# Training Large Neural Networks With Low-Dimensional Error Feedback

Maher Hanut and Jonathan Kadmon

Edmond and Lily Center for Brain Sciences,
The Hebrew University,
Jerusalem.
{maher.hanut, jonathan.kadmon}@mail.huji.ac.il

## Abstract

Training deep neural networks typically relies on backpropagating high-dimensional error signals—a computationally intensive process with little evidence supporting its implementation in the brain. However, since most tasks involve low-dimensional outputs, we propose that low-dimensional error signals may suffice for effective learning. To test this hypothesis, we introduce a novel local learning rule based on Feedback Alignment that leverages indirect, low-dimensional error feedback to train large networks. Our method decouples the backward pass from the forward pass, enabling precise control over error signal dimensionality while maintaining high-dimensional representations. We begin with a detailed theoretical derivation for linear networks, which forms the foundation of our learning framework, and extend our approach to nonlinear, convolutional, and transformer architectures. Remarkably, we demonstrate that even minimal error dimensionality—on the order of the task dimensionality—can achieve performance matching that of traditional backpropagation. Furthermore, our rule enables efficient training of convolutional networks, which have previously been resistant to Feedback Alignment methods, with minimal error. This breakthrough not only paves the way toward more biologically accurate models of learning but also challenges the conventional reliance on high-dimensional gradient signals in neural network training. Our findings suggest that low-dimensional error signals can be as effective as high-dimensional ones, prompting a reevaluation of gradient-based learning in high-dimensional systems. Ultimately, our work offers a fresh perspective on neural network optimization and contributes to understanding learning mechanisms in both artificial and biological systems.

# 1.  Introduction

The brain, despite its vast complexity, operates within the constraints of a low-dimensional world. Every movement—lifting a hand, turning a head—follows finite joint constraints. Navigation unfolds in modest spatial dimensions, and object recognition is drawn from a limited set of familiar categories. Even when possibilities seem vast, they are dwarfed by the sheer number of neurons and synapses in any brain region. Artificial neural networks reflect this paradox: the number of neurons in each layer often far exceeds the task-relevant variables. This contrast between task and internal representation dimensionality challenges us to rethink how we train such large networks and explore new learning paradigms.

Modern deep learning methods rely on high-dimensional gradients that propagate backward through layers, ensuring a precise credit assignment to each parameter [1, 2]. However, this process is computationally intensive and biologically implausible, given the lack of direct evidence for backpropagation in the brain [3, 4]. If the loss gradient does not propagate back through the network, it must be obtained through alternative routes [5]. In this case, there is no reason to assume that each neuron receives detailed information about the local gradient. Moreover, the brain likely operates with low-dimensional error signals, reflecting the inherently low-dimensional nature of tasks. This raises an essential question: Can large neural networks be trained efficiently with constrained, low-dimensional error signals? If error dimensionality is tied to task complexity rather than network size, it could revolutionize how we train large networks and improve our understanding of biological learning.

Despite the importance of this question, the effect of low-dimensional error feedback on learning and neural representations was never properly explored. The strong coupling of feedforward and backward passes in backpropagation, where both processes use the same synapses, has limited investigations into alternative training methods. While approaches such as Feedback Alignment have been proposed [5], the necessity of full gradients for effective learning has not been questioned.

In this study, we propose that high-dimensional error propagation is not strictly necessary for training deep, overparameterized networks. We introduce a novel learning rule based on Feedback Alignment that decouples forward and backward passes, enabling precise control over the dimensionality of error signals. Our method exploits the low-dimensional nature of real-world tasks to train large networks efficiently without compromising their representational capacity, achieving performance on par with backpropagation. Our theoretical analysis in linear settings highlights the potential pitfalls of naive low-dimensional error propagation and shows that a simple local online learning rule can recover backpropagation-level performance. Extending this approach to nonlinear, convolutional, and transformer-based networks, we demonstrate that learning with

low-dimensional error signals is both feasible and universal. Finally, we explore how error dimensionality influences neural representations by examining its effects on receptive fields in an early visual system model, underscoring the broader implications of our findings.

## 2. Background and related work

We consider a multilayered perceptron with $L$ layers, each layer $l$ computes its output as $\boldsymbol{h}_l = f(W_l \boldsymbol{h}_{l-1})$, where $W_l$ is the weight matrix and $f$ is an element-wise activation function. The input to the network is $\boldsymbol{h}_0 = \boldsymbol{x}$, and the final network output is $\hat{\boldsymbol{y}} = f_L(W_L \boldsymbol{h}_{L-1})$, which approximates the target $\boldsymbol{y}$. The *task dimensionality*, denoted $d$, is at most the number of components in $\boldsymbol{y}$ and $\hat{\boldsymbol{y}}$. Training the network involves minimizing a loss function $\mathcal{L}(\boldsymbol{y}, \hat{\boldsymbol{y}})$ by adjusting the weights $\{W_l\}$. The error signal at the output layer, $\delta_L = \frac{\partial \mathcal{L}}{\partial \hat{\boldsymbol{y}}}$, is a $d$-dimensional vector, typically much smaller than the number of neurons in the hidden layers.

**Backpropagation (BP)** [1] is the standard approach for training neural networks. It propagates the error backward through the network using $\delta_l = W_{l+1}^T \delta_{l+1} \odot f'(W_l \boldsymbol{h}_{l-1})$, and updates the weights using $\Delta W_l = -\eta \delta_l \boldsymbol{h}_{l-1}^T$, where $\eta$ is the learning rate. However, this method requires an exact transposition of the forward weights, $W_{l+1}^T$, which is biologically implausible [6, 7]. Moreover, BP tightly couples the error propagation with the forward pass, limiting the ability to explore how different properties of the error signal affect learning dynamics.

**Feedback Alignment (FA)** [5] was proposed to address the biological limitations of BP by replacing $W_{l+1}^T$ with a fixed random matrix $B_l$. The error is computed as:

$$\delta_l = B_l \delta_{l+1} \odot f'(W_l \boldsymbol{h}_{l-1}), \tag{1}$$

decoupling the forward and backward weights and providing a more biologically plausible mechanism. However, FA struggles to scale effectively in deep architectures, such as convolutional neural networks (CNNs), where it often fails to achieve competitive performance [8].

An extension of FA involves adapting $B_l$ by updating it alongside the forward weights $W_l$ to improve their alignment [9, 10]:

$$\Delta B_l = -\eta \boldsymbol{h}_{l-1} \delta_l^T - \lambda B_l, \quad \Delta W_l = -\eta \delta_l \boldsymbol{h}_{l-1}^T - \lambda W_l, \tag{2}$$

where $\lambda$ is a regularization parameter. Although this adaptive approach improves performance by better aligning forward and backward weights, it still requires high-dimensional error signals and struggles to match BP performance in complex architectures like CNNs. Furthermore, in Section 3, we show that this approach fails when the matrix $B$ is low-rank

3

and the dimensionality of the error is constrained.

Other studies have explored the use of *fixed* sparse feedback matrices to reduce the dimensionality of error propagation [11]. However, these approaches result in significantly lower performance and do not provide a systematic framework for studying how error constraints affect learning and representation formation.

Beyond FA-based methods, several studies have shown that weight updates using backpropagation can result in a low-dimensional weight update [12–14] and favor low-rank solutions [15]. These findings support our hypothesis that low-dimensional feedback is sufficient to train deep networks. However, no previous work has considered training with a constrained error pathway, and the effects of error dimensionality and training have not been systematically studied.

In this work, our aim is to systematically investigate how constraining the dimensionality of the error signal affects the training and performance of neural networks. To this end, we introduce a novel learning scheme, *Restricted Adaptive Feedback (RAF)*, that allows flexible control over the dimensionality of the errors (Figure 1).

**Our main contributions are**:

1. We present a novel learning rule, *Restricted Adaptive Feedback (RAF)*, which matches BP performance while requiring minimal error signals. We provide a detailed derivation of the learning dynamics in a simple linear case, establishing a foundational understanding of how RAF operates.

2. We demonstrate that nonlinear networks can efficiently learn nontrivial datasets using low-dimensional error signals, highlighting the versatility of RAF in practical scenarios.

3. We show that more complex yet highly useful architectures, such as convolutional networks and transformers, can also be effectively trained with low-dimensional feedback.

4. We reveal that error dimensionality shapes the receptive fields in a model of the ventral visual system, offering new insights into the relationship between learning mechanisms and biological neural representations.

In the final section, we discuss the broader implications of our results for both neuroscience and machine learning.
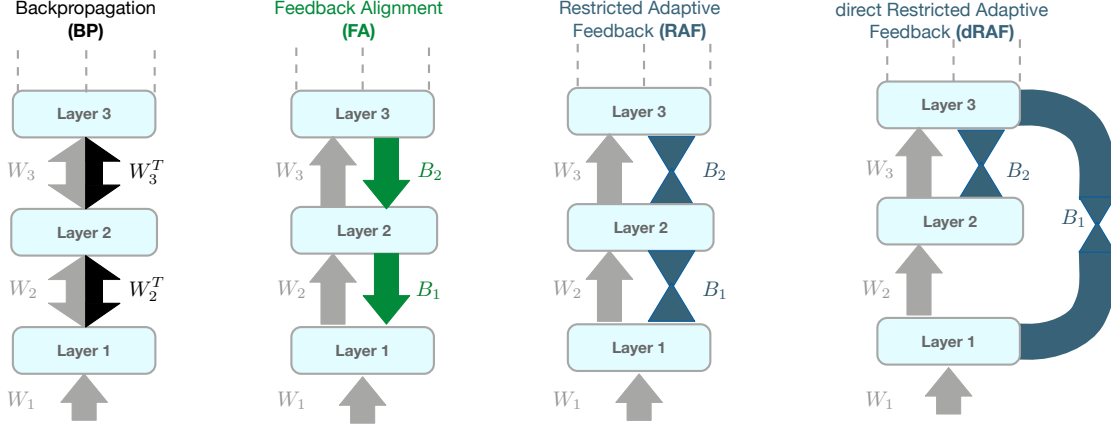
Figure 1: *Illustration of different approaches for propagating error to hidden layers.* From left to right: **Backpropagation (BP)** propagates error using the exact transpose of the forward weights. **Feedback Alignment (FA)** replaces the transposed weights with fixed random feedback matrices, which gradually align with the forward weights during training. **Restricted Adaptive Feedback (RAF)** constrains error propagation through low-rank feedback matrices, limiting error dimensionality and preventing an exact mirroring of BP. **Direct Restricted Adaptive Feedback (dRAF)** extends RAF by enabling error signals to bypass intermediate layers, propagating directly from the output layer or other non-adjacent layers.

