

---

# 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 (learning  
 62 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)

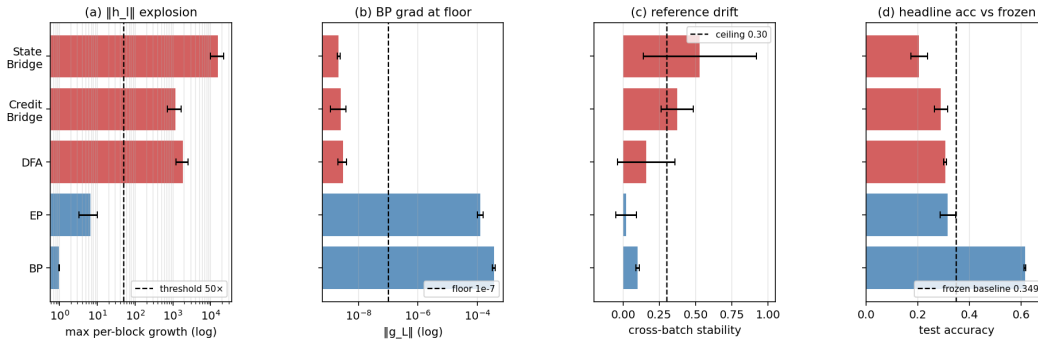


Figure 1: Five-method audit on the 4-block  $d=256$  pre-LayerNorm ResMLP: the field-standard pair looks superficially consistent across methods, but the diagnostic view separates trustworthy controls from walked-back methods.

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  
 96 objectives of the form  $\langle f_l(h_l), e_T B_l^T \rangle$  with no explicit scale constraint on  $f_l$ , so the residual stream  
 97 is free to inflate while still reducing the local loss [7]. In the same runs, each block’s  $w_1$  and  $w_2$   
 98 grows by roughly  $200\times$  in relative delta, their norm product reaches about  $5 \times 10^4$  per block, and  
 99 the terminal hidden-state norm  $\|h_L\|$  rises monotonically from about 9 at random initialization to  
 100 about  $4 \times 10^8$  by epoch 100 (Figure 2). Most of that growth appears immediately:  $\|h_L\|$  already  
 101 reaches about  $10^6$  by epoch 5. Once the residual stream reaches this regime, the backpropagation  
 102 reference vector no longer behaves like a healthy target.

103 The measurement failure occurs at the point where the hidden-layer BP gradient ceases to be a mean-  
 104 ingful reference direction. In terminal-LayerNorm architectures, the LayerNorm Jacobian scales  
 105 as  $\partial \text{LN}(h)/\partial h \propto 1/\|h\|$  in expectation, so the same residual-stream inflation is accompanied by  
 106 collapse of the hidden-layer BP reference norm: on DFA-trained ResMLP,  $\|g_L\|$  falls from about  
 107  $9.8 \times 10^{-4}$  at random initialization to about  $5 \times 10^{-10}$  by epoch 100, a six-order-of-magnitude drop,  
 108 while the reported cosine remains mathematically defined only because `F.cosine_similarity`  
 109 clamps the denominator at  $\varepsilon=10^{-8}$  (Table 1; Figure 1). At that endpoint the reference norm is about  
 110  $20\times$  below the clamp, so the quantity being reported is effectively  $(a \cdot b)/(\|a\| \max(\|b\|, 10^{-8}))$   
 111 rather than a comparison to an informative BP direction. At that point, reporting a cosine is no  
 112 longer evidence about credit quality.

