

---

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

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Standard evaluation of Feedback Alignment (FA) and related local-credit meth-  
2 ods on modern residual networks reports two numbers: headline accuracy and the  
3 cosine alignment  $\Gamma$  of the local credit signal with the true backpropagation  
4 gradient at hidden layers. We show, on standard pre-LayerNorm ResidualMLP and  
5 ViT-Mini architectures, that this evaluation is unreliable because it conflates two  
6 distinct failure modes: **(1) measurement degeneracy via terminal-LayerNorm**  
7 **gradient cancellation**, in which residual stream growth drives the BP gradient  
8 at hidden layers below the numerical floor and renders the cosine metric unin-  
9 terpretable; and **(2) low intrinsic credit-direction quality of random feedback**,  
10 which persists even when the BP gradient is in the meaningful regime and is invis-  
11 ible to the field-standard reporting pair.

12 We contribute a four-diagnostic protocol that detects both modes, a reference im-  
13 plementation, a calibrated scale for the new metrics, and a reproducible audit table  
14 on five methods (BP, DFA, State Bridge, Credit Bridge, EP) across three architec-  
15 ture families. The protocol walks back three of the five methods on the architec-  
16 tures we audit, where the field-standard reporting walks back none. A residual-  
17 stream penalty intervention partially alleviates both modes, and four independent  
18 control experiments—a null calibration with fresh random feedback, a hypothesis-  
19 disambiguation sweep on early-epoch vanilla checkpoints, a matched BP+penalty  
20 capacity-cost control, and a perturbation-correlation cross-metric triangulation—  
21 validate the two-mode separation. We release the protocol, the audit data, and a  
22 reporting template.

## 23 1 Introduction

24 Feedback Alignment (FA) and its variants [1–4] are routinely evaluated on modern residual architec-  
25 tures by reporting two numbers: the trained network’s test accuracy, and the cosine similarity  $\Gamma$  be-  
26 tween the method’s local credit signal and the true backpropagation gradient at hidden layers. A high  
27  $\Gamma$  is interpreted as evidence that the method is computing useful credit; an above-shallow accuracy is  
28 interpreted as evidence that the deep blocks are being trained. On a 4-block pre-LayerNorm Resid-  
29 ualMLP at  $d=256$  trained on CIFAR-10 under standard hyperparameters, DFA reports  $\Gamma \approx 0.10$   
30 and a test accuracy of 31%, both of which look reasonable to a reviewer who encounters them in  
31 isolation.

32 **Both numbers can silently mislead.** On the same architecture and seeds, an architecture-matched  
33 random-untrained-blocks baseline trained only at the embedding, terminal LayerNorm, and head  
34 reaches 34.9% test accuracy: the trainable-blocks DFA variant under-performs this baseline by

35 4 percentage points. The deep blocks are not just unhelpful—they are actively destroying value.  
36 Meanwhile, the BP gradient at the deepest hidden layer of the same trained DFA network has  
37  $\|g_L\| \approx 5 \times 10^{-10}$ , well below F.cosine\_similarity’s default  $\varepsilon=10^{-8}$  clamp and well below  
38 any reasonable numerical floor. The reported  $\Gamma \approx 0.10$  is a cosine to a noise-floor reference vector  
39 and is mathematically well defined but uninterpretable as “alignment quality.”

40 **Why both numbers fail together turns out to have a single source: the headline-accuracy and**  
41 **headline- $\Gamma$  pair conflates two distinct phenomena that the field treats as one.** This paper identi-  
42 fies the two phenomena, names them, and provides a protocol that separates them.

43 **The two failure modes (informal). Mode 1: measurement degeneracy via terminal-**  
44 **LayerNorm gradient cancellation.** In modern pre-LayerNorm residual networks with a terminal  
45 LN before the classification head, DFA-style local losses have no global constraint on residual-  
46 branch magnitude. Block parameters grow by  $\sim 95\times$  relative to initialization, the residual stream  
47  $\|h_L\|$  grows from  $\sim 9$  at random init to  $\sim 4 \times 10^8$  over 100 epochs, and the LayerNorm Jacobian  
48 rescaling drives the BP gradient at hidden layers from  $\sim 10^{-3}$  to  $\sim 10^{-10}$ . The cosine alignment met-  
49 ric is then computed against a numerical-floor reference vector and cannot meaningfully distinguish  
50 a useful credit signal from noise.

51 **Mode 2: low intrinsic credit-direction quality of random feedback.** Even at the very first epoch  
52 of vanilla DFA training, when  $\|g_L\|$  is still in the meaningful regime ( $\sim 10^{-6}$ , three orders above the  
53 floor), DFA’s local credit signal  $e_T B_l^\top$  has essentially zero alignment with the BP gradient on deep  
54 layers ( $\overline{\cos} = -0.008 \pm 0.013$  across three seeds). The deep-layer alignment is missing for a reason  
55 that has nothing to do with measurement: random feedback simply does not compute a useful credit  
56 direction at the block layers of pre-LN residual networks, and this would be visible if the metric  
57 were interpretable.

