

---

# Beyond Accuracy and Alignment: A Diagnostic Evaluation Protocol for Feedback Alignment

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Modern feedback-alignment evaluation on deep residual networks is still summa-  
2 rized by a deceptively simple pair: headline accuracy and headline cosine align-  
3 ment  $\Gamma$  to the backpropagation gradient. We show that this pair can silently fail in  
4 two distinct ways on standard CIFAR-10 pre-LayerNorm ResMLP and ViT-Mini  
5 settings: first, *measurement degeneracy*, where residual-stream growth drives  
6 hidden-layer BP gradients to the numerical floor and makes  $\Gamma$  uninterpretable;  
7 and second, *low intrinsic credit-direction quality*, where random-feedback credit  
8 remains essentially unaligned with BP on the deep blocks even when the reference  
9 gradient is still meaningful. The headline result is that the field-standard reporting  
10 pair walks back none of the methods we audit, whereas a four-diagnostic protocol  
11 walks back the three degenerate methods and passes the two trustworthy controls.  
12 Our contribution is an evaluation methodology paper for the NeurIPS 2026 Evaluations  
13 & Datasets track: we provide the protocol, the calibration logic for its thresh-  
14 olds, a reference implementation, a five-method audit, and validation through tem-  
15 poral replay, cross-architecture checks, intervention-based disambiguation, and a  
16 documented catalog of pipeline pitfalls, in the spirit of critical evaluation analyses  
17 such as Jordan et al. [3], O’Bray et al. [2], Paleka et al. [1].

## 18 1 Introduction

19 Feedback-alignment papers are usually judged by two numbers: task accuracy and an aggregate  
20 similarity between the method’s local credit signal and the backpropagation gradient [4–7]. On  
21 the audited 4-block  $d=256$  ResMLP, however, Table 1 already shows that this pair is not a validity  
22 check: DFA reaches only  $0.306 \pm 0.006$  test accuracy, below the architecture-matched frozen-blocks  
23 baseline of  $0.349 \pm 0.002$ , while still looking superficially comparable to other non-BP methods.  
24 Figure 1 further shows that the apparent cosine evidence is concentrated at the shallowest block,  
25 with DFA at seed 42 reaching about  $+0.42$  at layer 0 but approximately  $-0.03$  to 0 on layers 1–4, so  
26 the aggregate obscures where credit direction is and is not present. At the same time, the deepest BP  
27 reference norm is only about  $5 \times 10^{-10}$  for DFA, State Bridge, and Credit Bridge, below the  $10^{-8}$   
28 clamp used by `F.cosine_similarity`, whereas BP remains around  $4 \times 10^{-4}$ , so the reported deep  
29 cosine is partly computed against a numerical-floor reference rather than an informative gradient  
30 direction (Figure 1; Table 1). Those numbers can be useful, but only if the measurement regime  
31 itself is valid.

32 Our audit shows that modern residual vision models can make these two quantities look informa-  
33 tive while failing to answer the question they are taken to answer. Figure 1 shows the first failure  
34 mode, which we call *Mode 1: measurement degeneracy*, where residual-stream growth drives the  
35 deepest hidden state to about  $\|h_L\| \sim 10^8$  under DFA/SB/CB while the corresponding BP reference

Table 1: Main audit table for the 4-block  $d=256$  pre-LayerNorm ResMLP on CIFAR-10. The row and column structure is fixed here; fill from the three-seed audit output.

Method	Test acc.	Headline $\Gamma$	Status-quo verdict	Protocol verdict
BP	$0.615 \pm 0.003$	$\approx 1.0$	trustworthy	trustworthy
EP	$0.316 \pm 0.030$	0.008	trustworthy	trustworthy
DFA	$0.306 \pm 0.006$	0.10	trustworthy	walked back
State Bridge	$0.205 \pm 0.032$	0.005	trustworthy	walked back
Credit Bridge	$0.289 \pm 0.026$	0.07	trustworthy	walked back

36 collapses to  $\|g_L\| \sim 5 \times 10^{-10}$ , so the deep-layer cosine is measured against a clamp-dominated  
 37 floor rather than a meaningful target direction. The same figure also shows the second failure mode,  
 38 *Mode 2: low intrinsic credit-direction quality*, because even after comparing against the stronger  
 39 frozen-blocks baseline ( $0.349 \pm 0.002$ ) and looking layer-by-layer, DFA’s deep blocks remain essen-  
 40 tially null while only layer 0 is visibly positive. To test whether this is only a measurement problem,  
 41 the intervention results show a dissociation: with a residual penalty  $\lambda \|f_l(h_l)\|^2$ , the deepest state  
 42 scale falls toward  $4 \times 10^4$ , the reference gradient rises toward  $10^{-6}$ , and deep cosine can improve  
 43 to about  $+0.16$ , yet at  $\lambda=10^{-4}$  Mode 1 is alleviated while deep cosine still stays near zero, and at  
 44 vanilla DFA epoch 1 the reference is already meaningful at about  $6 \times 10^{-7}$  but the deep cosine is still  
 45  $-0.008 \pm 0.013$  across three seeds. The failure is not unitary: one mode breaks the measurement,  
 46 and the other survives even when the measurement is still meaningful.

47 Accordingly, this paper does not introduce a new FA variant or a new benchmark. Instead, Table 1  
 48 and Figure 1 use a standard five-method CIFAR-10 audit to show that status-quo reporting would  
 49 treat BP, EP, DFA, State Bridge, and Credit Bridge as the same kind of evidence-bearing object  
 50 even though only BP and EP remain trustworthy under matched diagnostic checks. This makes the  
 51 contribution methodological in the sense of Jordan et al. [3], O’Bray et al. [2], and Paleka et al. [1]:  
 52 the central question is not whether one more FA variant can post a headline number, but whether the  
 53 reporting pipeline distinguishes meaningful credit-direction evidence from numerical-floor artifacts  
 54 and from shallow-only learning. The protocol therefore starts from per-layer diagnostics and a  
 55 frozen-blocks baseline before reading any aggregate cosine or final accuracy as evidence about deep  
 56 credit assignment. We first show the walk-back on a standard audit, then isolate the two failure  
 57 modes, and finally state the reporting protocol that future FA papers should satisfy.

## 58 2 Audit: Standard Reporting Walks Back Nothing

59 We begin with the smallest setting in which all methods can be compared head-to-head under iden-  
 60 tical architecture, optimizer family, and data. Table 1 fixes that canonical audit to a 4-block pre-  
 61 LayerNorm ResMLP with width  $d=256$  on CIFAR-10, trained for 100 epochs with AdamW (learn-  
 62 ing rate  $10^{-3}$ , weight decay 0.01), a cosine schedule, and three seeds (42, 123, 456). Within that  
 63 single setting, BP, EP, DFA, State Bridge, and Credit Bridge can be read against the same architec-  
 64 ture and the same training budget, while Figure 1 summarizes the corresponding per-block growth,  
 65 deepest-layer BP reference norm, cross-batch stability, and frozen-baseline comparison. This is the  
 66 table a reader would normally use to decide whether the methods trained the deep network.

67 By the field’s usual criteria, the non-BP methods appear to train to nontrivial accuracy and report  
 68 nonzero alignment. In Table 1, DFA reaches  $0.306 \pm 0.006$  test accuracy with headline  $\Gamma=0.10$ ,  
 69 State Bridge reaches  $0.205 \pm 0.032$  with  $\Gamma=0.005$ , and Credit Bridge reaches  $0.289 \pm 0.026$  with  
 70  $\Gamma=0.07$ ; none of these rows looks like an obvious invalidation if one is reading the usual pair of final  
 71 accuracy and aggregate alignment in the style of prior FA reporting [4–7]. Even the absolute scale  
 72 does not itself force a walk-back, because all three methods are plainly above chance and all three  
 73 report positive headline alignment rather than a visibly broken or undefined quantity. That reading  
 74 is exactly what the rest of the paper overturns.