113 The simplest control is architectural, not theoretical. On the same ResMLP backbone, BP keeps  
 114  $\|h_L\|$  near 200 and  $\|g_L\|$  near  $4 \times 10^{-4}$  throughout training, while EP keeps  $\|h_L\|$  around  $5 \times 10^3$   
 115 and  $\|g_L\|$  around  $1.3 \times 10^{-4}$ , so hard optimization on CIFAR-10 by itself does not force hidden-

116 layer gradients to the numerical floor (Table 1; Figure 2). The broader cross-architecture pattern is  
 117 consistent with the same interpretation: StudentNet and the BatchNorm CNN, which lack terminal  
 118 LayerNorm, keep deepest BP gradients around  $10^{-4}$  and never trigger diagnostic (b), whereas ViT-  
 119 Mini with a terminal LN shows the same collapse pattern and triggers diagnostic (b) by epochs 2–3  
 120 (Figure 2). The pathology therefore belongs to the evaluated FA regime, not to CIFAR-10 or the  
 121 backbone alone.

122 The collapse is not a late-epoch curiosity. For vanilla DFA on the ResMLP temporal replay,  $\|g_L\|$   
 123 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  
 124 epoch 3, and  $6.7 \times 10^{-8}$  at epoch 4, so diagnostic (b) fires at epoch 3–4 across all three seeds, while  
 125 the max-per-block growth detector fires slightly later at epochs 8–11 (Figure 2). Both detectors  
 126 therefore fire in the first 11 epochs of a 100-epoch run, making the protocol actionable as an early-  
 127 stop criterion rather than a post hoc explanation. The practical point is reinforced by accuracy: DFA  
 128 is at 0.308 already at epoch 4 and ends at 0.306 by epoch 100, so the remaining training budget  
 129 adds essentially nothing to the headline result once the measurement has already degenerated. Once  
 130 measurement degeneracy is identified, the next question is whether poor deep credit remains even  
 131 before collapse.

## 132 4 Failure Mode 2: Low Intrinsic Credit-Direction Quality

133 The second failure mode appears even in the meaningful-measurement regime. At the earliest vanilla  
 134 DFA checkpoints on ResMLP, the hidden backpropagated gradient at the first deep block remains  
 135 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  
 136 seeds, all above the  $10^{-7}$  threshold used to distinguish measurable from collapsed gradients. Yet the  
 137 corresponding deep-layer cosine values are already essentially null: across layers 1–4, all seed-level  
 138 measurements at epoch 1 lie in  $[-0.04, +0.02]$ , with a three-seed mean of  $-0.008 \pm 0.013$ , and  
 139 by epoch 2 the deep mean is still only  $-0.018 \pm 0.018$  (Table 2). This is the observational pattern  
 140 predicted by low credit-direction quality rather than mere disappearance of signal: the gradient is  
 141 still present enough to measure, but the directions delivered to the deep network carry little agree-  
 142 ment with backpropagation, consistent with prior concerns that alternative feedback rules can fail by  
 143 supplying poor credit assignments even before full collapse [8, 9, 11? ]. This rules out the simplest  
 144 objection that the deep-layer null result is merely a byproduct of collapse.

145 A second metric with different numerical failure modes tells the same story. Cosine measures di-  
 146 rectional agreement with the BP gradient, whereas perturbation correlation  $\rho$  measures whether the  
 147 proposed update predicts the correct sign and relative magnitude of loss change under actual per-  
 148 turbations; their failure modes are therefore different, especially with respect to normalization and  
 149 small-denominator effects. In our controls,  $\rho$  behaves as expected, with a Taylor-ceiling positive  
 150 control near  $+0.997$  and a random-vector negative control near  $+0.006$  (Figure 3, Table 2). On  
 151 vanilla DFA, deep  $\rho$  is likewise null: for the early checkpoints where the gradients remain measur-  
 152 able, the deep average is  $-0.003 \pm 0.005$  across seeds and epochs, and in a floor-level checkpoint it is  
 153  $+0.002$ , again indistinguishable from noise. The agreement between cosine and  $\rho$  therefore rules out  
 154 the interpretation that the null deep result is an artifact of cosine’s  $\varepsilon$ -clamp or vector normalization.  
 155 The deep blocks are not just hard to measure; they are receiving weakly useful directions.