58 **Why the field hasn’t seen this before.** The two modes are normally entangled: Mode 1 makes  
59 Mode 2 invisible, and the field-standard (accuracy,  $\Gamma$ ) pair has no diagnostic for either. A reviewer  
60 reading “DFA reaches 31%,  $\Gamma \approx 0.10$ ” has no signal that the deep blocks are passive (Mode 2) or  
61 that the cosine is measured against the floor (Mode 1). The framing has stayed in place because the  
62 symptoms look like ordinary undertraining.

63 **Our contribution.** We propose a **four-diagnostic protocol** that detects both modes, together with  
64 a calibrated scale for each diagnostic, a reference implementation, and a five-method audit on three  
65 architecture families (pre-LN ResidualMLP, ViT-Mini, BatchNorm CNN). The protocol walks back  
66 DFA, State Bridge, and Credit Bridge on the modern residual architectures we audit, where the  
67 field-standard (accuracy,  $\Gamma$ ) pair walks back none. We additionally validate that the two modes  
68 are mechanistically distinct: a residual-stream penalty intervention restores the BP gradient to the  
69 meaningful regime (alleviating Mode 1) and *partially* restores deep-layer alignment from 0 to  $\sim$   
70 0.16 (alleviating Mode 2), but neither is fully fixed. Cross-metric triangulation with perturbation  
71 correlation, null calibration with fresh random feedback, and a matched BP+penalty capacity-cost  
72 control all confirm the separation.

73 The protocol, reference implementation, audit table, and reporting template are released as a com-  
74 munity artifact. Our goal is that future FA evaluations on modern architectures use the protocol or  
75 an equivalent calibrated reporting standard, instead of the present field-standard pair that silently  
76 conflates measurement degeneracy with credit quality.

## 77 2 Related work

78 **Feedback Alignment and local credit.** Random feedback alignment [1] demonstrated that back-  
79 ward weights need not match forward weights for shallow networks to learn. Direct Feedback Align-  
80 ment (DFA) [2] bypassed the symmetric backward pass entirely. Subsequent work [5, 6, 3] extended  
81 FA to deeper networks with mixed success. [4, 7] showed DFA can train modest CNNs and small  
82 Transformers, typically reporting  $\Gamma$  as evidence that the local signal is useful. [8] questioned whether  
83 FA-style methods can scale to ImageNet-class problems. State and credit bridges [10, 11] are recent  
84 attempts to learn explicit credit-prediction networks under similar constraints.

85 **FA evaluation.** The standard evaluation pair—test accuracy and the cosine  $\Gamma$  between local credit  
 86 and the true BP gradient at hidden layers—has been used in essentially all of the above work. To  
 87 our knowledge, no prior work questions whether  $\Gamma$  is measured in a meaningful regime on the  
 88 architectures it is reported on, or whether the deep blocks of the trained network actually contribute  
 89 over an architecture-matched random-untrained-blocks baseline. We call this combined oversight  
 90 the field-standard evaluation pair, and our paper identifies how it conflates two distinct phenomena.

91 **Evaluation as scientific object.** The NeurIPS 2026 Evaluations and Datasets track explicitly in-  
 92 vites critical analyses of existing evaluation practices and proposals for new evaluation protocols.  
 93 Adjacent work in deep learning evaluation has documented similar conflation issues: e.g., the well-  
 94 known “representation similarity is metric-dependent” literature, the “probing task validity” critique,  
 95 the LayerNorm-induced gradient pathology in pre-LN Transformers [9]. Our contribution is to iden-  
 96 tify the analogous conflation in FA evaluation specifically and to provide a protocol that resolves it  
 97 for the FA evaluation community.

### 98 3 The audit: standard FA evaluation walks back nothing

99 We apply the field-standard (accuracy,  $\Gamma$ ) reporting pair to five methods on the standard 4-block  
 100  $d=256$  pre-LayerNorm ResidualMLP on CIFAR-10 (Table 1, three seeds, 100 training epochs,  
 101 AdamW lr= $10^{-3}$ , wd=0.01, cosine schedule).

Table 1: Field-standard reporting on five methods, 4-block  $d=256$  ResidualMLP, CIFAR-10, three seeds. The headline pair gives no walk-back signal on any method.

method	test acc	headline $\Gamma$	status quo verdict	our verdict
BP	$0.609 \pm 0.004$	$\approx 1.0$	trustworthy	trustworthy
EP	$0.316 \pm 0.038$	0.008	trustworthy	trustworthy
DFA	$0.308 \pm 0.014$	0.10	trustworthy	<b>walked back</b>
Credit Bridge	$0.289 \pm 0.034$	0.07	trustworthy	<b>walked back</b>
State Bridge	$0.205 \pm 0.039$	0.005	trustworthy	<b>walked back</b>

102 A reviewer reading Table 1’s middle two columns has no signal that any of these methods is in a  
 103 degenerate regime: every (accuracy,  $\Gamma$ ) pair looks consistent with “DFA-style methods train deep  
 104 residual networks to roughly one-third of BP’s accuracy with a small but positive credit alignment.”  
 105 The status quo verdict treats all five methods as trustworthy.

