
Variance Collapse Predicts When Gate Density Diverges by Activation Class

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Whether the fraction of gradient-carrying units (*gate density*) rises or falls during
2 training is often treated as activation-specific folklore. We show it is governed by
3 a single, derivable quantity: whether the post-BatchNorm pre-activation variance
4 collapses during training. Under coupled-weight-decay optimizers (SGD, Adam),
5 variance collapses for every activation tested; whether this drives gate density
6 toward 0 or 1 depends only on each activation’s fixed threshold-crossing location,
7 computed once via automatic differentiation with no training involved. On an
8 architecture-fixed ablation (48 independent runs), ReLU declines while GELU,
9 SiLU, and Mish all rise instead, confirmed at the level of individual channels (12 of
10 12 activation \times seed cells, using a variance-normalized margin). Under AdamW’s
11 decoupled weight decay, variance no longer collapses — it is flat or growing,
12 not shrinking, for every activation; fed AdamW’s own measured statistics, the
13 identical predictor with no new free parameters correctly anticipates the resulting
14 uniform decline (12 of 12 cells), turning an apparent exception into a confirmation
15 of the mechanism’s scope rather than a separate, unexplained finding. The smooth-
16 activation rise further generalizes across scale (a 365-class real-photograph dataset)
17 and architecture family (two non-BatchNorm, non-convolutional models), while
18 the ReLU-specific decline and a separate representational-rank-rise observation
19 turn out to be specific to convolutional networks. We also report, rather than omit,
20 three honest negative results: gate density does not discriminate via representational
21 rank, is actively misleading for channel pruning on smooth activations, and does
22 not predict low-data fine-tuning accuracy.

23 1 Introduction

24 The local derivative of an activation function gates how much gradient signal a unit carries backward.
25 For ReLU specifically, individual units are known to lose this capacity permanently during training
26 — the dying-ReLU phenomenon [1]. We ask a different, comparative question: across ordinary,
27 successful training, does the *population-level fraction* of gradient-carrying units (*gate density*) move
28 in the same direction for every activation function, or does the direction depend on the activation’s
29 shape?

30 We answer this with a hook-based method that recovers a network’s exact pointwise gate $\Gamma(x) =$
31 $|f'(x)|$ for *any* elementwise activation, without modifying the architecture or knowing f' in closed
32 form (Section 3). Tracking gate density through standard supervised training, the direction is not
33 shared: ReLU-based CNNs decline monotonically while GELU-based networks rise, and holding
34 the architecture exactly fixed and swapping only the activation reproduces this split on identical
35 backbones, ruling out an architecture confound (Section 4).

36 We then derive, rather than only describe, why (Section 5). BatchNorm’s scale shrinks under
37 weight decay and a decaying learning rate — an established consequence of training scale-invariant
38 layers [10, 11] — collapsing the post-normalization pre-activation variance for every activation alike.
39 Whether a shrinking-variance gate concentrates toward 0 or 1 depends only on a fixed, activation-
40 specific threshold-crossing point, computable with no training at all. We verify this mechanism’s
41 central assumption directly, confirm it at the level of individual channels rather than only a population
42 average, and then ask the question this immediately raises: **is the mechanism only descriptive**
43 **within one optimizer, or is it actually predictive across optimizers?** Feeding the unmodified,
44 already-validated predictor AdamW’s own measured statistics — where decoupled weight decay
45 prevents the variance collapse the mechanism’s first step relies on — it correctly anticipates AdamW’s
46 qualitatively different outcome (Section 5.3). This is the paper’s central contribution: one mechanism,
47 one predictor, whose governing input (does variance collapse, set by the coupled-vs-decoupled
48 weight-decay axis) determines the output across every optimizer tested.

49 Finally, we test generalization deliberately rather than assuming it (Section 6), and report two non-
50 findings and one null result exactly as the data shows them (Sections 7, 8), including where this
51 complicates a clean narrative.

52 Contributions

- 53 1. **A general method for measuring the exact pointwise gradient gate** of any elementwise
54 activation without modifying the network, validated against the known ReLU-at-initialization
55 baseline.
- 56 2. **Evidence that gate density’s training-time direction is activation-class-dependent, not**
57 **architecture-dependent:** consistent decline for ReLU and rise for GELU/SiLU/Mish, con-
58 firmed via an architecture-matched ablation (48 independent runs) and robust to normalization
59 scheme and measurement threshold.
- 60 3. **A derived, cross-optimizer predictive mechanism**, verified at the population level, the individual-
61 channel level, and — decisively — by feeding it a different optimizer’s real measured dynamics
62 and confirming it anticipates the qualitatively different outcome with zero new free parameters.
- 63 4. **Generalization tested along three axes (optimizer, scale, architecture family)**, with honest,
64 characterized scope boundaries where the pattern is CNN-specific rather than general, alongside
65 three explicitly negative findings.

66 2 Related Work

67 **Dying ReLU** [1] documents individual ReLU units becoming permanently inactive, treated as
68 a pathology with established remedies. We measure a continuous, population-level statistic that
69 co-occurs with *rising* accuracy, and generalize the question to activations for which "death" is
70 not well-defined. **Activation sparsity:** trained ReLU networks develop increasing forward-pass
71 sparsity [2]; we measure backward-pass gate sparsity as an epoch-resolved trajectory. **Lottery Ticket**
72 **Hypothesis** [3] prunes via static weight masks; our gate-density statistic is continuous and non-
73 interventional, and we test a pruning use case directly with a negative result (Section 7). **Rank/Neural**
74 **Collapse** [4, 5] and **self-supervised representation collapse** [6, 7] use related vocabulary for
75 unrelated, in one case opposite-signed, phenomena; we flag the distinction explicitly (Appendix A.5).
76 **Gradient inversion attacks** [14, 8, 9] are reused unmodified in an exploratory probe (Appendix A.6).
77 **BatchNorm scale dynamics under weight decay** [10, 11] is the established result our mechanism
78 applies, not re-derives. The gate vocabulary itself originates in an earlier synthetic phase-transition
79 study (Appendix A.9), connected here via an exact mathematical correspondence (Section A.3).

80 3 Method

81 For an elementwise activation $y = f(x)$, the chain rule gives, exactly, $\text{grad_input} = \text{grad_output} \cdot$
82 $f'(x)$. A backward hook on any activation module observes both tensors; their elementwise ratio
83 recovers $\Gamma(x) = |f'(x)|$ exactly, for any f , without knowing its closed form. **Gate density** at
84 threshold θ is the fraction of elements with $\Gamma(x) > \theta$ on a fixed evaluation batch, re-measured every
85 epoch on the *same* batch.

86 We instrument CIFAR-native ResNet-18/50 and VGG-11 (a standard 3×3 -stride-1 stem, since the
 87 canonical ImageNet stem destroys a 32×32 image’s resolution immediately) and canonical ViT-B/16
 88 and ConvNeXt-Tiny. The method is validated against the only available ground truth — ReLU’s
 89 $\approx 50\%$ active-fraction at BatchNorm-normalized random initialization — across 20 seeds and three
 90 architectures (Appendix A.1).

91 To separate activation effects from architecture effects, we hold a ResNet-18/VGG-11 skeleton exactly
 92 fixed (identical depth, width, stem, normalization placement, parameter count) and vary only the
 93 activation module via a single constructor argument; the same mechanism extends to a BatchNorm-
 94 vs-GroupNorm ablation. **Statistical methodology:** every directional claim is computed at the level of
 95 the independent seed trajectory — one statistic (a correlation or a delta) per (architecture, activation,
 96 dataset, seed) — and tested via an exact binomial sign test across seeds; epochs within a single run
 97 are never pooled for a hypothesis test, since they are not independent observations.

98 4 The Activation-Class Direction Split

99 Training ResNet-18/50/VGG-11 with SGD on CIFAR-10/100 (18 independent trajectories: 3 archi-
 100 tectures \times 2 datasets \times 3 seeds), gate density falls and effective rank rises in *every* trajectory (sign
 101 test, $p = 3.8 \times 10^{-6}$ for both) while test accuracy rises throughout. In GELU-based ViT-B/16, gate
 102 density instead rises toward a ceiling (4 of 4 runs).

103 This ReLU-vs-ViT contrast confounds activation with architecture family, normalization, and other
 104 structural differences. Holding ResNet-18/VGG-11 exactly fixed and swapping only the activation
 105 (ReLU, GELU, SiLU, Mish; 48 independent runs) resolves this directly: every one of 12 ReLU runs
 106 declines, every one of 12 GELU/SiLU/Mish runs rises (sign test $p = 2.44 \times 10^{-4}$ per activation).
 107 **The direction flips with the activation, not the architecture family.** Effective rank rises in all 48 of
 108 48 runs regardless of activation ($p = 3.55 \times 10^{-15}$) — it does not track the activation-class split and
 109 is reported separately as a non-discriminating observation (Section 7).

110 The direction is robust to two further checks (full detail and figures in Appendix A.2). **Normalization:**
 111 replacing every BatchNorm layer with GroupNorm on the identical ResNet-18 skeleton (4 activations,
 112 3 seeds, 24 runs) leaves the direction unchanged — 3 of 3 ReLU runs decline and 9 of 9 smooth-
 113 activation runs rise under *both* normalization schemes — narrowing what could be driving the
 114 phenomenon to the activation function itself, although Section 5 shows BatchNorm’s presence
 115 specifically supplies the mechanism by which the direction is set. **Threshold:** sweeping the gate-
 116 density threshold $\theta \in \{0.001, \dots, 0.10\}$ (60 checks) gives zero sign flips, although the *magnitude* of
 117 the smooth-activation rise is small at the standard threshold (< 1 point against ReLU’s ~ 13 -point
 118 fall) and grows to $+4.7$ to $+5.5$ points at the strictest threshold tested — an asymmetry explained
 119 in Section 5 by quantile compression of the raw gate-value distribution (the gate-value distribution
 120 compresses from both tails toward its center for smooth activations, simultaneously raising the count
 121 above threshold while lowering the mean; ReLU’s distribution remains exactly bimodal at $\{0, 1\}$, so
 122 its mean and thresholded fraction must move together).

123 5 Mechanism

124 5.1 A derived, tested account

125 For a channel with pre-BatchNorm output u , BatchNorm computes

$$z = \gamma \cdot \frac{u - \text{mean}(u)}{\sqrt{\text{var}(u) + \epsilon}} + \beta \Rightarrow \text{mean}(z) \approx \beta, \quad \text{var}(z) \approx \gamma^2, \quad (1)$$

126 to first order, by BatchNorm’s own definition. A convolutional layer immediately followed by
 127 BatchNorm is scale-invariant, so under SGD with weight decay its effective scale is pulled toward an
 128 equilibrium that decreases as the learning rate decays [10, 11]; this is consistent with γ ’s observed
 129 monotonic shrinkage in every one of 12 trajectories tested ($p = 2.44 \times 10^{-4}$; -81% ReLU, -73 to
 130 -74% GELU/SiLU/Mish).

131 Define $z_{\text{low}}(\theta) = \inf\{z : |f'(z)| > \theta\}$, the point where a gate crosses θ on its negative tail — fixed
 132 and computable via automatic differentiation, no training involved. If $z \sim \mathcal{N}(\mu, \sigma^2)$ with $\sigma \rightarrow 0$ and
 133 μ drifts by a *similar* amount regardless of activation, then $\text{active_frac}(\theta) = P(|f'(z)| > \theta) \rightarrow 0$ or 1

134 depending only on which side of $z_{\text{low}}(\theta)$ the shared μ lands on. We tested the shared-drift assumption
135 directly (9 activations spanning the empirical sign transition, 3 seeds): μ flips from small-positive
136 at initialization to a narrow, shared negative band (-0.041 to -0.066) by epoch 24 regardless of
137 activation shape, and comparing this measured drift to each activation’s independently-computed
138 z_{low} predicts the observed direction in 9 of 9 cases (Appendix A.3), including correctly flagging in
139 advance which case (closest to the boundary) should be least reliable.

140 5.2 Direct per-channel verification

141 The test above is at the population level (one pooled μ per activation/seed/epoch). Logging per-
142 channel μ_c, σ_c and active-gradient fraction for every channel of ResNet-18 (4 activations, 3 seeds,
143 3,904 channels per run), a raw per-channel margin $\mu_c - z_{\text{low}}$ gives a confusing, activation-split corre-
144 lation against active-gradient fraction: strong for ReLU but weakly *negative* for GELU/SiLU/Mish.
145 Diagnosis: ReLU’s z_{low} sits within 0.03 standard deviations of its channels’ own mean (the threshold
146 lives at the live center of the distribution), while the smooth activations’ z_{low} sits 0.9–1.5 standard
147 deviations away; σ shrinks $\sim 3\times$ for every activation, while μ also drifts slightly negative for every
148 activation. For ReLU this μ -drift directly crosses the live threshold; for the smooth activations the
149 same small μ -drift instead pulls the *raw* margin down even as σ ’s much larger relative shrinkage
150 pulls active-gradient fraction *up* by thinning the sub-threshold tail. The raw margin conflates these
151 forces; it does not capture σ .

152 Replacing the raw margin with the σ -normalized z-score $(\mu_c - z_{\text{low}})/\sigma_c$ removes the confound: all
153 four activations now show strong, uniform, positive per-channel correlation (90–96% of channels per
154 seed; 12 of 12 activation \times seed cells majority-positive, sign test $p = 2.44 \times 10^{-4}$). **Thin-margin**
155 **activations (ReLU) are driven by μ crossing a threshold at their distribution’s center; thick-**
156 **margin activations (GELU, SiLU, Mish) are driven by σ -shrinkage thinning a distant tail —**
157 **two routes to the same threshold-crossing logic, unified by one σ -normalized quantity.**

158 5.3 A unified prediction across optimizers

159 Section 6.1 below reports that under AdamW the activation-class split collapses into a uniform decline.
160 We treat this as a falsifiable prediction problem: **does the exact, unmodified per-channel predictor**
161 **above, fed AdamW’s own measured $\mu_c(t)$ and $\sigma_c(t)$, anticipate AdamW’s outcome with no**
162 **new free parameters?** The per-channel correlation itself remains positive under AdamW (72–90%
163 of channels per seed) — expected, since the σ -normalized margin and active-gradient fraction are
164 related through the same monotonic threshold-crossing relationship regardless of which optimizer
165 produced the underlying μ, σ trajectory, so this step confirms the relationship is optimizer-invariant
166 rather than constituting new evidence on its own. The substantive test is whether the margin’s *trend*
167 matches AdamW’s observed decline: for 12 of 12 activation \times seed cells it does (68–89% channel-
168 level agreement per cell; population-pooled margin deltas negative for all four activations, matching
169 the observed active-gradient-fraction declines exactly in sign).

170 The mechanistic reason is direct, not inferred: per-channel σ shrinks ~ 67 – 77% under SGD for
171 every activation (driving the smooth-activation rise via tail-thinning), but under AdamW it is flat
172 for ReLU (-0.9%) and *grows* for every smooth activation ($+2.2$ to $+4.8\%$, Table 1). With σ flat-
173 or-growing and μ still drifting slightly negative, the Gaussian threshold-crossing account requires
174 the margin to decline for every activation — which is what is observed. **This confirms the account**
175 **is quantitatively faithful to AdamW’s actual dynamics, not that the predictor is an oracle**
176 **independent of how margin and active-gradient fraction are related:** the relationship between
177 them is close to definitional once $\mu, \sigma, z_{\text{low}}$ are fixed, and what is being tested is whether the *measured*
178 *inputs* (which are not fixed, and were not previously known for AdamW) take AdamW to the predicted
179 regime. They do, in every cell tested.

180 **The mechanism is therefore one predictor whose governing input — whether σ collapses, itself**
181 **set by coupled (SGD, Adam) vs. decoupled (AdamW) weight decay — determines the output**
182 **across every optimizer tested, rather than three disconnected, optimizer-specific observations.**
183 A first-principles account of *why* decoupling prevents σ -collapse remains open (Section 8).

184 A 13-condition continuous smoothness sweep (ReLU; LeakyReLU at four slopes; PReLU; Softplus at
185 $\beta \in \{50, 20, 10, 5\}$; GELU; SiLU; Mish; 3 seeds), whose Softplus(β) family has a derivative *exactly*
186 equal to the logistic sigmoid used by the originating synthetic theory (Appendix A.9), confirms a

Table 1: Per-channel σ change, epoch 0 to 24, by activation and optimizer. The mechanistic reason variance-driven activations stop rising under AdamW.

Activation	SGD	AdamW
ReLU	-77.0%	-0.9% (flat)
GELU	-67.5%	+2.2% (grows)
SiLU	-66.2%	+4.8% (grows)
Mish	-66.9%	+4.7% (grows)

187 monotonic dose-response between an initialization-time smoothness index ($\text{Var}[f'(x)]$ on the fixed
 188 evaluation batch) and the observed trend (seed-level Spearman $\rho = -0.71$, $p = 6 \times 10^{-7}$, $n = 37$;
 189 condition-level $\rho = -0.73$, $p = 0.0045$, $n = 12$). The sign transition occurs within the Softplus(β)
 190 family itself, between $\beta = 50$ (decline) and $\beta = 20$ (rise) — predicted exactly by the z_{low} values
 191 above. One disclosed exception (PReLU, whose learned slope changes its effective smoothness
 192 during training) and the full per-condition breakdown are in Appendix A.3.

193 6 Generalization

194 6.1 Optimizer

195 The architecture-fixed ablation replicates exactly under Adam, coupled weight decay matching SGD’s
 196 update structure (48 of 48 runs, $p = 4.88 \times 10^{-4}$ per activation, magnitudes matching SGD). Under
 197 AdamW, decoupled weight decay, all four activations decline instead of diverging by class (48 of
 198 48 runs, $p = 4.88 \times 10^{-4}$) — the pattern Section 5.3 shows the mechanism anticipates rather than
 199 merely accommodates.

200 6.2 Scale and architecture family

201 We test, rather than assume, two further axes (full detail and per-axis figures in Appendix A.4). **Scale:**
 202 on Tiny-ImageNet-200 (200 classes, 64×64), ReLU’s decline and rank-rise replicate exactly; the
 203 smooth-activation rise replicates in net displacement (9 of 9 / 8 of 9 seeds) though a whole-trajectory
 204 trend statistic loses power on this already-small effect once trajectories become non-monotonic at
 205 this scale. On a 365-class, real-photograph subsample of Places365-Standard (ResNet-50, 96×96 , 2
 206 seeds), the *full* pattern — large ReLU decline, small smooth-activation rise, universal rank-rise —
 207 replicates exactly, at a class count comparable to ImageNet-1k (untested directly: registration-gated,
 208 no anonymous scriptable download). **Architecture family:** on two non-convolutional, LayerNorm-
 209 only architectures (an 8-block MLP-Mixer [12] and a 6-block Transformer-Encoder [13]), the smooth-
 210 activation rise replicates cleanly (6 of 6 runs each, $p = 0.03125$) but **the ReLU decline does not** —
 211 it is small and architecture-dependent outside the CNN family, as is the separate universal-rank-rise
 212 observation.