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

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

## 166 5 Intervention and Cross-Architecture Evidence

167 The penalty intervention first matters as a rescue of the measurement regime. When we add a per-  
 168 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  
 169 epochs on CIFAR-10, the  $\lambda=10^{-2}$  setting contains the terminal hidden-state scale from  $\|h_L\| \sim$   
 170  $4.4 \times 10^8$  under vanilla DFA to  $\sim 4.0 \times 10^4$ , while lifting the deepest BP reference norm from  
 171  $\|g_L\| \sim 5 \times 10^{-10}$  to  $\sim 9.0 \times 10^{-7}$ , a roughly four-order-of-magnitude rescue on both quantities  
 172 (Figure 3; Table 2). At that setting, both diagnostic (a) and diagnostic (b) pass on penalized DFA,  
 173 and test accuracy rises to  $0.363 \pm 0.001$  from  $0.308 \pm 0.014$  for vanilla DFA. The key point is not  
 174 yet that the recovered network has good deep credit, but that the deep reference vector is again large  
 175 enough to function as a meaningful target direction rather than a clamp-dominated artifact. That  
 176 rescue makes the second question measurable rather than hypothetical.

177 Once the reference vector is meaningful again, the deep layers no longer sit exactly at null. At  
 178  $\lambda=10^{-2}$ , penalized DFA reaches a three-seed deep-layer mean cosine of  $+0.155 \pm 0.025$  and deep  
 179 perturbation correlation of  $+0.080 \pm 0.011$ , whereas vanilla DFA is essentially zero on both metrics  
 180 in the deep blocks, consistent with prior concerns that alternative feedback can fail by supplying  
 181 poor credit directions even before full collapse [8, 9, 11? ]. The null calibration rules out the inter-  
 182 pretation that this recovered signal is merely measurement noise: on the same penalized checkpoint,  
 183 replacing the training-time feedback matrices with 20 fresh random  $B_l$  draws gives a deep cosine  
 184 of only  $+0.002 \pm 0.022$ , with per-layer standard deviations of 0.013–0.023, all within noise of zero  
 185 (Table 2). The  $\lambda$  sweep sharpens the dissociation further: at  $\lambda=10^{-4}$ , Mode 1 is already alleviated,  
 186 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  
 187 rises to  $+0.165$  and deep  $\rho$  to  $+0.091$  (Figure 3). The improvement is real, but it is only partial.

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

198 The architecture comparison sharpens the scope of the critique. In the terminal-LN architectures we  
 199 audited, both diagnostics fire for DFA-trained ResMLP at  $d=256$ , the same pattern recurs at  $d=512$   
 200 with even larger max-per-block growth (about  $1.5 \times 10^4$ ), and ViT-Mini with a class token and  
 201 terminal LN shows diagnostic (a) by epoch 1 and diagnostic (b) by epochs 2–3 (Figure 2). In the non-  
 202 terminal-LN controls, the pattern is different: StudentNet shows diagnostic (a) only at epochs 14–25  
 203 while diagnostic (b) never fires across 100 epochs and three seeds, and the BatchNorm CNN on  
 204 CIFAR-10 likewise shows strong growth under DFA, with max-per-block growth up to  $237\times$ , but  
 205 keeps deepest BP gradients around  $\|g\| \sim 10^{-3}$  and never triggers diagnostic (b) (Figure 2). BP  
 206 never triggers either diagnostic in any audited architecture. This is an observational association  
 207 rather than a causal identification of terminal LayerNorm as the unique mechanism, but it is enough  
 208 to support a narrower claim: diagnostic (b) appears tied to the terminal-LN architectures audited  
 209 here, while diagnostic (a) remains useful more broadly. This lets the paper end with a reporting rule  
 210 rather than an overclaimed theory.

