ENABLING FINE-TUNING OF DIRECT FEEDBACK ALIGNMENT VIA FEEDBACK-WEIGHT MATCHING

Anonymous authors

004

006 007 008

009 010

011

012

013

014

015

016

017

018

019

021

022

024

025

026

027

028 029

030 031 Paper under double-blind review

ABSTRACT

In this paper, we introduce feedback-weight matching, a new method that facilitates reliable fine-tuning of fully connected neural networks using Direct Feedback Alignment (DFA). Although DFA has demonstrated potential by enabling efficient and parallel updates of weight parameters through direct propagation of the network's output error, its usage has been primarily restricted to training networks from scratch. We provide the first analysis showing that existing standard DFA struggles to fine-tune networks pre-trained via back-propagation. Through an analysis of weight alignment (WA) and gradient alignment (GA), we show that the proposed feedback-weight matching enhances DFA's ability and stability in fine-tuning pre-trained networks, providing insights into DFA's behavior and characteristics when applied to fine-tuning. In addition, we find that feedback-weight matching, when combined with weight decay, not only mitigates over-fitting but also further reduces the network output error, leading to improved learning performance during DFA-based fine-tuning. Our experimental results show that, for the first time, feedback-weight matching enables reliable and superior fine-tuning across various fine-tuning tasks compared to existing standard DFA, e.g., achieving 7.97% accuracy improvement on image classification tasks (i.e., 82.67% vs. 74.70%) and 0.66 higher correlation score on NLP tasks (i.e., 0.76 vs. 0.10). The code implementation is available at an anonymous GitHub repository¹.

1 INTRODUCTION

Recently, a new training mechanism called Direct Feedback Alignment (DFA) (Nøkland, 2016) has been proposed for deep neural networks to alleviate the weight transport problem (Grossberg, 1987; Crick, 1989). Based on the concept of Feedback Alignment (FA) (Lillicrap et al., 2016), DFA passes the error of the output layer directly to each layer of the network to update the weight parameters without back-propagation (Rumelhart et al., 1986). By using random feedback matrices, the weight gradient of each layer is independently approximated from the directly passed error, enabling efficient training of networks through the parallel update of multiple layers. This contrasts with back-propagation that propagates the network error sequentially from the last to the first layer.

Although Direct Feedback Alignment (DFA) (Nøkland, 2016) has shown its potential in training 040 primarily for fully connected networks (Garg & Vempala, 2022; Launay et al., 2020), its application 041 to fine-tuning (Devlin et al., 2018), i.e., adapting a pre-trained network to a new task, has been less 042 studied until today despite its practical usefulness. In fact, it has been known that fine-tuning net-043 works with DFA is challenging (Chu & Bacho, 2024); the performance of networks fine-tuned with 044 DFA is generally unreliable compared to that of those fine-tuned with back-propagation (Rumelhart 045 et al., 1986). Given that fine-tuning has become one of the practical and also effective ways of re-046 utilizing pre-trained networks for various downstream tasks (Church et al., 2021), investigating how 047 DFA can be applied to the fine-tuning mechanism both theoretically and empirically is necessary.

Enabling fine-tuning with Direct Feedback Alignment (DFA) (Nøkland, 2016) can not only broaden DFA's usability but also introduce an alternative approach to current back-propagation-based fine-tuning (Rumelhart et al., 1986; Church et al., 2021). Currently, DFA has not yet been established as a reliable stand-alone training method that can provide comparable performance to back-propagation (Launay et al., 2019; Crafton et al., 2019). Thus, taking a wide range of well-pre-trained

⁰⁵³

¹ https://anonymous.4open.science/r/Feedback-Weight-Matching-C7F0

models, such as Transformer-based foundation models (Kenton & Toutanova, 2019), as the starting point would be a practical strategy that can complement DFA's unstable and limited learning capabilities. Additionally, by incorporating DFA's unique advantages, such as being back-propagation-free and enabling parallel training, into the widely used fine-tuning scheme, we can explore new possibilities for re-utilizing pre-trained models in a more agile, efficient, and biologically plausible manner, in contrast to conventional back-propagation, which requires significantly more resources and time.

060 In this paper, we introduce a DFA-based fine-tuning method, which investigates the feasibility of 061 Direct Feedback Alignment (DFA) (Nøkland, 2016) for fine-tuning deep neural networks, with the 062 aim of extending the scope of DFA to embrace various pre-trained networks. We first analyze 063 the reasons why the existing standard DFA, which updates the pre-trained weights using random 064 feedback matrices, does not perform well in fine-tuning. This analysis is based on the weight alignment (WA) and gradient alignment (GA) (Refinetti et al., 2021), which are two measures proposed to 065 estimate the state and learning performance of DFA. From this analysis, we propose the feedback-066 weight matching, which first reconstructs the feedback matrices by decomposing the pre-trained 067 weights and then re-initializes the weights based on the reconstructed feedback matrices before 068 starting fine-tuning. Additionally, we prove that applying weight decay (Krogh & Hertz, 1991) on 069 top of feedback-weight matching considerably improves and stabilizes the fine-tuning performance of DFA, beyond the general regularization effect on weight parameters. Together with the simple 071 yet effective feedback-weight matching, weight decay acts as a key facilitator for fine-tuning fully 072 connected networks with DFA. To the best of our knowledge, this work is the first attempt to explore 073 the possibility of applying DFA to fine-tuning of fully connected networks via an in-depth study.

074 The experiments provide evaluation results consistent with our theoretical analysis; applying 075 feedback-weight matching enables more effective and reliable fine-tuning of fully connected 076 networks with DFA over various fine-tuning tasks, when compared to the existing standard 077 DFA (Nøkland, 2016) that does not apply the proposed feedback-weight matching. For instance, the 078 image classification accuracy of fully connected networks fine-tuned with feedback-weight match-079 ing reaches 82.67%, while that of standard DFA remains 74.70%. Also, it successfully fine-tunes 080 Transformer models (BERT) (Devlin et al., 2018) on NLP tasks, e.g., achieving 0.76 correlation 081 score, while the standard DFA barely conducts fine-tuning at all, i.e., achieving mere 0.10 correlation score. The results demonstrate the potential for extending DFA towards the widely used pre-training and fine-tuning strategy, moving beyond its limited usage in from-scratch training. 083

084 085

093 094

102 103

2 BACKGROUND AND RELATED WORK

DFA. It is common to train a neural network using the back-propagation algorithm (Rumelhart et al., 1986). Given a fully connected network, we denote W_l as the weight of *l*-th layer of the network, $\mathcal{L}(\hat{y}, y)$ as the loss function, where \hat{y} is the ground-truth output, and y is the network output, and $h_l = g(a_l)$ as the output of the *l*-th layer, where $g(\cdot)$ is activation function, and $a_l = W_l h_{l-1}$. To update the weight with the gradient descent algorithm (Ruder, 2016), the gradient of the loss \mathcal{L} w.r.t. the weight W_l is obtained using back-propagation (BP) as:

$$\delta \boldsymbol{W}_{l}^{BP} = -\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}_{l}} = -\left[\left(\boldsymbol{W}_{l+1}^{\top} \delta \boldsymbol{a}_{l+1}\right) \odot g'(\boldsymbol{a}_{l})\right] \boldsymbol{h}_{l-1}^{\top}, \ \delta \boldsymbol{a}_{l} = \frac{\partial \mathcal{L}}{\partial \boldsymbol{a}_{l}}$$
(1)

where \odot is the Hadamard product. However, back-propagation poses some challenges, specifically the weight transport (Grossberg, 1987; Crick, 1989) and backward locking problems (Lillicrap et al., 2020; Launay et al., 2019). Direct Feedback Alignment (DFA) (Nøkland, 2016) addresses the weight transport problem by employing random feedback and mitigates the backward locking problem by delivering the network's output error signal to each layer independently. Specifically, 1) the global error vector $e = \hat{y} - y$ is transmitted to each layer, and 2) the weight W_{l+1} at the *l*-th layer of the network is replaced with a random feedback matrix F_l , leading to the following weight gradient:

$$\delta \boldsymbol{W}_{l}^{DFA} = -\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}_{l}} = -\left[\left(\boldsymbol{F}_{l}\boldsymbol{e}\right) \odot \boldsymbol{g}'(\boldsymbol{a}_{l})\right] \boldsymbol{h}_{l-1}^{\top} - \lambda^{t} \boldsymbol{W}_{l}$$
(2)

where λ^t is the weight-decay hyperparameter at the step *t*. Equation (2) eliminates the necessity of sequential layer-wise gradient computations required by back-propagation (Rumelhart et al., 1986).