# 3.  Low-dimensional error feedback in linear networks

We begin our analysis by studying learning dynamics in multilayered linear networks. Although linear models may seem overly simplistic, they can exhibit rich learning dynamics due to the nonlinearity introduced by the loss function [16]. Additionally, imposing dimensional constraints on linear networks yields insightful results that extend beyond the linear case.

**A linear problem**  We consider a simple linear transformation problem with a low-dimensional structure, $\boldsymbol{y} = A\boldsymbol{x}$. Here, $\boldsymbol{x} \in \mathbb{R}^n$ represents the $n$-dimensional input, and $\boldsymbol{y} \in \mathbb{R}^m$ represents the target. The matrix $A$ is a rank-$d$ matrix defined as $A = \sum_{j=1}^{d} \boldsymbol{u}_j \boldsymbol{v}_j^T$, where $\boldsymbol{u}_j \in \mathbb{R}^n$ and $\boldsymbol{v}_j \in \mathbb{R}^m$ are random Gaussian vectors, and we assume $d \ll n$. Our data set consists of $p$ training samples $\{\boldsymbol{x}^\mu, \boldsymbol{y}^\mu\}_{\mu=1}^p$, with each input vector $\boldsymbol{x}^\mu$ being i.i.d. according to the standard normal distribution $\boldsymbol{x}^\mu \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$. The labels are given by $\boldsymbol{y}^\mu = A\boldsymbol{x}^\mu + \boldsymbol{\xi}^\mu$, where $\boldsymbol{\xi}^\mu$ is additive Gaussian noise with zero mean and unit variance.

The goal is to learn the low-dimensional structure of $A$ from the $p$ samples using a linear neural network. For simplicity, we assume that $p$ is sufficiently large to allow the network to fully recover the structure of $A$.

**A linear network model.**  To study the effects of restricted error pathways, we consider a simple linear network with three layers: an input layer $\boldsymbol{x} \in \mathbb{R}^n$, a hidden layer $\boldsymbol{h} \in \mathbb{R}^k$,

and an output layer $\boldsymbol{y} \in \mathbb{R}^m$. The input and hidden layers are connected by the weight matrix $W_1 \in \mathbb{R}^{k \times n}$, and the hidden and output layers are connected by the weight matrix $W_2 \in \mathbb{R}^{m \times k}$. The output of the network can be expressed as $\boldsymbol{y} = W_2 W_1 \boldsymbol{x}$ (Figure 2).

The network is trained to minimize the quadratic empirical loss function:

$$L = \frac{1}{p} \sum_\mu \|\boldsymbol{y}(\boldsymbol{x^\mu}) - W_2 W_1 \boldsymbol{x^\mu}\|^2. \tag{3}$$

We apply Feedback Alignment (FA) to update $W_1$, which does not have direct access to the loss gradient. Instead of backpropagating the error through $W_2^T$, we use a fixed low-rank feedback matrix $B$. This provides an alternative pathway for propagating the error signal to $W_1$.

For a given data point $\{\boldsymbol{x^\mu}, \boldsymbol{y^\mu}\}$, the weight updates, derived from the FA framework, are given by:

$$\begin{aligned}
\Delta W_1^\mu &= \eta B^T(\boldsymbol{y^\mu} - W_2 W_1 \boldsymbol{x^\mu})\boldsymbol{x}^{\mu T}, \text{ and} \\
\Delta W_2^\mu &= \eta(\boldsymbol{y^\mu} - W_2 W_1 \boldsymbol{x^\mu})\boldsymbol{x}^{\mu T} W_1^T.
\end{aligned} \tag{4}$$

Here, $\eta$ represents the learning rate, and the update for $W_1$ is computed using the indirect error feedback provided by $B$, while $W_2$ receives the full error signal directly from the output.

**Constraining error dimensionality with low-rank feedback** To control the dimensionality of the error feedback, we impose a low-rank constraint on the feedback matrix $B$. Rather than allowing full-dimensional feedback, we decompose $B$ as $B = QP$, where $Q \in \mathbb{R}^{k \times r}$ and $P \in \mathbb{R}^{r \times m}$. When $r < \min(k, m)$, $B$ is low rank, which means that it can project the error signal onto at most $r$ independent directions.

This low-rank structure introduces an "$r$-bottleneck," which limits the flow of error information. In linear settings, the problem is solvable using a single weight matrix, rendering the training of $W_1$ unnecessary. However, backpropagation dynamics still adjust these weights [16]. By controlling the value of $r$, we can systematically study how reducing the dimensionality of the error signal impacts learning. Initially, we follow the original Feedback Alignment framework, keeping $Q$ and $P$ as random matrices. However, as we will demonstrate, allowing $Q$ and $P$ to learn is crucial to high performance.

## 3.1 Learning dynamics

Our analysis extends the framework established by Saxe et. al. [16] to incorporate indirect feedback with constrained dimensionality. We begin by characterizing the task across the $p$ data points through the input-output covariance matrix, $\Sigma_{io} = \frac{1}{p} \sum_{\mu=1}^p \boldsymbol{y^\mu}(\boldsymbol{x^\mu})^T$, which captures the correlation between input vectors $\boldsymbol{x}$ and output vectors $\boldsymbol{y}$. Performing Singular Value Decomposition (SVD), we obtain $\Sigma_{io} = USV^T$, where $U \in \mathbb{R}^{m \times m}$ and

$V \in \mathbb{R}^{n \times n}$ contain the left and right singular vectors, respectively, and $S \in \mathbb{R}^{m \times n}$ is a rectangular diagonal matrix of singular values. For sufficiently large $p$, the first $d$ singular values in $S$ correspond to the prominent directions in the data (i.e., the singular values of $A$), while the remaining singular values are $O(1/\sqrt{p})$ and dominated by noise.

To track how training aligns the network weights with these prominent directions, we rotate the weight matrices $W_1$, $W_2$, and $B$ using the singular vectors of $\Sigma_{io}$. This transformation simplifies the analysis by aligning the network's weight dynamics with the key data directions:

$$W_1 = \bar{W}_1 V^T, \quad W_2 = U \bar{W}_2, \quad B = \bar{B} U^T,$$

where $\bar{W}_1$, $\bar{W}_2$, and $\bar{B}$ represent the transformed weight matrices. This rotation aligns the weight dynamics with the dominant singular vectors, allowing us to focus on how the network captures the important features of the data.

Since the inputs are uncorrelated, we can apply these transformations to the iterative weight-update equations derived from the FA learning rule. Assuming a small learning rate $\eta \ll 1$ with full-batch updates, we express the weight updates in continuous time:

$$\tau \frac{d\bar{W}_1}{dt} = \bar{B}^T(S - \bar{W}_2 \bar{W}_1), \quad \tau \frac{d\bar{W}_2}{dt} = (S - \bar{W}_2 \bar{W}_1)\bar{W}_1^T, \tag{5}$$

where $\eta = dt/\tau$. This continuous form captures the dynamics of the learning process, allowing us to study it from a dynamical systems perspective. By analyzing these equations, we can identify fixed points and evaluate their stability, providing insight into how the network converges and learns under constrained feedback. Figure 2, shows how the singular vectors of $W_2 W_1$ align to the corresponding singular vectors of $\Sigma_{io}$ .

**Stationary solutions for the training.** Training halts when the right-hand side of the weight update equations (5) vanishes, indicating that the dynamics have reached a stable fixed point. At this fixed point, the update equation for $\bar{W}_1$ leads to the condition:

$$\bar{B}(S - \bar{W}_2 \bar{W}_1) = 0 \implies S_{jj}\boldsymbol{B}_{:,j} = \sum_{i=1}^{m} \boldsymbol{B}_{:,i}(\bar{W}_2 \bar{W}_1)_{ij} \quad \forall j, \tag{6}$$

where $\boldsymbol{B}_{:,j} \in \mathbb{R}^k$ is the $j$-th column of $\bar{B}$, and $S_{jj}$ is the $j$-th singular value of $\Sigma_{io}$. This equation indicates that, at the stationary point, the weight products $\bar{W}_2 \bar{W}_1$ must align with the singular modes of the data.

However, since $B$ is of rank $r$, the feedback matrix $\bar{B}$ can only span at most $r$ independent directions. If $r = m$, the system has enough feedback dimensionality to align perfectly with the singular values in $S$, recovering the full structure of $\Sigma_{io}$ as demonstrated in Saxe et. al. [16]. In this case, the training successfully converges to a unique solution where $\bar{W}_2 \bar{W}_1 = S$ (Figure 2b).

Crucially, when $r < m$, the feedback matrix $\bar{B}$ lacks the sufficient rank to fully capture the
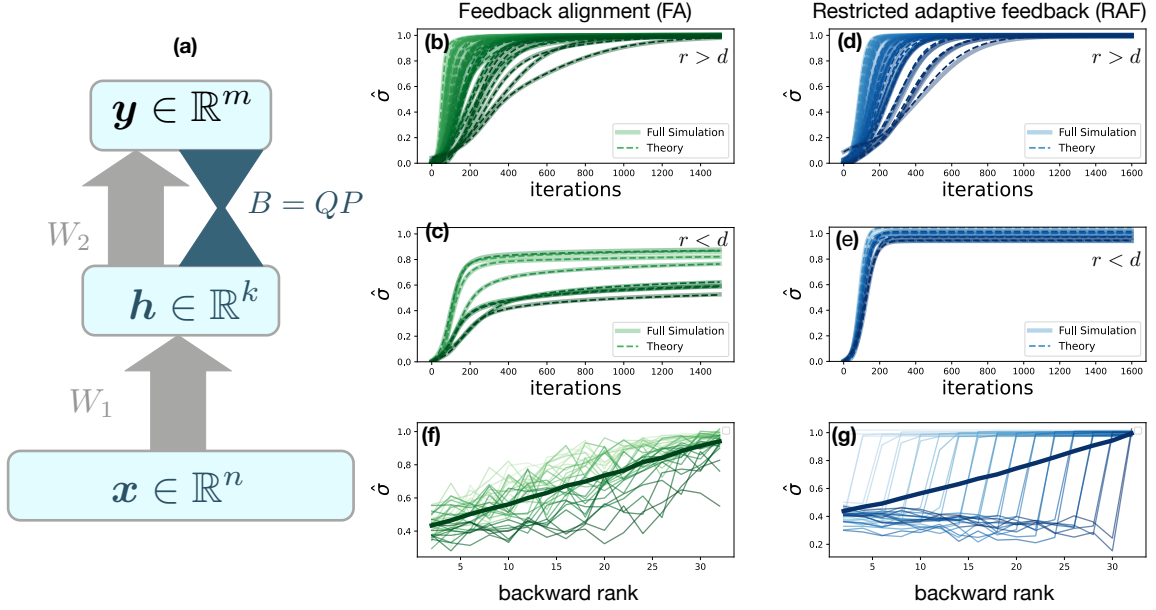
Figure 2: *Learning dynamics and component alignment in linear networks.* **(a)** Schematic of the network architecture with input dimension $n = 128$, hidden layer size $k = 64$, and output dimension $m = 64$. The feedback matrix $B$ is factorized as $B = QP$ and constrained to rank $r$. **(b, c)** Comparison of theoretical predictions (dashed) and numerical simulations (solid) for low-rank Feedback Alignment (FA), training only $Q$, with $r = 64$ and $r = 8$, respectively. The $y$-axis shows the alignment of singular vectors between $W_2 W_1$ and $\Sigma_{io}$. **(d, e)** Same as (b, c) but for Restricted Adaptive Feedback (RAF), where both $Q$ and $P$ are trained. **(f)** When $P$ is not trained, the singular modes align on average (bold), but the top $r$ components fail to recover fully. **(g)** Training $P$ to align with the error, as in RAF, ensures full recovery of the top $r$ singular components, improving convergence and alignment.

