
What Accuracy and Gradient Cosine Miss: Evaluating Feedback Alignment via Scale Stability, Reference Validity, and Depth Utility

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite the success of deep learning, training deep networks in biologically plausible and hardware-efficient ways remains an open challenge. Feedback alignment (FA) methods address this by replacing backpropagation’s symmetric backward weights with fixed random matrices, but their effectiveness depends critically on whether they can be accurately evaluated. The standard evaluation relies on two quantities: task accuracy and cosine similarity between the method’s credit signal and the backpropagation gradient. We show that this reporting pair is insufficient by identifying two independent failure modes, both silent under current reporting: (1) measurement degeneracy, where the BP reference gradient collapses to the numerical floor in terminal-LayerNorm residual architectures, rendering cosine uninterpretable; and (2) aggregation collapse, where the aggregate cosine masks layerwise heterogeneity that concentrates credit at one end of the network. To address these limitations, we propose a diagnostic evaluation protocol based on three checks—scale stability, reference validity, and depth utility—together with per-layer rather than aggregate cosine reporting. Across five architectures and 125 trained models, the standard reporting pair gives no signal of failure in any audited case, while our protocol identifies all failures with wide calibration margins. The two failure modes are causally independent: a per-block scale penalty alleviates Mode 1 (residual scale explosion driving reference collapse) without affecting Mode 2 (cosine ranking that contradicts every functional metric we measured).

21 1 Introduction

22 The backpropagation algorithm trains neural networks by propagating error gradients layer by layer from the output back to the input. Despite its effectiveness, BP relies on two mechanisms that are both biologically implausible and computationally constraining: weight symmetry between the forward and backward paths, and sequential layer-by-layer updates that prevent parallelization [4, 6, 9, 13, 19]. Feedback alignment (FA) methods address these constraints by replacing the backward weights with fixed random matrices [3, 5, 13, 15, 17, 20]. Direct feedback alignment (DFA) further removes the sequential dependence by projecting the output error directly to each hidden layer through independent random connections [16], and has been shown to train modern architectures including transformers and graph networks with performance approaching BP [1, 12]. These properties make the FA family a candidate for both biologically plausible learning and hardware-efficient training [8]. Yet as these methods scale to deeper residual architectures, a prior question becomes pressing: how should we evaluate whether they actually work?

34 The standard evidence for a feedback alignment method is two quantities: task accuracy, which measures whether the network learned, and cosine similarity between the method’s credit signal and the backpropagation gradient, which measures whether the learning signal points in roughly the right direction [3, 12, 13, 16, 17]. A method that reaches nontrivial accuracy and reports positive alignment

38 is typically interpreted as having trained the network with useful credit assignment. This pair has
 39 been the primary reporting convention across a decade of FA research, and the alignment angle in
 40 particular has been explicitly recommended as the basis for training best practices [11, 13, 17].

41 The reporting pair gives no signal when a
 42 method has not actually trained the network.
 43 On a representative residual setting (Figure 1),
 44 both FA and DFA report non-trivial accuracy
 45 and positive aggregate cosine—the kind of num-
 46 bers a reader would interpret as evidence of
 47 credit reaching the network. Yet both fall be-
 48 low an architecture-matched frozen-blocks base-
 49 line whose residual blocks were never trained
 50 at all: neither method’s trained network out-
 51 performs the same architecture with random
 52 blocks. The standard reporting pair shows no
 53 sign of this, for two independent reasons. First,
 54 the cosine measurement itself can be invalid:
 55 in terminal-LayerNorm residual architectures,
 56 residual-stream scale explosion compresses the
 57 BP reference gradient to the numerical floor
 58 through the LayerNorm Jacobian, so the re-
 59 ported cosine is computed against floating-point
 60 noise rather than a meaningful direction; and
 61 even when the reference is meaningful, an ag-
 62 gregate cosine can appear positive because the
 63 per-layer contributions are concentrated at one
 64 end of the network. Second, accuracy alone can-
 65 not distinguish whether deep blocks help or hurt
 66 the network’s prediction.

67 We demonstrate these failures through
 68 a controlled audit of three methods—
 69 backpropagation, feedback alignment [13], and direct feedback alignment [16]—on residual
 70 architectures. FA and DFA share the same local loss and differ only in how credit reaches the deep
 71 layers, providing a controlled comparison on which our protocol identifies both failures with wide
 72 calibration margins while the standard pair does not.

73 Contributions.

- 74 • We identify two independent failure modes of the standard FA evaluation pair (accuracy +
 75 aggregate cosine) on deep residual architectures: *measurement degeneracy*, where the BP
 76 reference gradient collapses to the numerical floor and makes cosine uninterpretable, and
 77 *low intrinsic credit-direction quality*, where deep-layer credit is essentially unaligned with
 78 BP even when the reference is meaningful.
- 79 • We provide a three-diagnostic evaluation protocol that detects both modes. On a repre-
 80 sentative setting where both FA and DFA report non-trivial accuracy and positive cosine
 81 alignment, the standard pair gives no signal that either method’s deep blocks are unused; the
 82 protocol identifies both as failures with wide calibration margins.
- 83 • We show that cosine alignment to the BP gradient, even when measured validly, does not
 84 predict depth utility: under matched conditions, three independent functional metrics rank
 85 the audited methods in an order that cosine contradicts.
- 86 • We validate the protocol across 5 architectures, 125 trained models, and 17 experimental
 87 settings, and release a reference implementation.

88 2 Motivating Observations

89 We fix a common training recipe across all primary-audit experiments: pre-LayerNorm [2, 21]
 90 ResMLPs [18] on CIFAR-10 [10], trained for 100 epochs with AdamW [14] (learning rate 10^{-3} ,
 91 weight decay 0.01, cosine schedule, batch size 128) across three seeds. Architecture depth L and

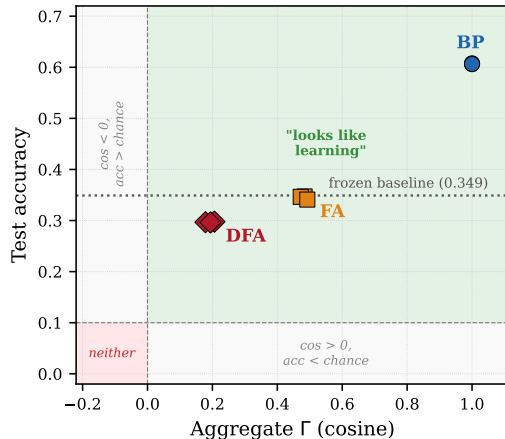


Figure 1: Standard reporting pair on the representative setting ($d=512$, $L=2$ ResMLP, three seeds): both FA and DFA report non-trivial accuracy and positive aggregate cosine.

Method	Test acc	Γ	vs Frozen	Verdict
BP	0.607	≈ 1.0	+25.7 pp	pass
FA	0.345	+0.48	-0.4 pp	fail
DFA	0.297	+0.19	-5.2 pp	fail

Table 1: Same data as Figure 1; trained blocks of FA and DFA fall below an architecture-matched frozen-random-blocks baseline.

width d are specified per experiment. Three methods are trained on this identical setup, differing only in how each layer’s credit signal a_l is computed.

Backpropagation (BP) computes the exact gradient via the chain rule: $a_l = \partial L / \partial h_l$.

Feedback Alignment (FA) [13] propagates credit sequentially through fixed random matrices: starting from the exact output-layer gradient $\partial L / \partial h_L$, each layer receives $a_l = B_l a_{l+1}$, where $B_l \in \mathbb{R}^{d \times d}$ is fixed at initialization.

Direct Feedback Alignment (DFA) [16] projects the output error directly to each layer, bypassing all intermediate layers: $a_l = B_l^\top e_T$, where $B_l \in \mathbb{R}^{C \times d}$ is a fixed random matrix and $e_T = \hat{y} - y$ is the output error.