GA and WA. To better elucidate the dynamics of DFA (Nøkland, 2016), the concept of gradient alignment (GA) is introduced (Lillicrap et al., 2016). GA quantitatively assesses the similarity be-

108 tween the weight gradients obtained through DFA and those derived via back-propagation (Rumel-109 hart et al., 1986). This is achieved by comparing the weight updates generated from the identically 110 initialized weights by both methods. It has been hypothesized that a stronger (higher) GA corre-111 sponds to enhanced learning performance in DFA. In addition, the concept of weight alignment 112 (WA) (Refinetti et al., 2021) has been introduced to evaluate the relationship between the weight and the feedback matrix in DFA, suggesting that strong WA is associated with strong GA. Although 113 GA and WA have been instrumental in analyzing the learning efficacy of DFA, prior research has 114 not explored their utility in the context of fine-tuning. In contrast, this paper pioneers the application 115 of GA and WA concepts to systematically investigate the fine-tuning process in DFA. 116

117 Applicability to Transformers and CNNs. Some studies (Launay et al., 2020) explore the applicability of DFA (Nøkland, 2016) to various fully connected networks, including neural radiance 118 fields (NeRF) (Mildenhall et al., 2021; Sitzmann et al., 2019), recommender systems (Guo et al., 119 2017), and NLP (Vaswani, 2017; Merity et al., 2016). While they show that DFA can train a wide 120 range of deep architectures, they also reveal a significant performance gap between DFA and back-121 propagation (Rumelhart et al., 1986), particularly in Transformer models (Vaswani, 2017). When 122 applied to models not based on fully connected networks, such as CNNs, the performance gap be-123 tween DFA and back-propagation is even more pronounced. For instance, VGG-16 (Simonyan & 124 Zisserman, 2014) on CIFAR-100 (Krizhevsky et al., 2009) trained with DFA achieves 1% top-1 ac-125 curacy (Launay et al., 2019), while back-propagation achieves 60%. Similarly, in ImageNet (Deng 126 et al., 2009), it is 6.2% vs. 53% (Crafton et al., 2019). Given that applying DFA to from-scratch 127 training scenarios 1) consistently underperforms relative to back-propagation, 2) takes a much longer 128 training time than fine-tuning, and 3) is limited to a narrower range of architectures, we argue that 129 utilizing DFA for fine-tuning would be a more effective, efficient, practical, and expedient approach. Thus, in this study, we investigate and analyze the potential of employing DFA in fine-tuning, which 130 is conducive to the widely-used pre-train-and-fine-tune strategy (Devlin et al., 2018). 131

132 Applying DFA to Back-Propagation Weights. As described above, in CNNs, DFA encounters 133 challenges in effectively learning the necessary spatial information (Crafton et al., 2019). Similarly, 134 in fully connected networks, DFA is known to produce feature representation clusters that deviate 135 from those learned via back-propagation (Nøkland, 2016). Moreover, although stable training can be achieved when transitioning from weights learned through DFA to back-propagation, the reverse 136 is not true; switching from back-propagation to DFA results in unstable training, and DFA fails to 137 fully recover its performance even after large training epochs (Chu & Bacho, 2024). These imply 138 the inherent difficulties in fine-tuning with DFA using weights pre-trained with back-propagation. 139

140 DFA with Weight Decay. In the study by Song et al. (2021), it is analyzed that weight decay (Krogh 141 & Hertz, 1991) can reduce the output error in fully connected networks when used with Feedback 142 Alignment (FA) (Lillicrap et al., 2016). Nevertheless, the analysis predominantly focuses on the training of networks from scratch using FA, rather than on the fine-tuning process with DFA. This 143 work, for the first time, examines the impact of weight decay in the context of fine-tuning with DFA. 144 Our findings indicate that weight decay can be beneficial in fine-tuning with DFA, as it reduces 145 network output error and over-fitting, thereby enhancing overall learning performance. 146

147 148

149

3 FEEDBACK-WEIGHT MATCHING

We first discuss why the existing standard DFA (Nøkland, 2016) does not behave stably in fine-150 tuning, based on weight alignment (WA) and gradient alignment (GA) (Refinetti et al., 2021). Then, 151 we introduce feedback-weight matching, which enables effective and reliable fine-tuning of DFA. 152

153 WHY DOES DFA PERFORM UNRELIABLY IN FINE-TUNING? 3.1 154

155 **Definition 3.1.** (Weak Weight Alignment) Given a L-layer fully connected linear network updated 156 (trained) with DFA (Nøkland, 2016), the weight of the *l*-th layer at the *t*-th training step, which is 157 denoted as $W_{1 \le 1 \le L}^t$, becomes (Refinetti et al., 2021) as follows:

158 159

$$W_{1}^{t} = F_{1}A_{1}^{t}, W_{1 < l < L}^{t} = F_{l}A_{l}^{t}F_{l-1}^{\top}, \text{ and } W_{L}^{t} = A_{L}^{t}F_{L-1}^{\top},$$

where $A_{1}^{t} = -\eta \sum_{t'=0}^{t-1} e^{t'}(x^{t'})^{\top}, \text{ and } A_{l \ge 2}^{t} = \eta^{2} \sum_{t'=0}^{t-1} \sum_{t''=0}^{t'-1} (B_{l}^{t'}x^{t''}) \cdot (B_{l}^{t''}x^{t''}) e^{t'}(e^{t''})^{\top}^{(3)}$

162 Here, F_l is the feedback matrix of the *l*-th layer, A_1^t and $A_{l>2}^t$ are the alignment matrices, and 163 $B_l = A_{l-2} \cdots A_0 \in \mathbb{R}^{n_L \times n_L}$ is defined recursively using the feedback matrices only, with $A_0 = I$ 164 (Refinetti et al., 2021). Equation (3) is referred to as weak weight alignment (WA) (Refinetti et al., 165 2021), representing the state where no particular relationship exists between $W_{1 < l < L}^{t}$ and $F_{l}F_{l-1}^{t}$ 166 and between W_L^t and F_{L-1}^{\top} . At the early stage of DFA training, weak WA is naturally induced since 167 $A_{l>2}^t$ in Equation (3) starts with arbitrary values. However, as the training proceeds, $A_{l\geq2}^t$ becomes 168 proportional to the identity matrix (Refinetti et al., 2021), i.e., $A_{l\geq 2}^t \propto I$, leading to another state 169 called strong weight alignment (WA), which is defined as follows. 170

Definition 3.2. (*Strong Weight Alignment*) If $A_{l\geq 2}^t \propto I$, Equation (3) becomes the state called *strong weight alignment (WA)*, which is defined as follows.

173 174

$$\boldsymbol{W}_{1

$$\tag{4}$$$$

It is known that the strong WA in Equation (4), given $F_l^{\top}F_l \equiv I$, implies *strong gradient alignment (GA)* (Refinetti et al., 2021) defined in Equation (9), causing the gradient direction of the DFA weight, $W_{1 < l \leq L}^t$, aligned to that of back-propagation (Rumelhart et al., 1986). Hence, strong WA leads the learning trajectory of DFA to be comparable to that of back-propagation with strong GA.

However, if the pre-trained weights are fine-tuned via existing standard DFA using arbitrary random feedback matrix F_l , it becomes difficult to achieve strong WA in Equation (4), as shown below, likely to result in sub-optimal fine-tuning performance by inducing weak GA from weak WA.

Proposition 3.3. If the pre-trained weight, denoted as W_l^0 , is updated using DFA with arbitrary random feedback matrices F_l , the strong WA condition in Equation (4) is unlikely to be satisfied.

183 184 185

186

187

188 189

190 191

182

$$W_{1 < l < L}^{t} \propto F_{l} F_{l-1}^{\top}, \ W_{L}^{t} \propto F_{L-1}^{\top}$$

$$\tag{5}$$

where W_l^t denotes the weight after t steps of training, starting from the pre-trained weight W_l^0 . Equation (5) shows that the weight trained from the backpropagation pre-trained weight does not satisfy the strong WA condition. The proof is detailed in Appendix A.

3.2 INDUCING STRONG WEIGHT ALIGNMENT

To enable fine-tuning with DFA by deriving strong GA from strong WA defined in Equation (4), we propose the *feedback-weight matching* method, which induces both strong WA and GA as follows.