8

$m$ independent directions in $S$. As a result, eq. (6) becomes under-determined, leading to potentially infinite solutions. This means that the trained weights $\bar{W}_2\bar{W}_1$ may not align with the true data structure encoded in $\Sigma_{io}$ (Figure 2c)

The implications of the uderdetrimed fixed point solutions are surprising. Naively, one would expect that if the bottleneck is not too narrow, the projections of the error would maintain the necessary structure to guide the learning. Specifically, the John-son–Lindenstrauss lemma [17] suggests that as long as $r > d \log m$, the pairwise correlation structure of the error signal would be maintained. Nevertheless, our analysis shows that the solutions weight are not guaranteed to converge to the correct solution. Thus, in the case of low-rank feedback, Feedback Alignment (FA) is likely to fail. The solution to this predicament, as we show next, is aligning the feedback weights with the data, ensuring that the network learns the correct representations.

## 3.2   Training the feedback weights

From the previous analysis, it is apparent that the feedback matrix $B$ must be trained for the learning dynamics to converge to the correct solution of Eq. (6). When $B$ is low-rank and fixed, the network lacks sufficient capacity to transmit the full error information, which can impede learning.

One approach is to adopt a learning rule inspired by Kolen-Pollack, as in eq. (2). When $B = QP$, the updates are given by

$$\Delta W_2^\mu = \eta \boldsymbol{\delta}^\mu \boldsymbol{x}^{\mu T} W_1^T - \lambda W_2 \quad \text{and} \quad \Delta Q^\mu = \eta W_1 \boldsymbol{x}^\mu \boldsymbol{\delta}^{\mu T} P^T - \lambda Q, \tag{7}$$

where $\boldsymbol{\delta}^\mu = \frac{\partial \mathcal{L}(\boldsymbol{x}^\mu, \boldsymbol{y}^\mu)}{\partial \boldsymbol{y}} = \boldsymbol{y}^\mu - W_2 W_1 \boldsymbol{x}^\mu$ is the error gradient for the $\mu$-th data point, $\eta$ is the learning rate, and $\lambda$ is the regularization parameter. Importantly, in this framework only the column space of $B$, or in our case $Q$ are updated.

However, adopting this learning framework is insufficient when $P$ remains fixed and randomly initialized (Figure 2c). Since $P$ can project onto at most $r$ unique directions of the error, it may not align with the relevant error subspace. In this case, the network may converge to an incorrect solution, as indicated by the non-uniqueness of solutions to (6) when $m > r$.

Ideally, we want $P$ to be an orthogonal matrix whose $r$ columns span the top $r$ principal directions of the output-output correlation matrix $\Sigma^{yy}$. This alignment ensures that the most significant components of the error are propagated back through the network.

In cases where the output correlations are unknown, we can update $P$ using a modified Oja learning rule [18]:

$$\Delta P^\mu = \eta P \boldsymbol{y}^\mu \boldsymbol{y}^{\mu T} (I - P^T P) - \lambda P, \tag{8}$$

where $I$ is the identity matrix. This rule adjusts $P$ incrementally so that its columns

converge to the top $r$ principal components of the outputs $\{\boldsymbol{y}^\mu\}$.

By training both $Q$ and $P$, we allow the feedback matrix $B = QP$ to adaptively align with the relevant error directions, enabling the network to learn the correct mappings even under constrained feedback dimensionality.

Repeating the linear analysis that led to (5), we extend the derivation to our case with adaptive feedback weights. We define the transformed feedback matrix as $P = \bar{P}U^T$. This transformation aligns the feedback matrix $P$ with the principal components of the data, simplifying the analysis.

Taking the continuous-time limit (with $\eta \to 0$ and $\eta p = \tau$), we obtain a set of differential equations that describe the learning dynamics of the forward and backward weights:

$$
\begin{aligned}
\tau \frac{d\bar{W}_1}{dt} &= \bar{B}^T(S - \bar{W}_2\bar{W}_1), \\
\tau \frac{d\bar{W}_2}{dt} &= (S - \bar{W}_2\bar{W}_1)\bar{W}_1^T - \lambda\bar{W}_2,
\end{aligned}
\tag{9}
$$

and

$$
\begin{aligned}
\tau \frac{dQ}{dt} &= \bar{W}_1(S - \bar{W}_2\bar{W}_1)^T\bar{P}^T - \lambda Q, \\
\tau \frac{d\bar{P}}{dt} &= \bar{P}SS^T(I - \bar{P}^T\bar{P}) - \lambda\bar{P}.
\end{aligned}
\tag{10}
$$

Here, $\bar{B} = Q\bar{P}$ and $\lambda$ is the regularization parameter. Note that $Q$ is not affected by rotation, as it does not come in contact with either the input or the output. The full derivation can be found in the Appendix.

Notably, the updates to the feedback weights are local and follow learning rules that were well-studied in theoretical neuroscience [18–21]. Local plasticity makes our framework an attractive alternative for backpropagation in models of brain circuits.

We refer to this learning framework as *Restricted Adaptive Feedback* (RAF). Figure 2 compares the learning dynamics of RAF with those of BP and FA where only the $Q$ matrix is learned. The results demonstrate that RAF effectively aligns the feedback weights, enabling the network to converge to the correct solution despite the constrained error dimensionality. Furthermore, RAF will learn the top $r$ components of $\Sigma$, even if $r < d$ (Figure 2e).

## 3.3 Restricted adaptive feedback in deep architectures

Our weight-update equations can be naturally extended to deeper networks. In the single hidden layer model above, the backward weights $P$ were updated using the true labels $\boldsymbol{y}^\mu$. However, in deeper models, hidden layers do not have ground-truth representations. Instead, each layer relies on the local error signal, which propagates through the network according to $\boldsymbol{\delta}_l^\mu = Q_l P_l \boldsymbol{\delta}_{l+1}^\mu$. This error signal $\boldsymbol{\delta}_l^\mu$ provides the necessary information for

learning at the layer $l$.

To update the feedback weights $P_l$ in the absence of ground-truth representations, we use local error signals $\boldsymbol{\delta}_{l+1}^{\mu}$. As in the single-layer model, we use Oja's rule to adjust $P_l$ to span the principal components of the error at the next layer, $\boldsymbol{\delta}_{l+1}^{\mu}$, ensuring efficient error propagation.

The complete update rules for layer $l$ are given by

$$\Delta W_l^{\mu} = \eta Q_l P_l \boldsymbol{\delta}_{l+1}^{\mu} \boldsymbol{h}_l^{\mu T} - \lambda W_l \qquad \text{and} \qquad \begin{aligned} \Delta Q_l^{\mu} &= \eta \boldsymbol{h}_l^{\mu} P_l \boldsymbol{\delta}_{l+1}^{\mu} - \lambda Q_l, \\ \Delta P_l^{\mu} &= \eta P_l \boldsymbol{\delta}_{l+1}^{\mu} \boldsymbol{\delta}_{l+1}^{\mu T} (I - P_l^T P_l) - \lambda P, \end{aligned} \tag{11}$$

where $\eta$ is the learning rate, $\lambda$ is the regularization parameter, $\boldsymbol{h}_l^{\mu}$ is the activation of layer $l$, and $\boldsymbol{\delta}_{l+1}^{\mu}$ is the error signal from the next layer.

By updating $P_l$ using the error signals, we ensure that the feedback weights of each layer are adapted to capture the most relevant directions in the error space, facilitating effective learning throughout the network. Notably, while we use the same learning rate $\eta$ and weight decay $\lambda$ for all components $\{W_l\}$, $\{Q_l\}$, and $\{P_l\}$, it can potentially differ.

In the interest of brevity, we omit simulations of deep linear networks, as the extension from the single-layer case is straightforward. Instead, we proceed directly to deep nonlinear networks, where the impact of constrained error feedback presents more complex and interesting dynamics.

## 4.  Nonlinear networks and complex architectures

Adapting our Restricted Adaptive Feedback (RAF) framework to nonlinear networks is straightforward because the core principles of local learning and constrained error feedback remain applicable. The local update rules in eq. (11) remain the same; the primary difference lies in the introduction of nonlinear activation functions during the propagation of signals and errors. Specifically, the forward and backward passes are modified as follows:

$$h_l = f(W_l h_{l-1}), \quad \delta_l = Q_l P_l \delta_{l+1} \odot f'(h_l), \tag{12}$$

where $f$ is the nonlinear activation function applied element-wise, and $f'$ is its derivative.

To test whether our learning rule extends effectively from linear to nonlinear models, we trained deep networks on the CIFAR-10 dataset. We used a simple nonlinear model with four fully connected layers of 512 ReLU neurons each. While not state-of-the-art, this model provides a suitable testbed for evaluating our theory's applicability to nonlinear architectures and complex data.

To isolate the impact of feedback dimensionality, we applied the RAF rule in eq. (11),

constraining the rank of feedback matrices in one specific layer at a time while leaving the others unrestricted. We then measured the network's test accuracy and compared it to backpropagation (BP) as a baseline. Figure 3a shows the accuracy as a function of the feedback rank $r_l$, with each curve representing a different constrained layer.