213 **The smooth-activation rise generalizes broadly across optimizer, scale, normalization, and**
 214 **architecture family; the ReLU decline and the universal-rank-rise observation are CNN-specific**
 215 — a scope narrowing discovered by deliberately testing outside the original setup, not assumed.

216 7 Honestly-Reported Non-Findings

217 Effective rank rises in every condition tested in this paper — 114+ independent runs regardless of
 218 activation, architecture, dataset, or threshold (pooled sign test $p < 10^{-14}$) — and has no power
 219 to discriminate anything this paper argues for; we do not list it as a contribution (full discussion,
 220 including its explicit distinction from Neural Collapse and self-supervised representation collapse, in
 221 Appendix A.5).

222 Testing whether end-of-training gate density predicts pruning importance (ResNet-18/CIFAR-10, 4
 223 activations, 3 seeds, global channel pruning across 1920 channels, ratios 10–70%, no retraining),
 224 gate-density and magnitude pruning are statistically indistinguishable for ReLU (paired- t $p = 0.55$,
 225 an informative equivalence) but gate-density pruning is dramatically and significantly *worse* than

226 magnitude pruning for every smooth activation ($p \leq 1.5 \times 10^{-4}$; 65–80 accuracy points lost at 10%
227 pruning vs. under 1.5 for magnitude pruning, Figure 7) — **gate density should not be used as a**
228 **general-purpose channel-importance proxy**; for three of four activations tested, doing so is actively
229 harmful (full detail in Appendix A.5).

230 An exploratory privacy probe and a pre-registered low-data fine-tuning test are reported in full in
231 Appendices A.6 and A.7: gate stiffness reduces one gradient-inversion attack’s convergence but not a
232 stronger one (attack-dependent, not a general privacy claim), and final gate density does not predict
233 low-data fine-tuning accuracy (a genuine, pre-registered null: pooled Pearson $r = 0.24$, $p = 0.26$;
234 Spearman $\rho = -0.33$, $p = 0.11$, the two not even agreeing in sign).

235 8 Limitations

236 The mechanism explains *direction*, not magnitude or ultimate root cause: we do not derive a closed
237 form for trend magnitude, and an attempted derivation of why the pre-activation mean drifts negative
238 (Appendix A.3) accounts for the order of magnitude of the shared drift but not its residual, statistically
239 real cross-activation variation. We do not have a first-principles account of why AdamW’s decoupled
240 weight decay specifically prevents σ -collapse, only that it empirically does and that the mechanism
241 correctly anticipates the consequence. ConvNeXt-Tiny is excluded from all claims (undertrained,
242 inconsistent rank sign across datasets; Appendix A.4). Architecture families beyond CNNs, MLP-
243 Mixer, and Transformer-Encoder remain untested, as does full-scale ImageNet-1k directly.

244 9 Conclusion

245 Gate density diverges by activation class during ordinary training under coupled-weight-decay
246 optimizers, and this divergence is derived, not asserted: BatchNorm’s scale shrinks under weight
247 decay, collapsing pre-activation variance for every activation alike, and whether this collapse drives
248 gate density toward 0 or 1 depends only on a fixed, activation-specific threshold-crossing point,
249 verified at the population level, the individual-channel level, and across a continuous smoothness
250 sweep. The decisive test is predictive rather than merely descriptive: fed a different optimizer’s
251 (AdamW’s) own measured statistics, the identical mechanism anticipates a qualitatively different,
252 uniform-decline outcome with zero new free parameters, because decoupled weight decay prevents
253 the variance collapse the mechanism’s first step requires. The smooth-activation half of this finding
254 generalizes broadly across optimizer, scale, and architecture family; the ReLU-specific decline and
255 a separate rank-rise observation are CNN-specific. We report, rather than omit, every place this
256 story does not extend cleanly. The result is a bounded, mechanistically-derived and mechanistically-
257 *predictive* account of activation-class-dependent gate dynamics — not a universal theory of gradient
258 dynamics, and we do not claim it as one.

259 References

260 References

- 261 [1] Lu, L., Shin, Y., Su, Y., and Karniadakis, G. E. (2019). Dying ReLU and initialization: theory and numerical
262 examples. *Communications in Computational Physics*, 28(5), 1671–1706.
- 263 [2] Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks.
264 *Proceedings of AISTATS*, 249–256.
- 265 [3] Frankle, J. and Carbin, M. (2019). The lottery ticket hypothesis: finding sparse, trainable neural networks.
266 *International Conference on Learning Representations (ICLR)*.
- 267 [4] Dong, Y., Cordonnier, J.-B., and Loukas, A. (2021). Attention is not all you need: pure attention loses rank
268 doubly exponentially with depth. *International Conference on Machine Learning (ICML)*.
- 269 [5] Pappan, V., Han, X. Y., and Donoho, D. L. (2020). Prevalence of neural collapse during the terminal phase
270 of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40), 24652–24663.
- 271 [6] Hua, T., Wang, W., Xue, Z., Ren, S., Wang, Y., and Zhao, H. (2021). On feature decorrelation in self-
272 supervised learning. *International Conference on Computer Vision (ICCV)*.

- 273 [7] Jing, L., Vincent, P., LeCun, Y., and Tian, Y. (2022). Understanding dimensional collapse in contrastive
274 self-supervised learning. *International Conference on Learning Representations (ICLR)*.
- 275 [8] Zhao, B., Mopuri, K. R., and Bilen, H. (2020). iDLG: improved deep leakage from gradients. *arXiv preprint*
276 *arXiv:2001.02610*.
- 277 [9] Geiping, J., Bauermeister, H., Dröge, H., and Moeller, M. (2020). Inverting gradients—how easy is it to
278 break privacy in federated learning? *Advances in Neural Information Processing Systems*, 33.
- 279 [10] van Laarhoven, T. (2017). L2 regularization versus batch and weight normalization. *arXiv preprint*
280 *arXiv:1706.05350*.
- 281 [11] Hoffer, E., Banner, R., Golan, I., and Soudry, D. (2018). Norm matters: efficient and accurate normalization
282 schemes in deep networks. *Advances in Neural Information Processing Systems*, 31.
- 283 [12] Tolstikhin, I., Houshy, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A.,
284 Keysers, D., Uszkoreit, J., Lucic, M., and Dosovitskiy, A. (2021). MLP-Mixer: an all-MLP architecture for
285 vision. *Advances in Neural Information Processing Systems*, 34.
- 286 [13] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I.
287 (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
- 288 [14] Zhu, L., Liu, Z., and Han, S. (2019). Deep leakage from gradients. *Advances in Neural Information*
289 *Processing Systems*, 32.
- 290 [15] Donoho, D. L. (2006). Compressed sensing. *IEEE Transactions on Information Theory*, 52(4), 1289–1306.

291 A Technical Appendices and Supplementary Material

292 This appendix is supplementary: the main text stands on its own. Every result summarized in the
293 main text is given in full here, with the exact numbers, figures, and experimental detail that did not fit
294 the 9-page limit.

295 A.1 Method detail: validation and architecture-fixed ablation

296 We validate the gate-recovery method against the only available ground truth — ReLU’s known
297 $\approx 50\%$ active-fraction at BatchNorm-normalized random initialization — across 20 seeds and three
298 architectures (Table 2).

Table 2: Active-gradient-fraction at random initialization, 20 seeds, $\theta = 0.01$. ReLU architectures cluster near the theoretical $\approx 50\%$ expectation; GELU-based architectures are near-ceiling; ViT-B/16 shows vanishing gradients at default initialization (resolved as an initialization artifact below, not a property of depth).

Architecture	Mean active-fraction	Std
ResNet-18	0.555	0.156
ResNet-50	0.529	0.116
VGG-11	0.618	0.141
ConvNeXt-Tiny	0.994	0.004
ViT-B/16	NaN (vanishing at init)	—

299 **The ViT-B/16 vanishing-gradient finding is an initialization artifact, not depth-related.** A
300 dedicated falsification sweep across depth (1–12 layers), initialization scheme (default vs. truncated-
301 normal std = 0.02), and batch size (4–64), replicated independently twice, finds gradients vanish in
302 100% of seeds at every depth under default initialization and in 0% of seeds at every depth under the
303 standard fix. This is unrelated to, and should not be conflated with, the training-time direction split
304 studied in the main paper.

305 **Architecture-fixed activation ablation, method.** To separate activation effects from architecture
 306 effects (Section 4), we hold the ResNet-18 and VGG-11 skeletons exactly fixed — identical depth,
 307 width, stem, BatchNorm placement, and parameter count, verified identical across activations — and
 308 vary only the activation module via a single constructor argument. The same mechanism extends to
 309 the BatchNorm-versus-GroupNorm ablation in Appendix A.2 (replacing every BatchNorm2d with
 310 GroupNorm, 32 groups or fewer for narrow layers) on the identical skeleton.

311 A.2 BatchNorm and threshold robustness

312 Replacing every BatchNorm layer with GroupNorm on the identical ResNet-18 skeleton (4 activations,
 313 3 seeds, 24 runs) leaves the direction unchanged: 3 of 3 ReLU runs decline and 9 of 9 smooth-
 314 activation runs rise under *both* normalization schemes (Figure 8). Magnitude is noisier under
 315 GroupNorm, particularly for Mish, and final accuracy is ~ 2 –4 points lower — both attributable to
 316 GroupNorm being unmatched in hyperparameters at this scale, not evidence about the underlying
 317 mechanism. BatchNorm is not a necessary condition for the phenomenon, narrowing what could be
 318 driving it to the activation function itself — though Section 5 shows BatchNorm’s presence *does*
 319 supply the specific mechanism by which the direction is set in the main experiments.

320 The active-gradient-fraction statistic depends on a threshold θ . Sweeping $\theta \in$
 321 $\{0.001, 0.005, 0.01, 0.05, 0.10\}$ on the architecture-fixed ResNet-18 ablation (4 activations, 3 seeds,
 322 60 checks) finds zero sign flips: ReLU declines and every smooth activation rises at every threshold
 323 tested. The magnitude of the effect is asymmetric: at the originally-used threshold ($\theta = 0.01$), the
 324 smooth-activation rise is small in absolute terms (< 1 percentage point) against ReLU’s ~ 13 -point
 325 fall; the effect becomes comparably large only at the strictest threshold tested ($\theta = 0.10$: $+4.7$ to
 326 $+5.5$ points).

327 The raw, threshold-free gate magnitude (`gate_mean`) declines for every activation tested, including
 328 the ones whose thresholded gate density rises (ReLU -25.5% , GELU -15.3% , SiLU -9.0% , Mish
 329 -8.4%). This is not a contradiction: per-channel quantile analysis shows the gate-value distribution
 330 compresses from both tails toward its center for smooth activations (e.g. GELU’s 5th percentile rises
 331 from 0.041 to 0.076 while its 95th percentile falls from 1.10 to 0.88 over training), simultaneously
 332 raising the count above a near-zero threshold while lowering the mean; ReLU’s distribution remains
 333 exactly bimodal at $\{0, 1\}$ throughout, so its mean and thresholded fraction are mathematically
 334 identical and must move together. Section 5 explains why this compression happens for every
 335 activation but produces opposite threshold-crossing outcomes.

336 A.3 Mechanism detail: threshold-crossing table, shared-drift test, smoothness sweep, and the 337 mu-drift derivation

338 Table 3 gives $z_{\text{low}}(\theta = 0.10)$, computed exactly via automatic differentiation, against the empirically
 339 observed gate-density trend.

Table 3: $z_{\text{low}}(\theta = 0.10)$, computed exactly via autograd, vs. the empirically observed gate-density trend.

Activation	z_{low}	Observed $\rho(\text{epoch, gate density})$	Direction
ReLU / LeakyReLU (all slopes)	0.000	-0.96 to -0.99	decline
Softplus $\beta = 50$	-0.044	-0.59	decline
Softplus $\beta = 20$	-0.110	$+0.73$	rise
Softplus $\beta = 10$	-0.219	$+1.00$	rise
Softplus $\beta = 5$	-0.439	$+1.00$	rise
GELU	-0.554	$+1.00$	rise
SiLU	-0.912	$+0.99$	rise
Mish	-0.891	$+0.97$	rise

340 We tested the shared-drift assumption rather than assuming it: logging signed pre-activation mean
 341 (forward-hook only, no backward pass) for 9 representative activations spanning the transition, 3
 342 seeds, 25 epochs. At initialization μ is small and positive for every activation (0.014–0.033); by
 343 epoch 6 it has flipped negative for every single activation tested and remains in a narrow, shared

344 band through epoch 24: -0.041 (SiLU) to -0.066 (Mish) — a common drift, not nine separate
 345 trajectories.

Table 4: Mechanism validation: predicted vs. observed direction, derived from the measured shared μ -drift and the independently-computed z_{low} . 9 of 9 predictions correct.

Activation	$\mu(\text{epoch } 24)$	z_{low}	Margin	Predicted	Match
ReLU	-0.0458	0.000	-0.0458	decline	✓
LeakyReLU(0.01)	-0.0478	0.000	-0.0478	decline	✓
Softplus $\beta = 50$	-0.0476	-0.044	-0.0036	decline	✓
Softplus $\beta = 20$	-0.0568	-0.110	$+0.0532$	rise	✓
Softplus $\beta = 10$	-0.0635	-0.219	$+0.1555$	rise	✓
Softplus $\beta = 5$	-0.0586	-0.439	$+0.3804$	rise	✓
GELU	-0.0508	-0.554	$+0.5032$	rise	✓
SiLU	-0.0413	-0.912	$+0.8707$	rise	✓
Mish	-0.0655	-0.891	$+0.8255$	rise	✓

346 **Smoothness sweep.** Softplus’s derivative is exactly the logistic sigmoid, $\text{softplus}'(x; \beta) =$
 347 $\text{sigmoid}(\beta x)$ — the same gate formalism used by the originating synthetic theory (Appendix A.9),
 348 with β playing the role of that theory’s stiffness parameter α ; this is an exact mathematical correspon-
 349 dence, not an analogy. A 13-condition sweep (ReLU; LeakyReLU at four slopes; PReLU; Softplus at
 350 $\beta \in \{50, 20, 10, 5\}$; GELU; SiLU; Mish; 3 seeds; 25 epochs) tests whether smoothness predicts the
 351 trend in a continuous dose-response, using $\text{Var}[f'(x)]$ at initialization on the fixed evaluation batch
 352 as a smoothness index (trajectories plotted in Figure 5 in the main text).

353 The relationship is significant at the seed-trajectory level for both trend shape ($\rho = -0.710$, $p =$
 354 6.1×10^{-7}) and magnitude ($\rho = -0.605$, $p = 5.7 \times 10^{-5}$), and at the condition level for trend
 355 shape ($\rho = -0.731$, $p = 0.0045$); the condition-level magnitude correlation does not individually
 356 reach $p < 0.05$ at $n = 12$ ($p = 0.122$) but its bootstrap 95% CI, $[-0.775, -0.401]$, excludes zero.
 357 PReLU is a genuine exception: its smoothness index at initialization (0.261) is the highest of all
 358 13 conditions, yet it rises. This reflects a limitation of measuring smoothness only at initialization:
 359 PReLU’s negative-slope parameter is learned and its effective smoothness changes substantially
 360 during training (gate variance falls from ~ 0.24 – 0.28 to ~ 0.08 – 0.12), unlike fixed-shape activations.

361 **The mu-drift derivation.** The mechanism above explains direction given that μ drifts negative; it
 362 does not explain *why*. This project’s optimizer applies weight decay to every parameter, including
 363 BatchNorm’s β . For a parameter evolving under plain SGD with weight decay, $\theta_{t+1} = \theta_t - \eta(g_t +$
 364 $\lambda\theta_t)$, where g_t is the task gradient and λ the decay coefficient: if g_t has a roughly-constant average
 365 value \bar{g} over a window where η is locally constant, the system relaxes toward a quasi-equilibrium
 366 $\theta^* \approx -\bar{g}/\lambda$. For BatchNorm’s β , $g_t = dL/d\beta = \sum_{\text{batch}} dL/dz$ — the gradient flowing into the
 367 pre-activation, summed over the batch. This is a generic fact about SGD with L2 regularization,
 368 not specific to BatchNorm’s scale-invariance (unlike the γ argument, which needs scale-invariance
 369 specifically). The falsifiable prediction: if \bar{g} is set mainly by the shared architecture/task rather than
 370 by each activation’s own shape, β ’s equilibrium value should be nearly identical across very different
 371 activation functions (same $\lambda = 5 \times 10^{-4}$ for all) and should not correlate with activation-specific
 372 shape properties like z_{low} or an initialization-time smoothness index.

373 Tested against the same 9-activation drift data above: μ at epoch 24 ranges only -0.065 to -0.041
 374 (spread 0.024) across activations whose z_{low} spans -0.00005 to -1.234 (four orders of magnitude)
 375 — the narrowness this account predicts. The cross-activation variation is nonetheless statistically
 376 real: one-way ANOVA across the 9 activations, $F = 29.7$, $p = 7.7 \times 10^{-9}$; between-activation
 377 standard deviation (0.0084) is $3.5 \times$ the within-activation (seed) standard deviation (0.0024). This
 378 residual variation does not correlate with z_{low} (Pearson $r = 0.26$, $p = 0.51$; Spearman $\rho = 0.28$,
 379 $p = 0.46$) nor with the initialization-time smoothness index (Pearson $r = 0.15$, $p = 0.71$; Spearman
 380 $\rho = 0.30$, $p = 0.43$). The equilibrium argument correctly predicts the order of magnitude of the
 381 shared drift band and gives it a mechanistic reason (shared λ , roughly shared \bar{g}) rather than leaving
 382 it as an unexplained empirical regularity; it does not explain the residual, statistically significant
 383 cross-activation variation within that band. We tested two natural candidates and stopped rather than
 384 searching for a third that fits.

385 A.4 Generalization detail: scale, architecture family, and ConvNeXt

386 **Tiny-ImageNet-200.** ResNet-18 (the existing CIFAR-native stem, which keeps an 8×8 feature
387 map before the global pool at 64×64 resolution rather than collapsing it prematurely), 200 classes,
388 64×64 , 100k training images, 4 activations, 3 seeds, 25 epochs. ReLU’s decline and effective-rank
389 rise replicate exactly (6 of 6 seed-level checks across both metrics, matching CIFAR in direction and
390 magnitude). The smooth activations’ smaller-magnitude rise also replicates in net displacement for
391 9 of 9 (active-gradient fraction) and 8 of 9 (effective rank) seeds, but the trajectories are no longer
392 monotonic at this scale — an early overshoot followed by partial relaxation — so a whole-trajectory
393 linear-trend statistic loses power on this already-small effect even though net direction is preserved;
394 we report both statistics rather than the more favorable one alone.

