

A Robust Backpropagation-Free Framework for Images

Anonymous authors

Paper under double-blind review

Abstract

While current deep learning algorithms have been successful for a wide variety of artificial intelligence (AI) tasks, including those involving structured image data, they present deep neurophysiological conceptual issues due to their reliance on the gradients, computed by backpropagation of errors (backprop). Gradients are required to obtain synaptic weight adjustments but require knowledge of feed-forward activities for the backward propagation, a biologically implausible process. This is known as the “weight transport problem”. Therefore, in this work, a more biologically plausible approach towards solving the weight transport problem for image data is presented. It accomplishes this by introducing an error-kernel driven activation alignment (EKDAA) algorithm to train convolutional neural networks (CNNs), using locally derived error transmission kernels and error maps. Like standard deep learning networks, EKDAA performs the standard forward process via weights and activation functions, but its backward error computation involves learning error kernels to propagate local error signals through the network. The efficacy of EKDAA is demonstrated by performing visual-recognition tasks on the Fashion MNIST, CIFAR-10 and SVHN benchmarks, along with demonstrating its ability to extract visual features from natural color images. Furthermore, to demonstrate its non-reliance on gradient computations, results are presented for an EDDAA-driven CNN trained using a non-differentiable activation function.

1 Introduction

One of the most daunting challenges still facing neuroscientists is the understanding of how the neurons in the complex, synaptic network of the brain work together and adjust their synapses in order to accomplish goals (Südhof & Malenka, 2008). While artificial neural networks (ANNs) trained by backpropagation of errors (backprop) present a practical, feasible implementation of learning by synaptic adjustment, it is largely regarded by neuroscientists as biologically implausible for various reasons, including the implausibility of the direct backwards propagation of error derivatives for synaptic updates. Further, it breaks a fundamental requirement for a bio-plausible learning mechanism, where backprop requires access to the feed-forward weights to create a continued error signal backwards to previous layers, a method highly regarded as impossible within the brain. This property, known as the *weight transport problem* (Grossberg, 1987) plagues backprop as a candidate for the base learning rule for building realistic bio-inspired neural modeling frameworks. For a bio-plausible learning rule, it is more likely that neural activity differences, driven by feedback connections, are used in locally effecting synaptic changes (Lillicrap et al., 2020). Such differences-based networks overcome some of backprop’s major implausibilities in a way that is more naturalistic and more compatible with the current understanding of how brain circuitry operates. Although a few classes of algorithms have been proposed to address the specific challenge of

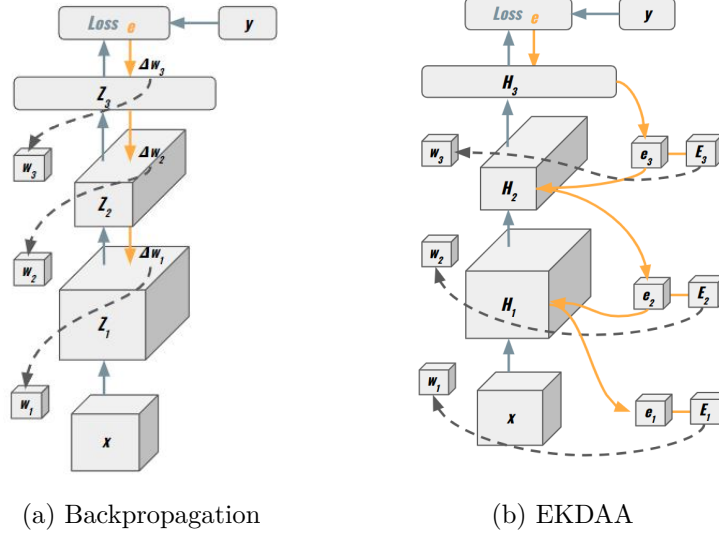


Figure 1: The signal flow for backprop and EKDAA. The forward pass (solid blue arrows), the backward pass (solid orange arrows), and the weight update (gray dashed lines). The pre-activation of a layer N is H_N , and has the corresponding post-activation Z_N . e is the convolutional error kernel, and E is used to transpose the error signal to the appropriate size to propagate backwards.

error gradient propagation in training ANNs, fewer still have been proposed to handle the highly structured data found in large-scale image datasets. Current-day convolutional neural networks (CNNs) and Visual Transformers (ViTs) continue to set the benchmark standards for difficult vision problems (Mnih et al., 2013; He et al., 2016; Dosovitskiy et al., 2020), and they do so using backprop, with symmetric weight matrices in both the feedforward and feedback pathways.

1.1 Bio-plausible Machine Learning

We create a learning rule for the weight transport problem for image data, leveraging spatial relationships with convolution. We introduce a more bio-plausible error synaptic feedback mechanism we deem the (learnable) *error-kernel*, that generates target activities for feature maps within a CNN to align to. We call this learning mechanism *error-kernel driven activation alignment (EKDAA)*.

The Weight Transport Problem In our learning scheme, the forward pathway relies on traditional weight matrices/tensors whereas the backward pathway focuses on error kernels and maps, thus eliminating two-way symmetric weight structure inherent to backprop-trained networks and resolving the weight transport problem.

While currently known bio-plausible methodologies have not reached the modeling performance of backprop and have yet to be scaled to large datasets such as ImageNet (Deng et al., 2009), we believe investigating bio-plausible learning rules are key in future neural modeling. EKDAA notably opens the door to a wider variety of neural structures, potentially enabling lateral neural connections, where forward/backward propagation no longer carry the traditional meaning. We successfully train a convolutional network on image data, using the signum function (which has a derivative of zero everywhere except at zero, i.e., its derivative is a Dirac delta function). Learning with the signum activator showcases how EKDAA allows for investigating bio-plausible activators. The signum function behaves similar to an action potential with a hard activation where a signal is either propagated or killed, abiding to Dale’s law for neurons (Eccles, 1976; Lillicrap et al., 2020).

Bio-plausibly and Convolution The convolution operator is a powerful mechanism to extract spatial features from images and video and has key properties for signal updates that can be integrated into a bio-plausible convolutional neural network. In a backprop model, deconvolution for layer-wise gradient updates, as well as the deconvolution of filter updates can be computed by taking the differentiation of each value of the filter with every element of the matrix that filter touches during the convolution pass. However, gradient updates can also be computed without the need to take direct derivatives on elements by using convolution of the same matrices that derivatives would be taken on. The update signals can be computed by:

$$\Delta \mathbf{W}_{\mathbf{m},\mathbf{n},:,}^{\ell} \leftarrow \mathbf{X}_{\mathbf{m},:,}^{\ell} * \Delta \mathbf{X}_{\mathbf{m},:,}^{\ell+1} \quad (1)$$

$$\Delta \mathbf{X}_{\mathbf{m},\mathbf{n},:,}^{\ell} \leftarrow \Delta \mathbf{X}_{\mathbf{m},\mathbf{n},:,}^{\ell+1} * \text{Flip}(\mathbf{W}_{\mathbf{m},\mathbf{n},:,}^{\ell}) \quad (2)$$

where, $\mathbf{W}_{\mathbf{m},\mathbf{n},:,}^{\ell}$ is the weight matrix for a layer l with feature map element $[m, n]$, $\mathbf{X}_{\mathbf{m},:,}^{\ell}$ is the corresponding model layer output. $\Delta \mathbf{W}_{\mathbf{m},\mathbf{n},:,}^{\ell}$ is the error signal weight matrix for layer l , $\Delta \mathbf{X}_{\mathbf{m},:,}^{\ell}$ is the error signal from the output of layer l , and Flip is the transpose of matrix with respect to both the x and y axis. The convolution operator is represented with $*$. Therefore, convolution works within a Hebbian model and can be used as a powerful feature extractor without the need for computed derivatives. The symmetrical process of convolution in the inference and training pass may also be considerable as a more bio-plausible mechanism over standard feed-forward models. Convolution and deconvolution also are agonistic mechanisms to whether a learning rule is or is not weight symmetric. There are no inherent rules that deconvolution must deconvolve on the exact same weight matrices that are convoluted on in the forward pass.

