spectrum.

Hence, our second challenge is, *how to measure EoS without estimation of  $\eta_*$ ?*

These challenges motivate our simulation-based approach, which directly probes the stability boundary without requiring explicit knowledge of  $\lambda_{\max}$  or  $\text{Cov}(\mathbf{M}_t)$ .

In the following subsections, we propose practical methods to detect each regime by simulating local trajectories, bypassing explicit eigendecomposition.

---

### 1134 H.1.2 EoC TEST VIA MINIMAL EIGENVALUE ESTIMATION

1135  
 1136 At a checkpoint  $\mathbf{w}_{t_0}$ , EoC corresponds to persistent negative population curvature, i.e.,  
 1137  $\lambda_{\min}(\nabla^2 \mathcal{L}(\mathbf{w}_{t_0})) < 0$ , under which the quadratic approximation is unbounded below and fails  
 1138 regardless of the LR. Our goal is to estimate  $\lambda_{\min}(\nabla^2 \mathcal{L}(\mathbf{w}_{t_0}))$  using only stochastic mini-batch  
 1139 access, without explicit eigendecomposition.

1140 **Stochastic Power Iteration Bisection.** To this end, we provide a stochastic power iteration bisection  
 1141 method that can estimate the smallest eigenvalue of a large matrix  $\mathbf{A}$  without costly eigendecom-  
 1142 position. Specifically, given a matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , and we get a sequence of independent random  
 1143 observations  $\{\mathbf{A}_t\}_{t=1}^T$  with  $\mathbb{E}\mathbf{A}_t = \mathbf{A}$  and  $\text{Cov}(\mathbf{A}_t) = \mathbf{\Sigma}$ . We first decide a large range for the  
 1144 minimal eigenvalue as  $|\lambda_{\min}(\mathbf{A})| \in [0, \Lambda]$ . We then consider the following power iteration

$$1145 \mathbf{x}_{t+1} = (\mathbf{I} - \eta(\mathbf{A}_t + \lambda \mathbf{I})) \mathbf{x}_t.$$

1146  
 1147 Intuitively, consider applying the above update for  $T$  steps. If  $\lambda < \lambda^*$ , the mean dynamics admit an  
 1148 expansive mode and  $\|\mathbf{x}_T\|$  explodes; if  $\lambda \geq \lambda^*$  and  $\eta$  is sufficiently small, the iteration becomes stable.  
 1149 Algorithm 1 formalizes this intuition via an EXPLODETEST predicate and performs  $N$  rounds of  
 1150 bisection to return the smallest stable  $\hat{\lambda}$ . The high-probability accuracy of this estimator is guaranteed  
 1151 by the following theorem.

1152 **Theorem H.6** (Error bound for stochastic power-iteration bisection). *Let  $\mathbf{A} \in \mathbb{S}^d$  be a symmetric*  
 1153 *matrix with minimal eigenvalue  $\lambda_{\min}$  and a unit eigenvector  $\mathbf{u}_{\min}$ . At each inner iteration  $t$ , we*  
 1154 *observe  $\mathbf{A}_t = \mathbf{A} + \Delta_t$ , where  $\Delta_t \in \mathbb{S}^d$  satisfies  $\mathbb{E}[\Delta_t \mid \mathcal{F}_{t-1}] = \mathbf{0}$ , where the filtration  $\mathcal{F}_t :=$*   
 1155  *$\sigma(\Delta_0, \dots, \Delta_{t-1})$  and  $\|\Delta_t\| \leq \sigma$  almost surely.<sup>1</sup> Assume the initialization  $\mathbf{x}_0$  satisfies*

$$1156 \frac{|\mathbf{x}_0^\top \mathbf{u}_{\min}|}{\|\mathbf{x}_0\|_2} \geq \varepsilon \tag{12}$$

1157  
 1158 for some absolute constant  $\varepsilon > 0$ .<sup>2</sup> Run Algorithm 1 with inner horizon  $T$ , bisection rounds  $N$ ,  
 1159 search range  $[0, \Lambda]$ , and inner stepsize  $\eta$  satisfying  $\eta = o(1/\sqrt{T})$ ,  $\eta = \omega(1/T)$ . Let  $\hat{\lambda}_{\text{shift}}$  be the  
 1160 estimated shift returned by Algorithm 1 and define  $\hat{\lambda}_{\min} := -\hat{\lambda}_{\text{shift}}$ . Then for any  $\delta \in (0, 1)$ , with  
 1161 probability at least  $1 - \delta$ ,

$$1162 \left| \hat{\lambda}_{\min} - \lambda_{\min} \right| = \tilde{\mathcal{O}} \left( \frac{1}{\eta T} + \frac{\sigma \sqrt{\log d}}{\sqrt{T}} \right) + \mathcal{O} \left( \frac{\Lambda}{2^{-N}} \right),$$

1163 where  $\tilde{\mathcal{O}}(\cdot)$  hides polylogarithmic factors in  $(d, 1/\delta)$  and  $\mathcal{O}(\cdot)$  hides absolute constants.

1164  
 1165 Theorem H.6 justifies Algorithm 1 as a consistent estimator of the minimal eigenvalue. The detailed  
 1166 proof of this theorem is provided in Appendix M.

1167  
 1168 **EoC Test Framework.** We now use the stochastic power iteration method to estimate the minimal  
 1169 eigenvalue of  $\mathbf{H}_{t_0}$  for any chosen checkpoint  $\mathbf{w}_{t_0}$ , and we observe the consistent negative eigenvalues  
 1170 during each stage of training across various training setups, as discussed in Appendix E. Also, as  
 1171 shown in Figure 1a, once we replace the instant loss function with a local quadratic approximation  
 1172 and run several iterations, we see that the approximated loss sharply diverges to  $-\infty$  while the real  
 1173 loss curve remains stable. Combining the above 2 steps for the detection of EoC gives an obvious  
 1174 takeaway: **EoS is always observed across the neural network training.** The detailed description  
 1175 for this framework is provided in Algorithm 2.

### 1179 H.1.3 EoS TEST VIA CONVEXITY-COMPENSATED LOCAL SIMULATION

1180  
 1181 Now we introduce our framework for the detection of EoS, which is more challenging compared with  
 1182 the straightforward detection of EoC. Detecting EoS requires projecting the parameter trajectory onto  
 1183 the subspace spanned by eigenvectors with positive eigenvalues, which is computationally prohibitive.  
 1184 However, we observe a significant gap between the maximal and minimal Hessian eigenvalues in  
 1185 practice:

1186 <sup>1</sup>The same statement holds under i.i.d. mean-zero noise as a special case.

1187 <sup>2</sup>Condition (12) can be achieved by a normal initialization  $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}_d)$  with high probability. We can  
 freely generalize this theorem to the case with normal initialization, shown at the end of Appendix M.

**Observation H.7.** For LLM pre-training, for all steps  $t$  after the warmup phase, we observe that

$$\lambda_{\max} \geq 10 \cdot |\lambda_{\min}|,$$

i.e., saddleness  $\nu \leq 0.1$ . See Appendix E for empirical validations.

This motivates us to add a convexity compensation term that eliminates negative curvature effects while preserving positive curvature instability. The idea is simple: by adding  $\frac{\lambda_P}{2} \|\boldsymbol{\theta}\|^2$  to the local loss, we shift all Hessian eigenvalues upward by  $\lambda_P$ . If  $\lambda_P \geq |\lambda_{\min}|$ , the effective Hessian becomes positive semi-definite, removing negative curvature effects. Figure 8 shows that an appropriate compensation can decouple EoS from the EoC. Furthermore, since  $\kappa \geq 10$  by the above observation, this shift only slightly perturbs the large positive eigenvalues that govern EoS.

To make the above idea rigorous, we first let  $\boldsymbol{\theta} := \mathbf{w} - \mathbf{w}_{t_0}$  denote the displacement from the current parameters. For each mini-batch  $\mathcal{B}_t$ , we define the compensated local loss as:

$$G_t(\boldsymbol{\theta}) := \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \boldsymbol{\theta}^\top \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \frac{1}{2} \boldsymbol{\theta}^\top \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) \boldsymbol{\theta} + \frac{\lambda_P}{2} \|\boldsymbol{\theta}\|^2, \quad (13)$$

where  $\lambda_P \geq 0$  is called the compensation coefficient. The derivation of this linear system and its corresponding stability analysis is provided in Appendix I and Appendix J.

Now we can extend the definition of critical LR in Definition H.3 with the compensation coefficient  $\lambda_P$  as

$$\eta_{\star}^{\lambda_P}(\mathbf{w}_{t_0}) := \sup\{\eta > 0 : \sup_t \mathbb{E} \|\Pi_{t_0}^{+\prime}(\mathbf{w}_t - \mathbf{w}_{\star})\|^2 < \infty\}, \quad (14)$$

where  $\Pi_{t_0}^{+\prime}$  denotes the orthogonal projection onto  $\mathbf{H}_{t_0}$ 's eigenspace spanned by the eigenvectors whose eigenvalue  $\mu$  satisfies  $\mu + \lambda_P > 0$ .

The next question is to decide how large the compensation coefficient should be. Ideally, we should set  $\lambda_P$  as small as possible and make  $G_t$  convex in expectation over time  $t$ , which gives  $\lambda_P = |\lambda_{\min}|$ .

In practice, we set the compensation coefficient using the estimation by Algorithm 1 proposed in Appendix H.1.2 as  $\lambda_P = \hat{\lambda}$  so that the compensated curvature  $\mathbf{H}_{t_0} + \lambda_P \mathbf{I}$  is (approximately) positive semidefinite, which suppresses the EoC mode while leaving the positive-curvature stability threshold nearly unchanged. Specifically, we can show that the perturbation of the critical LR with  $\lambda_P$ , denoted by  $\eta_{\star}^{\lambda_P}$  can be bounded by the condition number  $\kappa$ .

**Theorem H.8** (EoS Detection Error for SGD). *Given the Batch covariance matrix  $\text{Cov}(\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})) \preceq \sigma^2 \mathbf{I}$ , consider the vanilla SGD, then we have*

$$\frac{\eta_{\star}^{\lambda_P}}{\eta_{\star}^0} \in \left[ \frac{\lambda_{\max}(\lambda_{\max} + \lambda_P)}{(\lambda_{\max} + \lambda_P)^2 + \sigma^2}, \frac{\lambda_{\max}^2 + \sigma^2}{\lambda_{\max}(\lambda_{\max} + \lambda_P)} \right],$$

where  $\lambda_{\max}$ ,  $\lambda_{\min}$  represent the maximal and minimal eigenvalues of  $\nabla^2 \mathcal{L}(\mathbf{w}_{t_0})$ . Specifically, if  $\lambda_P = -\lambda_{\min}$ , then we have

$$\frac{\eta_{\star}^{\lambda_P}}{\eta_{\star}^0} \in \left[ \frac{\kappa^2 + \kappa}{(\kappa + 1)^2 + \alpha}, \frac{\kappa^2 + \alpha}{\kappa^2 + \kappa} \right]$$

where  $\alpha = \frac{\sigma^2}{|\lambda_{\min}|^2}$ ,  $\kappa := 1/\nu = \lambda_{\max}/|\lambda_{\min}|$ .

The proof of Theorem H.8 is given in Appendix K. Combining this theorem with Observation H.7, we provide a framework for detecting EoS decoupling from the impact of EoC.

**EoS Test Framework.** We summarize the above discussion into a unified framework for the detection of EoS:

- **Step1: remove the EoC mode with Convexity compensation.** When  $\lambda_{\min}(\nabla^2 \mathcal{L}(\mathbf{w}_{t_0})) < 0$ , the quadratic surrogate is unbounded below and diverges regardless of  $\eta$  (EoC), which would mask the EoS mode. Using the procedure in Algorithm 1, we obtain  $\lambda_P \approx -\lambda_{\min}$  and define the compensated mini-batch surrogate  $G_{\mathcal{B}_t}$  with curvature  $\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \lambda_P \mathbf{I}$ . Theorem H.8 together with Observation H.7 ensures that this compensation perturbs the positive-curvature stability threshold only by a relative factor  $1 + O(1/\kappa)$ .