Consistent with our findings for the linear model, constraining the feedback rank to $r = d = 10$, in any layer, match BP performance (Figure 3a), where $d$ is the number of classes. Interestingly, the shallower layers performed well even under tighter rank constraints $(r_l < d)$, suggesting that the deeper layers compensate for the limited feedback in the earlier layers by effectively adjusting their weights.

This compensatory effect is possible because no information is lost during the feedforward pass, unlike in a bottlenecked network. To demonstrate that the network still utilizes high-dimensional representations, we compared the performance of RAF-trained networks with constrained feedback to narrower networks without rank restrictions (Figure 3b). The results confirm that RAF-trained networks leverage their width to maintain high performance despite feedback constraints.

Moreover, as shown in Figure 3b, nonlinear networks trained with low-rank feedback in all layers using RAF can still match BP performance, demonstrating RAF's robustness even with dimensional constraints across the entire network.

**Task dimensionality determines the minimal rank**  Our linear analysis suggests that the error signal dimensionality needed for effective learning is tied to the loss gradient dimensionality, which depends on the number of classes in the data. To test this, we trained networks on subsets of CIFAR-100 with 50, 75, and 100 classes, constraining the ranks of all feedback matrices to $r$ (Figure 3c).

The results show that network performance matches BP when the feedback rank equals the number of classes $(r = d)$. It indicates that the minimal rank required for effective learning aligns with the task's complexity, defined by the number of output classes.

## Low-dimensional error signal from the top layer

Our theory shows how to propagate error signals from deeper layers to shallower ones using restricted adaptive feedback. However, error projections can also bypass intermediate layers entirely, leading to different variants of *Direct Restricted Adaptive Feedback* (dRAF), analogous to Direct Feedback Alignment [22]. For example, we can make direct connections from the output or the penultimate layer to earlier layers, training these connections using our algorithm (Figure 3d). As with RAF, this direct projection method matches the performance of BP when the rank of the feedback matrices satisfies $r \geq d$. This model is particularly important because it is more flexible and has greater potential to explain learning in the brain, where error signals may arrive from different pathways.
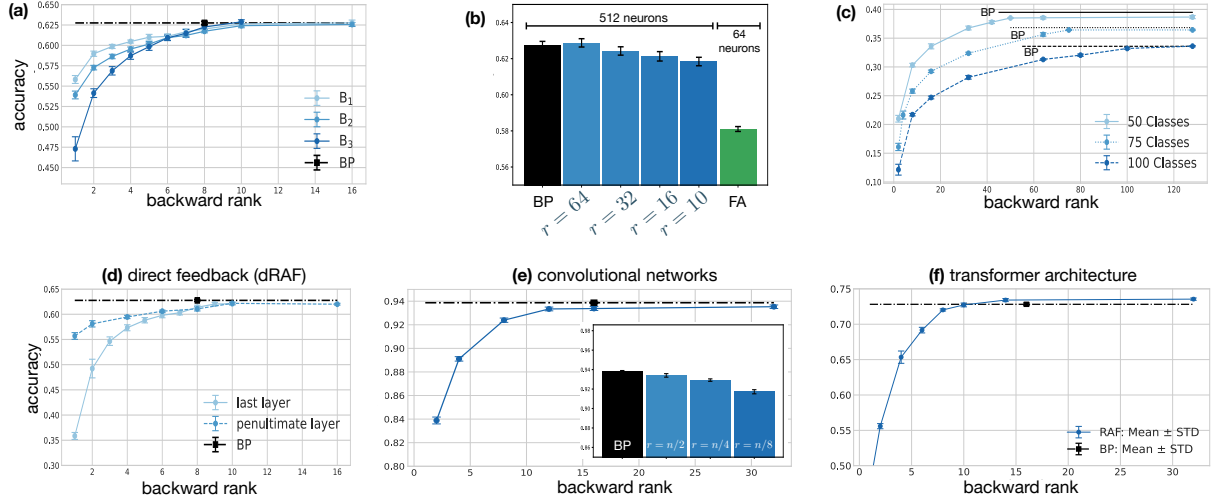
Figure 3: *Restricted Adaptive Feedback (RAF) efficiently trains nonlinear networks on CIFAR-10.* **(a)** Constraining feedback in any layer to $r = d = 10$ matches full BP performance. Different curves ($B_1$, $B_2$, and $B_3$) correspond to restricting feedback at different layers. **(b)** Reducing feedback dimensionality has minimal impact on performance, whereas reducing network width significantly degrades accuracy, indicating that high-dimensional representations are still utilized. **(c)** Minimal error dimensionality is sufficient for learning, as shown by subsampling classes from CIFAR-100. All feedback matrices are constrained to the same rank $r$. **(d)** Direct RAF: all layers receive a low-dimensional error either directly from the output layer (solid) or from the penultimate layer (dashed), both converging to BP-level performance. **(e)** Performance of a 4-block VGG-like convolutional network trained with RAF, constraining the layers with 512 channels. Inset: Training the same network with all layers constrained to a fraction of their size; the smallest layer has 64 channels. **(f)** A transformer-based architecture trained to classify images, demonstrating that RAF generalizes to more advanced network architectures. See Supplemental Information for details on all implementations.

## Convolutional neural networks

Our previous results demonstrate that fully connected networks can learn effectively from minimal error signals, even on complex datasets, matching backpropagation (BP) performance. Here, we extend this investigation to convolutional architectures.

Training convolutional networks with Feedback Alignment (FA) is notoriously difficult [8, 23]. Recent work has made progress by learning feedback weights within the FA framework [24], but our approach differs by restricting the error signal's dimensionality using Restricted Adaptive Feedback (RAF). We aim to determine whether convolutional networks can also benefit from low-dimensional error feedback.

We trained a VGG-like convolutional network [25] with four blocks and batch normalization on the CIFAR-10 dataset (see Appendix for details). Using RAF, we decoupled the error propagation from the feedforward pass in all layers. Initially, we constrained the feedback error only in the blocks containing 512 (Figure 3e). Consistent with our findings in fully connected networks, convolutional networks learn well with a feedback matrix with a rank similar to the number of classes, $d = 10$.

To test whether convolutional networks can train when constraining *all* feedback paths, we further restricted each block to have feedback matrices with ranks equal to 1/2, 1/4, or 1/8 of the block width (Figure 3e inset). Our results indicate that reducing the error dimensionality has a minimal impact on performance, except in the most extreme case. Specifically, constraining the feedback rank to 1/8 of the block width resulted in a noticeable drop in performance. This finding aligns with our previous results, as the layer with 64 channels received feedback with a rank of $r = 8$. Overall, our results demonstrate that convolutional networks can be efficiently trained using a minimal error signal.

## Transformers

Transformers have become a dominant architecture in natural language processing and, more recently, have achieved state-of-the-art performance in image classification tasks [26, 27]. Unlike convolutional and standard feedforward architectures, transformers rely on a self-attention mechanism that computes interactions between tokens through the multiplication of key, query, and value matrices [28]. These operations introduce complex, nonlinear dependencies that make it unclear whether low-dimensional error feedback, as proposed in our framework, can effectively support learning in transformer models.

To investigate this question, we applied our approach to visual transformer architectures trained for image recognition on the CIFAR-10 dataset. Each weight matrix in the transformer architecture, including key, query, and value matrices, was trained using low-dimensional error feedback from the subsequent layer. Feedback matrices were trained using the RAF framework (see the appendix for details). The network's performance as a function of the feedback rank, which was kept identical for all matrices, is shown in

Figure 3f.

Our results indicate that, consistent with our findings in both nonlinear and convolutional networks, an error signal with dimensionality as low as the task dimensionality (i.e. the number of output classes) is sufficient to train transformer architectures effectively (Figure 3f). Despite the intricate nature of self-attention computations, training with low-dimensional feedback yielded performance comparable to conventional backpropagation on CIFAR-10 with rank $r = 10$. Notably, training transformers with low-dimensional feedback using RAF not only matched but, in some cases, consistently outperformed backpropagation by a small margin. This improvement was observed even when hyperparameters were independently optimized for both BP and RAF. The improved generalization could be understood as a regularization effect achieved by constraining the feedback. However, since we observe it only in the transformer architecture, it is not a general property of a low-dimensional error signal, and we leave further investigation to a future study.

## 5. Error dimensionality shapes neural receptive fields

We have shown that neural networks can be efficiently trained using minimal error signals comparable to task dimensionality. Here, we investigate how error dimensionality affects neural representations, providing insights into receptive fields observed in the brain.

Lindsey et al. found that in convolutional models of the visual system, narrow feedforward bottlenecks between the retina and the brain led to center-surround receptive fields in the retinal layer, similar to mammals [29]. Wider bottlenecks resulted in orientation-selective receptive fields, as seen in salamanders. We hypothesize that these effects are due to constraints on the error signal reaching the retina, rather than the feedforward bottlenecks themselves.

To test this hypothesis, we trained a model similar to that of Lindsey et. al [29] but with full-width layers throughout. Instead of constraining the feedforward pathway, we constrained only the error signal using a low-rank feedback matrix trained with RAF (Figure 4c), thereby isolating the impact of error dimensionality on neural representations. We extracted the receptive fields of neurons in the retinal layer using visualization techniques [30] (Figure 4d).