2 Related Work

Although Hebbian learning (Hebb, 1949) is one of the earliest and simplest biologically plausible learning rules for addressing the credit assignment problem, extending them to the CNN has not yet been well-developed. Our proposed approach aims to fill this gap.

In addition to Hebbian-based learning, other related work includes alternative convolution-based learning schemes such as those found in (Akrouit et al., 2019) that base their work on the Kollen-Pollack (KP) method and have demonstrated promising results on larger, more extensive benchmarks without the need for weight transport. They create *weight mirrors* based on the KP method and incorporate it into the convolutional framework. Other training approaches are based on local losses (Nøkland & Eidnes, 2019; Grinberg et al., 2019; Guerguiev et al., 2017; Schmidhuber, 1990; Werbos, 1982; Linnainmaa, 1970); are methods that take the sign of the forward activities (Xiao et al., 2018); are schemes that utilize noise-based feedback modulation (Lansdell et al., 2019); or use synthetic gradients (Jaderberg et al., 2017) to stabilize learning for deeper networks. These approaches have shown better or comparable performance (to backprop) on challenging benchmarks using the convolution operator. However, there are significant dependencies on the network model’s forward activities used to guide the backward signal propagation (requiring weight transport), hence, these approaches belong to a different class of problems/algorithms than what we address in this work.

The Bottleneck Approach: One of the more recent efforts in this area, inspired by information theory, is the Hilbert-Schmidt independence criterion (HSIC) bottleneck training algorithm (Ma et al., 2019), based on the Information Bottleneck (IB) principle (Tishby et al., 2000). HSIC performs credit assignment locally and layer-wise, seeking hidden representations that have high mutuality with targets but less with the inputs (presented) to that layer (i.e., it is not driven by

the information propagated from the layer below). Approaches based on the bottleneck mechanism are considered the least bio-plausible. Other efforts (Salimans et al., 2017) use an evolutionary strategy to search for optimal weights without gradient descent. However, these approaches often struggle with slow convergence and require many iterations to find optimal solutions.

Feedback Alignment: Notably, an algorithm named *Random Feedback Alignment (RFA)* was proposed in (Lillicrap et al., 2016), where it was argued that the use of the transpose of the forward weights (\mathbf{W}^ℓ for any layer ℓ) in backprop, meant to carry backwards derivative information, was not required for learning. This work showed that network weights could be trained by replacing the transposed forward weights with fixed, random matrices of the same shape (\mathbf{B}^ℓ for layer ℓ), side-stepping the weight transport problem (Grossberg, 1987).

Direct Feedback Alignment (DFA) (Nøkland, 2016), and its variants (Han et al., 2020; Crafton et al., 2019; Chu et al., 2020), was inspired by RFA (Lillicrap et al., 2016), but in contrast to RFA, it directly propagates the error signal (at the output to individual layers directly, rather than layer-wise as is done in RFA). Across multiple neural architectures, it was observed that networks trained with DFA showed a steeper reduction in the classification error when compared to those trained with backprop. To compare these biologically-plausible feedback alignment-based training paradigms with EKDA, we extended the corresponding published works and implemented CNN versions of FA, DFA, and other related variants. Details of their performance on the benchmark image datasets we investigate are given in Section 4.

Target Propagation: Target propagation (target prop, or TP) (Lee et al., 2015) is another approach to credit assignment in deep neural networks, where the goal is to compute targets that are propagated backwards to each layer of the network. Target prop essentially designs each layer of the network as an auto-encoder, with the decoder portion attempting to learn the inverse of the encoder (modified by a linear correction to account for the imperfectness of the auto-encoders themselves). This corrected difference (between encoder and decoder) is then propagated throughout the network. This process allows difference target prop (DTP) (Lee et al., 2015) and variants (Bartunov et al., 2018; Ororbia & Mali, 2019) (e.g., DTP- σ) to side-step the vanishing/exploding gradient problem. However, TP approaches are expensive and can be unstable, requiring multiple forward/backward passes in each layer-wise encoder/decoder in order to produce useful targets.

Representation Alignment: *Local Representation Alignment (LRA)* (Ororbia et al., 2018) and recursive LRA (Ororbia et al., 2020) represent yet another class of credit assignment methods, inspired by predictive coding theory (Clark, 2015) and similar in spirit to target prop. Under LRA, each layer in the neural network has a target associated with it such that changing the synaptic weights in a particular layer will help move layer-wise activity towards better matching a target activity value. LRA was shown to perform competitively to other local learning rules for fully-connected models, but extending/applying it to vectorized natural images like CIFAR-10 resulted in significant performance degradation.

3 Error Kernel Credit Assignment

In implementing the EKDA algorithm, desirably, the forward pass in the CNN remains the same. However, the backward pass uses a form of Hebbian learning that creates targets for each layer (as shown in Figure 1) that allows for creating an error signal that is used to make weights adjustments. In convolutional layers, locally computed error kernels transmit signals by applying the convolution of the error kernel with the pre-activation of that layer aiming to align the forward activations accordingly. In the fully connected case, an error matrix is multiplied by the pre-

activation layer to generate targets allowing for computing pseudo-gradients to train the forward activations. While EKDAA is similar to TP in that it creates targets to optimize towards in an encoding/decoding fashion, EKDAA introduces the novel idea of encoding an error signal with a learned kernel. Similar to FA, EKDAA finds success with using random weights for projecting error signal changing layerwise dimensionality as weight matrices in the forward and backward pass are asymmetric. In contrast to TP and LRA approaches, with our proposed algorithm, EKDAA, the forward and backward activities share as minimal information as possible; this helps the model to better overcome poor initialization and stabilization issues. In this section, we describe the forward pass in our notation and then present the details of the EKDAA learning approach.

Notation: We denote standard convolution with the symbol $*$ and deconvolution with symbol \odot . Hadamard product is denoted by \odot while \cdot represents a matrix/vector multiplication. $()^T$ denotes the transpose operation. $\text{Flip}(\mathbf{X})$ is a function for flipping a tensor and is defined as taking the transpose of \mathbf{X} over both the x-axis and y-axis such that the value of an element $\mathbf{X}_{i,j}$ after flipping results in the location $\mathbf{X}_{n-i,n-j}$. $\text{Flatten}(\mathbf{z})$ means that the input tensor \mathbf{z} is converted to a column vector with a number of rows equal to the number of elements that it originally contained while $\text{UnFlatten}(\mathbf{z})$ is its inverse (i.e., it converts the vector back to its original tensor shape). We use the notation $:$ to indicate extracting a slice of a certain dimension in a tensor object, i.e., $\mathbf{V}_{j,:}$ means that we extract all scalar elements in the j th slice of the three dimensional tensor \mathbf{V} . Finally, $=$ denotes equality while \leftarrow denotes variable assignment.

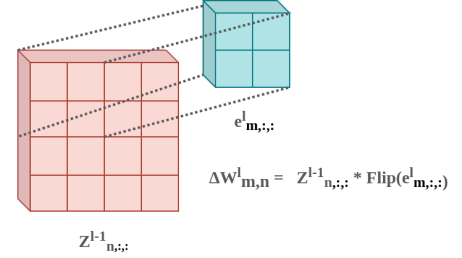


Figure 2: Kernel update to learn the filters $\mathbf{W}^{\ell}_{m,n,:}$ with EKDAA. $\mathbf{z}^{\ell-1}_{n,:}$, the n th post-activation of layer $\ell - 1$, is deconvolved on the m th error kernel $e^{\ell}_{m,:}$ of layer ℓ , propagating the error signal to update $\Delta \mathbf{W}^{\ell}_{m,n,:}$.