106 **The two diagnostics that should have fired.** The same trained networks have:

- 107 • **Per-block residual-stream growth** ( $\max_l \|h_{l+1}\|/\|h_l\|$ ) of 1.3 for BP, 2.4 for State Bridge,  
 108 11.6 for EP,  $96\times$  for Credit Bridge, and  $237\times$  for DFA. BP and EP are bounded; DFA, SB,  
 109 and CB show explosive per-block growth.
- 110 • **BP gradient at the deepest hidden layer** ( $\|g_L\|$ ) of  $\sim 4\times 10^{-4}$  for BP,  $\sim 2\times 10^{-4}$  for EP,  $\sim$   
 111  $10^{-9}$  for DFA, SB, and CB. The DFA/SB/CB values are below the `F.cosine_similarity`  
 112 default  $\varepsilon=10^{-8}$  clamp and several orders below any reasonable numerical floor for the  
 113 cosine metric to be interpretable.

114 Both diagnostics cleanly separate healthy methods from degenerate ones across three seeds: a sep-  
 115 aration gap of  $63\times$  for the per-block growth measure (healthy max 11, degenerate min 694) and  
 116  $24,338\times$  for the BP gradient floor measure (healthy min  $1.0\times 10^{-4}$ , degenerate max  $4.2\times 10^{-9}$ ).  
 117 Both gaps survive a sweep of the threshold value over an order of magnitude.

118 **The walked-back claim.** We report this finding as the primary audit result. Three of the five  
 119 methods we audit have claims that should be walked back, and the field-standard reporting pair does  
 120 not catch any of them.

121 **Walk-back: the deep blocks are not contributing.** Beyond the measurement-degeneracy diag-  
 122 nostics, an architecture-matched *frozen-random-blocks* baseline (training only the embedding, ter-  
 123 minal LN, and head while leaving the deep blocks at random initialization) reaches  $0.349 \pm 0.002$

124 on this architecture under all three of DFA, SB, and CB. The trainable-blocks variants reach 0.308,  
 125 0.205, and 0.289 respectively—*below* the random-untrained baseline. Training the deep blocks is  
 126 not just unhelpful; on this architecture and these seeds, it is actively destructive of accuracy.

127 **This is the central audit finding.** Three of five FA-style methods on a standard residual architecture  
 128 under standard hyperparameters do not beat their architecture’s frozen-random-blocks baseline. The  
 129 field-standard (accuracy,  $\Gamma$ ) reporting pair has no diagnostic for this.

## 130 4 The diagnostic protocol

131 We propose a four-diagnostic protocol that detects the audit findings of Section 3.

132 **Diagnostic (a): per-layer residual stream growth.** Compute  $\max_l \|h_{l+1}\|_2 / \|h_l\|_2$  over a fixed  
 133 evaluation batch. If the maximum per-block growth exceeds a calibrated threshold ( $50\times$  in our  
 134 default), the residual stream is in a regime incompatible with the original architectural intent. This  
 135 is the most direct test of Mode 1’s structural cause.

136 **Diagnostic (b): BP gradient at hidden layers.** Compute  $\|\partial L / \partial h_L\|_2$  on a fixed eval batch.  
 137 If this falls below a calibrated floor ( $10^{-7}$  in our default, well above fp32 subnormals and the  
 138 `F.cosine_similarity` clamp), the reference vector against which  $\Gamma$  is measured is at the nu-  
 139 merical floor and the metric is not interpretable as alignment quality. This is Mode 1’s symptom:  
 140 any cosine alignment reported in this regime is a cosine to noise.

141 **Diagnostic (c): cross-batch direction stability.** Compute the mean pairwise cosine of normalized  
 142 BP-grad direction across disjoint minibatches. A high value ( $> 0.30$  in our default) indicates the  
 143 reference vector is dominated by a sample-invariant global drift component, which means  $\Gamma$  mea-  
 144 sures alignment to drift rather than to per-sample credit. This is a sub-mode discriminator: it tells  
 145 you *how* Mode 1 has corrupted the reference, not whether (b) alone detects.

146 **Diagnostic (d): frozen-blocks baseline.** Train an architecture-matched variant with the deep  
 147 blocks frozen at random initialization. If the trainable-blocks variant fails to clear this baseline  
 148 by a calibrated margin (2 percentage points in our default), the deep blocks are not meaningfully  
 149 contributing. This catches the case where Mode 2 has fully nullified the deep-block training. Note  
 150 that this is a behavioral consequence and (as we discuss in Section 5) becomes ambiguous under  
 151 interventions that partially restore alignment.

152 **Calibrated thresholds.** Default thresholds ( $50\times$ ,  $10^{-7}$ , 0.30, 2pp) sit cleanly in the middle of  
 153 large separation gaps between healthy and degenerate networks: the per-block growth diagnostic  
 154 has a  $63\times$  gap, the BP gradient floor diagnostic has a  $24,338\times$  gap. Verdicts are robust to threshold  
 155 perturbations of a factor of two in either direction.

156 **Decision-utility ablation.** We compare seven reporting strategies on the five-method audit (Ta-  
 157 ble 2): the field-standard pair (S0: accuracy only, S1:  $+\Gamma$ ) walks back 0/5 methods. The full  
 158 protocol (S<sub>full</sub>: accuracy + (a) + (b) + (c) + (d)) walks back 3/5. Each of (a), (b), and (d) is inde-  
 159 pendently sufficient for binary detection of the three failing methods on this architecture; (c) is for  
 160 sub-mode discrimination, not primary detection.