395 **Places365-Standard.** ResNet-50 (same CIFAR-native stem family, 96×96 input, 12×12 feature
396 map before the global pool), a fixed, seeded 150-train/20-val-per-class subsample (365 scene classes,
397 real, uncurated photographs, comparable in class count to ImageNet-1k, which we did not test directly
398 — it requires registration and has no anonymous scriptable download), 4 activations, 2 seeds (disclosed
399 as 2, not the usual 3, for compute-budget reasons), 25 epochs. Every seed-level check agrees with
400 the established direction by both the trend and net-displacement statistics: ReLU declines (2 of 2
401 seeds, -6.7 to -8.3 points) while GELU, SiLU, and Mish all rise (2 of 2 seeds each, $+0.3$ to $+1.1$
402 points) — the same large-decline, small-rise asymmetry seen throughout this paper. Effective rank
403 rises for all four activations, consistent with ResNet-50 remaining a CNN. With only 2 seeds the sign
404 test cannot exceed $p = 0.5$; this experiment adds a real-photograph, larger-class-count replication,
405 not independent statistical power beyond what the other CNN experiments already establish.

406 **Architecture family.** Two CIFAR-native, activation-configurable, non-convolutional architectures
407 with no BatchNorm anywhere: an 8-block MLP-Mixer [12] and a 6-block, 4-head Transformer-
408 Encoder [13] (patch size 4, hidden dimension 128, $\sim 1\text{M}$ parameters each), both LayerNorm-only,
409 instrumented with the same gate-recovery hooks (every activation is an explicit submodule, unlike
410 `nn.TransformerEncoderLayer`’s functional activation, which the hook-based recovery cannot
411 see). Same 4 activations, both datasets, 3 seeds, 25 epochs. GELU, SiLU, and Mish all rise robustly
412 in both new architectures (6 of 6 independent runs each, sign test $p = 0.03125$, magnitude $+0.1$
413 to $+0.8$ points). The ReLU decline does not generalize: it rises in 5 of 6 MLP-Mixer runs, and is
414 small and dataset-dependent in the Transformer-Encoder (mildly declining on CIFAR-100, mixed
415 on CIFAR-10), with nothing resembling the large, robust CNN decline. Effective rank is the most
416 architecture-dependent statistic of all: MLP-Mixer’s rank *declines* for every activation including
417 ReLU, and the Transformer-Encoder’s rank splits by activation (ReLU rises, the smooth activations
418 decline).

419 **ConvNeXt-Tiny: excluded from all claims.** A single-seed, 15-epoch training-dynamics run on
420 CIFAR-10 (47.4% final test accuracy) showed effective rank rising (consistent with other smooth
421 activations); the corresponding CIFAR-100 run, badly undertrained at 18.9% final accuracy, showed
422 effective rank *falling*. We attribute the sign disagreement to undertraining, not a real architectural
423 effect, and exclude ConvNeXt-Tiny from every quantitative claim in this paper.

424 A.5 Representational rank and pruning, in full

425 **Representational rank: a null result.** Effective rank rises in every condition tested in this paper
426 — all 18 original trajectories, all 48 architecture-ablation runs, all 24 BatchNorm/GroupNorm runs,
427 regardless of activation, architecture, dataset, or threshold (114+ independent runs, sign test $p <$
428 10^{-14} pooled). A statistic that moves in the same direction in every condition tried has no power to
429 discriminate anything this paper argues for, and we do not list it as a contribution. We flag, but do not
430 resolve, a possible connection to BatchNorm-statistics dynamics (Section 5), and explicitly distinguish
431 it from Neural Collapse [5] (different layer, different training-phase window) and self-supervised
432 representation collapse [6, 7] (opposite sign, unrelated mechanism).

433 **Pruning: gate density is not a general-purpose importance signal.** We tested whether end-of-
434 training gate density predicts channel importance for pruning, using the same ResNet-18/CIFAR-10
435 checkpoints (4 activations, 3 seeds, global channel pruning across 1920 channels in 8 layers, ratios

436 10–70%, no retraining, compared against magnitude pruning and a random baseline; full per-ratio
437 results plotted in Figure 7 in the main text).

438 For ReLU, gate-density and magnitude pruning are statistically indistinguishable (paired- t $p = 0.55$)
439 — informative equivalence. For GELU, SiLU, and Mish, gate-density pruning is significantly worse
440 than magnitude pruning at every ratio ($p = 1.5 \times 10^{-4}$, 5.2×10^{-5} , 2.3×10^{-5} respectively): at 10%
441 pruning, removing the lowest-gate-density channels costs 65–80 accuracy points, while removing the
442 lowest-magnitude channels costs under 1.5 points. Gate density should not be presented or used as a
443 general-purpose proxy for channel importance; for three of the four activations studied, doing so is
444 actively harmful, not merely uninformative.

445 A.6 Exploratory privacy probe, in full

446 As one possible consequence of gate dynamics, we ask whether gate density affects gradient-inversion
447 attack success, using a controllable victim model (a small CNN with one tunable-stiffness sigmoid
448 gate, distinct from the real architectures studied above) against DLG [14], iDLG [8], and the stronger,
449 scale-invariant Inverting Gradients attack [9] (7 stiffness values, 20 seeds).

450 Increasing gate stiffness significantly reduces iDLG’s convergence probability (logistic regression
451 slope -0.608 , $p = 3.2 \times 10^{-5}$) but has no effect on Inverting Gradients, which converges regardless
452 (100%; DLG itself shows no significant trend, slope -0.142 , $p = 0.35$). We do not claim this
453 connects to the real-architecture gate dynamics studied in the main paper — this experiment uses
454 a different, minimal controlled model. The result we do claim: gate stiffness is not a general-
455 purpose privacy mitigation, and any claim that gradient sparsity protects privacy should specify which
456 adversary it was tested against.

457 A.7 Low-data fine-tuning payoff attempt, in full

458 Since z_{low} is computable before any training, we pre-registered (in writing, before computing any-
459 thing) a test of whether a checkpoint’s final gate density predicts its low-data fine-tuning accuracy: the
460 predictor is each existing checkpoint’s final population-level active-gradient fraction (24 checkpoints:
461 4 activations \times {SGD, AdamW} \times 3 seeds, all already trained 25 epochs on ResNet-18/CIFAR-10);
462 the outcome is test accuracy after fine-tuning the full pretrained network on a fixed, seeded 5%-
463 of-train-set subsample of CIFAR-10 (2,500 images, identical recipe for all 24 checkpoints: SGD,
464 lr= 0.01, momentum 0.9, weight decay 5×10^{-4} , 5 epochs, batch size 64); the pre-registered test is
465 Pearson and Spearman correlation, pooled and within each optimizer subset.

Table 5: Low-data fine-tuning payoff test, exactly as pre-registered.

Subset	n	Pearson r	p	Spearman ρ	p
All	24	0.241	0.258	-0.332	0.113
AdamW only	12	0.445	0.147	0.077	0.812
SGD only	12	0.246	0.441	-0.203	0.527

466 No statistically detectable correlation in any subset, and Pearson and Spearman do not even agree
467 in sign on the pooled set — consistent with no real relationship rather than a weak true effect being
468 underpowered. At the design tested (final population-level active-gradient fraction vs. post-fine-tune
469 accuracy on a fixed 5% CIFAR-10 subsample), gate density does not predict low-data fine-tuning
470 performance. We did not substitute a different outcome metric (e.g. accuracy drop instead of post-
471 fine-tune accuracy) after seeing this null; the metric was fixed in the pre-registration and that is what
472 is reported.

473 A.8 Reproducibility and hyperparameters

474 All training-dynamics experiments use SGD (momentum 0.9, weight decay 5×10^{-4} , Nesterov),
475 cosine learning-rate schedule from an initial rate of 0.1 (CIFAR-native architectures, batch size
476 128) or scaled proportionally for 224-resolution architectures, 25 epochs unless otherwise noted, 3
477 independent seeds per condition unless a training failure is disclosed. Gate-density and rank statistics
478 are computed on a single fixed 64-example held-out batch per run, re-measured every epoch. All

479 experiments ran on a single 80 GB data-center GPU. Code, raw result files, and figures will be
 480 released upon acceptance.

481 The optimizer-generalization experiment (Section 6.1) uses identical architecture/data/epoch settings
 482 with Adam and AdamW substituted for SGD, learning rate 10^{-3} , weight decay 5×10^{-4} (coupled for
 483 Adam, decoupled for AdamW), all other hyperparameters unchanged. The Tiny-ImageNet experiment
 484 uses the public Tiny-ImageNet-200 release (200 classes, 64×64 , 100k train/10k val images), the
 485 same CIFAR-native ResNet-18 stem and SGD configuration as above, ImageNet-standard per-channel
 486 normalization statistics, and the same random-crop/flip augmentation pattern used throughout this
 487 paper (crop padding scaled to the 64×64 input). The sequence-model experiment uses an 8-block
 488 MLP-Mixer (patch size 4, hidden dimension 128, token-mixing dimension 64, channel-mixing di-
 489 mension 512) and a 6-block, 4-head Transformer-Encoder (patch size 4, hidden dimension 128, MLP
 490 dimension 256, learnable positional embeddings, a prepended class token), both LayerNorm-only,
 491 $\sim 1\text{M}$ parameters each, otherwise the same SGD configuration and CIFAR-10/CIFAR-100 datasets as
 492 the architecture-fixed CNN ablation. The per-channel mechanism experiment uses the same ResNet-
 493 18/CIFAR-10/SGD configuration as the BatchNorm-vs-GroupNorm ablation, logging per-channel
 494 statistics at the same five epochs $\{0, 6, 12, 18, 24\}$ used by the shared-drift test. The Places365
 495 experiment uses the public Places365-Standard easyformat release, a fixed seeded subsample of 150
 496 train and 20 validation images per class (365 classes; the full 1.8M-image training set exceeds this
 497 project’s single-GPU budget), 96×96 input (the same CIFAR-native ResNet-50 stem, 12×12 feature
 498 map before the global pool), RandomResizedCrop (scale 0.7–1.0) and horizontal flip for training
 499 augmentation, ImageNet-standard normalization statistics, otherwise the same SGD configuration as
 500 the main experiments; 2 seeds (not the usual 3, disclosed plainly) for compute-budget reasons.

501 A.9 The synthetic phase-transition theory: motivating context

502 The gate vocabulary used throughout this paper ($\Gamma(x) = |f'(x)|$, active-gradient-fraction, gate
 503 collapse) originates in an earlier, fully synthetic study of a single fixed-kernel inverse-reconstruction
 504 problem with a tunable sigmoid stiffness α . That study is self-contained, uses bootstrap confidence
 505 intervals, Cramér–Rao bounds, and a depth-compounding theorem whose own follow-up experiment
 506 found its key independence assumption holds only asymptotically ($\alpha \geq 40$), which the original study
 507 reports candidly. We preserve it here in full as the originating context for this paper’s vocabulary and,
 508 specifically, for the exact mathematical correspondence between its sigmoid-stiffness parameter α
 509 and the Softplus- β family used in Appendix A.3. Its claims are scoped to the synthetic, fixed-kernel
 510 reconstruction setting described below and are not asserted to transfer to the real, trained-network
 511 setting of the main paper — the connection is established empirically above, not assumed.

512 A.9.1 Problem Setup

513 We consider the fixed single-layer convolutional neural network

$$f_{\theta}(x) = h(Ax), \quad (2)$$

514 where $A \in \mathbb{R}^{N \times N}$ is a convolution matrix corresponding to a fixed kernel θ , and $h(\cdot)$ is an element-
 515 wise sigmoid activation:

$$h(z) = \frac{1}{1 + e^{-\alpha(z-c)}}, \quad h'(z) = \alpha h(z)(1 - h(z)), \quad (3)$$

516 where $\alpha > 0$ is the stiffness parameter and $c > 0$ is a threshold. The input is constrained to $0 \leq x \leq 1$,
 517 the target is binary $y \in \{0, 1\}^N$, and the reconstruction objective is

$$J(x) = \|y - f_{\theta}(x)\|_2^2, \quad x^* = \arg \min_{0 \leq x \leq 1} J(x). \quad (4)$$

518 We define the **gradient gate** $\Gamma(x) = h'(Ax) \in \mathbb{R}^N$, with $\Gamma_i(x) = \alpha [f_{\theta}(x)]_i (1 - [f_{\theta}(x)]_i)$. Coordi-
 519 nate i is **active** if $|\Gamma_i(x)| > \varepsilon$. The **active gradient fraction** is $F_{\alpha}(x) = |\{i : |\Gamma_i(x)| > \varepsilon\}|/N$.

520 For the two-layer extension $f_2(x) = h(A_2 \cdot h(A_1 x))$, the gradient via chain rule is $\nabla J(x) =$
 521 $2A_1^T [\Gamma^{(1)} \odot A_2^T (\Gamma^{(2)} \odot (h_2 - y))]$, with effective gate $\Gamma_{\text{eff}} = \Gamma^{(1)} \odot A_2^T \Gamma^{(2)}$ —the source of
 522 compounding collapse.

523 **A.9.2 Landscape Analysis**

524 The gradient is

$$\nabla J(x) = 2A^\top ((h(Ax) - y) \odot h'(Ax)). \quad (5)$$

525 **Lemma A.1** (Gradient Bound). *The magnitude of the gradient $\nabla J(x)$ is bounded by the stiffness α :*

$$\|\nabla J(x)\|_2 \leq \frac{\alpha}{2} \|A^\top (h(Ax) - y)\|_2. \quad (6)$$

526 *Proof.* Since $h'(z) = \alpha h(z)(1 - h(z)) \leq \alpha/4$, applying the triangle inequality to the gradient
527 expression yields the result. \square

528 **Theorem A.2** (Gradient Gate Bound). *For any $\varepsilon \in (0, \alpha/4)$ and $x \in [0, 1]^N$:*

$$\mathbb{E}_x[F_\alpha(x)] \leq 1 - \frac{2\varepsilon}{\alpha}. \quad (7)$$

529 *As $\alpha \rightarrow \infty$ with ε fixed, $F_\alpha(x) \rightarrow 0$ for almost every x with $(Ax)_i \neq c$ for all i .*

530 *Proof.* The set $\{\Gamma_i > \varepsilon\}$ maps to a sigmoid-output interval (s_-, s_+) with $s_\pm = \frac{1}{2} \pm \frac{1}{2} \sqrt{1 - 4\varepsilon/\alpha}$.
531 Its width satisfies $s_+ - s_- = \sqrt{1 - 4\varepsilon/\alpha} \leq 1 - 2\varepsilon/\alpha$ (using $\sqrt{1-t} \leq 1 - t/2$). When $(Ax)_i$ is
532 approximately uniform, $P(\Gamma_i > \varepsilon) \leq 1 - 2\varepsilon/\alpha$. As $\alpha \rightarrow \infty$, the corresponding z -interval shrinks to
533 $\{c\}$, giving $F_\alpha(x) \rightarrow 0$ for any x not satisfying $(Ax)_i = c$ for all i . \square

534 **Theorem A.3** (Compounding Collapse — Scoped Validity). *Under an approximate gate-
535 independence assumption at convergence (empirically verified: mean $|\text{corr}(\Gamma^{(1)}, A_2^\top \Gamma^{(2)})| = 0.201$,
536 max < 0.50 at all tested configurations), for the two-layer model with equal stiffness α and kernel
537 A_2 :*

$$\frac{F_\alpha^{(2)}}{F_\alpha^{(1)}} \leq \frac{1}{1 + c_2 \alpha}, \quad c_2 = \frac{\varepsilon}{\|A_2\|_\infty}. \quad (8)$$

538 *For $\alpha \geq 16$, bootstrap $R^2 \geq 0.919$ (CI lower bound: 0.764). The rate c is not a universal constant; in
539 the valid regime it follows $c(\alpha) \approx 0.211 \cdot \alpha^{-0.674}$ ($R^2 = 0.973$), ranging from $c \approx 0.031$ at $\alpha = 16$
540 to $c \approx 0.012$ at $\alpha = 60$. Below $\alpha = 16$, gradients are insufficiently sparse for the multiplicative
541 independence assumption to hold. **This asymptotic scope ($\alpha \geq 16$, and an independence assumption
542 verified only weakly even in that regime—mean $|\text{corr}| = 0.201$, not strictly zero) is the reason this
543 theorem is not presented as a general law in the main paper.***

544 *Proof sketch.* Coordinate i is active in the two-layer model only if active in both $\Gamma^{(1)}$ and $A_2^\top \Gamma^{(2)}$.
545 Under approximate independence, $P(\text{active in 2L}) = P(\Gamma_i^{(1)} > \varepsilon) \cdot P((A_2^\top \Gamma^{(2)})_i > \varepsilon)$. Applying
546 Theorem A.2 to the second layer with effective threshold $\varepsilon/\|A_2\|_\infty$ gives the bound. \square

547 **Lemma A.4** (Adam Advantage from Momentum). *In the high- α regime, Adam’s first moment
548 $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$ preserves gradient signal from past active sets via exponential decay at
549 rate β_1^k . Oracle coordinate rescaling ($\beta_1 = 0$) does not recover this advantage.*

550 Verification (20 seeds): *Mean IoU — Adam: 0.950, Oracle: 0.794, PGD: 0.811. Wilcoxon $p =$
551 0.0078, Cohen’s $d > 9$ at all $\alpha \geq 5$. The Adam–Oracle gap is non-monotone, peaking at 0.332 at
552 $\alpha = 13$ with a local minimum of 0.246 at $\alpha = 25$.*

553 **Identifiability and Null Space Geometry.** Reconstruction is unique if and only if the effective
554 Jacobian $J(x^*) = \text{diag}(\Gamma(x^*)) \cdot A$ has trivial null space, $\sigma_{\min}(J(x^*)) > 0$.

555 $\sigma_{\min}(J(x^*))$ decreases monotonically with α , crossing zero near α^* . At $\alpha \geq 20$, the null space
556 dimension reaches 46/256 (18% of input space). The stable rank collapses monotonically from 46.8
557 at $\alpha = 1$ to 6.9 at $\alpha = 60$ (6.7 \times), computed via exact SVD of the effective Jacobian, confirming the
558 implicit dimensional reduction induced by the gradient gate.