Definition 3.4. (*Feedback Matching*) From the pre-trained weight W_l^0 , we set the feedback matrix \bar{F}_l such that:

195 196 197

207

212

213

194

$$\bar{F}_l \bar{F}_{l-1}^\top \equiv W_{1 < l < L}^0 \text{ and } \bar{F}_{L-1}^\top \equiv W_L^0.$$
(6)

Equation (6) requires us to decompose the pre-trained weight $W_{1 < l < L}^0$ into \bar{F}_l and \bar{F}_{l-1}^\top . It can be achieved either through traditional decomposition methods, such as SVD (Singular Value Decomposition) (Klema & Laub, 1980), or alternatively, by optimizing Equation (23) in Appendix B.

Once the feedback matrix \bar{F}_l is reconstructed as Equation (6), we proceed to the *weight matching* process to induce strong WA, as described below.

204 **Definition 3.5.** (*Weight Matching*) Given the reconstructed \bar{F}_l derived by feedback matching (Equation (6)), we re-initialize the pre-trained weight W_l^0 into \bar{W}_l^0 so that it matches \bar{F}_l such that:

$$\bar{W}_{1

$$\tag{7}$$$$

The following shows that Equation (6) and (7) lead to strong WA condition in Equation (4).

Proposition 3.6. If the re-initialized weight \bar{W}_l^0 in Equation (7) is updated using DFA with the feedback matrix \bar{F}_l derived by Equation (6), the strong WA condition in Equation (4) is induced.

$$\bar{\boldsymbol{W}}_{1 < l < L}^{t} \propto \bar{\boldsymbol{F}}_{l} \bar{\boldsymbol{F}}_{l-1}^{\top}, \ \bar{\boldsymbol{W}}_{L}^{t} \propto \bar{\boldsymbol{F}}_{L-1}^{\top}$$

$$\tag{8}$$

with \bar{W}_l^t is the weight at step t, initialized from \bar{W}_l^0 . Equation (8) indicates that the weight updated from the re-initialized weight satisfies the strong WA condition, the proof is detailed in Appendix A. Subsequently, strong WA, achieved through Equation (6) and Equation (7), leads to strong GA (Refinetti et al., 2021). By matching the feedback matrix to the pre-trained weights, as in Equation (6), it becomes possible to preserve the knowledge embedded in the pre-trained weights. Additionally, by re-initializing the pre-trained weights from the matched feedback matrices, as in Equation (7), it becomes possible to facilitate the attainment of strong WA through DFA in fine-tuning.

221 222

223

230

234

235

236

245 246

247

252

253

254 255 256

257

258 259

260

261

268

269

3.3 INDUCING STRONG GRADIENT ALIGNMENT

While the previous section (Section 3.2) shows that the proposed feedback-weight matching in Equation (6) and (7) promotes strong weight alignment (WA), naturally leading to strong gradient alignment (GA), we now show that feedback-weight matching also directly induces strong GA.

Definition 3.7. (*Gradient Alignment*) The gradient alignment (GA) is defined as the cosine similarity between the weight gradient obtained using DFA (Nøkland, 2016), denoted G_{DFA} , and the weight gradient of back-propagation (Rumelhart et al., 1986), denoted G_{BP} , which is given by:

$$\cos \angle (\boldsymbol{G}_{DFA}, \boldsymbol{G}_{BP}) = \boldsymbol{G}_{DFA} \cdot \boldsymbol{G}_{BP} / \|\boldsymbol{G}_{DFA}\| \|\boldsymbol{G}_{BP}\|.$$
(9)

We show that feedback-weight matching, i.e., Equation (6) and (7), also directly induce strong GA when fine-tuning the first layer of the two-layer fully connected linear network, as follows.

Proposition 3.8. Feedback-weight matching given in Equation (6) and (7) induces strong GA, i.e., a higher GA, in the first layer of a fully connected linear network.

 $\cos_{FWM} \angle (\boldsymbol{F}_1, \boldsymbol{W}_2^t) \ge \cos_{DFA} \angle (\boldsymbol{F}_1, \boldsymbol{W}_2^t)$ (10)

cos_{*FWM*} \angle (F_1 , W_2^t) refers to GA in the first layer using feedback-weight matching, while cos_{*DFA*} \angle (F_1 , W_2^t) refers to GA in the first layer without feedback-weight matching. Equation (10) shows the GA when feedback-weight matching is used and when it is not. The proof is detailed in Appendix A. Based on Proposition 3.8, which shows feedback-weight matching induces stronger GA, we provide the following conjecture, which generalizes it to an arbitrary *L*-layer fully connected linear network.

Conjecture 3.9. It is conjectured that Equation (6) and (7) induce strong gradient alignment (GA), i.e., a higher GA, for all $1 \le l \le L$ layers in a fully connected linear network, where L > 2.

4 WEIGHT DECAY

Similar to conventional training using back-propagation (Nøkland, 2016), weight decay (Krogh & Hertz, 1991) has been shown to mitigate over-fitting of DFA, though its effect in fine-tuning has not been studied. We discuss how the proposed feedback-weight matching helps weight decay to reduce the network error (i.e., improving learning performance) during fine-tuning when applied to DFA.

Lemma 4.1. Given the re-initialized weight $\overline{W}_{1 < l \leq L}^0$ in Equation (7) and the pre-trained weight $W_{1 < l < L}^0$, the following terms, $r_{1 < l < L}$ and r_L , are non-negative with high probability.

$$r_{1 < l < L} = \|\boldsymbol{W}_{l}^{t} - \boldsymbol{W}_{l}^{0}\| - \|\boldsymbol{W}_{l}^{t} - \bar{\boldsymbol{W}}_{l}^{0}\| = \|\bar{\boldsymbol{F}}_{l}\bar{\boldsymbol{F}}_{l-1}^{\top} - \boldsymbol{W}_{l}^{0}\| - |\boldsymbol{c}_{l}^{t} - 1|\|\bar{\boldsymbol{F}}_{l}\bar{\boldsymbol{F}}_{l-1}^{\top}\| \ge 0 \quad (11)$$

$$r_{L} = \|\boldsymbol{W}_{L}^{t} - \boldsymbol{W}_{L}^{0}\| - |\boldsymbol{W}_{L}^{t} - \bar{\boldsymbol{W}}_{L}^{0}\| = \|\bar{\boldsymbol{F}}_{L-1}^{\top} - \boldsymbol{W}_{L}^{0}\| - |c_{L}^{t} - 1|\|\bar{\boldsymbol{F}}_{L-1}^{\top}\| \ge 0,$$
(12)

implying that $\| \mathbf{W}_l^t - \mathbf{W}_l^0 \| \ge | \mathbf{W}_l^t - \bar{\mathbf{W}}_l^0 \|$ for all $1 < l \le L$.

Based on Lemma 4.1, we derive that feedback-weight matching reduces the network output error e^{t+1} over the train step t when combined with weight decay (Krogh & Hertz, 1991), as follows. The proof is detailed in Appendix A

Proposition 4.1. Let e^t denote the output error of a two-layer fully connected non-linear network (i.e., L=2) at the t-th training step, η is the learning rate, $\gamma \leq \lambda_{min}(\bar{G})$ is a positive constant, where $\bar{G} = \mathbb{E}_{w \sim \mathcal{N}(0, I_p)} \psi(w^\top x_i) \psi(w^\top x_j)$ with the number of neuron as p and a non-linear function $\psi(\cdot)$, λ^t is the weight-decay hyperparameter at the step t, and y is the output of the network. By applying feedback-weight matching in Equation (6) and (7), the following holds:

$$\|\boldsymbol{e}^{t+1}\| \le \left(1 - \frac{\eta\gamma}{4} - \eta\lambda^t\right)\|\boldsymbol{e}^t\| + \lambda^t\|\boldsymbol{y}\| - \alpha_2 r_2$$
(13)

for all $t \ge 0$ and some constants α_2 , with r_2 defined in (11).

It is shown (Song et al., 2021) that the inequality in Equation (13), i.e., $||e^{t+1}|| \le (1 - \frac{\eta\gamma}{4} - \eta\lambda^t)||e^t|| + \lambda^t ||\mathbf{y}||$, holds for a two-layer fully connected non-linear network when applying FA (Feedback Alignment) (Lillicrap et al., 2016) with weight decay (Krogh & Hertz, 1991). Specifically, the right-hand side of the inequality, i.e., $(1 - \frac{\eta\gamma}{4} - \eta\lambda^t)||e^t|| + \lambda^t ||\mathbf{y}||$, consists of the following term as a linear component in fine-tuning:

$$\|\boldsymbol{W}_2^t - \boldsymbol{W}_2^0\|$$
 s.t. $\boldsymbol{W}_2^0 \propto \boldsymbol{F}_1^ op$ (14)

where W_2^0 is the pre-trained weights. By assuming that W_2^0 is replaced with the re-initialized weights, \overline{W}_2^0 in Equation (7), $\|e^{t+1}\|$ in Equation (13) is decreased by $\alpha_2 r_2$ since $\|W_2^t - \overline{W}_2^0\| \ge \|W_2^t - \overline{W}_2^0\|$, as in Lemma 4.1.