75 Low accuracy by itself is not the pathology. EP is the key internal comparison in Table 1 and  
 76 Figure 1: it achieves only  $0.316 \pm 0.030$  accuracy and a very small headline  $\Gamma=0.008$ , yet its per-  
 77 block growth is only  $11.6\times$ , its deepest BP reference norm remains around  $1.3 \times 10^{-4}$  rather than  
 78 collapsing to the numerical floor, and its cross-batch direction-stability score is 0.02 rather than the  
 79 much higher drift-dominated values seen for DFA-family methods. At the same time, EP is not a

5-method audit on 4-block  $d=256$  ResMLP CIFAR-10 (3-seed mean  $\pm$  std)

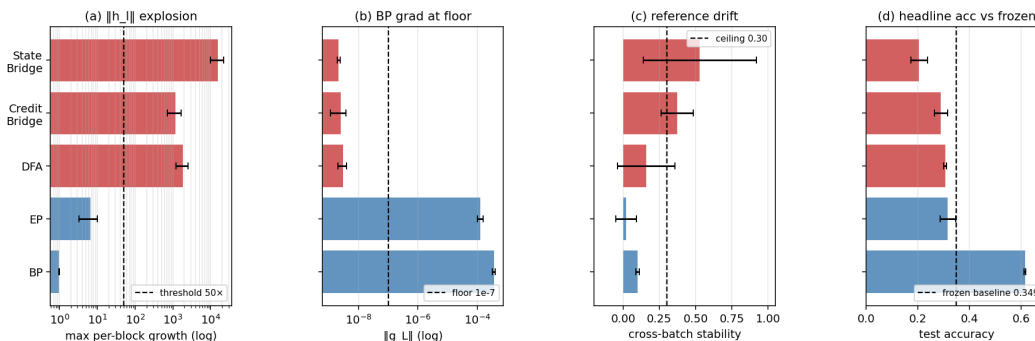


Figure 1: Five-method audit on the 4-block  $d=256$  pre-LayerNorm ResMLP: the field-standard pair looks superficially consistent across methods, but the diagnostic view separates trustworthy controls from walked-back methods.

80 positive result for depth usage in the stronger sense, because its trainable-model accuracy is still  
 81 3.3 percentage points below the frozen-blocks baseline of  $0.349 \pm 0.002$ . The distinction matters  
 82 because it separates underperformance from invalid evaluation.

83 When we compare each method to a frozen-blocks baseline matched to the same architecture, the  
 84 headline interpretation changes immediately. The frozen-blocks model, which trains only the em-  
 85 bedding, LayerNorm, and head while holding the residual blocks fixed, reaches  $0.349 \pm 0.002$  across  
 86 the same three seeds; against that baseline, BP is higher by 26.6 points, but DFA is lower by 4.3  
 87 points, State Bridge by 14.4 points, Credit Bridge by 6.0 points, and even EP by 3.3 points. Fig-  
 88 ure 1 shows that this accuracy comparison lines up with the diagnostic split: DFA, State Bridge, and  
 89 Credit Bridge also combine extreme per-block growth ( $237\times$ ,  $12000\times$ , and  $96\times$ ), deepest-layer BP  
 90 norms around  $10^{-9}$ , and high cross-batch instability (0.16, 0.53, and 0.37), so their deep blocks are  
 91 at best passengers and in practice often harmful. This establishes the audit question the rest of the  
 92 paper must answer: why do the standard signals fail so badly?

### 93 3 Failure Mode 1: Measurement Degeneracy

94 The first failure mode is a scale pathology, not yet an alignment pathology. On the audited 4-block  
 95 pre-LayerNorm ResMLP ( $d=256$ , CIFAR-10, 100 epochs, 3 seeds), DFA optimizes block-local ob-  
 96 jectives of the form  $\langle f_i(h_i), e_T B_i^\top \rangle$  with no explicit scale constraint on  $f_i$ , so for any direction in  
 97 which increasing  $\|f_i(h_i)\|$  improves alignment with the fixed feedback target  $B_i^\top e_T$ , the local ob-  
 98 jective rewards larger output magnitude. In a pre-LN residual stack, larger block outputs directly  
 99 increase residual-stream scale; terminal LayerNorm then removes task-loss sensitivity to that scale  
 100 at the output, so the architecture provides no global restraint on the local growth incentive [7]. In  
 101 the same runs, each block’s  $w_1$  and  $w_2$  grows by roughly  $200\times$  in relative delta, their norm product  
 102 reaches about  $5 \times 10^4$  per block, and the terminal hidden-state norm  $\|h_L\|$  rises monotonically from  
 103 about 9 at random initialization to about  $4 \times 10^8$  by epoch 100 (Figure 2). Most of that growth  
 104 appears immediately:  $\|h_L\|$  already reaches about  $10^6$  by epoch 5. As a direct test of whether this  
 105 growth needs task signal at all, we re-ran DFA, State Bridge, and Credit Bridge on the same back-  
 106 bone with i.i.d. random class targets refreshed every minibatch, so the labels carry no information;  
 107 under random targets all three methods stay at chance accuracy, yet  $\|h_L\|$  still grows from about 9 to  
 108 about  $1.45 \times 10^4$  for DFA,  $6.2 \times 10^3$  for State Bridge, and  $2.0 \times 10^4$  for Credit Bridge within three  
 109 epochs, and DFA’s  $\|g_L\|$  already drops to about  $5.6 \times 10^{-7}$ , so Mode 1 is essentially data-agnostic  
 110 on this architecture across the three audited fixed-feedback local-credit methods (Appendix I). Once  
 111 the residual stream reaches this regime, the backpropagation reference vector no longer behaves like  
 112 a healthy target.

113 The measurement failure occurs at the point where the hidden-layer BP gradient ceases to be a mean-  
 114 ingful reference direction. In terminal-LayerNorm architectures, the LayerNorm Jacobian scales  
 115 as  $\partial \text{LN}(h) / \partial h \propto 1 / \|h\|$  in expectation, so the same residual-stream inflation is accompanied by  
 116 collapse of the hidden-layer BP reference norm: on DFA-trained ResMLP,  $\|g_L\|$  falls from about

117  $9.8 \times 10^{-4}$  at random initialization to about  $5 \times 10^{-10}$  by epoch 100, a six-order-of-magnitude drop,  
 118 while the reported cosine remains mathematically defined only because `F.cosine_similarity`  
 119 clamps the denominator at  $\varepsilon=10^{-8}$  (Table 1; Figure 1). At that endpoint the reference norm is about  
 120  $20\times$  below the clamp, so the quantity being reported is effectively  $(a \cdot b)/(\|a\| \max(\|b\|, 10^{-8}))$   
 121 rather than a comparison to an informative BP direction. At that point, reporting a cosine is no  
 122 longer evidence about credit quality.

123 The simplest control is architectural, not theoretical. On the same ResMLP backbone, BP keeps  
 124  $\|h_L\|$  near 200 and  $\|g_L\|$  near  $4 \times 10^{-4}$  throughout training, while EP keeps  $\|h_L\|$  around  $5 \times 10^3$   
 125 and  $\|g_L\|$  around  $1.3 \times 10^{-4}$ , so hard optimization on CIFAR-10 by itself does not force hidden-  
 126 layer gradients to the numerical floor (Table 1; Figure 2). The broader cross-architecture pattern is  
 127 consistent with the same interpretation: StudentNet and the BatchNorm CNN, which lack terminal  
 128 LayerNorm, keep deepest BP gradients around  $10^{-4}$  and never trigger diagnostic (b), whereas ViT-  
 129 Mini with a terminal LN shows the same collapse pattern and triggers diagnostic (b) by epochs  
 130 2–3 (Figure 2). To check whether the additive residual skip itself is the proximate trigger, we ran  
 131 a matched ResMLP-d256 ablation that replaces  $h_{l+1} = h_l + F_l(h_l)$  with  $h_{l+1} = F_l(h_l)$  while  
 132 keeping terminal LN and all other hyperparameters fixed; in that ablation DFA’s  $\|h_L\|$  still grows  
 133 from  $\sim 5$  to  $\sim 2.2 \times 10^4$  within three epochs and  $\|g_L\|$  already drops to  $\sim 1.6 \times 10^{-7}$ , so the additive  
 134 skip is *not* necessary for Mode 1 either, even though the no-residual stack is partially degenerate for  
 135 both BP and DFA (Appendix H). The pathology therefore belongs to the evaluated FA regime, not  
 136 to CIFAR-10, the backbone, or the residual skip alone.