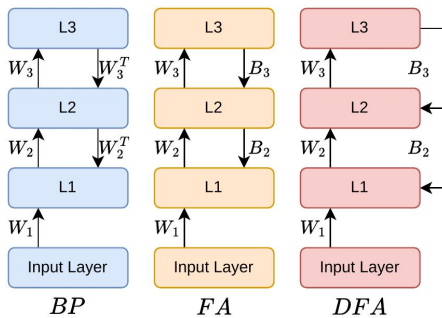


Figure 2: Backward architectures

FA and DFA share the same local loss: each block f_l is updated by reducing $-\langle f_l(h_l), a_l \rangle$. Neither loss contains a penalty on $\|f_l(h_l)\|$. The only difference between the two methods is how a_l is computed: FA preserves sequential structure in the credit path, DFA does not. A frozen-blocks baseline trains only the embedding, LayerNorm, and classification head while holding all residual blocks fixed at their random initialization; it reaches 0.349 ± 0.002 across the same three seeds.

The failure shown in Figure 1 has two independent sources internal to the standard reporting pair. To isolate them, we turn to the per-layer state of a related setting at $d = 256, L = 4$ (Figure 3). Two independent inconsistencies emerge from the per-layer data, each one severing the link between the standard reporting pair and the question it is meant to answer.

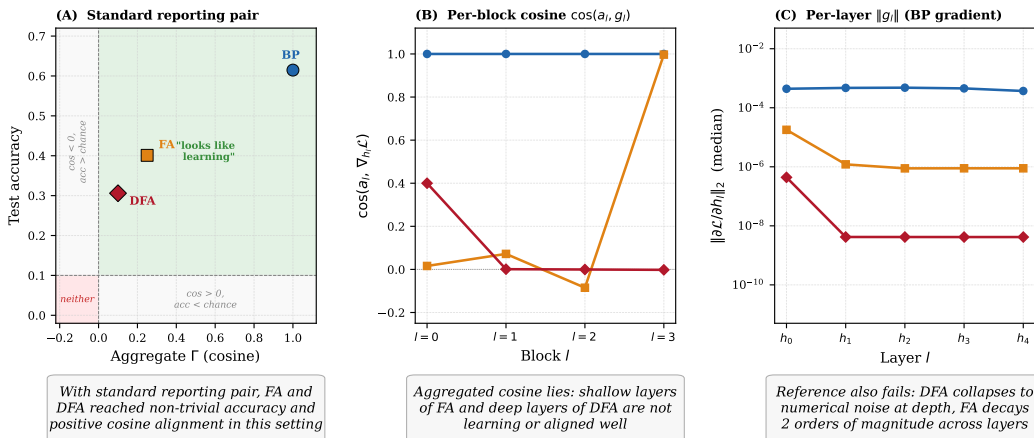


Figure 3: Per-layer state on the standard primary-audit setting (ResMLP $d = 256, L = 4$, CIFAR-10). (A) Standard reporting pair: all three methods report non-trivial accuracy and positive aggregate cosine. (B) Per-block cosine reveals that FA’s shallow blocks and DFA’s deep blocks contribute no aligned credit. (C) Per-layer BP gradient norm decays across FA’s depth and collapses to numerical noise in DFA, leaving cosine measurements at depth without a meaningful reference.

The aggregate cosine masks opposite per-layer patterns. On this setting, FA reports an aggregate cosine of $\Gamma \approx +0.23$ and DFA reports $\Gamma \approx +0.10$; both look like partial alignment to the BP gradient. Per-block cosine reveals that these aggregates are produced by inverse distributions (Figure 3B). FA’s shallowest block contributes near-zero cosine ($\approx +0.01$) while its deepest block contributes $+0.95$; the aggregate is carried almost entirely by the deep end of the network. DFA shows the opposite shape: its shallowest block contributes $+0.39$ while its deepest block contributes ≈ 0 ; the aggregate is carried almost entirely by the embedding. Two networks whose credit reaches opposite ends of the depth axis produce aggregate cosines that read as the same kind of evidence. Aggregation collapses this distinction, so the aggregate Γ cannot in principle answer whether credit reaches the deep layers.

The reference gradient can collapse below the cosine clamp. The BP gradient norm $\|g_l\|$ along depth tells the second story (Figure 3C). BP holds $\|g_l\|$ within an order of magnitude of $\sim 4 \times 10^{-4}$ across all layers, providing a reference of consistent scale. FA decays from $\sim 2 \times 10^{-5}$ at the embedding to $\sim 9 \times 10^{-7}$ at the deepest layer, but every layer remains above PyTorch’s default

127 cosine-similarity denominator floor of $\varepsilon = 10^{-8}$. DFA collapses by three orders of magnitude after
 128 the embedding—from $\sim 4 \times 10^{-7}$ at h_0 to $\sim 4 \times 10^{-9}$ at h_1 and below—and remains beneath ε for
 129 every deeper layer. When the denominator is supplied by the clamp rather than the reference vector,
 130 $\cos(a_l, g_l)$ is no longer comparing the credit direction to a meaningful BP direction; it is comparing
 131 it to a normalized floating-point residual.

132 **The two failures are independent.** On this setting, FA fails the first check (its aggregate is dominated
 133 by a single end of the network) but passes the second (every reference norm remains measurable);
 134 DFA fails both. The aggregate-dominance failure does not require reference collapse, and reference
 135 collapse does not require any particular layerwise distribution of cosine. Each one severs the standard
 136 reporting pair from its intended meaning on its own. Neither method universally avoids either
 137 mode—Section 4 shows that both modes can be triggered by either method in other settings, and that
 138 interventions on one mode do not resolve the other.

139 3 Failure Modes

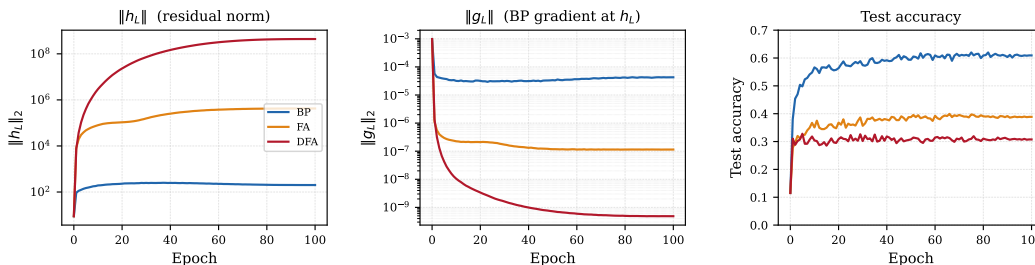


Figure 4: Temporal evolution of $\|h_L\|$, $\|g_L\|$, and test accuracy on ResMLP $d = 256$, $L = 4$, terminal LN. DFA’s residual norm grows four orders of magnitude and its BP reference gradient collapses to $\sim 10^{-9}$, the regime where cosine similarity becomes a comparison against floating-point noise. FA exhibits the same growth direction at attenuated magnitude, with $\|g_L\|$ stabilizing two orders of magnitude above this regime.

140 3.1 Mode 1: Residual scale explosion drives reference collapse

141 In Mode 1, the BP reference gradient at deep layers is compressed below the numerical floor used
 142 by the cosine implementation, so the cosine reported by the standard pair is computed against a
 143 normalized floating-point residual rather than against the BP direction. The mechanism is a chain of
 144 three steps connecting the local loss form to this measurement collapse.

145 The local loss $-\langle f_l(h_l), a_l \rangle$ is unbounded in $\|f_l(h_l)\|$: any direction in which the block output
 146 grows along a_l reduces the loss further. In a pre-LN residual block $h_{l+1} = h_l + f_l(h_l)$, growth
 147 in $\|f_l(h_l)\|$ accumulates directly into the residual stream, so $\|h_L\|$ increases along training. The
 148 terminal LayerNorm rescales h_L by its norm, giving a Jacobian whose spectrum scales as $1/\|h\|$ [21].
 149 The deepest BP reference gradient $\|g_L\| = \|\partial L / \partial h_L\|$ inherits this $1/\|h\|$ factor, so $\|g_L\|$ shrinks at
 150 the same rate that $\|h_L\|$ grows.

