

ZERO-ORDER OPTIMIZATION AT THE EDGE OF STABILITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Zeroth-order (ZO) methods are widely used when gradients are unavailable or prohibitively expensive, including black-box learning and memory-efficient fine-tuning of large models, yet their optimization dynamics in deep learning remain underexplored. In this work, we provide an explicit step size condition that exactly captures the (mean-square) linear stability of a family of ZO methods based on the standard two-point estimator. Our characterization reveals a sharp contrast with first-order (FO) methods: whereas FO stability is governed solely by the largest Hessian eigenvalue, mean-square stability of ZO methods depends on the entire Hessian spectrum. Since computing the full Hessian spectrum is infeasible in practical neural network training, we further derive tractable stability bounds that depend only on the largest eigenvalue and the Hessian trace. Empirically, we find that full-batch ZO methods operate at the *edge of stability*: ZO-GD, ZO-GDM, and ZO-Adam consistently stabilize near the predicted stability boundary across CNNs, ResNets, and Transformers on vision tasks. Our results highlight an implicit regularization effect specific to ZO methods, where large step sizes primarily regularize the Hessian trace, whereas in FO methods they regularize the top eigenvalue.

1 INTRODUCTION

Zeroth-order (ZO) optimization methods, which rely only on function evaluations, are widely used when gradients are unavailable, unreliable, or expensive to compute. Such a setting arises in black-box learning, derivative-free control, and increasingly in modern large-model pipelines, where memory and systems constraints can make backpropagation costly. Recent work shows that ZO methods based on two-point function evaluations can fine-tune large language models (LLMs) with competitive accuracy while substantially reducing memory usage and compute overhead (Malladi et al., 2023; Zhang et al., 2024b). Despite their growing practical relevance, the training dynamics of ZO methods in deep learning remain far less understood than those of first-order (FO) optimizers.

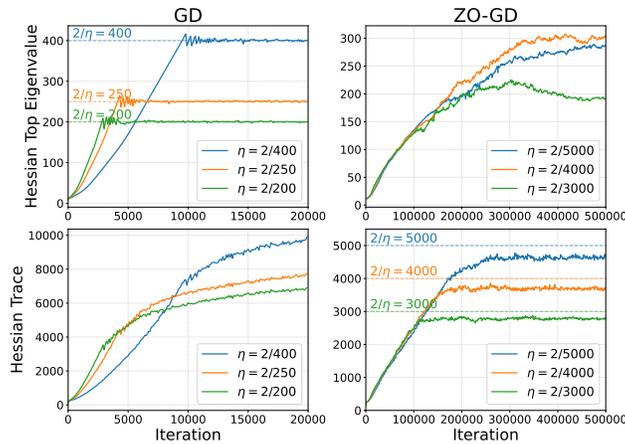


Figure 1: EoS behaviors of FO and ZO methods are captured by different spectral quantities of the Hessian. We train full-batch GD (left) and ZO-GD (right) with varying step sizes η on a CNN for CIFAR-10. For GD, the largest eigenvalue of the Hessian $\lambda_{\max}(\mathbf{H}_t)$ stabilizes near $2/\eta$. For ZO-GD, the trace of the Hessian $\text{Tr}(\mathbf{H}_t)$ instead stabilizes slightly below $2/\eta$.

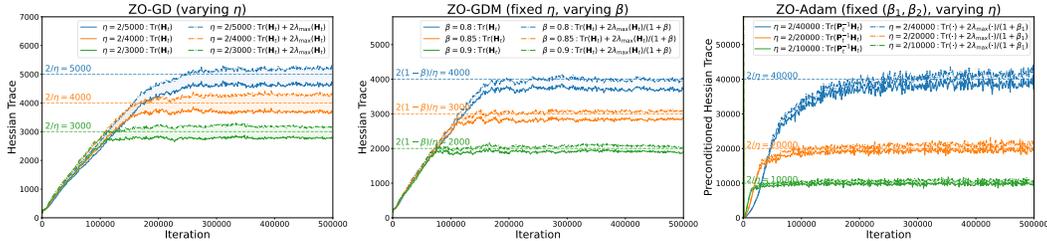


Figure 2: **Zeroth-order methods operate at the mean-square edge of stability.** We train full-batch ZO methods on a CNN for CIFAR-10 and track the curvature terms defining the mean-square stability interval from Appendix C.2. Across all panels, each color denotes one run; the solid curve is the lower-band term, the dash-dotted curve is the upper-band term, and the dashed line is the predicted stability threshold. **Left** (ZO-GD, varying η): lower $\text{Tr}(\mathbf{H}_t)$, upper $\text{Tr}(\mathbf{H}_t) + 2\lambda_{\max}(\mathbf{H}_t)$, threshold $2/\eta$ in Eq. (1). **Middle** (ZO-GDM, varying β with fixed $\eta = 10^{-4}$): lower $\text{Tr}(\mathbf{H}_t)$, upper $\text{Tr}(\mathbf{H}_t) + \frac{2}{1+\beta}\lambda_{\max}(\mathbf{H}_t)$, threshold $2(1-\beta)/\eta$ in Eq. (2). **Right** (ZO-Adam, varying η with fixed $(\beta_1, \beta_2) = (0.9, 0.999)$): lower $\text{Tr}(\mathbf{P}_t^{-1}\mathbf{H}_t)$, upper $\text{Tr}(\mathbf{P}_t^{-1}\mathbf{H}_t) + \frac{2}{1+\beta_1}\lambda_{\max}(\mathbf{P}_t^{-1}\mathbf{H}_t)$, threshold $2/\eta$ in Eq. (3). Across all methods, the threshold stays within (or very close to) the stability interval throughout training, indicating mean-square EoS behavior.

For FO optimization in deep learning, one prominent empirical phenomenon is the *edge of stability* (EoS). In full-batch gradient descent (GD) with step size η on a quadratic objective $f_{\text{quad}}(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{H}\mathbf{x}$, the iterates diverge when $\eta > 2/\lambda_{\max}(\mathbf{H})$. In neural network training, however, optimization often remains well-behaved even at step sizes beyond this local quadratic threshold. Along the training trajectory $\{\mathbf{x}_t\}$, the top eigenvalue of the Hessian $\lambda_{\max}(\mathbf{H}_t)$ increases early in training and then stabilizes near the threshold $2/\eta$ over long horizons (Cohen et al., 2021). This phenomenon has motivated a growing line of work aimed at understanding its mechanism and its connections to stability, curvature, and implicit regularization in deep learning optimization dynamics (Ahn et al., 2022; Arora et al., 2022; Damian et al., 2023; Cohen et al., 2025).

In this work, we ask whether a similar phenomenon arises in *zeroth-order* training. At first glance, this is far from obvious: ZO methods update parameters by drawing random search directions at every iteration, and their stability cannot be understood through the deterministic arguments typically used to explain FO dynamics and stability.

As a preliminary experiment, Figure 1 compares full-batch (first-order) GD and zeroth-order gradient descent (ZO-GD) with varying step sizes η when training a CNN on CIFAR-10. For GD, the top Hessian eigenvalue $\lambda_{\max}(\mathbf{H}_t)$ stabilizes near $2/\eta$, consistent with the behavior predicted by standard EoS theory. For ZO-GD, however, $\lambda_{\max}(\mathbf{H}_t)$ does not exhibit the same trend; instead, perhaps surprisingly, the Hessian trace $\text{Tr}(\mathbf{H}_t)$ stabilizes slightly below $2/\eta$. This suggests that ZO training may be governed by a different stability mechanism than FO training, and that ZO methods may exhibit their own EoS phenomenon in deep learning.

These observations motivate the following fundamental questions:

- (i) do ZO methods operate at the edge of stability in neural network training, and
- (ii) if so, which curvature-related quantity governs their stability?

We introduce a mean-square linear stability theory for ZO methods to investigate these questions. In Appendix C, we provide an exact step size characterization of mean-square linear stability under the linearized dynamics for a family of ZO optimizers, including ZO-GD and its momentum and preconditioned variants. Unlike FO methods, whose stability is governed solely by the largest Hessian eigenvalue, mean-square stability of ZO methods depends on the entire Hessian spectrum. Moreover, momentum affects stability in opposite ways: increasing β enlarges the stable regime for GD with momentum (GDM), but shrinks it for ZO-GD with momentum (ZO-GDM). Since computing the full spectrum is typically infeasible, we also derive tractable bounds that depend only on the Hessian trace and the largest eigenvalue (summarized in Table 1).

In Section 2, we empirically find that full-batch ZO methods operate at the *mean-square* edge of stability: across architectures (CNNs, ResNets, and Vision Transformers), ZO methods (ZO-GD, ZO-GDM, and ZO-Adam) consistently stabilize near the predicted mean-square stability boundary (see Figures 2, 6 and 7). Notably, this behavior is governed primarily by trace-based curvature quantities, providing new insight into how large step sizes implicitly bias ZO training toward solutions with small (preconditioned) Hessian trace.

Table 1: Summary of linear stability conditions for FO and ZO methods under the linearized dynamics (Definition 1). For FO methods of GD, GDM, and Adam, we report the exact critical step size η^* , and the mean and mean-square thresholds coincide. For ZO methods, we report lower and upper bounds on the mean-square critical step size η_{ms}^* ; the critical step size in the mean, η_{mean}^* , matches the corresponding FO threshold. For Adam and ZO-Adam, the reported conditions correspond to the *frozen-preconditioner* variants and are governed by the spectrum of the preconditioned Hessian $\mathbf{P}^{-1}\mathbf{H}$. See Appendix C for details.

Method	Linear Stability Condition
GD	$\eta_{\text{mean}}^* = \eta_{\text{ms}}^* = \frac{2}{\lambda_{\max}(\mathbf{H})}$
GDM (Cohen et al., 2021)	$\eta_{\text{mean}}^* = \eta_{\text{ms}}^* = \frac{2(1+\beta)}{\lambda_{\max}(\mathbf{H})}$
Adam (Cohen et al., 2022)	$\eta_{\text{mean}}^* = \eta_{\text{ms}}^* = \frac{2(1+\beta_1)}{(1-\beta_1)\lambda_{\max}(\mathbf{P}^{-1}\mathbf{H})}$
ZO-GD (Theorem 1)	$\frac{2}{\text{Tr}(\mathbf{H}) + 2\lambda_{\max}(\mathbf{H})} \leq \eta_{\text{ms}}^* \leq \frac{2}{\text{Tr}(\mathbf{H})}$
ZO-GDM (Theorem 2)	$\frac{2(1-\beta)}{\text{Tr}(\mathbf{H}) + \frac{2\lambda_{\max}(\mathbf{H})}{1+\beta}} \leq \eta_{\text{ms}}^* \leq \frac{2(1-\beta)}{\text{Tr}(\mathbf{H})}$
ZO-Adam (Theorem 3)	$\frac{2}{\text{Tr}(\mathbf{P}^{-1}\mathbf{H}) + \frac{2\lambda_{\max}(\mathbf{P}^{-1}\mathbf{H})}{1+\beta_1}} \leq \eta_{\text{ms}}^* \leq \frac{2}{\text{Tr}(\mathbf{P}^{-1}\mathbf{H})}$

2 ZERO-ORDER OPTIMIZATION OPERATES AT THE MEAN-SQUARE EOS

Building on the mean-square linear stability theory in Appendix C.2, we empirically show that full-batch ZO methods on neural networks operate at the *mean-square* edge of stability (EoS): the training dynamics stabilizes near the predicted mean-square linear stability boundary.

2.1 TRACKING MEAN-SQUARE STABILITY DURING TRAINING

Let $\mathbf{H}_t := \nabla^2 f(\mathbf{x}_t)$ denote the loss Hessian along the training trajectory. For each ZO optimizer, Appendix C.2 provides explicit mean-square stability conditions under the linearized dynamics, together with computable lower and upper bounds on the corresponding stability threshold that depend only on the trace and top eigenvalue of the relevant curvature matrix. In large-scale neural network training, however, evaluating the *exact* mean-square stability condition is typically infeasible, since it requires the full spectrum of the Hessian (or the preconditioned Hessian for Adam-style methods). We therefore estimate only the trace and top eigenvalue during training, and use them to form tractable lower and upper bounds. Concretely, for the purpose of visualizing these dynamic conditions between \mathbf{H}_t and the step size η , we present our results in the following format:

$$(\text{lower term}) \leq (\text{stability threshold}) \leq (\text{upper term}),$$

where the stability threshold depends on the step size η and does not change over iterations and the upper and lower terms depend on the spectrum of \mathbf{H}_t (e.g., Figure 2). We say the training operates near the mean-square EoS when the stability threshold remains within, or very close to, this interval for a sustained portion of training.

ZO-GD. From (4) in Theorem 1, we track mean-square stability via

$$\text{Tr}(\mathbf{H}_t) \leq \frac{2}{\eta} \leq \text{Tr}(\mathbf{H}_t) + 2\lambda_{\max}(\mathbf{H}_t). \quad (1)$$

ZO-GDM. From (5) in Theorem 2, we track mean-square stability via

$$\text{Tr}(\mathbf{H}_t) \leq \frac{2(1-\beta)}{\eta} \leq \text{Tr}(\mathbf{H}_t) + \frac{2\lambda_{\max}(\mathbf{H}_t)}{1+\beta}. \quad (2)$$

ZO-Adam. For ZO-Adam, the bounds depend on the preconditioner \mathbf{P}_t at iteration t . From (6) in Theorem 3, we track mean-square stability via

$$\text{Tr}(\mathbf{P}_t^{-1}\mathbf{H}_t) \leq \frac{2}{\eta} \leq \text{Tr}(\mathbf{P}_t^{-1}\mathbf{H}_t) + \frac{2\lambda_{\max}(\mathbf{P}_t^{-1}\mathbf{H}_t)}{1+\beta_1}. \quad (3)$$

2.2 EXPERIMENTAL SETUP

We consider an image classification task on a subset of CIFAR-10 and train ZO methods using the squared loss. We evaluate three representative vision architectures: a CNN, a ResNet, and a Vision Transformer (ViT). Unless stated otherwise, we use full-batch training and a constant step size to match the linearized stability theory in Appendix C. Experimental details are provided in Appendix G.

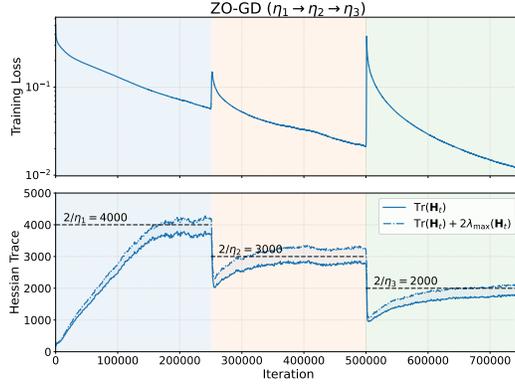


Figure 3: **Catapult dynamics in ZO-GD.** We train ZO-GD on CNN and increase the step size midway through training (from η_1 to η_2 and then to η_3). *Top*: the training loss exhibits a pronounced spike after each step size increase, consistent with catapult dynamics. *Bottom*: the Hessian trace $\text{Tr}(\mathbf{H}_t)$ drops sharply during the catapult phase and then rises again, re-equilibrating near the new stability threshold $2/\eta$.

2.3 MAIN EXPERIMENTS

In Figure 2, we train ZO methods on the CNN and track the stability intervals in (1), (2), and (3). Across all three optimizers, we observe a consistent mean-square EoS pattern: after an initial phase of progressive sharpening, the tracked curvature terms adjust and stabilize so that the stability threshold remains within, or very close to, the corresponding intervals. As shown in Figure 6 and Figure 7, we observe the same behavior when training a ResNet and a ViT.

Specifically, (i) for ZO-GD, the threshold $2/\eta$ remains close to the interval $[\text{Tr}(\mathbf{H}_t), \text{Tr}(\mathbf{H}_t) + 2\lambda_{\max}(\mathbf{H}_t)]$ across step sizes; (ii) for ZO-GDM, the threshold $2(1 - \beta)/\eta$ remains close to $[\text{Tr}(\mathbf{H}_t), \text{Tr}(\mathbf{H}_t) + \frac{2}{1+\beta}\lambda_{\max}(\mathbf{H}_t)]$ across momentum values; and (iii) for ZO-Adam, the threshold $2/\eta$ remains close to the corresponding preconditioned interval in (3). In all cases, the *trace* of the (preconditioned) Hessian provides the dominant stability signal throughout training.

2.4 ADDITIONAL EXPERIMENTS

Catapult dynamics. In Figure 3, we train ZO-GD and increase the step size midway through training (from η_1 to η_2 and then to η_3). We observe a pronounced spike in the training loss after each increase, consistent with the *catapult* dynamics (Lewkowycz et al., 2020). Immediately after the step size increase, the new step size temporarily exceeds the mean-square stability critical step size at the current iterate, so the dynamics become locally unstable and the loss increases sharply. At the same time, the Hessian trace drops rapidly below the new threshold and then rises again, re-equilibrating near the new threshold.

Effect of the smoothing parameter μ . In Figure 4, we train ZO-GD with fixed step size and vary the smoothing parameter μ in the two-point estimator. For moderate and small μ , the tracked stability terms increase early in training and then stabilize near the threshold $2/\eta$, consistent with mean-square EoS behavior. For larger μ , both $\text{Tr}(\mathbf{H}_t)$ and $\text{Tr}(\mathbf{H}_t) + 2\lambda_{\max}(\mathbf{H}_t)$ saturate at substantially smaller values and remain far below $2/\eta$, indicating that training does not approach the predicted mean-square stability boundary. Overall, mean-square EoS persists across a broad range of practically relevant smoothing levels, while overly large smoothing suppresses curvature growth. We attribute this phenomenon to implicit bias in Section D.

Beyond full-batch: mini-batch ZO-SGD. Although our main results focus on full-batch ZO methods, we include a preliminary mini-batch experiment in Figure 5. Compared to full-batch ZO-GD, mini-batch ZO-SGD converges to significantly flatter regimes.

3 CONCLUSION

In this work, we characterize exact mean-square stability thresholds for zeroth-order optimizers and empirically show that neural network training operates at the mean-square edge of stability. Due to space constraints, we defer detailed related work (Appendix A), the linear stability analysis and proofs (Appendix B, C, and F), additional experimental details and results (Appendix G and H), and extended discussion (Appendix D and E) to the appendix.

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

- Kwangjun Ahn, Jingzhao Zhang, and Suvrit Sra. Understanding the unstable convergence of gradient descent. In *International Conference on Machine Learning*, 2022.
- Arseniy Andreyev and Pierfrancesco Beneventano. Edge of stochastic stability: Revisiting the edge of stability for SGD. *arXiv preprint arXiv:2412.20553*, 2024.
- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: A query-efficient black-box adversarial attack via random search. In *European Conference on Computer Vision*, 2020.
- Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding gradient descent on the edge of stability in deep learning. In *International Conference on Machine Learning*, 2022.
- Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *ACM Workshop on Artificial Intelligence and Security*, 2017.
- Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021.
- Jeremy Cohen, Behrooz Ghorbani, Shankar Krishnan, Naman Agarwal, Sourabh Medapati, Michal Badura, Daniel Suo, David Cardoze, Zachary Nado, George E Dahl, et al. Adaptive gradient methods at the edge of stability. *arXiv preprint arXiv:2207.14484*, 2022.
- Jeremy Cohen, Alex Damian, Ameet Talwalkar, J Zico Kolter, and Jason D. Lee. Understanding optimization in deep learning with central flows. In *International Conference on Learning Representations*, 2025.
- Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of gradient descent at the edge of stability. In *International Conference on Learning Representations*, 2023.
- Aaron Defazio, Xingyu Yang, Ahmed Khaled, Konstantin Mishchenko, Harsh Mehta, and Ashok Cutkosky. The road less scheduled. In *Advances in Neural Information Processing Systems*, 2024.
- Klaus Deimling. *Nonlinear Functional Analysis*. Springer Berlin, Heidelberg, 1985.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- John C. Duchi, Michael I. Jordan, Martin J. Wainwright, and Andre Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.
- Wenzhi Fang, Ziyi Yu, Yuning Jiang, Yuanming Shi, Colin N. Jones, and Yong Zhou. Communication-efficient stochastic zeroth-order optimization for federated learning. *IEEE Transactions on Signal Processing*, 70:5058–5073, 2022.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021.
- Diego Granzio, Stefan Zohren, and Stephen Roberts. 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.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- Kevin G. Jamieson, Robert Nowak, and Ben Recht. Query complexity of derivative-free optimization. In *Advances in Neural Information Processing Systems*, 2012.

270 Mark Grigor’evich Krein and Moisei Aronovich Rutman. Linear operators leaving invariant a cone
271 in a banach space. *Uspekhi Matematicheskikh Nauk*, 3:3–95, 1948.

272
273 Sungyoon Lee and Cheongjae Jang. A new characterization of the edge of stability based on a
274 sharpness measure aware of batch gradient distribution. In *International Conference on Learning*
275 *Representations*, 2023.

276 Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large
277 learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*,
278 2020.

279 Philip M. Long and Peter L. Bartlett. Sharpness-aware minimization and the edge of stability. *Journal*
280 *of Machine Learning Research*, 25(179):1–20, 2024.

281
282 Chao Ma and Lexing Ying. On linear stability of SGD and input-smoothness of neural networks. In
283 *Advances in Neural Information Processing Systems*, 2021.

284
285 Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev
286 Arora. Fine-tuning language models with just forward passes. In *Advances in Neural Information*
287 *Processing Systems*, 2023.

288 Michael Muehlebach and Michael I. Jordan. Optimization with momentum: Dynamical, control-
289 theoretic, and symplectic perspectives. *Journal of Machine Learning Research*, 22(73):1–50,
290 2021.

291 Rotem Mulayoff and Tomer Michaeli. Exact mean square linear stability analysis for SGD. In
292 *Conference on Learning Theory*, 2024.