Table 2: Decision-utility ablation. “Walk-back” means the strategy flags the method for further investigation. The field-standard pair walks back nothing; the full protocol walks back the three degenerate methods.

method	S0	S1	+(a)	+(b)	+(c)	+(d)	full
BP	—	—	—	—	—	—	trust
EP	—	—	—	—	—	—	trust
DFA	—	—	WB	WB	—	WB	WB
State Bridge	—	—	WB	WB	WB	WB	WB
Credit Bridge	—	—	WB	WB	WB	WB	WB

161 **Cross-architecture validation.** We replicated the protocol on per-epoch training-time data for  
 162 three architecture families: 4-block pre-LN ResidualMLP, 4-block ViT-Mini, and a synthetic Stu-  
 163 dentNet without terminal LayerNorm, plus a five-method audit on a SmallCNN with BatchNorm  
 164 and no terminal LN. Across the  $3 \text{ archs} \times 3 \text{ seeds} \times 2 \text{ methods} = 18$  training trajectories of the  
 165 first three, the diagnostics fire on every DFA training run on the with-terminal-LN architectures  
 166 within 1–11 epochs (well before the headline accuracy stabilizes), and never fire on any BP run.  
 167 On the without-terminal-LN architectures (StudentNet, CNN), diagnostic (a) still fires on DFA but  
 168 diagnostic (b) does *not* fire on any of the methods we tested. This is consistent with diagnostic (b)  
 169 being specifically about LayerNorm-driven gradient cancellation rather than residual-stream growth  
 170 in general.

171 **Reference implementation.** We release `protocol/`, a  $\sim 200$ -line Python module that implements  
 172 the protocol on any model exposing a duck-typed interface (`model(x, return_hidden=True)`,  
 173 `model.embed` or `model.patch_embed`, `model.blocks`, and a terminal LN + head). The package  
 174 includes a smoke test that loads BP/DFA/EP checkpoints and verifies expected verdicts, a reporting  
 175 template, and a reproducible audit table.

## 176 5 Two distinct failure modes

177 The protocol of Section 4 catches the audit finding, but its main scientific interest is what it reveals  
 178 about *why* the field-standard pair fails. We argue that the failure is not a single phenomenon: it  
 179 conflates two distinct modes that respond differently to interventions and whose mechanisms are  
 180 separately measurable.

181 **Mode 1 (measurement degeneracy via terminal-LN gradient cancellation), in detail.** On the  
 182 standard 4-block  $d=256$  pre-LN ResidualMLP, DFA’s local block losses  $\langle f_l(h_l), e_T B_l^\top \rangle$  have no  
 183 scale constraint: the inner product can be increased indefinitely by inflating  $\|f_l(h_l)\|$ . Block param-  
 184 eters  $w_1, w_2$  inside each block grow by a factor of  $\sim 200\times$  during 100 epochs of training, and the  
 185 multiplicative product  $\|w_1\| \cdot \|w_2\|$  grows by  $\sim 5 \times 10^4$  per block. The residual stream  $\|h_L\|$  grows  
 186 from 9 at initialization to  $\sim 4 \times 10^8$  by epoch 100, with most of the growth happening in the first  
 187 10 epochs. Through the terminal LayerNorm Jacobian ( $\partial \text{LN}(h)/\partial h \propto 1/\|h\|$ ), this drives the BP  
 188 gradient at hidden layers from  $\sim 10^{-3}$  at random initialization to  $\sim 5 \times 10^{-10}$ . The cosine alignment  
 189 metric is then computed against a reference vector at the numerical floor: `F.cosine_similarity`  
 190 clamps the divisor at  $\varepsilon=10^{-8}$  rather than dividing by the true magnitude, scaling the reported value  
 191 by a factor of  $\sim 50\times$  in the wrong direction; the reported  $\Gamma \approx 0.10$  is not a “small alignment” but a  
 192 cosine to a degenerate reference.

193 **Causal validation: penalty intervention partially restores Mode 1.** Adding  $\lambda \|f_l(h_l)\|^2$  as a per-  
 194 block penalty to DFA’s local loss with  $\lambda=10^{-2}$  contains the residual stream:  $\|h_L\| : 4 \times 10^8 \rightarrow 4 \times 10^4$   
 195 (4 OOM rescue), and  $\|g_L\| : 5 \times 10^{-10} \rightarrow \sim 10^{-6}$  (4 OOM rescue, well into the meaningful regime).  
 196 Diagnostics (a) and (b) both pass on the penalized network. Three seeds:  $\|h_L\| = 4.0 \pm 0.1 \times 10^4$ ,  
 197  $\|g_L\| = 9.0 \pm 0.9 \times 10^{-7}$ .

198 **Mode 2 (low intrinsic credit-direction quality), in detail.** The penalty restores Mode 1, but the  
 199 test accuracy of penalized DFA only rises from 0.308 to 0.363 (3-seed mean  $0.363 \pm 0.001$ ). This  
 200 is +5.5pp over vanilla DFA but only +1.4pp over the architecture-matched random-blocks baseline  
 201 of 0.349. The deep blocks are still not meaningfully contributing.