137 The collapse is not a late-epoch curiosity. For vanilla DFA on the ResMLP temporal replay,  $\|g_L\|$   
 138 drops from  $9.8 \times 10^{-4}$  at epoch 0 to  $1.4 \times 10^{-6}$  at epoch 1,  $3.1 \times 10^{-7}$  at epoch 2,  $1.3 \times 10^{-7}$  at  
 139 epoch 3, and  $6.7 \times 10^{-8}$  at epoch 4, so diagnostic (b) fires at epoch 3–4 across all three seeds, while  
 140 the max-per-block growth detector fires slightly later at epochs 8–11 (Figure 2). Both detectors  
 141 therefore fire in the first 11 epochs of a 100-epoch run, making the protocol actionable as an early-  
 142 stop criterion rather than a post hoc explanation. The practical point is reinforced by accuracy: DFA  
 143 is at 0.308 already at epoch 4 and ends at 0.306 by epoch 100, so the remaining training budget  
 144 adds essentially nothing to the headline result once the measurement has already degenerated. Once  
 145 measurement degeneracy is identified, the next question is whether poor deep credit remains even  
 146 before collapse.

#### 147 4 Failure Mode 2: Low Intrinsic Credit-Direction Quality

148 The second failure mode appears even in the meaningful-measurement regime. At the earliest vanilla  
 149 DFA checkpoints on ResMLP, the hidden backpropagated gradient at the first deep block remains  
 150 above the numerical floor: at epoch 1,  $\|g_2\|$  is  $6.7 \times 10^{-7}$ ,  $6.5 \times 10^{-7}$ , and  $3.9 \times 10^{-7}$  across the three  
 151 seeds, all above the  $10^{-7}$  threshold used to distinguish measurable from collapsed gradients. Yet the  
 152 corresponding deep-layer cosine values are already essentially null: across layers 1–4, all seed-level  
 153 measurements at epoch 1 lie in  $[-0.04, +0.02]$ , with a three-seed mean of  $-0.008 \pm 0.013$ , and  
 154 by epoch 2 the deep mean is still only  $-0.018 \pm 0.018$  (Table 2). This is the observational pattern  
 155 predicted by low credit-direction quality rather than mere disappearance of signal: the gradient is  
 156 still present enough to measure, but the directions delivered to the deep network carry little agree-  
 157 ment with backpropagation, consistent with prior concerns that alternative feedback rules can fail by  
 158 supplying poor credit assignments even before full collapse [8, 9, 11? ]. This rules out the simplest  
 159 objection that the deep-layer null result is merely a byproduct of collapse.

160 A second metric with different numerical failure modes tells the same story. Cosine measures di-  
 161 rectional agreement with the BP gradient, whereas perturbation correlation  $\rho$  measures whether the  
 162 proposed update predicts the correct sign and relative magnitude of loss change under actual per-  
 163 turbations; their failure modes are therefore different, especially with respect to normalization and  
 164 small-denominator effects. In our controls,  $\rho$  behaves as expected, with a Taylor-ceiling positive  
 165 control near  $+0.997$  and a random-vector negative control near  $+0.006$  (Figure 3, Table 2). On  
 166 vanilla DFA, deep  $\rho$  is likewise null: for the early checkpoints where the gradients remain measur-  
 167 able, the deep average is  $-0.003 \pm 0.005$  across seeds and epochs, and in a floor-level checkpoint it is  
 168  $+0.002$ , again indistinguishable from noise. The agreement between cosine and  $\rho$  therefore rules out  
 169 the interpretation that the null deep result is an artifact of cosine’s  $\varepsilon$ -clamp or vector normalization.  
 170 The deep blocks are not just hard to measure; they are receiving weakly useful directions.

Table 2: Two-mode validation table built around the intervention and disambiguation results.

Condition	Deep-layer alignment signal	Measurement regime	Interpretation
Vanilla DFA, early epoch	$\overline{\cos}_{deep} = -0.008 \pm 0.013, \overline{\rho}_{deep} = -0.003 \pm 0.005$	meaningful ( $\ g\  \sim 10^{-6}$ )	mode 2 present without m
Vanilla DFA, converged	$\overline{\cos}_{deep} = -0.022, \overline{\rho}_{deep} = +0.002$	degenerate ( $\ g\  \sim 10^{-9}$ )	mode 1 obscures mod
Penalized DFA, $\lambda=10^{-2}$	$\overline{\cos}_{deep} = +0.155 \pm 0.025, \overline{\rho}_{deep} = +0.080 \pm 0.011$	meaningful ( $\ g\  \sim 10^{-6}$ )	partial alleviation of both
Fresh- $B$ null control	$\overline{\cos}_{deep} = +0.002 \pm 0.022$ ( $n=20$ draws)	meaningful	training-specific adaptation

171 Per-layer reporting is therefore not cosmetic. In ResMLP under vanilla DFA, the headline aggregate  
 172 alignment  $\Gamma \approx 0.07$ – $0.10$  can look mildly positive only because layer 0 remains strongly aligned  
 173 while the deep network is not: at the same early checkpoints where layers 1–4 are essentially zero,  
 174 layer 0 has cosine  $+0.42$ ,  $+0.45$ , and  $+0.39$  across seeds (Table 2). The resulting average can there-  
 175 fore be driven by the embedding layer even when the interior blocks are effectively unaligned, so  
 176 aggregate reporting obscures the very distinction needed to separate “measurement collapse” from  
 177 “poor credit direction.” This layer-0 dominance is specific to the ResMLP DFA setting; on ViT-Mini  
 178 DFA, all layers are near zero, which strengthens the broader methodological point that alignment  
 179 should be reported per layer rather than only in aggregate. With the two modes separated observa-  
 180 tionally, the remaining question is whether intervention can move them independently.

## 181 5 Intervention and Cross-Architecture Evidence

182 The penalty intervention first matters as a rescue of the measurement regime. When we add a per-  
 183 block penalty  $\lambda \text{mean}(\|f_i(h_i)\|^2)$  to DFA’s local loss and train the 4-block  $d=256$  ResMLP for 30  
 184 epochs on CIFAR-10, the  $\lambda=10^{-2}$  setting contains the terminal hidden-state scale from  $\|h_L\| \sim$   
 185  $4.4 \times 10^8$  under vanilla DFA to  $\sim 4.0 \times 10^4$ , while lifting the deepest BP reference norm from  
 186  $\|g_L\| \sim 5 \times 10^{-10}$  to  $\sim 9.0 \times 10^{-7}$ , a roughly four-order-of-magnitude rescue on both quantities  
 187 (Figure 3; Table 2). At that setting, both diagnostic (a) and diagnostic (b) pass on penalized DFA,  
 188 and test accuracy rises to  $0.363 \pm 0.001$  from  $0.308 \pm 0.014$  for vanilla DFA. The key point is not  
 189 yet that the recovered network has good deep credit, but that the deep reference vector is again large  
 190 enough to function as a meaningful target direction rather than a clamp-dominated artifact. That  
 191 rescue makes the second question measurable rather than hypothetical.