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

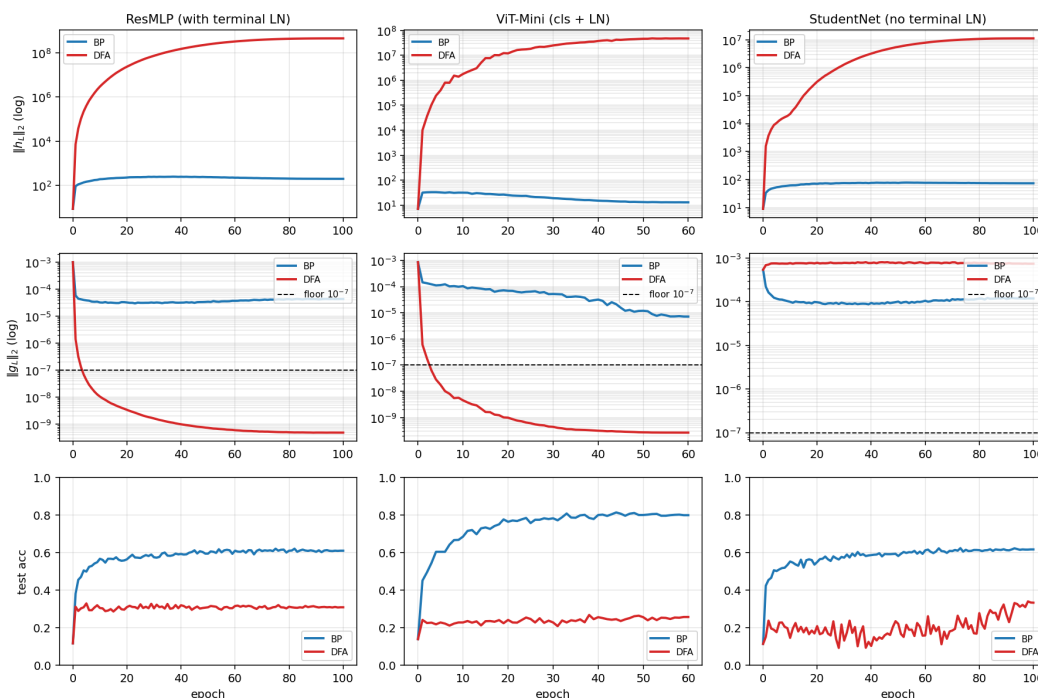


Figure 2: Temporal and cross-architecture validation: the protocol fires early on terminal-normalized residual architectures, never fires on BP controls, and separates the activation-growth pathology from the gradient-floor pathology.

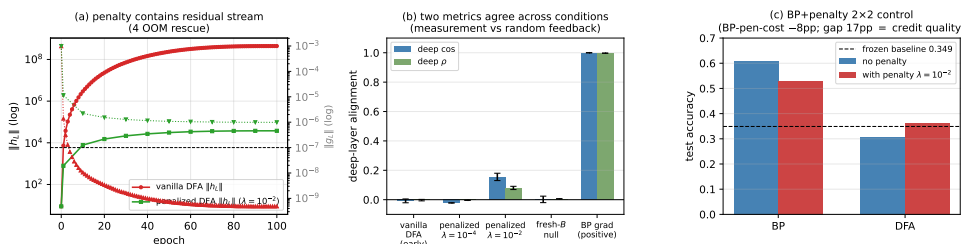


Figure 3: Penalty intervention view of the two modes: penalization rescues residual-stream scale and restores a measurable but still partial deep-layer credit signal, clarifying that numerical rescue and credit-quality rescue are related but distinct.

## 211 6 Recommended FA Evaluation Protocol

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

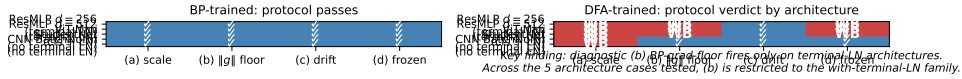


Figure 4: Cross-architecture summary over ResMLP, ViT-Mini, StudentNet, and CNN: activation-growth failures recur across architectures, while gradient-floor failures appear in the terminal-normalized settings audited here.