Consistent with our hypothesis, constraining the feedback rank led to the emergence of center-surround receptive fields in the retinal layer. To further validate our hypothesis, we trained a model that included the feedforward bottleneck, similar to [29] but used dRAF to train the retinal layer without restricting the backward pathway. In line with our expectations, the retinal receptive fields exhibited orientation selectivity (see the appendix).
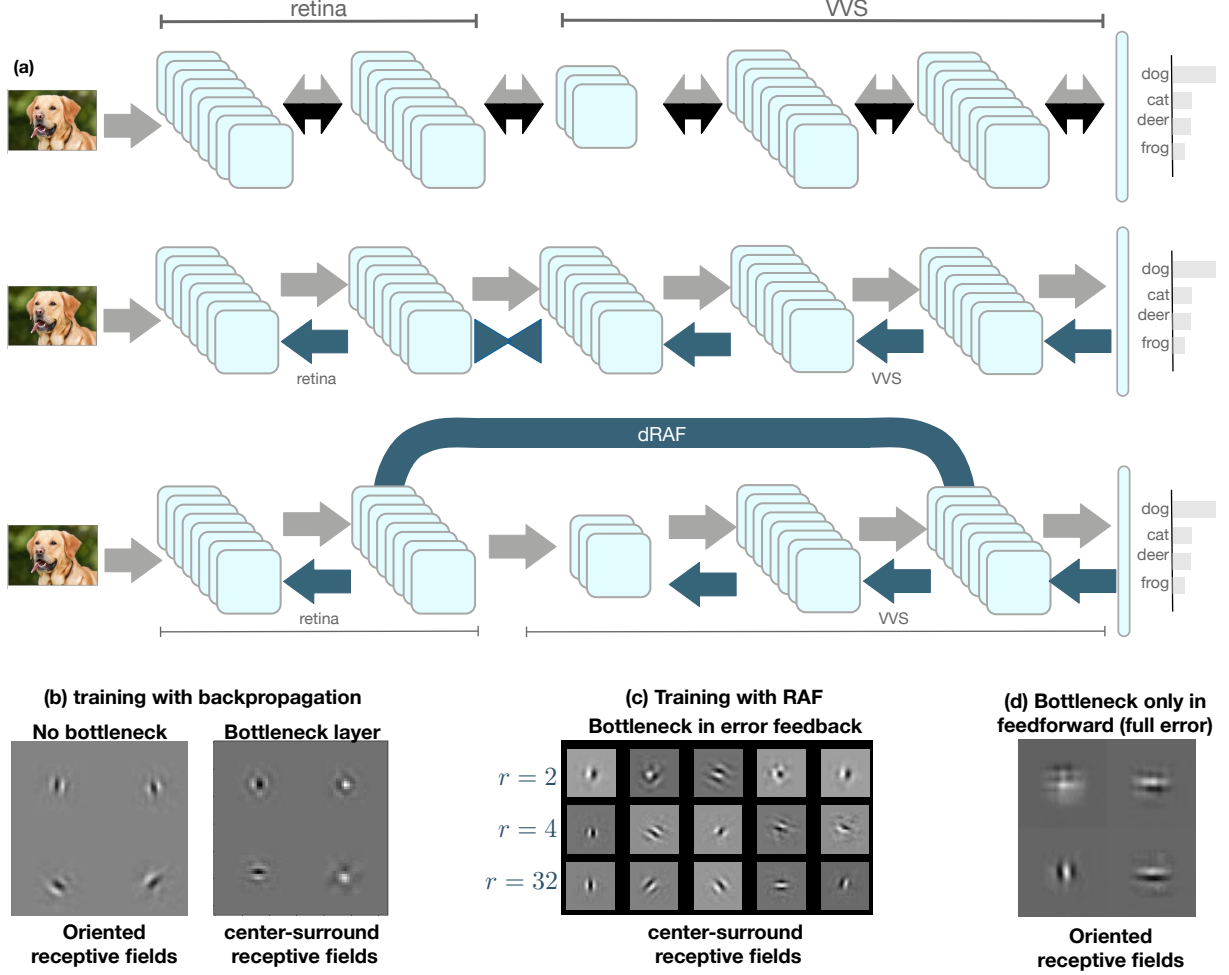
Figure 4: *Receptive fields (RFs) are shaped by feedback dimensionality rather than feedforward constraints.* **(a)** Model of the early ventral visual stream. Top: Retinal-to-cortical connectivity is modeled as a bottleneck in the network; adapted from [29]. Middle: A bottleneck is applied only to the feedback connections, restricting error dimensionality. Bottom: Using direct RAF (dRAF), retinal convolutional layers receive error signals from the penultimate layer, enabling full error feedback even with a feedforward bottleneck. **(b)** Receptive fields (RFs) in a model trained with BP, corresponding to the setup in (a) top. Left: No bottleneck. Right: Bottleneck at the retinal output, which induces center-surround RFs. **(c)** Training the retinal output with RAF. Restricting the rank of the feedback matrix produces center-surround RFs. **(d)** RFs in a network with a bottleneck at the retinal output but full-dimensional error feedback via dRAF. The presence of high-dimensional error signals results in oriented RFs.

This experiment demonstrates that error dimensionality influences neuronal tuning and neural representations. Specifically, lower-dimensional error signals promote higher symmetries in emergent receptive fields. Our findings underscore the importance of considering the dimensionality and pathways of error signals when studying neural computations in the brain.

# 6. Discussion

Our work demonstrates that neural networks can be trained effectively using minimal error signals constrained to the task's intrinsic dimensionality rather than the higher dimensionality of the network's representations. By adopting a factorized version of Feedback Alignment with low-rank matrices and training both left and right spaces of the feedback matrix, we showed that deep networks—linear, nonlinear, and convolutional—can match the performance of full backpropagation even under stringent error-dimensionality constraints. This finding highlights that the essential information required for learning is **tied to the complexity of the task**, measured by the number of output classes, and **does not scale with the number of parameters or size of the model**. Additionally, we revealed that constraining error dimensionality influences neural representations, providing a potential explanation for biological phenomena, such as center-surround receptive fields in the retina.

This work aims to explore potential mechanisms for implementing gradient descent in the brain. Recognizing that high-dimensional feedback is not necessary for effective learning is a significant step toward developing more flexible and biologically realistic models of learning. This insight suggests that the brain might utilize low-dimensional error signals to drive learning processes, aligning with anatomical and physiological constraints.

**Bridging biological constraints and gradient-based learning.** Our work aims to reconcile gradient-based learning with the anatomical constraints of cortical circuits. In the cortex, feedback connections are generally sparser, less structured, and more selective than their feedforward counterparts [31], suggesting that learning may rely on a restricted top-down error signal. Several alternative frameworks have been proposed to explain learning in hierarchical biological circuits. Predictive coding models suggest that the cortex continuously generates predictions about sensory input and computes discrepancies, or prediction errors, between expected and actual signals [32, 33] . Equilibrium propagation [34] offers another biologically inspired alternative, demonstrating how energy-based models can approximate backpropagation through local weight updates. While these approaches provide valuable insights into neural computation, they deviate from the conventional framework of gradient-based optimization. In contrast, our approach remains firmly within the gradient descent paradigm, allowing direct loss minimization. This key property enables its seamless application across diverse architectures, from fully

connected nonlinear networks to convolutional models and transformer-based attention mechanisms.

**Between reward-based learning and full gradient descent.**   By showing that deep networks can be trained with low-dimensional error signals, our framework bridges the gap between gradient-based learning and reward-based learning. Studies have suggested that dopaminergic neurons in the ventral tegmental area (VTA) encode reward signals that capture multiple task attributes rather than a single scalar value [35–37]. Our findings provide a theoretical foundation for moving beyond scalar reward-based learning toward a framework where low-dimensional reward signals carry richer information about the underlying loss (or gain). This perspective positions our model as an intermediary between reinforcement learning's scalar feedback and the full backpropagation of error gradients. More broadly, our results suggest that learning could emerge through the strategic compression of error signals into task-relevant representations, offering a biologically plausible alternative to high-dimensional gradient propagation.

**Aligning to the error space.**   On its surface, our novel learning rule is similar to previous Feedback Alignment (FA) schemes. However, it offers a conceptual novelty. In traditional FA, learning the feedback weights aims to align them with the feedforward weights [5, 10], mirroring full backpropagation. In contrast, by factorizing the feedback matrix $B = QP$, we also **align the row space of the feedback matrix with the source of the error**, thereby improving the quality of the error signal itself.

**Computational complexity.**   Despite the additional requirement of training the feedback matrices, the computational cost of using low-rank feedback does not increase significantly. The updates for the low-rank matrices are linear in the number of neurons or filters per layer, while the forward pass and the primary weight updates remain quadratic [38]. This efficiency makes our approach practical for large-scale networks and suggests that similar mechanisms could be feasible in biological neural networks.

**Gradient descent in high dimensions.**   Our findings invite a rethinking of gradient descent dynamics in large, overparameterized networks. Typically, the weight dynamics during training are high-dimensional [39]. However, when the error signal is low-dimensional, the weight updates in each layer are confined to a much lower-dimensional subspace. This constraint could have implications for understanding the generalization capabilities of neural networks, as it suggests that effective learning does not require exploring the full parameter space. Indeed, previous studies found that the Hessian of the loss function shows low-rank structure [40]. Exploring this connection could open new avenues for understanding the dynamics of gradient descent in high-dimensional loss spaces [41].

**Exploring the potential of low-dimensional feedback in more complex architectures.** This approach not only improves learning efficiency but may also act as a form of regularization, as can be seen from the improved performance of low-rank feedback matrices on visual transformers—the possibility of using restricted error signals as a form of regularization warrants further investigation.

In summary, our findings not only challenge the conventional reliance on high-dimensional error propagation but also pave the way for the development of biologically inspired, computationally efficient learning algorithms that hold promise for advancing both artificial intelligence and our understanding of brain function." This would leave the reader with a clear impression of both the scientific and translational impact of your work

## Acknowledgments

## References

[1] Rumelhart DE, Hinton GE, Williams RJ (1986) Learning representations by back-propagating errors. *nature* 323(6088):533–536.

[2] Chinta LV, Tweed DB (2012) Adaptive optimal control without weight transport. *Neural computation* 24(6):1487–1518.

[3] Stork (1989) Is backpropagation biologically plausible? in *International 1989 Joint Conference on Neural Networks*. (IEEE), pp. 241–246.

[4] Lillicrap TP, Santoro A, Marris L, Akerman CJ, Hinton G (2020) Backpropagation and the brain. *Nature Reviews Neuroscience* 21(6):335–346.

[5] Lillicrap TP, Cownden D, Tweed DB, Akerman CJ (2016) Random synaptic feedback weights support error backpropagation for deep learning. *Nature communications* 7(1):13276.

[6] Grossberg S (1987) Competitive learning: From interactive activation to adaptive resonance. *Cognitive science* 11(1):23–63.

[7] Crick F (1989) The recent excitement about neural networks. *Nature* 337(6203):129–132.

[8] Bartunov S, et al. (2018) Assessing the scalability of biologically-motivated deep learning algorithms and architectures. *Advances in neural information processing systems* 31.

[9] Kolen JF, Pollack JB (1994) Backpropagation without weight transport in *Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94)*. (IEEE), Vol. 3, pp. 1375–1380.

[10] Akrout M, Wilson C, Humphreys P, Lillicrap T, Tweed DB (2019) Deep learning without weight transport. *Advances in neural information processing systems* 32.

[11] Crafton B, Parihar A, Gebhardt E, Raychowdhury A (2019) Direct feedback alignment with sparse connections for local learning. *Frontiers in neuroscience* 13:450947.

[12] Liao Q, Leibo J, Poggio T (2016) How important is weight symmetry in backpropagation? in *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 30.