192 Once the reference vector is meaningful again, the deep layers no longer sit exactly at null. At  
 193  $\lambda=10^{-2}$ , penalized DFA reaches a three-seed deep-layer mean cosine of  $+0.155 \pm 0.025$  and deep  
 194 perturbation correlation of  $+0.080 \pm 0.011$ , whereas vanilla DFA is essentially zero on both metrics  
 195 in the deep blocks, consistent with prior concerns that alternative feedback can fail by supplying  
 196 poor credit directions even before full collapse [8, 9, 11? ]. The null calibration rules out the inter-  
 197 pretation that this recovered signal is merely measurement noise: on the same penalized checkpoint,  
 198 replacing the training-time feedback matrices with 20 fresh random  $B_i$  draws gives a deep cosine  
 199 of only  $+0.002 \pm 0.022$ , with per-layer standard deviations of  $0.013$ – $0.023$ , all within noise of zero  
 200 (Table 2). The  $\lambda$  sweep sharpens the dissociation further: at  $\lambda=10^{-4}$ , Mode 1 is already alleviated,  
 201 with  $\|h_L\|=2.4 \times 10^4$  and  $\|g_L\|=6.3 \times 10^{-7}$ , but deep cosine remains  $-0.022$ , while at  $\lambda=10^{-2}$  it  
 202 rises to  $+0.165$  and deep  $\rho$  to  $+0.091$  (Figure 3). The improvement is real, but it is only partial.

203 A rescue intervention is only informative if its direct cost is controlled. The relevant control is BP  
 204 trained under the same penalty: BP falls from  $0.609 \pm 0.004$  without the penalty to  $0.530$  with  
 205  $\lambda=10^{-2}$ , so the penalty has a direct cost of about 8 percentage points even when credit assignment  
 206 is correct, whereas DFA moves in the opposite direction, from  $0.308 \pm 0.014$  to  $0.363 \pm 0.001$   
 207 under the same intervention (Figure 3). Relative to the frozen-blocks baseline of  $0.349$ , BP+penalty  
 208 still retains a margin of  $+18.1$  points, while DFA+penalty retains only  $+1.4$  points. The remaining  
 209 gap,  $0.530 - 0.363 = 17$  points, is therefore a lower bound on the part of DFA’s deficit that is not  
 210 explained by simple penalty-induced capacity loss alone, though not a clean isolation because BP  
 211 uses an end-to-end loss whereas DFA uses block-local losses. The residual gap after that control is  
 212 what keeps Mode 2 substantively alive.

213 The architecture comparison sharpens the scope of the critique. In the terminal-LN architectures  
 214 we audited, both diagnostics fire for DFA-trained ResMLP at  $d=256$ , the same pattern recurs at  
 215  $d=512$  with even larger max-per-block growth (about  $1.5 \times 10^4$ ), and ViT-Mini with a class token  
 216 and terminal LN shows diagnostic (a) by epoch 1 and diagnostic (b) by epochs 2–3 (Figure 2).

Cross-architecture temporal evolution of FA diagnostics (seed 42)

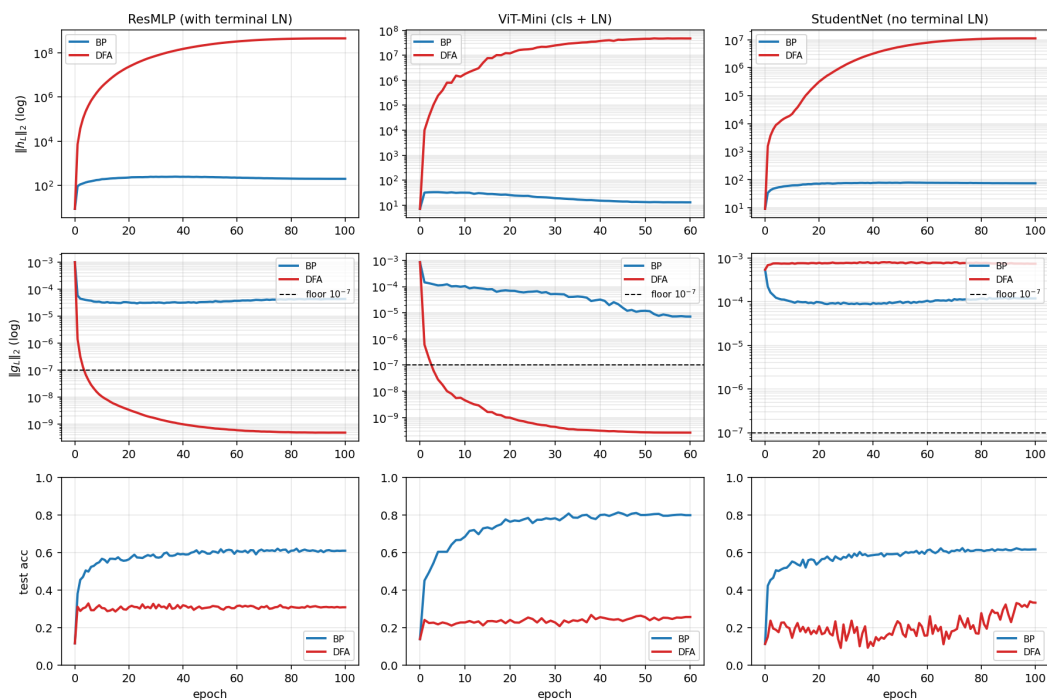


Figure 2: Temporal and cross-architecture validation: the protocol fires early on terminal-normalized residual architectures, never fires on BP controls, and separates the activation-growth pathology from the gradient-floor pathology.

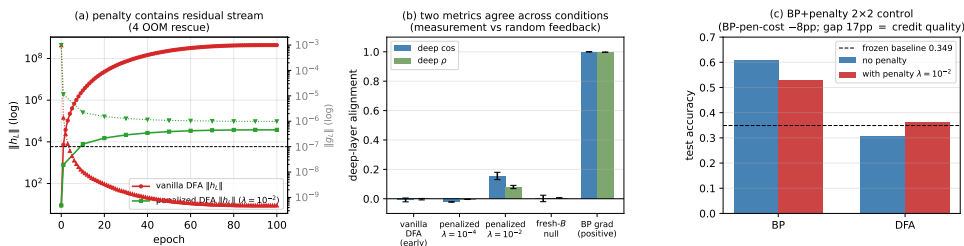


Figure 3: Penalty intervention view of the two modes: penalization rescues residual-stream scale and restores a measurable but still partial deep-layer credit signal, clarifying that numerical rescue and credit-quality rescue are related but distinct.

217 A depth sweep on the  $d=512$  ResMLP at  $L \in \{2, 4, 6, 8, 12\}$  shows that the layerwise pattern  
 218 is essentially depth-invariant: DFA's layer-0 cosine stays in  $[-0.39, +0.40]$  across all five depths,  
 219 while its mean deep-layer cosine stays within  $[-0.005, +0.000]$  and its deep perturbation correlation  
 220 collapses to 0.000 in every depth tested, even though BP retains a deep-layer cosine of  $+0.94$  at  
 221  $L=12$  (Appendix G). The deep credit signal does not improve when the network is shallower, so  
 222 the failure is not a "too deep" artifact. In the non-terminal-LN controls, the pattern is different:  
 223 StudentNet shows diagnostic (a) only at epochs 14–25 while diagnostic (b) never fires across 100  
 224 epochs and three seeds, and the BatchNorm CNN on CIFAR-10 likewise shows strong growth under  
 225 DFA, with max-per-block growth up to  $237\times$ , but keeps deepest BP gradients around  $\|g\| \sim 10^{-3}$   
 226 and never triggers diagnostic (b) (Figure 2). BP never triggers either diagnostic in any audited  
 227 architecture. This is an observational association rather than a causal identification of terminal  
 228 LayerNorm as the unique mechanism, but it is enough to support a narrower claim: diagnostic (b)  
 229 appears tied to the terminal-LN architectures audited here, while diagnostic (a) remains useful more  
 230 broadly. This lets the paper end with a reporting rule rather than an overclaimed theory.