559 **Uncertainty Quantification for α^* .** To avoid parameter collinearity in the Fisher information
 560 matrix ($\text{cond}(\mathcal{I}_4) = 4.02 \times 10^6$ for the 4-parameter fit, flagged numerically unreliable due to IoU_{\max} -
 561 δ collinearity, $r = 0.781$), we fix $\text{IoU}_{\max} = 1.0$ (justified by the data: Adam achieves $\text{IoU} = 1.000$
 562 at all $\alpha \leq 10$) and fit a 3-parameter model:

$$\overline{\text{IoU}}(\alpha) = \text{IoU}_{\min} + \frac{1 - \text{IoU}_{\min}}{1 + \exp((\alpha - \alpha^*)/\delta)}. \quad (9)$$

563 This reduces $\text{cond}(\mathcal{I}_3)$ to 3.0×10^4 ($135\times$ improvement). The primary uncertainty estimate is the
 564 distribution-free bootstrap 95% CI: $\alpha^* \in [27.92, 29.78]$ (1000 resamples, 20 seeds), corresponding
 565 to a CI width of 1.86. The 3-parameter CRLB gives $\sigma_{\text{CRLB}}(\alpha^*) = 0.344$.

566 A.9.3 Algorithmic Analysis

567 **Projected Gradient Descent (PGD).** Relies directly on gradient magnitude; struggles severely
 568 when gradients vanish. Applies uniform updates across all coordinates, failing to exploit the few
 569 regions where gradients remain informative.

570 **Momentum and Nesterov.** Improve upon PGD by accumulating gradients over time, partially
 571 compensating for sparse updates. Nesterov’s look-ahead correction yields faster convergence in
 572 smooth regimes, but the look-ahead step lands in an equally gradient-sparse region at high α .

573 **Adam.** Rescales gradients coordinate-wise based on historical gradient magnitudes. Lemma A.4
 574 shows that Adam’s advantage arises from momentum accumulation, not coordinate-wise rescaling as
 575 commonly assumed.

576 **L-BFGS-B.** Approximates curvature information and can accelerate convergence in smooth regimes.
 577 Performance degrades when gradients are sparse or unreliable.

Table 6: Qualitative and quantitative comparison of optimisers. IoU and active gradient fraction F_α at $\alpha = 10$ and $\alpha = 20$ (sobel_x, checkerboard, 5 seeds).

Optimiser	Robustness	Stability	$\alpha = 10$		$\alpha = 20$	
			IoU	F_α	IoU	F_α
Adam	High	High	1.000	1.000	.947	.470
L-BFGS-B	Medium	Variable	.999	.998	.883	.359
Momentum	Medium	Medium	1.000	1.000	.784	.328
Nesterov	Medium	Medium	1.000	1.000	.754	.278
PGD	Low	Low	.858	.896	.787	.270

578 A.9.4 Experiments

579 All experiments use 64×64 inputs with five fixed convolutional kernels: **identity-like** (triv-
 580 ial, separable), **avg_blur** (3×3 uniform, low-pass), **random_norm** (random 3×3 , unit Frobe-
 581 nius), **sobel_x** (horizontal edges, spectral norm ≈ 4.0), and **laplacian** (second-order, spectral
 582 norm ≈ 8.0). Threshold $c = 0.5$; stiffness $\alpha \in \{1, 2, 5, 10, 20, 40\}$ (core experiments) and
 583 $\alpha \in \{1, 1.5, 2, 3, 5, 7, 10, 13, 16, 20, 25, 30, 40, 60\}$ (phase transition characterisation, 20 seeds).
 584 Adam uses $\text{lr} = 0.03$, $\beta_1 = 0.9$, $\beta_2 = 0.999$; PGD and Momentum use $\text{lr} = 0.1$ with $\beta = 0.9$;
 585 Nesterov uses the same; L-BFGS-B uses its default line search. All methods run for 200 iterations.

586 **Metrics.** Final loss $J(x)$; Intersection-over-Union (IoU, binarised at 0.5); active gradient fraction
 587 F_α ($\varepsilon = 0.01$); relative loss reduction $(J(x_0) - J(x_T))/J(x_0)$; success rate $P[\text{IoU} > 0.7]$.

588 **Phase Transition Characterisation.** Mean reconstruction IoU follows a sigmoidal decay in α
 589 across a transition band $\alpha \in [13, 30]$. Three threshold values characterise the band:

Threshold	α	IoU at crossing
α_{10} : IoU first drops below 1.000	13	0.996
α_{50} : IoU drops below 0.900	30	0.875
α_{90} : IoU drops below 0.800	≥ 60	0.771

591 The 3-parameter sigmoid fit (IoU_{\max} fixed at 1.0) gives:

$$\text{IoU}_{\min} = 0.782, \quad \alpha^* = 28.97, \quad \delta = 6.29.$$

592 The bootstrap median is $\alpha^* = 28.83$ (95% CI [27.92, 29.78]). A sparse grid concentrated within
 593 the subcritical plateau ($\text{IoU} \equiv 1.0$) cannot resolve α^* reliably, since a bounded least-squares fit on
 594 plateau-only data is underdetermined; a 14-point grid extending to $\alpha = 60$ (20 independent seeds,
 595 1000-resample bootstrap) is required to observe the full sigmoid S-curve and obtain the estimate
 596 above.

597 **Effect of Sigmoid Stiffness.** Performance degrades consistently as α increases through the transi-
 598 tion band (Fig. 19). Adam sustains perfect reconstruction ($\text{IoU} = 1.000$) up to $\alpha = 10$, then degrades
 599 smoothly to 0.947 at $\alpha = 20$ and 0.830 at $\alpha = 40$. The non-adaptive optimisers begin degrading
 600 earlier, with PGD already at $\text{IoU} = 0.858$ at $\alpha = 10$.

Table 7: Performance across sigmoid stiffness values (Adam, `sobel_x`, 20 seeds).

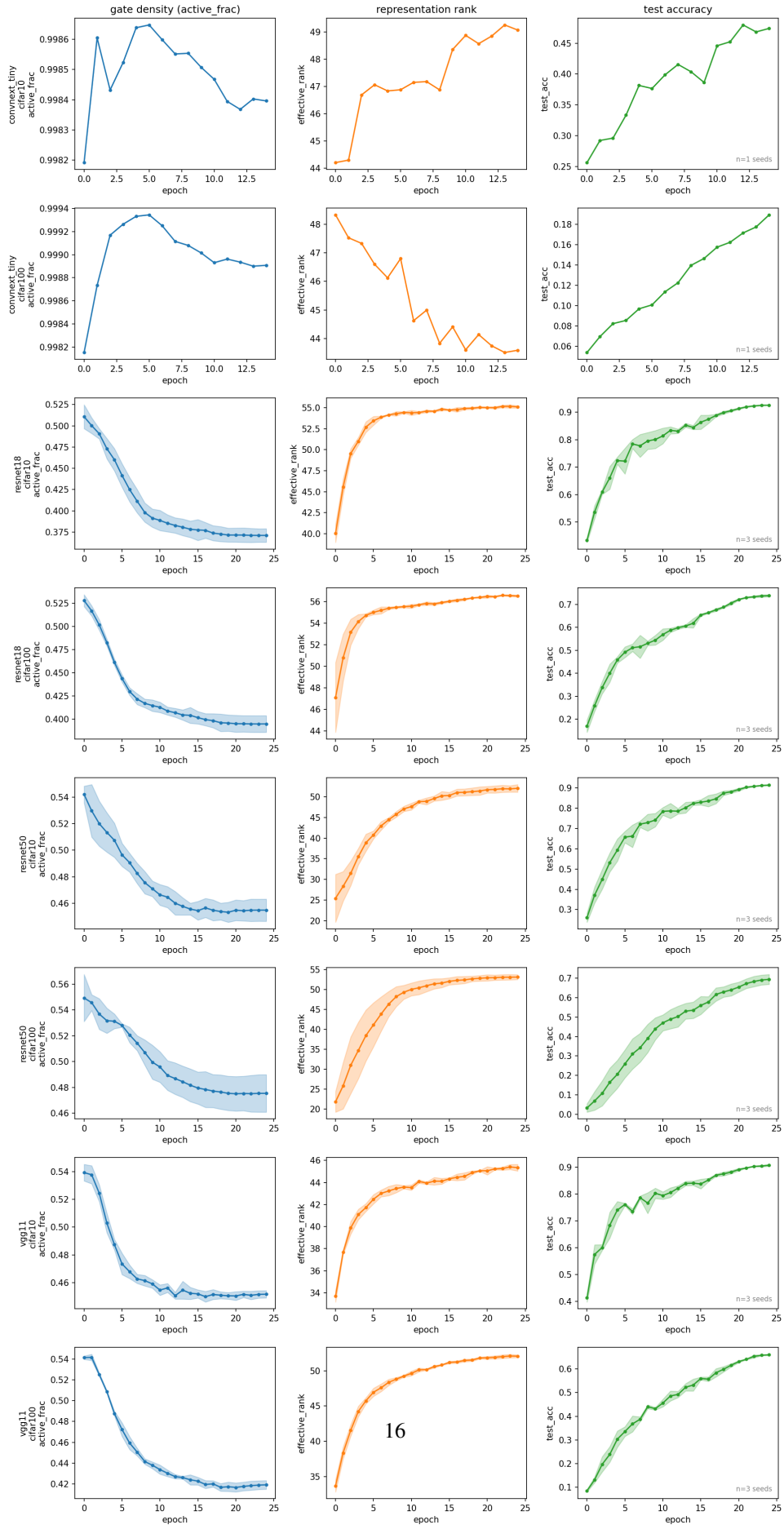
α	Loss (mean \pm std)	IoU (mean \pm std)
1.0	504.87 \pm 9.91	1.000 \pm 0.000
5.0	344.19 \pm 11.62	1.000 \pm 0.000
10.0	488.05 \pm 9.10	1.000 \pm 0.000
20.0	646.31 \pm 11.35	0.947 \pm 0.007
40.0	799.49 \pm 30.10	0.830 \pm 0.012
60.0	892.13 \pm 34.52	0.771 \pm 0.018

601 Impact of Cutoff Threshold c .

602 **Spatial Gradient Gate and Dimensional Collapse.** The gradient gate $\Gamma(x) = h'(Ax)$ collapses
 603 spatially as α increases: a dense, near-uniform field at low α narrows to isolated active coordinates
 604 at edge boundaries at high α , directly instantiating Theorem A.2. The quantitative version of
 605 this collapse, from `grad_sparsity_all_optimizers.csv`: at $\alpha = 1$, the gate is 100% active.
 606 The collapse becomes pronounced in the transition band: by $\alpha = 20$, Adam retains 47.0% active
 607 coordinates versus PGD’s 27.0%; by $\alpha = 40$, Adam retains 19.0% versus PGD’s 5.7%. Stable rank
 608 collapses monotonically from 46.8 at $\alpha = 1$ to 6.9 at $\alpha = 60$ ($6.7\times$), computed via exact SVD of the
 609 effective Jacobian.

610 Gradient Sparsity and Its Impact.

Gate collapse emerges during normal supervised training
(shaded = 95% CI across seeds)



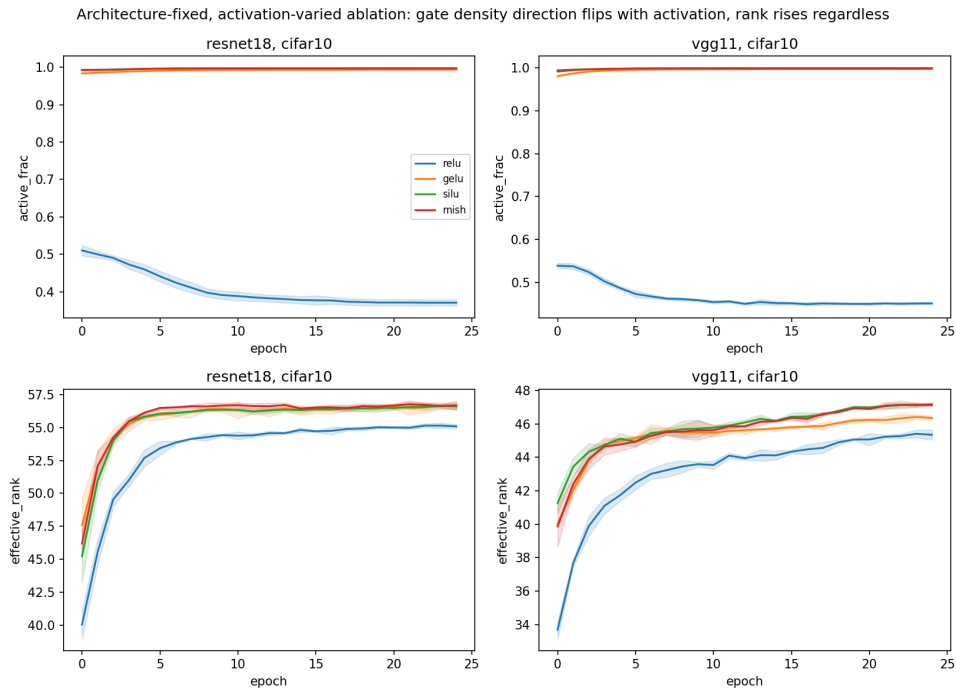


Figure 2: Architecture-fixed, activation-varied ablation. Gate density falls for ReLU and rises for GELU/SiLU/Mish on the *same* backbone; effective rank rises for every activation regardless of direction.

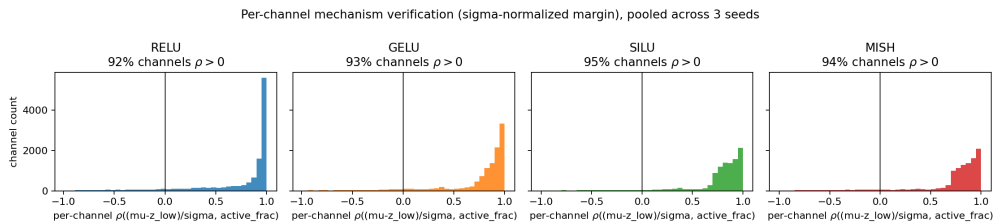


Figure 3: Per-channel verification. Histograms of the per-channel correlation between the σ -normalized margin and active-gradient fraction across 5 logged epochs, pooled over 3 seeds. All four activations show strong, uniform, positive correlation.

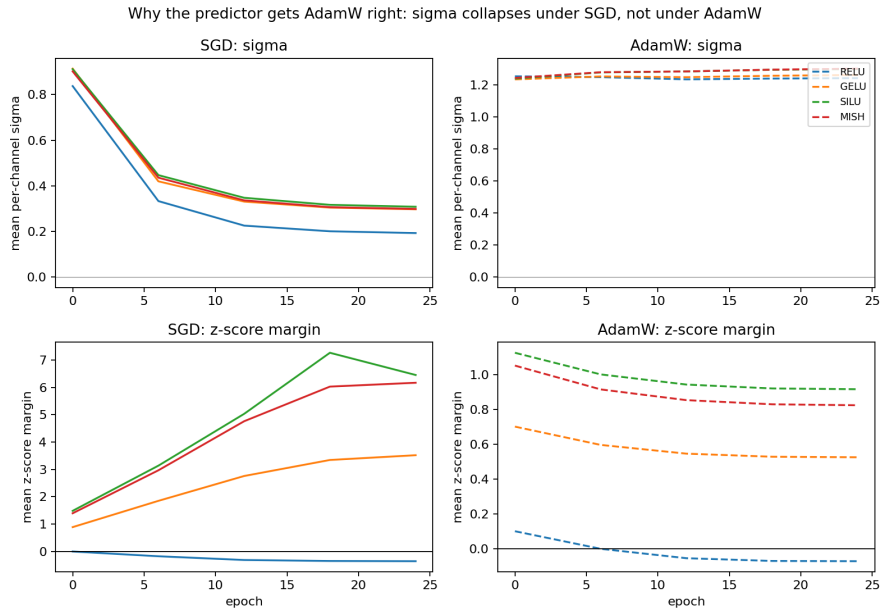


Figure 4: Why the predictor gets AdamW right. σ collapses for every activation under SGD but is flat (ReLU) or grows (GELU/SiLU/Mish) under AdamW; the resulting margin rises for smooth activations under SGD but declines for every activation under AdamW.

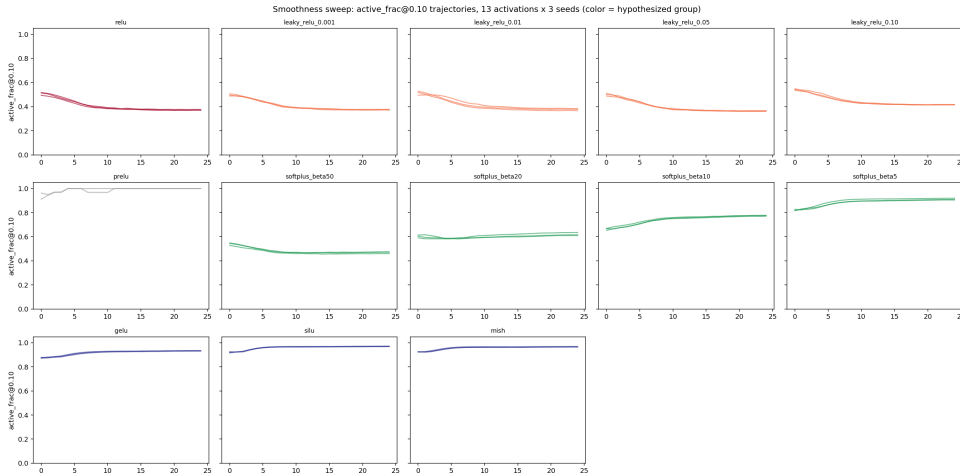


Figure 5: Smoothness sweep. Gate-density trajectories for 13 activation conditions. The sign transition occurs within the Softplus(β) family, predicted exactly by the z_{low} values computed independently of any training run.

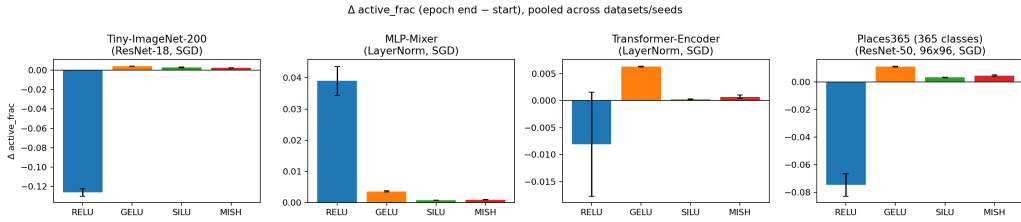


Figure 6: Generalization beyond CIFAR-scale CNNs. Mean Δ active-gradient-fraction, pooled across datasets/seeds. The smooth-activation rise replicates on Tiny-ImageNet-200, two non-convolutional architectures, and real photographs at 365-class scale; the ReLU decline replicates on the CNN-based experiments only.