Inference Dynamics: Given an input (color) image \mathbf{x} , inference in a basic CNN consists of running a feedforward pass through the underlying model, computing the activities (or nonlinear feature maps) for each layer/level ℓ , where the model contains L_C convolutional layers total. The CNN is parameterized by a set of synaptic tensors $\Theta = \{\mathbf{W}^1, \mathbf{W}^2, \dots, \mathbf{W}^L, \mathbf{W}_y\}$ where the last parameter \mathbf{W}_y is a two-dimensional tensor (or matrix) meant to be used in a softmax/maximum entropy classifier. All other tensors \mathbf{W}^{ℓ} , $\ell = 1, 2, \dots, L$ are four-dimensional and of shape $\mathbf{W}^{\ell} \in \mathcal{R}^{N_{\ell} \times N_{\ell-1} \times h_{\ell} \times w_{\ell}}$. This means that any tensor \mathbf{W}^{ℓ} houses N_{ℓ} sets of $N_{\ell-1}$ filters/kernels of shape $h_{\ell} \times w_{\ell}$. Note that the bottom tensor \mathbf{W}^0 , which takes in as input the source image, would be of shape $\mathbf{W}^0 \in \mathcal{R}^{N_1 \times N_0 \times h_0 \times w_0}$ where N_0 is the number of input color channels, e.g., three, for images of size $h_0 \times w_0$ pixels.

The m th feature map of any convolutional layer ℓ is calculated as a function of the $N_{\ell-1}$ features maps of the layer below ($n \in N_{\ell-1}$ – there are $N_{\ell-1}$ input channels to the m th channel of layer ℓ). This is done, with $\mathbf{h}^{\ell}_{:,m}$ initialized as $\mathbf{h}^{\ell}_{:,m} = \mathbf{0}$, in the following manner (bias omitted for clarity):

$$\mathbf{h}^{\ell}_{m,:} \leftarrow \mathbf{h}^{\ell}_{m,:} + \mathbf{W}^{\ell}_{m,n,:} * \mathbf{z}^{\ell-1}_{n,:}, \forall n \quad (3)$$

$$\mathbf{z}^{\ell}_{m,:} = \phi^{\ell}(\mathbf{h}^{\ell}_{m,:}) \quad (4)$$

where $\mathbf{W}^{\ell}_{m,n,:}$ denotes the specific filter/kernel that is applied to input channel n when computing values for the m th output channel/map. Note that ϕ^{ℓ} is the activation function applied to any output channel in layer ℓ , e.g., $\phi^{\ell}(v) = \max(0, v)$. Max (or average) pooling is typically applied directly after the nonlinear feature map/channel has been computed, i.e., $\mathbf{z}^{\ell}_{m,:} \leftarrow \Phi_{mp}(\mathbf{z}^{\ell}_{m,:})$.

Learning Dynamics: Once inference has been conducted, we may then compute the values needed to adjust the filters themselves. To calculate the updates for each filter in the CNN, EKDAA proceeds in two steps: 1) calculate target activity values for each feature map in each layer (shown in Figure 2) - this is then used to compute the error neurons (or error neuron maps), a type of neuron specialized for computing mismatch signals inspired by predictive processing brain theory (Clark, 2015), and, 2) calculate the updates/adjustments for each filter given the error neuron values. To do so, we introduce a specific set of filter parameters that we call the *error kernels*, each denoted as $\mathbf{E}_{m,n,:}^\ell$, for every map and layer in the CNN. This means that, if we include these error kernels as part of the parameter set of the CNN learned by EKDAA, $\Theta = \{\mathbf{W}^0, \mathbf{E}^0, \mathbf{W}^1, \mathbf{E}^1, \dots, \mathbf{W}^L, \mathbf{E}^L, \mathbf{W}_y, \mathbf{E}_y\}$. Each error filter/kernel is the same shape as its corresponding convolutional filter, i.e., $\mathbf{E}^\ell \in \mathcal{R}^{N_\ell \times N_{\ell-1} \times h_\ell \times w_\ell}$ (except \mathbf{E}_y , which is the same shapes as the transpose of \mathbf{W}_y).

Assuming the tensor target activity \mathbf{y}^ℓ is available to layer ℓ , we compute each channel’s error neuron map as $\mathbf{e}_{m,:}^\ell = -(\mathbf{y}_{m,:}^\ell - \mathbf{z}_{m,:}^\ell)$. Using this mismatch signal, we then work our way down to layer $\ell - 1$ by first convolving this error neuron map to project it downwards, using the appropriate error kernel. Once the projection is complete, if pooling has been applied to the output of each convolutional layer, we then up-sample the projection before computing the final target. This process proceeds formally as follows:

$$\mathbf{e}_{m,:}^\ell = -(\mathbf{y}_{m,:}^\ell - \mathbf{z}_{m,:}^\ell) \quad (5)$$

$$\mathbf{d}_{n,:}^{\ell-1} \leftarrow \mathbf{d}_{n,:}^{\ell-1} + \mathbf{E}_{m,n,:}^\ell \odot \mathbf{e}_{m,:}^\ell, \quad \forall m \in N_\ell \quad (6)$$

$$\mathbf{d}_{n,:}^{\ell-1} \leftarrow \Phi_{up}(\mathbf{d}_{n,:}^{\ell-1}) \quad // \text{ (If max-pooling used)} \quad (7)$$

$$\mathbf{y}_{n,:}^{\ell-1} = \phi^{\ell-1}(\mathbf{h}_{n,:}^{\ell-1} - \beta \mathbf{d}_{n,:}^{\ell-1}) \quad (8)$$

where we see that $\Phi_{up}()$ denotes the up-sampling operation (to recover the dimensionality of the map before max-pooling was applied). Note that if pooling was not used in layer $\ell - 1$, then Equation 7 is omitted in the calculation of layer $\ell - 1$ ’s target activity. Note that the update rule has a recursive nature, since it requires the existence of \mathbf{y}^ℓ which in turn would have been created by applying Equations 5-8 to the layer above, $\ell + 1$. Thus, the base case target activity \mathbf{y}^L , which would exist at the very top (or highest level) of the CNN, and, in the case of supervised classification, which is the focus of this paper, this would be the target label vector \mathbf{y} associated with input image \mathbf{x} .

Once targets have been computed for each convolutional layer, the adjustment for each filter/kernel in each requires a specialized local rule that entails convolving the post-activation maps of the level below with the error neuron map at ℓ . Formally, this means:

$$\Delta \mathbf{W}_{m,n,:}^\ell = \mathbf{z}_{n,:}^{\ell-1} * \text{Flip}(\mathbf{e}_{m,:}^\ell) \quad (9)$$

$$\Delta \mathbf{E}_{m,n,:}^\ell = -\gamma(\Delta \mathbf{W}_{m,n,:}^\ell)^T \quad (10)$$

which can then subsequently be treated as the gradient to be used in either a stochastic gradient descent update, i.e., $\mathbf{W}_{m,n,:}^\ell \leftarrow \mathbf{W}_{m,n,:}^\ell - \lambda \Delta \mathbf{W}_{m,n,:}^\ell$, or a more advanced rule such as Adam (Kingma & Ba, 2017) or RMSprop (Tieleman & Hinton, 2012).