293
294 Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions.
295 *Foundations of Computational Mathematics*, 17(2):527–566, 2017.

296
297 Zhen Qin, Daoyuan Chen, Bingchen Qian, Bolin Ding, Yaliang Li, and Shuiguang Deng. Federated
298 full-parameter tuning of billion-sized language models with communication cost under 18 kilobytes.
299 In *International Conference on Machine Learning*, 2024.

300 Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a
301 scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.

302
303 Ohad Shamir. An optimal algorithm for bandit and zero-order convex optimization with two-point
304 feedback. *Journal of Machine Learning Research*, 18(52):1–11, 2017.

305 Minhak Song and Chulhee Yun. Trajectory alignment: Understanding the edge of stability phe-
306 nomenon via bifurcation theory. In *Advances in Neural Information Processing Systems*, 2023.

307
308 Minhak Song, Beomhan Baek, Kwangjun Ahn, and Chulhee Yun. Through the river: Understand-
309 ing the benefit of schedule-free methods for language model training. In *Advances in Neural*
310 *Information Processing Systems*, 2025.

311 Maksim Velikanov, Denis Kuznedelev, and Dmitry Yarotsky. A view of mini-batch SGD via
312 generating functions: conditions of convergence, phase transitions, benefit from negative momenta.
313 In *International Conference on Learning Representations*, 2023.

314
315 Yining Wang, Simon Du, Sivaraman Balakrishnan, and Aarti Singh. Stochastic zeroth-order opti-
316 mization in high dimensions. In *International Conference on Artificial Intelligence and Statistics*,
317 2018.

318
319 Zixuan Wang, Zhouzi Li, and Jian Li. Analyzing sharpness along GD trajectory: Progressive
320 sharpening and edge of stability. In *Advances in Neural Information Processing Systems*, 2022.

321
322 Lei Wu, Chao Ma, and Weinan E. How SGD selects the global minima in over-parameterized learning:
323 A dynamical stability perspective. In *Advances in Neural Information Processing Systems*, 2018.

324 Mengwei Xu, Dongqi Cai, Yaozong Wu, Xiang Li, and Shangguang Wang. FwdLLM: Efficient
325 federated finetuning of large language models with perturbed inferences. In *USENIX Annual*
326 *Technical Conference*, 2024.

327
328 Geonhui Yoo, Minhak Song, and Chulhee Yun. Understanding sharpness dynamics in NN training
329 with a minimalist example: The effects of dataset difficulty, depth, stochasticity, and more. In
330 *International Conference on Machine Learning*, 2025.

331 Liang Zhang, Bingcong Li, Kiran Koshy Thekumparampil, Sewoong Oh, and Niao He. Dpzero:
332 Private fine-tuning of language models without backpropagation. In *International Conference on*
333 *Machine Learning*, 2024a.

334 Liang Zhang, Bingcong Li, Kiran Koshy Thekumparampil, Sewoong Oh, Michael Muehlebach,
335 and Niao He. Zeroth-order optimization finds flat minima. In *Advances in Neural Information*
336 *Processing Systems*, 2025a.

337
338 Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiaxiang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu
339 Chen, Jason D. Lee, Wotao Yin, Mingyi Hong, Zhangyang Wang, Sijia Liu, and Tianlong Chen.
340 Revisiting zeroth-order optimization for memory-efficient LLM fine-tuning: A benchmark. In
341 *International Conference on Machine Learning*, 2024b.

342 Yushun Zhang, Congliang Chen, Ziniu Li, Tian Ding, Chenwei Wu, Diederik P Kingma, Yinyu
343 Ye, Zhi-Quan Luo, and Ruoyu Sun. Adam-mini: Use fewer learning rates to gain more. In
344 *International Conference on Learning Representations*, 2025b.

345
346 Xingyu Zhu, Zixuan Wang, Xiang Wang, Mo Zhou, and Rong Ge. Understanding edge-of-stability
347 training dynamics with a minimalist example. In *International Conference on Learning Representations*,
348 2023.

349 Liu Ziyin, Botao Li, Tomer Galanti, and Masahito Ueda. Type-ii saddles and probabilistic stability of
350 stochastic gradient descent. *arXiv preprint arXiv:2303.13093*, 2023.

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

Appendix

378		
379		
380		
381		
382	A Related work	9
383		
384	B Preliminaries	9
385		
386	C Linear stability analysis	10
387		
388	C.1 Linear stability of first-order methods	11
389	C.2 Linear stability of zeroth-order methods	11
390		
391	D Discussion	12
392		
393		
394	E Summary and Conclusion	13
395		
396	F Proofs for the mean-square stability analysis of zeroth-order methods	14
397		
398	F.1 Second-moment recursion and covariance operator for ZO-GDM	14
399	F.2 Spectral analysis of the covariance operator	16
400	F.3 Proof of Theorem 2	22
401	F.4 Proof of Theorem 3	24
402	F.5 Spectral analysis of Frozen ZO-Adam covariance operator	27
403	F.6 Technical Lemmas	29
404		
405		
406	G Experimental details	30
407		
408		
409	H Additional experiments	31
410	H.1 Effect of β_1 in ZO-Adam	31
411	H.2 Empirical verification of the commutativity approximation in ZO-Adam	31
412		
413		
414		
415		
416		
417		
418		
419		
420		
421		
422		
423		
424		
425		
426		
427		
428		
429		
430		
431		

A RELATED WORK

Zeroth-order optimization. ZO methods optimize objectives using only function evaluations, typically via the standard two-point estimator along a random direction. Empirically, Malladi et al. (2023) report that ZO methods can fine-tune LLMs with performance comparable to FO methods while achieving up to $12\times$ memory and up to $2\times$ GPU-hour reduction, highlighting ZO optimization as a promising approach for memory-efficient fine-tuning. ZO methods have also been applied to adversarial robustness (Chen et al., 2017; Andriushchenko et al., 2020), reinforcement learning (Salimans et al., 2017), private fine-tuning (Zhang et al., 2024a) and distributed learning (Fang et al., 2022; Qin et al., 2024; Xu et al., 2024), where gradients may be noisy, unavailable, or expensive to compute or communicate.

On the theory side, most prior work studies ZO methods through the lens of classical (non)convex optimization (Jamieson et al., 2012; Duchi et al., 2015; Shamir, 2017; Wang et al., 2018; Malladi et al., 2023; Zhang et al., 2024a), emphasizing convergence under small step sizes and smoothness assumptions. Notably, Zhang et al. (2025a) study the implicit bias of ZO-GD under smooth convex objectives and show that it favors solutions with small Hessian trace.

Edge of stability. Recent empirical work shows that FO methods often train near instability. In particular, Cohen et al. (2021) identify the edge of stability (EoS) in full-batch GD, where $\lambda_{\max}(\mathbf{H}_t)$ grows early in training and then equilibrates near $2/\eta$. This observation has motivated extensive follow-up work on the mechanisms and implications of EoS (Ahn et al., 2022; Arora et al., 2022; Wang et al., 2022; Damian et al., 2023; Song & Yun, 2023; Zhu et al., 2023; Yoo et al., 2025; Cohen et al., 2025). EoS-type behavior has also been studied for momentum and adaptive optimizers (Cohen et al., 2022), mini-batch stochastic gradient descent (SGD) (Lee & Jang, 2023; Andreyev & Beneventano, 2024), and other families of FO optimizers, including sharpness-aware minimization (Foret et al., 2021; Long & Bartlett, 2024) and schedule-free methods (Defazio et al., 2024; Song et al., 2025).

Dynamical stability analysis of optimizers. A growing body of work studies optimization methods through the lens of dynamical stability. Recent work analyzes *linear stability* by examining the behavior of the linearized dynamics under a local quadratic approximation, extending beyond GD to momentum methods (Muehlebach & Jordan, 2021) and stochastic optimization. For SGD, prior analyses have characterized linear stability condition in the mean-square sense (Wu et al., 2018; Granzio et al., 2022; Velikanov et al., 2023), for higher moments (Ma & Ying, 2021), and in probability (Ziyin et al., 2023). Most closely related to our setting, Mulayoff & Michaeli (2024) derive the exact mean-square stability threshold of mini-batch SGD and show that it is monotonically non-decreasing in the batch size. Our work also studies mean-square linear stability, but for ZO methods, where stochasticity arises from the estimator directions even under full-batch training.

B PRELIMINARIES

We consider the optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}),$$

for a loss function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and a model parameter \mathbf{x} , and study the dynamics of zeroth-order (ZO) optimization methods. ZO methods iteratively update a sequence of iterates $\{\mathbf{x}_t\}_{t \geq 0}$ based solely on function evaluations without evaluating any gradients. We analyze the three most popular types of ZO methods: ZO-GD, ZO-GDM, and ZO-Adam.

ZO Gradient Descent (ZO-GD) replaces the gradient $\nabla f(\mathbf{x}_t)$ in the GD update with a gradient estimate $\hat{\nabla} f(\mathbf{x}_t)$:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \hat{\nabla} f(\mathbf{x}_t),$$

for a step size $\eta > 0$. We consider the standard two-point estimator (Nesterov & Spokoiny, 2017), defined as

$$\hat{\nabla} f(\mathbf{x}_t) := \frac{f(\mathbf{x}_t + \mu \mathbf{u}_t) - f(\mathbf{x}_t - \mu \mathbf{u}_t)}{2\mu} \cdot \mathbf{u}_t,$$

where $\mathbf{u}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\mu > 0$ is a smoothing parameter. In general, $\hat{\nabla} f(\mathbf{x}_t)$ is a biased estimator of $\nabla f(\mathbf{x}_t)$, with bias that vanishes as $\mu \rightarrow 0$.

ZO Gradient Descent with Momentum (ZO-GDM) combines Polyak momentum terms with gradient estimates:

$$\begin{aligned}\mathbf{m}_{t+1} &= \beta \mathbf{m}_t + \widehat{\nabla} f(\mathbf{x}_t), \\ \mathbf{x}_{t+1} &= \mathbf{x}_t - \eta \mathbf{m}_{t+1},\end{aligned}$$

for a step size $\eta > 0$ and a momentum parameter $\beta \in [0, 1)$.

ZO-Adam combines adaptive moment estimation (Adam) with gradient estimates:

$$\begin{aligned}\mathbf{m}_{t+1} &= \beta_1 \mathbf{m}_t + (1 - \beta_1) \widehat{\nabla} f(\mathbf{x}_t), \\ \mathbf{x}_{t+1} &= \mathbf{x}_t - \eta \mathbf{P}_{t+1}^{-1} \mathbf{m}_{t+1},\end{aligned}$$

for step size $\eta > 0$, momentum parameters $\beta_1, \beta_2 \in [0, 1)$, and $\epsilon > 0$, where \mathbf{P}_{t+1} is a preconditioner defined using an exponential moving average (EMA) of squared (estimated) gradients:

$$\begin{aligned}\boldsymbol{\nu}_{t+1} &= \beta_2 \boldsymbol{\nu}_t + (1 - \beta_2) \widehat{\nabla} f(\mathbf{x}_t) \odot \widehat{\nabla} f(\mathbf{x}_t), \\ \mathbf{P}_{t+1} &= (1 - \beta_1^{t+1}) \left[\text{diag} \left(\sqrt{\frac{\boldsymbol{\nu}_{t+1}}{1 - \beta_2^{t+1}}} \right) + \epsilon \mathbf{I} \right],\end{aligned}$$

and \odot denotes element-wise multiplication. This form is equivalent to the standard bias-corrected Adam update, written as a preconditioned momentum step.

In practical neural network training, the short-term stability behavior of Adam is often well approximated by a *frozen-preconditioner* variant, in which the preconditioner is held fixed at its current value (Cohen et al., 2022). Motivated by this, we also consider the corresponding ZO analogue.

Frozen ZO-Adam uses a fixed preconditioner $\mathbf{P} > 0$ with a step size $\eta > 0$:

$$\begin{aligned}\mathbf{m}_{t+1} &= \beta_1 \mathbf{m}_t + (1 - \beta_1) \widehat{\nabla} f(\mathbf{x}_t), \\ \mathbf{x}_{t+1} &= \mathbf{x}_t - \eta \mathbf{P}^{-1} \mathbf{m}_{t+1},\end{aligned}$$

and a momentum parameter $\beta_1 \in [0, 1)$. This optimizer is introduced purely for theoretical analysis and serves as the ZO analogue of *Frozen Adam*, which was used to study the *adaptive edge of stability* of Adam (Cohen et al., 2022).

C LINEAR STABILITY ANALYSIS

Directly analyzing the full dynamics of optimizers in deep learning is typically intractable. Instead, we adopt the standard dynamical systems approach of studying stability near a local minimizer via the *linearized dynamics*. Exponential stability of the linearized dynamics implies local stability of the corresponding nonlinear dynamics near the equilibrium, which justifies linear stability analysis as a principled tool.

Definition 1 (Linearized dynamics). Let \mathbf{x}^* be a twice differentiable local minimizer of f , and let $\mathbf{H} := \nabla^2 f(\mathbf{x}^*)$ be the Hessian at \mathbf{x}^* . The *linearized dynamics* of an optimizer around \mathbf{x}^* are the dynamics obtained by applying the optimizer to the quadratic Taylor approximation of f at \mathbf{x}^* ,

$$f_{\text{quad}}(\mathbf{x}) := f(\mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^\top \mathbf{H}(\mathbf{x} - \mathbf{x}^*).$$

Throughout, we assume $\mathbf{H} \geq 0$ and $\mathbf{H} \neq \mathbf{0}$. Let $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ denote the eigenvalues of \mathbf{H} , and define $\lambda_{\max}(\mathbf{H}) := \lambda_1$ and $\text{Tr}(\mathbf{H}) := \sum_{i=1}^d \lambda_i$.

We next formalize the notion of linear stability for the resulting (possibly stochastic) optimizer dynamics.

Definition 2 (Linear stability). Let \mathbf{x}^* be as in Definition 1, and let $\{\mathbf{x}_t\}_{t \geq 0}$ denote the iterates generated by an optimizer applied to f_{quad} .

We say the optimizer is *linearly stable in the mean* if

$$\sup_{t \geq 0} \mathbb{E}[\|\mathbf{x}_t - \mathbf{x}^*\|] < \infty \quad \text{for every initialization } \mathbf{x}_0 \in \mathbb{R}^d.$$

We say the optimizer is *mean-square linearly stable* if

$$\sup_{t \geq 0} \mathbb{E}[\|\mathbf{x}_t - \mathbf{x}^*\|^2] < \infty \quad \text{for every initialization } \mathbf{x}_0 \in \mathbb{R}^d.$$

Mean-square linear stability implies linear stability in the mean by Jensen's inequality. For deterministic optimizers, the expectation is redundant, and the two notions coincide.

Our goal is to theoretically characterize the *critical step size* that guarantees linear stability of ZO methods and empirically connect it to the edge of stability phenomena.

Definition 3 (Critical step size). Consider an optimizer with step size $\eta > 0$ applied to f_{quad} , producing iterates $\{\mathbf{x}_t\}_{t \geq 0}$. The *critical step size in the mean* is defined as

$$\eta_{\text{mean}}^* := \sup \left\{ \eta > 0 : \forall \mathbf{x}_0 \in \mathbb{R}^d, \sup_{t \geq 0} \|\mathbb{E}[\mathbf{x}_t - \mathbf{x}^*]\| < \infty \right\},$$

and the *critical step size in the mean-square* is defined as

$$\eta_{\text{ms}}^* := \sup \left\{ \eta > 0 : \forall \mathbf{x}_0 \in \mathbb{R}^d, \sup_{t \geq 0} \mathbb{E}[\|\mathbf{x}_t - \mathbf{x}^*\|^2] < \infty \right\}.$$

Next, we review known stability thresholds for FO methods (Appendix C.1) and present our main results for ZO methods (Appendix C.2). All proofs are deferred to Appendix F.

C.1 LINEAR STABILITY OF FIRST-ORDER METHODS

The first-order (FO) counterparts of the ZO optimizers in Section B, namely GD, GDM, and Frozen Adam, are deterministic, so their mean and mean-square linear stability conditions coincide. A key takeaway is that FO linear stability is governed solely by the top eigenvalue of the (preconditioned) Hessian, as summarized below.

Proposition 1 (Stability of GD). *The critical step size of GD is $\eta_{\text{mean}}^* = \eta_{\text{ms}}^* = 2/\lambda_{\max}(\mathbf{H})$.*

The stability condition of GD with Momentum (GDM) was established in Theorem 2 of Cohen et al. (2021).

Proposition 2 (Stability of GDM). *The critical step size of GDM is $\eta_{\text{mean}}^* = \eta_{\text{ms}}^* = 2(1 + \beta)/\lambda_{\max}(\mathbf{H})$.*

The stability condition of Frozen Adam was established in Lemma 2 and Proposition 1 of Cohen et al. (2022).

Proposition 3 (Stability of Frozen Adam). *The critical step size of Frozen Adam is*

$$\eta_{\text{mean}}^* = \eta_{\text{ms}}^* = \frac{2(1 + \beta_1)}{(1 - \beta_1)\lambda_{\max}(\mathbf{P}^{-1}\mathbf{H})}.$$

C.2 LINEAR STABILITY OF ZERO-ORDER METHODS

The inherent stochasticity in ZO updates fundamentally changes the linear stability as shown in Figure 1. However, the *mean* dynamics of a ZO method matches that of its FO counterpart, and η_{mean}^* coincides with the FO critical step size presented in Appendix C.1. This is due to the fact that under the quadratic model f_{quad} in Definition 1, the two-point estimator is *unbiased*: $\mathbb{E}[\widehat{\nabla} f_{\text{quad}}(\mathbf{x}_t)] = \nabla f_{\text{quad}}(\mathbf{x}_t)$. The intricacy of ZO linear stability is only captured by the *mean-square* stability. In particular, we show that η_{ms}^* for ZO-GD, ZO-GDM, and Frozen ZO-Adam, depends on the entire eigen spectrum of the (preconditioned) Hessian and is dominated by its trace value.

Theorem 1 (Stability of ZO-GD). *For ZO-GD, $\eta_{\text{mean}}^* = 2/\lambda_{\max}(\mathbf{H})$, and the mean-square critical step size η_{ms}^* is the unique $\eta > 0$ satisfying*

$$\eta \lambda_{\max}(\mathbf{H}) < 1 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \lambda_i}{2(1 - \eta \lambda_i)} = 1,$$

which admits the bounds

$$\frac{2}{\text{Tr}(\mathbf{H}) + 2\lambda_{\max}(\mathbf{H})} \leq \eta_{\text{ms}}^* \leq \frac{2}{\text{Tr}(\mathbf{H})}. \quad (4)$$

Remark 1 (Computing η_{ms}^*). If the full spectrum $\{\lambda_i\}_{i=1}^d$ were available, η_{ms}^* in Theorem 1 can be computed by solving $\sum_{i=1}^d (\eta \lambda_i / 2(1 - \eta \lambda_i)) = 1$ over $\eta \in (0, 1/\lambda_{\max}(\mathbf{H}))$. In practice, we instead track the bounds (4), which depend only on $\text{Tr}(\mathbf{H})$ and $\lambda_{\max}(\mathbf{H})$ and can be estimated efficiently during training, which is critical for large models.

Theorem 2 (Stability of ZO-GDM). For ZO-GDM, $\eta_{\text{mean}}^* = 2(1 + \beta)/\lambda_{\max}(\mathbf{H})$, and the mean-square critical step size η_{ms}^* is the unique $\eta > 0$ satisfying

$$\eta\lambda_{\max}(\mathbf{H}) < 1 - \beta^2 \text{ and } \sum_{i=1}^d \frac{\eta\lambda_i}{2(1 - \beta)(1 - \frac{\eta\lambda_i}{1 - \beta^2})} = 1,$$

which admits the bounds

$$\frac{2(1 - \beta)}{\text{Tr}(\mathbf{H}) + \frac{2\lambda_{\max}(\mathbf{H})}{1 + \beta}} \leq \eta_{\text{ms}}^* \leq \frac{2(1 - \beta)}{\text{Tr}(\mathbf{H})}. \quad (5)$$

Proof sketch. We analyze the recursions of the second-moment matrices $\mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top]$, $\mathbb{E}[\mathbf{x}_t \mathbf{m}_t^\top]$, and $\mathbb{E}[\mathbf{m}_t \mathbf{m}_t^\top]$. Using the Isserlis’ Theorem to evaluate Gaussian fourth moments, we obtain a linear recursion that can be expressed as a cone-preserving linear operator on a product cone. Mean-square stability is then equivalent to this operator having spectral radius smaller than one. We characterize this spectral radius using Theorem 4 in Appendix F.2, which leverages the Krein–Rutman Theorem and leads to the explicit mean-square stability condition. \square

Theorem 3 (Stability of Frozen ZO-Adam). Let $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_d \geq 0$ denote the eigenvalues of the preconditioned Hessian $\mathbf{P}^{-1}\mathbf{H}$. For Frozen ZO-Adam, $\eta_{\text{mean}}^* = 2(1 + \beta_1)/((1 - \beta_1)\lambda_{\max}(\mathbf{P}^{-1}\mathbf{H}))$. Assuming $\mathbf{P}\mathbf{H} = \mathbf{H}\mathbf{P}$, the mean-square critical step size η_{ms}^* is the unique $\eta > 0$ satisfying

$$\eta\lambda_{\max}(\mathbf{P}^{-1}\mathbf{H}) < 1 + \beta_1 \text{ and } \sum_{i=1}^d \frac{\eta\tilde{\lambda}_i}{2(1 - \frac{\eta\tilde{\lambda}_i}{1 + \beta_1})} = 1,$$

and it admits the bounds

$$\frac{2}{\text{Tr}(\mathbf{P}^{-1}\mathbf{H}) + \frac{2\lambda_{\max}(\mathbf{P}^{-1}\mathbf{H})}{1 + \beta_1}} \leq \eta_{\text{ms}}^* \leq \frac{2}{\text{Tr}(\mathbf{P}^{-1}\mathbf{H})}. \quad (6)$$

Remark 2 (Commutativity assumption). The condition $\mathbf{P}\mathbf{H} = \mathbf{H}\mathbf{P}$ ensures that $\mathbf{P}^{-1}\mathbf{H}$ is diagonalizable in the same eigenbasis as \mathbf{H} , which allows the mean-square dynamics to decouple across eigendirections and enables a tractable spectral analysis. Without commutativity, the second-moment recursion generally couples different eigenspaces, and obtaining an explicit stability characterization becomes substantially less tractable. Empirically, in neural network training with ZO-Adam, we observe that \mathbf{P}_t and \mathbf{H}_t are nearly commuting: the relative commutator Frobenius norm $\|\mathbf{P}_t\mathbf{H}_t - \mathbf{H}_t\mathbf{P}_t\|_F / \|\mathbf{P}_t\mathbf{H}_t\|_F$ decreases from 0.8–0.9 at initialization to below 0.05 and remains below 0.05 throughout training (see Appendix H.2).