202 **Direct measurement.** On the penalized DFA checkpoint, we directly compute the per-layer cosine  
 203 of the local credit signal  $e_T B_l^\top$  with the BP gradient at  $h_l$ , using the training-time random feedback  
 204 matrices  $B_l$  and no  $\varepsilon$  clamp. Three-seed result on deep layers ( $l = 1, 2, 3, 4$ ):  $\overline{\text{cos}} = +0.155 \pm 0.025$ .  
 205 This is *measurable, real, and small*: well above noise (see calibration below) but well below BP’s  
 206 self-cosine of 1.0. The deep blocks under the penalty are partially aligned with BP gradient but not  
 207 fully.

208 **Disambiguation: was the alignment always there, or did the penalty create it?** A reasonable  
 209 reading of the above would be: “the cosine was always there in vanilla DFA; the penalty just made  
 210 the measurement interpretable.” The disambiguation experiment falsifies this. We trained vanilla

211 DFA and saved checkpoints at every epoch from 1 to 5, where  $\|g_L\|$  is still in the meaningful regime  
 212 ( $1.4 \times 10^{-6}$  at epoch 1, well above the  $10^{-7}$  floor). Per-layer cosine on these vanilla checkpoints  
 213 (3 seeds, epochs 1 and 2): *deep-layer cosine*  $-0.008 \pm 0.013$  averaged over 24 measurements  
 214 (3 seeds  $\times$  2 epochs  $\times$  4 deep layers). The deep-layer alignment is essentially zero on vanilla DFA in  
 215 the meaningful regime; the  $+0.155$  on the penalized network is created by the penalty intervention,  
 216 not revealed by it.

217 **The penalty’s role.** The penalty does two things at once. It contains the residual stream (directly  
 218 addressing Mode 1), and it changes the training trajectory of the block parameters such that the  
 219 final  $f_l$  direction is partially aligned with the BP gradient direction (partially addressing Mode 2).  
 220 The second effect is non-obvious: the penalty does not directly optimize for alignment. A plausible  
 221 mechanism is that with no penalty, the local credit objective can be increased indefinitely by inflating  
 222  $\|f_l\|$ , so the optimizer follows directions uncorrelated with BP gradient; with the penalty,  $\|f_l\|$  is  
 223 constrained, so the optimizer must orient  $f_l$  more carefully, which incidentally yields better partial  
 224 alignment with BP gradient direction.

## 225 5.1 Calibration of the cosine measurement

226 A natural reviewer concern about the  $+0.155$  result is whether it is above or below noise. We anchor  
 227 it with explicit positive and negative controls.

228 **Positive control.** On a BP-trained network, using the BP gradient itself as the predicted credit signal,  
 229 the perturbation correlation  $\rho$  between  $\langle g_l, \varepsilon v \rangle$  and the true loss change  $L(h_l + \varepsilon v) - L(h_l)$  is  $+0.997$   
 230 at every layer (4-layer mean  $+0.9965$ ). This is the Taylor-expansion ceiling.

231 **Negative control.** On the same BP-trained network, using a random vector independent of the layer  
 232 as the credit signal,  $\rho$  is  $+0.006$  (4-layer mean), within statistical noise of zero.

### 233 Cross-metric triangulation on the test conditions.

Table 3: Two metrics, four conditions. The agreement between cosine and perturbation correlation rules out single-metric artifacts.

condition	deep cosine $\overline{\cos}$	deep $\overline{\rho}$
positive control (BP grad on BP net)	1.000	+0.997
negative control (random vector on BP net)	—	+0.006
vanilla DFA, ep 1 (3 seeds, meaningful regime)	$-0.008 \pm 0.013$	$-0.003 \pm 0.005$
penalized DFA, ep 30 (3 seeds, lam= $10^{-2}$ )	$+0.155 \pm 0.025$	$+0.080 \pm 0.011$

234 The penalized DFA’s  $+0.080$  perturbation correlation is  $\sim 13 \times$  the negative control and  $\sim 8\%$  of  
 235 the positive control. Both metrics agree on the vanilla-to-penalized transition: vanilla deep signal is  
 236 indistinguishable from random, penalized deep signal is small but well above noise. The agreement  
 237 across metrics rules out the possibility that cosine is capturing a directional artifact unrelated to  
 238 local-loss usefulness.

## 239 5.2 $\lambda$ sweep: independent dissociation of the two modes

240 The disambiguation experiment of Section 5 relied on vanilla DFA early-epoch checkpoints (epochs  
 241 1–2) to measure deep-layer cosine while  $\|g_L\|$  was still in the meaningful regime. A natural reviewer  
 242 concern is that early-epoch checkpoints are not at convergence and might be confounded by stochas-  
 243 tic initialization effects. We strengthen the disambiguation with an independent control: a sweep  
 244 over the penalty strength  $\lambda$  at convergence (30 epochs), with both metrics measured on each saved  
 245 checkpoint.