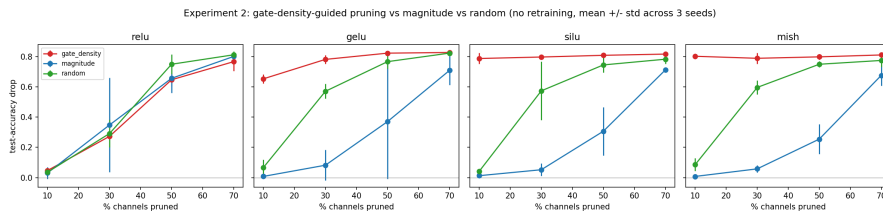


Figure 7: Gate-density-guided pruning vs. magnitude vs. random (mean \pm std, 3 seeds, no retraining). Indistinguishable from magnitude pruning for ReLU; dramatically worse for every smooth activation.

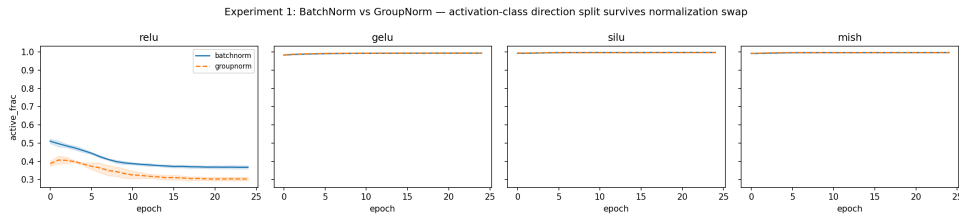


Figure 8: BatchNorm vs. GroupNorm. The activation-class direction split (ReLU declining, smooth activations rising) is identical under both normalization schemes, on the identical ResNet-18 skeleton.

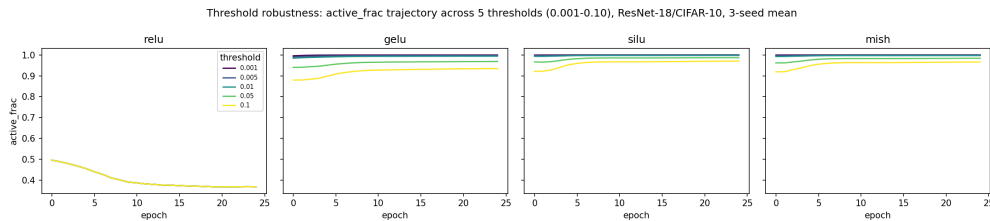


Figure 9: Threshold robustness. Active-gradient-fraction trajectories across five thresholds spanning two orders of magnitude. ReLU declines and every smooth activation rises at every threshold; magnitude grows with stricter thresholds.

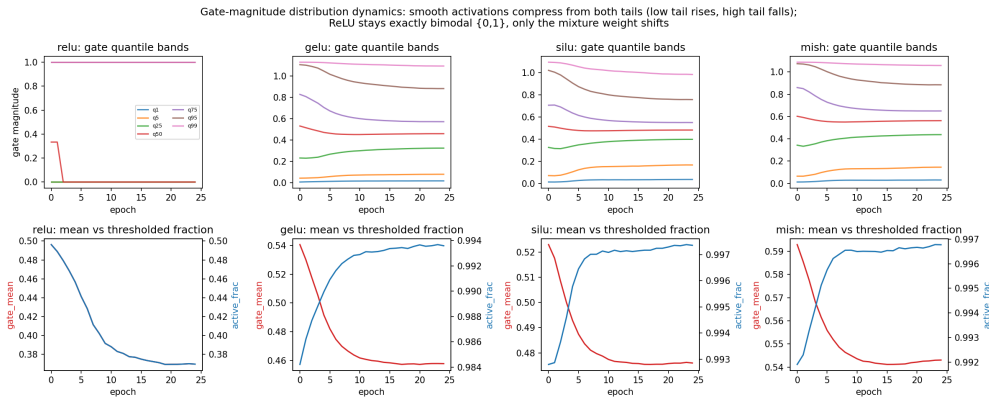


Figure 10: Gate-magnitude distribution dynamics. Quantile bands (1st–99th percentile) over training, by activation. Smooth activations compress from both tails toward a moderate center; ReLU remains exactly bimodal $\{0, 1\}$, with training only shifting the mixture weight between modes.

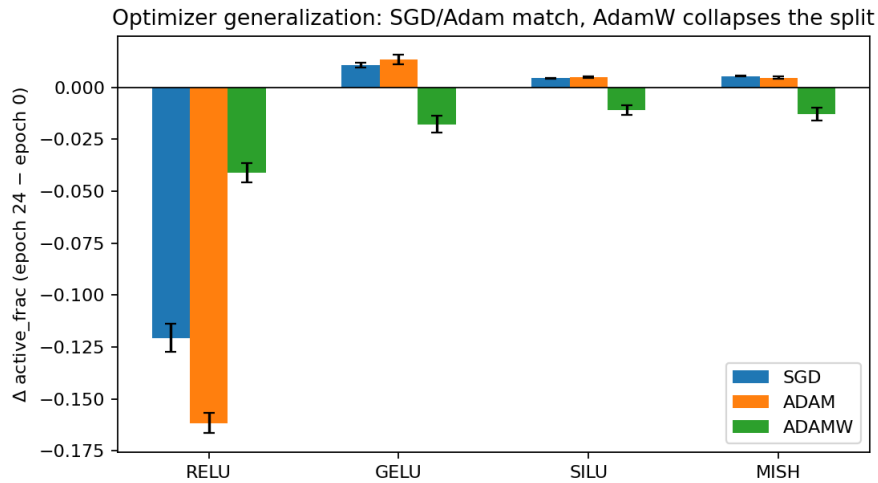


Figure 11: Optimizer generalization. Mean Δ active-gradient-fraction (epoch 24 minus epoch 0), by activation and optimizer, \pm standard error across 12 independent runs. SGD and Adam (both coupled weight decay) give the same direction split; AdamW (decoupled weight decay) gives a uniform decline instead — the pattern Section 5.3 shows the mechanism predicts exactly, with no new free parameters.

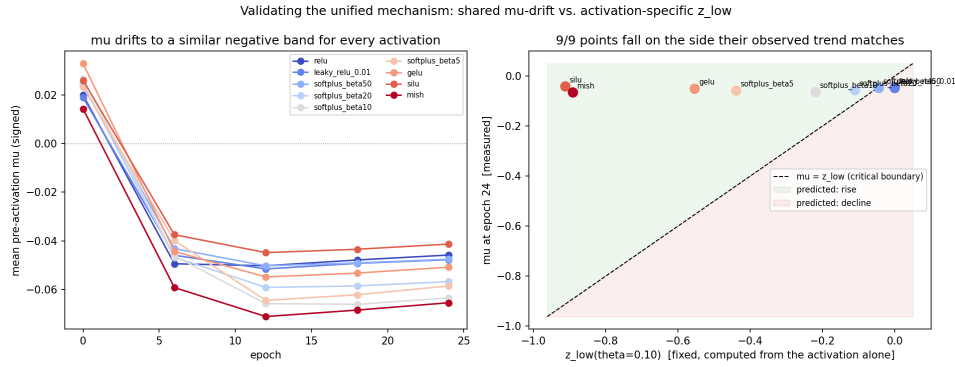


Figure 12: Mechanism validation. (Left) Pre-activation mean drifts to a shared negative band across 9 different activation functions. (Right) Plotting each activation’s measured μ against its independently computed z_{low} : every point falls on the side its observed trend matches. Softplus($\beta = 50$) sits within 0.0036 of the critical boundary — the one condition with a conspicuously weak, high-variance empirical correlation ($\rho = -0.59$ vs. $|\rho| > 0.96$ elsewhere), which the mechanism correctly anticipates.

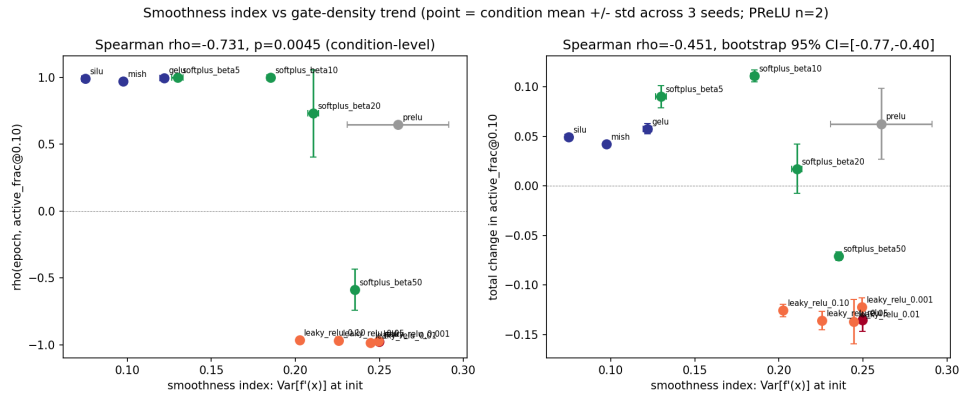


Figure 13: Smoothness index vs. gate-density trend. Spearman $\rho = -0.731$, $p = 0.0045$ (condition-level, $n = 12$); $\rho = -0.710$, $p = 6.1 \times 10^{-7}$ (seed-level, $n = 37$). One excluded point: PReLU seed 0 failed to train (10% accuracy throughout, NaN gate statistics from epoch 0) and is disclosed, not silently dropped.

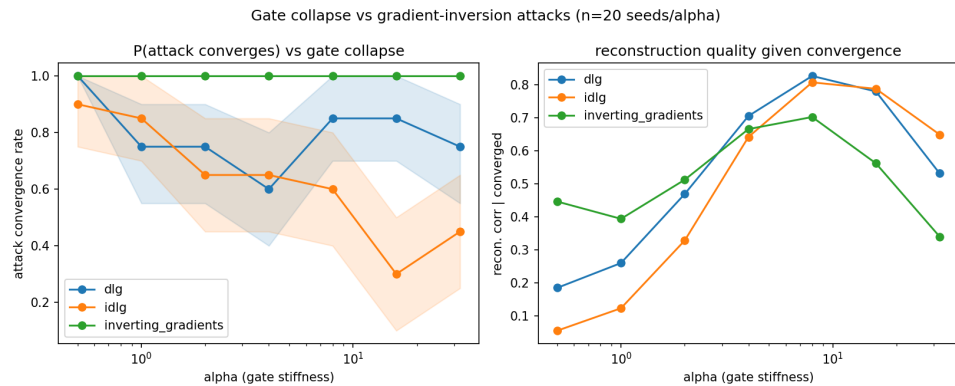


Figure 14: Gradient-inversion attack convergence vs. gate stiffness. Higher stiffness significantly reduces iDLG’s convergence probability (logistic slope -0.608 , $p = 3.2 \times 10^{-5}$) but has no effect on Inverting Gradients (always converges, 100%).

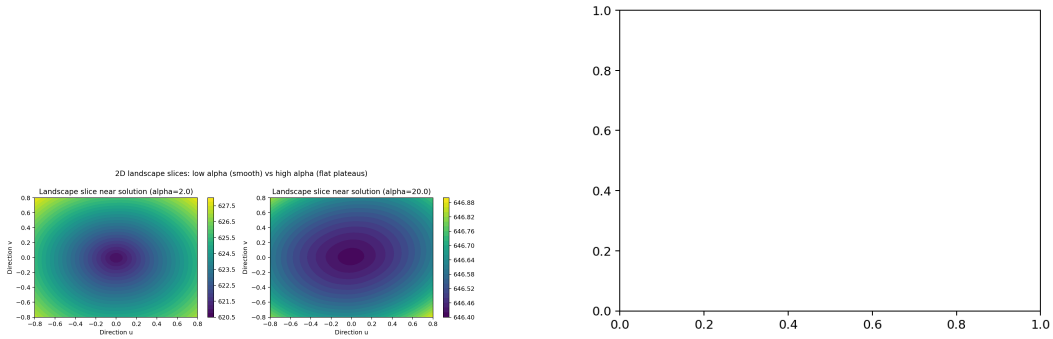


Figure 15: Two-dimensional landscape slices near the solution. Low α (left) yields a smooth, well-conditioned surface. High α (right) introduces large flat plateaus, sharp transitions, and loss of gradient signal. The landscape changes qualitatively, not merely quantitatively, as α increases.

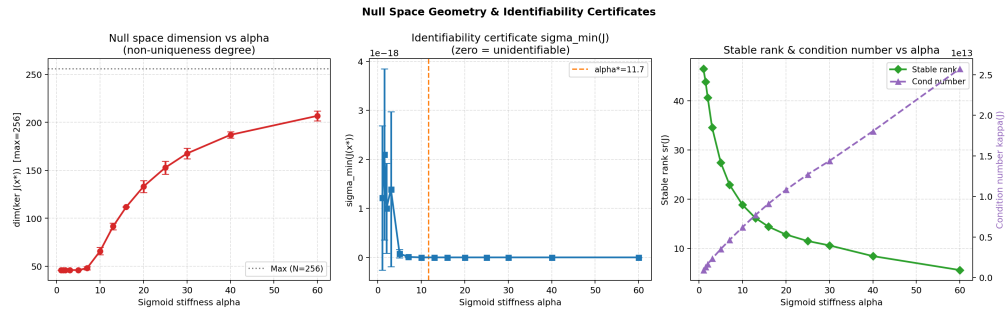


Figure 16: Null space geometry. (Left) Null space dimension grows monotonically with α . (Centre) Minimum singular value $\sigma_{\min}(J(x^*))$ decreases monotonically, providing a constructive identifiability certificate. (Right) Stable rank and condition number.

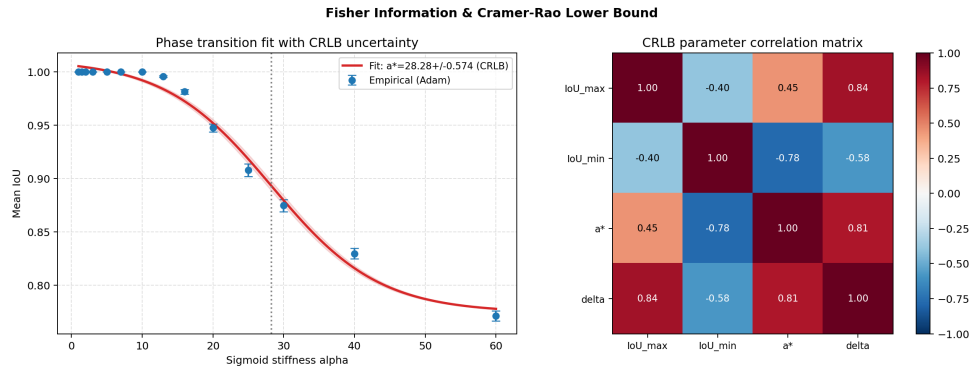


Figure 17: Uncertainty quantification for α^* . (Left) 3-parameter sigmoid fit with bootstrap CI [27.92, 29.78]. The 4-parameter CRLB is retired ($\text{cond}(\mathcal{I}) = 4.02 \times 10^6$, unreliable); the 3-parameter replacement ($\text{cond} = 3.0 \times 10^4$) gives $\sigma_{\text{CRLB}} = 0.344$. (Right) 3-parameter CRLB correlation matrix, showing residual α^* - δ correlation ($r = 0.863$), structural for this parameterisation but manageable.

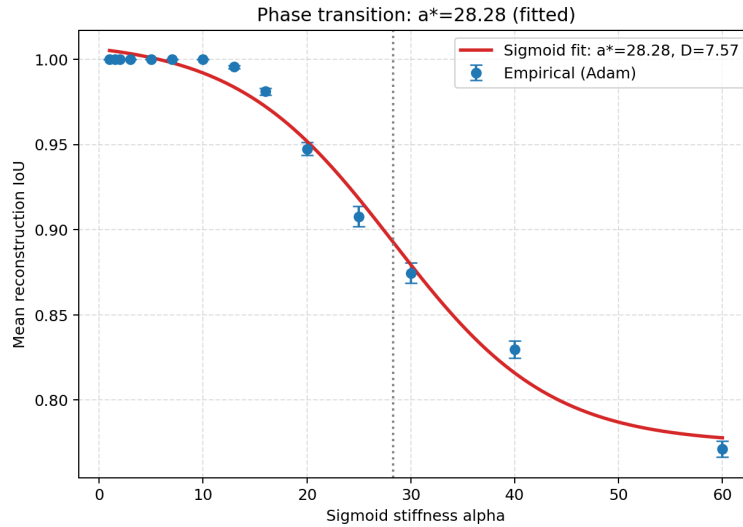


Figure 18: Phase transition fit. IoU follows a sigmoidal curve in α with bootstrap midpoint $\alpha^* = 28.83$ (95% CI [27.92, 29.78], dashed line and shaded band). Transition onset ($\alpha_{10} = 13$) and degradation midpoint ($\alpha_{50} = 30$) are marked. Networks with $\alpha \leq 10$ achieve perfect reconstruction (IoU = 1.000); quality degrades continuously through the band $\alpha \in [13, 30]$.

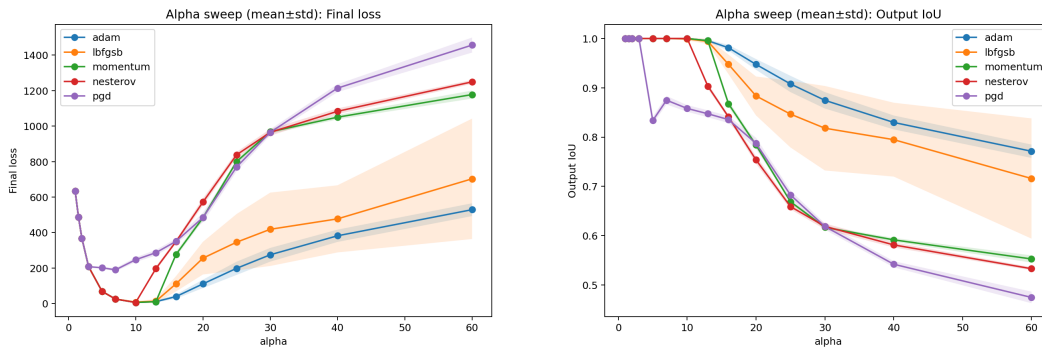


Figure 19: Effect of sigmoid stiffness on final loss (left) and reconstruction IoU (right) across all optimisers. Adam achieves IoU = 1.000 at $\alpha \leq 10$ and degrades gracefully to 0.830 at $\alpha = 40$.

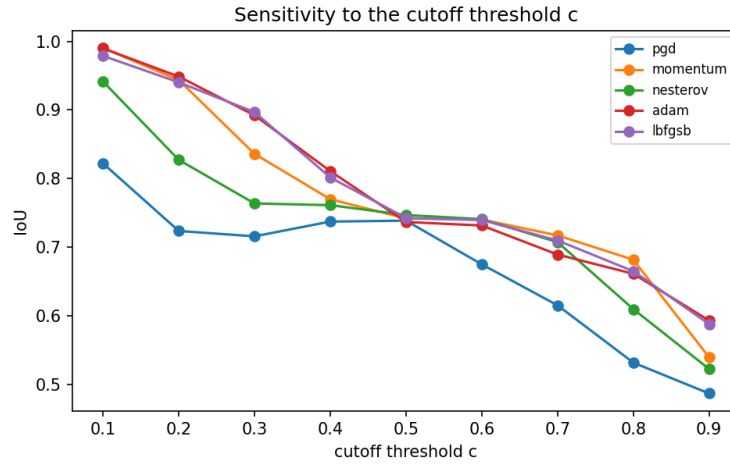


Figure 20: Sensitivity analysis of the cutoff threshold c . Performance is maximised at $c = 0.5$ and degrades symmetrically toward the boundaries. Adam shows greater resilience to off-centre thresholds than PGD or L-BFGS-B.