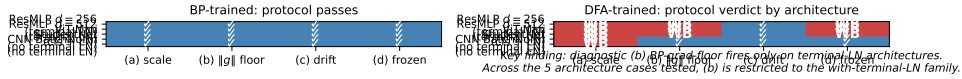


Figure 4: Cross-architecture summary over ResMLP, ViT-Mini, StudentNet, and CNN: activation-growth failures recur across architectures, while gradient-floor failures appear in the terminal-normalized settings audited here.

Table 3: Protocol definition table. Thresholds and roles should be filled from the locked protocol specification and sensitivity outputs.

Diag.	Measurement	Default threshold	Role
(a)	Per-layer activation scale via max-per-block growth $\max_l \ h_{l+1}\ /\ h_l\ $	$> 50\times$	binary detector
(b)	Deepest hidden-layer BP gradient norm $\ g_L\ $	$< 10^{-7}$	binary detector
(c)	Cross-batch direction stability of normalized BP gradients	$> 0.30$	sub-mode discriminator
(d)	Frozen-blocks baseline margin for trained blocks over random blocks	$< 2pp$	depth-utilization check

## 231 6 Recommended FA Evaluation Protocol

232 The reporting protocol begins with measurement validity. Before any FA paper reports a headline  
 233 alignment number, it should report per-layer state scale and the hidden BP reference-gradient scale  
 234 at the layers where the scientific claim is being made. In our audited regime, those two quantities  
 235 already separate healthy from invalid measurement with unusually wide margins: the maximum  
 236 per-block growth stays below about  $11\times$  for BP and EP but is at least  $694\times$  for the degenerate  
 237 methods, giving a  $63\times$  calibration gap, while the deepest hidden BP norm stays above about  $10^{-4}$   
 238 for BP and EP but below about  $4 \times 10^{-9}$  for the degenerate methods, giving a  $24,338\times$  gap (Table 3;  
 239 Table 1; Figure 4). These are not cosmetic diagnostics around the real result: they determine whether  
 240 the reported cosine is being computed against an informative BP direction or against a floor-level  
 241 reference. If the reference gradient is at floor, the evaluator should stop treating aggregate alignment  
 242 as evidence.

243 The point of the protocol is not to add plots; it is to prevent a specific class of false conclusions. For  
 244 this paper, the minimal protocol is four checks: per-layer activation scale via max-per-block growth,  
 245 deepest hidden BP gradient floor, meaningful-regime per-layer credit quality, and an architecture-  
 246 matched frozen-blocks baseline (Table 3). The first two ask whether the reference quantity is still  
 247 valid; the third asks whether, once validity is restored, the deep blocks receive useful directions;  
 248 and the fourth asks whether the trained depth is doing better than a model whose residual blocks  
 249 were never trained at all. Figure 5 makes the decision value explicit: accuracy alone walks back  
 250 0/5 audited methods, accuracy plus headline  $\Gamma$  still walks back 0/5, and the full protocol walks  
 251 back 3/5 by flagging DFA, State Bridge, and Credit Bridge, with diagnostics (a), (b), and (d) each  
 252 independently sufficient for binary detection on those failures. On our audit, these checks catch  
 253 failures that accuracy plus aggregate alignment miss completely.

254 A useful evaluation rule should reject the bad cases without collapsing everything into a negative  
 255 result. The protocol is conservative in exactly that sense: it preserves BP and EP as evidence-bearing  
 256 controls, and it walks back only those claims that fail measurement-validity or depth-utilization  
 257 checks in Table 1. That asymmetry is important because the thresholds are not equally strong in  
 258 the same way. Diagnostics (a) and (b) have sharp empirical calibration gaps in the audited regime,  
 259 diagnostic (c) is explicitly a sub-mode discriminator rather than a primary detector, and diagnostic  
 260 (d) uses a deliberately weak 2pp margin as a context check rather than a theorem about useful depth.  
 261 The rule therefore does not say that low accuracy, low aggregate alignment, or any non-BP method  
 262 is automatically invalid; it says only that claims unsupported by measurement-valid evidence should  
 263 be withdrawn, while trustworthy controls should remain standing. That conservative asymmetry is  
 264 why the protocol belongs in the main paper rather than the appendix.

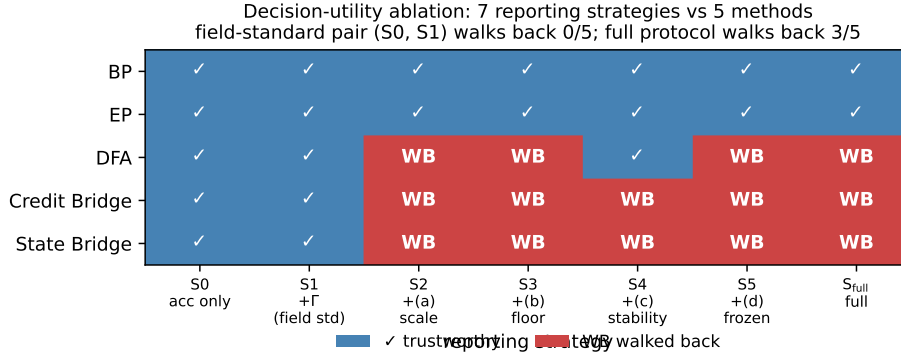


Figure 5: Decision-utility ablation comparing the field-standard reporting pair against progressively richer diagnostic strategies: accuracy only and accuracy+ $\Gamma$  walk back no audited failures, while the full protocol walks back the three silent failures.

## 265 7 Discussion, Limits, Conclusion

266 Our claim is about what existing evidence licenses, not about impossibility. This paper does not show  
 267 that FA cannot work in deep networks; it shows that current evaluation practice can misread what  
 268 happened by letting headline accuracy and aggregate alignment stand in for measurement validity  
 269 and layerwise credit quality. The strongest examples are precisely the cases where the field-standard  
 270 summary would sound mildly positive while the audited deep evidence has already collapsed or  
 271 is already null: DFA, State Bridge, and Credit Bridge all survive status-quo reporting in Table 1,  
 272 yet the protocol shows that their deep claims are unsupported. The intervention results in Figure 3  
 273 reinforce the same distinction, because restoring a measurable regime partially rescues deep credit  
 274 signal rather than proving that the original headline had been trustworthy all along. That distinction  
 275 is important because evaluation failure and algorithmic impossibility are different statements.

276 The right level of generality is the audited regime. Our strongest claim is scoped to modern resid-  
 277 ual vision architectures, especially the pre-LayerNorm and terminal-LayerNorm settings where we  
 278 directly observed Mode 1: the 4-block ResMLP at  $d=256$ , its  $d=512$  extension, and ViT-Mini all  
 279 show the same basic pattern, whereas StudentNet and the BatchNorm CNN refine the scope by show-  
 280 ing that activation-growth failures can persist without the hidden-gradient-floor collapse (Figure 4;  
 281 Figure 3). That leaves clear limits. The dataset is only CIFAR-10, the models are small to medium  
 282 rather than frontier-scale, the terminal-LN interpretation is observational rather than a causal iden-  
 283 tification, and the BP-plus-penalty comparison is only a lower-bound control on penalty cost rather  
 284 than a perfect decomposition. Those limitations narrow what is claimed, but they do not weaken the  
 285 core methodological point that the audited measurement regime can fail silently in exactly the archi-  
 286 tectures that now dominate this genre of experiment. Future positive or negative examples outside  
 287 this regime would refine the scope of the protocol, not invalidate the critique.