[13] Gunasekar S, Lee JD, Soudry D, Srebro N (2018) Implicit bias of gradient descent on linear convolutional networks. *Advances in neural information processing systems* 31.

[14] Caro JO, et al. (2024) Translational symmetry in convolutions with localized kernels causes an implicit bias toward high frequency adversarial examples. *Frontiers in Computational Neuroscience* 18.

[15] Patel N, Shwartz-Ziv R (2024) Learning to compress: Local rank and information compression in deep neural networks. *arXiv preprint arXiv:2410.07687*.

[16] Saxe AM, McClelland JL, Ganguli S (2013) Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*.

[17] Johnson WB (1984) Extensions of lipshitz mapping into hilbert space in *Conference modern analysis and probability, 1984*. pp. 189–206.

[18] Oja E (1982) Simplified neuron model as a principal component analyzer. *Journal of mathematical biology* 15:267–273.

[19] Clopath C, Büsing L, Vasilaki E, Gerstner W (2010) Connectivity reflects coding: a model of voltage-based stdp with homeostasis. *Nature neuroscience* 13(3):344–352.

[20] Turrigiano GG (2008) The self-tuning neuron: synaptic scaling of excitatory synapses. *Cell* 135(3):422–435.

[21] Pehlevan C, Hu T, Chklovskii DB (2015) A hebbian/anti-hebbian neural network for linear subspace learning: A derivation from multidimensional scaling of streaming data. *Neural computation* 27(7):1461–1495.

[22] Nøkland A (2016) Direct feedback alignment provides learning in deep neural networks. *Advances in neural information processing systems* 29.

[23] Launay J, Poli I, Krzakala F (2019) Principled training of neural networks with direct feedback alignment. *arXiv preprint arXiv:1906.04554.*

[24] Bacho F, Chu D (2024) Low-variance forward gradients using direct feedback alignment and momentum. *Neural Networks* 169:572–583.

[25] Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556.*

[26] Dosovitskiy A, et al. (2020) An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929.*

[27] Khan S, et al. (2022) Transformers in vision: A survey. *ACM computing surveys (CSUR)* 54(10s):1–41.

[28] Vaswani A, et al. (2017) Attention is all you need. *Advances in neural information processing systems* 30.

[29] Lindsey J, Ocko SA, Ganguli S, Deny S (2019) A unified theory of early visual representations from retina to cortex through anatomically constrained deep cnns. *arXiv preprint arXiv:1901.00945.*

[30] Erhan D, Bengio Y, Courville A, Vincent P (2009) Visualizing higher-layer features of a deep network. *University of Montreal* 1341(3):1.

[31] Markov NT, et al. (2014) Anatomy of hierarchy: feedforward and feedback pathways in macaque visual cortex. *Journal of comparative neurology* 522(1):225–259.

[32] Rao RP, Ballard DH (1999) Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature neuroscience* 2(1):79–87.

[33] Whittington JC, Bogacz R (2017) An approximation of the error backpropagation algorithm in a predictive coding network with local hebbian synaptic plasticity. *Neural computation* 29(5):1229–1262.

[34] Scellier B, Bengio Y (2017) Equilibrium propagation: Bridging the gap between energy-based models and backpropagation. *Frontiers in computational neuroscience* 11:24.

[35] Cohen JY, Haesler S, Vong L, Lowell BB, Uchida N (2012) Neuron-type-specific signals for reward and punishment in the ventral tegmental area. *nature* 482(7383):85–88.

[36] Parker NF, et al. (2016) Reward and choice encoding in terminals of midbrain dopamine neurons depends on striatal target. *Nature neuroscience* 19(6):845–854.

[37] Engelhard B, et al. (2019) Specialized coding of sensory, motor and cognitive variables in vta dopamine neurons. *Nature* 570(7762):509–513.

[38] Livni R, Shalev-Shwartz S, Shamir O (2014) On the computational efficiency of training neural networks. *Advances in neural information processing systems* 27.

[39] Jacot A, Gabriel F, Hongler C (2018) Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems* 31.

[40] Sagun L, Bottou L, LeCun Y (2017) Empirical analysis of the hessian of over-parameterized neural networks in *Proceedings of the International Conference on Learning Representations (ICLR)*.

[41] Zhang C, Bengio S, Hardt M, Recht B, Vinyals O (2017) Understanding deep learning requires rethinking generalization in *Proceedings of the International Conference on Learning Representations (ICLR)*.

# Appendix

## A. Full linear theory

in this section, we provide the analysis of the linear networks for both Feedback Alignment and RAF

### Detailed analysis of Feedback Alignment learning dynamics

From eq. (4) in the main text, the weight updates for a single sample $\mu$ are given by:

$$\begin{aligned}
\Delta W_1^\mu &= \eta B \left( \boldsymbol{y}^\mu - W_2 W_1 \boldsymbol{x}^\mu \right) \boldsymbol{x}^{\mu\top}, \\
\Delta W_2^\mu &= \eta \left( \boldsymbol{y}^\mu - W_2 W_1 \boldsymbol{x}^\mu \right) \boldsymbol{x}^{\mu\top} W_1^\top,
\end{aligned} \tag{A1}$$

where:

- $\eta$ is the learning rate,

- $\boldsymbol{x}^\mu \in \mathbb{R}^n$ is the input vector for sample $\mu$,

- $\boldsymbol{y}^\mu \in \mathbb{R}^m$ is the corresponding target output,

- $W_1 \in \mathbb{R}^{k \times n}$ and $W_2 \in \mathbb{R}^{m \times k}$ are the weight matrices,

- $B \in \mathbb{R}^{k \times m}$ is a predefined matrix (e.g., a feedback or scaling matrix).

We introduce the empirical covariance matrices:

$$\Sigma_{io} = \frac{1}{p} \sum_{\mu=1}^{p} \boldsymbol{y}^\mu \boldsymbol{x}^{\mu\top}, \tag{A2}$$

$$\Sigma_{oo} = \frac{1}{p} \sum_{\mu=1}^{p} \boldsymbol{y}^\mu \boldsymbol{y}^{\mu\top}, \tag{A3}$$

$$\Sigma_{ii} = \frac{1}{p} \sum_{\mu=1}^{p} \boldsymbol{x}^\mu \boldsymbol{x}^{\mu\top} = I, \tag{A4}$$

where $\Sigma_{ii} = I$ assumes that the input vectors are whitened (i.e., have unit covariance). Summing over all $p$ training examples, we obtain the average weight updates:

$$\begin{aligned}
\Delta W_1 &= \frac{\eta}{p} \sum_{\mu=1}^{p} \Delta W_1^\mu \\
&= \frac{\eta}{p} \sum_{\mu=1}^{p} B \left( \boldsymbol{y}^\mu - W_2 W_1 \boldsymbol{x}^\mu \right) \boldsymbol{x}^{\mu\top}.
\end{aligned} \tag{A5}$$

Using these definitions, the update for $W_1$ simplifies to:

$$\Delta W_1 = \eta B \left( \Sigma_{io} - W_2 W_1 \Sigma_{ii} \right) \quad (A6)$$

$$= \eta B \left( \Sigma_{io} - W_2 W_1 \right). \quad (A7)$$

Under the limit as $\eta \to 0$ with $\eta = \frac{dt}{\tau}$ (where $\tau$ is a time constant), we transition from discrete updates to continuous-time dynamics:

$$\tau \frac{dW_1}{dt} = B \left( \Sigma_{io} - W_2 W_1 \right), \quad (A8)$$

which matches the weight dynamics presented in eq. (5) of the main text.

Similarly, the update for $W_2$ becomes:

$$\Delta W_2 = \frac{\eta}{p} \sum_{\mu=1}^{p} \Delta W_2^{\mu} \quad (A9)$$

$$= \frac{\eta}{p} \sum_{\mu=1}^{p} \left( \boldsymbol{y}^{\mu} - W_2 W_1 \boldsymbol{x}^{\mu} \right) \boldsymbol{x}^{\mu \top} W_1^{\top} \quad (A10)$$

$$= \eta \left( \Sigma_{io} - W_2 W_1 \right) W_1^{\top}, \quad (A11)$$

which simplifies under the same limit to:

$$\tau \frac{dW_2}{dt} = \left( \Sigma_{io} - W_2 W_1 \right) W_1^{\top}. \quad (A12)$$

Eq. (A8) and eq. (A12) describe the continuous-time dynamics of the weights $W_1$ and $W_2$ under the given learning rule.

We consider the singular value decomposition (SVD) of the covariance matrix $\Sigma_{io}$:

$$\Sigma_{io} = U S V^{\top}, \quad (A13)$$

where:

- $U \in \mathbb{R}^{m \times d}$ and $V \in \mathbb{R}^{n \times d}$ are matrices with orthonormal columns,

- $S \in \mathbb{R}^{d \times d}$ is a diagonal matrix containing the singular values,

- $d$ is the rank of $\Sigma_{io}$.

We perform a rotation of the weight matrices and $B$ as follows:

$$\begin{aligned} W_1 &= \bar{W}_1 V^{\top}, \\ W_2 &= U \bar{W}_2, \\ B &= \bar{B} U^{\top}. \end{aligned} \quad (A14)$$

Substituting these into the previous weight dynamics, we have for $W_1$:

$$\tau \frac{dW_1}{dt} = B\left(\Sigma_{io} - W_2 W_1\right) \tag{A15}$$

$$= \bar{B} U^\top \left(USV^\top - U\bar{W}_2\bar{W}_1 V^\top\right). \tag{A16}$$

Since $U^\top U = I$ (due to orthonormal columns of $U$), we can simplify:

$$\tau \frac{dW_1}{dt} = \bar{B}\left(U^\top U\right)\left(S - \bar{W}_2\bar{W}_1\right)V^\top \tag{A17}$$

$$= \bar{B}\left(S - \bar{W}_2\bar{W}_1\right)V^\top. \tag{A18}$$

Recognizing that $W_1 = \bar{W}_1 V^\top$, we can write $\frac{dW_1}{dt} = \frac{d\bar{W}_1}{dt}V^\top$. Thus, multiplying both sides on the right by $V$ (since $V^\top V = I$):

$$\tau \frac{d\bar{W}_1}{dt} = \bar{B}\left(S - \bar{W}_2\bar{W}_1\right). \tag{A19}$$

Similarly, for $W_2$:

$$\tau \frac{dW_2}{dt} = \left(\Sigma_{io} - W_2 W_1\right)W_1^\top \tag{A20}$$

$$= \left(USV\top - U\bar{W}_2\bar{W}_1 V^\top\right)\left(\bar{W}_1 V^\top\right)^\top \tag{A21}$$

$$\tag{A22}$$

Using the fact that $V^\top V = I$ we simplify:

$$\tau \frac{dW_2}{dt} = \left(US - U\bar{W}_2\bar{W}_1\right)\bar{W}_1^\top \tag{A23}$$

$$\tag{A24}$$

Since $W_2 = U\bar{W}_2$, we have $\frac{dW_2}{dt} = U\frac{d\bar{W}_2}{dt}$. Multiplying both sides on the left by $U^\top$:

$$\tau \frac{d\bar{W}_2}{dt} = \left(S - \bar{W}_2\bar{W}_1\right)\bar{W}_1^\top. \tag{A25}$$

Eqs. (A19) and (A25) describe the dynamics of the rotated weights $\bar{W}_1$ and $\bar{W}_2$:

$$\begin{aligned} \tau \frac{d\bar{W}_1}{dt} &= \bar{B}\left(S - \bar{W}_2\bar{W}_1\right), \\ \tau \frac{d\bar{W}_2}{dt} &= \left(S - \bar{W}_2\bar{W}_1\right)\bar{W}_1^\top. \end{aligned} \tag{A26}$$

## Full derivation of Restricted Adaptive Feedback (RAF) learning dynamics

We consider an alternative algorithm where the matrix $B$ is replaced by $B = QP$, with $Q \in \mathbb{R}^{k \times r}$ and $P \in \mathbb{R}^{r \times m}$.