In Algorithm 1, we provide a full mathematical description of how EKDAA would be applied to a deep CNN specialized for classification. Note that while this paper focuses on feedforward classification, our approach is not dependent on the type of task that the CNN is required to solve. For example, one could readily employ our approach to learn convolutional autoencoders for the

Algorithm 1 EKDAA applied to a CNN with max-pooling and a fully-connected maximum entropy output.

```

// Feedforward inference
Input: sample  $(\mathbf{y}, \mathbf{x})$  and  $\Theta$ 
function INFER( $\mathbf{x}, \Theta$ )
    // Pass data thru convolution stack
    // Get image input channels
     $\mathbf{z}_{n,:}^0 = \mathbf{x}_{n,:}, \forall n \in N_0$ 
    for  $\ell = 1$  to  $L_C$  do
        // Calculate feature maps for layer  $\ell$ 
         $\mathbf{h}_{m,:}^\ell = \mathbf{0}, \forall m \in N_\ell$ 
        for  $m = 1$  to  $N_\ell$  do
             $\mathbf{h}_{m,:}^\ell \leftarrow \mathbf{h}_{m,:}^\ell + \mathbf{W}_{m,n,:}^\ell * \mathbf{z}_{n,:}^{\ell-1}, \forall n$ 
             $\mathbf{z}_{m,:}^\ell = \phi^\ell(\mathbf{h}_{m,:}^\ell)$ 
             $\mathbf{z}_{m,:}^\ell \leftarrow \Phi_{mp}(\mathbf{z}_{m,:}^\ell)$ 
         $\mathbf{h}_y = \mathbf{W}_y \cdot \text{Flatten}(\mathbf{z}^{L_C}), \mathbf{z}_y = \sigma(\mathbf{h}_y)$ 
         $\Lambda = \{(\mathbf{h}^1, \dots, \mathbf{h}^{L_C}, \mathbf{h}_y), (\mathbf{z}^0, \dots, \mathbf{z}^{L_C}, \mathbf{z}_y)\}$ 
    Return  $\Lambda$ 

// Calculate weight updates via EKDAA
Input: Statistics  $\Lambda$ , target  $\mathbf{y}$ ,  $\beta$ , and  $\Theta$ 
function CALCUPDATES( $\Lambda, \mathbf{y}, \Theta$ )
     $\mathbf{h}^1, \dots, \mathbf{h}^{L_C}, \mathbf{h}_y, \mathbf{z}^0, \dots, \mathbf{z}^{L_C}, \mathbf{z}_y \leftarrow \Lambda, \mathbf{y}^L = \mathbf{y}$ 
    // Compute softmax weight updates
     $\mathbf{e}_y = -(\mathbf{y} - \mathbf{z}_y)$ 
     $\Delta \mathbf{W}_y = \mathbf{e}_y \cdot (\text{Flatten}(\mathbf{z}^{L_C}))^T$ 
     $\Delta \mathbf{E}_y = -\gamma(\Delta \mathbf{W}_y)^T$ 
     $\mathbf{y}^{L_C} = \phi^{L_C}(\text{Flatten}(\mathbf{h}^{L_C}) - \beta(\mathbf{E} \cdot \mathbf{e}_y))$ 
    // Compute convolutional kernel updates
     $\mathbf{y}^{L_C} \leftarrow \text{UnFlatten}(\mathbf{y}^{L_C})$ 
    for  $\ell = L_C$  to 1 do
        for  $m = 1$  to  $N_\ell$  do
             $\mathbf{e}_{m,:}^\ell = -(\mathbf{y}_{m,:}^\ell - \mathbf{z}_{m,:}^\ell)$ 
         $\mathbf{d}^\ell = \mathbf{0}, \forall n \in N_{\ell-1}$ 
        for  $n = 1$  to  $N_{\ell-1}$  do
            for  $m = 1$  to  $N_\ell$  do
                 $\mathbf{d}_{n,:}^{\ell-1} \leftarrow \mathbf{d}_{n,:}^{\ell-1} + (\mathbf{E}_{m,n,:}^\ell \odot \mathbf{e}_{m,:}^\ell)$ 
             $\mathbf{y}_{n,:}^{\ell-1} = \phi^{\ell-1}(\mathbf{h}_{n,:}^{\ell-1} - \beta \Phi_{up}(\mathbf{d}_{n,:}^{\ell-1}))$ 
        for  $m = 1$  to  $N_\ell$  do
            for  $n = 1$  to  $N_{\ell-1}$  do
                 $\Delta \mathbf{W}_{m,n,:}^\ell = \mathbf{z}_{n,:}^{\ell-1} * \text{Flip}(\mathbf{e}_{m,:}^\ell)$ 
                 $\Delta \mathbf{E}_{m,n,:}^\ell = -\gamma(\Delta \mathbf{W}_{m,n,:}^\ell)^T$ 
     $\Delta = \{\Delta \mathbf{W}^0, \Delta \mathbf{E}^0, \dots, \mathbf{W}^{L_C}, \Delta \mathbf{E}^{L_C}, \mathbf{W}_y, \mathbf{E}_y\}$ 
    Return  $\Delta$ 

```

case of unsupervised learning, to craft alternative convolutional architectures that solve other types of computer vision problems, e.g., image segmentation, or to build more complex models such as those that model time series information, i.e., temporal/recurrent convolutional networks. One key advantage of the above approach is that the test-time inference of the CNN that is learned using the proposed EKDAA is no slower than a standard backprop-trained CNN given that the forward pass remains the same/untouched. EKDAA also benefits from model stabilization over BP, even if the incoming pre-activations are extreme values from poor initialization, it will still give a usable error signal to learn from, as the error signal is the difference in target versus actual (subtractive Hebbian learning) rather than differentiable values as in backprop.

Base Rule: The base rule of EKDAA is shown in Algorithm 2 where only fully connected layers are considered. EKDAA is able to converge to a loss minima like backprop because its error signal pseudo-gradients satisfy the condition that for a given weight update matrix, \mathbf{W}_n : $\Delta \mathbf{W}_n \leq |\Delta \tilde{\mathbf{W}}_n - 90^\circ|$, where $\tilde{\mathbf{W}}_n$ is the exact calculated derivative for the weight matrix \mathbf{W}_n . While the approximated gradients are almost never the exact gradient that backprop provides, they are always within 90 degrees of it, and thus they always tends towards the direction of backprop’s gradient descent. EKDAA’s “gradients” do not take quite as greedy steps towards loss minimization, but over many iterations of training, they do converge towards loss minimization, similar to backprop.

Algorithm 2 EKDAA for the fully-connected maximum entropy output case only.