- **Step2: probe the stability boundary by local simulation.** Fix a single mini-batch sequence  $\{\mathcal{B}_t\}_{t=0}^{T-1}$  and simulate, for all LR factor  $\alpha \in \{\alpha_1 < \dots < \alpha_K\}$ ,

$$\boldsymbol{\theta}_{t+1}^{(\alpha)} = \boldsymbol{\theta}_t^{(\alpha)} - \alpha\eta\nabla G_t(\boldsymbol{\theta}_t), \quad \boldsymbol{\theta}_0^{(\alpha)} = 0,$$

as summarized in Algorithm 3.

- **Step3: Estimate the critical LR and decide the existence of EoS.** We declare *explosion* if  $\max_{0 \leq t \leq T} \|\boldsymbol{\theta}_t^{(\alpha)}\|_2 \geq M$  and estimate the critical LR  $\eta_{\star}^{\lambda_P}$  by

$$\widehat{\eta}_{\star}^{\lambda_P}(\mathbf{w}_{t_0}) := \min\{\alpha_k\eta : \theta^{(\alpha_k)} \text{ explodes in } T \text{ steps}\},$$

with the convention that if no  $\alpha_k$  explodes, we only obtain the lower bound  $\widehat{\eta}_{\star}^{\lambda_P}(\mathbf{w}_{t_0}) > \alpha_K\eta$ . We then quantify proximity to EoS via the ratio  $r(\mathbf{w}_{t_0}) := \eta/\widehat{\eta}_{\star}^{\lambda_P}(\mathbf{w}_{t_0})$ : When  $r(\mathbf{w}_{t_0})$  is close to 1, the current LR is close to  $\widehat{\eta}_{\star}^{\lambda_P}$ , indicating proximity to the local EoS regime.  $r \ll 1$  indicates a stable margin, and  $r > 1$  indicates the compensated surrogate is already unstable at  $\eta$ .

**Remark H.9.** *Most training runs use a LR schedule, whereas our probe assumes a constant  $\eta$ . This is not a conflict: the local simulation spans only a small number of steps relative to training, during which  $\eta$  changes negligibly for commonly used schedules (e.g., WSD or cosine), so we treat it as approximately constant.*

## H.2 ADAPTIVE GRADIENT METHODS

Modern neural networks, such as transformer-based large language models, are usually trained with adaptive optimizers such as Adam (Kingma, 2014) and its variants (Loshchilov & Hutter, 2017). In this section, we extend the previous definition and detection framework of EoS to a class of Adaptive optimizers. Specifically, we consider a class of Adaptive Gradient Methods (AGMs) (Li et al., 2025) with the following general form:

$$\mathbf{m}_{t+1} = \beta_1 \cdot \mathbf{m}_t + (1 - \beta_1) \cdot \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_t) \quad (15)$$

$$\mathbf{P}_{t+1} = \Psi(\mathbf{P}_t, \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_t), \mathbf{w}_t) \quad (16)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{P}_{t+1}^{-1} \mathbf{m}_{t+1} - \eta \lambda_{\text{WD}} \mathbf{w}_t, \quad (17)$$

where momentum  $\mathbf{m}_t \in \mathbb{R}^d$  and the preconditioner mapping  $\Psi : \mathbb{S}_{++}^d \times \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{S}_{++}^d$ . Here  $\Psi$  is assumed to be a slowly-varying mapping, so we can treat  $\mathbf{P}_t$  as fixed over a short time window. Note that AdamW, SGDM, and SGD are special cases of this AGM class. For preconditioned optimizers, the relevant spectrum is that of the preconditioned Hessian  $\mathbf{P}_t^{-1} \nabla^2 \mathcal{L}(\mathbf{w}_t)$ ; we use  $\mu_{\max}$  and  $\mu_{\min}$  to denote its largest and smallest eigenvalues when the context is clear.

Similar to the case for SGD, we first build the linear dynamical system for AGMs when applying local quadratic approximation with a convexity compensation term  $Q_{\mathcal{B}_t}(\mathbf{w}_t)$  for the loss function  $\mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_t)$ , and then the following stability analysis would lead us to the definitions of EoS and EoC for AGMs.

### H.2.1 EoS AND EOC FOR AGMS

Fix a checkpoint  $\mathbf{w}_{t_0}$  and freeze the preconditioner over a short window as  $\mathbf{P} := \mathbf{P}_{t_0}$ . Under the same local quadratic approximation as in Appendix H.1, an AGM update can be written as a linear stochastic system in an augmented state  $\mathbf{z}_t$  (e.g., stacking the momentum  $\mathbf{m}$  and parameter error  $\mathbf{w}_t - \mathbf{w}_{\star}$ , where  $\mathbf{w}_{\star}$  is given in (10))

$$\mathbf{z}_{t+1} = \mathbf{M}_t(\eta) \mathbf{z}_t + \mathbf{q}_t.$$

Detailed definitions and the derivation of this system can be found in Appendix L. In this system, the relevant curvature for stability is the spectrum of the preconditioned Hessian  $\mathbf{P}^{-1} \mathbf{H}$  or equivalently  $\widetilde{\mathbf{H}}_{t_0} := \mathbf{P}^{-1/2} \mathbf{H}_{t_0} \mathbf{P}^{-1/2}$ , whose eigenvalues we denote by  $\mu_1, \dots, \mu_d$  (with  $\mu_{\max}, \mu_{\min}$  as extremes).

Compared with vanilla SGD, the critical LR  $\eta_{\star}$  is accordingly defined here and all stability statements in Appendix H.1 carry over after replacing the Hessian eigenvalues  $\{\lambda_i\}$  by the preconditioned eigenvalues  $\{\mu_i + \lambda_{\text{WD}}\}$  (incorporating weight decay as an additional isotropic shift):  $\mu_i + \lambda_{\text{WD}} \geq$

1296 0,  $\eta \leq \eta_*$ ,  $\forall i \in [d]$ . Correspondingly, we can define the saddleness with the weight decay factor as  
 1297  $\nu_{\lambda_{\text{WD}}}$ .

1298 In the same way, the above stable conditions imply the definition of EoC and EoS for AGMs:  
 1299

- 1300 • **Edge of Convexity (EoC) for AGMs.** When  $\mu_{\min} + \lambda_{\text{WD}} < 0$ , but  $\nu_{\lambda_{\text{WD}}} \approx 0$ , and the quadratic  
 1301 approximation is unbounded below, and the iterates diverge to  $-\infty$ .
- 1302 • **Edge of Stability (EoS) for AGMs.** When  $\eta_{t_0} \approx \eta_*(\mathbf{w}_{t_0})$ , a small positive perturbation of  $\eta$   
 1303 causes the local quadratic approximation  $Q_{\mathcal{B}_t}(\mathbf{w})$  going to  $+\infty$ .

1304 Again, by an analogous argument as in Appendix H.1.3, we add a convexity compensation term on  
 1305  $Q_{\mathcal{B}_t}$  to avoid the impact of negative curvature (those  $\mu$  such that  $\mu + \lambda_{\text{WD}} < 0$ ) as  
 1306

$$1307 G_t(\boldsymbol{\theta}) := \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \boldsymbol{\theta}^\top \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})$$

$$1308 + \frac{1}{2} \boldsymbol{\theta}^\top \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) \boldsymbol{\theta} + \frac{\lambda_P}{2} \|\boldsymbol{\theta}\|_{\mathcal{P}}^2,$$

1310 In Appendix H.1, the convexity compensation term  $\frac{\lambda_P}{2} \|\boldsymbol{\theta}\|_{\mathcal{P}}^2$  shifts the local curvature of  $\mathbf{H}$  by  $\lambda_P \mathbf{I}$   
 1311 and here we aim to shift the spectrum of  $\mathbf{P}^{-1} \mathbf{H}$ , thus we replace the convexity compensation term  
 1312 with  $\frac{\lambda_P}{2} \|\boldsymbol{\theta}\|_{\mathcal{P}}^2$ . And we use this proxy for the detection of EoS for AGMs.  
 1313

1314 In the same fashion of Equation (14), the critical LR with compensation  $\eta_*^{\lambda_P}$  is defined here. Then we  
 1315 have the following theorem, suggesting that the impact of  $\lambda_P$  is bounded by the spectral property of  
 1316 the preconditioned Hessian  $\mathbf{P}^{-1} \mathbf{H}$ .

1317 **Theorem H.10** (EoS Detection Error). *Given the Batch covariance matrix  $\text{Cov}(\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})) \preceq \sigma^2 \mathbf{I}$ ,  
 1318 consider AGMs in the form of Equations (15) to (17) with  $\beta_1 = 0$ , it holds that*

$$1319 \frac{\eta_*^{\lambda_P}}{\eta_*^0} \in \left[ \frac{\mu_{\max}(\mu_{\max} + \lambda_P)}{(\mu_{\max} + \lambda_P)^2 + \sigma^2}, \frac{\mu_{\max}^2 + \sigma^2}{\mu_{\max}(\mu_{\max} + \lambda_P)} \right],$$

1322 where  $\mu_{\max}, \mu_{\min}$  represent the maximal and minimal eigenvalues of  $\mathbf{P}^{-1} \nabla^2 \mathcal{L}(\mathbf{w}_{t_0}) + \lambda_{\text{WD}} \mathbf{I}$ .  
 1323 Specifically, if  $\lambda_P = -\mu_{\min}$ , then we have

$$1324 \frac{\eta_*^{\lambda_P}}{\eta_*^0} \in \left[ \frac{\kappa^2 + \kappa}{(\kappa + 1)^2 + \alpha}, \frac{\kappa^2 + \alpha}{\kappa^2 + \kappa} \right]$$

1327 where  $\alpha = \frac{\sigma^2}{|\mu_{\min}|}$ ,  $\kappa = \mu_{\max}/|\mu_{\min}|$ .

1329 The proof of Theorem H.10 is given in Appendix K, where Theorem K.1 itself is a restatement of  
 1330 the above theorem. Following the approach in Appendix H.1.3, we give the test of EoC and EoS for  
 1331 AGMs in Algorithm 4 and Algorithm 5.  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

---

## I PROOF OF THEOREM H.1

*Proof.* We give the characterization of a more general system with an extra weight decay term  $\frac{\lambda_P}{2} \|\mathbf{w} - \mathbf{w}_0\|_2^2$ , which is introduced in Appendix H.1.3. One can see that the result of Theorem H.1 is a special case when  $\lambda_P = 0$ .

We consider optimizing the following proximal loss function:

$$\begin{aligned} f_t(\mathbf{w}) &:= \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0) + (\mathbf{w} - \mathbf{w}_0)^\top \mathbf{g} \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0) \\ &\quad + \frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^\top \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0) (\mathbf{w} - \mathbf{w}_0) \\ &\quad + \frac{\lambda_P}{2} \|\mathbf{w} - \mathbf{w}_0\|_2^2, \end{aligned}$$

where we term  $\lambda_P$  the compensation coefficient. For notational simplicity, we define:

- (Reparameterization)  $\boldsymbol{\theta} := \mathbf{w} - \mathbf{w}_0$ ,  $g_t(\boldsymbol{\theta}) := f_t(\mathbf{w})$ ;
- (Gradient)  $\mathbf{g}_t := \mathbf{g} \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0)$ ,  $\mathbf{g} := \mathbb{E}_{\mathcal{B}_t}[\nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0)]$ ;
- (Hessian)  $\mathbf{H}_t := \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0)$ ,  $\mathbf{H} := \mathbb{E}_{\mathcal{B}_t}[\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0)]$ ;
- (Reference Loss)  $\mathcal{L}_t := \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0)$ ,  $\mathcal{L} := \mathbb{E}_{\mathcal{B}_t}[\mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_0)]$ ;
- (Effective Hessian)  $\mathbf{A}_t := \mathbf{H}_t + \lambda_P \cdot \mathbf{I}$ ,  $\mathbf{A} := \mathbb{E}_{\mathcal{B}_t}[\mathbf{H}_t + \lambda_P \cdot \mathbf{I}]$ ;
- (Contraction Map)  $\mathbf{C}_t := \mathbf{I} - \eta \mathbf{A}_t$ ,  $\mathbf{C} := \mathbb{E}_{\mathcal{B}_t}[\mathbf{I} - \eta \mathbf{A}_t]$ .