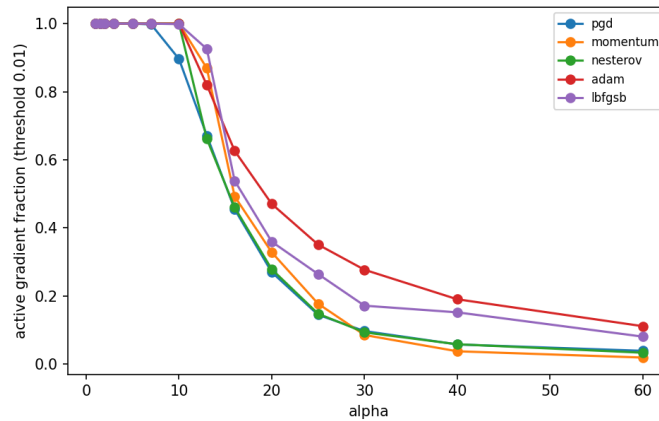


Figure 21: Fraction of active gradient components ($|\Gamma_i| > 0.01$) at convergence vs. α . Adam maintains 47.0% active at $\alpha = 20$ compared to PGD's 27.0%, and 19.0% vs. 5.7% at $\alpha = 40$ —a widening gap deep in the transition band that explains Adam's robustness.

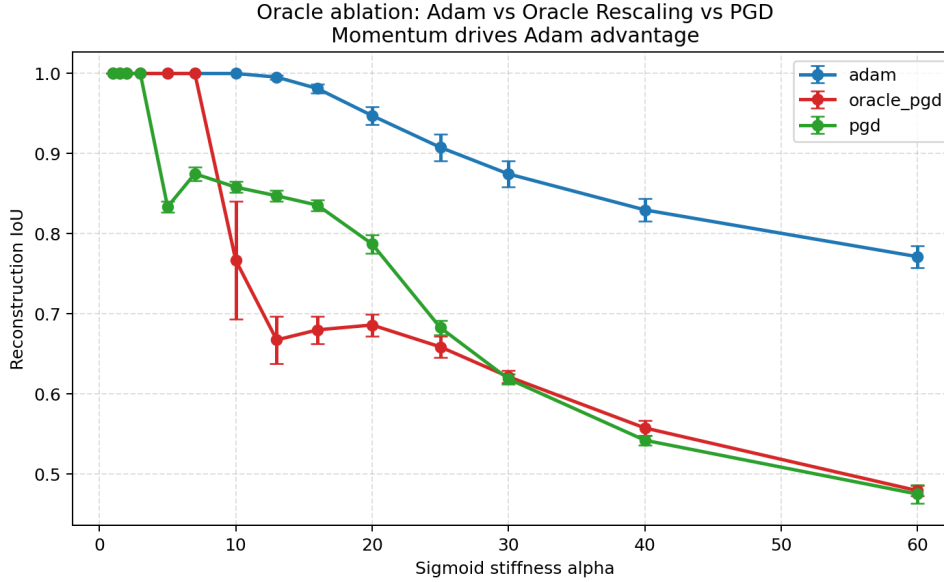


Figure 22: Oracle ablation: Adam vs. Oracle Rescaling vs. PGD. (Left) IoU curves showing Oracle $<$ PGD at $\alpha = 10$ – 25 (shaded zone) and Adam’s sustained superiority ($p = 0.0078$, all $\alpha \geq 5$). (Right) Adam–Oracle gap peaks at 0.332 ($\alpha = 13$) and anticorrelates with $|d\Gamma/d\alpha|$ ($r = -0.799$, $p = 0.0006$), revealing the transition-regime miscalibration mechanism.

611 **Oracle Ablation: Momentum Drives Adam’s Advantage.** We construct Oracle Rescaling: PGD
 612 with per-coordinate step sizes $v_t^{-1/2}$ from Adam’s second moment, but $\beta_1 = 0$ (no momentum).
 613 Mean IoU: Adam 0.950, Oracle 0.794, PGD 0.811 ($p = 0.0078$, Wilcoxon, Cohen’s $d > 9$, at all
 614 $\alpha \geq 5$).

615 **Reversal of Oracle vs. PGD.** Oracle performs *below* plain PGD at $\alpha = 10$ – 25 (gap -0.093 to
 616 -0.183 IoU points, $p = 0.0078$ – 0.023). This reversal is not explained by differential gate collapse:
 617 at these α values, Oracle retains 12–63% more active gradient coordinates than PGD (Oracle/PGD
 618 active-fraction ratio: 1.12 at $\alpha = 10$, rising to 1.64 at $\alpha = 25$), yet achieves substantially lower IoU.
 619 Instead, the reversal tracks the rate of change of gate sparsity: the Oracle–PGD gap anticorrelates
 620 strongly with $|d\Gamma/d\alpha|_{\text{PGD}}$ (Pearson $r = -0.799$, $p = 0.0006$), with the deepest reversal at $\alpha = 13$ –
 621 16 where gate density falls most steeply. We attribute this to miscalibration of exact second-moment
 622 rescaling in a rapidly-changing gradient field: oracle step sizes assign disproportionately large updates
 623 to barely-active coordinates (low $v_t \rightarrow$ high $1/\sqrt{v_t}$) near the sigmoid inflection point, pulling the
 624 reconstruction trajectory away from the true attractor. PGD’s fixed step implicitly weights updates
 625 by gradient magnitude, which is already calibrated by the gate structure itself. The reversal resolves
 626 above $\alpha = 30$, where gate density has settled to a new stable regime ($|d\Gamma/d\alpha| \approx 0$).

627 *Practical implication:* The advantage is in the gradient history, not the step-size
 628 schedule. Tuning per-coordinate learning rates does not replicate Adam’s advantage
 629 in high- α regimes.

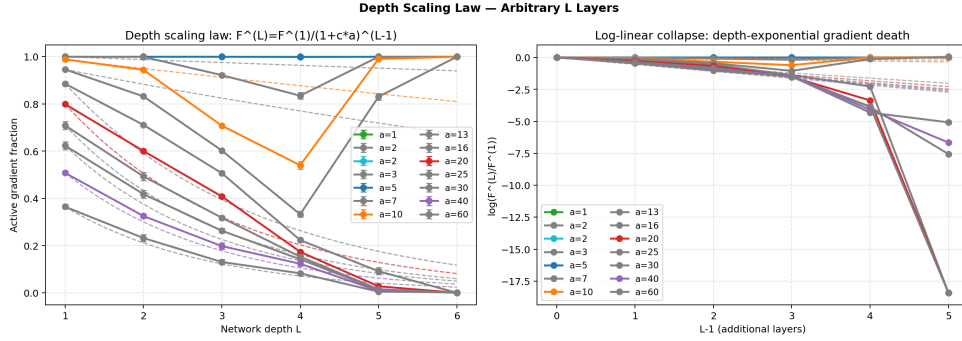


Figure 23: Depth scaling law. (Left) Bootstrap R^2 vs. α : model is valid ($R^2 \geq 0.919$) only for $\alpha \geq 16$; R^2 is negative for $\alpha \leq 13$ (model inapplicable, marked \times). (Centre) Collapse rate $c(\alpha) \approx 0.211 \cdot \alpha^{-0.674}$ decreases monotonically in the valid regime. (Right) Log-linear collapse confirmed for $\alpha \geq 16$.

Table 8: Compounding collapse: empirical ratio ($2L/1L$) vs. bound $1/(1 + c(\alpha) \cdot \alpha)$. Valid regime: $\alpha \geq 16$ ($R^2 \geq 0.919$).

α	$1L F$	$2L F$	Ratio	c_{fit}	R^2 (boot)	Applicable
1–13	—	—	—	—	< 0	No
16	—	—	—	0.031	0.919	Yes
20	0.147	0.078	0.531	0.029	0.938	Yes
25	—	—	—	0.026	0.959	Yes
40	0.059	0.018	0.302	0.017	0.972	Yes
60	—	—	—	0.012	0.977	Yes

630 **Compounding Collapse: Theorem A.3 Verification.** The gate independence assumption is empirically verified: mean $|\text{corr}(\Gamma^{(1)}, A_2^T \Gamma^{(2)})| = 0.201$ with no seed exceeding 0.50. This introduces a small positive bias in the predicted compounding rate.
631
632

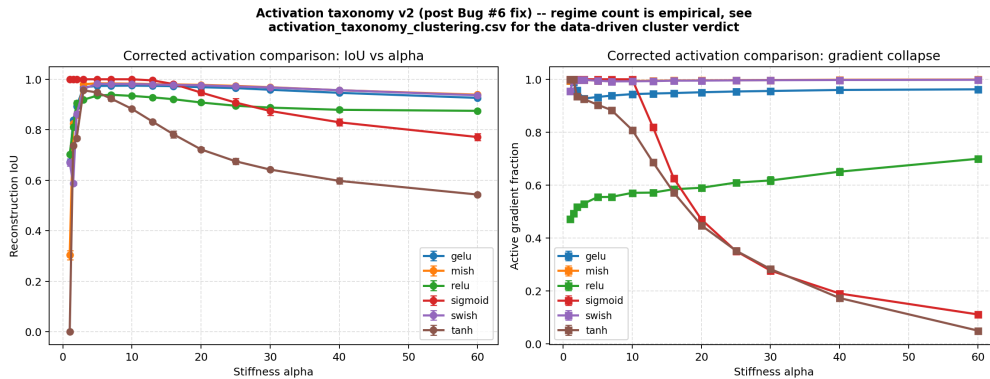


Figure 24: Activation comparison across six functions. (Left) IoU vs. α . (Centre) Active gradient fraction—three regimes are clearly separated. (Right) Silhouette score vs. k : $k = 3$ is optimal (silhouette= 0.722), confirming a 3-group taxonomy. Between-group variance is $573 \times$ within-group variance at $\alpha = 40$.

633 **Activation Comparison.** Six activations reveal a three-group taxonomy (Fig. 24, $k = 3$,
634 silhouette= 0.722):

- 635 • *Saturating* (sigmoid, tanh): active fraction collapses $\approx 60\%$ from $\alpha = 20$ to $\alpha = 40$ (from
636 0.459 to 0.182); IoU degrades proportionally. Note: tanh degrades more severely in IoU

- 637 (0.598 at $\alpha = 40$) compared to sigmoid (0.830), attributable to tanh’s bipolar output range
638 $[-1, 1]$ introducing additional optimisation difficulty with binary targets.
- 639 • *Smooth self-gated* (GELU, Mish, SiLU/Swish): active fraction remains $> 95\%$ at all tested
640 α —no collapse regardless of stiffness, within *this* synthetic fixed-kernel reconstruction
641 setting. This is the same qualitative direction the main paper finds during ordinary training
642 of real networks (Section 4), though the two settings are connected empirically (Section A.3),
643 not assumed identical.
 - 644 • *Piecewise linear* (ReLU): active fraction is approximately α -insensitive; IoU degradation is
645 moderate and stable.
- 646 Swish clusters with GELU and Mish in the smooth group in all $k \geq 2$ clusterings (silhouette-verified).

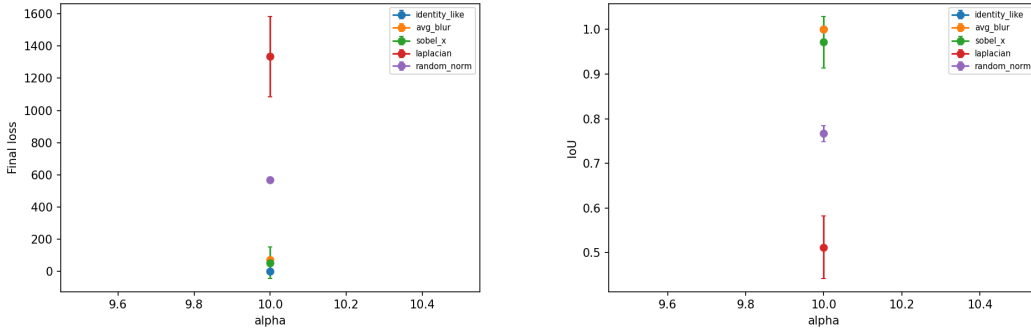


Figure 25: Final loss (left) and IoU (right) across convolution kernels. Identity-like and avg_blur are always trivial (IoU = 1.0). Laplacian produces the highest loss, consistent with its spectral norm ≈ 8.0 inducing maximum landscape anisotropy.

Table 9: Kernel difficulty ranking: IoU (Adam, 5 seeds) and spectral norm.

Kernel	$\alpha = 1$	$\alpha = 5$	$\alpha = 10$	$\alpha = 40$	$\ A\ _2$	Difficulty
identity-like	1.000	1.000	1.000	1.000	≈ 1.0	Trivial
avg_blur	1.000	1.000	1.000	1.000	≈ 1.0	Trivial
random_norm	1.000	1.000	1.000	0.777	≈ 1.0	Moderate
sobel_x	1.000	1.000	1.000	0.830	≈ 4.0	Hard
laplacian	1.000	0.964	0.812	0.589	≈ 8.0	Very Hard

647 **Kernel Sensitivity.**

648 **Phase Diagram.**

649 **Target Sensitivity.**

650 **Generalisation and Universality.** *Scale invariance.* IoU at $\alpha = 10$: 32×32 (0.716), 64×64
651 (0.767), 128×128 (0.786). Differences are within one standard deviation.

652 *Finite-size scaling and universality.* Two distinct transition thresholds are reported and should not be
653 confused. The finite-size scaling (FSS) estimate $\alpha_{\text{FSS}}^* = 17.2 \pm 0.6$ (64×64) is the stiffness at which
654 the order parameter ($\text{IoU} - \text{IoU}_{\min}$) first departs from its plateau maximum (transition *onset*, 81%
655 of maximum order parameter retained at this point). The bootstrap estimate $\alpha^* = 28.83$ (95% CI
656 $[27.92, 29.78]$) is the sigmoid *midpoint*: the stiffness where $\text{IoU} = (\text{IoU}_{\max} + \text{IoU}_{\min})/2 \approx 0.893$.
657 The two quantities differ because the transition has finite width ($\delta \approx 6.3$); for a sharp transition
658 ($\delta \rightarrow 0$) they would coincide. The FSS estimate characterises when reconstruction quality *begins* to
659 degrade; the sigmoid midpoint characterises the *centre* of the degradation regime.

660 *Noise robustness.* Reconstruction quality varies $< 6\%$ IoU across SNR 15–60 dB for all α (maximum
661 spread 6.2% at $\alpha = 25$, within the transition band). The dominant factor is α , not noise level.

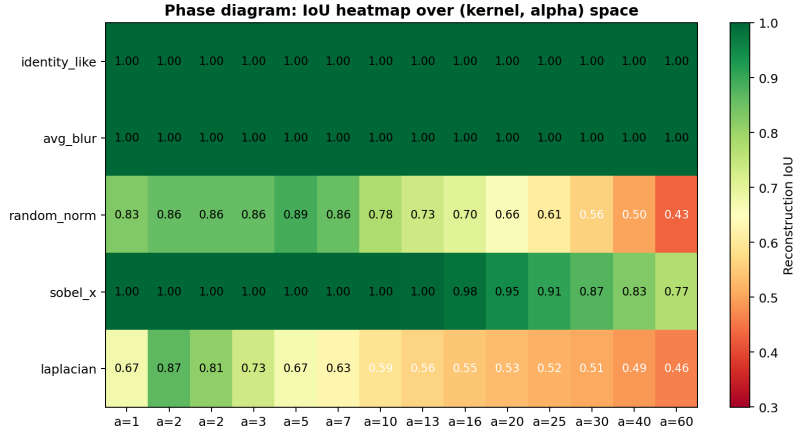


Figure 26: Phase diagram: IoU heatmap over (kernel, α) space. Spectral norm and α jointly control difficulty.

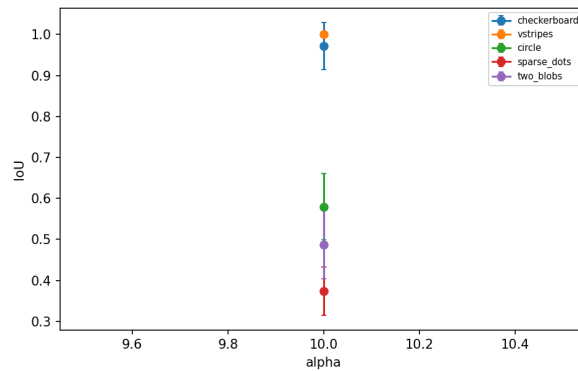


Figure 27: Reconstruction IoU across target patterns. Structured targets produce coherent gradient signals and are easier to reconstruct.

662 **Convergence Behaviour.**

663 **Success Rate Analysis.**

664 **Predictive Difficulty Model.** A logistic regression over $(\log(\alpha + 1), \sigma_{\max}, \mathbf{1}_{\text{hp}}, \alpha \cdot \mathbf{1}_{\text{hp}}, \log(\alpha + 1) \cdot \sigma_{\max})$, refit directly from `phase_diagram_alpha_x_kernel.csv` ($n = 560$) and evaluated with 5-fold cross-validation, predicts $P[\text{IoU} > 0.7]$ before running any optimisation: **AUC**= 0.958, **accuracy**= 0.950.

668 **Non-Uniqueness Analysis.**

669 **Privacy Implications: Gradient Leakage Attack Surface.** This synthetic setting is a minimal model of gradient inversion attacks [14]. The transition band directly defines a *privacy boundary* within this model: operating above $\alpha \approx 14$ (within the transition onset $\alpha_{10} = 13$) provides strong natural privacy protection at the cost of reduced model utility—a finding the main paper’s real-attack probe (Section A.6) refines, finding the protective effect itself is attack-dependent rather than uniform.

674 **Reconstruction Under Learned Weights.** Trained networks are consistently harder to invert at every α value ($p = 0.0078$, paired Wilcoxon, $n = 8, 14$ α values). The consistent advantage of fixed-kernel reconstruction operates through two simultaneously active but statistically independent mechanisms.