151 Figure 4 shows this chain in the DFA trajectory on the primary-audit setting: $\|h_L\|$ grows from $\sim 10^2$
 152 at initialization to $\sim 10^8$ by epoch 100, while $\|g_L\|$ collapses from $\sim 10^{-3}$ to $\sim 10^{-9}$ over the same
 153 window—four orders of magnitude in each direction, in tandem. By epoch 20, $\|g_L\|$ is already below
 154 PyTorch’s default cosine-similarity denominator floor of $\varepsilon = 10^{-8}$, and remains there for the rest of
 155 training.

156 FA shares DFA’s local loss form but exhibits a markedly attenuated chain on the same architecture
 157 (Figure 4): $\|h_L\|$ grows to $\sim 10^5$ rather than $\sim 10^8$, and $\|g_L\|$ stabilizes at $\sim 10^{-7}$, within the
 158 regime where cosine is computed against a meaningful reference. The local loss form is therefore
 159 not sufficient for Mode 1 to become a measurement failure; the credit propagation rule determines
 160 whether the chain reaches the cosine clamp. FA is not generally safe from Mode 1—the same
 161 mechanism produces full reference collapse on terminal-LN architectures with different geometry
 162 (Section 4).

163 **3.2 Mode 2: Aggregation collapses layerwise heterogeneity**

164 In Mode 2, the per-layer cosine values along depth follow distributions that concentrate at one end
 165 of the network, and the aggregate cosine averages over this concentration; the standard pair’s Γ no
 166 longer indicates whether credit reaches the deep layers, only that some layers are aligned somewhere
 167 along the depth. Section 2 already showed two such concentrations on the same architecture: FA’s
 168 per-block cosine rises from $\approx +0.01$ at the embedding to $+0.95$ at the deepest block, while DFA’s
 169 drops from $+0.39$ at the embedding to ≈ 0 at the deepest block. Both produce aggregate cosines that
 170 read as partial alignment.

171 The two concentrations follow from the credit propagation rule. DFA’s $a_l = B_l^\top e_T$ is a fixed random
 172 projection of the output error and contains no information about the forward state at layer l . The
 173 embedding sits one block away from the output and receives credit whose correlation with the BP
 174 gradient is preserved by the projection’s geometry, but every deeper block receives a signal generated
 175 independently of its forward state, and the cosine reflects that independence—near zero in expectation.
 176 FA’s $a_l = B_l a_{l+1}$ inherits the deepest block’s exact gradient $\partial L/\partial h_L$ and degrades upstream by an
 177 additional random matrix product per layer. The two rules produce alignment patterns that are inverse
 178 but symmetric in their effect on the aggregate.

179 Per-layer cosine reporting [13, 17] addresses this specific aggregation failure but does not recover the
 180 standard pair: it carries no information about whether each layer’s reference gradient is numerically
 181 valid (Mode 1), and it does not indicate whether the trained depth as a whole contributes to the
 182 network’s prediction (Section 5). Section 4.2 establishes that Mode 1 and Mode 2 are causally
 183 independent failure modes, both observationally and under a penalty intervention that alleviates Mode
 184 1 without affecting Mode 2.

185 **4 Validation**

186 **4.1 Scope of Mode 1 across architectures**

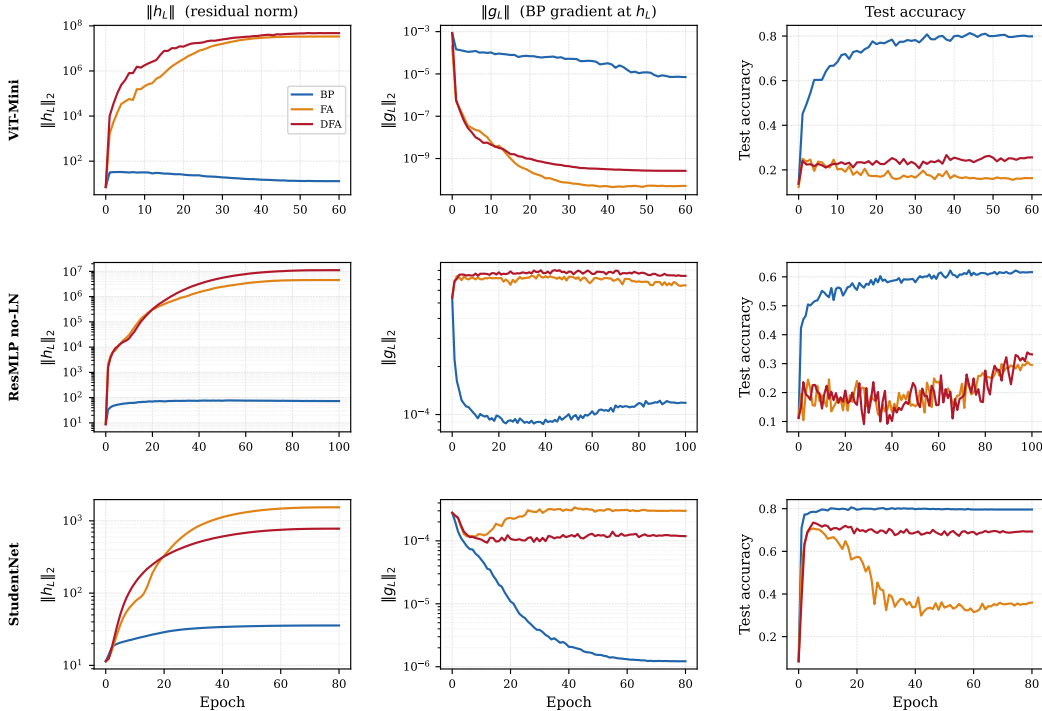


Figure 5: Temporal evolution of $\|h_L\|$, $\|g_L\|$, and test accuracy across three architectures: ViT-Mini (terminal LN), ResMLP $d = 256$, $L = 4$ with terminal LN removed, and StudentNet (no LN). The $\|g_L\|$ y-axis ranges differ across rows to expose within-row dynamics; the qualitative split is between collapsing trajectories on ViT-Mini and stable trajectories on the two no-LN architectures.

187 The two failure modes documented in Section 3 on ResMLP have different cross-architecture
 188 scopes. We examine three additional settings that vary terminal LayerNorm presence while holding

189 the other components of the architecture roughly fixed (Figure 5): a vision transformer ViT-Mini
 190 ($d = 128, L = 4$, cls token + terminal LN), the same ResMLP $d = 256, L = 4$ used in the primary
 191 audit but with the terminal LayerNorm removed, and StudentNet ($d = 128, L = 4$, no LN). Across
 192 these three settings, Mode 1a (residual scale growth under FA and DFA) appears in all three; Mode
 193 1b (reference gradient collapse) appears only in the architecture with terminal LayerNorm.

194 On ViT-Mini, both FA and DFA reproduce the full Mode 1 chain seen on ResMLP: $\|h_L\|$ grows to
 195 $\sim 10^7$ and $\|g_L\|$ collapses to $\sim 10^{-9}$ within 30 epochs, well below the regime where cosine similarity
 196 computes against a meaningful reference. The transformer architecture provides no protection from
 197 Mode 1 once terminal LN is in place. On ResMLP with terminal LN removed and on StudentNet (no
 198 LN), the residual stream still grows under FA and DFA— $\|h_L\|$ reaches $\sim 10^7$ on no-LN ResMLP
 199 and $\sim 10^3$ on StudentNet—but $\|g_L\|$ stabilizes at or above $\sim 10^{-4}$ on both, in the regime where
 200 cosine remains a valid measurement. The unbounded growth incentive of the local loss is architecture-
 201 agnostic, but its translation into measurement collapse requires terminal LayerNorm to compress the
 202 BP reference gradient.