246 **The killer row is  $\lambda=10^{-4}$ .** At this penalty strength, the residual stream is already contained  
 247 ( $\|h_L\| = 2.4 \times 10^4$ , four orders below vanilla), and the BP gradient at the deepest hidden layer  
 248 is at  $6.3 \times 10^{-7}$  (well above the  $10^{-7}$  floor and in the meaningful measurement regime). Diag-  
 249 nostics (a) and (b) both pass: **Mode 1 is fully alleviated.** But the deep-layer cosine ( $-0.022$ ) and  
 250 perturbation correlation ( $-0.004$ ) are essentially zero, on both metrics independently. **Mode 2 is**  
 251 **not alleviated at all.**

Table 4:  $\lambda$  sweep on the penalty strength, all 30 epochs, seed 42. The deep-layer cosine and perturbation correlation rise from essentially zero at  $\lambda=10^{-4}$  to small-but-positive at  $\lambda=10^{-2}$ , even though diagnostics (a) and (b) already pass at  $\lambda=10^{-4}$ .

$\lambda$	test acc	$\ h_L\ $	$\ g_L\ $	deep $\overline{\text{cos}}$	deep $\overline{\rho}$
0	0.308	$4.4 \times 10^8$	$5 \times 10^{-10}$	(degenerate)	(degenerate)
$10^{-4}$	0.359	$2.4 \times 10^4$	$6.3 \times 10^{-7}$	-0.022	-0.004
$10^{-2}$	0.363	$4.0 \times 10^4$	$9.0 \times 10^{-7}$	+0.165	+0.091
$10^{-1}$	0.349	$1.2 \times 10^4$	$1.6 \times 10^{-6}$	+0.131	+0.067

252 This is direct evidence that the two modes are mechanistically distinct: they do not even respond  
 253 to the same intervention strength. There exists a regime ( $\lambda=10^{-4}$ , 30 epochs of training) in which  
 254 Mode 1 is fully alleviated and Mode 2 is unchanged from vanilla, with both metrics agreeing.

255 The threshold for Mode 2 alleviation is somewhere between  $\lambda=10^{-4}$  and  $\lambda=10^{-2}$ . At  $\lambda=10^{-2}$  the  
 256 penalty is strong enough to alter the optimization trajectory of the block parameters (constraining  
 257  $\|f_i\|$  tightly enough that the direction of  $f_i$  has to be coordinated more carefully with the local  
 258 credit signal), and the deep-layer alignment rises to  $\sim +0.16$ . At  $\lambda=10^{-1}$  the penalty starts to over-  
 259 constrain and the alignment is slightly lower ( $\sim +0.13$ ), giving an inverted-U relationship between  
 260  $\lambda$  and deep alignment.

### 261 5.3 Capacity-cost control

262 A second reviewer concern is whether the 0.36  $\rightarrow$  0.61 accuracy gap between penalized DFA and  
 263 BP-trainable is due to credit quality (Mode 2) or simply to the penalty’s capacity-regularization cost.  
 264 We disambiguate with a  $2 \times 2$  matched control.

Table 5:  $2 \times 2$  capacity-cost control. The penalty is the same in both the BP and DFA conditions. BP+penalty still clears the random-blocks baseline by 18.1pp; DFA+penalty clears it by only 1.4pp.

	no penalty	with penalty
BP	0.609	0.530
DFA	0.308	0.363
$\Delta$	-8.0pp	+5.5pp

265 Two observations make this control informative. First, the penalty’s effect on BP is -8pp (a small  
 266 capacity loss), which is one order of magnitude smaller than the residual gap between BP+penalty  
 267 and DFA+penalty ( $0.530 - 0.363 = 17\text{pp}$ ). The 17pp residual gap is consistent with credit-quality  
 268 cost, not with capacity regularization. Second, the penalty has *opposite* effects on the two methods:  
 269 it hurts BP by 8pp while helping DFA by 5.5pp, the opposite pattern expected from a generally  
 270 beneficial regime shift.

271 **The clean phrasing.** The  $2 \times 2$  control identifies a residual performance gap under matched architec-  
 272 ture, data, optimizer family, and matched penalty, after accounting for the penalty’s direct capacity  
 273 cost on BP. It is not a perfect isolation of “credit quality” in a vacuum (BP uses end-to-end loss while  
 274 DFA uses local block losses, and the two trainers may differ in non-capacity ways), but it is a strong  
 275 lower bound on the non-capacity penalty-unexplained gap.

### 276 5.4 Summary: five validations of the two-mode separation

277 Together, the disambiguation experiment, the  $\lambda$  sweep, the cross-metric triangulation, the capacity-  
 278 cost control, and the threshold robustness analysis provide five independent lines of evidence that the  
 279 failure of standard FA evaluation is not a single phenomenon. Mode 1 (measurement degeneracy)  
 280 is detected by diagnostic (b), is causally controlled by the residual-stream penalty at any  $\lambda \geq 10^{-4}$ ,  
 281 and is specifically associated with terminal- LayerNorm architectures in our audits. Mode 2 (low  
 282 intrinsic credit quality) persists after Mode 1 is alleviated at weak penalty strengths ( $\lambda=10^{-4}$ ), is  
 283 detected by direct per-layer cosine in the meaningful regime, and rises only when the penalty is

284 strong enough to alter the optimization trajectory of the deep blocks ( $\lambda \geq 10^{-2}$ ). The fact that the  
285 two modes have different intervention thresholds is the strongest single piece of evidence that they  
286 are mechanistically distinct.