Conjecture 4.2. Given an L-layer fully connected non-linear network, suppose that the right-hand side of the inequality in Equation (13), i.e., $(1 - \frac{\eta\gamma}{4} - \eta\lambda^t) \|\mathbf{e}^t\| + \lambda^t \|\mathbf{y}\|$, contains $\|\mathbf{W}_l^t - \mathbf{W}_l^0\|$ as linear components for some $1 < l \leq L$. Then, based on Proposition 4.1 and Lemma 4.1, it is conjectured that Equation (13) can be generalized into:

$$\|\boldsymbol{e}^{t+1}\| \leq \left(1 - \frac{\eta\gamma}{4} - \eta\lambda^{t}\right)\|\boldsymbol{e}^{t}\| + \lambda^{t}\|\boldsymbol{y}\| - \sum_{l=2}^{L} \alpha_{l}r_{l}$$
(15)

with constants α_l , and r_l defined in (11) and (12) for some $1 < l \le L$ and all $t \ge 0$. Error Equation (12) and subsequently Equation (15) it can be seen that feedback w

From Equation (13), and subsequently Equation (15), it can be seen that feedback-weight matching preserves the weight decay effect by decreasing the network error $\|\boldsymbol{e}^{t+1}\|$ by the quantity $\eta\lambda^t\|\boldsymbol{e}^t\| - \lambda^t\|\boldsymbol{y}\|$. It is achieved by $\sum_{l=2}^{L} \alpha_l r_l$, which effectively counteracts the adverse impact of weight decay, namely, the increase in error $\|\boldsymbol{e}^{t+1}\|$ when $\eta\|\boldsymbol{e}^t\| \leq \|\boldsymbol{y}\|$, if $\sum_{l=2}^{L} \alpha_l r_l \geq \lambda^t\|\boldsymbol{y}\| - \eta\lambda^t\|\boldsymbol{e}^t\|$.

5 EXPERIMENT

295 We evaluate the proposed feedback-weight matching on two types of fine-tuning tasks. First, it is 296 applied to image classification tasks using two fully connected networks with four and six hidden 297 layers, respectively. These networks are pre-trained with CIFAR-100 (Krizhevsky et al., 2009) and 298 TinyImageNet (Le & Yang, 2015) using back-propagation (Rumelhart et al., 1986), and then finetuned on CIFAR-10 (Krizhevsky et al., 2009), SVHN (Netzer et al., 2011), and STL-10 (Coates 299 et al., 2011) through DFA applying feedback-weight matching. Next, we apply it to natural lan-300 guage processing (NLP) tasks using Transformer models, i.e., BERT-Tiny and Small (Kenton & 301 Toutanova, 2019; Turc et al., 2019), pre-trained on BookCorpus (Zhu et al., 2015) & Wikipedia, 302 and then fine-tuned with a set of GLUE tasks (Wang, 2018). For fine-tuning of BERT, feedback-303 weight matching is applied to the attention, intermediate, and block outputs of the encoder layers in 304 a similar way to previous works (Launay et al., 2020) that attempt to apply DFA to Transformer's 305 attention architectures (Vaswani, 2017). It is important to highlight that even standard DFA has 306 rarely been applied to Transformer models for from-scratch training due to its inherent challenges 307 and difficulties (Launay et al., 2020). Our experiment is the first attempt to apply DFA fine-tuning 308 to Transformer models (i.e., BERT), which is considered more challenging than from-scratch DFA 309 training. The detailed experimental setups are provided in Appendix E.

310 311

312

276

284 285 286

289

5.1 FINE-TUNING PERFORMANCE

Table 1 summarizes the fine-tuning performance on image classification tasks (i.e., test classification 313 accuracy) of the proposed feedback-weight matching compared against 1) back-propagation-based 314 fine-tuning and 2) standard DFA fine-tuning that does not apply feedback-weight matching. As 315 shown in the table, the proposed feedback-weight matching enables reliable fine-tuning for vari-316 ous network architectures and tasks, which consistently outperforms standard DFA with an aver-317 age of 2.16% accuracy gap, while underperforming when compared to back-propagation with an 318 average of 2.32%. For instance, feedback-weight matching achieves 82.67% accuracy when fine-319 tuning the 6-layer network from CIFAR-100 to SVHN, which is 7.97% higher than standard DFA 320 that achieves 74.70%, but 1.67% lower than back-propagation. It also indicates that the proposed 321 feedback-weight matching maintains more robust performance over network depths, whereas the performance of standard DFA deteriorates with deeper networks. For instance, in the case of fine-322 tuning from CIFAR-100 to SVHN, the accuracy drop between 4-layer and 6-layer networks is only 323 0.20% with feedback-weight matching, which is 24x smaller than the case not applying it (4.85%).

Table 1: Image Classification Tasks. The fine-tuning performance of feedback-weight matching (DFA_{ours}) on the 4 and 6-layer fully connected networks, compared with standard DFA fine-tuning (DFA_{fine}) and back-propagation fine-tuning (BP_{fine}). The pre-trained weights are obtained through back-propagation (BP). For reference, we also present the from-scratch-training results of back-propagation (BP_{scratch}) and DFA (DFA_{scratch}).
 The bold indicates better performance in DFA fine-tuning.

		Target Data	Source Data										
	Model		Scratch			CIFAR-10	0		TinyImageNet				
_			BPscratch	DFA scratch	BP _{fine}	DFA _{fine}	DFA ours	BP _{fine}	DFA _{fine}	DFA _{ours}			
_	4 layers	CIFAR-10	55.48	52.78	57.16	53.79	55.38	57.66	56.75	55.51			
		SVHN	85.10	82.93	84.32	79.55	82.87	84.69	80.31	83.16			
		STL-10	43.15	42.20	47.73	44.83	45.30	50.29	50.62	45.61			
_		CIFAR-10	54.93	51.94	58.85	53.04	55.39	55.97	51.08	55.54			
	6 layers	SVHN	85.10	81.89	84.34	74.70	82.67	84.72	76.03	81.39			
		STL-10	43.10	40.48	47.78	43.42	45.28	47.63	43.33	45.21			

Table 2 presents the evaluation results of feedback-weight matching applied to BERT-Tiny and BERT-Small, fine-tuned for NLP tasks, using the same experimental setup in image classification tasks (Table 1). Similar to image classification tasks, feedback-weight matching enables DFA to fine-tune BERT for various tasks of the GLUE dataset in a more robust manner compared to standard DFA. For example, on CoLA, feedback-weight matching achieves a Matthews correlation of 0.53 in BERT-Small, compared to 0.06 with standard DFA. Similarly, on STSB, BERT-Small achieves a Pearson correlation of 0.76 with feedback-weight matching, while standard DFA yields only 0.10, demonstrating a significant gap in both learning performance and reliability. In the worst case, standard DFA fails to learn from the fine-tuning data entirely, achieving 0.00 Matthews correlation for CoLA with BERT-Tiny, whereas feedback-weight matching achieves 0.29.

Table 2: NLP Tasks. The fine-tuning performance of feedback-weight matching (DFA_{ours}) on Transformer
 architectures (i.e., BERT-Tiny and BERT-Small), compared with standard DFA fine-tuning (DFA_{fine}) and back propagation-based fine-tuning (BP_{fine}). The pre-trained weights are obtained via back-propagation (BP). For
 reference, we also present the from-scratch-training results of back-propagation (BP_{scratch}) and DFA (DFA_{scratch}).
 The bold indicates better performance in DFA fine-tuning.