Thus, the objective function can be rewritten as:

$$g_t(\boldsymbol{\theta}) = \mathcal{L}_t + \boldsymbol{\theta}^\top \mathbf{g}_t + \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{H}_t \boldsymbol{\theta} + \frac{\lambda_P}{2} \|\boldsymbol{\theta}\|_2^2.$$

We analyze the dynamics of Stochastic Gradient Descent (SGD) given by:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta(\mathbf{g}_t + \mathbf{A}_t \boldsymbol{\theta}_t). \quad (18)$$

In this notation, we restate Condition (10) as: we assume the existence of an optimal parameter vector  $\mathbf{x}_*$  such that:

$$\mathbb{E}[\mathbf{g}_t + \mathbf{A}_t \mathbf{x}_*] = \mathbf{g} + \mathbf{A} \mathbf{x}_* = 0.$$

We now examine the deviation of the current parameter from the optimal one, defined as:

$$\boldsymbol{\epsilon}_t := \boldsymbol{\theta}_t - \mathbf{x}_*.$$

The SGD recurrence relation can be expressed in terms of the error  $\boldsymbol{\epsilon}_t$  as:

$$\boldsymbol{\epsilon}_{t+1} = \mathbf{C}_t \cdot \boldsymbol{\epsilon}_t - \eta(\mathbf{g}_t + \mathbf{A}_t \mathbf{x}_*).$$

We are interested in the evolution of the first and second moments of the error  $\boldsymbol{\epsilon}_t$ . We introduce:

$$\begin{aligned} \mathbf{m}^{(1)} &:= \mathbb{E}[\boldsymbol{\epsilon}_t], \quad \mathbf{M}^{(2)} := \mathbb{E}[\boldsymbol{\epsilon}_t \boldsymbol{\epsilon}_t^\top], \\ \mathbf{m}^{(2)} &:= \text{vec}(\mathbf{M}). \end{aligned}$$

For the first moment, we derive:

$$\begin{aligned} \mathbf{m}_{t+1}^{(1)} &= \mathbb{E}[\boldsymbol{\epsilon}_{t+1}] \\ &= \mathbb{E}[\boldsymbol{\theta}_{t+1} - \mathbf{x}_*] \\ &= \mathbb{E}[\boldsymbol{\theta}_t - \eta(\mathbf{g}_t + \mathbf{A}_t \boldsymbol{\theta}_t) - \mathbf{x}_*] \\ &= \mathbb{E}[\boldsymbol{\epsilon}_t - \eta \mathbf{A}_t \boldsymbol{\epsilon}_t - \eta(\mathbf{g}_t + \mathbf{A}_t \mathbf{x}_*)] \\ &= \mathbb{E}[\boldsymbol{\epsilon}_t - \eta \mathbf{A}_t \boldsymbol{\epsilon}_t] - \mathbb{E}[\eta(\mathbf{g}_t + \mathbf{A}_t \mathbf{x}_*)] \\ &= (\mathbf{I} - \eta \mathbf{A}) \mathbf{m}_t^{(1)} \\ &= \mathbf{C} \cdot \mathbf{m}_t^{(1)}. \end{aligned}$$

1404 Next, we derive the update rule for the second moment. Let  $\mathbf{r}_t := \mathbf{g}_t + \mathbf{A}_t \mathbf{x}_*$ . Observe that:

$$\begin{aligned}
1405 & \\
1406 & \epsilon_{t+1} \epsilon_{t+1}^\top = [\mathbf{C}_t \epsilon_t - \eta \mathbf{r}_t] [\mathbf{C}_t \epsilon_t - \eta \mathbf{r}_t]^\top \\
1407 & = \mathbf{C}_t \epsilon_t \epsilon_t^\top \mathbf{C}_t - \eta \mathbf{r}_t \epsilon_t^\top \mathbf{C}_t \\
1408 & \quad - \eta \mathbf{C}_t \epsilon_t \mathbf{r}_t^\top + \eta^2 \mathbf{r}_t \mathbf{r}_t^\top.
\end{aligned}$$

1410 Recall that  $\text{vec}(\mathbf{A} \mathbf{X} \mathbf{m}) = (\mathbf{m}^\top \otimes \mathbf{A}) \cdot \text{vec}(\mathbf{X})$ , where  $\otimes$  denotes the Kronecker product. Additionally,  
1411 we define:

$$\begin{aligned}
1412 & \mathbf{b} := \text{vec}\left(\mathbb{E}[\mathbf{r}_t \mathbf{r}_t^\top]\right), \\
1413 & \mathbf{S} := \mathbb{E}[\mathbf{C}_t \otimes \mathbf{r}_t + \mathbf{r}_t \otimes \mathbf{C}_t].
\end{aligned}$$

1416 Consequently, the update for the vectorized second moment is:

$$\begin{aligned}
1417 & \\
1418 & \mathbf{m}_{t+1}^{(2)} = \text{vec}(\mathbf{M}_{t+1}) \\
1419 & = \text{vec}(\mathbb{E}[\epsilon_{t+1} \epsilon_{t+1}^\top]) \\
1420 & = \mathbb{E}[\mathbf{C}_t \otimes \mathbf{C}_t] \cdot \mathbf{m}_t^{(2)} + \eta \mathbf{S} \cdot \mathbf{m}_t^{(1)} + \eta^2 \mathbf{b}.
\end{aligned}$$

1422 The first and second moments evolve according to the following linear dynamical system:

$$\begin{aligned}
1423 & \\
1424 & \begin{bmatrix} \mathbf{m}_{t+1}^{(1)} \\ \mathbf{m}_{t+1}^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ \eta \mathbf{S} & \mathbb{E}[\mathbf{C}_t \otimes \mathbf{C}_t] \end{bmatrix} \begin{bmatrix} \mathbf{m}_t^{(1)} \\ \mathbf{m}_t^{(2)} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \eta^2 \mathbf{b} \end{bmatrix}. \quad (19)
\end{aligned}$$

1427 Equation (19) completes the proof.  $\square$

## 1429 J PROOF OF COROLLARY H.4 AND COROLLARY H.2

1431 The stability of the linear system (19) is governed by the spectral radius (maximal absolute value  
1432 of eigenvalues) of the transition matrix. Since the matrix is block lower triangular, we only need to  
1433 consider the maximal and minimal eigenvalues of the diagonal blocks  $\mathbf{C}$  and  $\mathbb{E}[\mathbf{C}_t \otimes \mathbf{C}_t]$ . In other  
1434 words, the stability is equivalent to the following condition:

$$\max\{\rho(\mathbf{C}), \rho(\mathbb{E}[\mathbf{C}_t \otimes \mathbf{C}_t])\} \leq 1,$$

1437 where the operator  $\rho(\cdot)$  denotes the spectrum norm when applied to a matrix.

1438 For  $\mathbf{C}$ , we have  $\mathbf{C} = \mathbf{I} - \eta \mathbf{A} = \mathbf{I} - \eta(\mathbf{H} + \lambda_P \mathbf{I})$ . Assuming the eigenvalues of  $\mathbf{H}$  are  $\lambda_1, \lambda_2, \dots, \lambda_d$ ,  
1439 the maximal eigenvalue of  $\mathbf{C}$  is denoted by  $c_{\max} = \max_i(1 - \eta\lambda_i - \eta\lambda_P)$ , and the minimal eigenvalue  
1440 is denoted by  $c_{\min} = \min_i(1 - \eta\lambda_i - \eta\lambda_P)$ . By the condition  $\rho(\mathbf{C}) \leq 1$  We have

$$1441 \quad c_{\max} \leq 1, \quad c_{\min} \geq -1,$$

1443 which gives

$$1444 \quad \eta \leq \frac{2}{\lambda_{\max}(\mathbf{A}) + \lambda_P}, \quad \lambda_i + \lambda_P \geq 0.$$

1448 For  $\mathbb{E}[\mathbf{C}_t \otimes \mathbf{C}_t]$ , we have:

$$\begin{aligned}
1449 & \mathbb{E}[\mathbf{C}_t \otimes \mathbf{C}_t] = \mathbb{E}[(\mathbf{I} - \eta \mathbf{A}_t) \otimes (\mathbf{I} - \eta \mathbf{A}_t)] \\
1450 & = \mathbf{I} - \eta \mathbf{A} \otimes \mathbf{I} - \eta \mathbf{I} \otimes \mathbf{A} + \mathbb{E}[\eta^2 \mathbf{A}_t \otimes \mathbf{A}_t] \\
1451 & = (\mathbf{I} - \eta \mathbf{A}) \otimes (\mathbf{I} - \eta \mathbf{A}) + \eta^2 \text{Cov}(\mathbf{A}_t).
\end{aligned}$$

1454 We observe that  $\text{Cov}(\mathbf{A}_t)$  is positive semi-definite (PSD). The spectral radius of the first term is  
1455  $\max_i(1 - \eta\lambda_i - \eta\lambda_P)^2$ . The stability similarly gives two conditions:

$$\begin{aligned}
1456 & \lambda_i + \lambda_P \geq 0 \\
1457 & \eta \leq \eta_U,
\end{aligned}$$

where  $\eta_U$  is some constant dependent on  $\sigma, \lambda_{\max}(\mathbf{C})$ . The second inequality for  $\eta$  comes from a crucial observation: when  $\eta \leq \max_{i \in [d]} \left\{ \frac{\lambda_i}{\lambda_i + \sigma^2 + \sigma} \right\}$ , where  $\sigma^2 := \rho(\text{Cov}(\mathbf{A}_t))$  then

$$\rho((\mathbf{I} - \eta \mathbf{A}) \otimes (\mathbf{I} - \eta \mathbf{A}) + \eta^2 \text{Cov}(\mathbf{A}_t)) \leq \max \{ (1 - \eta(\lambda_i + \lambda_P))^2 + \eta^2 \sigma^2 \}.$$

By the upper bound of  $\eta$ , one can directly verify that for any  $i \in [d]$ ,

$$1 - \eta(\lambda_i + \lambda_P))^2 + \eta^2 \sigma^2 \leq 1,$$

Thus we conclude that when  $\eta \leq \min_{i \in [d]} \left\{ \frac{2}{\lambda_{\max}(\mathbf{A}) + \lambda_P}, \max_{i \in [d]} \left\{ \frac{\lambda_i}{\lambda_i + \sigma^2 + \sigma} \right\} \right\}$ , the system is stable, thus there exists an supremum  $\eta_*$  for any  $\eta \leq \eta_*$ , the system is stable. Taking  $\lambda_P = 0$  completes the proof of the original Corollary H.4 and Corollary H.2. And actually, we have finished the stability analysis for the system with  $\lambda > 0$ , which is introduced in Appendix H.1.3.  $\square$

## K EOS DETECTION ERROR FOR AGMS

For the proof of Theorem H.8 and Theorem H.10, we provide a unified theorem covering both these two as shown below.