Taken together, Theorems 1 to 3 provide exact mean-square stability characterizations under the linearized dynamics. In the next section, we use the corresponding upper and lower bounds in (4), (5), and (6) to empirically test whether ZO methods operate near the mean-square edge of stability during neural network training.

D DISCUSSION

In this section, we discuss implications of our mean-square stability theory and empirical mean-square EoS results, and highlight several directions for future work.

Why mean-square stability is the relevant notion for ZO dynamics. In ZO optimization, randomness persists even in full-batch training due to random perturbation directions. As a result, stability cannot be assessed solely through the mean trajectory; $\mathbb{E}[\mathbf{x}_t]$ may remain bounded even when fluctuations grow and dominate the behavior of the iterates. Mean-square stability captures this effect by directly controlling the second moment $\mathbb{E}\|\mathbf{x}_t - \mathbf{x}^*\|^2$, while remaining analyzable under the linearized dynamics and yielding explicit step size conditions. Other notions of stability are also meaningful, such as stability of higher moments or tail-probability bounds. Nevertheless, our experiments suggest that the curvature quantities appearing in the mean-square stability conditions closely track the stability behavior observed during ZO neural network training.

Curvature quantities that govern ZO stability. For FO methods under the linearized dynamics, stability depends only on the largest eigenvalue of the (preconditioned) Hessian. In contrast, our results show that mean-square stability of ZO methods depends on the full Hessian spectrum. This

dependence appears explicitly in the exact stability conditions, and it is reflected in the computable bounds through both trace and top-eigenvalue terms. A practically important regime is when $\text{Tr}(\mathbf{H})$ dominates $\lambda_{\max}(\mathbf{H})$, in which case these bounds become tight and the trace term largely sets the stability scale. Empirically, this matches our empirical observations: the Hessian trace (or the preconditioned trace for ZO-Adam) closely tracks the relevant stability threshold throughout training, while $\lambda_{\max}(\mathbf{H}_t)$ can be less informative for ZO dynamics.

Momentum reshapes ZO stability differently from FO stability. Momentum provides a concrete example showing that ZO stability is not a direct analogue of FO stability. For GD with momentum (GDM), the linearized stability threshold increases with β , corresponding to stable training at sharper curvature levels, with $\lambda_{\max}(\mathbf{H}_t) \approx 2(1 + \beta)/\eta$ at the EoS. For ZO-GDM, our mean-square conditions imply the opposite dependence: increasing β shrinks the stable regime and reduces the corresponding stability scale. Empirically, this is consistent with ZO-GDM operating in lower-curvature regimes as β increases, with $\text{Tr}(\mathbf{H}_t) \approx 2(1 - \beta)/\eta$ at the mean-square EoS.

This qualitative difference suggests that, in ZO training, momentum not only affects optimization speed but also changes how estimator noise accumulates across iterations, thereby reshaping the stability constraint. A related contrast appears for adaptive methods. For (frozen) Adam, increasing β_1 increases the stability threshold, with $\lambda_{\max}(\mathbf{P}_t^{-1}\mathbf{H}_t) \approx 2(1 + \beta_1)/((1 - \beta_1)\eta)$. For ZO-Adam, by comparison, our experiments suggest that $\text{Tr}(\mathbf{P}_t^{-1}\mathbf{H}_t) \approx 2/\eta$ at the mean-square EoS, which is independent of β_1 (see Figure 8 and Appendix H.1). Understanding how these effects translate into practical benefits (or tradeoffs) of momentum in ZO training is an interesting open question.

Effect of the smoothing parameter μ and trace-related implicit bias. The two-point estimator introduces smoothing, controlled by μ , that changes both the bias and the noise structure of the ZO update. Zhang et al. (2025a) connect such ZO optimization to an implicit preference for small-trace regions, formalized as approximately minimizing

$$f_\mu(\mathbf{x}) := f(\mathbf{x}) + \frac{\mu^2}{2} \text{Tr}(\nabla^2 f(\mathbf{x})),$$

up to higher-order terms. Under this perspective, μ directly modulates a curvature-dependent bias in the effective objective, and large μ can prevent the dynamics from approaching the mean-square stability threshold predicted by the unsmoothed linearized model. A systematic theory that jointly captures (i) the mean-square stability constraint and (ii) the effect of smoothing bias on the effective landscape is an interesting direction for future work.

Beyond full-batch: mini-batch ZO methods. Our analysis focuses on full-batch ZO dynamics, where the only randomness comes from the estimator directions. In practical settings, mini-batching introduces an additional noise source through stochastic sampling of data points. A complete stability theory for mini-batch ZO training would need to incorporate both estimator noise and sampling noise, and quantify how these two sources interact in the second-moment recursion. Recent work gives sharp mean-square stability thresholds for mini-batch SGD (Mulayoff & Michaeli, 2024). Deriving analogous results for mini-batch ZO methods would clarify whether mean-square EoS persists under data subsampling, and which curvature quantities control stability in that regime.

E SUMMARY AND CONCLUSION

We developed a mean-square linear stability theory for zeroth-order (ZO) optimization methods based on the standard two-point estimator, including ZO-GD, ZO-GDM, and Adam-style preconditioned variants. We derived exact step size characterizations for mean-square stability via linearization, showing that ZO stability depends on the full spectrum of the (preconditioned) Hessian and admits computable bounds in terms of the trace and top eigenvalue.

Guided by these results, we empirically studied full-batch ZO training on standard neural network architectures (CNN, ResNet, and ViT) and found consistent evidence that ZO methods operate at the mean-square edge of stability: the curvature quantities governing the theoretical stability threshold adapt during training and stabilize near the predicted boundary. Across methods, the stability behavior is driven primarily by trace-based curvature terms, providing a concrete mechanism through which large step sizes implicitly regularize ZO training dynamics.

Our results position mean-square stability as a principled framework for analyzing and predicting ZO optimization behavior in deep learning. An important next step is to extend the theory beyond full-batch settings to account for additional sources of stochasticity, such as mini-batch sampling noise, and to understand how stability constraints interact with optimization efficiency and generalization in practice.

F PROOFS FOR THE MEAN-SQUARE STABILITY ANALYSIS OF ZERO-ORDER METHODS

In this section, we provide the proofs of the mean-square linear stability results in Appendix C.2. Our analysis treats zeroth-order gradient descent with momentum (ZO-GDM) as the canonical zeroth-order method. The remaining algorithms studied in the paper are handled as special cases or extensions: ZO-GD corresponds to the specialization $\beta = 0$, while Frozen ZO-Adam is analyzed by modifying the covariance recursion under an additional commutativity assumption.

The core of the analysis proceeds in two steps. First, we derive an explicit linear recursion for the second-moment quantities of the ZO-GDM iterates under the linearized dynamics (Appendix F.1). This recursion induces a linear operator acting on a product space of covariance matrices. Second, we characterize mean-square linear stability by analyzing the spectral radius of this operator (Appendix F.2). A key technical ingredient is that the operator preserves a natural cone and includes a rank-one global coupling term, which allows its spectral radius to be characterized via the Krein–Rutman Theorem.

We first derive the second-moment recursion for ZO-GDM and express it in operator form. We then develop the spectral analysis of the resulting covariance operator and use it to prove Theorem 2. Theorem 1 follows immediately as the special case $\beta = 0$, and Theorem 3 is proved by adapting the analogous argument to the frozen preconditioned setting.

The proof of Theorem 2 is provided in Appendix F.3, and the proof of Theorem 3 is provided in Appendix F.4.

F.1 SECOND-MOMENT RECURSION AND COVARIANCE OPERATOR FOR ZO-GDM

In this subsection, we derive the linear recursion governing the second-moment dynamics of ZO-GDM under the linearized dynamics. Following Definition 1, we consider applying ZO-GDM to the quadratic objective

$$f_{\text{quad}}(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x},$$

where $\mathbf{H} \geq 0$ and $\mathbf{H} \neq \mathbf{0}$.

Recall that the ZO-GDM updates are given by

$$\mathbf{m}_{t+1} = \beta \mathbf{m}_t + \widehat{\nabla} f_{\text{quad}}(\mathbf{x}_t),$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathbf{m}_{t+1},$$

where $\beta \in [0, 1)$ with initialization $\mathbf{m}_0 = \mathbf{0}$, and $\widehat{\nabla} f_{\text{quad}}(\mathbf{x}_t)$ denotes the standard two-point estimator

$$\widehat{\nabla} f_{\text{quad}}(\mathbf{x}_t) = \frac{f_{\text{quad}}(\mathbf{x}_t + \mu \mathbf{u}_t) - f_{\text{quad}}(\mathbf{x}_t - \mu \mathbf{u}_t)}{2\mu} \mathbf{u}_t, \quad \mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, I_d).$$

Under the quadratic model, this estimator admits the explicit form

$$\widehat{\nabla} f_{\text{quad}}(\mathbf{x}_t) = (\mathbf{u}_t \mathbf{u}_t^\top) \mathbf{H} \mathbf{x}_t,$$

and is unbiased, i.e., $\mathbb{E}[\widehat{\nabla} f_{\text{quad}}(\mathbf{x}_t)] = \mathbf{H} \mathbf{x}_t$.

Let $\mathbf{H} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ with $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d)$ and orthogonal \mathbf{U} . Defining rotated variables $\bar{\mathbf{x}}_t = \mathbf{U}^\top \mathbf{x}_t$, $\bar{\mathbf{m}}_t = \mathbf{U}^\top \mathbf{m}_t$, and $\bar{\mathbf{u}}_t = \mathbf{U}^\top \mathbf{u}_t$, rotational invariance implies $\bar{\mathbf{u}}_t \sim \mathcal{N}(\mathbf{0}, I_d)$ i.i.d., and the dynamics in the rotated coordinates take the same form with \mathbf{H} replaced by $\mathbf{\Lambda}$. Hence, without loss of generality, we assume $\mathbf{H} = \mathbf{\Lambda}$ in what follows.

We now state a lemma that reduces the ZO-GDM second-moment dynamics to a linear covariance operator, which will be the basis for the spectral analysis in subsequent subsections.

Lemma 1 (Second-moment recursion and covariance operator for ZO-GDM). *Consider ZO-GDM applied to the quadratic model $f_{\text{quad}}(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{\Lambda} \mathbf{x}$ with $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d) \neq \mathbf{0}$, step size $\eta > 0$, and momentum $\beta \in [0, 1)$. For each iteration t and coordinate $i = 1, \dots, d$, define the 2×2 covariance block*

$$\mathbf{W}_{i,t} := \begin{bmatrix} \mathbb{E}[x_{i,t}^2] & \mathbb{E}[\eta x_{i,t} m_{i,t}] \\ \mathbb{E}[\eta x_{i,t} m_{i,t}] & \mathbb{E}[\eta^2 m_{i,t}^2] \end{bmatrix},$$

where $\mathbf{x}_t = (x_{1,t}, \dots, x_{d,t}) \in \mathbb{R}^d$ and $\mathbf{m}_t = (m_{1,t}, \dots, m_{d,t}) \in \mathbb{R}^d$. Then the covariance blocks satisfy the recursion

$$\mathbf{W}_{i,t+1} = \mathbf{A}_i \mathbf{W}_{i,t} \mathbf{A}_i^\top + \eta^2 \left(\lambda_i^2 (\mathbf{W}_{i,t})_{11} + \sum_{j=1}^d \lambda_j^2 (\mathbf{W}_{j,t})_{11} \right) \mathbf{Q},$$

where

$$\mathbf{A}_i := \begin{bmatrix} 1 - \eta \lambda_i & -\beta \\ \eta \lambda_i & \beta \end{bmatrix}, \quad \mathbf{Q} := \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Proof of Lemma 1. We consider with the augmented state $(\mathbf{x}_t, \eta \mathbf{m}_t) \in \mathbb{R}^{2d}$, for which the ZO-GDM update can be written as

$$\begin{bmatrix} \mathbf{x}_{t+1} \\ \eta \mathbf{m}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \eta \mathbf{u}_t \mathbf{u}_t^\top \mathbf{\Lambda} & -\beta \mathbf{I} \\ \eta \mathbf{u}_t \mathbf{u}_t^\top \mathbf{\Lambda} & \beta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x}_t \\ \eta \mathbf{m}_t \end{bmatrix}.$$

To track second moments, define

$$\mathbf{X}_t := \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top], \quad \mathbf{Y}_t := \mathbb{E}[\eta^2 \mathbf{m}_t \mathbf{m}_t^\top], \quad \mathbf{C}_t := \mathbb{E}[\eta \mathbf{x}_t \mathbf{m}_t^\top].$$

Using independence of \mathbf{u}_t from $(\mathbf{x}_t, \mathbf{m}_t)$ and applying Lemma 4 (Isserlis's theorem) to evaluate the Gaussian fourth moments, we obtain the following closed recursion:

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta(\mathbf{\Lambda} \mathbf{X}_t + \mathbf{X}_t \mathbf{\Lambda}) + \eta^2(2\mathbf{\Lambda} \mathbf{X}_t \mathbf{\Lambda} + \text{Tr}(\mathbf{\Lambda} \mathbf{X}_t \mathbf{\Lambda}) \mathbf{I}) - \beta(\mathbf{C}_t + \mathbf{C}_t^\top) + \beta \eta(\mathbf{\Lambda} \mathbf{C}_t + \mathbf{C}_t^\top \mathbf{\Lambda}) + \beta^2 \mathbf{Y}_t,$$

$$\mathbf{Y}_{t+1} = \eta^2(2\mathbf{\Lambda} \mathbf{X}_t \mathbf{\Lambda} + \text{Tr}(\mathbf{\Lambda} \mathbf{X}_t \mathbf{\Lambda}) \mathbf{I}) + \eta \beta(\mathbf{\Lambda} \mathbf{C}_t + \mathbf{C}_t^\top \mathbf{\Lambda}) + \beta^2 \mathbf{Y}_t,$$

$$\mathbf{C}_{t+1} = \eta \mathbf{X}_t \mathbf{\Lambda} - \eta^2(2\mathbf{\Lambda} \mathbf{X}_t \mathbf{\Lambda} + \text{Tr}(\mathbf{\Lambda} \mathbf{X}_t \mathbf{\Lambda}) \mathbf{I}) + \beta \mathbf{C}_t - \beta \eta(\mathbf{\Lambda} \mathbf{C}_t + \mathbf{C}_t^\top \mathbf{\Lambda}) - \beta^2 \mathbf{Y}_t.$$

Since $\mathbf{\Lambda}$ is diagonal, the coordinates decouple except through the scalar coupling term

$$\text{Tr}(\mathbf{\Lambda} \mathbf{X}_t \mathbf{\Lambda}) = \sum_{j=1}^d \lambda_j^2 (\mathbf{X}_t)_{jj}.$$

Taking diagonal entries, for each coordinate $i = 1, \dots, d$, we obtain

$$(\mathbf{X}_{t+1})_{ii} = (1 - 2\eta \lambda_i + 2\eta^2 \lambda_i^2) (\mathbf{X}_t)_{ii} + \beta^2 (\mathbf{Y}_t)_{ii} + 2\beta(-1 + \eta \lambda_i) (\mathbf{C}_t)_{ii} + \eta^2 \sum_{j=1}^d \lambda_j^2 (\mathbf{X}_t)_{jj},$$

$$(\mathbf{Y}_{t+1})_{ii} = 2\eta^2 \lambda_i^2 (\mathbf{X}_t)_{ii} + \beta^2 (\mathbf{Y}_t)_{ii} + 2\beta \eta \lambda_i (\mathbf{C}_t)_{ii} + \eta^2 \sum_{j=1}^d \lambda_j^2 (\mathbf{X}_t)_{jj},$$

$$(\mathbf{C}_{t+1})_{ii} = (\eta \lambda_i - 2\eta^2 \lambda_i^2) (\mathbf{X}_t)_{ii} - \beta^2 (\mathbf{Y}_t)_{ii} + \beta(1 - 2\eta \lambda_i) (\mathbf{C}_t)_{ii} - \eta^2 \sum_{j=1}^d \lambda_j^2 (\mathbf{X}_t)_{jj}.$$

Recalling that

$$(\mathbf{X}_t)_{ii} = \mathbb{E}[x_{i,t}^2], \quad (\mathbf{Y}_t)_{ii} = \mathbb{E}[\eta^2 m_{i,t}^2], \quad (\mathbf{C}_t)_{ii} = \mathbb{E}[\eta x_{i,t} m_{i,t}],$$

the above relations can be grouped into the 2×2 covariance block $\mathbf{W}_{i,t}$ defined in the statement of the lemma. A direct calculation then gives

$$\mathbf{W}_{i,t+1} = \mathbf{A}_i \mathbf{W}_{i,t} \mathbf{A}_i^\top + \eta^2 \left(\lambda_i^2 (\mathbf{W}_{i,t})_{11} + \sum_{j=1}^d \lambda_j^2 (\mathbf{W}_{j,t})_{11} \right) \mathbf{Q},$$

with

$$\mathbf{A}_i := \begin{bmatrix} 1 - \eta \lambda_i & -\beta \\ \eta \lambda_i & \beta \end{bmatrix}, \quad \mathbf{Q} := \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

This completes the proof. \square

Lemma 1 shows that the ZO-GDM second-moment dynamics are governed by a closed linear recursion over the covariance blocks $\{\mathbf{W}_{i,t}\}_{i=1}^d$. In particular,

$$\mathbb{E}\|\mathbf{x}_t\|^2 = \sum_{i=1}^d (\mathbf{W}_{i,t})_{11},$$

so mean-square linear stability is equivalent to uniform boundedness of the quantity $\sum_{i=1}^d (\mathbf{W}_{i,t})_{11}$ over iterations.

Collecting the blocks into the product space of 2×2 symmetric matrices $(\mathbb{S}^2)^d$, the mapping $\{\mathbf{W}_{i,t}\}_{i=1}^d \mapsto \{\mathbf{W}_{i,t+1}\}_{i=1}^d$ can be viewed as the action of a linear operator $\mathcal{T} : (\mathbb{S}^2)^d \rightarrow (\mathbb{S}^2)^d$, parameterized by the step size η , momentum β , and the eigenvalues $\lambda_1, \dots, \lambda_d$. In the next subsection, we analyze the spectral properties of this covariance operator and derive explicit conditions under which its spectral radius is smaller than one.

F.2 SPECTRAL ANALYSIS OF THE COVARIANCE OPERATOR

In this subsection, we analyze the spectral properties of the covariance operator induced by the ZO-GDM second-moment recursion derived in Lemma 1 (Appendix F.1). That recursion is governed by a linear operator $\mathcal{T} : (\mathbb{S}^2)^d \rightarrow (\mathbb{S}^2)^d$ acting on the covariance blocks $\{\mathbf{W}_{i,t}\}_{i=1}^d$, and in particular controls $\mathbb{E}\|\mathbf{x}_t\|^2 = \sum_{i=1}^d (\mathbf{W}_{i,t})_{11}$. We therefore analyze the spectral radius of \mathcal{T} .

The operator \mathcal{T} preserves the cone $(\mathbb{S}_+^2)^d$ and contains a rank-one coupling term across coordinates, allowing us to invoke the Krein–Rutman theorem. The following theorem gives an exact characterization of the spectral radius of \mathcal{T} , which forms the technical core of the mean-square stability analysis for ZO-GDM.

Theorem 4 (Spectral characterization of ZO-GDM covariance operator). *Let $\mathcal{X} := (\mathbb{S}^2)^d$ and $\mathcal{K} := (\mathbb{S}_+^2)^d$. Fix $\eta > 0$, $\beta \in [0, 1)$, and $\lambda_{\max} = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$. Define the linear operator $\mathcal{T} : \mathcal{X} \rightarrow \mathcal{X}$ blockwise by*

$$(\mathcal{T}(\mathbf{W}))_i = \mathbf{A}_i \mathbf{W}_i \mathbf{A}_i^\top + \eta^2 \left(\lambda_i^2 (\mathbf{W}_i)_{11} + \sum_{j=1}^d \lambda_j^2 (\mathbf{W}_j)_{11} \right) \mathbf{Q}, \quad i = 1, \dots, d,$$

where

$$\mathbf{A}_i := \begin{bmatrix} 1 - \eta\lambda_i & -\beta \\ \eta\lambda_i & \beta \end{bmatrix}, \quad \mathbf{Q} := \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Define the scalar function

$$S(\eta, \beta) := \sum_{i=1}^d \frac{\eta\lambda_i}{2(1-\beta) \left(1 - \frac{\eta\lambda_i}{1-\beta^2}\right)}.$$

Then the operator \mathcal{T} satisfies the following properties.