	Model	Training	CoLA	SST-2	MRPC	QQP	MNLI	QNLI	STSB	RTE	WNLI
		manning	(mat-cor)	(acc)	(acc)	(acc)	(acc)	(acc)	(pearson)	(acc)	(acc)
	BERT-Tiny	BPscratch	0.07	96.3	67.4	82.8	63.4	89.2	-0.19	64.1	50.0
		BP _{fine}	0.00	93.5	70.7	86.9	73.8	88.2	-0.25	60.3	52.6
		DFAscratch	0.00	95.2	67.4	81.2	59.2	84.2	-0.11	50.2	50.0
		DFA _{fine}	0.00	92.4	67.4	80.6	60.0	80.2	-0.17	51.2	51.0
		DFA _{ours}	0.29	95.9	69.7	82.3	60.2	84.3	0.36	55.5	52.6
		BPscratch	0.55	96.3	95.4	91.3	75.3	93.4	0.67	89.8	51.9
		BP _{fine}	0.87	98.9	96.7	98.0	93.0	99.1	0.90	94.0	53.3
	BERT-Small	DFAscratch	0.19	96.5	75.2	86.7	67.4	80.9	0.05	60.0	50.3
		DFA _{fine}	0.06	95.6	70.9	86.0	67.0	85.3	0.10	59.0	49.3
		DFA _{ours}	0.53	97.3	92.5	86.9	65.8	87.2	0.76	59.0	51.0

5.2 WEIGHT ALIGNMENT (WA) AND GRADIENT ALIGNMENT (GA)

Figure 1a and 1b plot the weight alignment (WA) and the gradient alignment (GA), along with the train and test accuracy, for some fine-tuning setups. As shown in the figures, the proposed feedback-weight matching (green) induces strong weight alignment (WA) from the outset, subsequently strong gradient alignment (GA) as analyzed in Section 3.2 and 3.3, leading to both enhanced train and test accuracy across all experiments with faster and stable convergence. In contrast, standard DFA (yel-low), not applying feedback-weight matching, achieves significantly lower WA and GA. While they gradually increase over fine-tuning epochs in some cases, the initially low WA and GA impede effective fine-tuning. As a result, the train and test accuracy of standard DFA do not improve sub-stantially from the pre-trained weight parameters, especially for BERT-Small. This suggests that standard DFA struggles to adapt to the target dataset during fine-tuning, likely due to the mismatch between its random feedback matrices and the pre-trained weights. In other words, it overly relies on pre-trained weights in the hope that they will fit and perform well on new target fine-tuning data.

397

398

399 400

401

402

412 413



Figure 1: WA, GA, train accuracy, and test accuracy. The green graph indicates DFA fine-tuning with feedback-weight matching (ours), the yellow indicates DFA fine-tuning without feedback-weight matching, the blue indicates DFA trained from scratch, and the gray indicates fine-tuning with back-propagation.

Table 3: Ablation experiment. The fine-tuning performance when removing weight matching (DFA_{weight*}), feedback matching (DFA_{feed}*), and weight decay (DFA_{decay}*). 'DFA_{ours}' denotes applying all of them.

		Source Data										
Model	Target Data		CIFAF	R-100			TinyIm	ageNet				
		DFA _{weight*}	DFA _{feed*}	DFA _{decay*}	DFA _{ours}	DFA _{weight*}	DFA _{feed*}	DFA _{decay*}	DFA _{ours}			
	CIFAR-10	53.92	55.23	48.82	55.38	53.73	55.05	48.66	55.51			
4 layers	SVHN	80.65	81.34	77.99	82.87	79.77	83.13	77.63	83.16			
	STL-10	44.25	45.20	40.00	45.30	44.05	45.42	40.47	45.61			
	CIFAR-10	53.47	55.03	46.21	55.39	53.50	55.03	45.77	55.54			
6 layers	SVHN	79.70	82.76	76.71	82.67	79.77	82.76	76.76	82.72			
	STL-10	43.86	45.42	39.17	45.28	43.78	45.43	40.23	45.21			

5.3 ABLATION STUDY: FEEDBACK MATCHING, WEIGHT MATCHING, AND WEIGHT DECAY

Table 3 presents the impact of feedback matching, weight matching, and weight decay on fine-tuning 414 with DFA. To assess their effectiveness, we remove each of them in isolation. Removing feedback 415 matching results in a marginal performance decline, such as a reduction from 55.54% to 55.03% 416 when the 6-layer network is fine-tuned from TinyImageNet to CIFAR-10. This marginal drop occurs 417 because bypassing feedback matching applies random feedback matrices to the re-initialized weights 418 that are amenable to arbitrary random feedback matrices, resulting in a reasonable level of WA 419 and GA. In contrast, omitting weight matching leads to a relatively bigger performance drop, e.g., 420 classification accuracy decreases from 83.16% to 79.77% when fine-tuning the 4-layer network from 421 TinyImageNet to SVHN. Similarly, the correlation score drops from 0.76 to -0.06 when fine-tuning 422 BERT-Small to STSB as shown in Table 5 (Appendix D). It is presumed that excluding weight matching causes the pre-trained weights obtained by back-propagation, not by DFA, to be fine-tuned 423 with mismatched feedback matrices, thereby resulting in weak WA and GA. 424

425 When weight decay is not applied, the fine-tuning of feedback-weight matching performance also 426 exhibits some declines, e.g., classification accuracy decreases from 55.38% to 48.82% when fine-427 tuning the 4-layer network from CIFAR-100 to CIFAR-10. It should be noted that weight decay 428 appears to have minimal impact on fine-tuning of standard DFA when feedback-weight matching is 429 not applied; in our experiment, the classification accuracy even increases, such as from 54.38% to 56.75% when fine-tuning the 4-layer network from TinyImageNet to CIFAR-10. This demonstrates 430 the synergistic effect of feedback-weight matching and weight decay, i.e., reducing network output 431 error as shown in Section 4.

1

432 5.4FEEDBACK-WEIGHT MATCHING AND WEIGHT DECAY 433

434 To evaluate the impact of feedback-weight matching on weight decay, we measure the fine-tuning 435 performance with weight decay, with and without applying feedback-weight matching, which is 436 shown in Table 4. The results indicate that weight decay enhances fine-tuning accuracy (reducing network output error) when used in conjunction with feedback-weight matching, with an average im-437 provement of 8.35%. This demonstrates that feedback-weight matching facilitates weight decay in 438 reducing network output error, thereby improving fine-tuning accuracy, as provided in Equation (15). 439 In contrast, weight decay is less likely to improve fine-tuning performance without feedback-weight 440 matching. In fact, when applied to the standard DFA (not applying feedback-weight matching), 441 weight decay results in fine-tuning accuracy with minimal variation (providing similar accuracy). 442

Table 4: Feedback-weight matching and weight decay . 'DFA_{fine}' applies weight decay without feedbackweight matching, compared with 'DFA_{ours}' applying both weight decay and feedback-weight matching.

A lower

			4	+ layers			0 layers					
Tar	Target Data		CIFAR-100		inyImage	Net	CIFA	R-100	TinyIı	TinyImageNet		
		DFA	fine DFA _{ot}	irs DF/	A _{fine} DI	Aours	DFA _{fine}	DFA _{ours}	DFA _{fine}	DFA _{fine} DFA _{our}		
CI	CIFAR-10		39 55.3 8	3 54	54.38 55		54.08	55.39	53.50	55.54		
S	SVHN		77 82.8 7	7 80	80.74 8		78.73	82.67	79.57	82.7	2	
S	STL-10		00 45.30) 50	.40 4	5.61	43.56	45.28	45.28	45.2	1	
(b) Fine-tuning NLP tasks (BERT)												
Mode	1 1	raining	(mat-cor)	(acc)	(acc)	(acc)	(acc)	(acc)	(pearson)	(acc)	(acc)	
BERT-T	iny D D	OFA _{fine} OFA _{ours}	0.00 0.29	92.4 95.9	67.4 69.7	80.6 82.3	60.0 60.2	80.2 84.3	-0.17 0.36	51.2 55.5	51.0 52.6	
BERT-Sn	nall D	OFA _{fine} OFA _{ours}	0.06 0.53	95.6 97.3	70.9 92.5	86.0 86.9	67.0 65.8	85.3 87.2	0.10 0.76	59.0 59.0	49.3 51.0	

(a) Fine-tuning image classification tasks (fully connected networks)

Т

6 101100

Figure 2 plots the weight alignment (WA), gradient alignment (GA), training accuracy, and test accuracy across varying strengths of weight decay during the fine-tuning of 4-layer network from CIFAR-100 to CIFAR-10. The proposed feedback-weight matching ensures strong WA and GA as discussed in Section 3.2 and 3.3 from the beginning, which helps mitigate alignment degradation (Song et al., 2021), while exhibiting varying behaviors depending on different levels of weight decay. In the absence of weight decay (black curve), GA declines and exhibits significant oscillations, ultimately causing a decrease in test accuracy. Conversely, when a strong weight decay is applied (blue curve), both WA and GA decrease sharply, followed by substantial reductions in both training and test accuracy. These observations suggest that an appropriate weight decay strength is crucial for effective fine-tuning (green curve) when applied with feedback-weight matching.



Figure 2: WA, GA, train accuracy, and test accuracy over different weight decays (0, 5e-4, 1e-3, and 1e-2). A 4-layer fully connected network is fine-tuned from CIFAR-100 to CIFAR-10 by feedback-weight matching.