**Theorem K.1** (EoS Detection Error for RMSProp). *Given the Batch covariance matrix  $\text{Cov}(\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})) \preceq \sigma \mathbf{I}$ , consider AGMs in the form of Equations (15) to (17) with  $a = 0$ , it holds that*

$$\frac{\eta_*^{\lambda_P}}{\eta_*^0} \in \left[ \frac{\mu_{\max}(\mu_{\max} + \lambda_P)}{(\mu_{\max} + \lambda_P)^2 + \sigma^2}, \frac{\mu_{\max}^2 + \sigma^2}{\mu_{\max}(\mu_{\max} + \lambda_P)} \right],$$

where  $\mu_{\max}, \mu_{\min}$  represent the maximal and minimal eigenvalues of  $\mathbf{P}^{-1} \nabla^2 \mathcal{L}(\mathbf{w}_{t_0}) + \lambda_{\text{WD}} \mathbf{I}$ . Specifically, if  $\lambda_P = -\mu_{\min}$ , then we have

$$\frac{\eta_*^{\lambda_P}}{\eta_*^0} \in \left[ \frac{\kappa^2 + \kappa}{(\kappa + 1)^2 + \alpha}, \frac{\kappa^2 + \alpha^2}{\kappa^2 + \kappa} \right]$$

where  $\alpha = \frac{\sigma^2}{|\mu_{\min}|^2}$ ,  $\kappa = \mu_{\max}/|\mu_{\min}|$ .

This theorem is basically a replication of Theorem H.10. One can see that Theorem H.8 is a direct corollary of Theorem K.1 by setting  $\mathbf{P} = \mathbf{I}$ .

*Proof of Theorem K.1.* We first write out the update rule of RMSProp with a weight decay parameter  $\lambda_{\text{WD}}$ :

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{P}^{-1}(\mathbf{g}_t + \mathbf{A}_t \boldsymbol{\theta}_t) - \eta \lambda_{\text{WD}}(\boldsymbol{\theta}_t + \mathbf{w}_{t_0}).$$

Let  $\epsilon_t = \boldsymbol{\theta}_t - \mathbf{x}_*$ , then we have

$$\begin{bmatrix} \mathbf{m}_{t+1}^{(1)} \\ \mathbf{m}_{t+1}^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{C}, & \mathbf{0} \\ \eta \mathbf{S}, & \mathbb{E}[\mathbf{C}_t \otimes \mathbf{C}_t] \end{bmatrix} \begin{bmatrix} \mathbf{m}_t^{(1)} \\ \mathbf{m}_t^{(2)} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \eta^2 \mathbf{b} \end{bmatrix},$$

where  $\mathbf{C} = (\mathbf{I} - \eta(\mathbf{P}^{-1} \mathbf{H} + \lambda_P \mathbf{I} + \lambda_{\text{WD}} \mathbf{I}))$ . Recall that our goal is to get the bound for  $\frac{\eta_*^{\lambda_P}}{\eta_*^0}$ . The assumption  $\preceq \text{Cov}(\mathbf{H}_t) \preceq \sigma^2 \mathbf{I}$  gives the upper bound and lower bound for  $\eta_*^{\lambda_P}$ :

**Upper bound:**  $\eta_*^{\lambda_P} \leq \eta_U^{\lambda_P}$ , where  $\eta_U^{\lambda_P} := \sup \{ \eta : \rho(M(\eta; \lambda_P)) < 1 \}$ .

**Lower bound:**  $\eta_*^{\lambda_P} \geq \eta_L^{\lambda_P}$ , where  $\eta_L^{\lambda_P} := \sup \{ \eta : \|M(\eta; \lambda_P)\|^2 + \eta^2 \sigma^2 < 1 \}$ .

So we have

$$\eta_L^{\lambda_P} \leq \eta_*^{\lambda_P} \leq \eta_U^{\lambda_P}, \quad \eta_L^0 \leq \eta_*^0 \leq \eta_U^0.$$

And

$$\frac{\eta_L^{\lambda_P}}{\eta_U^0} \leq \frac{\eta_*^{\lambda_P}}{\eta_*^0} \leq \frac{\eta_U^{\lambda_P}}{\eta_L^0}.$$

1512 We denote the eigenvalues of  $\mathbf{P}^{-1}\mathbf{H} + \lambda_{\text{WD}}\mathbf{I}$  by  $\mu_1, \mu_2, \dots, \mu_d$ ,

$$1513 \eta_U^{\lambda_{\text{P}}} = \frac{2}{\mu_{\max} + \lambda_{\text{P}}}$$

1514 and

$$1515 \eta_L^{\lambda_{\text{P}}} = \frac{2(\mu_{\max} + \lambda_{\text{P}})}{(\mu_{\max} + \lambda_{\text{P}})^2 + \sigma_2^2}$$

$$1516 \frac{\eta_L^{\lambda_{\text{P}}}}{\eta_U^0} = \frac{\mu_{\max}(\mu_{\max} + \lambda_{\text{P}})}{(\mu_{\max} + \lambda_{\text{P}})^2 + \sigma^2}$$

1517 Taking  $\lambda_{\text{P}} = |\mu_{\min}|$  gives

$$1518 \frac{\eta_L^{\lambda_{\text{P}}}}{\eta_U^0} = \frac{\kappa^2 + \kappa}{(\kappa + 1)^2 + \alpha^2},$$

1519 where  $\kappa = \frac{\mu_{\max}}{|\mu_{\min}|}$  and  $\alpha = \frac{\sigma}{|\mu_{\min}|}$ .

1520 Similarly we have

$$1521 \frac{\eta_U^{\lambda_{\text{P}}}}{\eta_L^0} = \frac{\mu_{\max}^2 + \sigma^2}{\mu_{\max}(\mu_{\max} + \lambda_{\text{P}})}.$$

1522 Taking  $\lambda_{\text{P}} = |\mu_{\min}|$  gives

$$1523 \frac{\eta_U^{\lambda_{\text{P}}}}{\eta_L^0} = \frac{\kappa^2 + \alpha^2}{\kappa^2 + \kappa}.$$

1524 Thus we have  $\frac{\eta_U^{\lambda_{\text{P}}}}{\eta_L^0} = 1(1 + O(\frac{1}{\kappa}))$ . □

## 1525 L LINEAR SYSTEM DERIVATION FOR AGMS

1526 For the local simulation of AGMs, we write out the update rule as

$$1527 \mathbf{m}_{t+1} = a \cdot \mathbf{m}_t + (1 - a) \cdot (\mathbf{g}_t + \mathbf{A}_t \boldsymbol{\theta}_t)$$

$$1528 \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{P}^{-1} \mathbf{m}_{t+1} - \eta \lambda_{\text{WD}} (\boldsymbol{\theta}_t + \mathbf{w}_{t_0}).$$

1529 To simplify notations, we denote  $a = \beta_1$  and  $b = 1 - \beta_1$ . Let  $\boldsymbol{\epsilon}_t = \boldsymbol{\theta}_t - \mathbf{x}_*$  and organizing the above equations give

$$1530 \eta \mathbf{P}^{-1} \mathbf{m}_{t+1} = a \eta \mathbf{P}^{-1} \cdot \mathbf{m}_t + b \eta \mathbf{P}^{-1} \mathbf{A}_t \boldsymbol{\epsilon}_t + b \eta \mathbf{P}^{-1} (\mathbf{g}_t + \mathbf{A}_t \mathbf{x}_*).$$

$$1531 \boldsymbol{\epsilon}_{t+1} = (\mathbf{I} - \eta b \mathbf{P}^{-1} \mathbf{A}_t - \eta \lambda_{\text{WD}}) \boldsymbol{\epsilon}_t - \eta a \mathbf{P}^{-1} \mathbf{m}_t$$

$$1532 - \eta (b \mathbf{g}_t + b \mathbf{A}_t \mathbf{x}_* + \lambda_{\text{WD}} \mathbf{w}_{t_0} + \lambda_{\text{WD}} \mathbf{x}_*).$$

1533 Define  $\mathbf{D}_t = \mathbf{P}^{-1} \mathbf{H}_t$  and let

$$1534 \mathbf{M}_t = \begin{bmatrix} a \mathbf{I} & b \eta \mathbf{D}_t \\ -a \mathbf{P}^{-1} & \mathbf{I} - \eta b \mathbf{D}_t - \eta \lambda_{\text{WD}} \mathbf{I} \end{bmatrix},$$

1535 and

$$1536 \mathbf{q}_t = \begin{bmatrix} b \eta \mathbf{P}^{-1} (\mathbf{g}_t + \mathbf{A}_t \mathbf{x}_*) \\ -\eta (b \mathbf{g}_t + b \mathbf{A}_t \mathbf{x}_* + \lambda_{\text{WD}} \mathbf{w}_{t_0} + \lambda_{\text{WD}} \mathbf{x}_*). \end{bmatrix}$$

1537 Let  $\mathbf{z}_t = ((\eta \mathbf{P}^{-1} \mathbf{m}_t)^\top, \boldsymbol{\epsilon}_t^\top)^\top$ , then we have the linear system

$$1538 \mathbf{z}_{t+1} = \mathbf{M}_t \mathbf{z}_t + \mathbf{q}_t.$$

1539 We denote  $\mathbf{m}_t^{(1)} = \mathbb{E}[\mathbf{z}_t]$ , thus we have

$$1540 \mathbf{m}_{t+1}^{(1)} = \mathbf{M} \mathbf{m}_t^{(1)},$$

1566 where

$$1567 \quad M = \begin{bmatrix} \beta_1 \mathbf{I} & b\eta \mathbf{D} \\ -a\mathbf{P}^{-1} & \mathbf{I} - \eta b \mathbf{D} - \eta \lambda_{\text{WD}} \mathbf{I} \end{bmatrix}$$

1569 We define the second moment matrix  $\Sigma_t := \mathbb{E}[\mathbf{z}_t \mathbf{z}_t^\top]$  and the second moment vector  $\mathbf{m}_t^{(2)} = \text{vec}(\Sigma_t)$ . Then we write out the second moment evolution

$$1572 \quad \mathbf{m}_{t+1}^{(2)} = \mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t] \mathbf{m}_t^{(2)} + \mathbf{S} \mathbf{m}_t^{(1)} + \mathbf{b},$$

1573 where

$$1574 \quad \mathbf{S} := \mathbb{E}[\mathbf{q}_t \otimes \mathbf{M}_t + \mathbf{M}_t \otimes \mathbf{q}_t], \quad \mathbf{b} := \text{vec}(\mathbb{E}[\mathbf{q}_t \mathbf{q}_t^\top]).$$

1576 So the evolution for this linear system can be written as

$$1578 \quad \begin{bmatrix} \mathbf{m}_{t+1}^{(1)} \\ \mathbf{m}_{t+1}^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{M} & 0 \\ \mathbf{S} & \mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t] \end{bmatrix} \begin{bmatrix} \mathbf{m}_t^{(1)} \\ \mathbf{m}_t^{(2)} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix}.$$

1581 We denote  $\rho(\mathbf{M})$  the spectral radius of  $\mathbf{M}$ . Since the transition matrix is block lower triangular, we have

$$1583 \quad \rho \left( \begin{bmatrix} \mathbf{M} & 0 \\ \mathbf{S} & \mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t] \end{bmatrix} \right) = \max\{\rho(\mathbf{M}), \rho(\mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t])\}.$$

1585 In fact, we know that

$$1586 \quad \mathbb{E}[\mathbf{M}_t \otimes \mathbf{M}_t] = \mathbf{M} \otimes \mathbf{M} + \text{Cov}(\mathbf{M}_t).$$

## 1588 M MINIMAL EIGENVALUE ESTIMATION WITH STOCHASTIC POWER 1589 ITERATION

### 1591 M.1 NOTATION AND PROBLEM SETUPS

1593 Before we start proving Theorem H.6, we first recall some necessary notations and clarify our problem setups.

1595 Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$  be symmetric with eigenvalues  $\lambda_1(\mathbf{A}) \geq \dots \geq \lambda_d(\mathbf{A}) =: \lambda_{\min}(\mathbf{A})$ . Define the compensation threshold

$$1598 \quad \lambda_\star := \inf\{\lambda \geq 0 : \mathbf{A} + \lambda \mathbf{I} \succeq \mathbf{0}\} = \max\{0, -\lambda_{\min}(\mathbf{A})\}.$$

1599 For each step  $t$ , we observe symmetric matrices

$$1600 \quad \mathbf{A}_t = \mathbf{A} + \Delta_t, \quad \mathbb{E}[\Delta_t \mid \mathcal{F}_{t-1}] = \mathbf{0}, \quad \|\Delta_t\| \leq \sigma, \quad \text{a.s.} \quad (20)$$