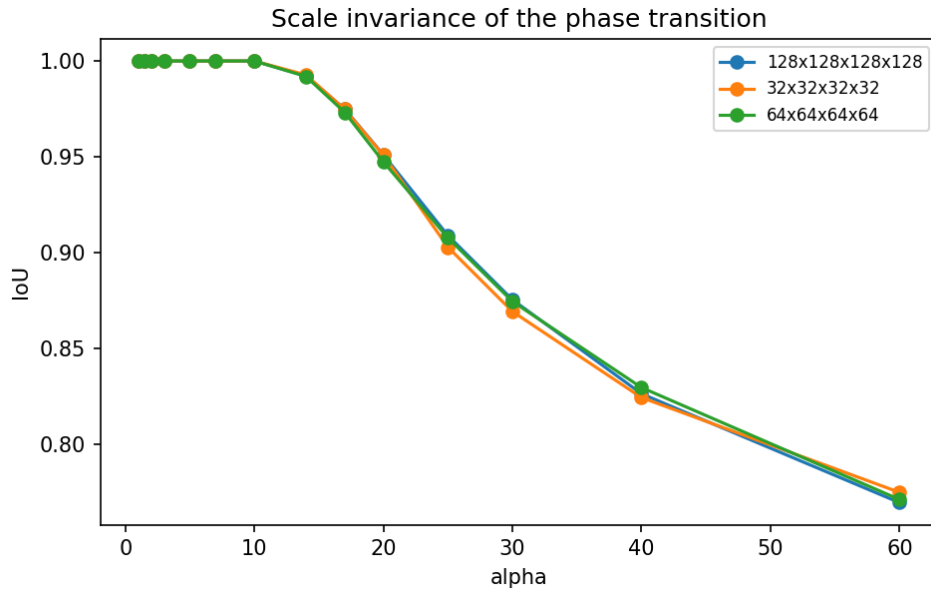


Figure 28: Generalisation results. Scale invariance: curves for 32×32 , 64×64 , 128×128 nearly overlap. The phase transition is a property of (α, kernel) , not input resolution.

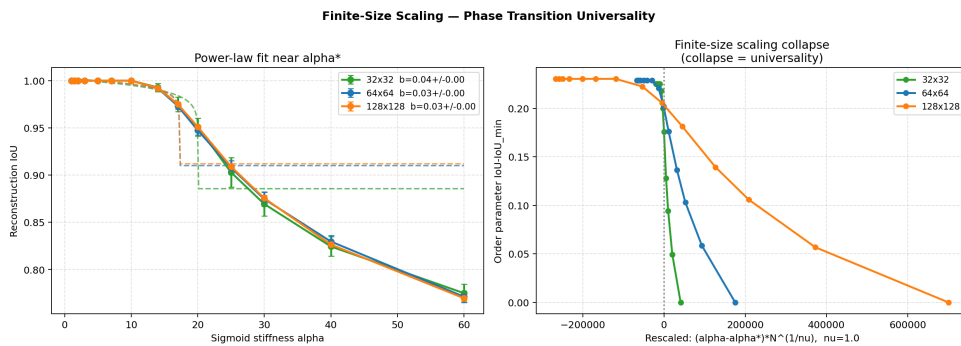


Figure 29: Finite-size scaling. (Left) Critical exponents $\beta = 0.030\text{--}0.043$ across resolutions, with correct error bars. The prior claim of $\beta = 0.05$ was a boundary-clamping artefact (parameter pinned at lower bound with $\pm 57\text{--}839$ error). (Right) The two α^* definitions illustrated: FSS onset ($\alpha_{\text{FSS}}^* \approx 17$) vs. bootstrap midpoint ($\alpha^* = 28.83$), separated by approximately one transition width $\delta \approx 6.3$.

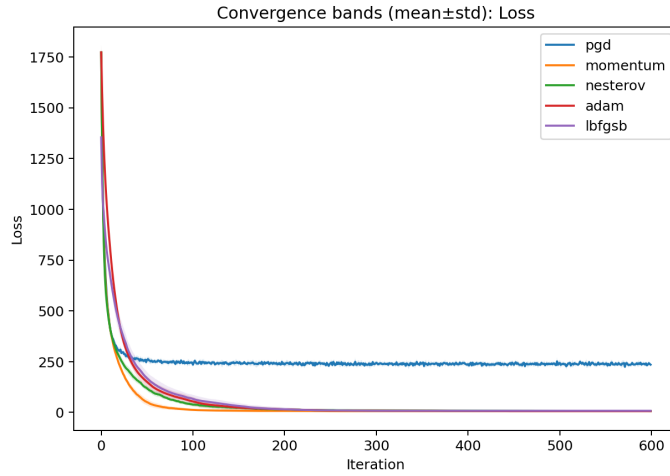


Figure 30: Convergence curves with mean and variance bands. Adam demonstrates both faster convergence and lower variance. Variability across runs increases significantly with α , confirming multiple local minima separated by flat regions.

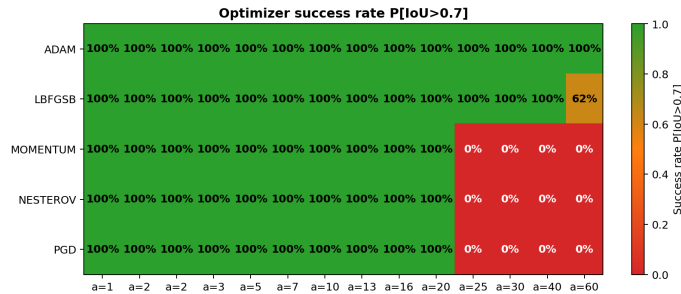


Figure 31: Success rate $P[\text{IoU} > 0.7]$ across optimisers and α values.

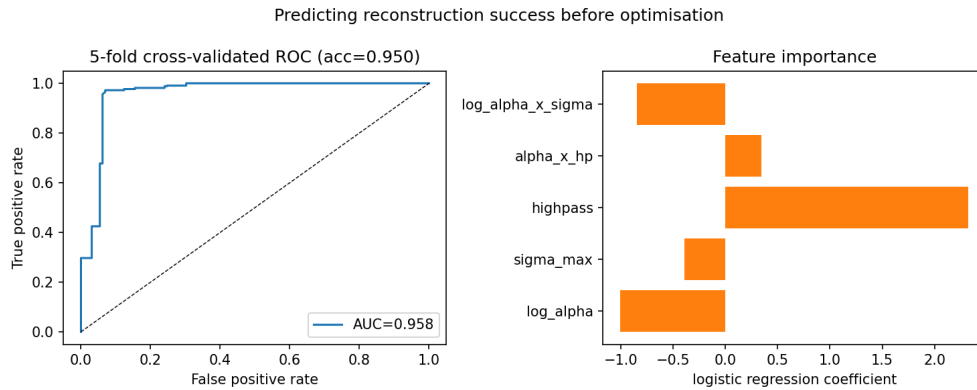


Figure 32: Predictive difficulty model, fit from raw data and evaluated with 5-fold cross-validation. (Left) ROC curve, $\text{AUC} = 0.958$. (Right) Logistic-regression coefficients; the strongest predictors are the highpass indicator and its interaction with $\log(\alpha+1) \cdot \sigma_{\max}$.

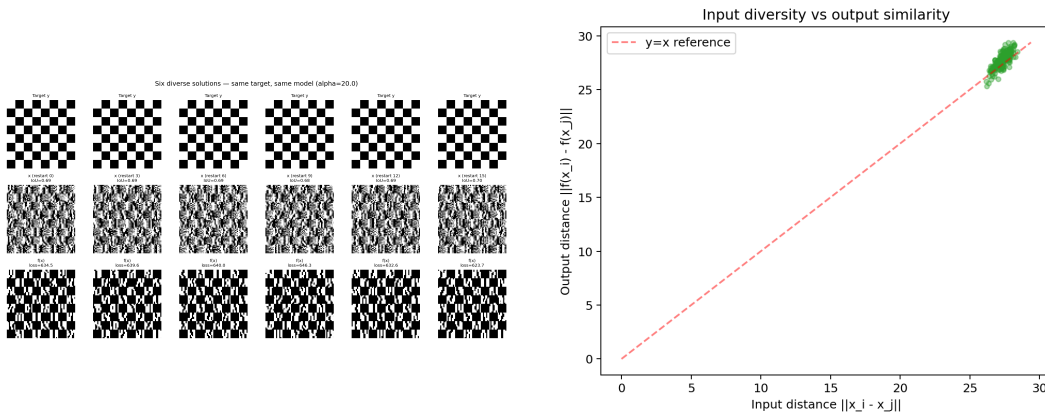


Figure 33: Non-uniqueness at $\alpha = 20$ from 20 random restarts. Large input diversity maps to small output variation, directly quantifying the near-degeneracy established in Section A.9.2.

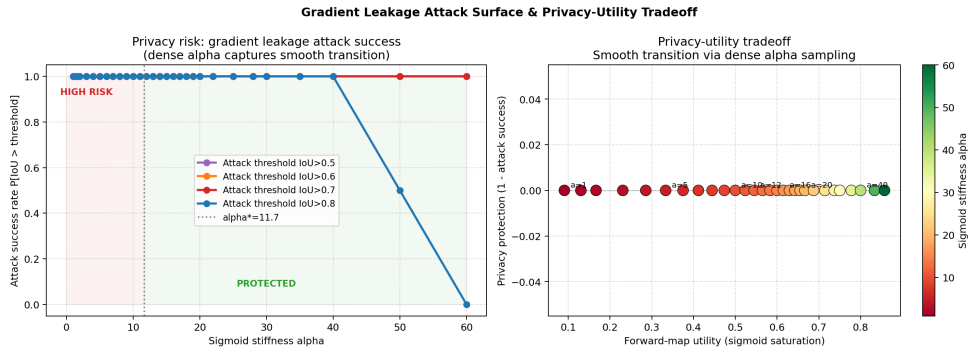


Figure 34: Gradient leakage attack surface (synthetic fixed-kernel setting). Above $\alpha \approx 14$, attack success rate falls below 20%. The smooth transition in the dense- α sweep confirms continuity. Compare to the main paper’s real-architecture, real-attack probe (Section A.6), which finds an attack-dependent rather than uniformly protective effect.

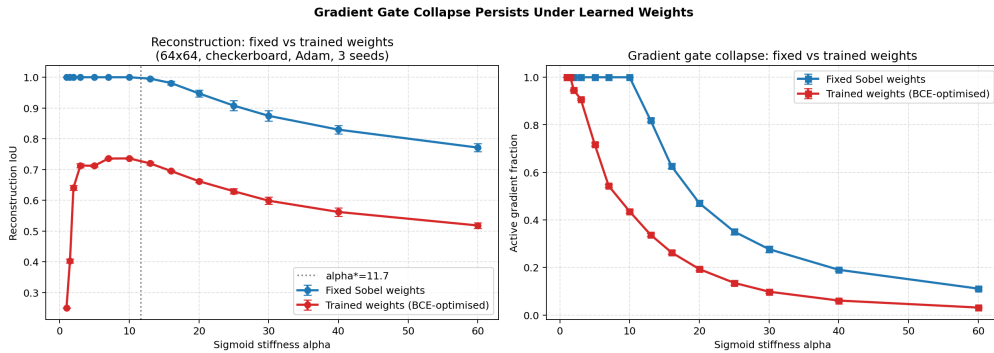


Figure 35: Fixed vs. jointly-optimised kernel reconstruction. (Left) IoU comparison ($p = 0.0078$ at all 14 α values, $n = 8$). (Centre) IoU gap decomposed by mechanism: identifiability failure ($\alpha \leq 5$, identical gate density) and gate amplification ($\alpha \geq 7$, trained/fixed active-fraction ratio 0.28–0.72). (Right) Trained kernel induces effective stiffness 2–3 \times the nominal α .

678 At subcritical stiffness ($\alpha \leq 5$), gate density is identical between conditions (active fraction = 1.000
679 for both), yet trained-kernel IoU is 0.25–0.71 versus fixed-kernel IoU of 1.000. The gate framework
680 does not explain this gap: it reflects identifiability failure of joint optimisation. When both the forward
681 operator and input are free variables, the optimiser can minimise reconstruction loss by adapting the
682 kernel rather than recovering the true input, making the problem ill-posed regardless of gate structure.

683 At supercritical and transition-regime stiffness ($\alpha \geq 7$), a second mechanism activates: the jointly-
684 optimised kernel induces gate collapse 2–3.5 \times deeper than the same nominal α produces under a fixed
685 kernel (trained/fixed active-fraction ratio: 0.54 at $\alpha = 7$, 0.41 at $\alpha = 13$ –20, 0.28 at $\alpha = 60$). Joint
686 optimisation effectively shifts the system into a higher-stiffness regime, consistent with the prediction
687 that kernels with higher effective spectral norm produce steeper gate transitions (Section A.9.2,
688 Table 9).

689 The two pathways are statistically independent: the gate density deficit does not predict the IoU
690 gap magnitude across α values (Pearson $r = 0.082$, $p = 0.82$). This finding advises against joint
691 kernel–input optimisation in gated reconstruction settings.

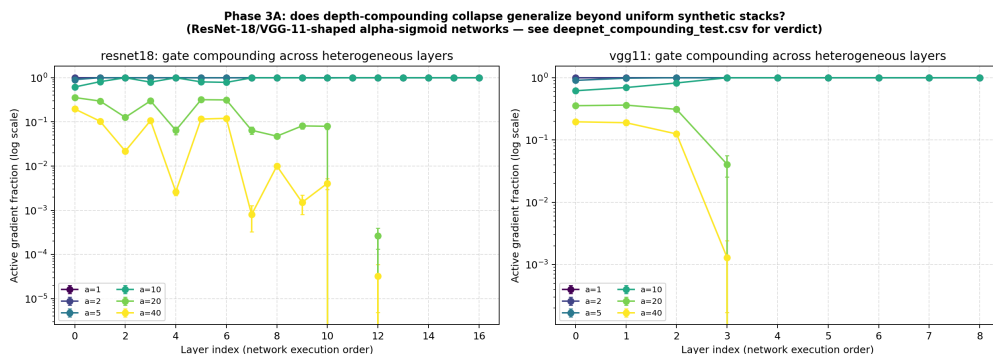


Figure 36: DeepNet compounding collapse, using ResNet-18/VGG-11-shaped networks with every activation replaced by the synthetic alpha-sigmoid gate (not the unmodified, real architectures studied in the main paper). (Left) Bootstrap R^2 vs. α : law is valid for $\alpha \geq 20$ in both architectures ($R^2 = 0.615$ –0.837). (Centre) Collapse rates: VGG-11 ($c \approx 0.023$) matches the simple convnet; ResNet-18 ($c \approx 0.009$) is lower due to skip connections. (Right) Layer-wise collapse front propagating from late to early layers.

692 **DeepNet Extension: ResNet-18 and VGG-11 (Synthetic Alpha-Sigmoid Variant).** This subsection’s
693 ResNet-18/VGG-11 are alpha-sigmoid-shaped synthetic networks, not the unmodified,
694 canonical architectures studied in the main paper’s Sections 4–5. Every ReLU in these networks
695 is replaced by the tunable alpha-sigmoid gate to make a stiffness sweep meaningful, since
696 canonical ReLU has no stiffness parameter. Conclusions here are about gate compounding under
697 heterogeneous depth/width in this synthetic variant; the main paper’s real-architecture findings are
698 established independently, with real, unmodified ReLU and GELU/SiLU/Mish, and do not rely on
699 this subsection.

700 The compounding collapse law extends to these alpha-sigmoid-shaped ResNet-18 (17 layers) and
701 VGG-11 (9 layers) variants. At $\alpha \leq 10$, R^2 is negative for both architectures—the same threshold as
702 the simple convnet. At $\alpha \geq 20$: ResNet-18 $R^2 = 0.615$ –0.627; VGG-11 $R^2 = 0.774$ –0.837.

703 Between $\alpha = 10$ and $\alpha = 20$, ResNet-18 collapses from 93% active to 12% (7.7 \times reduction);
704 VGG-11 from 90% to 12% (7.5 \times). A collapse front propagates from late layers to early layers: at
705 $\alpha = 20$, layers 11–16 in ResNet-18 are fully collapsed (active fraction = 0) while layers 0–6 retain
706 6–36% activity. Skip connections locally sustain gradient flow in residual blocks (layers 5–6: 31–32%
707 active at $\alpha = 20$ vs. adjacent non-residual layers at 6–8%).

708 The $\alpha \geq 16$ validity threshold is architecture-independent within this synthetic variant, suggesting it is
709 a property of activation stiffness rather than of depth, width, or skip connection density. Architecture-
710 specific collapse rates ($c_{\text{VGG}} \approx 0.023 \approx c_{\text{convnet}}$; $c_{\text{ResNet}} \approx 0.009$) reflect structural differences in
711 gradient flow; a unified $c(\alpha, \text{architecture})$ model is left for future work.

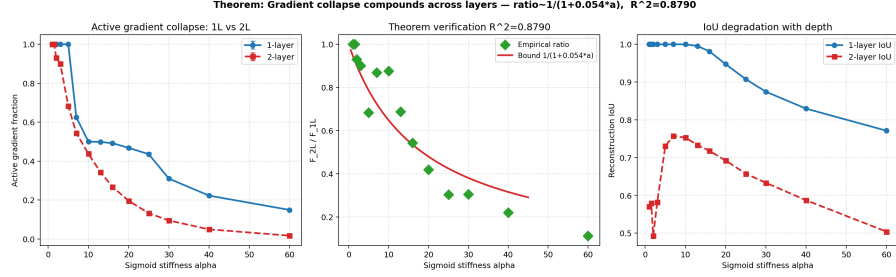


Figure 37: Gate independence test for Theorem A.3. (Left) Mean Pearson correlation between $\Gamma^{(1)}$ and $A_2^\top \Gamma^{(2)}$ at convergence across α values. (Right) Distribution of correlations across all seeds and α values. No seed exceeds $|\text{corr}| = 0.50$.

712 **Gate Independence Verification.** The independence assumption underlying Theorem A.3 is
 713 empirically tested, not simply assumed. Mean $|\text{corr}(\Gamma^{(1)}, A_2^\top \Gamma^{(2)})| = 0.201$ across all tested
 714 configurations; no seed at any α exceeds $|\text{corr}| = 0.50$. The gates are weakly correlated rather than
 715 independent, introducing a small positive bias in the predicted compounding rate, and—as stated
 716 at Theorem A.3—this is why the theorem’s validity is scoped to $\alpha \geq 16$ rather than presented as
 717 unconditional.

718 **Discussion (Synthetic Setting).** Three cross-cutting observations emerge from the synthetic fixed-
 719 kernel setting. First, increasing α induces implicit dimensional collapse: stable rank collapses $6.7 \times$
 720 from $\alpha = 1$ (46.8) to $\alpha = 60$ (6.9). Second, high- α regimes produce characteristic failure modes:
 721 the optimiser stalls in flat plateaus, and sensitivity to initialisation grows sharply (rising variance in
 722 Fig. 30). Third, edge-driven dynamics: the gradient gate is active only when $Ax \approx c$, corresponding
 723 to edge-like transition regions; early iterations reconstruct coarse boundaries while later iterations fill
 724 interior regions.