## 287 6 Limitations

288 Our audit covers a specific slice of the FA literature: pre-LayerNorm ResidualMLP, ViT-Mini, and  
289 SmallCNN architectures on CIFAR-10, evaluated under standard hyperparameters. We do not claim  
290 that FA evaluation is broken everywhere; we identify a specific evaluation failure mode on mod-  
291 ern pre-LN residual networks with terminal LayerNorm, and we explicitly observe that diagnostic  
292 (b) does not fire on architectures without a terminal LN (StudentNet, CNN with BN). This is ob-  
293 servational association, not a causal identification of LayerNorm per se: a future non-terminal-LN  
294 architecture where (b) fires would refine the claim. Section 2 cites the classical FA literature where  
295 non-terminal-LN architectures dominate; our central claim concerns the modern with-terminal-LN  
296 residual case.

297 The Mode 2 measurement in Section 5 relies on direct cosine and perturbation correlation in the  
298 meaningful regime, which is only accessible after a Mode 1 intervention. We cannot directly ob-  
299 serve Mode 2 on a vanilla DFA-trained network at convergence, because by then  $\|g_L\|$  has crashed  
300 below the floor. The disambiguation experiment (early-epoch vanilla checkpoints) addresses this by  
301 measuring at epochs where  $\|g_L\|$  is still meaningful, but those checkpoints are not at convergence.

302 The matched-penalty  $2 \times 2$  control disambiguates capacity loss from credit quality but does not ac-  
303 count for non-capacity differences between end-to-end BP and local DFA training. The 17pp resid-  
304 ual gap is therefore a lower bound on the credit-quality cost rather than a clean isolation.

## 305 7 Broader impacts

306 This paper does not introduce a new training method, dataset, or generative model. It identifies a  
307 measurement problem in the evaluation of an existing class of training methods. Its primary impact  
308 is on the scientific record of the FA literature: future evaluations on modern residual architectures  
309 should use the protocol or an equivalent calibrated reporting standard, and existing claims about FA  
310 performance on these architectures should be re-evaluated under the protocol where possible. We  
311 are not aware of any negative downstream applications of this work.

## 312 8 Conclusion

313 We have shown that standard Feedback Alignment evaluation on modern residual networks is un-  
314 reliable because it conflates two distinct failure modes: measurement degeneracy via terminal-  
315 LayerNorm gradient cancellation, and low intrinsic credit-direction quality of random feedback.  
316 We provide a four-diagnostic protocol that detects both modes, a calibrated scale anchored by posi-  
317 tive and negative controls, a five-method audit on three architecture families, and four independent  
318 control experiments validating the two-mode separation. The protocol, audit data, and reporting  
319 template are released as a community artifact for the FA evaluation community.

## 320 References

- 321 [1] T. P. Lillicrap, D. Cownden, D. B. Tweed, and C. J. Akerman. Random synaptic feedback  
322 weights support error backpropagation for deep learning. *Nature Communications*, 7:13276,  
323 2016.
- 324 [2] A. Nøkland. Direct feedback alignment provides learning in deep neural networks. In *NeurIPS*,  
325 2016.
- 326 [3] M. Akrouf, C. Wilson, P. Humphreys, T. Lillicrap, and D. B. Tweed. Deep learning without  
327 weight transport. In *NeurIPS*, 2019.
- 328 [4] J. Launay, I. Poli, F. Boniface, and F. Krzakala. Direct feedback alignment scales to modern  
329 deep learning tasks and architectures. In *NeurIPS*, 2020.

- 330 [5] T. H. Moskowitz, A. Litwin-Kumar, and L. F. Abbott. Feedback alignment in deep convolu-  
331 tional networks. *arXiv:1812.06488*, 2018.
- 332 [6] M. Refinetti, S. d’Ascoli, R. Ohana, and S. Goldt. Align, then memorise: the dynamics of  
333 learning with feedback alignment. In *ICML*, 2021.
- 334 [7] B. Crafton, A. Parihar, E. Gebhardt, and A. Raychowdhury. Direct feedback alignment with  
335 sparse connections for local learning. *Frontiers in Neuroscience*, 13:525, 2019.
- 336 [8] S. Bartunov, A. Santoro, B. Richards, L. Marris, G. Hinton, and T. Lillicrap. Assessing the  
337 scalability of biologically-motivated deep learning algorithms and architectures. In *NeurIPS*,  
338 2018.
- 339 [9] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T. Liu.  
340 On layer normalization in the transformer architecture. In *ICML*, 2020.
- 341 [10] Anonymous. State Bridge: terminal-conditioned predictor for credit assignment. *Anonymous*  
342 *in-progress reference, 2024-2026*.
- 343 [11] Anonymous. Credit Bridge: value-field local credit without hidden BP. *Anonymous in-progress*  
344 *reference, 2024-2026*.

## 345 A Reproducibility

346 All experiments use PyTorch  $\geq 2.0$  on a single NVIDIA A6000 GPU. Source for the proto-  
347 col package is in `protocol/`; experimental scripts are in `experiments/`. Random seeds are  
348 42, 123, 456 for all 3-seed measurements, with additional seeds (789, 1024, 2048) used where  
349 reported. CIFAR-10 is loaded via `torchvision` with the standard normalization  $(\mu, \sigma) =$   
350  $((0.4914, 0.4822, 0.4465), (0.2470, 0.2435, 0.2616))$ .