1602 where filtration  $\mathcal{F}_t = \sigma(\Delta_0, \dots, \Delta_{t-1})$ . Given a candidate compensation coefficient  $\lambda \in [0, \Lambda]$  and step size  $\eta = \left[ \omega(\frac{1}{T}), \frac{1}{\sqrt{T}} \right]$ , define

$$1604 \quad \mathbf{Y}_t(\lambda) := \mathbf{I} - \eta(\mathbf{A}_t + \lambda \mathbf{I}), \quad \mathbf{M}(\lambda) := \mathbb{E} \mathbf{Y}_t(\lambda) = \mathbf{I} - \eta(\mathbf{A} + \lambda \mathbf{I}),$$

1606 and the product

$$1607 \quad \mathbf{Z}_t(\lambda) := \mathbf{Y}_{t-1}(\lambda) \cdots \mathbf{Y}_0(\lambda), \quad \mathbf{Z}_0(\lambda) = \mathbf{I}.$$

1608 With initialization  $\mathbf{x}_0 \neq \mathbf{0}$ , the iterates satisfy  $\mathbf{x}_t(\lambda) = \mathbf{Z}_t(\lambda) \mathbf{x}_0$ . Also, one important property implied by the range of  $\eta$  is

$$1610 \quad \eta(\lambda_{\max}(\mathbf{A}) + \Lambda) \leq 1 \quad \implies \quad \sup_{\lambda \in [\lambda_\star, \Lambda]} \|\mathbf{M}(\lambda)\| \leq 1. \quad (21)$$

1612 And this property is frequently used in our following proof.

1613 **Explosion test and  $N$ -round bisection.** Fix a threshold  $B > 1$ . For a given  $\lambda$ , we run  $T$  steps and declare

$$1615 \quad \text{Explode}(\lambda) = 1 \iff \max_{0 \leq t \leq T} \|\mathbf{x}_t(\lambda)\|_2 > B \|\mathbf{x}_0\|_2.$$

1617 We run  $N$  rounds of bisection over  $[0, \Lambda]$  using this predicate and output  $\hat{\lambda}$  as the final upper endpoint.

1618 We assume each bisection query uses an independent fresh length- $T$  trajectory  $\{\mathbf{A}_t\}_{t=0}^{T-1}$ .

1619 Now we are ready to give the following complete statement of Theorem H.6.

---

1620 M.2 RESTATEMENT OF THEOREM H.6

1621  
 1622 **Theorem M.1** (High-probability accuracy of  $\hat{\lambda}$ ). *Given LR  $\eta = \left[\omega\left(\frac{1}{T}\right), o\left(\frac{1}{\sqrt{T}}\right)\right]$ . Run  $N$  rounds*  
 1623 *of bisection on  $[0, \Lambda]$  with threshold  $B \geq 2$ , where each query  $\lambda$  uses an independent length- $T$*   
 1624 *trajectory of samples. Given initialization  $\mathbf{x}_0$  such that there exists some absolute constant  $\varepsilon > 0$ ,*  
 1625

$$1626 \frac{|\mathbf{x}_0^\top \mathbf{u}_{\min}|}{\|\mathbf{x}_0\|} \geq \varepsilon > 0,$$

1627  
 1628 *where  $\mathbf{u}_{\min}$  is the corresponding eigenvector of  $\mathbf{A}$ 's minimal eigenvalue  $\lambda_{\min}$ . Given  $\lambda_{\min}(A) < 0$ ,*  
 1629 *then with probability at least  $1 - \delta$ ,*  
 1630

$$1631 \left| \hat{\lambda} - \lambda_* \right| \leq \tilde{\mathcal{O}}\left(\frac{1}{\eta T}\right) + \tilde{\mathcal{O}}\left(\frac{\sigma\sqrt{\log d}}{\sqrt{T}}\right) + \frac{\Lambda}{2^N}, \quad (22)$$

1632  
 1633 *In particular,*

$$1634 \left| \hat{\lambda} - \lambda_* \right| = \tilde{\mathcal{O}}\left(\frac{1}{\eta T}\right) + \tilde{\mathcal{O}}\left(\frac{\sigma\sqrt{\log d}}{\sqrt{T}}\right) + \mathcal{O}(2^{-N}). \quad (\text{w.h.p.})$$

1635  
 1636 To prove the above theorem, we consider two scenarios: (1)  $\lambda \geq \lambda_*$ , (2)  $\lambda < \lambda_*$ . We give the  
 1637 following two lemmas, which imply Theorem M.1.

1638  
 1639 **Lemma M.2** (stable when  $\lambda \geq \lambda_*$ ). *Assume  $\eta = \left[\omega\left(\frac{1}{T}\right), o\left(\frac{1}{\sqrt{T}}\right)\right]$  and take  $B \geq 2$ . Fix a query*  
 1640  *$\lambda \in [\lambda_*, \Lambda]$ . Then, with probability at least  $1 - \delta/N$ ,*

$$1641 \max_{0 \leq t \leq T} \|\mathbf{x}_t(\lambda)\|_2 \leq B \|\mathbf{x}_0\|_2.$$

1642  
 1643 **Lemma M.3** (Explosion when  $\lambda < \lambda_*$ ). *Assume  $\eta = \left[\omega\left(\frac{1}{T}\right), o\left(\frac{1}{\sqrt{T}}\right)\right]$  and take  $B \geq 2$ . Fix one*  
 1644 *query  $\lambda = \lambda_* - \epsilon$ ,  $\epsilon > 0$  with initialization  $\mathbf{x}_0$  such that there exists some positive constant  $\varepsilon > 0$*   
 1645 *independent of  $T, \eta$ ,*

$$1646 \frac{|\mathbf{x}_0^\top \mathbf{u}_{\min}|}{\|\mathbf{x}_0\|} \geq \varepsilon > 0,$$

1647  
 1648 *define  $\rho := \|\mathbf{M}(\lambda)\|$  and for  $\delta \in (0, 1)$  define*

$$1649 \Gamma_T(\delta) := C \eta \sigma \left( \sqrt{T \log \frac{4dTN}{\delta}} + \log \frac{4dTN}{\delta} \right). \quad (23)$$

1650  
 1651 *where  $C > 0$  is a sufficiently large absolute constant independent of  $T, \eta$ . Fix  $\delta \in (0, 1)$ . Then, with*  
 1652 *probability at least  $1 - \delta/N$ ,*

$$1653 \frac{\|\mathbf{x}_T(\lambda)\|_2}{\|\mathbf{x}_0\|_2} \geq \rho^T \varepsilon ((1 - 2\Gamma_T(\delta))).$$

1654  
 1655 *Consequently, if*

$$1656 \rho^T \geq \frac{B}{\varepsilon(1 - 2\Gamma_T(\delta))},$$

1657  
 1658 *then  $\text{Explode}(\lambda) = 1$ . In particular, since  $\rho \geq 1 + \eta\epsilon$ , it suffices that*

$$1659 (1 + \eta\epsilon)^T \geq \frac{B}{\varepsilon(1 - 2\Gamma_T(\delta))}.$$

1660  
 1661 *Equivalently, defining the (random-init) gray-zone width*

$$1662 \epsilon_T(\delta) := \frac{1}{\eta T} \left( \log \left( \frac{B}{\varepsilon} \right) + 2\Gamma_T(\delta) \right) = \mathcal{O}\left(\frac{1}{\eta T}\right) + \tilde{\mathcal{O}}\left(\frac{\sigma\sqrt{\log d}}{\sqrt{T}}\right), \quad (24)$$

1663  
 1664 *we have: if  $\epsilon \geq \epsilon_T(\delta)$  then  $\text{Explode}(\lambda) = 1$  with probability at least  $1 - \delta/N$ .*

1665  
 1666 Now, assuming the correctness of the above Lemmas, which we will prove later, we use these two  
 1667 Lemmas to prove Theorem M.1.

1674 *Proof of Theorem M.1.* By Lemma M.2, for any queried  $\lambda \geq \lambda_*$  the procedure is stable (no ex-  
 1675 plosion) with probability at least  $1 - \delta/N$ . By Lemma M.3, for any queried  $\lambda \leq \lambda_* - \varepsilon_T(\delta)$  the  
 1676 procedure explodes on the same high-probability event. Therefore, on an event of probability at  
 1677 least  $1 - \delta$  (union bound over  $N$  queries), the explode/stable predicate is correct outside the interval  
 1678  $[\lambda_* - \varepsilon_T(\delta), \lambda_* + \varepsilon_T(\delta)]$ , which directly implies the output  $\hat{\lambda}$  lies within  $\varepsilon_T(\delta)$  of  $\lambda_*$  up to the  
 1679 bisection resolution  $\Lambda/2^N$ , yielding Equation (22).  $\square$

1680

### 1681 M.3 PROOF OF LEMMA M.2

1682

1683 Fix  $B > 1$  and define the operator-norm stopping time

$$1684 \tau(\lambda) := \inf\{t \in \{0, 1, \dots, T\} : \|\mathbf{Z}_t(\lambda)\| > B\} \wedge T. \quad (25)$$

1685

1686 By definition,

$$1687 \|\mathbf{Z}_s(\lambda)\| \leq B \quad \text{for all integers } s \in \{0, 1, \dots, \tau(\lambda) - 1\}. \quad (26)$$

1688

1688 We then give an important lemma below.

1689

1689 **Lemma M.4** (Stopping-time matrix-Freedman concentration). *Assume (20) and (21) hold. Fix a*  
 1690 *query  $\lambda \in [0, \Lambda]$  and define  $\tau = \tau(\lambda)$  as in (25). Then for any  $\delta \in (0, 1)$ , with probability at least*  
 1691  *$1 - \delta/N$ ,*

1692

$$1693 \|\mathbf{Z}_t(\lambda) - \mathbf{M}(\lambda)^t\| \leq B \Gamma_T(\delta) \quad \text{for all } t \in \{0, 1, \dots, T\} \text{ such that } t \leq \tau(\lambda). \quad (27)$$

1694

1694 *Consequently, by a union bound over the  $N$  bisection queries (each using an independent trajectory),*  
 1695 *Equation (27) holds simultaneously for all  $N$  queries with probability at least  $1 - \delta$ .*

1696

1697 To prove the above lemma, we first give two technical lemmas

1698

1698 **Lemma M.5** (Self-adjoint dilation). *For any (real) matrix  $X \in \mathbb{R}^{d \times d}$  define its self-adjoint dilation*

1699

$$1700 \mathcal{D}(X) := \begin{pmatrix} 0 & X \\ X^\top & 0 \end{pmatrix} \in \mathbb{R}^{2d \times 2d}.$$

1701

1702 *Then  $\mathcal{D}(X)$  is symmetric and*

1703

$$1704 \|\mathcal{D}(X)\| = \|X\|, \quad \mathcal{D}(X)^2 = \begin{pmatrix} XX^\top & 0 \\ 0 & X^\top X \end{pmatrix}.$$

1705

1706 **Theorem M.6** (Matrix Freedman (Theorem 1.2 in Tropp (2011))). *Let  $\{\mathbf{S}_k\}_{k \geq 0}$  be a self-adjoint*  
 1707 *matrix martingale of dimension  $m$  with difference sequence  $\mathbf{X}_k := \mathbf{S}_k - \mathbf{S}_{k-1}$  adapted to a filtration*  
 1708  *$\{\mathcal{F}_k\}$ . Assume  $\mathbb{E}[\mathbf{X}_k | \mathcal{F}_{k-1}] = 0$  and  $\lambda_{\max}(\mathbf{X}_k) \leq R$  almost surely for all  $k$ . Define the predictable*  
 1709 *quadratic variation*

1710

$$1711 V_n := \sum_{k=1}^n \mathbb{E}[\mathbf{X}_k^2 | \mathcal{F}_{k-1}].$$

1712

1713 *Then for all  $u \geq 0$ ,*

1714

$$1715 \mathbb{P}\left\{\lambda_{\max}(\mathbf{S}_n - \mathbf{S}_0) \geq u \text{ and } \|V_n\| \leq v\right\} \leq m \cdot \exp\left(-\frac{u^2}{2(v + Ru/3)}\right).$$

1716

1717 *In particular, there exists an absolute constant  $C_0 > 0$  such that for any  $\rho \in (0, 1)$ , with probability*  
 1718 *at least  $1 - \rho$ ,*

1719

$$1720 \|\mathbf{S}_n - \mathbf{S}_0\| \leq C_0 \left( \sqrt{v \log \frac{m}{\rho}} + R \log \frac{m}{\rho} \right). \quad (28)$$

1721

1722 *Proof of Lemma M.4.* Fix one query  $\lambda \in [0, \Lambda]$  and abbreviate  $\mathbf{Y}_t := \mathbf{Y}_t(\lambda)$ ,  $\mathbf{M} := \mathbf{M}(\lambda)$ ,  $\mathbf{Z}_t :=$   
 1723  $\mathbf{Z}_t(\lambda)$ ,  $\tau := \tau(\lambda)$ . Let  $\mathcal{F}_s := \sigma(\Delta_0, \dots, \Delta_{s-1}) = \sigma(\mathbf{Y}_0, \dots, \mathbf{Y}_{s-1})$  be the natural filtration.