203 Table 2 (✓=passes, ✗=fails, FA/DFA) shows
 204 Mode 1a failing universally—the unbounded
 205 growth incentive is architecture-agnostic—
 206 while Mode 1b tracks terminal LN presence.
 207 The same-backbone control isolates this: re-
 208 moving terminal LN from the primary-audit
 209 ResMLP flips the Mode 1b verdict for both
 210 methods without changing Mode 1a. The depth-
 211 utility consequence on ViT-Mini is severe—the frozen-blocks baseline reaches 0.570 ± 0.003 while
 212 FA and DFA reach 0.163 and 0.256, trained blocks underperforming random blocks by more than 30
 213 percentage points.

Architecture	LN	Scale stable	Reference valid
ResMLP	yes	✗/✗	✓/✗
ResMLP, no LN	no	✗/✗	✓/✓
ViT-Mini	yes	✗/✗	✗/✗
StudentNet	no	✗/✗	✓/✓

Table 2: Mode 1 verdict per architecture (FA/DFA).

214 4.2 Penalty intervention separates the two modes

215 The two failure modes can be intervened on independently. Adding a per-block scale penalty
 216 $\lambda \cdot \|f_l(h_l)\|^2$ to the local loss on each residual block—without modifying the credit rule, the optimizer,
 217 or any other component—suppresses Mode 1 in a dose-response manner on methods where Mode 1
 218 is severe; at the dose where Mode 1 is fully alleviated but Mode 2 is not yet engaged, the deep-layer
 219 cosine remains at the vanilla value, separating the two modes causally.

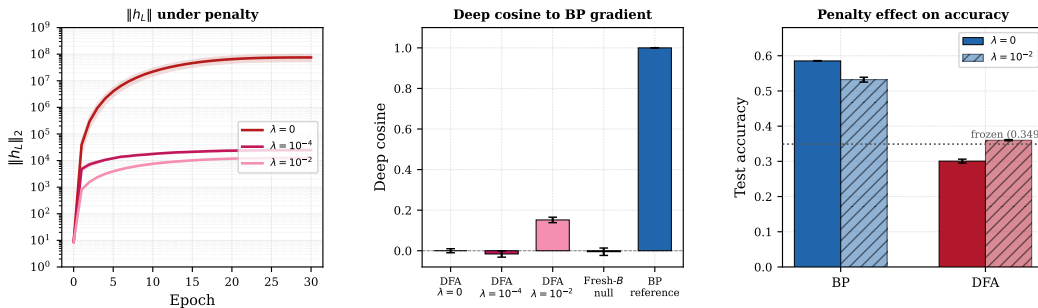


Figure 6: Penalty intervention separates Mode 1 from Mode 2 on DFA, ResMLP $d = 256, L = 4$, three seeds, 30 epochs. **Left:** residual norm under $\lambda \in \{0, 10^{-4}, 10^{-2}\}$ shows dose-response containment. **Middle:** deep cosine across conditions. At $\lambda = 10^{-4}$ the penalty contains $\|h_L\|$ and lifts $\|g_L\|$ above the cosine clamp, but deep cosine remains at the vanilla value; at $\lambda = 10^{-2}$ deep cosine recovers partially to $+0.152 \pm 0.013$. **Right:** BP and DFA accuracy with and without $\lambda = 10^{-2}$ penalty against the frozen-blocks baseline. BP loses 5.3 pp under penalty (capacity cost) but stays 18.3 pp above frozen; DFA gains 5.9 pp but only reaches 1.1 pp above frozen.

220 We sweep $\lambda \in \{0, 10^{-4}, 10^{-2}\}$ on DFA, ResMLP $d = 256, L = 4$, three seeds, 30 epochs (Figure 6,
 221 left and middle panels). Vanilla DFA reaches $\|h_L\| \sim 4 \times 10^8$ and $\|g_L\| \sim 5 \times 10^{-10}$, with deep
 222 cosine at the vanilla zero. At $\lambda = 10^{-4}$, the residual stream is contained ($\|h_L\| \sim 2.3 \times 10^4$) and
 223 $\|g_L\| \sim 5 \times 10^{-6}$ recovers well above the clamp—Mode 1 is alleviated—while deep cosine remains
 224 at -0.016 ± 0.016 , indistinguishable from zero. At $\lambda = 10^{-2}$, $\|h_L\|$ is further contained and deep
 225 cosine rises to $+0.152 \pm 0.013$, indicating partial recovery of Mode 2. The two modes respond at

226 different intervention strengths, and Mode 2 does not follow Mode 1’s recovery passively. FA on
 227 ResMLP exhibits only mild Mode 1 and produces no dose response (Appendix A).

228 Two control experiments rule out alternative explanations of the deep-cosine recovery: a BP control
 229 shows the penalty does not artificially lift cosine on a method whose credit is already exact (Ap-
 230 pendix C.1); a fresh- B null calibration shows the recovered cosine is specific to the matrices used
 231 during training rather than trivial adaptation (Appendix C.2).

232 Two observational checks corroborate the intervention. Across methods on ResMLP $d = 256, L = 4$,
 233 FA’s $\|g_l\|$ remains above the cosine clamp at every layer (Figure 4, middle panel), so its standard
 234 pair fails only through Mode 2; DFA on the same architecture fails both. Along training time, DFA’s
 235 deep-layer cosine is already -0.008 ± 0.013 at epoch 1 while $\|g_L\| \sim 10^{-7}$ remains above ε : Mode
 236 2 is present before Mode 1 develops. Together with the penalty intervention, this establishes Mode 1
 237 and Mode 2 as causally independent failure modes the protocol must diagnose separately. Neither
 238 method universally avoids either mode—FA on ViT-Mini [7] exhibits both despite passing Mode
 239 1b on ResMLP—so the modes are properties of (architecture, credit-rule) pairings rather than of
 240 methods in isolation.

241 4.3 Cosine cannot predict depth utility

242 Mode 1 and Mode 2 together explain why the standard reporting pair fails on settings where the
 243 measurement is invalid or the aggregate masks layerwise heterogeneity. A natural follow-up is
 244 whether cosine alignment, when measured validly and reported per-layer, suffices on its own—that is,
 245 whether the depth-utility check (the frozen-blocks comparison) is redundant once Mode 1 and Mode
 246 2 are addressed. We show it is not: even on settings where Mode 1 has been alleviated and the cosine
 247 measurement is valid, deep-layer cosine to the BP gradient does not predict whether depth is being
 248 used.

249 To establish this, we introduce two diagnostic probes designed to vary credit quality in the cosine
 250 space:

251 **State Bridge (SB)** learns a state predictor $G_\psi(h_l, t_l)$ that estimates the deepest hidden state h_L from
 252 the current layer’s state; the credit signal is the gradient of a cross-entropy loss computed by passing
 253 $G_\psi(h_l, t_l)$ through the classification head, $a_l = \nabla_{h_l} \text{CE}(\text{head}(G_\psi(h_l, t_l)), y)$.

254 **Credit Bridge (CB)** learns a value network $V_\phi(h_l, t_l)$ that directly estimates a scalar value at the
 255 current layer; the credit signal is its input gradient, $a_l = \nabla_{h_l} V_\phi(h_l, t_l)$.

256 SB and CB are not proposed as competitive FA methods; they exist to produce credit signals whose
 257 deep-layer cosine to the BP gradient differs systematically from DFA’s, so that we can ask whether
 258 these cosine differences correspond to functional differences in credit quality. Both probes train under
 259 the same per-block penalty $\lambda = 10^{-2}$ used in Section 4.2 to alleviate Mode 1, the predictor and value
 260 network themselves are unpenalized.