<pre> // Feedforward inference Input: sample (\mathbf{y}, \mathbf{x}) and Θ function INFER(\mathbf{x}, Θ) $\mathbf{z}^0 = \mathbf{x}$ for $\ell = 1$ to L_C do $\mathbf{h}^\ell \leftarrow \mathbf{h}^\ell + \mathbf{W}^\ell \cdot \mathbf{z}^{\ell-1}$ $\mathbf{z}^\ell = \phi^\ell(\mathbf{h}^\ell)$, $\mathbf{h}_y = \mathbf{W}_y \cdot (\mathbf{z}^{L_C})$, $\mathbf{z}_y = \sigma(\mathbf{h}_y)$ $\Lambda = \{(\mathbf{h}^1, \dots, \mathbf{h}^{L_C}, \mathbf{h}_y), (\mathbf{z}^0, \dots, \mathbf{z}^{L_C}, \mathbf{z}_y)\}$ Return Λ </pre>	<pre> // Calculate weight updates via EKDAA Input: Statistics Λ, target \mathbf{y}, β, and Θ function CALCUPDATES($\Lambda, \mathbf{y}, \Theta$) $\mathbf{h}^1, \dots, \mathbf{h}^{L_C}, \mathbf{h}_y, \mathbf{z}^0, \dots, \mathbf{z}^{L_C}, \mathbf{z}_y \leftarrow \Lambda$, $\mathbf{y}^L = \mathbf{y}$ // Compute softmax weight updates $\mathbf{e}_y = -(\mathbf{y} - \mathbf{z}_y)$ $\Delta \mathbf{W}_y = \mathbf{e}_y \cdot (\mathbf{z}^{L_C})^T$, $\Delta \mathbf{E}_y = -\gamma(\Delta \mathbf{W}_y)^T$ $\mathbf{y}^{L_C} = \phi^{L_C}((\mathbf{h}^{L_C}) - \beta(\mathbf{E} \cdot \mathbf{e}_y))$ for $\ell = L_C$ to 1 do $\mathbf{e}^\ell = -(\mathbf{y}^\ell - \mathbf{z}^\ell)$ $\mathbf{d}^\ell = \mathbf{0}$ $\mathbf{d}^{\ell-1} \leftarrow \mathbf{d}^{\ell-1} + (\mathbf{E}^\ell \cdot \mathbf{e}^\ell)$ $\mathbf{y}^{\ell-1} = \phi^{\ell-1}(\mathbf{h}^{\ell-1} - \beta(\mathbf{d}^{\ell-1}))$ $\Delta \mathbf{W}^\ell = \mathbf{z}^{\ell-1} \cdot (\mathbf{e}^\ell)$ $\Delta \mathbf{E}^\ell = -\gamma(\Delta \mathbf{W}^\ell)^T$ $\Delta = \{\Delta \mathbf{W}^0, \Delta \mathbf{E}^0, \dots, \mathbf{W}^{L_C}, \Delta \mathbf{E}^{L_C}, \mathbf{W}_y, \mathbf{E}_y\}$ Return Δ </pre>
---	---

4 Experimental Setup and Results

4.1 Datasets and Experimental Tasks

To understand learning capacity for fitting and generalization under EKDAA, we design and train several models and test them against three standard datasets, Fashion MNIST (Xiao et al., 2017) (FMNIST), CIFAR-10 (Krizhevsky et al., 2014), and SVHN (Netzer et al., 2011).

Fashion MNIST, while only being a single channel (gray-scale) image dataset at a $[28 \times 28]$ resolution, has a more complicated pixel input space than MNIST, facilitating a better analysis of the convolutional contribution to overall network performance. SVHN and CIFAR-10 images are of shape $[32 \times 32]$ pixels and are represented in three color channels, having more complicated patterns and motifs compared to the gray scale Fashion MNIST (SVHN even contains many images with distractors at the sides of the character the image is centered on). Fully-connected layers are not strong enough to learn the spatial relationships between pixels, and, as a result, convolutional filters are needed in order to learn the key patterns within the data that would help in distinguishing between the varying classes.

Taken together, the FMNIST, CIFAR-10, and SVHN datasets allow us to investigate not only how well EKDAA learns filters when engaged in the process of data fitting but also how effective it is in creating models that generalize on visual datasets. Additionally, we show that our networks can be trained using non-differentiable activations, such as the signum function, or, more formally: $\text{signum}(x) = 1$ if $x > 0$, 0 if $x = 0$, -1 if $x < 0$. In this case, we train all convolutional and fully-connected layers of our model using signum as the activation function, except for the softmax output layer, in order to investigate how well an EKDAA-driven network handles non-differentiable activity without specific tuning.

Dataset/Algorithm	Memory Required (MB)	Computation Time (Sec)
FMNIST/BP	1517.29	96.36 +/- 1.79
FMNIST/EKDAA	2574.25	96.39 +/- 9.16
SVHN/BP	4723.84	316.84 +/- 19.91
SVHN/EKDAA	4728.03	268.95 +/- 54.56
CIFAR-10/BP	4723.84	319.72 +/- 13.59
CIFAR-10/EKDAA	4728.03	271.99 +/- 11.58

Table 1: Memory and computation time required for backprop and EKDAA models. GPU memory and average computation time with variance is recorded per 1000 updates with 10 trials run.

4.2 Technical Implementation

We design CNNs for the Fashion MNIST, CIFAR-10, and SVHN datasets. The Fashion MNIST CNN consists of three convolutional layers before flattening and propagating through one fully-connected layer followed by a softmax. The filter size is $[3 \times 3]$ for all convolutional layers with the first layer starting with one channel, expanding to 32 to 64 and finishing at 128 filters. The fully connected layers start after flattening the filters which are then propagated through 128 fully-connected nodes before finishing at 10 output nodes (one per image class). Max-pooling with a kernel of $[2 \times 2]$ and a stride of 2 was used at the end of the first and second layers of convolution.

The CIFAR-10 and SVHN models use six layers of convolution and is inspired by the blocks of convolution and pooling layers used in the VGG family of networks (Simonyan & Zisserman, 2014). First, two convolutional layers are used before finally passing through a max-pooling layer with a kernel of $[2 \times 2]$ and a stride of two. Three of these mini-blocks of two convolution layers, followed by a max-pooling layer, are used to build the final network. The first three layers of convolution use 64 filters while the last three layers use 128 filters. All layers use a filter size of $[3 \times 3]$. After traversing through the last convolutional layer, the final neural activities are flattened and propagated through a single 128-node, fully-connected layer before shrinking down to 10 output nodes (which are run through the softmax nonlinearity). Both Fashion MNIST, SVHN, and CIFAR-10 models use a very small amount of fully-connected nodes and instead use multiple large filter layers to learn/extract distributed representations (see Appendix for details).

Each model was tuned for optimal performance and several hyper-parameters were adjusted including: batch size, learning rate, filter size, number of filters per layer, number of fully-connected nodes per layer, weight initialization, optimizer choice, and dropout rate. Additional details can be found in Appendix. This exact architecture was used for the EKDAA model as well as for the other learning mechanisms that it was compared against. Models were optimized using stochastic gradient descent with momentum and Pascanu (Pascanu et al., 2013) re-scaling was applied to all layer-wise weight updates.

All models were trained on the original datasets at the original resolutions without any augmentation or pre-training. Unlike backprop, EKDAA did not benefit from extensive heuristic knowledge on optimal parameterization, making a grid search for parameters ineffective as the tuning limits would be significantly wider than with backprop. As a result, for tuning, the learning rate was tuned from the range $1e-1$ to $1e-4$, number of filters were tuned from the range 32 to 128, dropout was tuned from 0 to 0.5, and the activation function was evaluated to be either the hyperbolic tangent (tanh) or the linear rectifier (relu). The final meta-parameter setting for EKDAA was a learning rate of $0.5e-3$, 0.9 momentum rate, tanh for activation, and a dropout rate of 0.1 for filters and 0.3 for fully-connected layers (complete model specifications are in the Appendix). All model weights were randomly initialized with system time as a seed. For EKDAA, the error

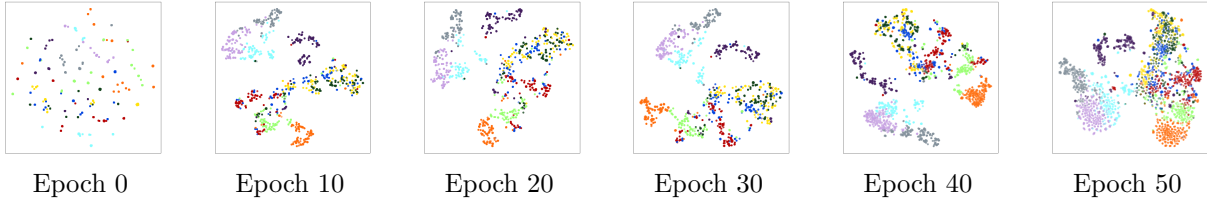


Figure 3: t-SNE visualization depicting the learned representations of EKDAA, shown for Epoch 0 (initial weights) through the final training epoch for FMNIST.

kernels were randomly initialized with the Glorot uniform weight initialization scheme (Glorot & Bengio, 2010). Furthermore, all models were trained on a single Tesla P4 GPU with 8GB of GPU RAM and ran on Linux Ubuntu 18.04.5 LTS, using Tensorflow 2.1.0. The code for this work has been designed in a novel library that allows for defining convolutional and fully-connected models with the ability to quickly change the learning mechanism between BP and EKDAA. The library also allows for defining new custom learning rules for analysis. While this codebase offers an advantage for analyzing learning mechanisms, it has been custom written without the optimization techniques that common libraries have implemented. This codebase takes advantage of Tensorflow tensors when possible but has a custom defined forward and backward pass that is not nearly as memory or computationally efficient as it can be. Our library implementation can be found at: https://anonymous.4open.science/r/EKDAA_release-5613/.