1724

1724 Fix an integer  $t \in \{1, \dots, T\}$ . Define, for  $s = 0, 1, \dots, t$ ,

1725

$$1726 \mathbf{H}_s^{(t)} := \mathbb{E}[\mathbf{Z}_t | \mathcal{F}_s]. \quad (29)$$

1727

1727 Then  $\{\mathbf{H}_s^{(t)}\}_{s=0}^t$  is a matrix martingale w.r.t.  $\{\mathcal{F}_s\}$  with  $H_0^{(t)} = \mathbb{E}\mathbf{Z}_t$  and  $H_t^{(t)} = \mathbf{Z}_t$ .

We claim that  $\mathbf{H}_s^{(t)}$  admits the explicit expression

$$\mathbf{H}_s^{(t)} = \mathbf{M}^{t-s} \mathbf{Z}_s, \quad s = 0, 1, \dots, t. \quad (30)$$

Indeed, by independence of the future factors  $\{\mathbf{Y}_s, \dots, \mathbf{Y}_{t-1}\}$  from  $\mathcal{F}_s$  and  $\mathbb{E}\mathbf{Y}_j = M$  for each  $j$ ,

$$\mathbb{E}[\mathbf{Z}_t | \mathcal{F}_s] = \mathbb{E}[\mathbf{Y}_{t-1} \cdots \mathbf{Y}_s] \mathbf{Z}_s = \mathbf{M}^{t-s} \mathbf{Z}_s,$$

proving (30). In particular,  $\mathbf{H}_0^{(t)} = \mathbf{M}^t$  and  $\mathbf{H}_t^{(t)} = \mathbf{Z}_t$ .

Define the stopped martingale

$$\widetilde{\mathbf{H}}_s^{(t)} := \mathbf{H}_{s \wedge \tau \wedge t}^{(t)}, \quad s = 0, 1, \dots, t,$$

and its differences  $\widetilde{\mathbf{D}}_s^{(t)} := \widetilde{\mathbf{H}}_s^{(t)} - \widetilde{\mathbf{H}}_{s-1}^{(t)}$  for  $s = 1, \dots, t$ . Then  $\{\widetilde{\mathbf{H}}_s^{(t)}\}$  is a martingale and  $\sum_{s=1}^t \widetilde{\mathbf{D}}_s^{(t)} = \widetilde{\mathbf{H}}_t^{(t)} - \widetilde{\mathbf{H}}_0^{(t)}$ .

Using (30) and  $\mathbf{Z}_s = \mathbf{Y}_{s-1} \mathbf{Z}_{s-1}$ , the (un-stopped) increment is

$$\mathbf{H}_s^{(t)} - \mathbf{H}_{s-1}^{(t)} = \mathbf{M}^{t-s} \mathbf{Z}_s - \mathbf{M}^{t-(s-1)} \mathbf{Z}_{s-1} = \mathbf{M}^{t-s} (\mathbf{Y}_{s-1} - M) \mathbf{Z}_{s-1}. \quad (31)$$

From (20) we have

$$\mathbf{Y}_{s-1} - M = -\eta \Delta_{s-1}, \quad \|\mathbf{Y}_{s-1} - M\| \leq \eta \sigma \quad \text{a.s.} \quad (32)$$

Moreover, by (21),  $\|M\| \leq 1$  and hence  $\|\mathbf{M}^{t-s}\| \leq 1$ . Finally, by the stopping rule (25)–(26), whenever  $s-1 < \tau$  we have  $\|\mathbf{Z}_{s-1}\| \leq B$ . Therefore, combining (31)–(32),

$$\|\mathbf{H}_s^{(t)} - \mathbf{H}_{s-1}^{(t)}\| \leq \|\mathbf{M}^{t-s}\| \|\mathbf{Y}_{s-1} - M\| \|\mathbf{Z}_{s-1}\| \leq \eta \sigma B \quad \text{on the event } \{s-1 < \tau\}.$$

On the complementary event  $\{s-1 \geq \tau\}$  the stopped increment  $\widetilde{\mathbf{D}}_s^{(t)}$  is identically zero. Hence, we obtain the uniform almost sure bound

$$\|\widetilde{\mathbf{D}}_s^{(t)}\| \leq R := \eta \sigma B \quad \text{for all } s = 1, \dots, t, \text{ a.s.} \quad (33)$$

We bound the conditional second moment. On the event  $\{s-1 < \tau\}$ . We use (31):

$$(\mathbf{H}_s^{(t)} - \mathbf{H}_{s-1}^{(t)})(\mathbf{H}_s^{(t)} - \mathbf{H}_{s-1}^{(t)})^\top = \mathbf{M}^{t-s} (\mathbf{Y}_{s-1} - M) \mathbf{Z}_{s-1} \mathbf{Z}_{s-1}^\top (\mathbf{Y}_{s-1} - M)^\top \mathbf{M}^{t-s \top}.$$

Conditioning on  $\mathcal{F}_{s-1}$  makes  $\mathbf{Z}_{s-1}$  measurable while  $(\mathbf{Y}_{s-1} - M) = -\eta \Delta_{s-1}$  is independent of  $\mathcal{F}_{s-1}$ . Using  $\mathbf{Z}_{s-1} \mathbf{Z}_{s-1}^\top \preceq \|\mathbf{Z}_{s-1}\|^2 I \preceq B^2 I$  on  $\{s-1 < \tau\}$  yields

$$\mathbb{E}\left[(\mathbf{H}_s^{(t)} - \mathbf{H}_{s-1}^{(t)})(\mathbf{H}_s^{(t)} - \mathbf{H}_{s-1}^{(t)})^\top \mid \mathcal{F}_{s-1}\right] \preceq B^2 \mathbf{M}^{t-s} \mathbb{E}[(\mathbf{Y}_{s-1} - M)(\mathbf{Y}_{s-1} - M)^\top] \mathbf{M}^{t-s \top}.$$

By (32),

$$\mathbb{E}[(\mathbf{Y}_{s-1} - M)(\mathbf{Y}_{s-1} - M)^\top] = \eta^2 \mathbb{E}[\Delta_{s-1}^2] \preceq \eta^2 \sigma^2 I,$$

where the last step uses  $\Delta_{s-1}^2 \preceq \|\Delta_{s-1}\|^2 I \preceq \sigma^2 I$  a.s. Using again  $\|\mathbf{M}^{t-s}\| \leq 1$ , we conclude that on  $\{s-1 < \tau\}$ ,

$$\left\| \mathbb{E}\left[(\widetilde{\mathbf{H}}_s^{(t)} - \widetilde{\mathbf{H}}_{s-1}^{(t)})(\widetilde{\mathbf{H}}_s^{(t)} - \widetilde{\mathbf{H}}_{s-1}^{(t)})^\top \mid \mathcal{F}_{s-1}\right] \right\| \leq \eta^2 \sigma^2 B^2.$$

On  $\{s-1 \geq \tau\}$  the stopped increment is 0, so the same bound holds trivially. Therefore,

$$\left\| \sum_{s=1}^t \mathbb{E}\left[\widetilde{\mathbf{D}}_s^{(t)} \widetilde{\mathbf{D}}_s^{(t) \top} \mid \mathcal{F}_{s-1}\right] \right\| \leq t \eta^2 \sigma^2 B^2. \quad (34)$$

An identical argument gives the same bound for  $\sum \mathbb{E}[\widetilde{\mathbf{D}}_s^{(t) \top} \widetilde{\mathbf{D}}_s^{(t)} \mid \mathcal{F}_{s-1}]$ .

The process  $\widetilde{\mathbf{H}}_s^{(t)}$  need not be self-adjoint, so we apply Theorem M.6 to the self-adjoint dilation  $\mathcal{D}(\widetilde{\mathbf{H}}_s^{(t)})$ . Let  $\widetilde{S}_s := \mathcal{D}(\widetilde{\mathbf{H}}_s^{(t)})$  and  $\widetilde{X}_s := \widetilde{S}_s - \widetilde{S}_{s-1} = \mathcal{D}(\widetilde{\mathbf{D}}_s^{(t)})$ . By Lemma M.5,  $\|\widetilde{X}_s\| = \|\widetilde{\mathbf{D}}_s^{(t)}\| \leq R$  and the predictable quadratic variation satisfies

$$\left\| \sum_{s=1}^t \mathbb{E}[\widetilde{X}_s^2 \mid \mathcal{F}_{s-1}] \right\| = \max \left\{ \left\| \sum_{s=1}^t \mathbb{E}[\widetilde{\mathbf{D}}_s^{(t)} \widetilde{\mathbf{D}}_s^{(t) \top} \mid \mathcal{F}_{s-1}] \right\|, \left\| \sum_{s=1}^t \mathbb{E}[\widetilde{\mathbf{D}}_s^{(t) \top} \widetilde{\mathbf{D}}_s^{(t)} \mid \mathcal{F}_{s-1}] \right\| \right\} \leq t \eta^2 \sigma^2 B^2.$$

Thus Theorem M.6 (in the inverted form (28)) with dimension  $m = 2d$ , variance proxy  $v = t\eta^2\sigma^2B^2$ , and increment bound  $R = \eta\sigma B$  yields: with probability at least  $1 - \rho$ ,

$$\begin{aligned} \|\widetilde{\mathbf{H}}_{t\wedge\tau}^{(t)} - \widetilde{\mathbf{H}}_0^{(t)}\| &= \|\widetilde{\mathbf{H}}_{t\wedge\tau}^{(t)} - \mathbf{H}_0^{(t)}\| \leq C_0 \left( \sqrt{t\eta^2\sigma^2B^2 \log \frac{2d}{\rho}} + \eta\sigma B \log \frac{2d}{\rho} \right) \\ &= C_0 \eta\sigma B \left( \sqrt{t \log \frac{2d}{\rho}} + \log \frac{2d}{\rho} \right). \end{aligned} \quad (35)$$

Recall  $H_0^{(t)} = M^t$ . Also, if  $\tau \geq t$ , then  $t \wedge \tau = t$  and  $H_t^{(t)} = \mathbf{Z}_t$ , hence  $\widetilde{\mathbf{H}}_{t\wedge\tau}^{(t)} = \mathbf{Z}_t$ . Therefore, on the event  $\{\tau \geq t\}$ ,

$$\|\widetilde{\mathbf{H}}_{t\wedge\tau}^{(t)} - \mathbf{H}_0^{(t)}\| = \|\mathbf{Z}_t - M^t\|.$$

Consequently,

$$\mathbb{P}\left(\tau \geq t \text{ and } \|\mathbf{Z}_t - M^t\| > u\right) \leq \mathbb{P}\left(\|\widetilde{\mathbf{H}}_{t\wedge\tau}^{(t)} - \mathbf{H}_0^{(t)}\| > u\right).$$

Combining this with (35) and choosing  $\rho = \delta/(TN)$  gives that with probability at least  $1 - \delta/(TN)$ ,

$$t \leq \tau \implies \|\mathbf{Z}_t - M^t\| \leq C_0 \eta\sigma B \left( \sqrt{t \log \frac{4dTN}{\delta}} + \log \frac{4dTN}{\delta} \right).$$

Since  $t \leq T$  implies  $\sqrt{t} \leq \sqrt{T}$ , we further obtain the uniform implication

$$t \leq \tau \implies \|\mathbf{Z}_t - M^t\| \leq C_0 \eta\sigma B \left( \sqrt{T \log \frac{4dTN}{\delta}} + \log \frac{4dTN}{\delta} \right). \quad (36)$$

Let  $\mathcal{E}_t$  be the event that (36) holds for this particular  $t$ . We have  $\mathbb{P}(\mathcal{E}_t) \geq 1 - \delta/(TN)$ . By a union bound over  $t = 1, \dots, T$ , with probability at least  $1 - \delta/N$ , (36) holds simultaneously for all  $t \in \{1, \dots, T\}$ . This is exactly (27) after enlarging  $C$  in the definition of  $\Gamma_T(\delta)$  in (23). (Also  $t = 0$  is trivial because  $\mathbf{Z}_0 = M^0 = I$ .)