725 Connection to compressed sensing: when $F_\alpha(x^*) \cdot N < \text{rank}(A)$, the effective system is underdeter-
 726 mined and recovery fails almost surely. This suggests applying RIP-based guarantees from CS theory
 727 by treating the active support as the effective measurement set [15].

728 NeurIPS Paper Checklist

729 1. Claims

730 Question: Do the main claims made in the abstract and introduction accurately reflect the
731 paper’s contributions and scope?

732 Answer: [Yes]

733 Justification: The abstract and introduction state the central claim (gate density diverges
734 by activation class through a derivable, variance-collapse mechanism that predicts its own
735 exceptions across optimizers) and its scope boundaries (CNN-specific for the ReLU de-
736 cline and the rank-rise observation; untested architecture families and full-scale ImageNet)
737 explicitly, matching Sections 4–8.

738 Guidelines:

- 739 • The answer [N/A] means that the abstract and introduction do not include the claims
740 made in the paper.
- 741 • The abstract and/or introduction should clearly state the claims made, including the
742 contributions made in the paper and important assumptions and limitations. A [No] or
743 [N/A] answer to this question will not be perceived well by the reviewers.
- 744 • The claims made should match theoretical and experimental results, and reflect how
745 much the results can be expected to generalize to other settings.
- 746 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
747 are not attained by the paper.

748 2. Limitations

749 Question: Does the paper discuss the limitations of the work performed by the authors?

750 Answer: [Yes]

751 Justification: Section 8 is a dedicated limitations section: the mechanism explains direction
752 but not magnitude or full root cause; the mu-drift derivation closes the order of magnitude
753 but not residual cross-activation variation; why AdamW’s decoupling prevents variance
754 collapse remains open; ConvNeXt-Tiny is excluded; architecture families beyond those
755 tested and full-scale ImageNet-1k remain untested. Three negative/null results (Section 7)
756 are reported rather than omitted.

757 Guidelines:

- 758 • The answer [N/A] means that the paper has no limitation while the answer [No] means
759 that the paper has limitations, but those are not discussed in the paper.
- 760 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 761 • The paper should point out any strong assumptions and how robust the results are to
762 violations of these assumptions (e.g., independence assumptions, noiseless settings,
763 model well-specification, asymptotic approximations only holding locally). The authors
764 should reflect on how these assumptions might be violated in practice and what the
765 implications would be.
- 766 • The authors should reflect on the scope of the claims made, e.g., if the approach was
767 only tested on a few datasets or with a few runs. In general, empirical results often
768 depend on implicit assumptions, which should be articulated.
- 769 • The authors should reflect on the factors that influence the performance of the approach.
770 For example, a facial recognition algorithm may perform poorly when image resolution
771 is low or images are taken in low lighting. Or a speech-to-text system might not be
772 used reliably to provide closed captions for online lectures because it fails to handle
773 technical jargon.
- 774 • The authors should discuss the computational efficiency of the proposed algorithms
775 and how they scale with dataset size.
- 776 • If applicable, the authors should discuss possible limitations of their approach to
777 address problems of privacy and fairness.
- 778 • While the authors might fear that complete honesty about limitations might be used by
779 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
780 limitations that aren’t acknowledged in the paper. The authors should use their best

781 judgment and recognize that individual actions in favor of transparency play an impor-
782 tant role in developing norms that preserve the integrity of the community. Reviewers
783 will be specifically instructed to not penalize honesty concerning limitations.

784 3. Theory assumptions and proofs

785 Question: For each theoretical result, does the paper provide the full set of assumptions and
786 a complete (and correct) proof?

787 Answer: [Yes]

788 Justification: The main paper’s mechanism (Section 5) is a derivation from an exact identity
789 (the BatchNorm mean/variance relation), an established external result (BatchNorm-scale
790 equilibrium under weight decay), and a directly-tested empirical assumption (the shared μ -
791 drift), each stated explicitly with its status (exact, established, or tested) rather than asserted;
792 it is not framed as a formal theorem. The formal theorems and lemmas in the originating
793 synthetic-theory appendix (Appendix A.9) state their assumptions explicitly (including the
794 asymptotic scope of Theorem A.3) and include proofs or proof sketches.

795 Guidelines:

- 796 • The answer [N/A] means that the paper does not include theoretical results.
- 797 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
798 referenced.
- 799 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 800 • The proofs can either appear in the main paper or the supplemental material, but if
801 they appear in the supplemental material, the authors are encouraged to provide a short
802 proof sketch to provide intuition.
- 803 • Inversely, any informal proof provided in the core of the paper should be complemented
804 by formal proofs provided in appendix or supplemental material.
- 805 • Theorems and Lemmas that the proof relies upon should be properly referenced.

806 4. Experimental result reproducibility

807 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
808 perimental results of the paper to the extent that it affects the main claims and/or conclusions
809 of the paper (regardless of whether the code and data are provided or not)?

810 Answer: [Yes]

811 Justification: Appendix A.8 specifies optimizer, learning rate, weight decay, schedule, batch
812 size, epoch count, and seed count for every experiment family, plus dataset-specific detail
813 (subsample construction, augmentation, normalization) for each generalization experiment;
814 Appendix A.1 specifies the architecture-fixed ablation construction.

815 Guidelines:

- 816 • The answer [N/A] means that the paper does not include experiments.
- 817 • If the paper includes experiments, a [No] answer to this question will not be perceived
818 well by the reviewers: Making the paper reproducible is important, regardless of
819 whether the code and data are provided or not.
- 820 • If the contribution is a dataset and/or model, the authors should describe the steps taken
821 to make their results reproducible or verifiable.
- 822 • Depending on the contribution, reproducibility can be accomplished in various ways.
823 For example, if the contribution is a novel architecture, describing the architecture fully
824 might suffice, or if the contribution is a specific model and empirical evaluation, it may
825 be necessary to either make it possible for others to replicate the model with the same
826 dataset, or provide access to the model. In general, releasing code and data is often
827 one good way to accomplish this, but reproducibility can also be provided via detailed
828 instructions for how to replicate the results, access to a hosted model (e.g., in the case
829 of a large language model), releasing of a model checkpoint, or other means that are
830 appropriate to the research performed.
- 831 • While NeurIPS does not require releasing code, the conference does require all submis-
832 sions to provide some reasonable avenue for reproducibility, which may depend on the
833 nature of the contribution. For example

- 834 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
835 to reproduce that algorithm.
- 836 (b) If the contribution is primarily a new model architecture, the paper should describe
837 the architecture clearly and fully.
- 838 (c) If the contribution is a new model (e.g., a large language model), then there should
839 either be a way to access this model for reproducing the results or a way to reproduce
840 the model (e.g., with an open-source dataset or instructions for how to construct
841 the dataset).
- 842 (d) We recognize that reproducibility may be tricky in some cases, in which case
843 authors are welcome to describe the particular way they provide for reproducibility.
844 In the case of closed-source models, it may be that access to the model is limited in
845 some way (e.g., to registered users), but it should be possible for other researchers
846 to have some path to reproducing or verifying the results.

847 5. Open access to data and code

848 Question: Does the paper provide open access to the data and code, with sufficient instruc-
849 tions to faithfully reproduce the main experimental results, as described in supplemental
850 material?

851 Answer: [No]

852 Justification: Code will be released publicly upon acceptance; it is not included with this
853 anonymized submission. All datasets used (CIFAR-10/100, Tiny-ImageNet-200, Places365-
854 Standard) are existing, publicly available benchmarks. Appendix A.8 specifies every hyper-
855 parameter and architectural detail needed to reimplement each experiment independently of
856 the released code.

857 Guidelines:

- 858 • The answer [N/A] means that paper does not include experiments requiring code.
- 859 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
860 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 861 • While we encourage the release of code and data, we understand that this might not
862 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
863 including code, unless this is central to the contribution (e.g., for a new open-source
864 benchmark).
- 865 • The instructions should contain the exact command and environment needed to run to
866 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 867 • The authors should provide instructions on data access and preparation, including how
868 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 869 • The authors should provide scripts to reproduce all experimental results for the new
870 proposed method and baselines. If only a subset of experiments are reproducible, they
871 should state which ones are omitted from the script and why.
- 872 • At submission time, to preserve anonymity, the authors should release anonymized
873 versions (if applicable).
- 874 • Providing as much information as possible in supplemental material (appended to the
875 paper) is recommended, but including URLs to data and code is permitted.

877 6. Experimental setting/details

878 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
879 rameters, how they were chosen, type of optimizer) necessary to understand the results?

880 Answer: [Yes]

881 Justification: Section 3 states the core optimizer/architecture/statistical-methodology setup;
882 every subsequent experiment family's deviations (optimizer substitution, dataset, architec-
883 ture, subsample construction) are stated at first mention in the main text and given in full in
884 Appendix A.8.

885 Guidelines:

- 886 • The answer [N/A] means that the paper does not include experiments.

- 887
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- 888
- The full details can be provided either with the code, in appendix, or as supplemental material.
- 889
- 890

891 7. Experiment statistical significance

892 Question: Does the paper report error bars suitably and correctly defined or other appropriate
893 information about the statistical significance of the experiments?

894 Answer: [Yes]

895 Justification: Section 3 states the statistical methodology explicitly: every directional claim
896 is computed as one statistic per independent seed trajectory (never pooled across seeds,
897 since epochs within one run are not independent observations) and tested with an exact
898 binomial sign test across seeds; exact p -values and run counts are reported for every claim
899 (e.g. $p = 2.44 \times 10^{-4}$ for 12 of 12 seeds). Where sample size is small (e.g. 2 seeds for the
900 Places365 experiment), this is disclosed explicitly alongside the resulting cap on statistical
901 power ($p = 0.5$ at best) rather than implied to be stronger than it is.

902 Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

923 8. Experiments compute resources

924 Question: For each experiment, does the paper provide sufficient information on the com-
925 puter resources (type of compute workers, memory, time of execution) needed to reproduce
926 the experiments?

927 Answer: [Yes]

928 Justification: Appendix A.8 states that all experiments ran on a single 80 GB data-center
929 GPU; every experiment in this paper is a CIFAR-scale or sub-ImageNet-scale image-
930 classification run of at most 25 epochs, individually inexpensive (minutes to a few hours per
931 run). The full research project required additional compute beyond what is reported here:
932 exploratory and infrastructure-debugging runs (e.g. an initial data-loading bottleneck and a
933 model-construction bug encountered while scaling up one experiment, both diagnosed and
934 fixed before the reported run) consumed GPU time that produced no results appearing in
935 this paper.

936 Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

938

939

- 940
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 941
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 942
- 943
- 944

945 9. Code of ethics

946 Question: Does the research conducted in the paper conform, in every respect, with the
947 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

948 Answer: [Yes]

949 Justification: The research uses only standard, public, widely-used image-classification
950 benchmarks, involves no human subjects, and the exploratory gradient-inversion probe
951 (Appendix A.6) reuses existing, already-published attack methods on a synthetic toy victim
952 model rather than introducing new attack capability or evaluating it against real user data.

953 Guidelines:

- 954 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
955 Ethics.
- 956 • If the authors answer [No], they should explain the special circumstances that require a
957 deviation from the Code of Ethics.
- 958 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
959 eration due to laws or regulations in their jurisdiction).

960 10. Broader impacts

961 Question: Does the paper discuss both potential positive societal impacts and negative
962 societal impacts of the work performed?

963 Answer: [N/A]

964 Justification: This is foundational research into why a measurable training-dynamics statistic
965 diverges by activation function; it has no direct application or deployment path. The one
966 place this work touches a dual-use-adjacent topic (Appendix A.6's exploratory probe of
967 gradient-inversion attack convergence) studies existing, already-published attacks on a toy
968 model and explicitly does not claim a generalizable privacy mitigation, so it introduces no
969 new attack capability.

970 Guidelines:

- 971 • The answer [N/A] means that there is no societal impact of the work performed.
- 972 • If the authors answer [N/A] or [No], they should explain why their work has no societal
973 impact or why the paper does not address societal impact.
- 974 • Examples of negative societal impacts include potential malicious or unintended uses
975 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
976 (e.g., deployment of technologies that could make decisions that unfairly impact specific
977 groups), privacy considerations, and security considerations.
- 978 • The conference expects that many papers will be foundational research and not tied
979 to particular applications, let alone deployments. However, if there is a direct path to
980 any negative applications, the authors should point it out. For example, it is legitimate
981 to point out that an improvement in the quality of generative models could be used to
982 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
983 that a generic algorithm for optimizing neural networks could enable people to train
984 models that generate Deepfakes faster.
- 985 • The authors should consider possible harms that could arise when the technology is
986 being used as intended and functioning correctly, harms that could arise when the
987 technology is being used as intended but gives incorrect results, and harms following
988 from (intentional or unintentional) misuse of the technology.
- 989 • If there are negative societal impacts, the authors could also discuss possible mitigation
990 strategies (e.g., gated release of models, providing defenses in addition to attacks,
991 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
992 feedback over time, improving the efficiency and accessibility of ML).

993 **11. Safeguards**

994 Question: Does the paper describe safeguards that have been put in place for responsible
995 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
996 image generators, or scraped datasets)?

997 Answer: [N/A]

998 Justification: The paper releases no pretrained models, generative systems, or new datasets; it
999 trains small classifiers on standard public benchmarks for diagnostic measurement purposes
1000 only, posing no misuse risk beyond that of the original benchmarks themselves.

1001 Guidelines:

- 1002 • The answer [N/A] means that the paper poses no such risks.
- 1003 • Released models that have a high risk for misuse or dual-use should be released with
1004 necessary safeguards to allow for controlled use of the model, for example by requiring
1005 that users adhere to usage guidelines or restrictions to access the model or implementing
1006 safety filters.
- 1007 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1008 should describe how they avoided releasing unsafe images.
- 1009 • We recognize that providing effective safeguards is challenging, and many papers do
1010 not require this, but we encourage authors to take this into account and make a best
1011 faith effort.

1012 **12. Licenses for existing assets**

1013 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1014 the paper, properly credited and are the license and terms of use explicitly mentioned and
1015 properly respected?

1016 Answer: [No]

1017 Justification: The paper names and uses CIFAR-10/100, Tiny-ImageNet-200, and Places365-
1018 Standard, all standard, long-established public research benchmarks, downloaded from their
1019 official public release locations and used only for research measurement (no redistribution),
1020 but does not separately cite each dataset’s originating publication or state its specific license
1021 terms in the text.

1022 Guidelines:

- 1023 • The answer [N/A] means that the paper does not use existing assets.
- 1024 • The authors should cite the original paper that produced the code package or dataset.
- 1025 • The authors should state which version of the asset is used and, if possible, include a
1026 URL.
- 1027 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1028 • For scraped data from a particular source (e.g., website), the copyright and terms of
1029 service of that source should be provided.
- 1030 • If assets are released, the license, copyright information, and terms of use in the
1031 package should be provided. For popular datasets, `paperswithcode.com/datasets`
1032 has curated licenses for some datasets. Their licensing guide can help determine the
1033 license of a dataset.
- 1034 • For existing datasets that are re-packaged, both the original license and the license of
1035 the derived asset (if it has changed) should be provided.
- 1036 • If this information is not available online, the authors are encouraged to reach out to
1037 the asset’s creators.

1038 **13. New assets**

1039 Question: Are new assets introduced in the paper well documented and is the documentation
1040 provided alongside the assets?

1041 Answer: [N/A]

1042 Justification: The paper introduces no new public dataset or pretrained model as a contribu-
1043 tion; the small architectures built for the generalization experiments (the MLP-Mixer and
1044 Transformer-Encoder variants) are fully specified in Appendix A.8 but are not released as
1045 standalone assets distinct from the analysis code.

1046 Guidelines:

1047 • The answer [N/A] means that the paper does not release new assets.

1048 • Researchers should communicate the details of the dataset/code/model as part of their

1049 submissions via structured templates. This includes details about training, license,

1050 limitations, etc.

1051 • The paper should discuss whether and how consent was obtained from people whose

1052 asset is used.

1053 • At submission time, remember to anonymize your assets (if applicable). You can either

1054 create an anonymized URL or include an anonymized zip file.

1055 **14. Crowdsourcing and research with human subjects**

1056 Question: For crowdsourcing experiments and research with human subjects, does the paper

1057 include the full text of instructions given to participants and screenshots, if applicable, as

1058 well as details about compensation (if any)?

1059 Answer: [N/A]

1060 Justification: The paper involves no crowdsourcing or human subjects research.

1061 Guidelines:

1062 • The answer [N/A] means that the paper does not involve crowdsourcing nor research

1063 with human subjects.

1064 • Including this information in the supplemental material is fine, but if the main contribu-

1065 tion of the paper involves human subjects, then as much detail as possible should be

1066 included in the main paper.

1067 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,

1068 or other labor should be paid at least the minimum wage in the country of the data

1069 collector.

1070 **15. Institutional review board (IRB) approvals or equivalent for research with human**

1071 **subjects**

1072 Question: Does the paper describe potential risks incurred by study participants, whether

1073 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

1074 approvals (or an equivalent approval/review based on the requirements of your country or

1075 institution) were obtained?

1076 Answer: [N/A]

1077 Justification: The paper involves no human subjects research.

1078 Guidelines:

1079 • The answer [N/A] means that the paper does not involve crowdsourcing nor research

1080 with human subjects.

1081 • Depending on the country in which research is conducted, IRB approval (or equivalent)

1082 may be required for any human subjects research. If you obtained IRB approval, you

1083 should clearly state this in the paper.

1084 • We recognize that the procedures for this may vary significantly between institutions

1085 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the

1086 guidelines for their institution.

1087 • For initial submissions, do not include any information that would break anonymity (if

1088 applicable), such as the institution conducting the review.

1089 **16. Declaration of LLM usage**

1090 Question: Does the paper describe the usage of LLMs if it is an important, original, or

1091 non-standard component of the core methods in this research? Note that if the LLM is used

1092 only for writing, editing, or formatting purposes and does *not* impact the core methodology,

1093 scientific rigor, or originality of the research, declaration is not required.

1094 Answer: [Yes]

1095 Justification: [AUTHORS: fill in to match your actual workflow before submission.] An

1096 LLM-based coding agent was used substantially beyond writing/editing/formatting in this

1097 project: under author direction and review, it implemented experiment code, ran and

1098 debugged training jobs, performed the statistical analyses, and drafted text. This goes
1099 beyond the declaration-exempt case in the question's note and should be disclosed per
1100 current LLM-usage policy, with the specific scope of that involvement and the authors'
1101 verification process described here.

1102 Guidelines:

- 1103 • The answer [N/A] means that the core method development in this research does not
1104 involve LLMs as any important, original, or non-standard components.
- 1105 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1106 be described.