In addition to train and test accuracy, we compare the GPU memory usage and computation time required to train EKDAA and backprop. Results for this are shown in Table 1. Overall, EKDAA shows a slight improvement in computation time as models become larger, and has a minor increase in required memory to run the same model than backprop. Overall, EKDAA’s architecture does not exhibit worse computational requirements to those of backprop.

	FMNIST		SVHN		CIFAR-10	
	Train Acc	Test Acc	Train Acc	Test Acc	Train Acc	Test Acc
BP	95.31 \pm 0.18	89.97 \pm 0.14	90.98 \pm 0.23	88.52 \pm 0.10	83.33 \pm 0.22	71.08 \pm 0.08
BP (FC)	92.91 \pm 0.39	87.02 \pm 0.41	84.36 \pm 0.12	79.81 \pm 0.22	57.05 \pm 0.34	55.03 \pm 0.29
LRA-E (FC)	93.59 \pm 0.26	87.58 \pm 0.33	80.17 \pm 0.08	73.24 \pm 0.19	58.10 \pm 0.28	55.51 \pm 0.42
EKDAA	95.83 \pm 0.33	90.01 \pm 0.11	84.31 \pm 0.21	82.27 \pm 0.19	75.05 \pm 0.27	63.38 \pm 0.12
EKDAA, Sig.	94.00 \pm 0.14	88.69 \pm 0.06	79.43 \pm 0.09	76.87 \pm 0.13	64.22 \pm 0.12	59.71 \pm 0.08
HSIC (Ma et al., 2019)		88.30 \pm —				59.50 \pm —
FA	95.30 \pm 0.51	89.10 \pm 0.18	79.18 \pm 0.22	76.50 \pm 0.18	77.50 \pm 0.25	58.80 \pm 0.11
DFA	93.99 \pm 0.32	88.90 \pm 0.10	82.50 \pm 0.24	80.30 \pm 0.21	79.50 \pm 0.20	60.50 \pm 0.08
SDFA	94.10 \pm 0.28	89.00 \pm 0.10	84.50 \pm 0.23	81.40 \pm 0.19	80.00 \pm 0.19	59.60 \pm 0.06
DRTP	93.50 \pm 0.40	87.99 \pm 0.15	85.21 \pm 0.21	81.90 \pm 0.20	79.50 \pm 0.22	58.20 \pm 0.14

Table 2: Train and test accuracy on the selected datasets. Mean and standard deviation over 10 trials reported. **Note:** the signum (sig.) function is included only to demonstrate that we are able to successfully train a non-differentiable activation function with EKDAA, and obtain reasonable performance. BP results are shown to serve as a best case scenario optimal gradients when training.

4.3 Results and Discussion

We analyzed EKDAA by comparing it to several bio-inspired learning rules such as HSIC (Ma et al., 2019), RFA (Lillicrap et al., 2016), DFA (Nøkland, 2016), sparse direct feedback alignment (SDFA) (Crafton et al., 2019), and direct random target projection (DRTP) (Frenkel et al., 2021) (see Appendix for baseline details). The results are presented in Table 2 (we also added two fully-

connected baselines – an MLP trained by backprop and one trained by LRA-E (Ororbias & Mali, 2019)). Comparable BP results are shown only to provide intuition on how well the constructed models could perform if trained with precise gradients. We find EKDAA performs competitively with the other algorithms and exhibits both increased train and test accuracy on the natural color images with SVHN and CIFAR-10. Additionally, when testing EKDAA with the signum activation, we find the resulting CNN can operate with the non-differentiable function successfully on Fashion MNIST.

Figure 3 shows the t-SNE plots of the outputs of the last layer of a CNN across epochs, to demonstrate how EKDAA disentangles the feature space of Fashion MNIST. The t-SNE plots were generated every 10th epoch and visualized using default t-SNE parameters (i.e. no tuning) with a perplexity value of 30 and the maximum number of iterations set to 1000. Qualitatively, we find EKDAA successfully learns to group features together with primarily convolutional layers, indicating the error kernels learned are indeed benefiting the CNN model. Figure 4 shows the training and test learning curves of backprop and EKDAA (plotting the accuracy as a function of epoch for the best configuration of the model using each algorithm).

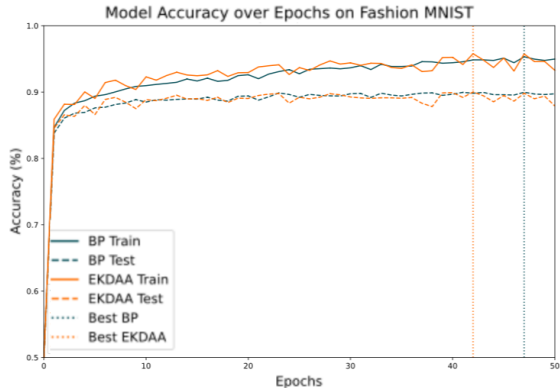


Figure 4: FMNIST train and test accuracy curves for BP and EKDAA.

While many biologically-plausible alternatives have also been developed to learn models of natural images, many of them incorporate error derivatives as part of the process and the architectures that they are generally applied to have been designed with multiple, large fully-connected layers with only a few convolutional layers. We argue that adding many fully connected layers corrupts the original input signal such that the neural model is engaged in a greedy optimization process that results in fitting to noise rather than extracting useful features from natural images. The role of convolution in such models is still debatable and, as a result, it is difficult to determine if model generalization/performance is coming from the bio-plausible learning mechanism or from the fully-connected layers. In contrast, EKDAA emphasizes the role of convolutional filters in extracting useful image features while reducing the amount of fully-connected elements. Our results on the three datasets examined above validate that this approach still yields models that generalize well.

Limitations: The main limitation of the proposed EKDAA algorithm is currently its scalability to massive datasets, especially when compared with highly optimized tensor computations implemented in standard deep learning libraries that support backprop-based convolution/deconvolution operations. Due to the current lack of advanced optimizations when compared to frameworks such as TensorFlow and PyTorch, supported by the large tech companies, EKDAA is not as efficient, thus requiring more computational resources than the established frameworks. Based on our practical experience with our custom software library that implements EKDAA, it does not scale easily to very large networks with many filters in each layer given our constrained computational budget and hardware. Specifically, we used one Ubuntu 18.04 server with 8GB Tesla P4, an Intel Xeon CPU E5-2650, and 256GB of RAM. Our primary future work will be to further modularize our EKDAA software library focusing on improving the algorithm’s ability to scale to training on much larger datasets such as ImageNet, despite the limited resources. (Deng et al., 2009).

Broader Impacts: Backprop generally works well on a carefully parameterized network, but it has a few drawbacks. For example, it requires computing gradients layer-by-layer, thus enforcing strict requirements on propagation flow, i.e. needing to propagate in only a feed-forward flow (i.e. lateral connections make no sense for backprop). Local learning mechanisms do not have these restrictions and allow for exploration into novel network structures and propagation flow.