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

223 The point of the protocol is not to add plots; it is to prevent a specific class of false conclusions. For  
 224 this paper, the minimal protocol is four checks: per-layer activation scale via max-per-block growth,  
 225 deepest hidden BP gradient floor, meaningful-regime per-layer credit quality, and an architecture-  
 226 matched frozen-blocks baseline (Table 3). The first two ask whether the reference quantity is still  
 227 valid; the third asks whether, once validity is restored, the deep blocks receive useful directions;  
 228 and the fourth asks whether the trained depth is doing better than a model whose residual blocks  
 229 were never trained at all. Figure 5 makes the decision value explicit: accuracy alone walks back  
 230 0/5 audited methods, accuracy plus headline  $\Gamma$  still walks back 0/5, and the full protocol walks  
 231 back 3/5 by flagging DFA, State Bridge, and Credit Bridge, with diagnostics (a), (b), and (d) each  
 232 independently sufficient for binary detection on those failures. On our audit, these checks catch  
 233 failures that accuracy plus aggregate alignment miss completely.

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

## 245 7 Discussion, Limits, Conclusion

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

256 The right level of generality is the audited regime. Our strongest claim is scoped to modern resid-  
 257 ual vision architectures, especially the pre-LayerNorm and terminal-LayerNorm settings where we

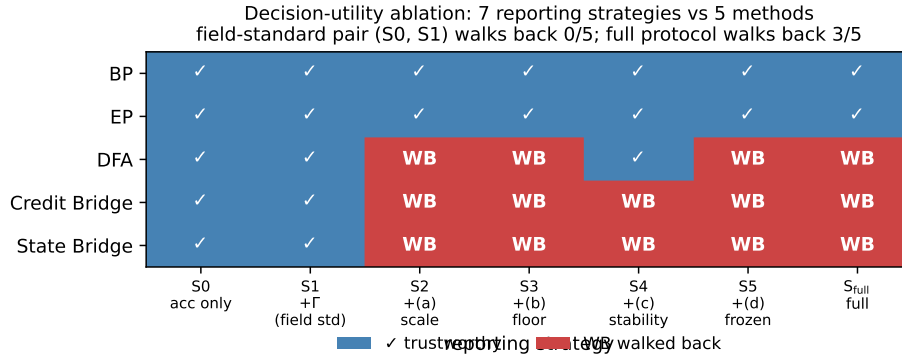


Figure 5: Decision-utility ablation comparing the field-standard reporting pair against progressively richer diagnostic strategies: accuracy only and accuracy+ $\Gamma$  walk back no audited failures, while the full protocol walks back the three silent failures.

258 directly observed Mode 1: the 4-block ResMLP at  $d=256$ , its  $d=512$  extension, and ViT-Mini all  
 259 show the same basic pattern, whereas StudentNet and the BatchNorm CNN refine the scope by show-  
 260 ing that activation-growth failures can persist without the hidden-gradient-floor collapse (Figure 4;  
 261 Figure 3). That leaves clear limits. The dataset is only CIFAR-10, the models are small to medium  
 262 rather than frontier-scale, the terminal-LN interpretation is observational rather than a causal iden-  
 263 tification, and the BP-plus-penalty comparison is only a lower-bound control on penalty cost rather  
 264 than a perfect decomposition. Those limitations narrow what is claimed, but they do not weaken the  
 265 core methodological point that the audited measurement regime can fail silently in exactly the archi-  
 266 tectures that now dominate this genre of experiment. Future positive or negative examples outside  
 267 this regime would refine the scope of the protocol, not invalidate the critique.