LIMITATIONS AND FUTURE WORKS 6

We discuss the limitations and future works of this paper in Appendix C.

483

484 485

443

444

445

446

458

459 460

461

462

463

486 7 CONCLUSION

We propose feedback-weight matching, a method that enhances the fine-tuning capability and stability of Direct Feedback Alignment (DFA) for pre-trained networks. While standard DFA struggles in fine-tuning networks trained via back-propagation, the proposed feedback-weight matching
improves weight and gradient alignment, boosting stability and performance of DFA fine-tuning.
Combined with weight decay, it also reduces over-fitting and network errors. Our experiments show
significant improvements in both image classification and NLP tasks compared to standard DFA.

494 495

496

501

502 503

504

8 REPRODUCIBILITY STATEMENT

For reproduction of the experimental results presented in this paper, we provide access to an anony mous GitHub repository¹ containing the code implementation and reproduction instructions. The
 detailed experimental setups are provided in Appendix E.

References

- AF Agarap. Deep learning using rectified linear units (relu). arXiv preprint arXiv:1803.08375, 2018. 15
- Dominique Chu and Florian Bacho. Random feedback alignment algorithms to train neural networks: why do they align? *Machine Learning: Science and Technology*, 5(2):025023, 2024. 1, 3
- Kenneth Ward Church, Zeyu Chen, and Yanjun Ma. Emerging trends: A gentle introduction to fine-tuning. *Natural Language Engineering*, 27(6):763–778, 2021. 1
- Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 215–223. JMLR Workshop and Conference Proceedings, 2011. 6, 15
- Brian Crafton, Abhinav Parihar, Evan Gebhardt, and Arijit Raychowdhury. Direct feedback alignment with sparse connections for local learning. *Frontiers in neuroscience*, 13:525, 2019. 1, 3
- Francis Crick. The recent excitement about neural networks. *Nature*, 337(6203):129–132, 1989. 1,
 2
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi erarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009. 3
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
 1, 2, 3
- Shivam Garg and Santosh Vempala. How and when random feedback works: A case study of low-rank matrix factorization. In *International Conference on Artificial Intelligence and Statistics*, pp. 4070–4108. PMLR, 2022. 1
- Stephen Grossberg. Competitive learning: From interactive activation to adaptive resonance. *Cognitive science*, 11(1):23–63, 1987. 1, 2
- Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: a factorizationmachine based neural network for ctr prediction. *arXiv preprint arXiv:1703.04247*, 2017. 3
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015. 15
- 539 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415, 2016. 16

540 541 542	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of naacL-HLT</i> , volume 1, pp. 2, 2019. 2, 6, 15, 16
543 544 545	Diederik P Kingma. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> , 2014. 15
546 547	Virginia Klema and Alan Laub. The singular value decomposition: Its computation and some appli- cations. <i>IEEE Transactions on automatic control</i> , 25(2):164–176, 1980. 4, 14
549 550	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 3, 6, 15
551 552	Anders Krogh and John Hertz. A simple weight decay can improve generalization. Advances in neural information processing systems, 4, 1991. 2, 3, 5, 6
555 555	Julien Launay, Iacopo Poli, and Florent Krzakala. Principled training of neural networks with direct feedback alignment. <i>arXiv preprint arXiv:1906.04554</i> , 2019. 1, 2, 3
556 557 558	Julien Launay, Iacopo Poli, François Boniface, and Florent Krzakala. Direct feedback alignment scales to modern deep learning tasks and architectures. <i>Advances in neural information processing systems</i> , 33:9346–9360, 2020. 1, 3, 6
559 560	Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015. 6, 15
561 562 563 564	Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Akerman. Random synaptic feedback weights support error backpropagation for deep learning. <i>Nature communications</i> , 7(1): 13276, 2016. 1, 2, 3, 6, 14
565 566	Timothy P Lillicrap, Adam Santoro, Luke Marris, Colin J Akerman, and Geoffrey Hinton. Back- propagation and the brain. <i>Nature Reviews Neuroscience</i> , 21(6):335–346, 2020. 2
567 568	I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 16
569 570	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. <i>arXiv preprint arXiv:1609.07843</i> , 2016. 3
571 572 573 574	Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. <i>Communications</i> <i>of the ACM</i> , 65(1):99–106, 2021. 3
575 576 577	Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In <i>NIPS workshop on deep</i> <i>learning and unsupervised feature learning</i> , volume 2011, pp. 4. Granada, 2011. 6, 15
578 579 580	Arild Nøkland. Direct feedback alignment provides learning in deep neural networks. Advances in neural information processing systems, 29, 2016. 1, 2, 3, 5, 15, 16
581 582 583	Maria Refinetti, Stéphane d'Ascoli, Ruben Ohana, and Sebastian Goldt. Align, then memorise: the dynamics of learning with feedback alignment. In <i>International Conference on Machine Learning</i> , pp. 8925–8935. PMLR, 2021. 2, 3, 4, 5, 13
584 585	Sebastian Ruder. An overview of gradient descent optimization algorithms. <i>arXiv preprint</i> arXiv:1609.04747, 2016. 2
587 588	David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back- propagating errors. <i>nature</i> , 323(6088):533–536, 1986. 1, 2, 3, 4, 5, 6, 14, 15, 16
589 590 591	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014. 3
592 593	Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Con- tinuous 3d-structure-aware neural scene representations. <i>Advances in Neural Information Pro-</i> <i>cessing Systems</i> , 32, 2019. 3

594 595 596 597 598 599 600	 Ganlin Song, Ruitu Xu, and John Lafferty. Convergence and alignment of gradient descent with random backpropagation weights. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 19888–19898. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/a576eafbce762079f7d1f77fcalc5cc2-Paper.pdf. 3, 6, 9, 13 Julia Tura Ming Wei Chang Kanten Lag and Kristing Touteneur. Well read students lagra battern.
601 602 603	On the importance of pre-training compact models. <i>arXiv preprint arXiv:1908.08962</i> , 2019. 6, 15, 16
604	Ashish Vaswani. Attention is all you need. arXiv preprint arXiv:1706.03762, 2017. 3, 6
605 606 607	Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understand- ing. <i>arXiv preprint arXiv:1804.07461</i> , 2018. 6, 16
608 609 610 611 612	Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In <i>The IEEE International Conference on Computer Vision (ICCV)</i> , December 2015. 6
613 614	
615	
616	
617	
618	
619	
620	
621	
622	
623	
624	
625	
626	
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	
638	
039	
04U 641	
04 I 642	
6/3	
647	
645	
646	
647	

A PROOF

A.1 PROOF OF PROPOSITION 3.3

Proof. We prove Proposition 3.3 for $W_{1 < l < L}^t$ in Equation (16), and the same reasoning applies to W_L^t in (17). Since $A_{l \ge 2}^t$ in Equation (3) becomes such that $A_{l \ge 2}^t \propto I$ as the training proceeds (Refinetti et al., 2021), the weight newly updated with DFA, which is denoted as $\bar{W}_{1 < l < L}^t$, comes to satisfy Equation (4), i.e., $\bar{W}_{1 < l < L}^t = c_l^t F_l F_{l-1}^\top$ with some constant c_l^t . Given that we take the pretrained weight $W_{1 < l < L}^0$ as the initial point in our fine-tuning, the overall weight $W_{1 < l < L}^t$ obtained by DFA is expressed as the sum of $W_{1 < l < L}^0$ and $\bar{W}_{1 < l < L}^t$, which is given by:

$$\boldsymbol{W}_{1(16)$$

$$\boldsymbol{W}_{L}^{t} = \boldsymbol{W}_{L}^{0} + \bar{\boldsymbol{W}}_{L}^{t} = \boldsymbol{W}_{L}^{0} + c_{L}^{t} \boldsymbol{F}_{L-1}^{\top} \not\propto \boldsymbol{F}_{L-1}^{\top}$$
(17)

where $c_{1<l\leq L}^t$ is a constant. In Equation (16), since $W_{1<l< L}^0$ is unlikely to be proportional to $F_l F_{l-1}^{\top}$, i.e., $W_{1<l< L}^0 \propto F_l F_{l-1}^{\top}$, the overall weight $W_{1<l< L}^t$, which includes $W_{1<l< L}^0$, is also unlikely to be proportional to $F_l F_{l-1}^{\top}$, i.e., $W_{1<l< L}^t \propto F_l F_{l-1}^{\top}$, though $\bar{W}_{1<l< L}^t = c_l^t F_l F_{l-1}^{\top} \propto F_l F_{l-1}^{\top}$. Hence, Equation (16) can hardly induce strong WA in Equation (4).