We introduce a novel framework for training images without the need for back-propagation. While our proposed work is limited with scale and speed compared to the highly optimized tensor computations implemented in standard deep learning libraries, this work serves as a foundation for exploring local learning on image data. We introduce EKDAA, an algorithm that learns error kernels from local layer convolutional signals to better represent image data in a backprop-free manner. While several novel bio-plausible methodologies have been developed in recent years they their learning rules tend to strictly focus on learning standard feed-forward linear layers. Some of these methods apply convolution before linear layers, but often fix randomly initialized filters or use backprop to train those layers. By utilizing learnable error kernels, we introduce a way to transfer error signals through convolutional layers of the model all while not requiring gradient information like in backprop, and learn in a local way that does not impose the weight transport problem. Continued work in this area may have profound impact in future model development by ablating the severe restrictions of backprop and still provide the ability to model highly complex and high dimensional data such as natural color images.

5 Conclusion

We have presented an initial exploration of a back-propagation-free convolutional neural network learning algorithm. We implemented a local feedback mechanism that transmits information across layers in order to compute target activity values and relevant error neuron maps (independent of activation function type), resulting in Hebbian-like update rules for the convolutional filters of a CNN. Specifically, this credit assignment process was made possible through the introduction of a mechanism that we call the error kernel, which provides a means to reverse filter error neuron activity signals and complements the normal filters used to extract features in convolutional models. We refer to our proposed process as the *error-kernel driven activation alignment (EKDAA)* method.

We compared various learning algorithms in training a small CNN and find that EKDAA outperforms other bio-inspired alternatives on natural color images for the task of image classification for the datasets tested. Notably, our method offers several benefits. It resolves the major bio-implausibility of the weight transport problem, works with non-differentiable activities, and is relatively computationally efficient since it can operate in a layer-wise parallel/asynchronous fashion.

Our experiments demonstrate that EKDAA learns “good” representations during training and, furthermore, we found that an EKDAA-trained CNN acquires latent representations that improve over time (epochs) as training evolves. Additionally, we find that EKDAA has similar computational and memory requirements as backprop, and shows similar loss convergence as well. While there is still much to explore in the future work, we have successfully presented an analysis of the novel EKDAA algorithm, yielding promising evidence that it is capable of training convolutional networks without backprop. Future directions to expand this work involve expanding to larger and more complex architectures and imaging data, as well as a study of the properties a Hebbian-based algorithm like EKDAA offer over backprop. The implication of this could have far-reaching effects in expanding the future designs of CNN architectures.

References

- Mohamed Akrouf, Collin Wilson, Peter C Humphreys, Timothy Lillicrap, and Douglas Tweed. Deep learning without weight transport. *arXiv preprint arXiv:1904.05391*, 2019.
- Sergey Bartunov, Adam Santoro, Blake Richards, Luke Marris, Geoffrey E Hinton, and Timothy Lillicrap. Assessing the scalability of biologically-motivated deep learning algorithms and architectures. In *Advances in Neural Information Processing Systems*, pp. 9368–9378, 2018.
- Tien Chu, Kamil Mykitiuk, Miron Szewczyk, Adam Wiktor, and Zbigniew Wojna. Training dnns in $\mathcal{O}(1)$ memory with mem-dfa using random matrices. *arXiv preprint arXiv:2012.11745*, 2020.
- Andy Clark. *Surfing uncertainty: Prediction, action, and the embodied mind*. Oxford University Press, 2015.
- Brian Crafton, Abhinav Parihar, Evan Gebhardt, and Arijit Raychowdhury. Direct feedback alignment with sparse connections for local learning. *Frontiers in neuroscience*, 13:525, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. IEEE, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- John Carew Eccles. From electrical to chemical transmission in the central nervous system: the closing address of the sir henry dale centennial symposium cambridge, 19 september 1975. *Notes and records of the Royal Society of London*, 30(2):219–230, 1976.
- Charlotte Frenkel, Martin Lefebvre, and David Bol. Learning without feedback: Fixed random learning signals allow for feedforward training of deep neural networks. *Frontiers in Neuroscience*, 15:20, 2021.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS’10)*. Society for Artificial Intelligence and Statistics, 2010.
- Leopold Grinberg, John Hopfield, and Dmitry Krotov. Local unsupervised learning for image analysis. *arXiv preprint arXiv:1908.08993*, 2019.
- Stephen Grossberg. Competitive learning: From interactive activation to adaptive resonance. *Cognitive Science*, 11(1):23–63, 1987. doi: <https://doi.org/10.1111/j.1551-6708.1987.tb00862.x>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1551-6708.1987.tb00862.x>.
- Jordan Guerguiev, Timothy P Lillicrap, and Blake A Richards. Towards deep learning with segregated dendrites. *Elife*, 6:e22901, 2017.
- Donghyeon Han, Gwangtae Park, Junha Ryu, and Hoi-jun Yoo. Extension of direct feedback alignment to convolutional and recurrent neural network for bio-plausible deep learning. *arXiv preprint arXiv:2006.12830*, 2020.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Donald Olding Hebb. The organization of behavior; a neuropsychological theory. *A Wiley Book in Clinical Psychology*, 62:78, 1949.
- Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, David Silver, and Koray Kavukcuoglu. Decoupled neural interfaces using synthetic gradients. In *International Conference on Machine Learning*, pp. 1627–1635. PMLR, 2017.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. The cifar-10 dataset. *online: <http://www.cs.toronto.edu/kriz/cifar.html>*, 55, 2014.
- BJ Lansdell, PR Prakash, and KP Kording. Learning to solve the credit assignment problem. *arXiv preprint arXiv:1906.00889*, 2019.
- Dong-Hyun Lee, Saizheng Zhang, Asja Fischer, and Yoshua Bengio. Difference target propagation, 2015.
- Timothy P Lillicrap, Adam Santoro, Luke Marris, Colin J Akerman, and Geoffrey Hinton. Back-propagation and the brain. *Nature Reviews Neuroscience*, 21(6):335–346, 2020.
- TP Lillicrap, D Cownden, DB Tweed, and C Akerman. Random synaptic feedback weights support error backpropagation for deep learning. *Nature Communications*, 7(13276):1–10, 2016.
- Seppo Linnainmaa. The representation of the cumulative rounding error of an algorithm as a taylor expansion of the local rounding errors. *Master’s Thesis (in Finnish), Univ. Helsinki*, pp. 6–7, 1970.
- Wan-Duo Kurt Ma, J. P. Lewis, and W. Bastiaan Kleijn. The hsic bottleneck: Deep learning without back-propagation, 2019.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011. URL http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf.
- Arild Nøkland and Lars Hiller Eidnes. Training neural networks with local error signals. In *International Conference on Machine Learning*, pp. 4839–4850. PMLR, 2019.
- Arild Nøkland. Direct feedback alignment provides learning in deep neural networks, 2016.
- Alexander Ororbia, Ankur Mali, Daniel Kifer, and C Lee Giles. Reducing the computational burden of deep learning with recursive local representation alignment. *arXiv preprint arXiv:2002.03911*, 2020.

- Alexander G Ororbia and Ankur Mali. Biologically motivated algorithms for propagating local target representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 4651–4658, 2019.
- Alexander G. Ororbia, Ankur Mali, Daniel Kifer, and C. Lee Giles. Conducting credit assignment by aligning local representations, 2018.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pp. 1310–1318. PMLR, 2013.
- Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- Jürgen Schmidhuber. Networks adjusting networks. In *Proceedings of "Distributed Adaptive Neural Information Processing"*, pp. 197–208, 1990.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Thomas C Südhof and Robert C Malenka. Understanding synapses: past, present, and future. *Neuron*, 60(3):469–476, 2008.
- T. Tieleman and G. Hinton. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.
- Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.
- Paul J Werbos. Applications of advances in nonlinear sensitivity analysis. In *System modeling and optimization*, pp. 762–770. Springer, 1982.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- Will Xiao, Honglin Chen, Qianli Liao, and Tomaso Poggio. Biologically-plausible learning algorithms can scale to large datasets. *arXiv preprint arXiv:1811.03567*, 2018.