- (a) (**Leading eigenvalue in the cone**) The operator \mathcal{T} preserves the cone \mathcal{K} , i.e., $\mathcal{T}(\mathcal{K}) \subseteq \mathcal{K}$. Moreover, \mathcal{T} has an eigenvalue equal to its spectral radius $\rho(\mathcal{T})$, with an associated eigenvector $\mathbf{W}^* = (\mathbf{W}_1^*, \dots, \mathbf{W}_d^*) \in \mathcal{K} \setminus \{\mathbf{0}\}$ satisfying

$$\rho(\mathcal{T}) \geq \eta^2 \sum_{i=1}^d \lambda_i^2 > 0 \quad \text{and} \quad \sum_{i=1}^d (\mathbf{W}_i^*)_{11} > 0.$$

- (b) (**Critical case**)

$$\rho(\mathcal{T}) = 1 \iff \eta\lambda_{\max} < 1 - \beta^2 \text{ and } S(\eta, \beta) = 1.$$

- (c) (**Subcritical case**)

$$\rho(\mathcal{T}) < 1 \iff \eta\lambda_{\max} < 1 - \beta^2 \text{ and } S(\eta, \beta) < 1.$$

Proof of Theorem 4. We work on the product space $\mathcal{X} := (\mathbb{S}^2)^d$ equipped with the product cone $\mathcal{K} := (\mathbb{S}_+^2)^d$ and the induced partial order \leq .

Proof of (a). Let $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_d) \in \mathcal{K}$, i.e., $\mathbf{W}_i \geq 0$ for all $i = 1, \dots, d$. Then $\mathbf{A}_i \mathbf{W}_i \mathbf{A}_i^\top \geq 0$ and $(\mathbf{W}_i)_{11} \geq 0$, so $\sum_{j=1}^d \lambda_j^2 (\mathbf{W}_j)_{11} \geq 0$. Hence $(\mathcal{T}(\mathbf{W}))_i \geq 0$ for all i , and therefore $\mathcal{T}(\mathcal{K}) \subseteq \mathcal{K}$.

Since \mathcal{X} is finite dimensional and \mathcal{K} is a closed, convex, pointed cone with nonempty interior, the Krein–Rutman theorem (Lemma 6) implies that \mathcal{T} has an eigenvalue equal to its spectral radius $\rho(\mathcal{T})$ with a corresponding eigenvector $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$.

To lower bound $\rho(\mathcal{T})$, define $\bar{\mathbf{Q}} := (\mathbf{Q}, \dots, \mathbf{Q}) \in \mathcal{K} \setminus \{\mathbf{0}\}$. For each i ,

$$(\mathcal{T}(\bar{\mathbf{Q}}))_i = \mathbf{A}_i \mathbf{Q} \mathbf{A}_i^\top + \eta^2 \left(\lambda_i^2 + \sum_{j=1}^d \lambda_j^2 \right) \mathbf{Q} \geq \eta^2 \left(\sum_{j=1}^d \lambda_j^2 \right) \mathbf{Q}.$$

Thus $\mathcal{T}(\bar{\mathbf{Q}}) \geq c \bar{\mathbf{Q}}$ with $c := \eta^2 \sum_{j=1}^d \lambda_j^2 > 0$. Iterating gives $\mathcal{T}^t(\bar{\mathbf{Q}}) \geq c^t \bar{\mathbf{Q}}$ for all $t \geq 1$. Fixing any norm on \mathcal{X} and applying Gelfand’s formula, we obtain

$$\rho(\mathcal{T}) \geq c = \eta^2 \sum_{j=1}^d \lambda_j^2.$$

Finally, we show $\sum_i (\mathbf{W}_i^*)_{11} > 0$. If $(\mathbf{W}_i^*)_{11} = 0$ for all i , then $\mathbf{W}_i^* = \begin{bmatrix} 0 & 0 \\ 0 & y_i \end{bmatrix}$ with $y_i \geq 0$. But then

$$\rho(\mathcal{T}) \mathbf{W}_i^* = (\mathcal{T}(\mathbf{W}^*))_i = \mathbf{A}_i \mathbf{W}_i^* \mathbf{A}_i^\top = \beta^2 y_i \mathbf{Q}.$$

Note that $\beta^2 y_i \mathbf{Q}$ has the same value at the (1, 1)-entry and (2, 2)-entry. Consequently, $\rho(\mathcal{T}) \mathbf{W}_i^*$ should have the same value at the (1, 1)-entry and (2, 2)-entry. Thus $\rho(\mathcal{T}) y_i = 0$, forcing $y_i = 0$ since $\rho(\mathcal{T}) > 0$. Hence $\mathbf{W}^* = \mathbf{0}$, a contradiction. This concludes the proof of (a).

Local-global decomposition. For each $i = 1, \dots, d$, we define the local linear map $\mathcal{M}_i : \mathbb{S}^2 \rightarrow \mathbb{S}^2$ by

$$\mathcal{M}_i(\mathbf{X}) := \mathbf{A}_i \mathbf{X} \mathbf{A}_i^\top + \eta^2 \lambda_i^2 (\mathbf{X})_{11} \mathbf{Q}.$$

Define the global coupling scalar function $s : \mathcal{X} \rightarrow \mathbb{R}$ by

$$s(\mathbf{W}) := \sum_{j=1}^d \eta^2 \lambda_j^2 (\mathbf{W}_j)_{11}.$$

Then, the linear operator $\mathcal{T} : \mathcal{X} \rightarrow \mathcal{X}$ can be decomposed as below:

$$(\mathcal{T}(\mathbf{W}))_i = \mathcal{M}_i(\mathbf{W}_i) + s(\mathbf{W}) \mathbf{Q}, \quad i = 1, \dots, d. \quad (7)$$

Note that each linear map \mathcal{M}_i is \mathbb{S}_+^2 -preserving, i.e., $\mathcal{M}_i(\mathbb{S}_+^2) \subseteq \mathbb{S}_+^2$. Define the block-diagonal operator $\mathcal{M} : \mathcal{X} \rightarrow \mathcal{X}$ and $\bar{\mathbf{Q}} \in \mathcal{K}$ by

$$(\mathcal{M}(\mathbf{W}))_i := \mathcal{M}_i(\mathbf{W}_i), \quad \text{and} \quad \bar{\mathbf{Q}} := (\mathbf{Q}, \dots, \mathbf{Q}) \in \mathcal{K}. \quad (8)$$

Then (7) can be written as

$$\mathcal{T}(\mathbf{W}) = \mathcal{M}(\mathbf{W}) + s(\mathbf{W}) \bar{\mathbf{Q}}.$$

Positivity of $s(\mathbf{W}^*)$. Let $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ denote the leading eigenvector of \mathcal{T} satisfying $\mathcal{T}(\mathbf{W}^*) = \rho(\mathcal{T}) \mathbf{W}^*$ and $\sum_{i=1}^d (\mathbf{W}_i^*)_{11} > 0$, which exists by Theorem 4(a). Then, $s(\mathbf{W}^*) = \sum_{j=1}^d \eta^2 \lambda_j^2 (\mathbf{W}_j^*)_{11} > 0$.

Key lemmas. We use the following key lemmas to prove (b) and (c): Lemma 2 and Lemma 3.

Lemma 2 shows that $\rho(\mathcal{M}_i) < 1$ if and only if $\eta \lambda_i < 1 - \beta^2$. Moreover, if $\eta \lambda_i < 1 - \beta^2$, then $\text{Id} - \mathcal{M}_i$ is invertible, $(\text{Id} - \mathcal{M}_i)^{-1}(\mathbb{S}_+^2) \subseteq \mathbb{S}_+^2$, and $\mathbf{Y}_i := (\text{Id} - \mathcal{M}_i)^{-1}(\mathbf{Q}) \in \mathbb{S}_+^2$ satisfies

$$\gamma_i := (\mathbf{Y}_i)_{11} = \frac{1 + \beta}{2\eta \lambda_i (1 - \beta^2 - \eta \lambda_i)}.$$

Lemma 3 shows that if $\eta \lambda_i \geq 1 - \beta^2$, then for every $\alpha > 0$, there does not exist $\mathbf{W} \geq \mathbf{0}$ such that $(\text{Id} - \mathcal{M}_i)(\mathbf{W}) \geq \alpha \mathbf{Q}$.

Proof of (b). We prove

$$\rho(\mathcal{T}) = 1 \iff \eta \lambda_{\max} < 1 - \beta^2 \text{ and } S(\eta, \beta) = 1.$$

(\implies) Assume $\rho(\mathcal{T}) = 1$. By Theorem 4(a), there exists $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ with $\mathcal{T}(\mathbf{W}^*) = \mathbf{W}^*$. Let $s^* := s(\mathbf{W}^*)$. Then, $s^* > 0$ by the positivity of $s(\mathbf{W}^*)$. From (7), for each i ,

$$\mathbf{W}_i^* = (\mathcal{T}(\mathbf{W}^*))_i = \mathcal{M}_i(\mathbf{W}_i^*) + s^* \mathbf{Q},$$

or equivalently,

$$(\text{Id} - \mathcal{M}_i)(\mathbf{W}_i^*) = s^* \mathbf{Q}.$$

918 Since $s^* > 0$, Lemma 3 implies $\rho(\mathcal{M}_i) < 1$ for every i , and by Lemma 2(a),
 919 $\eta\lambda_i < 1 - \beta^2$ for all i , hence $\eta\lambda_{\max} < 1 - \beta^2$.

920 By Lemma 2(b), $\text{Id} - \mathcal{M}_i$ is invertible and $(\text{Id} - \mathcal{M}_i)^{-1}$ is \mathbb{S}_+^2 -preserving, so

$$921 \mathbf{W}_i^* = s^*(\text{Id} - \mathcal{M}_i)^{-1}\mathbf{Q} = s^*\mathbf{Y}_i,$$

922 where $\mathbf{Y}_i \geq 0$ and $\gamma_i := (\mathbf{Y}_i)_{11}$ is given by Lemma 2(c). Taking $(1, 1)$ -entries and plugging into the
 923 definition of s^* , we get

$$924 s^* = \sum_{i=1}^d \eta^2 \lambda_i^2 (\mathbf{W}_i^*)_{11} = \sum_{i=1}^d \eta^2 \lambda_i^2 \gamma_i s^*.$$

925 Cancelling $s^* > 0$ gives

$$926 1 = \sum_{i=1}^d \eta^2 \lambda_i^2 \gamma_i.$$

927 Using Lemma 2(c),

$$928 \eta^2 \lambda_i^2 \gamma_i = \eta^2 \lambda_i^2 \cdot \frac{1 + \beta}{2\eta\lambda_i(1 - \beta^2 - \eta\lambda_i)} = \frac{\eta\lambda_i}{2(1 - \beta) \left(1 - \frac{\eta\lambda_i}{1 - \beta^2}\right)}.$$

929 Therefore $\sum_{i=1}^d \eta^2 \lambda_i^2 \gamma_i = S(\eta, \beta)$, and thus we conclude that $S(\eta, \beta) = 1$.

930 (\Leftarrow) Assume $\eta\lambda_{\max} < 1 - \beta^2$ and $S(\eta, \beta) = 1$. Then, $\eta\lambda_i < 1 - \beta^2$ for all i .

931 By Lemma 2(b), the matrices $\mathbf{Y}_i := (\text{Id} - \mathcal{M}_i)^{-1}\mathbf{Q}$ are well-defined and satisfy $\mathbf{Y}_i \geq 0$. Set
 932 $\mathbf{Y} := (\mathbf{Y}_1, \dots, \mathbf{Y}_d) \in \mathcal{K} \setminus \{\mathbf{0}\}$. For each i , $(\text{Id} - \mathcal{M}_i)(\mathbf{Y}_i) = \mathbf{Q}$, i.e. $\mathcal{M}_i(\mathbf{Y}_i) = \mathbf{Y}_i - \mathbf{Q}$. Moreover,

$$933 s(\mathbf{Y}) = \sum_{i=1}^d \eta^2 \lambda_i^2 (\mathbf{Y}_i)_{11} = \sum_{i=1}^d \eta^2 \lambda_i^2 \gamma_i = S(\eta, \beta) = 1.$$

934 Thus (7) gives, for each i ,

$$935 (\mathcal{T}(\mathbf{Y}))_i = \mathcal{M}_i(\mathbf{Y}_i) + s(\mathbf{Y})\mathbf{Q} = (\mathbf{Y}_i - \mathbf{Q}) + \mathbf{Q} = \mathbf{Y}_i,$$

936 so $\mathcal{T}(\mathbf{Y}) = \mathbf{Y}$. Hence, 1 is an eigenvalue of \mathcal{T} with eigenvector \mathbf{Y} , and therefore $\rho(\mathcal{T}) \geq 1$. Set
 937 $r := \rho(\mathcal{T}) \geq 1$.

938 By Theorem 4(a), there exists $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ with $\mathcal{T}(\mathbf{W}^*) = \rho(\mathcal{T})\mathbf{W}^*$. Set $s^* := s(\mathbf{W}^*)$. By the
 939 positivity of $s(\mathbf{W}^*)$, we have $s^* > 0$. From $\mathcal{T}(\mathbf{W}^*) = r\mathbf{W}^*$ and the decomposition (8), we obtain

$$940 r\mathbf{W}^* = \mathcal{M}(\mathbf{W}^*) + s^*\bar{\mathbf{Q}}, \quad \text{so that} \quad \left(\text{Id} - \frac{1}{r}\mathcal{M}\right)\mathbf{W}^* = \frac{s^*}{r}\bar{\mathbf{Q}}.$$

941 Since $\eta\lambda_{\max} < 1 - \beta^2$, Lemma 2(a) implies that $\rho(\mathcal{M}_i) < 1$ for all i . Consequently, $\rho(\mathcal{M}) =$
 942 $\max_i \rho(\mathcal{M}_i) < 1$. Since $r \geq 1$, we have $\rho(\frac{1}{r}\mathcal{M}) \leq \rho(\mathcal{M}) < 1$ and hence $\text{Id} - \frac{1}{r}\mathcal{M}$ is invertible
 943 with

$$944 \left(\text{Id} - \frac{1}{r}\mathcal{M}\right)^{-1} = \sum_{k=0}^{\infty} \left(\frac{1}{r}\mathcal{M}\right)^k,$$

945 where the series converges in operator norm. Applying this inverse gives

$$946 \mathbf{W}^* = \frac{s^*}{r} \left(\text{Id} - \frac{1}{r}\mathcal{M}\right)^{-1} \bar{\mathbf{Q}}.$$

947 Applying the nonnegative functional $s(\cdot)$ to both sides and cancelling $s^* > 0$, we obtain

$$948 r = s\left(\left(\text{Id} - \frac{1}{r}\mathcal{M}\right)^{-1} \bar{\mathbf{Q}}\right). \quad (9)$$

949 Using the Neumann series, we have

$$950 s\left(\left(\text{Id} - \frac{1}{r}\mathcal{M}\right)^{-1} \bar{\mathbf{Q}}\right) = \sum_{k=0}^{\infty} r^{-k} s(\mathcal{M}^k(\bar{\mathbf{Q}})) \leq \sum_{k=0}^{\infty} s(\mathcal{M}^k(\bar{\mathbf{Q}})) = s\left((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}}\right).$$

951 Because \mathcal{M} is block-diagonal, $(\text{Id} - \mathcal{M})^{-1}$ is block-diagonal with blocks $(\text{Id} - \mathcal{M}_i)^{-1}$, and thus

$$952 (\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}} = ((\text{Id} - \mathcal{M}_1)^{-1}\mathbf{Q}, \dots, (\text{Id} - \mathcal{M}_d)^{-1}\mathbf{Q}) = (\mathbf{Y}_1, \dots, \mathbf{Y}_d) = \mathbf{Y}.$$

953 Therefore,

$$954 s\left((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}}\right) = s(\mathbf{Y}) = S(\eta, \beta) = 1.$$

955 Combining with (9) shows that $r \leq 1$, and hence $r = 1$. Therefore, we conclude $\rho(\mathcal{T}) = 1$. This
 956 finishes the proof of (b).

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

Proof of (c). We prove

$$\rho(\mathcal{T}) < 1 \iff \eta\lambda_{\max} < 1 - \beta^2 \text{ and } S(\eta, \beta) < 1.$$

(\implies) Assume $\rho(\mathcal{T}) < 1$. Set $r := \rho(\mathcal{T}) < 1$. By Theorem 4(a), there exists $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ with $\mathcal{T}(\mathbf{W}^*) = r\mathbf{W}^*$. Let $s^* := s(\mathbf{W}^*)$. Then, $s^* > 0$ by the positivity of $s(\mathbf{W}^*)$. From (7), for each i ,

$$r\mathbf{W}_i^* = (\mathcal{T}(\mathbf{W}^*))_i = \mathcal{M}_i(\mathbf{W}_i^*) + s^*\mathbf{Q},$$

or equivalently,

$$(\text{Id} - \frac{1}{r}\mathcal{M}_i)(\mathbf{W}_i^*) = \frac{s^*}{r}\mathbf{Q}.$$

Since $\frac{s^*}{r} > 0$ and $\mathbf{W}_i^* \geq 0$, Lemma 3 applied to the map $\frac{1}{r}\mathcal{M}_i$ implies $\rho(\frac{1}{r}\mathcal{M}_i) < 1$, and thus $\rho(\mathcal{M}_i) < r < 1$. By Lemma 2(a),

$$\eta\lambda_i < 1 - \beta^2 \text{ for all } i, \quad \text{hence} \quad \eta\lambda_{\max} < 1 - \beta^2.$$

Moreover, since $\eta\lambda_i < 1 - \beta^2$, Lemma 2(b) implies $(\text{Id} - \mathcal{M}_i)^{-1}$ exists and is \mathbb{S}_+^2 -preserving, so

$$\mathbf{W}_i^* = \frac{s^*}{r} \left(\text{Id} - \frac{1}{r}\mathcal{M}_i \right)^{-1} \mathbf{Q} \geq 0.$$

Applying $s(\cdot)$ to the identity $\mathbf{W}^* = \frac{s^*}{r}(\text{Id} - \frac{1}{r}\mathcal{M})^{-1}\bar{\mathbf{Q}}$ gives, after cancelling $s^* > 0$,

$$r = s \left(\left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}} \right). \quad (10)$$

Using the Neumann series and $r < 1$, we have

$$s \left(\left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}} \right) = \sum_{k=0}^{\infty} r^{-k} s(\mathcal{M}^k(\bar{\mathbf{Q}})) \geq \sum_{k=0}^{\infty} s(\mathcal{M}^k(\bar{\mathbf{Q}})) = s \left((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}} \right).$$

Then from (10), we obtain

$$r \geq s \left((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}} \right) = s(\mathbf{Y}) = S(\eta, \beta).$$

Since $r < 1$, we conclude that $S(\eta, \beta) < 1$.

(\impliedby) Assume $\eta\lambda_{\max} < 1 - \beta^2$ and $S(\eta, \beta) < 1$. We prove $\rho(\mathcal{T}) < 1$ by contradiction. Assume $\rho(\mathcal{T}) \geq 1$, and let $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ satisfy $\mathcal{T}(\mathbf{W}^*) = r\mathbf{W}^*$ with $r := \rho(\mathcal{T}) \geq 1$. Set $s^* := s(\mathbf{W}^*) > 0$. Recall from (8) that

$$r\mathbf{W}^* = \mathcal{M}(\mathbf{W}^*) + s^*\bar{\mathbf{Q}}, \quad \text{so} \quad \left(\text{Id} - \frac{1}{r}\mathcal{M} \right) \mathbf{W}^* = \frac{s^*}{r}\bar{\mathbf{Q}}.$$

Since $\eta\lambda_{\max} < 1 - \beta^2$, Lemma 2(a) gives $\rho(\mathcal{M}_i) < 1$ for all i , hence $\rho(\mathcal{M}) < 1$. Since $r \geq 1$, we have $\rho(\frac{1}{r}\mathcal{M}) \leq \rho(\mathcal{M}) < 1$, so $\text{Id} - \frac{1}{r}\mathcal{M}$ is invertible and

$$\mathbf{W}^* = \frac{s^*}{r} \left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}}.$$

Applying $s(\cdot)$ and cancelling $s^* > 0$, we obtain

$$r = s \left(\left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}} \right). \quad (11)$$

Using the Neumann series, we have

$$s \left(\left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}} \right) = \sum_{k=0}^{\infty} r^{-k} s(\mathcal{M}^k(\bar{\mathbf{Q}})) \leq \sum_{k=0}^{\infty} s(\mathcal{M}^k(\bar{\mathbf{Q}})) = s \left((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}} \right).$$

As in (b), $(\text{Id} - \mathcal{M})^{-1}\bar{\mathbf{Q}} = \mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_d)$, hence

$$s \left((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}} \right) = s(\mathbf{Y}) = S(\eta, \beta).$$

Combining with (11) yields

$$r \leq S(\eta, \beta) < 1,$$

contradicting $r \geq 1$. Therefore $\rho(\mathcal{T}) < 1$. This completes the proof of (c). \square

Lemma 2. For each $i = 1, \dots, d$, define $\mathcal{M}_i : \mathbb{S}^2 \rightarrow \mathbb{S}^2$ by

$$\mathcal{M}_i(\mathbf{X}) := \mathbf{A}_i \mathbf{X} \mathbf{A}_i^\top + \eta^2 \lambda_i^2(\mathbf{X})_{11} \mathbf{Q},$$

where \mathbf{A}_i and \mathbf{Q} are defined in Theorem 4, and $\eta\lambda_i > 0$. It holds that

(a) $\rho(\mathcal{M}_i) < 1$ if and only if $\eta\lambda_i < 1 - \beta^2$

(b) If $\eta\lambda_i < 1 - \beta^2$, then $\text{Id} - \mathcal{M}_i$ is invertible and $(\text{Id} - \mathcal{M}_i)^{-1}$ is \mathbb{S}_+^2 -preserving.