A.2 PROOF OF PROPOSITION 3.6

Proof. Similar to (16) and (17), the overall weight W_l^t obtained by DFA is the sum of W_l^0 and \bar{W}_l^t . Specifically, now that $\bar{W}_{1< l< L}^0 = \bar{F}_l \bar{F}_{l-1}^\top$ and $\bar{W}_L^0 = \bar{F}_{L-1}^\top$, these become proportional to $\bar{F}_l \bar{F}_{l-1}^\top$ and \bar{F}_{L-1} , respectively, as follows:

$$W_{1 < l < L}^{t} = \bar{W}_{1 < l < L}^{0} + \bar{W}_{1 < l < L}^{t} = \bar{F}_{l}\bar{F}_{l-1}^{\top} + c_{l}^{t}\bar{F}_{l}\bar{F}_{l-1}^{\top} = (1 + c_{l}^{t})\bar{F}_{l}\bar{F}_{l-1}^{\top} \propto \bar{F}_{l}\bar{F}_{l-1}^{\top}$$
(18)

$$\boldsymbol{W}_{L}^{t} = \bar{\boldsymbol{W}}_{L}^{0} + \bar{\boldsymbol{W}}_{L}^{t} = \bar{\boldsymbol{F}}_{L-1}^{\top} + c_{L}^{t} \bar{\boldsymbol{F}}_{L-1}^{\top} = (1 + c_{L}^{t}) \bar{\boldsymbol{F}}_{L-1}^{\top} \propto \bar{\boldsymbol{F}}_{L-1}^{\top}$$
(19)

with constants $c_{1 < l \le L}^t$, which aligns with the strong WA condition in Equation (4).

A.3 PROOF OF PROPOSITION 3.8

Proof. The weight at the second layer of the network, W_2^t , can be expressed with the pre-trained weight, W_2^0 , with the learning rate η , the number of neurons as $p, F_1 \in \mathbb{R}^p$, and $W_2^t \in \mathbb{R}^p$ as follows (Song et al., 2021).

$$\boldsymbol{W}_{2}^{t} = \boldsymbol{W}_{2}^{t-1} - \eta \frac{1}{\sqrt{p}} \boldsymbol{W}_{1}^{t-1} \boldsymbol{X}^{\top} \boldsymbol{e}^{t-1} = \boldsymbol{W}_{2}^{0} - \frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} \boldsymbol{W}_{1}^{t} \boldsymbol{X}^{\top} \boldsymbol{e}^{t}$$
(20)

For the standard DFA that does not apply feedback-weight matching in Equation (6) and (7), we have $G_{DFA} = F_1$ and $G_{BP} = W_2^t$. By using Equation (20), the gradient alignment (GA) defined in Equation (9) between them, which is denoted as $\cos_{DFA} \angle (F_1, W_2^t)$, is at least as follows.

$$\cos_{DFA} \angle (F_1, W_2^t) = \frac{F_1^\top W_2^t}{\|F_1\| \|W_2^t\|} = \frac{F_1^\top}{\|F_1\|} W_2^t = \frac{F_1^\top}{\|F_1\|} (W_2^0 - \frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t) \\ \|W_2^0 - \frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t\| \\ \ge \frac{F_1^\top}{\|F_1\|} (W_2^0 - \frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t) \\ \|W_2^0\| + \|\frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t\|$$
(21)

Conversely, when applying feedback-weight matching in Equation (6) and (7), we have $F_1 = W_2^0$ for L=2. Using Equation (20) again, GA between them, $\cos_{FWM} \angle (F_1, W_2^t)$, is at least as follows.

$$\cos_{FWM} \angle (F_1, W_2^t) = \frac{\frac{F_1^\top}{\|F_1\|} (W_2^0 - \frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t)}{\|W_2^0 - \frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t\|} \ge \frac{\frac{F_1^\top}{\|F_1\|} (F_1 - \frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t)}{\|F_1\| + \|\frac{\eta}{\sqrt{p}} \sum_{t=0}^{t'-1} W_1^t X^\top e^t\|}.$$
(22)

700 If we assume that both F_1 and W_2^0 follow the standard Gaussian distribution, we have $||F_1^\top W_2^0|| \le ||F_1||^2$ (Song et al., 2021). Thus, $\cos_{FWM} \angle (F_1, W_2^t)$ exhibits a higher lower bound compared to $\cos_{DFA} \angle (F_1, W_2^t)$, i.e., $\cos_{FWM} \angle (F_1, W_2^t) \ge \cos_{DFA} \angle (F_1, W_2^t)$, implying a higher GA. \Box

702 A.4 PROOF OF LEMMA 4.1

704 *Proof.* We show that $r_{1 \le l < L} \ge 0$ in Equation (11), and the same reasoning extends to r_L in (12). 705 Given that $\bar{W}_l^0 = \bar{F}_l \bar{F}_{l-1}^\top \propto W_l^t = c_l^t \bar{F}_l \bar{F}_{l-1}^\top$, we can interpret W_l^t as a scaled version of \bar{W}_l^0 , 706 which implies that $||W_l^t - \bar{W}_l^0||$ is small. Conversely, since W_l^0 is not proportional to W_l^t , i.e., 707 $W_l^0 \propto W_l^t = c_l^t \bar{F}_l \bar{F}_{l-1}^\top$, it follows that $||W_l^t - W_l^0||$ is generally larger than $||W_l^t - \bar{W}_l^0||$. Therefore, 708 $||W_l^t - \bar{W}_l^0||$ is likely smaller than $||W_l^t - W_l^0||$.

709 710

711 712

713

714 715 716

B DECOMPOSITION OF WEIGHT INTO FEEDBACK MATRICES

One way of finding feedback matrices \bar{F}_l and \bar{F}_{l-1}^{\top} in Equation (6) from $W_{1 < l < L}^0$, other than SVD (Singular Value Decomposition) (Klema & Laub, 1980), is to optimize the following objective \mathcal{L}_{FM} .

$$\mathcal{L}_{FM} = \frac{1}{2} \sum_{l=2}^{L-1} (\boldsymbol{W}_{l}^{0} \boldsymbol{h}_{l-1} - \bar{\boldsymbol{F}}_{l} \bar{\boldsymbol{F}}_{l-1}^{\top} \boldsymbol{h}_{l-1})^{2} + \frac{1}{2} (\boldsymbol{W}_{L}^{0} \boldsymbol{h}_{L-1} - \bar{\boldsymbol{F}}_{L-1} \boldsymbol{h}_{L-1})^{2} + \frac{1}{2} \sum_{l=1}^{L-1} (\boldsymbol{I} - \bar{\boldsymbol{F}}_{l}^{\top} \bar{\boldsymbol{F}}_{l})^{2}$$
(23)

⁷¹⁷ Here, \mathcal{L}_{FM} is minimized to ensure that the layer output, when replaced by the feedback matrix ⁷¹⁸ $\bar{F}_l \bar{F}_{l-1}^{\top} h_{l-1}$, matches the output obtained using the pre-trained weight $W_l^0 h_{l-1}$, while \bar{F}_l is to be ⁷¹⁹ orthogonal to itself in accordance with the regular DFA condition (Lillicrap et al., 2016).

720 721

722

C LIMITATIONS AND FUTURE WORKS

Extending to Different Architectures. Although this study presents the significant potential of fine tuning with DFA, its current application is restricted to fully connected networks. This limitation
 arises because, at present, DFA is predominantly effective for fully connected architectures, and
 further research is needed to extend its applicability to other network types. In our future work, we
 plan to explore the application of DFA fine-tuning to various network architectures, such as CNNs.
 Meanwhile, we anticipate the development of more generalized methods that will enable DFA to be
 applied across a broader range of network types, thereby enhancing the applicability of our work.

Improving Learning Performance. The learning performance of the proposed feedback-weight
 matching is shown to surpass both 1) training networks with DFA from scratch and 2) fine-tuning
 networks with DFA using random feedback matrices. While fine-tuning with DFA applying the
 proposed method achieves superior and more stable performance compared to them, it still falls
 short of the performance achieved with fine-tuning using back-propagation (Rumelhart et al., 1986).
 We plan to explore how to achieve fine-tuning performance comparable to that of back-propagation
 by investigating DFA from its fundamental mechanism, along with the proposed method.