The updates for $Q$ and $P$ for a single sample $\mu$ are given by:

$$\begin{aligned}
\Delta Q^\mu &= \eta W_1 \boldsymbol{x}^\mu P \left( \boldsymbol{y}^\mu - W_2 W_1 \boldsymbol{x}^\mu \right), \\
\Delta P^\mu &= \eta P \boldsymbol{y}^\mu \boldsymbol{y}^{\mu\top} \left( I - P^\top P \right).
\end{aligned} \tag{A27}$$

Summing over all $p$ training examples, we obtain the average updates:

$$\begin{aligned}
\Delta Q &= \frac{\eta}{p} \sum_{\mu=1}^{p} \Delta Q^\mu = \eta W_1 \left( \frac{1}{p} \sum_{\mu=1}^{p} \boldsymbol{x}^\mu \left( \boldsymbol{y}^\mu - W_2 W_1 \boldsymbol{x}^\mu \right)^\top P^\top \right), \\
\Delta P &= \frac{\eta}{p} \sum_{\mu=1}^{p} \Delta P^\mu = \eta P \left( \frac{1}{p} \sum_{\mu=1}^{p} \boldsymbol{y}^\mu \boldsymbol{y}^{\mu\top} \right) \left( I - P^\top P \right).
\end{aligned} \tag{A28}$$

We simplify the update for $Q$:

$$\Delta Q = \eta W_1 \left( \left( \frac{1}{p} \sum_{\mu=1}^{p} \boldsymbol{x}^\mu \left( \boldsymbol{y}^\mu - W_2 W_1 \boldsymbol{x}^\mu \right)^\top \right) P^\top \right) \tag{A29}$$

$$= \eta W_1 \left( \left( \Sigma_{io}^\top - W_1^\top W_2^\top \Sigma_{ii} \right) P^\top \right) \tag{A30}$$

$$= \eta W_1 \left( \left( \Sigma_{io} - W_2 W_1 \right)^\top P^\top \right) \tag{A31}$$

Similarly, the update for $P$ simplifies to:

$$\Delta P = \eta P \Sigma_{oo} \left( I - P^\top P \right). \tag{A32}$$

Thus, the updates for $Q$ and $P$ become:

$$\begin{aligned}
\Delta Q &= \eta W_1 \left( \Sigma_{io} - W_2 W_1 \right)^\top P^\top, \\
\Delta P &= \eta P \Sigma_{oo} \left( I - P^\top P \right).
\end{aligned} \tag{A33}$$

Under the continuous-time assumption, where $\eta \to 0$ with $\eta = \frac{dt}{\tau}$, the updates for $Q$ and $P$ become differential equations:

$$\tau \frac{dQ}{dt} = W_1 \left( \Sigma_{io} - W_2 W_1 \right)^\top P^\top, \tag{A34}$$

and

$$\tau \frac{dP}{dt} = P \Sigma_{oo} \left( I - P^\top P \right). \tag{A35}$$

Finally, substituting $B = QP$ into the update for $W_1$ from eq. A8, we find that the dynamics for $W_1$ become:

$$\tau \frac{dW_1}{dt} = QP\left(\Sigma_{io} - W_2 W_1\right) \tag{A36}$$

We perform rotations similar to before:

$$\begin{aligned}
W_1 &= \bar{W}_1 V^\top, \\
W_2 &= U\bar{W}_2, \\
P &= \bar{P} U^\top.
\end{aligned} \tag{A37}$$

Substituting for $Q$ we get:

$$\tau \frac{dQ}{dt} = \bar{W}_1 V_1^\top \left(USV^\top - U\bar{W}_2 \bar{W}_1 V^\top\right)^\top (PU^\top)^\top \tag{A38}$$

$$= \bar{W}_1 V_1^\top V \left(S - \bar{W}_2 \bar{W}_1\right)^T U^\top U \bar{P}^\top \tag{A39}$$

$$= \bar{W}_1 \left(S - \bar{W}_2 \bar{W}_1\right)^\top \bar{P}^\top \tag{A40}$$

Substituting these rotations into the previous derivations, we update the dynamics.

For $W_1$, the update equation is:

$$\tau \frac{dW_1}{dt} = B\left(\Sigma_{io} - W_2 W_1\right). \tag{A41}$$

Since $B = QP = Q\bar{P}U^\top$, $W_1 = \bar{W}_1 V^\top$, $W_2 = U\bar{W}_2$, and $\Sigma_{io} = USV^\top$, we have:

$$\tau \frac{d(\bar{W}_1 V^\top)}{dt} = Q\bar{P}U^\top \left(USV^\top - U\bar{W}_2 \bar{W}_1 V^\top\right) \tag{A42}$$

$$= Q\bar{P}\left(SV^\top - \bar{W}_2 \bar{W}_1 V^\top\right). \tag{A43}$$

Since $V^\top$ is constant, we can write:

$$\tau \frac{d\bar{W}_1}{dt} V^\top = Q\bar{P}\left(SV^\top - \bar{W}_2 \bar{W}_1 V^\top\right). \tag{A44}$$

Multiplying both sides on the right by $V$ (using $V^\top V = I$):

$$\tau \frac{d\bar{W}_1}{dt} = Q\bar{P}\left(S - \bar{W}_2 \bar{W}_1\right). \tag{A45}$$

Substituting $P = \bar{P}U^\top$ and $\Sigma_{oo} = US^2U^\top$ in (A35), we get:

$$\tau \frac{d(\bar{P}U^\top)}{dt} = \bar{P}U^\top \left(US^2U^\top\right) \left(I - U\bar{P}^\top \bar{P}U^\top\right) \tag{A46}$$

$$= \bar{P}S^2 \left(I - \bar{P}^\top \bar{P}\right) U^\top. \tag{A47}$$

Multiplying both sides on the right by $U$ (since $U^\top U = I$):

$$\tau \frac{d\bar{P}}{dt} = \bar{P}S^2 \left(I - \bar{P}^\top \bar{P}\right). \tag{A48}$$

In summary, under the rotations, the updated dynamics are:

$$
\begin{aligned}
\tau \frac{d\bar{W}_1}{dt} &= Q\bar{P}\left(S - \bar{W}_2\bar{W}_1\right), \\
\tau \frac{d\bar{W}_2}{dt} &= \left(S - \bar{W}_2\bar{W}_1\right)\bar{W}_1^\top, \\
\tau \frac{d\bar{P}}{dt} &= \bar{P}S^2\left(I - \bar{P}^\top\bar{P}\right). \\
\tau \frac{d\bar{Q}}{dt} &= \bar{W}_1\left(S - \bar{W}_2\bar{W}_1\right)^\top \bar{P}^\top
\end{aligned}
\tag{A49}
$$

## B.  Learning rule implementation

Our implementation and optimization were conducted entirely using PyTorch. For the RAF layers, we initialized all weights using Kaiming uniform initialization. We modified the backward pass by adjusting the gradients with respect to the input, ensuring they align with our proposed update rule. In the output layer, we learn the projection matrix $\boldsymbol{P}$ to capture the principal directions of the target labels $\boldsymbol{y}$. In the hidden layers, $\boldsymbol{P}$ is learned to project onto the principal directions of the error signal from the subsequent layer, represented as $\boldsymbol{\delta}_{l+1}$. Specifically, for a layer $l$ with input dimension $n$, output dimension $m$, and rank constraint $r$, we proceed as follows:

---
**Algorithm 1:** Modified Backward Pass for RAF Layer
---

**Input:** Error signal $\boldsymbol{\delta}_{l+1} \in \mathbb{R}^{b \times m}$, activations $\boldsymbol{h}_l \in \mathbb{R}^{b \times n}$, matrices $Q_l \in \mathbb{R}^{n \times r}$,
$\qquad$ $P_l \in \mathbb{R}^{r \times m}$, weight matrix $W_l \in \mathbb{R}^{m \times n}$

**Output:** Gradient w.r.t. input grad_input, updates $\Delta Q_l$, $\Delta P_l$, $\Delta W_l$

**Compute covariance matrix:**

**if** *use_targets* **then**
$\qquad$ $C \leftarrow (\boldsymbol{y} - \bar{\boldsymbol{y}})^{\top} (\boldsymbol{y} - \bar{\boldsymbol{y}})$;

**else**
$\qquad$ $C \leftarrow (\boldsymbol{\delta_{l+1}} - \boldsymbol{\delta_{l+1}^-})^{\top} (\boldsymbol{\delta_{l+1}} - \boldsymbol{\delta_{l+1}^-})$;

$\gamma \leftarrow \max (\operatorname{diag} (C))$;

**Compute 'gradient' w.r.t. input:**
grad_input $\leftarrow \left(Q_l P_l \boldsymbol{\delta}_{l+1}^{\top}\right)^{\top} \in \mathbb{R}^{b \times n}$;

**Compute updates for $Q_l$ and $P_l$:**
$\Delta Q_l \leftarrow \boldsymbol{h}_l^{\top} \left(P_l \boldsymbol{\delta}_{l+1}^{\top}\right)^{\top}$;

$\Delta P_l \leftarrow P_l \left(\dfrac{C}{\gamma}\right) (I - P_l^{\top} P_l)$;

**Compute update for $W_l$:**
$\Delta W_l \leftarrow \boldsymbol{\delta}_{l+1}^{\top} \boldsymbol{h}_l$;
---

The derivation of the dRAF (directed restricted adaptive feedback) follows the same procedure, with the key difference being that the error signal $\boldsymbol{\delta}$ does not necessarily originate from the next layer; instead, it can come from any subsequent layer.