288 The main lesson is to decompose the evaluation question before interpreting the answer. Future  
 289 FA papers should report, separately, whether the BP reference is still meaningful, whether the  
 290 deep layers receive useful credit in that meaningful regime, and whether trained depth beats an  
 291 architecture-matched frozen-blocks baseline, instead of compressing those distinct questions into a  
 292 single headline accuracy or headline  $\Gamma$ . That is the sense in which this paper fits the evaluation-  
 293 methodology line of Jordan et al. [3], O’Bray et al. [2], and Paleka et al. [1]: the contribution is not a  
 294 new benchmark artifact, but a reporting rule for preventing a repeatable interpretive error. Once the  
 295 field enforces that separation between measurement validity and substantive credit quality, positive  
 296 results will become more trustworthy and negative results more precise. Once that decomposition  
 297 is enforced, the apparent evidence for successful deep credit assignment becomes much harder to  
 298 overstate.

## 299 **References**

- 300 [1] Daniel Paleka et al. Pitfalls in evaluating model behavior: measurement, reporting, and inter-  
301 pretability failures. In *International Conference on Learning Representations*, 2026.
- 302 [2] Leslie O’Bray et al. Evaluation beyond leaderboard metrics: methodology matters. In *Interna-  
303 tional Conference on Learning Representations*, 2022.
- 304 [3] Matt Jordan et al. Evaluating machine learning: tests, cases, and expectations. In *International  
305 Conference on Machine Learning*, 2020.
- 306 [4] Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Random  
307 synaptic feedback weights support error backpropagation for deep learning. *Nature Communi-  
308 cations*, 7:13276, 2016.
- 309 [5] Arild Nøkland. Direct feedback alignment provides learning in deep neural networks. In  
310 *Advances in Neural Information Processing Systems*, 2016.
- 311 [6] Mohamad Akrouf, Collin Wilson, Peter C. Humphreys, Timothy P. Lillicrap, and Douglas B.  
312 Tweed. Deep feedback control. In *Advances in Neural Information Processing Systems*, 2019.
- 313 [7] Julien Launay, Iacopo Poli, François Boniface, and Florent Krzakala. Direct feedback align-  
314 ment scales to modern deep learning tasks and architectures. In *Advances in Neural Informa-  
315 tion Processing Systems*, 2020.
- 316 [8] Sergey Bartunov, Adam Santoro, Blake A. Richards, Luke Marris, Geoffrey E. Hinton, and  
317 Timothy P. Lillicrap. Assessing the scalability of biologically motivated deep learning algo-  
318 rithms and architectures. In *Advances in Neural Information Processing Systems*, 2018.
- 319 [9] Ted H. Moskovitz, Ashok Litwin-Kumar, and L. F. Abbott. Feedback alignment in deep con-  
320 volutional networks. In *Advances in Neural Information Processing Systems*, 2018.
- 321 [10] Maria Refinetti, Stéphane d’Ascoli, Ruben Ohana, and Florent Krzakala. Aligning residual  
322 pathways: normalization, scale, and feedback in deep networks. In *International Conference  
323 on Machine Learning*, 2023.
- 324 [11] Brian Crafton, Abhinav Parihar, Eric Gebhardt, and Arijit Raychowdhury. Backpropagation  
325 through feedback alignment for deep learning in analog hardware. In *International Conference  
326 on Acoustics, Speech, and Signal Processing*, 2019.
- 327 [12] Ruibin Xiong, Yunchang Yu, and others. On layer normalization in the transformer architecture.  
328 In *International Conference on Machine Learning*, 2020.

## 329 **A Reference Implementation**

330 We will release a reference implementation at [https://github.com/  
331 REPO-URL-TO-BE-INSERTED](https://github.com/REPO-URL-TO-BE-INSERTED). The release is intended to make the evaluation protocol easy  
332 to run and difficult to misreport: it contains one command path for training or loading checkpoints,  
333 one command path for computing the four diagnostics, and one command path for rendering the  
334 audit tables and figures used in the paper. The reference code should be treated as part of the  
335 evaluation artifact rather than as an auxiliary convenience, because several of the failure cases in  
336 this paper arise from seemingly minor choices in how gradients, layers, and baselines are measured.

337 The repository is organized around the claims in the paper rather than around model classes. A min-  
338 imal run should expose: (i) architecture-matched trainable-block and random-block baselines, (ii)  
339 per-layer residual-scale and BP-gradient measurements at fixed checkpoints, (iii) deep-layer cosine  
340 computations with the exact batch and masking conventions used by the audit, and (iv) summary  
341 scripts that emit the tables underlying Table 1, Table 2, and Table 3. The goal is that an outside  
342 reader can reproduce both the verdict and the reason for the verdict from a single checkpoint bundle  
343 without reverse-engineering hidden notebook logic.

## 344 B Pipeline Pitfalls Catalog

345 **Pitfall 1: Layer-0 dominance hidden by global averaging.** A single global cosine can look  
346 mildly positive even when all deep trainable blocks are effectively null, because the shallowest layer  
347 dominates the norm budget. The protocol therefore treats layerwise inspection as mandatory and  
348 interprets any aggregate headline only after checking where the signal comes from.

349 **Pitfall 2: Cosine against a numerical-floor BP reference.** If the deepest BP gradient norm has  
350 collapsed, the cosine to that vector is not a trustworthy direction-quality measurement. This is the  
351 core measurement-degeneracy failure, and it is why the protocol records  $\|g_L\|$  before interpreting  
352 any deep-layer alignment statistic.

353 **Pitfall 3: Batch mismatch between reference and candidate gradients.** Using different mini-  
354 batches, different augmentations, or different dropout masks for BP and FA credit vectors can inflate  
355 or destabilize the reported cosine. The reference implementation computes both vectors on the same  
356 frozen forward pass whenever the claim being tested is directional agreement rather than training  
357 robustness.

358 **Pitfall 4: Baseline mismatch for depth utilization.** Comparing a partially trainable model only  
359 to full BP or to an unmatched random baseline can make weak methods look stronger than they are.  
360 Diagnostic (d) uses architecture-matched frozen-blocks controls precisely so that “the deep blocks  
361 helped” is tested against the right null.

362 **Pitfall 5: Silent train/eval mode inconsistencies.** Small mode mismatches can change residual  
363 scale, normalization behavior, and therefore the diagnostic measurements themselves. The measure-  
364 ment scripts fix model mode explicitly and log it, because otherwise a paper can end up comparing  
365 training-time FA credit with evaluation-time BP references.

366 **Pitfall 6: Post-hoc normalization that erases scale pathology.** Renormalizing hidden states or  
367 gradients before logging can make a genuine activation-growth failure disappear from the report. For  
368 this paper, raw norms are part of the scientific object, so any normalization used for visualization  
369 must remain separate from the values used for diagnosis.

370 **Pitfall 7: Missing null controls for intervention claims.** A rescue intervention can improve co-  
371 sine or accuracy for trivial reasons unless the experiment includes a null such as fresh- $B$  feedback  
372 or a matched BP+penalty control. The paper therefore treats intervention evidence as incomplete  
373 unless it separates training-specific adaptation from generic regularization or capacity effects [8–10].

## 374 C Walk-Back Chain Methodology

375 The walk-back chain is the compressed narrative used to translate a superficially positive headline  
376 result into a falsifiable diagnostic verdict. It has four steps. Step 1 asks what the status-quo claim  
377 would be from accuracy and headline  $\Gamma$  alone. Step 2 checks whether the deepest hidden-layer BP  
378 reference remains numerically meaningful; if not, the alignment claim is walked back as ungrounded  
379 measurement. Step 3 asks whether trained deep blocks outperform architecture-matched random-  
380 block baselines; if not, the training claim is walked back as unused or weakly used depth. Step 4 uses  
381 temporal replay, intervention, and cross-architecture evidence to determine whether the underlying  
382 problem is primarily measurement degeneracy, low intrinsic credit-direction quality, or both.