1026 (c) If $\eta\lambda_i < 1 - \beta^2$, define $\mathbf{Y}_i := (\text{Id} - \mathcal{M}_i)^{-1}\mathbf{Q} \in \mathbb{S}_+^2$ and $\gamma_i := (\mathbf{Y}_i)_{11}$. Then,

$$1027 \quad \gamma_i = \frac{1 + \beta}{2\eta\lambda_i(1 - \beta^2 - \eta\lambda_i)}.$$

1028

1029

1030

1031

1032

Proof of Lemma 2. Let $\Phi : \mathbb{S}^2 \rightarrow \mathbb{R}^3$ be the isomorphism defined by

1033

1034

1035

$$\Phi \left(\begin{bmatrix} x_{11} & x_{12} \\ x_{12} & x_{22} \end{bmatrix} \right) = \begin{bmatrix} x_{11} \\ x_{12} \\ x_{22} \end{bmatrix} \in \mathbb{R}^3.$$

1036

Then, for each $i = 1, \dots, d$, it holds that

1037

1038

1039

$$\forall \mathbf{X} \in \mathbb{S}^2 \quad \Phi(\mathcal{M}_i(\mathbf{X})) = \mathbf{K}_i\Phi(\mathbf{X}), \quad \mathbf{K}_i = \begin{bmatrix} 1 - 2\eta\lambda_i + 2\eta^2\lambda_i^2 & 2\beta(-1 + \eta\lambda_i) & \beta^2 \\ \eta\lambda_i - 2\eta^2\lambda_i^2 & \beta(1 - 2\eta\lambda_i) & -\beta^2 \\ 2\eta^2\lambda_i^2 & 2\beta\eta\lambda_i & \beta^2 \end{bmatrix}.$$

1040

1041

Since Φ is an isomorphism, $\rho(\mathcal{M}_i) = \rho(\mathbf{K}_i)$.

1042

1043

Proof of (a): Exact condition for $\rho(\mathbf{K}_i) < 1$. Fix i and write $x := \eta\lambda_i > 0$ and $b := \beta \in [0, 1)$. Under the identification $\Phi : \mathbb{S}^2 \rightarrow \mathbb{R}^3$ in the proof, $\rho(\mathcal{M}_i) = \rho(\mathbf{K}_i)$. Let $p(r) := \det(r\mathbf{I} - \mathbf{K}_i) = r^3 + c_1r^2 + c_2r + c_3$. A direct expansion gives

1044

1045

1046

1047

1048

1049

$$c_1 = -b^2 + 2bx - b - 2x^2 + 2x - 1, \quad (12)$$

$$c_2 = b(b^2 + b + 1 - 2(b+1)x), \quad (13)$$

$$c_3 = -b^3. \quad (14)$$

1050

We use the strict Jury criterion for a monic cubic with real coefficients: all roots of $p(r) = r^3 + c_1r^2 + c_2r + c_3$ lie in $\{|r| < 1\}$ if and only if

1051

1052

1053

$$|c_3| < 1, \quad p(1) > 0, \quad 1 - c_1 + c_2 - c_3 > 0, \quad 1 - c_3^2 > |c_2 - c_1c_3|. \quad (15)$$

1054

We verify (15) and show that they hold if and only if $0 < x < 1 - b^2$.

1055

1056

First, from (14), $|c_3| = b^3 < 1$, so the first condition in (15) holds.

1057

A direct substitution of (12),(13), and (14) gives

1058

$$p(1) = 1 + c_1 + c_2 + c_3 = 2x(1 - b^2 - x).$$

1059

Hence,

1060

$$p(1) > 0 \iff 0 < x < 1 - b^2.$$

1061

Next, substitution of (12),(13), and (14) gives

1062

$$1 - c_1 + c_2 - c_3 = 2\phi(x), \quad \phi(x) := x^2 - (b+1)^2x + (b^3 + b^2 + b + 1).$$

1063

The discriminant of ϕ is

1064

1065

1066

$$\Delta_\phi = (b+1)^4 - 4(b^3 + b^2 + b + 1) = b^4 + 2b^2 - 3 = (b^2 - 1)(b^2 + 3) < 0 \quad (b \in [0, 1)),$$

and ϕ has leading coefficient $1 > 0$, so $\phi(x) > 0$ for all $x \in \mathbb{R}$. Therefore $1 - c_1 + c_2 - c_3 > 0$ holds automatically for all $x \geq 0$ and $b \in [0, 1)$.

1067

It remains to check the last condition in (15): $1 - c_3^2 > |c_2 - c_1c_3|$. Again, by substituting (12),(13), and (14) gives

1068

1069

1070

$$(1 - c_3^2) - (c_2 - c_1c_3) = 2b^3x^2 + 2b(1+b)^2(1-b)x + (1+b)(1-b)^3(1+b+b^2) > 0$$

for any $x \geq 0$ and $b \in [0, 1)$. Moreover, we have

1071

1072

1073

1074

1075

$$(1 - c_3^2) + (c_2 - c_1c_3) = -2b^3x^2 - 2b(1+b)^2(1-b)x + (1+b)(1-b)(1+b^2)(1+b+b^2) := \psi(x).$$

If $b = 0$, then $\psi(x) = 1 > 0$. Otherwise, $b \in (0, 1)$, and ψ is a concave quadratic and has a global maximum at $x < 0$. Hence, if $0 < x < 1 - b^2$, then $\psi(x) > \min\{\psi(0), \psi(1 - b^2)\}$. Note that

$$\psi(0) = (1+b)(1-b)(1+b^2)(1+b+b^2) > 0,$$

1076

$$\psi(1 - b^2) = (1+b)(1-b)(1+b+b^2)[(1-b)^2 + 2b^3] > 0,$$

1077

for all $b \in (0, 1)$. Therefore we conclude that all the strict Jury conditions (15) hold if and only if $0 < x < 1 - b^2$. Putting $x = \eta\lambda_i$ and $b = \beta$, we conclude

1078

$$\rho(\mathcal{M}_i) = \rho(\mathbf{K}_i) < 1 \iff 0 < \eta\lambda_i < 1 - \beta^2,$$

1079

which proves part (a).

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Proof of (b). Assume that $\eta\lambda_i < 1 - \beta^2$. Note that

$$\det(\mathbf{I} - \mathbf{K}_i) = 2\eta\lambda_i(1 - \beta^2 - \eta\lambda_i) \neq 0,$$
so $\mathbf{I} - \mathbf{K}_i$ is invertible, and hence $\text{Id} - \mathcal{M}_i$ is also invertible. According to (a), we have $\rho(\mathcal{M}_i) < 1$, and

$$(\text{Id} - \mathcal{M}_i)^{-1} = \sum_{k=0}^{\infty} \mathcal{M}_i^k$$

converges in operator norm (finite-dimensional Neumann series). Since $\mathcal{M}_i(\mathbb{S}_+^2) \subseteq \mathbb{S}_+^2$, we have $\mathcal{M}_i^k(\mathbb{S}_+^2) \subseteq \mathbb{S}_+^2$ for any $k \geq 0$, and thus

$$(\text{Id} - \mathcal{M}_i)^{-1}(\mathbb{S}_+^2) \subseteq \mathbb{S}_+^2.$$

This concludes the proof of (b).

Proof of (c). Assume that $\eta\lambda_i < 1 - \beta^2$ and define $\mathbf{Y}_i := (\text{Id} - \mathcal{M}_i)^{-1}\mathbf{Q} \in \mathbb{S}_+^2$. In \mathbb{R}^3 , we have

$$\begin{bmatrix} (\mathbf{Y}_i)_{11} \\ (\mathbf{Y}_i)_{12} \\ (\mathbf{Y}_i)_{22} \end{bmatrix} = \Phi(\mathbf{Y}_i) = (\mathbf{I} - \mathbf{K}_i)^{-1}\Phi(\mathbf{Q}) = \begin{bmatrix} 2\eta\lambda_i - 2\eta^2\lambda_i^2 & 2\beta(1 - \eta\lambda_i) & -\beta^2 \\ -\eta\lambda_i + 2\eta^2\lambda_i^2 & 1 - \beta(1 - 2\eta\lambda_i) & \beta^2 \\ -2\eta^2\lambda_i^2 & -2\beta\eta\lambda_i & 1 - \beta^2 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}.$$

Solving this 3×3 system and extracting the first coordinate gives

$$\gamma_i := (\mathbf{Y}_i)_{11} = \frac{1 + \beta}{2\eta\lambda_i(1 - \beta^2 - \eta\lambda_i)}.$$

This concludes the proof of (c). \square

Lemma 3. Let \mathcal{M}_i denote a linear map defined as Lemma 2. If $\rho(\mathcal{M}_i) \geq 1$, then for every $\alpha > 0$, there does not exist $\mathbf{W} \geq 0$ such that

$$(\text{Id} - \mathcal{M}_i)(\mathbf{W}) \geq \alpha\mathbf{Q},$$

where \mathbf{Q} is defined in Theorem 4.

Proof of Lemma 3. Let $\mathcal{M}_i^* : \mathbb{S}^2 \rightarrow \mathbb{S}^2$ denote the adjoint of \mathcal{M}_i . Since \mathcal{M}_i^* is a linear operator with $\mathcal{M}_i^*(\mathbb{S}_+^2) \subseteq \mathbb{S}_+^2$ and $\mathbb{S}_+^2 \subset \mathbb{S}^2$ is a closed, convex, pointed cone with nonempty interior, the Krein–Rutman Theorem (see also Lemma 6) implies that there exists $\mathbf{Y} \geq 0$ with $\mathbf{Y} \neq \mathbf{0}$ such that

$$\mathcal{M}_i^*(\mathbf{Y}) = \rho(\mathcal{M}_i)\mathbf{Y}.$$

Define $\mathcal{L}_i := \text{Id} - \mathcal{M}_i$. Then, its adjoint $\mathcal{L}_i^* = \text{Id} - \mathcal{M}_i^*$ satisfies

$$\mathcal{L}_i^*(\mathbf{Y}) = (1 - \rho(\mathcal{M}_i))\mathbf{Y} \leq 0,$$

since $\rho(\mathcal{M}_i) \geq 1$. Since $\mathbf{Y} \geq 0$ and $\mathbf{Q} \geq 0$, we have $\langle \mathbf{Y}, \mathbf{Q} \rangle \geq 0$.

Now, we show that $\langle \mathbf{Y}, \mathbf{Q} \rangle > 0$. Assume the contrary that $\langle \mathbf{Y}, \mathbf{Q} \rangle = 0$. Then,

$$\langle \mathbf{Y}, \mathbf{Q} \rangle = \text{Tr}(\mathbf{Y}\mathbf{Q}) = \text{Tr}(\mathbf{Y}^{1/2}\mathbf{Q}\mathbf{Y}^{1/2}) = 0.$$

Since $\mathbf{Y}^{1/2}\mathbf{Q}\mathbf{Y}^{1/2} \geq 0$ and $\text{Tr}(\mathbf{Y}^{1/2}\mathbf{Q}\mathbf{Y}^{1/2}) = 0$, we have $\mathbf{Y}^{1/2}\mathbf{Q}\mathbf{Y}^{1/2} = \mathbf{0}$ and thus $\mathbf{Q}\mathbf{Y} = \mathbf{0}$.

Note that $\mathbf{Q} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}^\top \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}$ has null space $\text{Null}(\mathbf{Q}) = \text{span}(\begin{bmatrix} 1 & 1 \end{bmatrix}^\top)$. Since $\mathbf{Q}\mathbf{Y} = \mathbf{0}$ and $\mathbf{Y} \neq \mathbf{0}$, we have

$$\text{Range}(\mathbf{Y}) \subseteq \text{Null}(\mathbf{Q}) = \text{span}\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right), \quad \text{and hence} \quad \mathbf{Y} = \alpha \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} = \alpha \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \text{ for some } \alpha > 0.$$

Then, it holds that

$$\begin{aligned} \mathcal{M}_i^*(\mathbf{Y}) &= \mathbf{A}_i^\top \mathbf{Y} \mathbf{A}_i + \eta^2 \lambda_i^2 \langle \mathbf{Y}, \mathbf{Q} \rangle \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \\ &= \mathbf{A}_i^\top \mathbf{Y} \mathbf{A}_i \\ &= \alpha \left(\mathbf{A}_i^\top \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) \left(\mathbf{A}_i^\top \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^\top \\ &= \alpha \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \end{aligned}$$

However,

$$\alpha \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} = \mathcal{M}_i^*(\mathbf{Y}) = \rho(\mathcal{M}_i)\mathbf{Y} = \rho(\mathcal{M}_i)\alpha \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

gives a contradiction. Hence, $\langle \mathbf{Y}, \mathbf{Q} \rangle > 0$.

1134 Finally, assume for contradiction that there exists $\mathbf{W} \geq 0$ such that

$$1135 (\text{Id} - \mathcal{M}_i)(\mathbf{W}) \geq \alpha \mathbf{Q}.$$

1136 Then,

$$1137 \langle \mathbf{Y}, (\text{Id} - \mathcal{M}_i)(\mathbf{W}) \rangle \geq \alpha \langle \mathbf{Y}, \mathbf{Q} \rangle > 0.$$

1138 However, since $\mathbf{W} \geq 0$ and $\mathcal{L}_i^*(\mathbf{Y}) \leq 0$, it holds that

$$1139 \langle \mathbf{Y}, (\text{Id} - \mathcal{M}_i)(\mathbf{W}) \rangle = \langle (\text{Id} - \mathcal{M}_i^*)(\mathbf{Y}), \mathbf{W} \rangle = \langle \mathcal{L}_i^*(\mathbf{Y}), \mathbf{W} \rangle \leq 0,$$

1140 a contradiction. This concludes the proof of Lemma 3. \square

1141

1142

1143 F.3 PROOF OF THEOREM 2

1144 We prove Theorem 2 by characterizing mean-square linear stability of the linearized ZO-GDM
1145 dynamics. Specifically, if

$$1146 \eta \lambda_{\max}(\mathbf{H}) < 1 - \beta^2 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \lambda_i}{2(1 - \beta) \left(1 - \frac{\eta \lambda_i}{1 - \beta^2}\right)} < 1, \quad (16)$$

1147 then the linearized dynamics is mean-square stable. Conversely, mean-square linear stability implies

$$1148 \eta \lambda_{\max}(\mathbf{H}) < 1 - \beta^2 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \lambda_i}{2(1 - \beta) \left(1 - \frac{\eta \lambda_i}{1 - \beta^2}\right)} \leq 1. \quad (17)$$

1149 By definition, the mean-square critical step size η_{ms}^* is therefore the unique $\eta > 0$ satisfying

$$1150 \eta \lambda_{\max}(\mathbf{H}) < 1 - \beta^2 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \lambda_i}{2(1 - \beta) \left(1 - \frac{\eta \lambda_i}{1 - \beta^2}\right)} = 1. \quad (18)$$

1151 Moreover, η_{ms}^* obeys the bounds

$$1152 \frac{2(1 - \beta)}{\text{Tr}(\mathbf{H}) + \frac{2\lambda_{\max}(\mathbf{H})}{1 + \beta}} \leq \eta_{\text{ms}}^* \leq \frac{2(1 - \beta)}{\text{Tr}(\mathbf{H})}. \quad (19)$$

1153

1154

1155 *Proof of Theorem 2.* Without loss of generality, we assume $\mathbf{x}^* = \mathbf{0}$ and diagonalize the Hessian as
1156 $\mathbf{H} = \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d)$ with $\lambda_1 \geq \dots \geq \lambda_d \geq 0$.

1157 By Lemma 1, the covariance blocks evolve according to $\mathcal{T}(\mathbf{W}_{1,t}, \dots, \mathbf{W}_{d,t}) =$
1158 $(\mathbf{W}_{1,t+1}, \dots, \mathbf{W}_{d,t+1})$, and

$$1159 \mathbb{E} \|\mathbf{x}_t\|^2 = \sum_{i=1}^d (\mathbf{W}_{i,t})_{11}.$$

1160 Thus mean-square linear stability is equivalent to uniform boundedness of $\sum_{i=1}^d (\mathbf{W}_{i,t})_{11}$ over time,
1161 for every initialization.

1162

1163 **Reduction to strictly positive eigenvalues.** Let $P := \{i : \lambda_i > 0\}$ and $Z := \{i : \lambda_i = 0\}$.
1164 The coordinates indexed by P form an autonomous subsystem. Define the restricted operator
1165 $\mathcal{T}_P : (\mathbb{S}^2)^{|P|} \rightarrow (\mathbb{S}^2)^{|P|}$ by

$$1166 (\mathcal{T}_P(\mathbf{W}))_i = \mathbf{A}_i \mathbf{W}_i \mathbf{A}_i^\top + \eta^2 \left(\lambda_i^2 (\mathbf{W}_i)_{11} + \sum_{j \in P} \lambda_j^2 (\mathbf{W}_j)_{11} \right) \mathbf{Q}, \quad i \in P. \quad (20)$$

1167

1168 Since all $\lambda_i > 0$ for $i \in P$, Theorem 4 applies directly to \mathcal{T}_P . In particular,

$$1169 \rho(\mathcal{T}_P) < 1 \iff \eta \lambda_{\max}(\mathbf{H}) < 1 - \beta^2 \quad \text{and} \quad \sum_{i \in P} \frac{\eta \lambda_i}{2(1 - \beta) \left(1 - \frac{\eta \lambda_i}{1 - \beta^2}\right)} < 1,$$

1170 with the corresponding statement for $\rho(\mathcal{T}_P) \leq 1$.

1171

1172 It therefore suffices to relate $\rho(\mathcal{T}_P)$ to uniform boundedness of $\sum_{i=1}^d (\mathbf{W}_{i,t})_{11}$ over time, for every
1173 initialization. In particular, it suffices to prove that:

$$1174 \rho(\mathcal{T}_P) < 1 \implies \sup_{t \geq 0} \sum_{i=1}^d (\mathbf{W}_{i,t})_{11} < \infty \quad \text{for every initialization,}$$

1175

1188 and

$$1189 \sup_{t \geq 0} \sum_{i=1}^d (\mathbf{W}_{i,t})_{11} < \infty \text{ for every initialization} \implies \rho(\mathcal{T}_P) \leq 1.$$

1192 **(i)** Proof of $(\sup_{t \geq 0} \sum_{i=1}^d (\mathbf{W}_{i,t})_{11} < \infty \text{ for every initialization} \implies \rho(\mathcal{T}_P) \leq 1)$. We prove the
 1193 contrapositive. Assume $\rho(\mathcal{T}_P) > 1$ and set $r := \rho(\mathcal{T}_P)$. Recall that by Theorem 4(a), there exists
 1194 $\mathbf{W}^* \in \mathcal{K}_P \setminus \{\mathbf{0}\}$ such that

$$1195 \mathcal{T}_P(\mathbf{W}^*) = r\mathbf{W}^* \quad \text{and} \quad \sum_{i \in P} (\mathbf{W}_i^*)_{11} > 0. \quad (21)$$

1197 Extend \mathbf{W}^* to $\bar{\mathbf{W}}^* \in \mathcal{K}$ by setting $(\bar{\mathbf{W}}^*)_i = \mathbf{W}_i^*$ for $i \in P$ and $(\bar{\mathbf{W}}^*)_i = \mathbf{0}$ for $i \in Z$. Since the
 1198 P -subsystem is autonomous, by (21) it holds that

$$1199 (\mathcal{T}^t(\bar{\mathbf{W}}^*))_i = r^t \mathbf{W}_i^* \\
 1200 \text{ for all } t \geq 0 \text{ and coordinate } i \in P. \text{ Consequently, we have} \\
 1201 \mathcal{T}^t(\bar{\mathbf{W}}^*) \geq r^t \bar{\mathbf{W}}^* \quad (22)$$

1202 for all $t \geq 0$.

1203 Initialize $\mathbf{m}_0 = \mathbf{0}$ and $\mathbf{x}_0 = (x_{0,0}, x_{1,0}, \dots, x_{d,0}) \in \mathbb{R}^d$ such that

$$1204 \mathbf{W}_{i,0} = \begin{bmatrix} x_{i,0}^2 & 0 \\ 0 & 0 \end{bmatrix}, \quad x_{i,0}^2 := (\bar{\mathbf{W}}_i^*)_{11}.$$

1207 Denote $s_0 := \sum_{i=1}^d \eta^2 \lambda_i^2(\mathbf{W}_{i,0})_{11} = \sum_{i \in P} \eta^2 \lambda_i^2(\mathbf{W}_i^*)_{11}$. Note that $s_0 > 0$ by (21). Moreover, it
 1208 holds that

$$1209 \mathbf{W}_{i,1} = \mathbf{A}_i \mathbf{W}_{i,0} \mathbf{A}_i^\top + \eta^2 \lambda_i^2(\mathbf{W}_{i,0})_{11} \mathbf{Q} + s_0 \mathbf{Q} \geq s_0 \mathbf{Q}, \quad i = 1, \dots, d.$$

1210 Equivalently,

$$1211 \mathbf{W}_1 \geq s_0 \bar{\mathbf{Q}}, \quad \bar{\mathbf{Q}} := (\mathbf{Q}, \dots, \mathbf{Q}) \in \mathcal{K}. \quad (23)$$

1212 Since $\bar{\mathbf{W}}^* \in \mathcal{K}$ and each block \mathbf{Q} is positive definite on its support, there exists $\alpha > 0$ such that

$$1213 \alpha \bar{\mathbf{W}}^* \leq \bar{\mathbf{Q}}. \quad (24)$$

1214 Combining (23) and (24), we have

$$1215 \mathbf{W}_1 \geq s_0 \alpha \bar{\mathbf{W}}^*.$$

1216 Applying \mathcal{T}^{t-1} and using positivity and linearity of \mathcal{T} with (22), we obtain for all $t \geq 1$,

$$1217 \mathbf{W}_t = \mathcal{T}^{t-1}(\mathbf{W}_1) \geq s_0 \alpha \mathcal{T}^{t-1}(\bar{\mathbf{W}}^*) \geq s_0 \alpha r^{t-1} \bar{\mathbf{W}}^*.$$

1218 Define $\ell(\mathbf{W}) := \sum_{i=1}^d (\mathbf{W}_i)_{11}$, which is monotone on \mathcal{K} . Applying ℓ gives

$$1219 \sum_{i=1}^d (\mathbf{W}_{i,t})_{11} \geq s_0 \alpha r^{t-1} \sum_{i=1}^d (\bar{\mathbf{W}}_i^*)_{11} \xrightarrow{t \rightarrow \infty} \infty,$$

1220 since $r > 1$. This contradicts boundedness of $\sum_{i=1}^d (\mathbf{W}_{i,t})_{11}$ for all initializations. Therefore
 1221 $\rho(\mathcal{T}_P) \leq 1$.