A Appendix

B Experimental Setup Details

We performed a grid search for all of the models investigated in this work in order to find optimal meta-parameters and extract optimal behavior for each. Primarily, tuned hyper-parameters included: batch size, learning rate, filter size, number of filters per layer, number of fully-connected nodes/units per layer, weight initialization, choice of optimizer, and the dropout rate. Note that this work does not aim to obtain state-of-the-art image classification results. Rather, its intent is to present a method that efficiently tackles the credit assignment issues in a convolution neural network (CNN) by effectively operating with our proposed error kernel mechanism. Furthermore, our method offers additional flexibility in design choices (such as permitting the use of non-differentiable activation functions).

Meta-parameter Tuning: We report our grid search ranges for each model’s meta-parameters in Tables 2 (EKDAA), 6 (DFA), 5 (FA), 8 (RDFA), and 7 (SDFA), respectively. Furthermore, in the “Best” column, we report the final values selected/used for the models reported in the main paper.

Architecture Design: In Table 3, we present the architectures used across the learning algorithms investigated in this paper, i.e., the proposed EKDAA, feedback alignment (FA, also referred to as RFA in the main paper), direct feedback alignment (DFA), sparse direct feedback alignment (SDFA), and random direct feedback alignment (RDFA). We built the models for Fashion MNIST, SVHN, and CIFAR-10 to include several layers of convolution (conv), with a sizeable amount of filters, and only small (in terms of dimensionality) fully-connected (fc) layers to focus the learning process on adapting/using the model kernels/filters to extract useful features from the input image signals. In particular, the model for SVHN and CIFAR-10 had multiple layers with 128 filters per layer and, before flattening the activities for the fully-connected layers, the image size was reduced using three max pooling layers in order to propagate forward the image to obtain a $[4 \times 4]$ resolution.

General Comments/Discussion: With respect to the main paper’s results, what is significant about our findings is that EKDAA demonstrates that adjusting the synaptic weight parameters of a CNN is possible using recurrent error synapses formulated as error kernels themselves. This means that the target feature map values (and the error neuron maps that calculate the distance between the original feature maps and these targets) inherent to our backprop-free computational process provide useful teaching signals that facilitate the learning of useful neural vision architectures. The main results of our paper provide promising initial evidence that EKDAA can serve as a potentially useful bio-inspired alternative to backprop for training CNNs on natural images.

C Asset Usage

We build our codebase on top of TensorFlow 2.0 for fundamental functionality. TensorFlow is open-source with an Apache license. In addition, for analysis we use the publicly available Fashion-MNIST, SVHN, and CIFAR-10 datasets, all of which have licenses permitting unlimited use and modification. In addition, none of the datasets used in this study entail any data that could be considered sensitive (thus not requiring data consent) or offensive.

D Model and Training Specifications

Layer	Fashion MNIST	Fashion MNIST Output	SVHN/CIFAR-10	SVHN/CIFAR-10 Output
L0	Input	[:, 28, 28, 1]	Input	[:, 32, 32, 3]
L1	Conv1 (1, 32) $[3 \times 3]$	[:, 28, 28, 32]	Conv1 (3, 64) $[3 \times 3]$	[:, 32, 32, 64]
L2	MaxP1 (2, 2)	[:, 14, 14, 32]	Conv2 (64, 64) $[3 \times 3]$	[:, 32, 32, 64]
L3	Conv2 (32, 64) $[3 \times 3]$	[:, 14, 14, 64]	MaxP1 (2, 2)	[:, 16, 16, 64]
L4	MaxP2 (2, 2)	[:, 7, 7, 64]	Conv3 (64, 64) $[3 \times 3]$	[:, 16, 16, 64]
L5	Conv3 (64, 128) $[3 \times 3]$	[:, 7, 7, 128]	Conv4 (64, 128) $[3 \times 3]$	[:, 16, 16, 128]
L6	Flatten()	[:, 6272]	MaxP2 (2, 2)	[:, 8, 8, 128]
L7	FC1 (6272, 128)	[:, 128]	Conv5 (128, 128) $[3 \times 3]$	[:, 8, 8, 128]
L8	Softmax (128, 10)	[:, 10]	Conv6 (128, 128) $[3 \times 3]$	[:, 8, 8, 128]
L9	-	-	MaxP3 (2, 2)	[:, 4, 4, 128]
L10	-	-	Flatten()	[:, 2048]
L11	-	-	FC1 (2048, 128)	[:, 128]
L12	-	-	Softmax (128, 10)	[:, 10]

Table 3: Model architectures that were trained on Fashion MNIST, SVHN, and CIFAR-10. The layers of each model are defined as well as the outputs from each layer.

Parameter	Range Min	Range Max	Interval	Activation Functions	Best
batch_size	50	250	50	-	50
learning_rate	1e-5	1e-2	0.5	-	5e-4
filter_size	3	7	2	-	3
num_filters	32	256	32	-	-
fc_per_layer	128	128	-	-	128
weight_init	-	-	-	glorot_uniform, glorot_normal	glorot_uniform
optimizer	-	-	-	tanh, relu, signum	tanh
dropout	0.0	0.5	0.1	-	0.1 conv, 0.3 fc

Table 4: Hyper-parameter tuning ranges and best found parameters for EKDA.

Parameter	Range Min	Range Max	Increment	Activation Functions	Best
batch_size	32	256	64	-	64
learning_rate	5e-5	3e-2	0.5	-	5e-4
filter_size	3	7	2	-	3
num_filters	32	256	32	-	-
fc_per_layer	128	128	-	-	128
weight_init	-	-	-	glorot_uniform, glorot_normal	glorot_normal
optimizer	-	-	-	tanh, relu	relu
dropout	0.0	0.5	0.1	-	0.1 conv, 0.3 fc

Table 5: Hyper-parameter tuning ranges and best found parameters for FA.

Parameter	Range Min	Range Max	Increment	Activation Functions	Best
batch_size	32	256	64	-	64
learning_rate	5e-5	3e-2	0.5	-	5e-3
filter_size	3	7	2	-	3
num_filters	32	256	32	-	-
fc_per_layer	128	128	-	-	128
weight_init	-	-	-	glorot_uniform, glorot_normal	glorot_uniform
optimizer	-	-	-	tanh, relu	relu
dropout	0.0	0.5	0.1	-	0.1 conv, 0.2 fc

Table 6: Hyper-parameter tuning ranges and best found parameters for DFA.

Parameter	Range Min	Range Max	Increment	functions	Best
batch_size	32	256	64	-	64
learning_rate	5e-5	3e-2	0.5	-	3e-3
filter_size	3	7	2	-	3
num_filters	32	256	32	-	-
fc_per_layer	128	128	-	-	128
weight_init	-	-	-	glorot_uniform, glorot_normal	glorot_uniform
optimizer	-	-	-	tanh, relu	tanh
dropout	0.0	0.5	0.1	-	0.1 conv, 0.1 fc

Table 7: Hyper-parameter tuning ranges and best found parameters for SDFA.

Parameter	Range Min	Range Max	Increment	functions	Best
batch_size	32	256	32	-	32
learning_rate	5e-5	3e-2	0.5	-	4e-3
filter_size	3	7	2	-	3
num_filters	32	256	32	-	-
fc_per_layer	128	128	-	-	128
weight_init	-	-	-	glorot_uniform, glorot_normal	glorot_normal
optimizer	-	-	-	tanh, relu	relu
dropout	0.0	0.5	0.1	-	0.2 conv, 0.3 fc

Table 8: Hyper-parameter tuning ranges and best found parameters for RDFA.