Table 3: Four methods under matched penalty rescue ($\lambda = 10^{-2}$, ResMLP $d = 256, L = 4$, three seeds, 30 epochs). Three functional metrics agree on SB \gg rest; deep cosine ranks the methods in a different order.

Method	Test acc	Deep cosine	Nudging ($\eta = 0.01$)	Train loss Δ
SB + pen	0.453 \pm 0.003	+0.322 \pm 0.008	-1.93×10^{-3}	-0.447
FA + pen	0.369 \pm 0.003	+0.423 \pm 0.006	-2.09×10^{-4}	-0.128
CB + pen	0.360 \pm 0.004	+0.679 \pm 0.010	-4.26×10^{-4}	-0.121
DFA + pen	0.360 \pm 0.002	+0.152 \pm 0.013	-4.98×10^{-5}	-0.095

261 Table 3 reports four methods under matched penalty rescue. The functional metrics agree on a single
 262 ranking: SB \gg FA \approx CB \approx DFA. SB reaches 0.453 accuracy, 9 percentage points above the others;
 263 its nudging effect (the change in test loss after a small step in the credit direction) is 4 to 40 times
 264 larger than the other three; its training loss decreases by -0.447 over 30 epochs, 3 to 5 times more
 265 than the others. Deep cosine, by contrast, ranks the methods CB $>$ FA $>$ SB $>$ DFA. CB has the
 266 highest deep cosine ($+0.679$), more than double SB’s ($+0.322$), yet matches DFA’s accuracy (0.360).
 267 FA has higher deep cosine than SB ($+0.423$ vs $+0.322$) but 9 percentage points lower accuracy and
 268 an order-of-magnitude smaller nudging effect. The cosine ranking does not track the functional
 269 ranking on any of the three functional metrics.

270 The dissociation rules out the interpretation that deep-layer cosine to the BP gradient measures
271 whether the credit signal is useful for training depth. What cosine measures is angular agreement:
272 the credit direction a_l matches the BP gradient direction g_l to varying degrees. CB is designed to
273 produce a credit signal close to g_l in angle, and it succeeds—its deep cosine is the highest of the four.
274 But angular agreement at a single point does not guarantee that an update in that direction reduces
275 loss across the relevant region of state space, nor that the resulting weight updates compose into a
276 working forward computation. SB’s credit signal has lower angular agreement with the BP gradient
277 but produces larger functional effects per step (nudging) and over training (loss decrease, accuracy).
278 The two are different quantities, and cosine measures only the first.

279 This has a direct consequence for the protocol. The depth-utility check in Section 5 cannot be replaced
280 by a per-layer cosine check, even one that conditions on Mode 1 being alleviated and Mode 2 being
281 addressed by per-layer reporting. Whether the trained depth contributes to the network’s prediction is
282 a question about functional behavior, not angular agreement, and it requires its own measurement:
283 the architecture-matched frozen-blocks baseline. The standard reporting pair fails at the cosine axis
284 through Modes 1 and 2; it also fails at the accuracy axis by giving no signal that the trained network
285 underperforms an architecture whose deep blocks were never trained at all. The protocol therefore
286 requires three separate diagnostics, not a refinement of the cosine alone.

287 5 Recommended evaluation protocol

288 The audit and the failure-mode analysis converge on three checks that together address how the
289 standard reporting pair fails. We state them as a protocol that an evaluator can run on any FA training
290 run; the design rationale follows the structure of Sections 3 and 4.

291 5.1 Three diagnostic checks

292 **Diagnostic 1 (Scale stability).** Compute the maximum per-block residual growth $\rho =$
293 $\max_l \|h_{l+1}\|/\|h_l\|$ at the end of training. Flag the run if $\rho > 50$. This detects the residual-stream
294 explosion that drives Mode 1a. The threshold is calibrated against the audit: BP and FA on no-LN
295 architectures stay below 10, while DFA on every audited architecture exceeds 90, giving a calibration
296 gap of nearly an order of magnitude.

297 **Diagnostic 2 (Reference validity).** Compute the deepest BP reference gradient norm $\|g_L\|$ at the
298 end of training. Flag the run if $\|g_L\| < 10 \cdot \varepsilon$, where ε is the cosine implementation’s denominator
299 floor for the training dtype (PyTorch’s default is $\varepsilon = 10^{-8}$ for fp32, giving a threshold of 10^{-7} in our
300 setup). At this threshold, the floor contributes more than $\sim 10\%$ of the cosine denominator, making
301 the reported angle a mixture of the true reference direction and the implementation’s clamp. In the
302 audit, BP and FA on ResMLP $d = 256, L = 4$ stay above 10^{-6} , while DFA collapses to $\sim 5 \times 10^{-10}$,
303 well below the threshold.

304 **Diagnostic 3 (Depth utility).** Compute the gap between the trained model’s test accuracy and
305 an architecture-matched frozen-blocks baseline (residual blocks frozen at random initialization,
306 embedding/LN/head trained). Flag the run if the gap is below 2 percentage points. This detects
307 whether the trained depth contributes to the network’s prediction, the question that Section 4.3 showed
308 cosine cannot answer even when measured validly.

309 **Verdict logic.** A run fails the protocol if either Mode 1 (Diagnostics 1 and 2 both flag) or Depth
310 Utility (Diagnostic 3 flags) is triggered. Diagnostic 1 alone does not constitute a Mode 1 failure—
311 scale growth without reference collapse leaves cosine measurable, as observed for FA on ResMLP.
312 Diagnostic 2 alone is unlikely without Diagnostic 1 by the mechanism in Section 3.1. Diagnostic
313 3 stands on its own: a run can pass both Mode 1 diagnostics and still fail at depth utility, which is
314 precisely what Section 4.3 establishes.

315 5.2 Decision-utility ablation

316 To verify that each diagnostic group contributes to detecting failures the previous strategies miss,
317 we ablate the protocol on two representative cases. Each case is a method-architecture pair where
318 the standard reporting pair gives a misleading verdict, and we ask which reporting strategy first
319 identifies the failure. Case 1 is DFA on ResMLP $d = 256, L = 4$ (the primary audit), where the
320 failure is driven by Mode 1 measurement degeneracy. Case 2 is FA on ResMLP $d = 512, L = 2$ (the
321 representative setting from Section 2), where Mode 1 is not triggered but the trained depth fails to

322 outperform the frozen baseline. Together they exercise both the Mode 1 check and the depth utility
 323 check independently.

Table 4: Decision-utility ablation. Each cell shows the reporting strategy’s verdict on the case; **X** marks a missed failure. Both cases are genuine failures (the trained network underperforms or matches a frozen-blocks baseline).

Reporting strategy	Case 1: DFA, $d = 256, L = 4$	Case 2: FA, $d = 512, L = 2$
S0: acc only	X pass (acc 0.306)	X pass (acc 0.345)
S1: acc + aggregate Γ	X pass ($\Gamma = +0.10$)	X pass ($\Gamma = +0.48$)
S2: + Mode 1 check (D1 & D2)	✓ fail	X pass (D1, D2 both ok)
S3: + Depth utility (D3)	✓ fail	✓ fail

324 Table 4 reports the verdicts. The standard reporting pair (S0, S1) misses both failures: DFA reports
 325 above-chance accuracy and a positive aggregate cosine, FA does the same on the shallower setting;
 326 nothing in the standard pair indicates that either method’s trained residual blocks underperform a
 327 frozen baseline. Adding the Mode 1 check (S2) catches Case 1 but not Case 2: DFA’s Mode 1 is fully
 328 triggered ($\|h_L\| \sim 5 \times 10^8, \|g_L\| \sim 5 \times 10^{-10}$, both diagnostics flag), while FA’s Mode 1 on the
 329 shallower setting is not triggered ($\|h_L\|$ stays moderate, $\|g_L\|$ stays above the clamp), so S2 reports
 330 FA as passing. The depth utility check (S3) catches Case 2: FA’s accuracy of 0.345 is below the
 331 frozen baseline of 0.349, flagging the run regardless of Mode 1 status. The full protocol catches both
 332 failures.