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

## 279 References

- 280 [1] Daniel Paleka et al. Pitfalls in evaluating model behavior: measurement, reporting, and inter-  
 281 pretability failures. In *International Conference on Learning Representations*, 2026.
- 282 [2] Leslie O’Bray et al. Evaluation beyond leaderboard metrics: methodology matters. In *Interna-  
 283 tional Conference on Learning Representations*, 2022.
- 284 [3] Matt Jordan et al. Evaluating machine learning: tests, cases, and expectations. In *International  
 285 Conference on Machine Learning*, 2020.
- 286 [4] Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Random  
 287 synaptic feedback weights support error backpropagation for deep learning. *Nature Communi-  
 288 cations*, 7:13276, 2016.
- 289 [5] Arild Nøkland. Direct feedback alignment provides learning in deep neural networks. In  
 290 *Advances in Neural Information Processing Systems*, 2016.
- 291 [6] Mohamad Akrouf, Collin Wilson, Peter C. Humphreys, Timothy P. Lillicrap, and Douglas B.  
 292 Tweed. Deep feedback control. In *Advances in Neural Information Processing Systems*, 2019.

- 293 [7] Julien Launay, Iacopo Poli, François Boniface, and Florent Krzakala. Direct feedback align-  
294 ment scales to modern deep learning tasks and architectures. In *Advances in Neural Informa-  
295 tion Processing Systems*, 2020.
- 296 [8] Sergey Bartunov, Adam Santoro, Blake A. Richards, Luke Marris, Geoffrey E. Hinton, and  
297 Timothy P. Lillicrap. Assessing the scalability of biologically motivated deep learning algo-  
298 rithms and architectures. In *Advances in Neural Information Processing Systems*, 2018.
- 299 [9] Ted H. Moskovitz, Ashok Litwin-Kumar, and L. F. Abbott. Feedback alignment in deep con-  
300 volutional networks. In *Advances in Neural Information Processing Systems*, 2018.
- 301 [10] Maria Refinetti, Stéphane d’Ascoli, Ruben Ohana, and Florent Krzakala. Aligning residual  
302 pathways: normalization, scale, and feedback in deep networks. In *International Conference  
303 on Machine Learning*, 2023.
- 304 [11] Brian Crafton, Abhinav Parihar, Eric Gebhardt, and Arijit Raychowdhury. Backpropagation  
305 through feedback alignment for deep learning in analog hardware. In *International Conference  
306 on Acoustics, Speech, and Signal Processing*, 2019.
- 307 [12] Ruibin Xiong, Yunchang Yu, and others. On layer normalization in the transformer architecture.  
308 In *International Conference on Machine Learning*, 2020.

## 309 A Reference Implementation

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

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

## 324 B Pipeline Pitfalls Catalog

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

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

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

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

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

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

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

## 354 C Walk-Back Chain Methodology

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

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

## 368 D All Seven Validations

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

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

## 379 E Threshold Sensitivity Full Sweep

380 The sensitivity sweep is intentionally small because the paper does not claim that all four thresholds  
381 are equally canonical. The important result is qualitative stability for diagnostics (a) and (b): over a  
382 reasonable range of nearby cutoffs, the same methods are flagged on the same audited settings, and  
383 the same controls remain unflagged. This is the strongest calibration evidence in the paper because

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

384 these two diagnostics track the physical quantities most directly tied to the measurement-degeneracy  
 385 story.

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

## 391 **F Per-Architecture Detailed Audits**

392 The per-architecture appendix should be short and comparative. On pre-LayerNorm ResMLP and  
 393 ViT-Mini, the key pattern is the same as in the main text: residual-scale growth can become large  
 394 enough that the deepest BP reference becomes numerically weak, and the status-quo pair of accuracy  
 395 plus headline  $\Gamma$  fails to expose that. These are the settings where both failure modes matter and  
 396 where the full protocol is most necessary.

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

## 403 **G Reproducibility**

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

410 Hyperparameters should be listed exactly as run, not reconstructed from memory after the fact. For  
 411 intervention experiments, the appendix should report the penalty coefficient, where in the network

412 the penalty is applied, and which control runs share the same added objective. For diagnostic scripts,  
413 reproducibility requires logging the model mode, minibatch identity, and layer-index convention  
414 used for per-layer statistics. The point of this appendix is simple: because the paper's claims hinge  
415 on how evaluation is performed, measurement configuration is part of the result and must be repro-  
416 ducible with the same care as training configuration.