1222 **(ii)** Proof of $(\rho(\mathcal{T}_P) < 1 \implies \sup_{t \geq 0} \sum_{i=1}^d (\mathbf{W}_{i,t})_{11} < \infty \text{ for every initialization})$. Assume $\rho(\mathcal{T}_P) <$
 1223 1 . Fix any initial condition $(\mathbf{W}_{i,0})_{i \in [d]} \in \mathcal{K}$ and let $(\mathbf{W}_{i,t})_{i \in [d]} := \mathcal{T}^t(\mathbf{W}_{1,0}, \dots, \mathbf{W}_{d,0})$. Define the
 1224 coupling scalar

$$1225 s_t := \sum_{j=1}^d \lambda_j^2(\mathbf{W}_{j,t})_{11} = \sum_{j \in P} \lambda_j^2(\mathbf{W}_{j,t})_{11}.$$

1232 Then the P -coordinates evolve autonomously according to (20) with initial condition $(\mathbf{W}_{i,0})_{i \in P}$,
 1233 hence

$$1234 (\mathbf{W}_{i,t})_{i \in P} = \mathcal{T}_P^t((\mathbf{W}_{i,0})_{i \in P}).$$

1235 Since $\rho(\mathcal{T}_P) < 1$, there exist $C < \infty$ and $\rho \in (0, 1)$ such that

$$1236 \sum_{i \in P} \|\mathbf{W}_{i,t}\| \leq C \rho^t \sum_{i \in P} \|\mathbf{W}_{i,0}\|, \quad \forall t \geq 0, \quad (25)$$

1237 and in particular $s_t \leq \lambda_{\max}^2 \sum_{i \in P} \|\mathbf{W}_{i,t}\|$ is summable:

$$1238 \sum_{t=0}^{\infty} s_t < \infty. \quad (26)$$

1241

For $i \in Z$, we have $\lambda_i = 0$ and $\mathbf{A}_i = \mathbf{A}_0 := \begin{bmatrix} 1 & -\beta \\ 0 & \beta \end{bmatrix}$, so

$$\mathbf{W}_{i,t+1} = \mathbf{A}_0 \mathbf{W}_{i,t} \mathbf{A}_0^\top + \eta^2 s_t \mathbf{Q}.$$

Unrolling gives

$$\mathbf{W}_{i,t} = \mathbf{A}_0^t \mathbf{W}_{i,0} (\mathbf{A}_0^\top)^t + \eta^2 \sum_{k=0}^{t-1} \mathbf{A}_0^{t-1-k} (s_k \mathbf{Q}) (\mathbf{A}_0^\top)^{t-1-k}. \quad (27)$$

Since $\beta \in [0, 1)$, $\sup_{t \geq 0} \|\mathbf{A}_0^t\| < \infty$ and $\sup_{r \geq 0} \|\mathbf{A}_0^r \mathbf{Q} (\mathbf{A}_0^\top)^r\| < \infty$. Combining these bounds with (26) and (27) yields $\sup_{t \geq 0} \|\mathbf{W}_{i,t}\| < \infty$ for each $i \in Z$. Together with (25) and finiteness of Z , we obtain $\sup_{t \geq 0} \|\mathcal{T}^t(\mathbf{W}_{1,0}, \dots, \mathbf{W}_{d,0})\| < \infty$ for every initial condition. Therefore, $\sup_{t \geq 0} \sum_{i=1}^d (\mathbf{W}_{i,t})_{11}$ is bounded for all initializations.

This proves the explicit mean-square stability conditions (16) and (17).

Critical step size and bounds. Define

$$S(\eta, \beta) := \sum_{i=1}^d \frac{\eta \lambda_i}{2(1-\beta) \left(1 - \frac{\eta \lambda_i}{1-\beta^2}\right)}.$$

On $\eta \in \left(0, \frac{1-\beta^2}{\lambda_{\max}(\mathbf{H})}\right)$, each summand is strictly increasing in η , hence so is $S(\eta, \beta)$. Therefore η_{ms}^* is the unique $\eta > 0$ satisfying (18).

For the upper bound in (19), using $1 - \frac{\eta_{\text{ms}}^* \lambda_i}{1-\beta^2} \leq 1$ gives

$$1 = \sum_{i=1}^d \frac{\eta_{\text{ms}}^* \lambda_i}{2(1-\beta) \left(1 - \frac{\eta_{\text{ms}}^* \lambda_i}{1-\beta^2}\right)} \geq \sum_{i=1}^d \frac{\eta_{\text{ms}}^* \lambda_i}{2(1-\beta)} = \frac{\eta_{\text{ms}}^* \text{Tr}(\mathbf{H})}{2(1-\beta)},$$

so $\eta_{\text{ms}}^* \leq \frac{2(1-\beta)}{\text{Tr}(\mathbf{H})}$.

For the lower bound, since $\lambda_i \leq \lambda_{\max}(\mathbf{H})$,

$$1 - \frac{\eta_{\text{ms}}^* \lambda_i}{1-\beta^2} \geq 1 - \frac{\eta_{\text{ms}}^* \lambda_{\max}(\mathbf{H})}{1-\beta^2},$$

and hence

$$1 = \sum_{i=1}^d \frac{\eta_{\text{ms}}^* \lambda_i}{2(1-\beta) \left(1 - \frac{\eta_{\text{ms}}^* \lambda_i}{1-\beta^2}\right)} \leq \frac{\eta_{\text{ms}}^* \text{Tr}(\mathbf{H})}{2(1-\beta) \left(1 - \frac{\eta_{\text{ms}}^* \lambda_{\max}(\mathbf{H})}{1-\beta^2}\right)}.$$

Rearranging, we obtain

$$\eta_{\text{ms}}^* \geq \frac{2(1-\beta)}{\text{Tr}(\mathbf{H}) + \frac{2\lambda_{\max}(\mathbf{H})}{1+\beta}},$$

which completes the proof. \square

F.4 PROOF OF THEOREM 3

We prove Theorem 3 by characterizing mean-square linear stability of the linearized Frozen ZO-Adam dynamics. Specifically, if We prove the following implication for mean-square stability: if

$$\eta \lambda_{\max}(\mathbf{P}^{-1} \mathbf{H}) < 1 + \beta_1 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \tilde{\lambda}_i}{2 \left(1 - \frac{\eta \tilde{\lambda}_i}{1+\beta_1}\right)} < 1,$$

then the linearized dynamics of frozen ZO-Adam is mean-square linearly stable. Conversely, if the linearized dynamics is mean-square linearly stable, then

$$\eta \lambda_{\max}(\mathbf{P}^{-1} \mathbf{H}) < 1 + \beta_1 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \tilde{\lambda}_i}{2 \left(1 - \frac{\eta \tilde{\lambda}_i}{1+\beta_1}\right)} \leq 1.$$

By definition of the mean-square critical step size η_{ms}^* , it follows that η_{ms}^* is the unique $\eta > 0$ satisfying

$$\eta \lambda_{\max}(\mathbf{P}^{-1}\mathbf{H}) < 1 + \beta_1 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \tilde{\lambda}_i}{2\left(1 - \frac{\eta \tilde{\lambda}_i}{1 + \beta_1}\right)} = 1.$$

Moreover, η_{ms}^* satisfies

$$\frac{2}{\text{Tr}(\mathbf{P}^{-1}\mathbf{H}) + \frac{2\lambda_{\max}(\mathbf{P}^{-1}\mathbf{H})}{1 + \beta_1}} \leq \eta_{\text{ms}}^* \leq \frac{2}{\text{Tr}(\mathbf{P}^{-1}\mathbf{H})}.$$

Why the commutativity assumption $\mathbf{P}\mathbf{H} = \mathbf{H}\mathbf{P}$ is needed. The mean-square analysis for ZO methods closes on the coordinatewise 2×2 covariance blocks only after rotating to a basis where the quadratic is diagonal: the global coupling term is $\text{Tr}(\mathbf{H}\mathbf{X}_t\mathbf{H}) = \sum_i \lambda_i^2 (\mathbf{X}_t)_{ii}$, and this structure is what yields the rank-one coupling operator in Theorem 4. For frozen ZO-Adam, the update involves both \mathbf{H} and the fixed preconditioner \mathbf{P}^{-1} ; to reduce the dynamics to independent eigencoordinates with a single scalar coupling, we need a basis that diagonalizes \mathbf{H} and \mathbf{P} simultaneously. The condition $\mathbf{P}\mathbf{H} = \mathbf{H}\mathbf{P}$ (with $\mathbf{P} > 0$ and $\mathbf{H} \geq 0$) guarantees such a simultaneous orthogonal diagonalization, so that the recursion depends only on the eigenvalues $\{\tilde{\lambda}_i\}$ of $\mathbf{P}^{-1}\mathbf{H}$ and matches the ZO-GDM covariance-operator form after a change of variables. Without commutativity, $\mathbf{P}^{-1/2}\mathbf{H}\mathbf{P}^{-1/2}$ need not be diagonalizable in the same basis as the random rank-one factor, and the diagonal-block closure used in Theorem 4 generally fails.

Proof of Theorem 3. We assume $\mathbf{x}^* = \mathbf{0}$ without loss of generality, so $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{H}\mathbf{x}$ and $\widehat{\nabla}f(\mathbf{x}_t) = \mathbf{u}_t \mathbf{u}_t^\top \mathbf{H}\mathbf{x}_t$ with $\mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ i.i.d. Frozen ZO-Adam is

$$\mathbf{m}_{t+1} = \beta_1 \mathbf{m}_t + (1 - \beta_1) \mathbf{u}_t \mathbf{u}_t^\top \mathbf{H}\mathbf{x}_t, \quad (28)$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathbf{P}^{-1} \mathbf{m}_{t+1}. \quad (29)$$

Preconditioned coordinates. Define $\tilde{\eta} := (1 - \beta_1)\eta$ and the change of variables

$$\tilde{\mathbf{x}}_t := \mathbf{P}^{1/2} \mathbf{x}_t, \quad \tilde{\mathbf{m}}_t := \mathbf{P}^{-1/2} \mathbf{m}_t, \quad \tilde{\mathbf{H}} := \mathbf{P}^{-1/2} \mathbf{H} \mathbf{P}^{-1/2}.$$

Multiplying (28) by $\mathbf{P}^{-1/2}$ and (29) by $\mathbf{P}^{1/2}$ gives

$$\tilde{\mathbf{m}}_{t+1} = \beta_1 \tilde{\mathbf{m}}_t + (1 - \beta_1) \mathbf{P}^{-1/2} \mathbf{u}_t \mathbf{u}_t^\top \mathbf{P}^{1/2} \tilde{\mathbf{H}} \tilde{\mathbf{x}}_t,$$

$$\tilde{\mathbf{x}}_{t+1} = \tilde{\mathbf{x}}_t - \tilde{\eta} \tilde{\mathbf{m}}_{t+1},$$

where we introduced the scaled momentum $\hat{\mathbf{m}}_t := \tilde{\mathbf{m}}_t / (1 - \beta_1)$, so that

$$\hat{\mathbf{m}}_{t+1} = \beta_1 \hat{\mathbf{m}}_t + \mathbf{P}^{-1/2} \mathbf{u}_t \mathbf{u}_t^\top \mathbf{P}^{1/2} \tilde{\mathbf{H}} \tilde{\mathbf{x}}_t, \quad \tilde{\mathbf{x}}_{t+1} = \tilde{\mathbf{x}}_t - \tilde{\eta} \hat{\mathbf{m}}_{t+1}. \quad (30)$$

The linear map $(\mathbf{x}_t, \mathbf{m}_t) \leftrightarrow (\tilde{\mathbf{x}}_t, \hat{\mathbf{m}}_t)$ is invertible, hence $\sup_t \mathbb{E}\|\mathbf{x}_t\|^2 < \infty$ is equivalent to $\sup_t \mathbb{E}\|\tilde{\mathbf{x}}_t\|^2 < \infty$.

Simultaneous diagonalization. Since $\mathbf{P}\mathbf{H} = \mathbf{H}\mathbf{P}$ with $\mathbf{P} > 0$ and $\mathbf{H} \geq 0$ symmetric, there exists an orthogonal \mathbf{U} and diagonal matrices $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_d)$ with $\sigma_i > 0$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ with $\lambda_i \geq 0$ such that

$$\mathbf{P} = \mathbf{U}\Sigma\mathbf{U}^\top, \quad \mathbf{H} = \mathbf{U}\Lambda\mathbf{U}^\top.$$

Consequently,

$$\tilde{\mathbf{H}} = \mathbf{P}^{-1/2} \mathbf{H} \mathbf{P}^{-1/2} = \mathbf{U} \text{diag}\left(\frac{\lambda_1}{\sigma_1}, \dots, \frac{\lambda_d}{\sigma_d}\right) \mathbf{U}^\top =: \mathbf{U}\tilde{\Lambda}\mathbf{U}^\top,$$

so $\tilde{\lambda}_i := \lambda_i / \sigma_i$ are the eigenvalues of $\mathbf{P}^{-1}\mathbf{H}$.

Let $\bar{\mathbf{x}}_t := \mathbf{U}^\top \tilde{\mathbf{x}}_t$, $\bar{\mathbf{m}}_t := \mathbf{U}^\top \hat{\mathbf{m}}_t$, and $\bar{\mathbf{u}}_t := \mathbf{U}^\top \mathbf{u}_t$. Then $\bar{\mathbf{u}}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ i.i.d., and (30) becomes

$$\bar{\mathbf{m}}_{t+1} = \beta_1 \bar{\mathbf{m}}_t + \Sigma^{-1/2} \bar{\mathbf{u}}_t \bar{\mathbf{u}}_t^\top \Sigma^{1/2} \tilde{\Lambda} \bar{\mathbf{x}}_t, \quad \bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \tilde{\eta} \bar{\mathbf{m}}_{t+1}.$$

Since $\|\bar{\mathbf{x}}_t\| = \|\tilde{\mathbf{x}}_t\|$ and $\mathbb{E}\|\mathbf{x}_t\|^2 \leq \lambda_{\min}(\mathbf{P})^{-1} \mathbb{E}\|\tilde{\mathbf{x}}_t\|^2$, it suffices to analyze mean-square boundedness of $\bar{\mathbf{x}}_t$.

Second-moment recursion and covariance operator. Define

$$\mathbf{X}_t := \mathbb{E}[\bar{\mathbf{x}}_t \bar{\mathbf{x}}_t^\top], \quad \mathbf{Y}_t := \mathbb{E}[\tilde{\eta}^2 \bar{\mathbf{m}}_t \bar{\mathbf{m}}_t^\top], \quad \mathbf{C}_t := \mathbb{E}[\tilde{\eta} \bar{\mathbf{x}}_t \bar{\mathbf{m}}_t^\top].$$

Applying Lemma 5 with $\mathbf{P} = \Sigma$ and using independence of $\bar{\mathbf{u}}_t$ from $(\bar{\mathbf{x}}_t, \bar{\mathbf{m}}_t)$ yields a closed recursion on the diagonal entries. Equivalently, for each i , define the diagonal 2×2 covariance block

$$\mathbf{W}_{i,t} := \begin{bmatrix} (\mathbf{X}_t)_{ii} & (\mathbf{C}_t)_{ii} \\ (\mathbf{C}_t)_{ii} & (\mathbf{Y}_t)_{ii} \end{bmatrix} \in \mathbb{S}_+^2.$$

Let $\mathcal{X} := (\mathbb{S}^2)^d$ and $\mathcal{K} := (\mathbb{S}_+^2)^d$. Then there exists a linear operator $\mathcal{T}_{\text{adam}} : \mathcal{X} \rightarrow \mathcal{X}$ such that

$$(\mathbf{W}_{1,t+1}, \dots, \mathbf{W}_{d,t+1}) = \mathcal{T}_{\text{adam}}(\mathbf{W}_{1,t}, \dots, \mathbf{W}_{d,t}),$$

and $\mathcal{T}_{\text{adam}}$ is exactly the operator in Theorem 5, i.e., for $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_d) \in \mathcal{X}$,

$$(\mathcal{T}_{\text{adam}}(\mathbf{W}))_i = \mathbf{A}_i \mathbf{W}_i \mathbf{A}_i^\top + \tilde{\eta}^2 \left(\tilde{\lambda}_i^2(\mathbf{W}_i)_{11} + \sigma_i^{-1} \sum_{j=1}^d \sigma_j \tilde{\lambda}_j^2(\mathbf{W}_j)_{11} \right) \mathbf{Q},$$

where

$$\mathbf{A}_i := \begin{bmatrix} 1 - \tilde{\eta} \tilde{\lambda}_i & -\beta_1 \\ \tilde{\eta} \tilde{\lambda}_i & \beta_1 \end{bmatrix}, \quad \mathbf{Q} := \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Moreover,

$$\mathbb{E} \|\bar{\mathbf{x}}_t\|^2 = \text{Tr}(\mathbf{X}_t) = \sum_{i=1}^d (\mathbf{W}_{i,t})_{11},$$

so mean-square stability reduces to boundedness of $\sum_i (\mathbf{W}_{i,t})_{11}$.

In Theorem 5, we show that the analogous statement to Theorem 4 also holds for $\mathcal{T}_{\text{adam}}$. Hence, the remainder of the proof follows by the same argument as Theorem 2. After reduction to strictly positive eigenvalues as in Theorem 2 and applying Theorem 5, we obtain that: if

$$\eta \lambda_{\max}(\mathbf{P}^{-1} \mathbf{H}) < 1 + \beta_1 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \tilde{\lambda}_i}{2 \left(1 - \frac{\eta \tilde{\lambda}_i}{1 + \beta_1}\right)} < 1,$$

then the linearized dynamics of frozen ZO-Adam is mean-square linearly stable. Conversely, if the linearized dynamics is mean-square linearly stable, then

$$\eta \lambda_{\max}(\mathbf{P}^{-1} \mathbf{H}) < 1 + \beta_1 \quad \text{and} \quad \sum_{i=1}^d \frac{\eta \tilde{\lambda}_i}{2 \left(1 - \frac{\eta \tilde{\lambda}_i}{1 + \beta_1}\right)} \leq 1.$$

Hence, the mean-square critical step size η_{ms}^* is the unique $\eta > 0$ satisfying

$$\eta \tilde{\lambda}_{\max} < 1 + \beta_1, \quad \sum_{i=1}^d \frac{\eta \tilde{\lambda}_i}{2 \left(1 - \frac{\eta \tilde{\lambda}_i}{1 + \beta_1}\right)} = 1,$$

since the left-hand side is strictly increasing in η on $(0, \frac{1 + \beta_1}{\tilde{\lambda}_{\max}})$.