333 Each of the three checks catches at least one case the others miss. Appendix B extends this ablation
 334 to additional architectures and methods, with consistent results.

335 5.3 Scope, recommendations, and discussion

336 **Scope.** The thresholds are calibrated on supervised image classification with small to medium residual
 337 architectures (ResMLP, ViT-Mini, StudentNet) on CIFAR-10 in fp32. On other architecture families,
 338 datasets, or precisions (fp16/bf16 raise the effective denominator floor), the threshold values may
 339 need recalibration; the measurements themselves remain well-defined wherever a forward pass, a
 340 BP reference gradient, and a frozen-blocks baseline can be computed. The Mode 1b mechanism
 341 is similarly scoped to terminal-LayerNorm architectures; on architectures without it, Diagnostic 2
 342 functions as a defensive check rather than a primary failure detector. Extending the audit to RL,
 343 generative training, or larger-scale NLP is open work.

344 **Cross-architecture footprint.** Diagnostic 1 applies anywhere the local loss form $-\langle f_l(h_l), a_l \rangle$ is
 345 used without a norm penalty. Diagnostic 2 is informative only on architectures where the BP reference
 346 can mechanistically collapse; on no-LN architectures it adds little signal but costs nothing. Diagnostic
 347 3 is architecture-agnostic.

348 **Recommendations.** (i) Report raw diagnostic values, not only pass/fail, to allow recalibration.
 349 (ii) Report cosine alignment per-layer rather than aggregate. (iii) Treat the protocol as a minimum
 350 bar: passing rules out the failures we audited but does not validate the method as a BP alternative.

351 6 Conclusion

352 The standard reporting pair for feedback alignment—task accuracy and aggregate cosine alignment
 353 to the BP gradient—can fail to indicate that a method has trained the network. We identified two
 354 independent ways the cosine half can fail: the BP reference gradient can collapse below the cosine
 355 implementation’s denominator floor in terminal-LayerNorm residual architectures (Mode 1), and the
 356 aggregate can mask layerwise heterogeneity that concentrates credit at a single end of the network
 357 (Mode 2). We further showed that even when cosine is measured validly, it does not predict whether
 358 trained depth contributes to the network’s prediction; this is a separate failure of the accuracy half,
 359 requiring an architecture-matched frozen-blocks baseline to detect. The two failures are causally
 360 independent, both observationally and under a penalty intervention that alleviates Mode 1 without
 361 affecting Mode 2.

362 The recommended protocol consists of three diagnostic checks—scale stability, reference validity,
 363 and depth utility—together with per-layer rather than aggregate cosine reporting. Across the audited
 364 architectures and methods, the standard reporting pair gives no signal of failure where our protocol
 365 identifies failure. The protocol does not certify a method as universally effective; it rules out the

366 specific class of silent failures the audit revealed, and leaves the standard pair’s interpretation intact
367 when all three diagnostics pass.

368 **References**

- 369 [1] Mohamed Akrouf, Collin Wilson, Peter C. Humphreys, Timothy P. Lillicrap, and Douglas B. Tweed. Deep
370 learning without weight transport. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
371 arXiv:1904.05391.
- 372 [2] Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- 373 [3] Sergey Bartunov, Adam Santoro, Blake A. Richards, Geoffrey E. Hinton, and Timothy P. Lillicrap.
374 Assessing the scalability of biologically-motivated deep learning algorithms and architectures. In *Advances*
375 *in Neural Information Processing Systems (NeurIPS)*, 2018. arXiv:1807.04587.
- 376 [4] Eugene Belilovsky, Michael Eickenberg, and Edouard Oyallon. Greedy layerwise learning can scale to
377 imagenet. In *International Conference on Machine Learning (ICML)*, 2018. arXiv:1812.11446.
- 378 [5] Brian Crafton, Abhinav Parihar, Evan Gebhardt, and Arijit Raychowdhury. Direct feedback alignment with
379 sparse connections for local learning. *Frontiers in Neuroscience*, 2019. doi: 10.3389/fnins.2019.00525.
380 arXiv:1903.02083.
- 381 [6] Giorgia DellaFerrera and Gabriel Kreiman. Error-driven input modulation: Solving the credit assignment
382 problem without a backward pass. In *International Conference on Machine Learning (ICML)*, 2022.
- 383 [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
384 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
385 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In
386 *International Conference on Learning Representations (ICLR)*, 2021. arXiv:2010.11929.
- 387 [8] Charlotte Frenkel, Martin Lefebvre, and David Bol. Learning without feedback: Fixed random learning
388 signals allow for feedforward training of deep neural networks. *Frontiers in Neuroscience*, 2021. doi:
389 10.3389/fnins.2021.629892.
- 390 [9] Geoffrey E. Hinton. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint*
391 *arXiv:2212.13345*, 2022. doi: 10.48550/arXiv.2212.13345.
- 392 [10] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of
393 Toronto, 2009.
- 394 [11] Julien Launay, Iacopo Poli, and Florent Krzakala. Principled training of neural networks with direct
395 feedback alignment, 2019.
- 396 [12] Julien Launay, Iacopo Poli, François Boniface, and Florent Krzakala. Direct feedback alignment scales
397 to modern deep learning tasks and architectures. In *Advances in Neural Information Processing Systems*
398 *(NeurIPS)*, 2020.
- 399 [13] Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Random synaptic
400 feedback weights support error backpropagation for deep learning. *Nature Communications*, 7(1):13276,
401 2016. doi: 10.1038/ncomms13276.
- 402 [14] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on*
403 *Learning Representations (ICLR)*, 2019.
- 404 [15] Theodore H. Moskovitz, Ashok Litwin-Kumar, and L. F. Abbott. Feedback alignment in deep convolutional
405 networks, 2018.
- 406 [16] Arild Nøkland. Direct feedback alignment provides learning in deep neural networks. In *Advances in*
407 *Neural Information Processing Systems (NeurIPS)*, 2016.
- 408 [17] Maria Refinetti, Stéphane d’Ascoli, Ruben Ohana, and Sebastian Goldt. Align, then memorise: The
409 dynamics of learning with feedback alignment. In *International Conference on Machine Learning (ICML)*,
410 2021.
- 411 [18] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby, Edouard Grave,
412 Gautier Izacard, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, and Hervé Jégou. Resmlp: Feedforward
413 networks for image classification with data-efficient training. *IEEE Transactions on Pattern Analysis and*
414 *Machine Intelligence*, 2021. doi: 10.1109/TPAMI.2022.3206148. arXiv:2105.03404.

- 415 [19] James C. R. Whittington and Rafal Bogacz. An approximation of the error backpropagation algorithm
 416 in a predictive coding network with local hebbian synaptic plasticity. *Neural Computation*, 2017. doi:
 417 10.1162/NECO_a_00949.
- 418 [20] Will Xiao, Honglin Chen, Qianli Liao, and Tomaso Poggio. Biologically-plausible learning algorithms
 419 can scale to large datasets. In *International Conference on Learning Representations (ICLR)*, 2018.
 420 arXiv:1811.03567.
- 421 [21] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan
 422 Lan, Liwei Wang, and Tie-Yan Liu. On layer normalization in the transformer architecture. In *International
 423 Conference on Machine Learning (ICML)*, 2020. arXiv:2002.04745.