# C. Numerical experiments

## Layer-Wise Constraints

In Figure [3.a], we present the results of training a network with four hidden layers, each containing 512 neurons, and an output layer with 10 neurons on the CIFAR-10 dataset. The network was trained using RAF, as described earlier, without rank constraints, except for one layer at a time. For comparison, we also trained the network using standard backpropagation as a baseline. All networks were trained with a batch size of 32, a learning rate of $6 \times 10^{-4}$, and weight decay of $4 \times 10^{-4}$. The Adam optimizer with AMSGrad was used, with training conducted for 160 epochs and an exponential learning rate decay factor of 0.975. Each experiment was repeated 10 times.

## Constraining All Layers

For the results shown in Figure [3.b], we trained the same network with four hidden layers, each containing 512 neurons. This time, we applied rank constraints to all layers simultaneously, with rank values $r = 64, 32, 16, 10$.

To further demonstrate that the network still utilizes high-rank representations, we also trained a variant of the network with 64 neurons in each hidden layer, without applying any rank constraints. All training was conducted with a batch size of 32, a learning rate of $6 \times 10^{-4}$, and weight decay of $4 \times 10^{-4}$. The Adam optimizer with AMSGrad was used, with training carried out for 160 epochs and an exponential learning rate decay factor of 0.975. Each experiment was repeated 10 times.

## CIFAR-100 Sub-sampling

For the results shown in Figure [3.c], we trained the same model on the CIFAR-100 dataset, sampling different numbers of classes $d$, with $d = 50, 75, 100$. For each sub-sample, we trained the model while applying rank constraints to all layers, using various rank values. This was done to demonstrate that the dimensionality of the error signal depends on the task dimensionality $d$. All optimizations were performed with a batch size of 32, a learning rate of $6 \times 10^{-4}$, and weight decay of $4 \times 10^{-4}$. The Adam optimizer with AMSGrad was used, with training conducted for 160 epochs and an exponential learning rate decay factor of 0.975. Each training run was repeated 5 times.

## dRAF

For the results shown in Figure [3.d], we trained the same model using dRAF. In one experiment, we propagated the error signal from the last layer to all preceding layers. In a separate experiment, we propagated the error signal directly from the penultimate layer to all earlier layers, applying rank constraints to these layers while keeping the last layer at full rank. All optimizations were performed with a batch size of 32, a learning rate of $6 \times 10^{-4}$, and weight decay of $4 \times 10^{-4}$. The Adam optimizer with AMSGrad was used, training for 160 epochs with an exponential learning rate decay factor of 0.975. Each training run was repeated 5 times

## Convolutional Neural Networks

To extend our RAF algorithm to convolutional layers, we apply the rank constraint to the number of channels in the error signal. This effectively constrains the dimensionality of the error signal across the spatial dimensions. Similar to the implementation for fully connected layers in Algorithm 1, we modify the backward pass in the same way for convolutional layers.

The key difference in the convolutional context is how the projection matrix $\boldsymbol{P}$ operates on the error signals. In convolutional layers, $\boldsymbol{P}$ functions as a $1 \times 1$ convolutional filter, projecting the error signal at each spatial location from $m$ channels down to $r$ channels. This reduces the error signal's dimensionality to $r$ per pixel, adhering to the rank constraint.

As in the fully connected case, the matrix $\boldsymbol{P}$ is learned using Oja's rule. However, the covariance matrix $\boldsymbol{C}$ is computed over both the batch and spatial dimensions—that is, across all pixels in all images within the batch. This approach captures the covariance structure of pixel representations more effectively, enabling $\boldsymbol{P}$ to project the error signals appropriately in the convolutional setting.

## VGG-like Architecture

We trained a VGG-like convolutional neural network comprising four convolutional blocks, as detailed in Table C1. For the experiments presented in Figure 4(a), a rank constraint was applied exclusively to the final convolutional block, which consists of 512 channels. In contrast, for the results shown in Figure 4(b), rank constraints were imposed on all layers, reducing the rank to 1/2, 1/4, and 1/8 of the original channel dimensions.

Each network was trained using the Adam optimizer with a learning rate of $5 \times 10^{-4}$, a weight decay of $5 \times 10^{-5}$, and an exponential learning rate decay factor of 0.98. Training was conducted for 250 epochs, with each experiment repeated five times to ensure robustness.

## Vision Transformer (ViT)

We implemented a Vision Transformer (ViT) model that consists of seven self-attention blocks, each with an embedding dimension of 384 and four attention heads. The input images were partitioned into nonoverlapping patches, which were linearly projected into a 384-dimensional feature space. These patch embeddings were augmented with learned positional encodings and processed through a stack of seven transformer encoder layers. Each encoder block employed multi-head self-attention with four heads, followed by pre-norm layer normalization and dropout for regularization. The final representation was obtained from a learned class token, which was subsequently processed by a fully connected layer for classification.

Training was carried out using a batch size of 32, a learning rate of $3 \times 10^{-4}$, and a weight loss of $1 \times 10^{-4}$. We used the Adam optimizer with AMSGrad, training for 150 epochs with an exponential learning rate decay factor of 0.98. Each experiment was repeated five times to ensure robustness.

## Neural Receptive Fields

To analyze neural receptive fields, we trained the same convolutional network architecture as used in Lindsey et. al. [29]. The model consists of two convolutional layers with ReLU activations to simulate retinal processing, followed by three additional convolutional layers with ReLU activations to model the ventral visual stream (VVS). Fully connected layers were used for classification, with the complete architecture detailed in Table C2. A

bottleneck was introduced between the retinal and VVS layers to constrain the information flow.

Training was performed using the Restricted Adaptive Feedback (RAF) method, where rank constraints were applied to feedback sent to the retina with ranks $r = 2, 4, 32$. The network was trained with hyperparameters similar to those in Linsdey et. al. [29], using the RMSProp optimizer with a learning rate of $1 \times 10^{-4}$, a weight decay of $1 \times 10^{-5}$, and an exponential learning rate decay factor of 0.985. Each model was trained for 120 epochs, and all experiments were repeated five times to ensure robustness.

Additionally, we conducted experiments with dynamic Restricted Adaptive Feedback (dRAF), where feedback to the retina originated from higher visual layers while feedforward connections within the retina were constrained to four channels.

Table C1: CNN Architecture

| Layer(s) | Output Size | Details |
|---|---|---|
| Input | $32 \times 32 \times 3$ | |
| **Convolutional Block 1** | | |
| Conv2D + ReLU | $32 \times 32 \times 64$ | $3 \times 3$ conv, 64 filters, padding=1 |
| BatchNorm2D | $32 \times 32 \times 64$ | |
| Conv2D + ReLU | $32 \times 32 \times 64$ | $3 \times 3$ conv, 64 filters, padding=1 |
| BatchNorm2D | $32 \times 32 \times 64$ | |
| MaxPool2D | $16 \times 16 \times 64$ | $2 \times 2$ max pool, stride=2 |
| **Convolutional Block 2** | | |
| Conv2D + ReLU | $16 \times 16 \times 128$ | $3 \times 3$ conv, 128 filters, padding=1 |
| BatchNorm2D | $16 \times 16 \times 128$ | |
| Conv2D + ReLU | $16 \times 16 \times 128$ | $3 \times 3$ conv, 128 filters, padding=1 |
| BatchNorm2D | $16 \times 16 \times 128$ | |
| MaxPool2D | $8 \times 8 \times 128$ | $2 \times 2$ max pool, stride=2 |
| **Convolutional Block 3** | | |
| Conv2D + ReLU | $8 \times 8 \times 256$ | $3 \times 3$ conv, 256 filters, padding=1 |
| BatchNorm2D | $8 \times 8 \times 256$ | |
| Conv2D + ReLU | $8 \times 8 \times 256$ | $3 \times 3$ conv, 256 filters, padding=1 |
| BatchNorm2D | $8 \times 8 \times 256$ | |
| MaxPool2D | $4 \times 4 \times 256$ | $2 \times 2$ max pool, stride=2 |
| **Convolutional Block 4** | | |
| Conv2D + ReLU | $4 \times 4 \times 512$ | $3 \times 3$ conv, 512 filters, padding=1 |
| BatchNorm2D | $4 \times 4 \times 512$ | |
| Conv2D + ReLU | $4 \times 4 \times 512$ | $3 \times 3$ conv, 512 filters, padding=1 |
| BatchNorm2D | $4 \times 4 \times 512$ | |
| AdaptiveAvgPool2D | $1 \times 1 \times 512$ | Output size $(1, 1)$ |
| Flatten | 512 | Flatten to vector |
| **Classifier** | | |
| Fully Connected + ReLU | 256 | Linear layer, $512 \rightarrow 256$ |
| Dropout | 256 | Dropout probability $p = 0.4$ |
| Fully Connected | $C$ | Linear layer, $256 \rightarrow C$ |

Table C2: Retina model Architecture

| Layer(s) | Output Size | Details |
|---|---|---|
| Input | $32 \times 32 \times 1$ | Grayscale input |
| **Retina** | | |
| Conv2D + ReLU | $32 \times 32 \times 32$ | $9 \times 9$ conv, 32 filters, padding=4 |
| Conv2D + ReLU | $32 \times 32 \times 32$ | $9 \times 9$ conv, 32 filters, padding=4 |
| **VVS** | | |
| Conv2D + ReLU | $32 \times 32 \times 32$ | $9 \times 9$ conv, 32 filters, padding=4 |
| Conv2D + ReLU | $32 \times 32 \times 32$ | $9 \times 9$ conv, 32 filters, padding=4 |
| Conv2D + ReLU | $32 \times 32 \times 32$ | $9 \times 9$ conv, 32 filters, padding=4 |
| Flatten | $32,768$ | Flatten to vector |
| **Classifier** | | |
| Fully Connected + ReLU | $1,024$ | Linear layer, $32,768 \rightarrow 1,024$ |
| Dropout | $1,024$ | Dropout probability $p = 0.5$ |
| Fully Connected | $C$ | Linear layer, $1,024 \rightarrow C$ |