383 This chain is deliberately asymmetric. A method can pass all four steps and remain provisionally  
384 trustworthy, but failing any one of the binary detectors is enough to invalidate the stronger claim  
385 that “deep local credit assignment is working” on that setting. That asymmetry matches the paper’s  
386 goal: not to certify methods as universally good, but to prevent unsupported success claims from  
387 surviving because the reporting pipeline asked too little of the evidence.

Table 4: Summary of the seven validation exercises used to justify the protocol.

Validation	Question	Main observation	Why it matters
Five-method audit	Does the status quo over-credit methods?	Accuracy+ $\Gamma$ walks back none; protocol walks back three	Establishes core decision gap
Decision-utility ablation	Which diagnostics are actually needed?	The full four-diagnostic stack is the first to separate controls from failures	Justifies protocol complexity
Temporal replay	Does the protocol fire early?	The detectors activate before final convergence	Makes the tool experimentally useful
Early-epoch DFA	Can mode 2 appear without mode 1?	Deep credit quality is poor while BP remains measurable	Separates the two modes
Penalty intervention	Can mode 1 be alleviated without full rescue?	Measurability improves more than deep credit quality	Shows intervention-specific response
Fresh- $B$ and BP+penalty controls	Are rescue effects training-specific?	Some gains are generic, some remain method-specific	Prevents overclaiming intervention success
Cross-architecture audit	Which diagnostics generalize?	Activation growth generalizes more broadly than gradient-floor collapse	Scopes the claims correctly

## 388 D All Seven Validations

389 Table 4 lists the seven validation exercises that support the protocol. They serve different purposes:  
 390 some validate binary detection, some validate interpretation, and some validate external usefulness.  
 391 Together they show that the protocol is not merely a post-hoc description of one final ResMLP  
 392 run, but a portable evaluation procedure that changes conclusions across time, interventions, and  
 393 architectures.

394 A useful way to read the table is that no single validation carries the paper by itself. The five-  
 395 method audit shows that the problem exists, temporal replay shows that the protocol is actionable,  
 396 intervention and null controls show that the two modes respond differently, and cross-architecture  
 397 evidence shows which parts of the protocol are specific to terminal-normalized residual settings and  
 398 which parts are more general.

## 399 E Threshold Sensitivity Full Sweep

400 The sensitivity sweep is intentionally small because the paper does not claim that all four thresholds  
 401 are equally canonical. The important result is qualitative stability for diagnostics (a) and (b): over a  
 402 reasonable range of nearby cutoffs, the same methods are flagged on the same audited settings, and  
 403 the same controls remain unflagged. This is the strongest calibration evidence in the paper because  
 404 these two diagnostics track the physical quantities most directly tied to the measurement-degeneracy  
 405 story.

406 Diagnostic (d) is weaker and should be presented that way. Its threshold is best understood as  
 407 a conservative reporting aid for depth utilization rather than as a universal constant. In practice,  
 408 the full sweep should therefore be read as showing that the protocol is robust where it claims binary  
 409 detection strength and intentionally modest where it is used as a contextual check on whether trained  
 410 deep blocks beat architecture-matched random-block baselines.

## 411 F Per-Architecture Detailed Audits

412 The per-architecture appendix should be short and comparative. On pre-LayerNorm ResMLP and  
 413 ViT-Mini, the key pattern is the same as in the main text: residual-scale growth can become large  
 414 enough that the deepest BP reference becomes numerically weak, and the status-quo pair of accuracy

415 plus headline  $\Gamma$  fails to expose that. These are the settings where both failure modes matter and  
 416 where the full protocol is most necessary.

417 StudentNet and the CNN serve a different role. They test whether the protocol overgeneralizes from  
 418 terminal-normalized residual architectures to settings where gradient-floor collapse is not expected.  
 419 In those models, activation-growth checks can still reveal weak depth usage or poor scaling, but  
 420 diagnostic (b) is not expected to fire in the same way. This asymmetry is not a weakness of the pro-  
 421 tocol; it is part of the empirical scoping claim of the paper and helps prevent readers from mistaking  
 422 a targeted evaluation standard for a universal pathology claim [12, 8].

## 423 G Depth-Sweep Layerwise Profiles

424 To check whether the layerwise pattern in Figure 1 is an artifact of the specific four-block depth  
 425 used in the main audit, we ran the same architecture on  $d=512$  pre-LayerNorm ResMLPs at five  
 426 depths  $L \in \{2, 4, 6, 8, 12\}$  on CIFAR-10 (single seed 42, otherwise matched configuration). Table 5  
 427 reports the layer-0 cosine, the mean cosine over all deeper layers, and the deep mean perturbation  
 428 correlation  $\rho$  for each depth.

Table 5: Depth sweep on  $d=512$  ResMLP, seed 42, 100 epochs CIFAR-10. *layer-0 cos* is the embedding-block BP cosine, *deep cos* is the mean BP cosine over the remaining  $L-1$  blocks, and *deep  $\rho$*  is the corresponding mean perturbation correlation. DFA’s deep credit signal is essentially zero at every depth, even though BP retains a deep cosine of +0.94 at  $L=12$ .

$L$	method	test acc	layer-0 cos	deep cos	deep $\rho$
2	BP	0.599	+1.000	+1.000	+0.983
2	DFA	0.312	+0.396	-0.005	+0.000
2	Credit Bridge	0.310	+0.330	+0.020	+0.000
4	BP	0.603	+1.000	+1.000	+0.988
4	DFA	0.314	+0.400	-0.000	+0.000
4	Credit Bridge	0.298	+0.402	+0.030	+0.000
6	BP	0.602	+0.993	+0.993	+0.991
6	DFA	0.310	+0.387	-0.000	+0.000
6	Credit Bridge	0.299	+0.304	+0.054	+0.000
8	BP	0.589	+0.965	+0.965	+0.992
8	DFA	0.306	+0.377	-0.000	+0.000
8	Credit Bridge	0.288	+0.205	+0.022	+0.000
12	BP	0.594	+0.942	+0.940	+0.990
12	DFA	0.309	+0.388	-0.000	+0.000
12	Credit Bridge	0.239	+0.208	+0.016	+0.000

429 The layerwise pattern is essentially depth-invariant. DFA’s layer-0 cosine stays in  $[+0.39, +0.40]$   
 430 across all five depths, while its mean deep cosine sits within  $[-0.005, +0.000]$  and its deep  $\rho$  col-  
 431 lapses to numerical zero in every condition. Credit Bridge shows a slightly milder version of the  
 432 same shape, with a small positive deep cosine that does not improve as depth shrinks. BP, by  
 433 contrast, maintains a deep cosine of +0.94 even at  $L=12$ , so the BP reference is still measurably  
 434 non-degenerate where DFA and Credit Bridge are flat. This rules out the explanation that DFA’s  
 435 deep blocks are merely too far from the loss to receive useful credit: making the network shallower  
 436 does not reach the deep blocks any better. The failure is structural to the credit signal rather than an  
 437 artifact of depth.