424 A FA + penalty on ResMLP: no dose response

425 For completeness, we apply the same penalty intervention to FA on ResMLP $d = 256$, $L = 4$, three
 426 seeds, 30 epochs:

	λ	Test acc	Deep cos	$\ h_L\ $	$\ g_L\ $
427	0	0.372 ± 0.007	$+0.325 \pm 0.015$	$\sim 9 \times 10^4$	$\sim 2 \times 10^{-6}$
	10^{-4}	0.377 ± 0.006	$+0.298 \pm 0.031$	$\sim 9 \times 10^3$	$\sim 1 \times 10^{-5}$
	10^{-2}	0.369 ± 0.003	$+0.423 \pm 0.006$	$\sim 1 \times 10^4$	$\sim 2 \times 10^{-5}$

428 FA on this architecture has mild Mode 1 even at $\lambda = 0$: the residual norm $\|h_L\| \sim 9 \times 10^4$ is well
 429 below DFA’s $\sim 4 \times 10^8$, and the BP reference gradient $\|g_L\| \sim 2 \times 10^{-6}$ is already two orders of
 430 magnitude above the cosine clamp. The penalty has no Mode 1 to alleviate; accordingly, accuracy
 431 and deep cosine remain in narrow windows (0.37–0.38 and +0.30–+0.42 respectively) across all
 432 three λ values. This is consistent with the cross-architecture finding (Section 4.1) that FA’s Mode 1
 433 severity depends on the architecture, and the intervention is most informative where Mode 1 is fully
 434 developed.

435 B Extended decision-utility ablation

436 B.1 Protocol verdict on CIFAR-10

437 We extend the decision-utility ablation to additional architecture-method pairs to verify that the
 438 Section 5.2 pattern (each diagnostic group catching at least one case the others miss) is not specific to
 439 the two cases shown in the main text.

Table 5: Extended decision-utility ablation. Each cell shows the reporting strategy’s verdict; \times marks a missed failure (silently passing a method whose trained blocks underperform a frozen-blocks baseline). Healthy reference cases (BP) appear with \checkmark throughout. Cases include the two main-text cases (Section 5.2) plus three additional method-architecture pairs.

Reporting strategy	DFA, RM-d256-L4	FA, RM-d512-L2	DFA, ViT-Mini	FA, ViT-Mini	BP, RM-d256-L4
S0: acc only	\times	\times	\times	\times	\checkmark
S1: acc + Γ	\times	\times	\times	\times	\checkmark
S2: + Mode 1 (D1 & D2)	\checkmark	\times	\checkmark	\checkmark	\checkmark
S3: + Depth utility (D3)	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark

440 The pattern from the main text holds: the standard pair (S0, S1) silently passes all four genuine
 441 failure cases; the Mode 1 check (S2) catches the three cases where measurement is degenerate but
 442 does not catch the case where measurement is valid (FA on $d = 512$, $L = 2$); the depth utility
 443 check (S3) closes the remaining gap. The healthy case (BP on the primary audit) passes all reporting
 444 strategies, confirming that the protocol does not flag working methods. ViT-Mini provides the
 445 strongest depth-utility evidence in the extended set: the frozen-blocks baseline reaches 0.570 ± 0.003
 446 on this architecture, while DFA and FA reach 0.256 and 0.163 respectively—trained transformer
 447 blocks underperform random blocks by more than 30 percentage points.

448 StudentNet, the no-LN architecture used in Section 4.1 as a Mode 1b control, is excluded from
 449 this ablation: its frozen-blocks baseline reaches 0.908 ± 0.009 on the synthetic task, exceeding
 450 even fully trained BP (0.796). The task is solvable by a linear probe on random features, which
 451 makes Diagnostic 3 uninformative regardless of the training method. We use StudentNet only as a
 452 mechanism control for Mode 1b’s architectural scope and not as a depth-utility test case.

453 **B.2 Protocol verdict on CIFAR-100**

454 To check whether the protocol’s verdict pattern depends on the audited dataset, we replicate the
 455 primary-audit setup on CIFAR-100: ResMLP $d = 256, L = 4$, three seeds, 100 epochs, with the
 456 same optimizer, schedule, and training implementation (Section 2). We also compute an architecture-
 457 matched frozen-blocks baseline on CIFAR-100, which reaches 0.178 ± 0.001 —close to the linear-
 458 probe ceiling on CIFAR-100 pixels (0.177).

Table 6: Protocol verdict on CIFAR-100 (ResMLP $d = 256, L = 4$, three seeds, 100 epochs). The pattern from the primary CIFAR-10 audit is reproduced: BP passes; DFA fails Mode 1 and Depth Utility; FA passes Mode 1 (growth $15.6\times$, $\|g_L\| = 1.3 \times 10^{-6}$ above clamp) but fails Depth Utility (-4.5 pp vs frozen).

Method	Test acc	D1 (Growth)	D2 ($\ g_L\ $)	D3 (vs Frozen)	Verdict
BP	0.321 ± 0.002	$1.0\times \checkmark$	$9.5 \times 10^{-4} \checkmark$	$+14.3$ pp \checkmark	pass
FA	0.133 ± 0.013	$15.6\times \checkmark$	$1.3 \times 10^{-6} \checkmark$	-4.5 pp \times	fail (D3)
DFA	0.088 ± 0.001	$1004\times \times$	$9.1 \times 10^{-9} \times$	-9.0 pp \times	fail (D1+D2+D3)

459 The CIFAR-100 verdict pattern matches CIFAR-10 exactly. BP passes all three diagnostics; DFA
 460 triggers Mode 1 (both D1 and D2) and depth utility, with the same residual-scale-explosion-to-
 461 reference-collapse mechanism developed in Section 3; FA’s Mode 1 stays mild (growth below $50\times$,
 462 $\|g_L\|$ above the clamp) but its trained depth fails to outperform random blocks. FA on CIFAR-100 is
 463 a second instance of the depth-utility-only failure pattern (FA on $d = 512, L = 2$ ResMLP being the
 464 first in Section 5.2), strengthening the case that Diagnostic 3 catches a class of failures distinct from
 465 those caught by Mode 1 alone. The protocol verdict is dataset-invariant on this architecture.

466 **C Controls for the penalty intervention**

467 **C.1 BP control: penalty does not artificially lift cosine**

468 A potential alternative explanation for the deep-cosine recovery on DFA at $\lambda = 10^{-2}$ (Section 4.2) is
 469 that the penalty itself shifts cosine measurement statistics independent of any change in credit quality.
 470 We rule this out with a BP control: BP under the same penalty schedule reaches 0.585 ± 0.001
 471 accuracy at $\lambda = 0$ and 0.532 ± 0.007 at $\lambda = 10^{-2}$, a capacity cost of 5.3 percentage points (Figure 6,
 472 right panel). BP’s deep cosine remains ≈ 1.0 throughout: the penalty does not artificially lift cosine
 473 on a method whose credit is already exact. The penalty is not a free intervention—it costs accuracy on
 474 every method we tested—but BP+penalty still exceeds the frozen-blocks baseline by 18.3 percentage
 475 points, so the penalty does not turn working methods into broken ones.

476 **C.2 Fresh- B null calibration**

477 A second alternative is that the recovered deep cosine reflects trivial adaptation to the fixed feedback
 478 matrices B_l rather than an improvement in credit quality. We rule this out with a fresh- B null
 479 calibration: holding the trained DFA network at the $\lambda = 10^{-2}$ checkpoint (seed 42) fixed and
 480 replacing each B_l with a freshly drawn random matrix gives deep cosine -0.005 ± 0.018 over 20
 481 draws, against the trained- B value of $+0.166$ on the same checkpoint—a 9.3σ separation. The
 482 recovered cosine is well above the noise floor and is specific to the matrices used during training.

483 **D Related work**

484 **NeurIPS Paper Checklist**

485 The checklist is designed to encourage best practices for responsible machine learning research,
486 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
487 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
488 follow the references and follow the (optional) supplemental material. The checklist does NOT count
489 towards the page limit.