Bounds on η_{ms}^* . At $\eta = \eta_{\text{ms}}^*$,

$$1 = \sum_{i=1}^d \frac{\eta_{\text{ms}}^* \tilde{\lambda}_i}{2 \left(1 - \frac{\eta_{\text{ms}}^* \tilde{\lambda}_i}{1 + \beta_1}\right)} \geq \sum_{i=1}^d \frac{\eta_{\text{ms}}^* \tilde{\lambda}_i}{2} = \frac{\eta_{\text{ms}}^* \text{Tr}(\mathbf{P}^{-1} \mathbf{H})}{2},$$

which gives $\eta_{\text{ms}}^* \leq \frac{2}{\text{Tr}(\mathbf{P}^{-1} \mathbf{H})}$. For the lower bound, use $\tilde{\lambda}_i \leq \tilde{\lambda}_{\max}$ to get

$$1 - \frac{\eta_{\text{ms}}^* \tilde{\lambda}_i}{1 + \beta_1} \geq 1 - \frac{\eta_{\text{ms}}^* \tilde{\lambda}_{\max}}{1 + \beta_1},$$

hence

$$1 = \sum_{i=1}^d \frac{\eta_{\text{ms}}^* \tilde{\lambda}_i}{2 \left(1 - \frac{\eta_{\text{ms}}^* \tilde{\lambda}_i}{1 + \beta_1}\right)} \leq \frac{\eta_{\text{ms}}^* \text{Tr}(\mathbf{P}^{-1} \mathbf{H})}{2 \left(1 - \frac{\eta_{\text{ms}}^* \tilde{\lambda}_{\max}}{1 + \beta_1}\right)}.$$

Rearranging, we obtain

$$\eta_{\text{ms}}^* \geq \frac{2}{\text{Tr}(\mathbf{P}^{-1} \mathbf{H}) + \frac{2 \tilde{\lambda}_{\max}}{1 + \beta_1}}.$$

□

F.5 SPECTRAL ANALYSIS OF FROZEN ZO-ADAM COVARIANCE OPERATOR

Theorem 5 (Spectral characterization of the frozen ZO-Adam covariance operator). *Assume $\mathbf{P} > 0$ and $\mathbf{H} \geq 0$ are symmetric and commute: $\mathbf{P}\mathbf{H} = \mathbf{H}\mathbf{P}$. Let $\tilde{\lambda}_1 \geq \dots \geq \tilde{\lambda}_d > 0$ be the eigenvalues of $\mathbf{P}^{-1}\mathbf{H}$. Fix $\eta > 0$ and $\beta_1 \in [0, 1)$, and define $\tilde{\eta} := (1 - \beta_1)\eta$. Let $\mathcal{X} := (\mathbb{S}^2)^d$ and $\mathcal{K} := (\mathbb{S}_+^2)^d$. In the simultaneous diagonalization basis $\mathbf{P} = \text{diag}(\sigma_1, \dots, \sigma_d)$ and $\tilde{\mathbf{H}} = \mathbf{P}^{-1/2}\mathbf{H}\mathbf{P}^{-1/2} = \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_d)$, define $\mathcal{T}_{\text{adam}} : \mathcal{X} \rightarrow \mathcal{X}$ blockwise by*

$$(\mathcal{T}_{\text{adam}}(\mathbf{W}))_i = \mathbf{A}_i \mathbf{W}_i \mathbf{A}_i^\top + \tilde{\eta}^2 \left(\tilde{\lambda}_i^2 (\mathbf{W}_i)_{11} + \sigma_i^{-1} \sum_{j=1}^d \sigma_j \tilde{\lambda}_j^2 (\mathbf{W}_j)_{11} \right) \mathbf{Q}, \quad i = 1, \dots, d \quad (31)$$

$$\mathbf{A}_i := \begin{bmatrix} 1 - \tilde{\eta} \tilde{\lambda}_i & -\beta_1 \\ \tilde{\eta} \tilde{\lambda}_i & \beta_1 \end{bmatrix}, \quad \mathbf{Q} := \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Define

$$S_{\text{adam}}(\eta, \beta_1) := \sum_{i=1}^d \frac{\eta \tilde{\lambda}_i}{2 \left(1 - \frac{\eta \tilde{\lambda}_i}{1 + \beta_1} \right)}.$$

Then $\mathcal{T}_{\text{adam}}$ satisfies:

(a) (**Leading eigenvalue in the cone**) $\mathcal{T}_{\text{adam}}(\mathcal{K}) \subseteq \mathcal{K}$. Moreover, $\mathcal{T}_{\text{adam}}$ has an eigenvalue equal to its spectral radius $\rho(\mathcal{T}_{\text{adam}})$, with an associated eigenvector $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ satisfying

$$\rho(\mathcal{T}_{\text{adam}}) \geq \tilde{\eta}^2 \sum_{i=1}^d \tilde{\lambda}_i^2 > 0 \quad \text{and} \quad \sum_{i=1}^d (\mathbf{W}_i^*)_{11} > 0.$$

(b) (**Critical case**)

$$\rho(\mathcal{T}_{\text{adam}}) = 1 \iff \eta \tilde{\lambda}_1 < 1 + \beta_1 \text{ and } S_{\text{adam}}(\eta, \beta_1) = 1.$$

(c) (**Subcritical case**)

$$\rho(\mathcal{T}_{\text{adam}}) < 1 \iff \eta \tilde{\lambda}_1 < 1 + \beta_1 \text{ and } S_{\text{adam}}(\eta, \beta_1) < 1.$$

Proof. We work on $\mathcal{X} = (\mathbb{S}^2)^d$ with cone $\mathcal{K} = (\mathbb{S}_+^2)^d$ and order \leq .

Proof of (a). If $\mathbf{W} \in \mathcal{K}$, then $\mathbf{A}_i \mathbf{W}_i \mathbf{A}_i^\top \geq 0$, $(\mathbf{W}_i)_{11} \geq 0$, and $\sigma_i^{-1} \sum_j \sigma_j \tilde{\lambda}_j^2 (\mathbf{W}_j)_{11} \geq 0$. Since $\mathbf{Q} \geq 0$, (31) implies $(\mathcal{T}_{\text{adam}}(\mathbf{W}))_i \geq 0$ for all i , hence $\mathcal{T}_{\text{adam}}(\mathcal{K}) \subseteq \mathcal{K}$.

Since \mathcal{X} is finite dimensional and \mathcal{K} is a closed, convex, pointed cone with nonempty interior, the Krein–Rutman theorem implies that $\mathcal{T}_{\text{adam}}$ admits an eigenvalue equal to $\rho(\mathcal{T}_{\text{adam}})$ with an eigenvector $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$.

To lower bound $\rho(\mathcal{T}_{\text{adam}})$, define $\bar{\mathbf{Q}}_\Sigma := (\sigma_1^{-1} \mathbf{Q}, \dots, \sigma_d^{-1} \mathbf{Q}) \in \mathcal{K} \setminus \{\mathbf{0}\}$. For each i ,

$$\begin{aligned} (\mathcal{T}_{\text{adam}}(\bar{\mathbf{Q}}_\Sigma))_i &= \mathbf{A}_i (\sigma_i^{-1} \mathbf{Q}) \mathbf{A}_i^\top + \tilde{\eta}^2 \left(\tilde{\lambda}_i^2 (\sigma_i^{-1} \mathbf{Q})_{11} + \sigma_i^{-1} \sum_{j=1}^d \sigma_j \tilde{\lambda}_j^2 (\sigma_j^{-1} \mathbf{Q})_{11} \right) \mathbf{Q} \\ &\geq \tilde{\eta}^2 \left(\sigma_i^{-1} \sum_{j=1}^d \tilde{\lambda}_j^2 \right) \mathbf{Q} = c (\bar{\mathbf{Q}}_\Sigma)_i, \end{aligned}$$

where $c := \tilde{\eta}^2 \sum_{j=1}^d \tilde{\lambda}_j^2$. Thus $\mathcal{T}_{\text{adam}}(\bar{\mathbf{Q}}_\Sigma) \geq c \bar{\mathbf{Q}}_\Sigma$, so $\mathcal{T}_{\text{adam}}^t(\bar{\mathbf{Q}}_\Sigma) \geq c^t \bar{\mathbf{Q}}_\Sigma$ for all $t \geq 1$. By Gelfand's formula, $\rho(\mathcal{T}_{\text{adam}}) \geq c$.

Finally, if $(\mathbf{W}_i^*)_{11} = 0$ for all i , then each $\mathbf{W}_i^* = \begin{bmatrix} 0 & 0 \\ 0 & y_i \end{bmatrix}$ with $y_i \geq 0$. In (31), the coupling term vanishes, and $(\mathcal{T}_{\text{adam}}(\mathbf{W}^*))_i = \mathbf{A}_i \mathbf{W}_i^* \mathbf{A}_i^\top$ has $(1, 1)$ -entry equal to $\beta_1^2 y_i$. Since $\rho(\mathcal{T}_{\text{adam}}) > 0$ and $\mathcal{T}_{\text{adam}}(\mathbf{W}^*) = \rho(\mathcal{T}_{\text{adam}}) \mathbf{W}^*$, we get $\beta_1^2 y_i = 0$ for all i , hence $y_i = 0$ and $\mathbf{W}^* = \mathbf{0}$, a contradiction. Therefore $\sum_i (\mathbf{W}_i^*)_{11} > 0$.

Local–global decomposition. For each i , define $\mathcal{M}_i : \mathbb{S}^2 \rightarrow \mathbb{S}^2$ by

$$\mathcal{M}_i(\mathbf{X}) := \mathbf{A}_i \mathbf{X} \mathbf{A}_i^\top + \tilde{\eta}^2 \tilde{\lambda}_i^2 (\mathbf{X})_{11} \mathbf{Q},$$

and define the block-diagonal map $\mathcal{M} : \mathcal{X} \rightarrow \mathcal{X}$ by $(\mathcal{M}(\mathbf{W}))_i := \mathcal{M}_i(\mathbf{W}_i)$. Define the nonnegative linear functional $s_\Sigma : \mathcal{X} \rightarrow \mathbb{R}$ and the injection direction $\bar{\mathbf{Q}}_\Sigma \in \mathcal{K}$ by

$$s_\Sigma(\mathbf{W}) := \sum_{j=1}^d \sigma_j \tilde{\lambda}_j^2(\mathbf{W}_j)_{11}, \quad \bar{\mathbf{Q}}_\Sigma := (\sigma_1^{-1} \mathbf{Q}, \dots, \sigma_d^{-1} \mathbf{Q}).$$

Then (31) is equivalently

$$\mathcal{T}_{\text{adam}}(\mathbf{W}) = \mathcal{M}(\mathbf{W}) + \tilde{\eta}^2 s_\Sigma(\mathbf{W}) \bar{\mathbf{Q}}_\Sigma. \quad (32)$$

Key local facts. All statements below use the ZO-GDM local lemmas with the substitutions $(\eta, \beta, \lambda_i) \leftarrow (\tilde{\eta}, \beta_1, \tilde{\lambda}_i)$:

- Lemma 2(a): $\rho(\mathcal{M}_i) < 1$ if and only if $\tilde{\eta} \tilde{\lambda}_i < 1 - \beta_1^2$.
- Lemma 2(b),(c): Under $\tilde{\eta} \tilde{\lambda}_i < 1 - \beta_1^2$, $(\text{Id} - \mathcal{M}_i)^{-1}$ exists, is \mathbb{S}_+^2 -preserving, and $\mathbf{Y}_i := (\text{Id} - \mathcal{M}_i)^{-1} \mathbf{Q} \geq 0$ satisfies

$$\gamma_i := (\mathbf{Y}_i)_{11} = \frac{1 + \beta_1}{2\tilde{\eta} \tilde{\lambda}_i (1 - \beta_1^2 - \tilde{\eta} \tilde{\lambda}_i)}. \quad (33)$$

- Lemma 3: If $\tilde{\eta} \tilde{\lambda}_i \geq 1 - \beta_1^2$, then for any $\alpha > 0$ there is no $\mathbf{W} \geq 0$ with $(\text{Id} - \mathcal{M}_i)(\mathbf{W}) \geq \alpha \mathbf{Q}$.

Proof of (b). (\Rightarrow) Assume $\rho(\mathcal{T}_{\text{adam}}) = 1$. Let $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ satisfy $\mathcal{T}_{\text{adam}}(\mathbf{W}^*) = \mathbf{W}^*$, and set $s^* := s_\Sigma(\mathbf{W}^*) > 0$. From (32), for each i ,

$$(\text{Id} - \mathcal{M}_i)(\mathbf{W}_i^*) = \tilde{\eta}^2 \frac{s^*}{\sigma_i} \mathbf{Q}.$$

Since $\tilde{\eta}^2 s^* / \sigma_i > 0$, Lemma 3 forces $\rho(\mathcal{M}_i) < 1$ for all i and hence $\tilde{\eta} \tilde{\lambda}_i < 1 - \beta_1^2$, i.e. $\tilde{\eta} \tilde{\lambda}_i < 1 + \beta_1$.

Moreover, $\mathbf{W}_i^* = \tilde{\eta}^2 (s^* / \sigma_i) \mathbf{Y}_i$, so $(\mathbf{W}_i^*)_{11} = \tilde{\eta}^2 (s^* / \sigma_i) \gamma_i$. Plugging into $s^* = s_\Sigma(\mathbf{W}^*)$ yields

$$s^* = \sum_{i=1}^d \sigma_i \tilde{\lambda}_i^2 (\mathbf{W}_i^*)_{11} = \tilde{\eta}^2 s^* \sum_{i=1}^d \tilde{\lambda}_i^2 \gamma_i,$$

and cancelling $s^* > 0$ gives

$$1 = \tilde{\eta}^2 \sum_{i=1}^d \tilde{\lambda}_i^2 \gamma_i. \quad (34)$$

Using (33),

$$\tilde{\eta}^2 \tilde{\lambda}_i^2 \gamma_i = \frac{\tilde{\eta} \tilde{\lambda}_i}{2(1 - \beta_1) \left(1 - \frac{\tilde{\eta} \tilde{\lambda}_i}{1 - \beta_1^2}\right)}.$$

Since $\tilde{\eta} = (1 - \beta_1)\eta$ and $1 - \beta_1^2 = (1 - \beta_1)(1 + \beta_1)$, the right-hand side becomes $\frac{\eta \tilde{\lambda}_i}{2 \left(1 - \frac{\eta \tilde{\lambda}_i}{1 + \beta_1}\right)}$. Thus

(34) is equivalent to $S_{\text{adam}}(\eta, \beta_1) = 1$.

(\Leftarrow) Assume $\eta \tilde{\lambda}_i < 1 + \beta_1$ and $S_{\text{adam}}(\eta, \beta_1) = 1$. Equivalently, $\tilde{\eta} \tilde{\lambda}_i < 1 - \beta_1^2$ and $\sum_i \tilde{\eta}^2 \tilde{\lambda}_i^2 \gamma_i = 1$. Define $\mathbf{Y}_i := (\text{Id} - \mathcal{M}_i)^{-1} \mathbf{Q} \geq 0$ and set $\bar{\mathbf{Y}} := (\mathbf{Y}_1 / \sigma_1, \dots, \mathbf{Y}_d / \sigma_d) \in \mathcal{K} \setminus \{\mathbf{0}\}$. Then

$$(\mathcal{T}_{\text{adam}}(\bar{\mathbf{Y}}))_i = \mathcal{M}_i(\mathbf{Y}_i / \sigma_i) + \tilde{\eta}^2 s_\Sigma(\bar{\mathbf{Y}}) \sigma_i^{-1} \mathbf{Q}.$$

But

$$s_\Sigma(\bar{\mathbf{Y}}) = \sum_{i=1}^d \sigma_i \tilde{\lambda}_i^2 \sigma_i^{-1} (\mathbf{Y}_i)_{11} = \sum_{i=1}^d \tilde{\lambda}_i^2 \gamma_i, \quad \text{and} \quad \tilde{\eta}^2 \sum_{i=1}^d \tilde{\lambda}_i^2 \gamma_i = 1,$$

so that $\tilde{\eta}^2 s_\Sigma(\bar{\mathbf{Y}}) = 1$. Also, by definition, $\mathcal{M}_i(\mathbf{Y}_i / \sigma_i) = \sigma_i^{-1} (\mathbf{Y}_i - \mathbf{Q})$. Hence $(\mathcal{T}_{\text{adam}}(\bar{\mathbf{Y}}))_i = (\mathbf{Y}_i - \mathbf{Q}) / \sigma_i + \sigma_i^{-1} \mathbf{Q} = \mathbf{Y}_i / \sigma_i = \bar{\mathbf{Y}}_i$, i.e. $\mathcal{T}_{\text{adam}}(\bar{\mathbf{Y}}) = \bar{\mathbf{Y}}$. Therefore 1 is an eigenvalue, so $\rho(\mathcal{T}_{\text{adam}}) \geq 1$.

To show $\rho(\mathcal{T}_{\text{adam}}) \leq 1$, let $r := \rho(\mathcal{T}_{\text{adam}}) \geq 1$ and take $\mathbf{W}^* \in \mathcal{K} \setminus \{\mathbf{0}\}$ with $\mathcal{T}_{\text{adam}}(\mathbf{W}^*) = r \mathbf{W}^*$ and $s^* := s_\Sigma(\mathbf{W}^*) > 0$. From (32),

$$\left(\text{Id} - \frac{1}{r} \mathcal{M}\right) \mathbf{W}^* = \frac{\tilde{\eta}^2 s^*}{r} \bar{\mathbf{Q}}_\Sigma.$$

1512 Since $\tilde{\eta} \tilde{\lambda}_1 < 1 - \beta_1^2$, we have $\rho(\mathcal{M}) < 1$, hence $\text{Id} - \frac{1}{r}\mathcal{M}$ is invertible and

$$1513 \mathbf{W}^* = \frac{\tilde{\eta}^2 s^*}{r} \left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}}_\Sigma.$$

1514 Applying $s_\Sigma(\cdot)$ and cancelling $s^* > 0$ gives the fixed-point identity

$$1515 r = \tilde{\eta}^2 s_\Sigma \left(\left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}}_\Sigma \right). \quad (35)$$

1516 Using the Neumann series and $r \geq 1$,

$$1517 s_\Sigma \left(\left(\text{Id} - \frac{1}{r}\mathcal{M} \right)^{-1} \bar{\mathbf{Q}}_\Sigma \right) = \sum_{k=0}^{\infty} r^{-k} s_\Sigma(\mathcal{M}^k(\bar{\mathbf{Q}}_\Sigma)) \leq \sum_{k=0}^{\infty} s_\Sigma(\mathcal{M}^k(\bar{\mathbf{Q}}_\Sigma)) = s_\Sigma((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}}_\Sigma).$$

1518 Because $(\text{Id} - \mathcal{M})^{-1}$ is block-diagonal with blocks $(\text{Id} - \mathcal{M}_i)^{-1}$ and $(\bar{\mathbf{Q}}_\Sigma)_i = \sigma_i^{-1} \mathbf{Q}$, we have

$$1519 ((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}}_\Sigma)_i = \sigma_i^{-1} (\text{Id} - \mathcal{M}_i)^{-1} \mathbf{Q} = \sigma_i^{-1} \mathbf{Y}_i,$$

1520 and therefore

$$1521 s_\Sigma((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}}_\Sigma) = \sum_{i=1}^d \sigma_i \tilde{\lambda}_i^2 (\sigma_i^{-1} \mathbf{Y}_i)_{11} = \sum_{i=1}^d \tilde{\lambda}_i^2 \gamma_i.$$

1522 Plugging into (35) yields

$$1523 r \leq \tilde{\eta}^2 \sum_{i=1}^d \tilde{\lambda}_i^2 \gamma_i = 1,$$

1524 so $r \leq 1$ and hence $\rho(\mathcal{T}_{\text{adam}}) = 1$.

1525 **Proof of (c).** The proof is identical to Theorem 4(c) with the decomposition (32) and the functional s_Σ :

- 1526 • If $\rho(\mathcal{T}_{\text{adam}}) < 1$, the eigenpair in \mathcal{K} and Lemma 3 force $\tilde{\eta} \tilde{\lambda}_1 < 1 - \beta_1^2$ and the fixed-point identity (35) with $r < 1$ implies $\tilde{\eta}^2 \sum_i \tilde{\lambda}_i^2 \gamma_i < 1$, equivalently $S_{\text{adam}}(\eta, \beta_1) < 1$.
- 1527 • Conversely, if $\tilde{\eta} \tilde{\lambda}_1 < 1 - \beta_1^2$ and $S_{\text{adam}}(\eta, \beta_1) < 1$, assuming $\rho(\mathcal{T}_{\text{adam}}) \geq 1$ contradicts (35) by the same comparison with $\tilde{\eta}^2 s_\Sigma((\text{Id} - \mathcal{M})^{-1} \bar{\mathbf{Q}}_\Sigma) = \tilde{\eta}^2 \sum_i \tilde{\lambda}_i^2 \gamma_i < 1$.

1528 This completes (c) and the proof. \square

1529 F.6 TECHNICAL LEMMAS

1530 **Lemma 4.** Let $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then, the following identities hold by the Isserlis' Theorem.

1531 (a) $\mathbb{E}[\mathbf{u}\mathbf{u}^\top] = \mathbf{I}$.

1532 (b) $\mathbb{E}[\mathbf{u}\mathbf{u}^\top \mathbf{A} \mathbf{u}\mathbf{u}^\top] = \mathbf{A} + \mathbf{A}^\top + \text{Tr}(\mathbf{A})\mathbf{I}$ for any matrix \mathbf{A} .

1533 *Proof.* Let $\mathbf{u} = (u_1, \dots, u_d)^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

1534 (a) Let $[d] := \{1, 2, \dots, d\}$. For $i, j \in [d]$, $\mathbb{E}[u_i u_j] = \delta_{ij}$ since the coordinates are centered, independent, and $\text{Var}(u_i) = 1$. Thus $\mathbb{E}[\mathbf{u}\mathbf{u}^\top] = \mathbf{I}$.

1535 (b) Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be arbitrary. For $p, q \in [d]$, the (p, q) -entry of $\mathbb{E}[\mathbf{u}\mathbf{u}^\top \mathbf{A} \mathbf{u}\mathbf{u}^\top]$ equals

$$1536 \mathbb{E}[u_p (\mathbf{u}^\top \mathbf{A} \mathbf{u}) u_q] = \sum_{i,j=1}^d A_{ij} \mathbb{E}[u_p u_i u_j u_q].$$

1537 By the Isserlis' Theorem (Wick's formula) for centered Gaussians,

$$1538 \mathbb{E}[u_p u_i u_j u_q] = \mathbb{E}[u_p u_i] \mathbb{E}[u_j u_q] + \mathbb{E}[u_p u_j] \mathbb{E}[u_i u_q] + \mathbb{E}[u_p u_q] \mathbb{E}[u_i u_j] = \delta_{pi} \delta_{jq} + \delta_{pj} \delta_{iq} + \delta_{pq} \delta_{ij}.$$

1539 Substituting back gives

$$1540 \sum_{i,j} A_{ij} \delta_{pi} \delta_{jq} = A_{pq}, \quad \sum_{i,j} A_{ij} \delta_{pj} \delta_{iq} = A_{qp}, \quad \sum_{i,j} A_{ij} \delta_{pq} \delta_{ij} = \delta_{pq} \text{Tr}(\mathbf{A}).$$

1541 Hence,

$$1542 (\mathbb{E}[\mathbf{u}\mathbf{u}^\top \mathbf{A} \mathbf{u}\mathbf{u}^\top])_{pq} = A_{pq} + A_{qp} + \text{Tr}(\mathbf{A}) \delta_{pq},$$

1566 which is exactly $\mathbf{A} + \mathbf{A}^\top + \text{Tr}(\mathbf{A})\mathbf{I}$. □

1567
1568 **Lemma 5.** *Let $\Sigma > 0$ be a fixed diagonal matrix and let $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then, the following identities*

1569 *hold by the Isserlis' Theorem.*
1570
1571 (a) $\mathbb{E}[\Sigma^{-1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{1/2}] = \mathbf{I}$.

1572 (b) $\mathbb{E}[\Sigma^{-1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{1/2}\mathbf{A}\Sigma^{1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{-1/2}] = \mathbf{A} + \mathbf{A}^\top + \text{Tr}(\Sigma\mathbf{A})\Sigma^{-1}$ for any matrix \mathbf{A} .