737 **Proving Hypotheses.** This work provides some hypotheses regarding fine-tuning and weight decay in the context of DFA. For example, Conjecture 3.9 suggests that applying the proposed feedback-738 weight matching can achieve strong weight alignment (WA) for fully connected networks of arbi-739 trary depth. Additionally, Conjecture 4.2 posits that applying the proposed method to weight decay 740 enhances fine-tuning performance of DFA for fully connected networks of arbitrary layers. How-741 ever, formal proofs are necessary to substantiate these hypotheses and validate the efficacy of the 742 proposed approach. In future research, we intend to generalize the propositions presented in this 743 study to encompass various types of fully connected network architectures. 744

744 745 746

747

748

749

750

D ABLATION EXPERIMENT ON BERT

Table 5 presents the fine-tuning performance of BERT models when weight matching, feedback matching, and weight decay are individually removed. It is important to note that DFA is not applied to all fully connected layers in BERT, which limits the ability to properly assess the effectiveness of feedback-weight matching. Thus, this experimental setup may not provide an accurate evaluation.

751 752 753

754

E EXPERIMENTAL SETUPS

In this section, we offer an explanation of the experimental setup utilized throughout our research. Appendix E.1 outlines the training details of the feedback matrix used for feedback matching in all

758	Madal	Tasining	CoLA	SST-2	MRPC	QQP	MNLI	QNLI	STSB	RTE	WNLI
750	Model	Training	(mat-cor)	(acc)	(acc)	(acc)	(acc)	(acc)	(pearson)	(acc)	(acc)
759		DFA _{weight*}	0.00	94.7	67.4	81.4	59.2	88.4	-0.15	50.3	50.9
760	DEDT Time	DFA _{feed*}	0.00	95.8	68.9	82.4	60.8	86.9	0.35	55.5	50.0
761	BERT-Tiny	DFA _{decay*}	0.31	95.9	71.4	81.9	61.0	83.3	0.36	53.3	51.9
762		DFA _{ours}	0.29	95.9	69.7	82.3	60.2	84.3	0.36	50.8	52.6
763		DFA _{weight*}	0.08	96.0	75.1	85.0	66,7	79.7	-0.06	61.8	50.1
764	DEDT Small	DFA _{feed*}	0.54	97.0	91.5	87.4	65.2	85.3	0.75	62.0	50.2
704	DERI-Sillali	DFA _{decay*}	0.53	97.2	91.2	87.1	64.7	85.4	0.78	68.7	50.9
765		DFA _{ours}	0.53	97.3	92.5	86.9	65.8	87.2	0.76	59.0	51.0
766		I	I								

Table 5: Ablation experiment. The fine-tuning performance when removing weight matching (DFAweight*), feedback matching (DFA_{feed}*), and weight decay (DFA_{decay}*). 'DFA_{ours}' denotes applying all of them.

models. Appendix E.2 covers the configuration settings required for the fully connected network experiments. Appendix E.3 describes the setup necessary for experiments involving BERT, which employs a transformer architecture. To ensure the robustness of our findings, we report the average results over three different random seeds.

773 E.1 FEEDBACK MATRIX 774

775 We train feedback matrices to reconstruct pre-trained weights that were trained using backpropagation (Rumelhart et al., 1986). The loss function, in Equation (23), is used to guide the 776 feedback matching process. The two learned feedbacks are then combined and re-initialized into a 777 single weight matrix for each layer. We use the Adam optimizer (Kingma, 2014) without weight 778 decay or any scheduler. In fully connected networks, a learning rate of 1e-5 is applied, while in 779 transformers (BERT) (Kenton & Toutanova, 2019; Turc et al., 2019), a learning rate of 1e-3 is used. For all experiments on the model and dataset, training is conducted for 3 epochs with a batch size of 781 64

782 783

797

756

767 768

769

770

771

772

E.2 FULLY CONNECTED NETWORKS 784

785 We pre-train two fully connected networks with four and six layers on the CIFAR-100 (Krizhevsky 786 et al., 2009) and TinyImageNet (Le & Yang, 2015) datasets utilizing weights obtained through 787 back-propagation (BP). These pre-trained weights are subsequently fine-tuned on the CIFAR-788 10 (Krizhevsky et al., 2009), SVHN (Netzer et al., 2011), and STL-10 (Coates et al., 2011) datasets. During the pre-processing phase, we apply image resizing and normalization, without any augmen-789 tations. For Dynamic Feedback Alignment (DFA) (Nøkland, 2016), the weights are initialized with 790 a uniform distribution within the range of (-0.01, 0.01). Conversely, for back-propagation (Rumel-791 hart et al., 1986), we employ the He initialization (He et al., 2015). The optimization process is 792 carried out using Stochastic Gradient Descent, and ReLU (Agarap, 2018) is employed as the activa-793 tion function. The hyperparameters for both the 4-layer and 6-layer architectures remain consistent. 794 A comprehensive description of each hyperparameter under various training conditions is presented 795 in Table 6. 796

Table 6: Hyperparameters for fully connected networks training.

798									
799	Target Data	Hyperparmeters	BPscratch	BP_{fine}	DFAscratch	DFA _{fine}	DFA _{feed}	DFAweight	DFAours
		Learning Rate	1e-3	1e-3	1e-3	1e-3	1e-3	1e-3	1e-3
800		Batch size	64	64	64	64	64	64	64
801		Hidden Dim	1000	1000	1000	1000	1000	1000	1000
802		Input size	3072	3072	3072	3072	3072	3072	3072
803		Epochs	5000	5000	5000	5000	5000	5000	5000
000	CIFAR-10	Weight Decay	5e-4	5e-4	0	0	5e-4	5e-4	5e-4
804		Dropout	0.1	0.1	0	0	0	0	0
805		Epochs	5000	5000	5000	5000	5000	5000	5000
806	SVHN	Weight Decay	5e-4	5e-4	0	0	5e-4	5e-4	5e-4
807		Dropout	0.1	0.1	0	0	0	0	0
000		Epochs	5000	5000	5000	5000	30000	30000	30000
808	STL-10	Weight Decay	5e-4	5e-4	0	0	1e-3	1e-3	1e-3
809		Dropout	0.1	0.1	0	0	0.1	0.1	0.1

810 E.3 BERT

We train BERT-Tiny and Small models (Kenton & Toutanova, 2019; Turc et al., 2019) on the GLUE (Wang, 2018) dataset using the AdamW (Loshchilov, 2017) optimizer with a fixed learning rate and no scheduler. We apply weight decay and dropout techniques. GeLU (Hendrycks & Gimpel, 2016) is used for the activation function, which is commonly employed in BERT. Layers such as the encoder block outputs, intermediate outputs, and attention outputs are optimized using Dynamic Feedback Alignment (DFA) (Nøkland, 2016), while the projection layers for key, query, and value are trained using back-propagation (BP) (Rumelhart et al., 1986). The weights are initial-ized using a uniform distribution, and the feedback matrix is specifically designed to satisfy the left orthogonality condition. A comprehensive description of the hyperparameter values is presented in Table 7.

Table 7: Hyperparameters for BERT training.

823	Model	Hyperparmeters	Target Data	BPcaratah	BPfina	DFA	DFAfina	DFAfad	DFAwaight	DFA
824	moder	Batch size	Turger Duiu	64	64	64	64	64	64	64
825		Dropout		0.1	0.1	0.1	0.1	0.1	0.1	0.1
826		Weight Decay		0.01	0.01	0.01	0.01	0.01	0.01	0.01
827		Epochs		6	6	6	6	6	6	6
021		Max length		512	512	512	512	512	512	512
828		Num of heads		2	2	2	2	2	2	2
829		Num of layers		2	2	2	2	2	2	2
830		Hidden dim		128	128	128	128	128	128	128
831		Intermediate dim	C-LA	512	512	512	512	512	512	512
001			COLA SST-2	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
832	BERT-Tiny		MRPC	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
833			OOP	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
834		Learning Rate	MNLI	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
835			QNLI	5e-5	5e-5	5e-5	5e-5	5e-5	5e-5	5e-5
926			STSB	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
030			RTE	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
837			WNLI	5e-5	5e-5	5e-5	5e-5	5e-5	5e-5	5e-5
838		Num of heads		8	8	8	8	8	8	8
839		Num of layers		4	4	4	4	4	4	4
840		Hidden of dim		512	512	512 2048	512	512	512	512
0.4.4			CoLA	2040	2046	2040	2040	2040	2046	2046
841			SST-2	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
842	BERT-Small		MRPC	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
843	BERT-Sillali		OOP	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
844		Learning Rate	MNLI	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
0.1.5			QNLI	5e-5	5e-5	5e-5	5e-5	5e-5	5e-5	5e-5
043			STSB	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
846			RTE	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5
847			WNLI	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5	1e-5

ooo