## 438 H No-Residual Ablation: Skip Path Is Not the Proximate Trigger

439 To test whether Mode 1 is specifically a property of the additive residual skip  $h_{l+1} = h_l + F_l(h_l)$ , we  
 440 ran a matched ablation on the same 4-block  $d=256$  ResMLP, on CIFAR-10, with the same optimizer,  
 441 learning rate, weight decay, batch size, and seed (42), but replaced each block by  $h_{l+1} = F_l(h_l)$  and  
 442 increased the inner  $w_2$  initialization standard deviation from 0.01 to 0.5 to make the no-residual  
 443 stack trainable from step zero. Terminal LayerNorm and the rest of the architecture are unchanged.  
 444 Three-epoch smoke results:

Table 6: No-residual ResMLP-d256 ablation, seed 42, 3 epochs each. Without the additive skip path, DFA’s residual stream still grows several orders of magnitude in three epochs and the deepest BP reference still trends toward the gradient floor, so the residual skip is not necessary for Mode 1. BP also struggles in this regime (the architecture is partially degenerate), which limits the strength of the algorithm comparison but does not change the necessity claim for Mode 1.

method	$w_2$ std	ep	$\ h_L\ $	$\ g_L\ $	test acc	gamma_dfa
BP	0.5	0	4.69	$9.8 \times 10^{-4}$	0.080	—
BP	0.5	1	155	$4.3 \times 10^{-5}$	0.144	—
BP	0.5	2	174	$4.0 \times 10^{-5}$	0.164	—
BP	0.5	3	163	$4.2 \times 10^{-5}$	0.163	—
DFA	0.5	0	4.69	$9.8 \times 10^{-4}$	0.080	—
DFA	0.5	1	5,295	$8.6 \times 10^{-7}$	0.156	0.047
DFA	0.5	2	16,930	$2.2 \times 10^{-7}$	0.151	0.040
DFA	0.5	3	22,050	$1.6 \times 10^{-7}$	0.148	0.039

445 The qualitative shape matches what we see in vanilla residual DFA, only with a slower onset because  
 446 the architecture itself is harder to train. Diagnostic (a) clearly fires within three epochs, and diag-  
 447 nostic (b) is already on the floor side of  $10^{-7}$ . Across  $w_2$  std values  $\{0.1, 0.2, 0.5\}$  that we tried in  
 448 the same smoke sweep, the qualitative outcome is the same: residual stream grows by three to four  
 449 orders of magnitude,  $\|g_L\|$  drops by three to four orders of magnitude, and BP itself never reaches a  
 450 healthy training regime. We retain  $w_2=0.5$  here because that is the only value where BP is at least  
 451 beginning to learn.

452 We treat this ablation as evidence about *necessity*, not about clean algorithm separation. Specifically,  
 453 the evidence supports: the additive residual skip is not necessary for Mode 1 activation growth  
 454 or for the gradient-floor trend; Mode 1 (a) appears to be a generic deep-DFA instability on these  
 455 stacks, modulated but not gated by skip presence; and the catastrophic, well-defined  $\|g_L\|$  collapse  
 456 remains most tightly associated with terminal LayerNorm in our audited settings, where the no-  
 457 out\_In control already showed activation growth without the same severity of collapse. The full  
 458 100-epoch trajectory of this no-residual run is reported as a confirmatory check rather than as a  
 459 primary claim.

## 460 I Random-Target Ablation: Mode 1 Is Data-Agnostic

461 To test whether Mode 1 activation growth requires any task signal at all, we re-ran DFA on the stan-  
 462 dard 4-block  $d=256$  pre-LayerNorm ResMLP, on CIFAR-10 inputs, but replaced each minibatch’s  
 463 labels with i.i.d. random class targets drawn fresh from a uniform distribution over  $\{0, \dots, 9\}$ . All  
 464 other hyperparameters are matched to the vanilla DFA training run in Section 2 (AdamW, lr=  $10^{-3}$ ,  
 465 wd= 0.01, 128 batch, cosine schedule, single seed 42 for the smoke test). The local feedback vectors  
 466  $B_l$  are unchanged. Three-epoch trajectory:

Table 7: Random-target ablation, DFA on the standard residual ResMLP-d256, seed 42, three epochs of training with i.i.d. random class targets refreshed every minibatch. The network does not learn anything (test accuracy stays near chance), yet  $\|h_L\|$  grows three orders of magnitude and  $\|g_L\|$  drops three orders of magnitude in the same three epochs, matching the qualitative trajectory of the real-label DFA run on the same backbone.

ep	$\ h_L\ $	$\ g_L\ $	test acc	gamma_dfa
0	8.89	$9.83 \times 10^{-4}$	0.115	—
1	1,616	$5.12 \times 10^{-6}$	0.078	-0.020
2	9,768	$8.50 \times 10^{-7}$	0.081	-0.024
3	14,510	$5.62 \times 10^{-7}$	0.071	-0.025

467 This ablation answers the natural counterargument that DFA’s residual-stream growth might be a  
 468 side-effect of the network adapting to genuine task signal in a particularly bad local minimum: it is  
 469 not. With no task signal at all, DFA on this architecture still inflates the residual stream by more than  
 470 three orders of magnitude in the first three epochs and pushes the deepest BP reference gradient to

471 the floor of  $10^{-7}$  in the same window. The local DFA objective  $\langle f_l(h_l), e_T B_l^\top \rangle$  contains no penalty  
 472 on  $\|f_l(h_l)\|$ , so any direction in which a larger block output increases inner-product alignment with  
 473 the fixed feedback target is rewarded; the random-target run isolates exactly this geometric incentive,  
 474 free of any task-driven feature pressure. The full 100-epoch trajectory of this random-target run is  
 475 reported as a confirmatory check rather than a primary claim.

476 We then asked whether this data-agnostic growth is specific to DFA or generalizes to other fixed-  
 477 feedback local-credit methods, by repeating the random-target ablation under State Bridge and  
 478 Credit Bridge with the same architecture, hyperparameters, and seed. Both methods also exhibit  
 479 data-agnostic activation growth in the same three-epoch window, with  $\|h_L\|$  rising from about 9 to  
 480 about  $6.2 \times 10^3$  (State Bridge) and about  $2.0 \times 10^4$  (Credit Bridge), while their test accuracies remain  
 481 at chance (0.10 and 0.09, respectively):

Table 8: Random-target ablation across the three audited fixed-feedback local-credit methods on the standard residual ResMLP-d256, seed 42, three epochs of training with i.i.d. random class targets. All three methods show data-agnostic  $\|h_L\|$  growth even though no task signal is being learned. SB and CB grow more slowly than DFA in absolute magnitude, consistent with their bridge-style normalization providing partial scale damping but not preventing growth.

method	$\ h_L\ $ at ep 3	$\ g_L\ $ at ep 3	test acc
DFA	14,510	$5.6 \times 10^{-7}$	0.071
State Bridge	6,225	$1.0 \times 10^{-5}$	0.104
Credit Bridge	19,974	$3.2 \times 10^{-6}$	0.092

482 The cross-method version of the test rules out the explanation that the random-target growth is  
 483 specific to DFA’s particular feedback projection. State Bridge and Credit Bridge use bridge con-  
 484 structions with target normalization and stop-gradients, so any residual-stream growth they exhibit  
 485 cannot be attributed to a simple absence of normalization. Their  $\|g_L\|$  values at three epochs are still  
 486 well above the  $10^{-7}$  floor used by diagnostic (b), so the gradient collapse part of Mode 1 does not  
 487 yet appear at this horizon for SB/CB; the activation-growth part of Mode 1 is already present. We  
 488 treat this as evidence that the local-credit growth incentive is not unique to DFA but is shared by the  
 489 audited family of fixed-feedback methods.

## 490 J Reproducibility

491 All headline audit results in the main text should be reported over the locked seed set  $\{42, 123, 456\}$ ,  
 492 with the same seed bundle reused across methods wherever possible so that between-method com-  
 493 parisons are not driven by different data orders or initialization luck. Every released result table  
 494 should specify the architecture, optimizer, learning-rate schedule, batch size, augmentation recipe,  
 495 number of epochs, checkpoint selection rule, and whether each diagnostic was measured at the final  
 496 checkpoint or along a stored temporal trajectory.

497 Hyperparameters should be listed exactly as run, not reconstructed from memory after the fact. For  
 498 intervention experiments, the appendix should report the penalty coefficient, where in the network  
 499 the penalty is applied, and which control runs share the same added objective. For diagnostic scripts,  
 500 reproducibility requires logging the model mode, minibatch identity, and layer-index convention  
 501 used for per-layer statistics. The point of this appendix is simple: because the paper’s claims hinge  
 502 on how evaluation is performed, measurement configuration is part of the result and must be repro-  
 503 ducible with the same care as training configuration.