Also, for each query, we apply the above argument with failure probability  $\delta/N$  (implemented by taking  $\rho = \delta/(TN)$  inside the proof). A final union bound over the  $N$  queries yields a simultaneous success event of probability at least  $1 - \delta$ .  $\square$

For fixed  $(d, N, \delta, \sigma)$ ,

$$\Gamma_T(\delta) = C\eta\sigma \left( \sqrt{T \log(4dTN/\delta)} + \log(4dTN/\delta) \right) \rightarrow 0 \quad \text{whenever } \eta = o(1/\sqrt{T}).$$

Hence, on the non-explosion segment ( $t \leq \tau(\lambda)$ ), the deviation  $\|\mathbf{Z}_t(\lambda) - M(\lambda)^t\|$  vanishes in high probability as  $T \rightarrow \infty$ . So that we can assume that  $\Gamma_T(\delta) \leq 1/4$  for the rest of our proof, which is implied naturally by the big-O notation.

*Proof of Lemma M.2.* On the event of Lemma M.4, apply Equation (27) at  $t = \tau(\lambda)$  (which satisfies  $t \leq \tau(\lambda)$  trivially):

$$\|\mathbf{Z}_\tau - M(\lambda)^\tau\| \leq B\Gamma_T(\delta).$$

Because  $\lambda \geq \lambda_*$  and Equation (21) holds,  $\|M(\lambda)^\tau\| \leq 1$ . Thus

$$\|\mathbf{Z}_\tau\| \leq \|M(\lambda)^\tau\| + \|\mathbf{Z}_\tau - M(\lambda)^\tau\| \leq 1 + B\Gamma_T(\delta) \leq 1 + \frac{B}{4} \leq B,$$

where we used  $\Gamma_T(\delta) \leq 1/4$  and  $B \geq 2$ . This contradicts the definition of  $\tau(\lambda)$  unless  $\tau(\lambda) = T$ . Finally,  $\|x_t\| \leq \|\mathbf{Z}_t\| \|\mathbf{x}_0\| \leq B\|\mathbf{x}_0\|$  for all  $t \leq T$ .  $\square$

1836 M.4 PROOF OF LEMMA M.3

1837

1838 Define the normalized matrices

1839

1840

$$\tilde{\mathbf{Y}}_t := \rho^{-1} \mathbf{Y}_t, \quad \tilde{\mathbf{M}} := \rho^{-1} \mathbf{M}, \quad \tilde{\mathbf{Z}}_t := \tilde{\mathbf{Y}}_{t-1} \cdots \tilde{\mathbf{Y}}_0 = \rho^{-t} \mathbf{Z}_t.$$

1841

1842

Then  $\|\tilde{\mathbf{M}}\| = 1$ . Moreover, by the condition  $\mathbf{A}_t = \mathbf{A} + \Delta_t$  with  $\|\Delta_t\| \leq \sigma$ ,

1843

1844

$$\|\tilde{\mathbf{Y}}_t - \tilde{\mathbf{M}}\| = \rho^{-1} \|\mathbf{Y}_t - \mathbf{M}\| = \rho^{-1} \eta \|\Delta_t\| \leq \frac{\eta \sigma}{\rho} \quad \text{a.s.}$$

1845

Introduce the (analysis) stopping time

1846

1847

$$\tilde{\tau} := \inf\{t \in \{0, 1, \dots, T\} : \|\tilde{\mathbf{Z}}_t\| > 2\} \wedge T.$$

1848

1849

1850

Now apply the same stopping-time matrix-Freedman argument as in Lemma M.4 to the normalized sequence  $\{\tilde{\mathbf{Y}}_t\}$  (whose mean operator norm is  $\leq 1$ ) with threshold 2. This yields the existence of an event  $\mathcal{E}$  with  $\mathbb{P}(\mathcal{E}) \geq 1 - \delta/N$  such that

1851

1852

1853

$$\|\tilde{\mathbf{Z}}_t - \tilde{\mathbf{M}}^t\| \leq 2 \cdot \frac{\Gamma_T(\delta)}{\rho} \quad \text{for all } t \leq \tilde{\tau}. \quad (37)$$

1854

1855

Given  $\Gamma_T(\delta) \leq \frac{1}{4}$ , and since  $\rho \geq 1$  we have  $2\Gamma_T(\delta)/\rho \leq \frac{1}{2}$ . Thus on  $\mathcal{E}$ ,

1856

1857

$$\|\tilde{\mathbf{Z}}_{\tilde{\tau}}\| \leq \|\tilde{\mathbf{M}}^{\tilde{\tau}}\| + \|\tilde{\mathbf{Z}}_{\tilde{\tau}} - \tilde{\mathbf{M}}^{\tilde{\tau}}\| \leq 1 + \frac{1}{2} < 2,$$

1858

1859

which contradicts the definition of  $\tilde{\tau}$  unless  $\tilde{\tau} = T$ . Hence on  $\mathcal{E}$ , (37) holds at  $t = T$ , and

1860

1861

$$\|\tilde{\mathbf{Z}}_T\| \geq \|\tilde{\mathbf{M}}^T\| - \|\tilde{\mathbf{Z}}_T - \tilde{\mathbf{M}}^T\| \geq 1 - \frac{2\Gamma_T(\delta)}{\rho} \geq 1 - 2\Gamma_T(\delta).$$

1862

Rescaling back gives

1863

1864

$$\|\mathbf{Z}_T\| = \rho^T \|\tilde{\mathbf{Z}}_T\| \geq \rho^T (1 - 2\Gamma_T(\delta)). \quad (38)$$

1865

1866

1867

Condition on  $\mathbf{Z}_t$ . Let  $u$  be a unit top right singular vector of  $\mathbf{Z}_t$ , so that  $\|\mathbf{Z}_t u\|_2 = \|\mathbf{Z}_t\|$ . Let  $s := x_0 / \|x_0\|_2$ . By  $x_0 \sim \mathcal{N}(0, I_d)$ , the direction  $s$  is uniform on  $\mathbb{S}^{d-1}$  and independent of  $\mathbf{Z}_t$  (and hence independent of  $u$ ). Then

1868

1869

1870

$$\frac{\|\mathbf{x}_T\|_2}{\|\mathbf{x}_0\|_2} = \|\mathbf{Z}_T s\|_2 \geq \|\mathbf{Z}_T\| |\mathbf{u}^\top s| \geq \rho^T \varepsilon ((1 - 2\Gamma_T(\delta)))$$

1871

The above inequality completes the proof.  $\square$

1872

1873

1874

**Extension with Gaussian Initialization.** Further we can extend the theorem with initialization  $\mathbf{x}_0 \sim \mathcal{N}(0, I_d)$ . We use a standard spherical-cap bound: for  $s \sim \text{Unif}(\mathbb{S}^{d-1})$  and any fixed unit  $u$ ,

1875

1876

1877

$$\mathbb{P}(|\mathbf{u}^\top s| \leq t) \leq 2t\sqrt{d}, \quad \forall t \in (0, 1).$$

1878

1879

Taking  $t = \delta_0 / (2\sqrt{d})$  yields

$$\mathbb{P}\left(|\mathbf{u}^\top s| \geq \frac{\delta_0}{2\sqrt{d}}\right) \geq 1 - \delta_0.$$

1880

1881

Combine this with (38): on the intersection of the event  $\mathcal{E}$  and the above anti-concentration event (probability at least  $1 - \delta/N - \delta_0$ ),

1882

1883

1884

$$\frac{\|\mathbf{x}_T\|_2}{\|\mathbf{x}_0\|_2} \geq \rho^T (1 - 2\Gamma_T(\delta)) \cdot \frac{\delta_0}{2\sqrt{d}}.$$

1885

1886

1887

Therefore, if  $\rho^T \geq 4B\sqrt{d}/\delta_0$ , then  $\|\mathbf{x}_T\|_2 > B\|\mathbf{x}_0\|_2$ , which implies  $\text{Explode}(\lambda) = 1$  (since the test checks  $\max_{0 \leq t \leq T} \|x_t\|_2$ ). The same argument derives a whp bound for Gaussian initialization.  $\square$

1888

1889

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

## N ALGORITHMS

---

### Algorithm 1 Minimal Eigenvalue Approximation via Stochastic Power Iteration + Bisection

---

**Require:** Stochastic matrix access via i.i.d. samples  $\mathbf{A}_t$ , step size  $\eta$ , inner horizon  $T$ , bisection rounds  $N$ , explosion threshold  $B > 1$ , initial search range  $[0, \Lambda]$  (with  $\Lambda$  large enough), and initialization  $\mathbf{x}_0 \sim \mathcal{N}(0, I)$ .

**Ensure:** Estimated shift  $\hat{\lambda}$  and minimal eigenvalue estimate  $\hat{\lambda}_{\min} = -\hat{\lambda}$ .

1: **Subroutine** EXPLODETEST( $\lambda$ ):

2:     Sample  $\mathbf{x}_0 \sim \mathcal{N}(0, I)$  and normalize; set  $\mathbf{x} \leftarrow \mathbf{x}_0$ .

3: **for**  $t = 0, 1, \dots, T - 1$  **do**

4:     Draw an independent sample  $\mathbf{A}_t$ .

5:      $\mathbf{x} \leftarrow \mathbf{x} - \eta(\mathbf{A}_t + \lambda I)\mathbf{x}$ .

6:     **if**  $\|\mathbf{x}\|_2 > B$  **then**

7:         **return** false {explode}

8:     **end if**

9: **end for**

10: **return** true {stable}

11: **Main procedure:**

12:  $\lambda_{\text{lo}} \leftarrow 0, \quad \lambda_{\text{hi}} \leftarrow \Lambda$ .

13: **while** EXPLODETEST( $\lambda_{\text{hi}}$ ) = false **do**

14:      $\lambda_{\text{hi}} \leftarrow 2\lambda_{\text{hi}}$  {find a stable upper bound}

15: **end while**

16: **for**  $k = 1, 2, \dots, N$  **do**

17:      $\lambda_{\text{mid}} \leftarrow (\lambda_{\text{lo}} + \lambda_{\text{hi}})/2$ .

18:     **if** EXPLODETEST( $\lambda_{\text{mid}}$ ) = true **then**

19:          $\lambda_{\text{hi}} \leftarrow \lambda_{\text{mid}}$  {stable}

20:     **else**

21:          $\lambda_{\text{lo}} \leftarrow \lambda_{\text{mid}}$  {explode}

22:     **end if**

23: **end for**

24:  $\hat{\lambda} \leftarrow \lambda_{\text{hi}}, \quad \hat{\lambda}_{\min} \leftarrow -\hat{\lambda}$ .