490 Please read the checklist guidelines carefully for information on how to answer these questions. For
491 each question in the checklist:

- 492 • You should answer [Yes], [No], or [N/A].
- 493 • [N/A] means either that the question is Not Applicable for that particular paper or the
494 relevant information is Not Available.
- 495 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

496 **The checklist answers are an integral part of your paper submission.** They are visible to the
497 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it
498 (after eventual revisions) with the final version of your paper, and its final version will be published
499 with the paper.

500 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
501 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a
502 proper justification is given (e.g., error bars are not reported because it would be too computationally
503 expensive” or “we were unable to find the license for the dataset we used”). In general, answering
504 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we
505 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
506 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
507 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
508 please point to the section(s) where related material for the question can be found.

509 **IMPORTANT, please:**

- 510 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 511 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 512 • **Do not modify the questions and only use the provided macros for your answers.**

513 **1. Claims**

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

516 Answer: **[TODO]**

517 Justification: **[TODO]**

518 Guidelines:

- 519 • The answer [N/A] means that the abstract and introduction do not include the claims
520 made in the paper.
- 521 • The abstract and/or introduction should clearly state the claims made, including the
522 contributions made in the paper and important assumptions and limitations. A [No] or
523 [N/A] answer to this question will not be perceived well by the reviewers.
- 524 • The claims made should match theoretical and experimental results, and reflect how
525 much the results can be expected to generalize to other settings.
- 526 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
527 are not attained by the paper.

528 **2. Limitations**

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

530 Answer: **[TODO]**

531 Justification: **[TODO]**

532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582

Guidelines:

- The answer [N/A] means that the paper has no limitation while the answer [No] means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate “Limitations” section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren’t acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

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

Answer: [TODO]

Justification: [TODO]

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [TODO]

Justification: [TODO]

Guidelines:

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

- 583 • If the paper includes experiments, a [No] answer to this question will not be perceived
584 well by the reviewers: Making the paper reproducible is important, regardless of
585 whether the code and data are provided or not.
- 586 • If the contribution is a dataset and/or model, the authors should describe the steps taken
587 to make their results reproducible or verifiable.
- 588 • Depending on the contribution, reproducibility can be accomplished in various ways.
589 For example, if the contribution is a novel architecture, describing the architecture fully
590 might suffice, or if the contribution is a specific model and empirical evaluation, it may
591 be necessary to either make it possible for others to replicate the model with the same
592 dataset, or provide access to the model. In general, releasing code and data is often
593 one good way to accomplish this, but reproducibility can also be provided via detailed
594 instructions for how to replicate the results, access to a hosted model (e.g., in the case
595 of a large language model), releasing of a model checkpoint, or other means that are
596 appropriate to the research performed.
- 597 • While NeurIPS does not require releasing code, the conference does require all submis-
598 sions to provide some reasonable avenue for reproducibility, which may depend on the
599 nature of the contribution. For example
 - 600 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
601 to reproduce that algorithm.
 - 602 (b) If the contribution is primarily a new model architecture, the paper should describe
603 the architecture clearly and fully.
 - 604 (c) If the contribution is a new model (e.g., a large language model), then there should
605 either be a way to access this model for reproducing the results or a way to reproduce
606 the model (e.g., with an open-source dataset or instructions for how to construct
607 the dataset).
 - 608 (d) We recognize that reproducibility may be tricky in some cases, in which case
609 authors are welcome to describe the particular way they provide for reproducibility.
610 In the case of closed-source models, it may be that access to the model is limited in
611 some way (e.g., to registered users), but it should be possible for other researchers
612 to have some path to reproducing or verifying the results.

613 5. Open access to data and code

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

617 Answer: [TODO]

618 Justification: [TODO]

619 Guidelines:

- 620 • The answer [N/A] means that paper does not include experiments requiring code.
- 621 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
622 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 623 • While we encourage the release of code and data, we understand that this might not
624 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
625 including code, unless this is central to the contribution (e.g., for a new open-source
626 benchmark).
- 627 • The instructions should contain the exact command and environment needed to run to
628 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 629 • The authors should provide instructions on data access and preparation, including how
630 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 631 • The authors should provide scripts to reproduce all experimental results for the new
632 proposed method and baselines. If only a subset of experiments are reproducible, they
633 should state which ones are omitted from the script and why.
- 634 • At submission time, to preserve anonymity, the authors should release anonymized
635 versions (if applicable).
636

- 637 • Providing as much information as possible in supplemental material (appended to the
638 paper) is recommended, but including URLs to data and code is permitted.

639 **6. Experimental setting/details**

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

642 Answer: **[TODO]**

643 Justification: **[TODO]**

644 Guidelines:

- 645 • The answer [N/A] means that the paper does not include experiments.
- 646 • The experimental setting should be presented in the core of the paper to a level of detail
647 that is necessary to appreciate the results and make sense of them.
- 648 • The full details can be provided either with the code, in appendix, or as supplemental
649 material.

650 **7. Experiment statistical significance**

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

653 Answer: **[TODO]**

654 Justification: **[TODO]**

655 Guidelines:

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

676 **8. Experiments compute resources**

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

680 Answer: **[TODO]**

681 Justification: **[TODO]**

682 Guidelines:

- 683 • The answer [N/A] means that the paper does not include experiments.
- 684 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
685 or cloud provider, including relevant memory and storage.
- 686 • The paper should provide the amount of compute required for each of the individual
687 experimental runs as well as estimate the total compute.

- 688 • The paper should disclose whether the full research project required more compute
689 than the experiments reported in the paper (e.g., preliminary or failed experiments that
690 didn't make it into the paper).

691 9. Code of ethics

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

694 Answer: **[TODO]**

695 Justification: **[TODO]**

696 Guidelines:

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

703 10. Broader impacts

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

706 Answer: **[TODO]**

707 Justification: **[TODO]**

708 Guidelines:

- 709 • The answer [N/A] means that there is no societal impact of the work performed.
- 710 • If the authors answer [N/A] or [No], they should explain why their work has no societal
711 impact or why the paper does not address societal impact.
- 712 • Examples of negative societal impacts include potential malicious or unintended uses
713 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
714 (e.g., deployment of technologies that could make decisions that unfairly impact specific
715 groups), privacy considerations, and security considerations.
- 716 • The conference expects that many papers will be foundational research and not tied
717 to particular applications, let alone deployments. However, if there is a direct path to
718 any negative applications, the authors should point it out. For example, it is legitimate
719 to point out that an improvement in the quality of generative models could be used to
720 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
721 that a generic algorithm for optimizing neural networks could enable people to train
722 models that generate Deepfakes faster.
- 723 • The authors should consider possible harms that could arise when the technology is
724 being used as intended and functioning correctly, harms that could arise when the
725 technology is being used as intended but gives incorrect results, and harms following
726 from (intentional or unintentional) misuse of the technology.
- 727 • If there are negative societal impacts, the authors could also discuss possible mitigation
728 strategies (e.g., gated release of models, providing defenses in addition to attacks,
729 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
730 feedback over time, improving the efficiency and accessibility of ML).

731 11. Safeguards

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

735 Answer: **[TODO]**

736 Justification: **[TODO]**

737 Guidelines:

- 738 • The answer [N/A] means that the paper poses no such risks.

- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

748 12. Licenses for existing assets

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

752 Answer: **[TODO]**

753 Justification: **[TODO]**

754 Guidelines:

- 755
- 756
- 757
- 758
- 759
- 760
- 761
- 762
- 763
- 764
- 765
- 766
- 767
- 768
- 769
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

770 13. New assets

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

773 Answer: **[TODO]**

774 Justification: **[TODO]**

775 Guidelines:

- 776
- 777
- 778
- 779
- 780
- 781
- 782
- 783
- The answer [N/A] means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

784 14. Crowdsourcing and research with human subjects

785 Question: For crowdsourcing experiments and research with human subjects, does the paper
786 include the full text of instructions given to participants and screenshots, if applicable, as
787 well as details about compensation (if any)?

788 Answer: **[TODO]**

789 Justification: **[TODO]**

790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

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