## 351 B Pipeline pitfalls catalog

352 Beyond the four diagnostics, we found seven evaluation-pipeline bugs in our own dogfood  
353 codebase that silently corrupt FA evaluation results. Each has a standalone reproducer in  
354 `protocol/examples/verify_pitfalls*.py`.

- 355 1. `tensor.norm(-1)` is the  $L_{-1}$  “norm” of the entire flattened tensor, not the per-row  $L_2$   
356 norm. The correct call is `tensor.norm(dim=-1)`. This bug invalidated several months of  
357 our gradient-norm measurements.
- 358 2. `F.cosine_similarity(a, b)` divides by  $\max(\|a\| \|b\|, \epsilon)$  with  $\epsilon=10^{-8}$  by default.  
359 When  $\|b\| \sim 10^{-10}$  (the regime of the BP gradient on degenerate DFA-trained pre-LN  
360 networks), the divisor becomes  $\|a\| \cdot 10^{-8}$  instead of  $\|a\| \cdot 10^{-10}$ , scaling the reported  
361 cosine by a factor of  $\sim 100\times$  in the wrong direction.
- 362 3. fp16 mixed precision underflows BP gradients at  $\sim 5 \times 10^{-10}$ , below fp16’s smallest sub-  
363 normal of  $\sim 6 \times 10^{-8}$ . bf16 works because it shares fp32’s exponent range.
- 364 4. Random feedback  $B_l$  matrices are training-specific. DFA reports  $\Gamma \approx 0.106$  with the  
365 training-time  $B_l$ ; with 20 fresh random  $B_l$  draws on the same trained network,  $\Gamma \approx 0 \pm$   
366  $0.005$ . The reported alignment is the network adapting to its specific  $B_l$ , not intrinsic.
- 367 5. Aggregation strategy across (layers, samples, batches) is rarely specified but determines  
368 the headline number. Same DFA seed-42 gives  $\Gamma \in [-0.028, +0.074]$  across four valid  
369 aggregation strategies (a  $3.45\times$  ratio, with sign flip).
- 370 6. Per-layer  $\Gamma$  structure is hidden by aggregation. On the 4-block ResMLP, DFA’s headline  
371  $\Gamma \approx 0.10$  is driven almost entirely by the embedding layer ( $\Gamma_{l_0} \approx +0.43$ ); deeper layers  
372 have  $\Gamma \approx 0$ . The pattern is architecture- specific: on ViT-Mini all layers are uniformly near  
373 zero.

374 7. Auxiliary networks (random feedback  $B_l$ , bridge predictors) not saved alongside model  
375 checkpoints can cause post-hoc  $\Gamma$  scripts to silently fall back to  $\cos(\text{BP\_grad}, \text{BP\_grad}) =$   
376  $1.0$  and report “perfect alignment.” We discovered this in our own pipeline during the  
377 protocol development. Check that auxiliary networks are persisted before reporting any  $\Gamma$   
378 value.

## 379 **C Methodology: walk-back chain**

380 The framing of this paper underwent several corrections during the development of the protocol. We  
381 document the four-step progression explicitly as part of the methodology, not as narrative drama:

- 382 1. Initial metric ( $\Gamma \approx 0.10$  for DFA) suggested the method was learning useful credit on  
383 modern residuals.
- 384 2. Diagnostic showed the metric was measured against a numerical-floor reference vector  
385 ( $\|g_L\| \sim 10^{-10}$ ); the headline number was not interpretable.
- 386 3. Revised control (the residual-stream penalty) restored the reference but only partially  
387 closed the accuracy gap to BP, identifying a residual phenomenon.
- 388 4. Final interpretation (this paper) separates measurement failure (Mode 1) from genuine  
389 credit-quality cost (Mode 2), validated by the four control experiments of Section 5.

## 390 **D Six independent validations of the two-mode separation**

391 For completeness we list all six independent validation experiments, beyond the four reported in the  
392 main text:

- 393 1. Direct deep-layer cosine on penalized DFA (3 seeds): deep mean  $+0.155 \pm 0.025$ .
- 394 2. Null calibration with 20 fresh random  $B_l$ : deep cosine  $+0.002 \pm 0.022$  (within noise).
- 395 3. Hypothesis-disambiguation sweep: vanilla DFA early-epoch deep cosine  $-0.008 \pm 0.013$   
396 across 3 seeds at epoch 1.
- 397 4. BP+penalty matched-control: 8pp BP capacity cost vs 17pp residual gap at  $\lambda=10^{-2}$ .
- 398 5. Multi-seed lock-in: 24 measurements (3 seeds  $\times$  2 epochs  $\times$  4 deep layers) all in  
399  $[-0.04, +0.02]$  on vanilla.
- 400 6. Cross-metric triangulation via perturbation correlation: vanilla  $+0.002$ , penalized  $+0.080$   
401 (3 seeds), positive control (BP grad)  $+0.997$ , negative control (random vector)  $+0.006$ .