25: **return**  $\hat{\lambda}, \hat{\lambda}_{\min}$ .

---

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

---

### Algorithm 2 EoC Test for SGD

---

**Require:** Checkpoint  $\mathbf{w}_{t_0}$ ; weight decay  $\lambda_{\text{WD}} \geq 0$ ; tolerance  $\nu_{\text{tol}} > 0$ ; Algorithm 1 parameters  $(\eta_{\text{PI}}, T_{\text{PI}}, N_{\text{PI}}, B_{\text{PI}}, \Lambda_{\text{PI}})$ ; Lanczos iterations  $T_{\text{LZ}}$ .  
**Ensure:**  $\text{EOC} \in \{\text{true}, \text{false}\}$ ; estimates  $\hat{\lambda}_{\text{min}}$ ,  $\hat{\lambda}_{\text{max}}$ , and  $\hat{\nu}$ .  
1: Define symmetric stochastic samples  $\mathbf{A}_t := \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \lambda_{\text{WD}} \mathbf{I}$ .  
2:  $(\hat{\lambda}_{\text{shift}}, \hat{\lambda}_{\text{min}}) \leftarrow \text{Algorithm 1}(\mathbf{A}_t; \eta_{\text{PI}}, T_{\text{PI}}, N_{\text{PI}}, B_{\text{PI}}, \Lambda_{\text{PI}})$ .  
3:  $\hat{\lambda}_{\text{max}} \leftarrow \text{LANCZOSMAXEIG}(\mathbf{A}_t; T_{\text{LZ}})$ .  
4:  $\hat{\nu} \leftarrow \frac{\min\{0, -\hat{\lambda}_{\text{min}}\}}{\hat{\lambda}_{\text{max}}}$ .  
5:  $\text{EOC} \leftarrow (\hat{\lambda}_{\text{min}} < 0) \wedge (\hat{\nu} < \nu_{\text{tol}})$ .  
6: **return**  $\text{EOC}, \hat{\lambda}_{\text{min}}, \hat{\lambda}_{\text{max}}, \hat{\nu}$ .

---

### Algorithm 3 EoS Test for SGD

---

**Require:** Checkpoint  $\mathbf{w}_{t_0}$ ; training LR  $\eta$ ; horizon  $T$ ; explosion threshold  $M$ ; seed  $s$ ; increasing scaling factors  $\{\alpha_k\}_{k=1}^K$ ; EoS tolerance  $\epsilon_{\text{eos}} > 0$ . Also require Algorithm 1 parameters  $(\eta_{\text{PI}}, T_{\text{PI}}, N_{\text{PI}}, B_{\text{PI}}, \Lambda_{\text{PI}})$ .  
**Ensure:** Estimated critical LR  $\hat{\eta}_{\lambda_P}^*(\mathbf{w}_{t_0})$ , ratio  $r(\mathbf{w}_{t_0})$ , and EOS flag.  
1: (**Step 1: estimate convexity compensation.**) Define stochastic samples  $\mathbf{A}_t := \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})$  via Hessian–vector products. Run Algorithm 1 on  $\mathbf{A}_t$  to obtain  $\hat{\lambda}_{\text{shift}}$ . Set  $\lambda_P \leftarrow \hat{\lambda}_{\text{shift}}$ .  
2: **Subroutine** EXPLODETEST( $\alpha$ ):  
3:     Set RNG seed to  $s$ ;  $\boldsymbol{\theta} \leftarrow 0$ ;  $R \leftarrow 0$ .  
4: **for**  $t = 0, 1, \dots, T - 1$  **do**  
5:     Sample mini-batch  $\mathcal{B}_t$ .  
6:      $\mathbf{g} \leftarrow \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) \boldsymbol{\theta} + \lambda_P \boldsymbol{\theta}$ .  
7:      $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \eta \mathbf{g}$ .  
8:      $R \leftarrow \max\{R, \|\boldsymbol{\theta}\|_2\}$ .  
9:     **if**  $R \geq M$  **then**  
10:         **return** true {explode}  
11:     **end if**  
12: **end for**  
13:     **return** false {stable}  
14: (**Step 2: search over**  $\{\alpha_k\}$ .) Find  $k^* := \min\{k : \text{EXPLODETEST}(\alpha_k) = \text{true}\}$ .  
15: **if** no such  $k^*$  exists **then**  
16:      $\hat{\eta}_{\lambda_P}^* \leftarrow \alpha_K \eta$ ;  $r \leftarrow 1/\alpha_K$ ; EOS  $\leftarrow$  false.  
17:     **return**  $\hat{\eta}_{\lambda_P}^*, r, \text{EOS}$ .  
18: **end if**  
19: Set  $\alpha_{\text{lo}} \leftarrow (\alpha_{k^*-1}$  if  $k^* > 1$  else 0) and  $\alpha_{\text{hi}} \leftarrow \alpha_{k^*}$ .  
20:  $\hat{\eta}_{\lambda_P}^*(\mathbf{w}_{t_0}) \leftarrow \alpha_{\text{hi}} \eta$ .  
21:  $r(\mathbf{w}_{t_0}) \leftarrow 1/\alpha_{\text{hi}}$ .  
22: EOS  $\leftarrow (r(\mathbf{w}_{t_0}) \in [1 - \epsilon_{\text{eos}}, 1 + \epsilon_{\text{eos}}])$ .  
23: **return**  $\hat{\eta}_{\lambda_P}^*, r, \text{EOS}$ .

---

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

---

#### Algorithm 4 EoC Test for AGMs

---

**Require:** Checkpoint  $\mathbf{w}_{t_0}$ ; frozen preconditioner  $\mathbf{P}_{t_0} \succ 0$ ; weight decay  $\lambda_{\text{WD}} \geq 0$ ; tolerance  $\nu_{\text{tol}} > 0$ ; Algorithm 1 parameters  $(\eta_{\text{PI}}, T_{\text{PI}}, N_{\text{PI}}, B_{\text{PI}}, \Lambda_{\text{PI}})$ ; Lanczos iterations  $T_{\text{LZ}}$ .  
**Ensure:**  $\text{EOC} \in \{\text{true}, \text{false}\}$ ; estimates  $\hat{\mu}_{\min}$ ,  $\hat{\mu}_{\max}$ , and  $\hat{\nu}$ .

- 1: Define symmetric stochastic samples  $\tilde{\mathbf{A}}_t := \mathbf{P}_{t_0}^{-1/2} (\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})) \mathbf{P}_{t_0}^{-1/2} + \lambda_{\text{WD}} \mathbf{I}$ .
- 2:  $(\hat{\lambda}_{\text{shift}}, \hat{\mu}_{\min}) \leftarrow \text{Algorithm 1}(\tilde{\mathbf{A}}_t; \eta_{\text{PI}}, T_{\text{PI}}, N_{\text{PI}}, B_{\text{PI}}, \Lambda_{\text{PI}})$ .
- 3:  $\hat{\mu}_{\max} \leftarrow \text{LANCZOSMAXEIG}(\tilde{\mathbf{A}}_t; T_{\text{LZ}})$ .
- 4:  $\hat{\nu} \leftarrow \frac{\min\{0, -\hat{\mu}_{\min}\}}{\hat{\mu}_{\max}}$ .
- 5:  $\text{EOC} \leftarrow (\hat{\mu}_{\min} < 0) \wedge (\hat{\nu} < \nu_{\text{tol}})$ .
- 6: **return**  $\text{EOC}$ ,  $\hat{\mu}_{\min}$ ,  $\hat{\mu}_{\max}$ ,  $\hat{\nu}$ .

---

#### Algorithm 5 EoS Test for AGMs

---

**Require:** Checkpoint  $\mathbf{w}_{t_0}$ ; momentum state  $\mathbf{m}_{t_0}$ ; frozen preconditioner  $\mathbf{P}_{t_0} \succ 0$ ; weight decay  $\lambda_{\text{WD}} \geq 0$ ; training LR  $\eta$ ; momentum parameters  $\beta_1 \in [0, 1)$ ; horizon  $T$ ; explosion threshold  $M$ ; seed  $s$ ; increasing scaling factors  $\{\alpha_k\}_{k=1}^K$ ; EoS tolerance  $\epsilon_{\text{eos}} > 0$ . Also require Algorithm 1 parameters  $(\eta_{\text{PI}}, T_{\text{PI}}, N_{\text{PI}}, B_{\text{PI}}, \Lambda_{\text{PI}})$ .  
**Ensure:** Estimated critical LR  $\hat{\eta}_{\lambda_P}^*$  ( $\mathbf{w}_{t_0}$ ), ratio  $r(\mathbf{w}_{t_0})$ , and EOS flag.

- 1: **(Step 1: estimate convexity compensation.)** Define symmetric stochastic samples  $\tilde{\mathbf{A}}_t := \mathbf{P}_{t_0}^{-1/2} (\nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0})) \mathbf{P}_{t_0}^{-1/2} + \lambda_{\text{WD}} \mathbf{I}$ . Run Algorithm 1 on  $\tilde{\mathbf{A}}_t$  to obtain  $\hat{\lambda}_{\text{shift}}$ . Set  $\lambda_P \leftarrow \hat{\lambda}_{\text{shift}}$ .
- 2: **Subroutine** EXPLODETEST( $\alpha$ ):
- 3:     Set RNG seed to  $s$ ;  $\boldsymbol{\theta} \leftarrow 0$ ;  $m \leftarrow m_{t_0}$ ;  $R \leftarrow 0$ .
- 4: **for**  $t = 0, 1, \dots, T - 1$  **do**
- 5:     Sample mini-batch  $\mathcal{B}_t$ .
- 6:      $\mathbf{g} \leftarrow \nabla \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) + \nabla^2 \mathcal{L}_{\mathcal{B}_t}(\mathbf{w}_{t_0}) \boldsymbol{\theta} + \lambda_P \mathbf{P}_{t_0} \boldsymbol{\theta}$ .
- 7:      $\mathbf{m} \leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \mathbf{g}$ .
- 8:      $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \eta \mathbf{P}_{t_0}^{-1} \mathbf{m} - \alpha \eta \lambda_{\text{WD}} (\boldsymbol{\theta} + \mathbf{w}_{t_0})$ .
- 9:      $R \leftarrow \max\{R, \|\boldsymbol{\theta}\|_{\mathbf{P}_{t_0}}\}$ .
- 10:    **if**  $R \geq M$  **then**
- 11:     **return** `true` {explode}
- 12:    **end if**
- 13: **end for**
- 14:     **return** `false` {stable}
- 15: **(Step 2: search over  $\{\alpha_k\}$ .)** Find  $k^* := \min\{k : \text{EXPLODETEST}(\alpha_k) = \text{true}\}$ .
- 16: **if** no such  $k^*$  exists **then**
- 17:      $\hat{\eta}_{\lambda_P}^* \leftarrow \alpha_K \eta$ ;  $r \leftarrow 1/\alpha_K$ ;  $\text{EOS} \leftarrow \text{false}$ .
- 18:     **return**  $\hat{\eta}_{\lambda_P}^*$ ,  $r$ ,  $\text{EOS}$ .
- 19: **end if**
- 20: Set  $\alpha_{1_0} \leftarrow (\alpha_{k^* - 1}$  if  $k^* > 1$  else 0) and  $\alpha_{\text{hi}} \leftarrow \alpha_{k^*}$ .
- 21:  $\hat{\eta}_{\lambda_P}^*(\mathbf{w}_{t_0}) \leftarrow \alpha_{\text{hi}} \eta$ .
- 22:  $r(\mathbf{w}_{t_0}) \leftarrow 1/\alpha_{\text{hi}}$ .
- 23:  $\text{EOS} \leftarrow (r(\mathbf{w}_{t_0}) \in [1 - \epsilon_{\text{eos}}, 1 + \epsilon_{\text{eos}}])$ .
- 24: **return**  $\hat{\eta}_{\lambda_P}^*$ ,  $r$ ,  $\text{EOS}$ .

---