1573
1574 *Proof.* Let $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and let $\Sigma > 0$ be diagonal. Write $\mathbf{D} := \Sigma^{1/2}$, so \mathbf{D} is invertible and
1575 $\Sigma^{-1/2} = \mathbf{D}^{-1}$.

1576
1577 (a) We have

$$\mathbb{E}[\Sigma^{-1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{1/2}] = \mathbf{D}^{-1} \mathbb{E}[\mathbf{u}\mathbf{u}^\top] \mathbf{D} = \mathbf{D}^{-1}\mathbf{I}\mathbf{D} = \mathbf{I},$$

1578 where we used Lemma 4(a).

1579 (b) Let \mathbf{A} be arbitrary and define $\tilde{\mathbf{A}} := \mathbf{D}\mathbf{A}\mathbf{D}$. Then

$$\Sigma^{-1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{1/2}\mathbf{A}\Sigma^{1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{-1/2} = \mathbf{D}^{-1}\mathbf{u}\mathbf{u}^\top\tilde{\mathbf{A}}\mathbf{u}\mathbf{u}^\top\mathbf{D}^{-1}.$$

1581 Taking expectation and applying Lemma 4(b) to $\tilde{\mathbf{A}}$ yields

$$\mathbb{E}[\mathbf{u}\mathbf{u}^\top\tilde{\mathbf{A}}\mathbf{u}\mathbf{u}^\top] = \tilde{\mathbf{A}} + \tilde{\mathbf{A}}^\top + \text{Tr}(\tilde{\mathbf{A}})\mathbf{I}.$$

1582 Therefore,

$$\mathbb{E}[\mathbf{D}^{-1}\mathbf{u}\mathbf{u}^\top\tilde{\mathbf{A}}\mathbf{u}\mathbf{u}^\top\mathbf{D}^{-1}] = \mathbf{D}^{-1}\tilde{\mathbf{A}}\mathbf{D}^{-1} + \mathbf{D}^{-1}\tilde{\mathbf{A}}^\top\mathbf{D}^{-1} + \text{Tr}(\tilde{\mathbf{A}})\mathbf{D}^{-1}\mathbf{I}\mathbf{D}^{-1}.$$

1583 Using $\mathbf{D}^{-1}\tilde{\mathbf{A}}\mathbf{D}^{-1} = \mathbf{A}$, $\mathbf{D}^{-1}\tilde{\mathbf{A}}^\top\mathbf{D}^{-1} = \mathbf{A}^\top$, and

$$\text{Tr}(\tilde{\mathbf{A}}) = \text{Tr}(\mathbf{D}\mathbf{A}\mathbf{D}) = \text{Tr}(\mathbf{A}\mathbf{D}^2) = \text{Tr}(\mathbf{A}\Sigma), \quad \mathbf{D}^{-1}\mathbf{I}\mathbf{D}^{-1} = \mathbf{D}^{-2} = \Sigma^{-1},$$

1584 we conclude that

$$\mathbb{E}[\Sigma^{-1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{1/2}\mathbf{A}\Sigma^{1/2}\mathbf{u}\mathbf{u}^\top\Sigma^{-1/2}] = \mathbf{A} + \mathbf{A}^\top + \text{Tr}(\mathbf{A}\Sigma)\Sigma^{-1}.$$

1585 □

1586
1587 **Lemma 6** (The Krein–Rutman Theorem (Krein & Rutman, 1948); see also Theorem 19.2 in Deimling (1985)). *Let \mathcal{X} be a finite-dimensional vector space and $\mathcal{K} \subset \mathcal{X}$ be a closed, convex, pointed cone with nonempty interior. Let $\mathcal{T} : \mathcal{X} \rightarrow \mathcal{X}$ be a linear operator with $\mathcal{T}(\mathcal{K}) \subseteq \mathcal{K}$. Then, there exists $\mathbf{x} \in \mathcal{K} \setminus \{\mathbf{0}\}$ such that*

$$\mathcal{T}(\mathbf{x}) = \rho(\mathcal{T})\mathbf{x}.$$

1588 This is a standard corollary of the Krein–Rutman Theorem in finite dimension.

1602 G EXPERIMENTAL DETAILS

1603
1604 In this section, we provide additional experimental details. Our dataset construction, preprocessing,
1605 and model architectures follow the setup of Cohen et al. (2025).

1606
1607 **Dataset.** We train on a subset of CIFAR-10 consisting of 1,000 training examples drawn from the
1608 first four CIFAR-10 classes. We apply standard preprocessing by subtracting the dataset-wide channel-
1609 wise mean and dividing by the dataset-wide channel-wise standard deviation. For the squared-loss
1610 objective, we use one-hot targets: the ground-truth class is encoded as 1 and the remaining classes
1611 are encoded as 0.

1612 **Architectures.** We evaluate three representative vision architectures:

- 1613 • **CNN.** A four-layer convolutional network with initial channel width 32 and 3×3 convolutional
1614 kernels. We use GeLU activations, average pooling, and a linear readout layer.
- 1615 • **ResNet.** A 20-layer ResNet (He et al., 2016) with GeLU activations and GroupNorm (Wu & He,
1616 2018).
- 1617 • **Vision Transformer (ViT).** A Vision Transformer (Dosovitskiy et al., 2021) with depth 3, embed-
1618 ding dimension 64, 8 attention heads, MLP dimension 256, and patch size 4.

1620 **Top eigenvalue and trace estimation.** During training, we log curvature statistics every 1,000
 1621 iterations. Specifically, we compute the largest eigenvalue and trace of the Hessian (or the pre-
 1622 conditioned Hessian $\mathbf{P}_t^{-1}\mathbf{H}_t$ for ZO-Adam) using matrix-free procedures, without explicitly forming
 1623 the Hessian. We estimate the top eigenvalue via power iteration with 50 iterations, and estimate the
 1624 trace via Hutchinson’s method with 500 probe vectors. Both computations rely on Hessian–vector
 1625 products.

1626
 1627 **Compute.** All experiments are run on a single NVIDIA H100 GPU.

1628 H ADDITIONAL EXPERIMENTS

1629 H.1 EFFECT OF β_1 IN ZO-ADAM

1630
 1631 Figure 8 reports full-batch ZO-Adam training of a CNN on CIFAR-10 while sweeping both β_1 and
 1632 the step size η with fixed $\beta_2 = 0.999$. Across $\beta_1 \in \{0.1, 0.5, 0.9\}$, we observe that the preconditioned
 1633 trace $\text{Tr}(\mathbf{P}_t^{-1}\mathbf{H}_t)$ stabilizes near the stability threshold $2/\eta$, suggesting that the effective stability
 1634 boundary is primarily controlled by η and is largely insensitive to β_1 .

1635 H.2 EMPIRICAL VERIFICATION OF THE COMMUTATIVITY APPROXIMATION IN ZO-ADAM

1636
 1637 We empirically assess the commutativity approximation $\mathbf{P}_t\mathbf{H}_t \approx \mathbf{H}_t\mathbf{P}_t$ used in Theorem 3.

1638
 1639 **Metric and estimator.** We measure the *relative commutator Frobenius norm*

$$1640 \text{RelComm}_F(\mathbf{P}_t, \mathbf{H}_t) := \frac{\|[\mathbf{P}_t, \mathbf{H}_t]\|_F}{\|\mathbf{P}_t\mathbf{H}_t\|_F}, \quad [\mathbf{P}_t, \mathbf{H}_t] = \mathbf{P}_t\mathbf{H}_t - \mathbf{H}_t\mathbf{P}_t.$$

1641
 1642 Using the identity $\|\mathbf{A}\|_F^2 = \mathbb{E}\|\mathbf{Az}\|_2^2$ for \mathbf{z} with i.i.d. Rademacher entries, we estimate both
 1643 $\|[\mathbf{P}_t, \mathbf{H}_t]\|_F$ and $\|\mathbf{P}_t\mathbf{H}_t\|_F$ via Hutchinson probes using only Hessian–vector products. Each probe
 1644 requires two HVPS, namely $\mathbf{H}_t\mathbf{z}$ and $\mathbf{H}_t(\mathbf{P}_t\mathbf{z})$, together with a diagonal scaling by \mathbf{P}_t .

1645
 1646 **Setup.** We use the same CNN and CIFAR-10 setup as in Figure 8 and log $\text{RelComm}_F(\mathbf{P}_t, \mathbf{H}_t)$ at
 1647 the same frequency as the curvature statistics. We use 50 probes per checkpoint.

1648
 1649 **Results.** Across all tested (β_1, η) , $\text{RelComm}_F(\mathbf{P}_t, \mathbf{H}_t)$ decreases rapidly from the value 0.8–0.9
 1650 at initialization to below 0.05 and remains below 0.05 throughout training (bottom row of Figure 8),
 1651 supporting $\mathbf{P}_t\mathbf{H}_t \approx \mathbf{H}_t\mathbf{P}_t$ as a reasonable approximation in this setting.

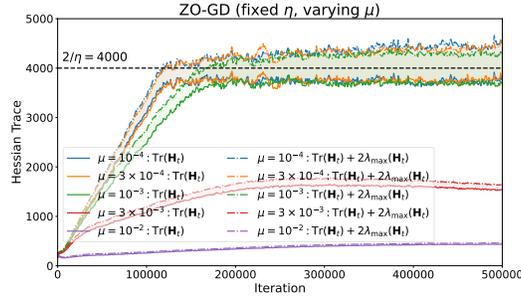
1652
 1653 **Why do \mathbf{P}_t and \mathbf{H}_t approximately commute in deep learning?** For ZO-Adam, \mathbf{P}_t is diagonal in
 1654 the parameter basis. Writing entries explicitly,

$$1655 ([\mathbf{P}_t, \mathbf{H}_t])_{ij} = (p_{t,i} - p_{t,j})(\mathbf{H}_t)_{ij},$$

1656
 1657 so non-commutativity is driven by off-diagonal Hessian couplings between coordinates that receive
 1658 different preconditioning. Consequently, \mathbf{P}_t and \mathbf{H}_t are expected to approximately commute when
 1659 (i) \mathbf{H}_t is close to block-diagonal under a natural parameter partition (e.g., by layer blocks), so
 1660 cross-block couplings are small, and (ii) the preconditioner varies primarily across such blocks
 1661 (or is slowly varying within blocks). This heuristic is consistent with empirical observations that
 1662 neural network Hessians exhibit a near-block-diagonal structure, and with blockwise second-moment
 1663 approximations studied in recent work (e.g., Zhang et al. (2025b)).

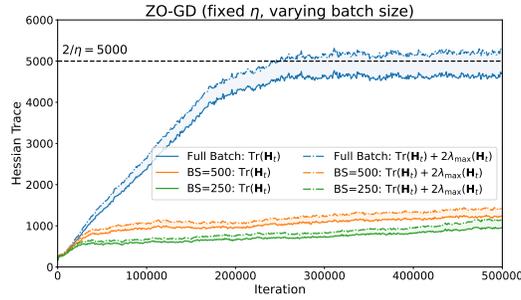
1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684



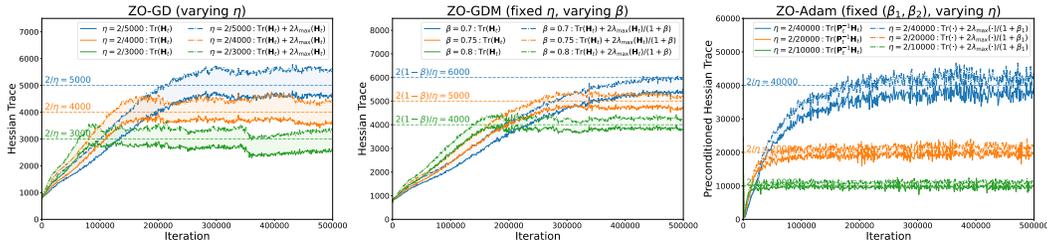
1685 **Figure 4: Effect of the smoothing parameter μ .** We train ZO-GD on a CNN with a fixed step size and vary
1686 the smoothing parameter μ in the two-point estimator. For moderate and small smoothing ($\mu \leq 10^{-3}$), ZO-GD
1687 operates at the mean-square EoS. For larger smoothing ($\mu \geq 3 \times 10^{-3}$), ZO-GD no longer reaches the EoS
1688 threshold and instead trains in a lower-curvature regime with a smaller Hessian trace.

1690
1691
1692
1693
1694
1695
1696
1697
1698
1699



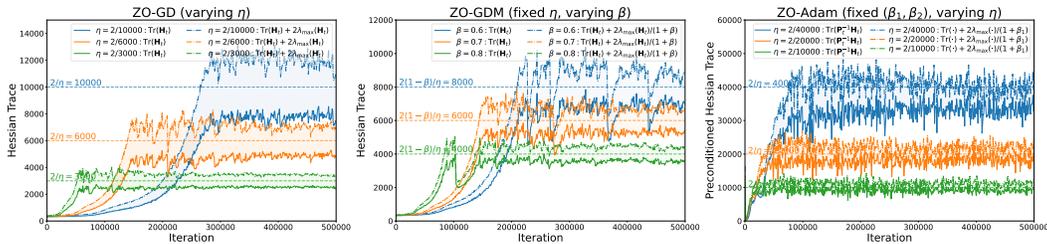
1700 **Figure 5: Effect of batch size in mini-batch ZO-SGD.** We train mini-batch ZO-SGD on a CNN with a fixed
1701 step size and vary the batch size. Compared to full-batch ZO-GD, mini-batch ZO-SGD trains in a lower-curvature
1702 regime with a smaller Hessian trace.

1703
1704
1705
1706
1707
1708
1709
1710
1711
1712



1713 **Figure 6: Mean-square EoS for full-batch ZO methods on ResNet.** We train full-batch ZO-GD, ZO-GDM,
1714 and ZO-Adam on ResNet20 for CIFAR-10 and track the corresponding mean-square stability bounds and
1715 threshold in Section 2.1.

1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727



1728 **Figure 7: Mean-square EoS for full-batch ZO methods on Vision Transformer.** We train full-batch ZO-GD,
1729 ZO-GDM, and ZO-Adam on a Vision Transformer for CIFAR-10 and track the corresponding mean-square
1730 stability bounds and threshold in Section 2.1.

1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781

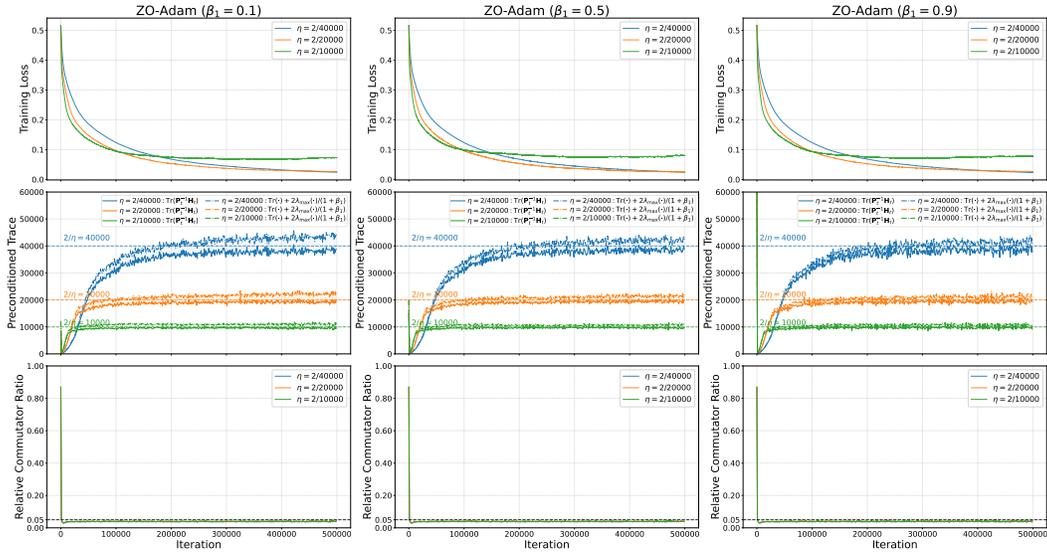


Figure 8: **ZO-Adam sweep over β_1 and η .** Full-batch ZO-Adam on a CNN trained on CIFAR-10 with $\beta_1 \in \{0.1, 0.5, 0.9\}$ (left to right) and multiple step sizes η per setting. *Top*: training loss. *Middle*: preconditioned curvature statistics $\text{Tr}(\mathbf{P}_t^{-1} \mathbf{H}_t)$ and $\text{Tr}(\mathbf{P}_t^{-1} \mathbf{H}_t) + \frac{2}{1+\beta_1} \lambda_{\max}(\mathbf{P}_t^{-1} \mathbf{H}_t)$. *Bottom*: relative commutator ratio $\|[\mathbf{P}_t, \mathbf{H}_t]\|_F / \|\mathbf{P}_t \mathbf{H}_t\|_F$ (cf. Appendix H.2).

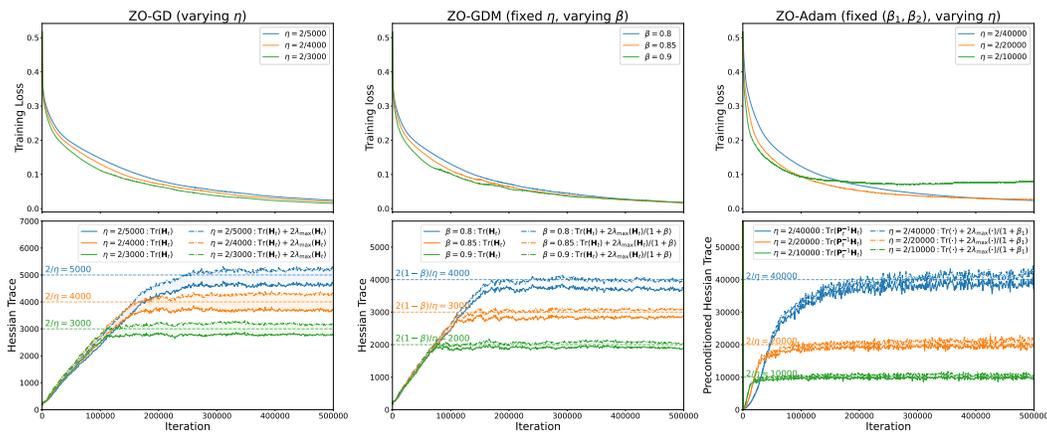


Figure 9: **CNN: same experiments as Figure 2, with training loss.** *Top*: training loss. *Bottom*: the stability-band plots from Figure 2 for ZO-GD (left), ZO-GDM (middle), and ZO-Adam (right).

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835

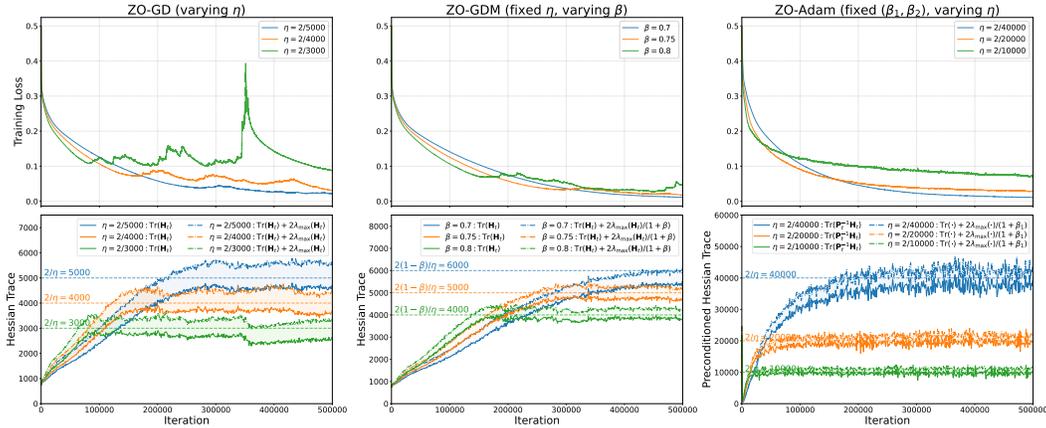


Figure 10: ResNet: same experiments as Figure 6, with training loss. Top: training loss. Bottom: the stability-band plots from Figure 6 for ZO-GD (left), ZO-GDM (middle), and ZO-Adam (right).

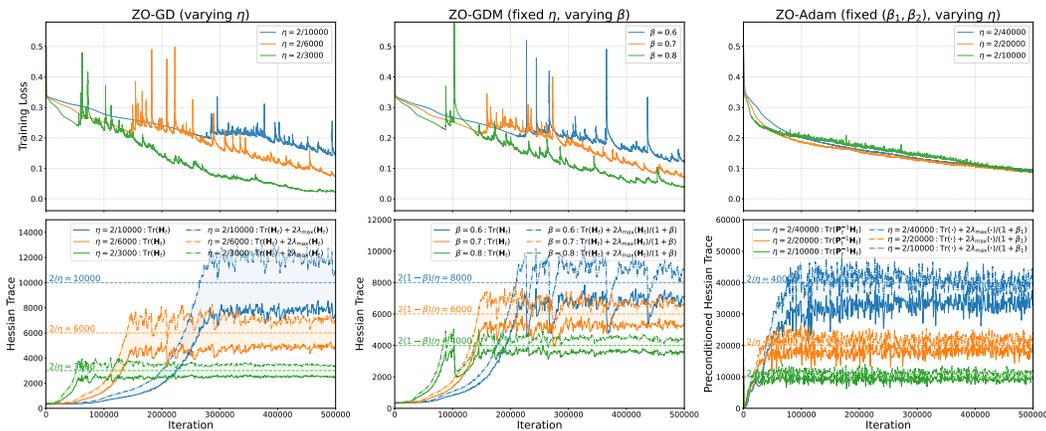


Figure 11: Vision Transformer: same experiments as Figure 7, with training loss. Top: training loss. Bottom: the stability-band plots from Figure 7 for ZO-GD (left), ZO-GDM (middle), and ZO-Adam (right).