

# EFFICIENT GRADIENT-BASED ALGORITHM FOR TRAINING DEEP LEARNING MODELS WITH MANY NONLINEAR ACTIVATIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This research paper presents a novel algorithm for training deep neural networks with many nonlinear layers (e.g., 30). The method is based on backpropagation of an approximated gradient, averaged over the range of a weight update. Unlike the gradient, the average gradient of a loss function is proven within this research to provide more accurate information on the change in loss caused by the associated parameter update of a model. Therefore, it may be utilized to improve learning. In our implementation, the efficiently approximated average gradient is paired with RMSProp and compared to the typical gradient-based approach. For the tested deep model with numerous stacked fully-connected layers featuring nonlinear activations on MNIST and Fashion MNIST, the presented algorithm: (a) generalizes better, at least in a reasonable epoch count, (b) in the case of optimal implementation, learning would require less computation time than the gradient-based RMSProp, with the memory requirement of the Adam optimizer, (c) performs well on a broader range of learning rates, therefore it may bring time and energy savings from reduced hyperparameter searches, (d) improves sample efficiency about three times according to median training losses. On the other hand, for a deep sequential convolutional model trained on the IMDB dataset, sample efficiency is improved by about 55%. However, in the case of the tested shallow model, the method performs approximately the same as the gradient-based RMSProp in terms of both training and test loss. The source code is provided at [...].

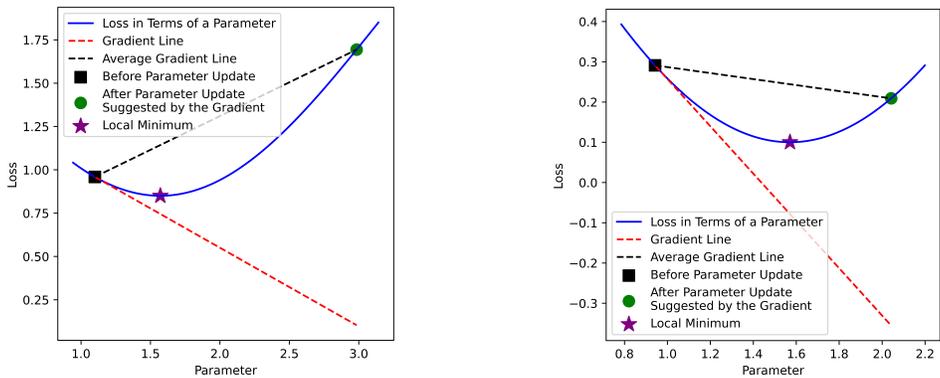
## 1 INTRODUCTION

### 1.1 AVERAGE GRADIENT

In this research, we focus on solving deep learning problems by calculating the precise influence of potential updates on the loss for each model parameter separately. Our goal is to obtain more precise information about the influence of each parameter on loss than what the gradient provides. Each potential update of a model parameter influences the locally-optimal direction of other model parameters during the same weight update, highlighting the complexity of the problem. The average gradient (defined in Appendix A), unlike the gradient, stores the accurate contribution of each model parameter to the loss delta related to a given weight update (Fig. 1; Eq. 14). Therefore, the average gradient can be utilized to efficiently minimize the loss. In this research, we propose a very fast algorithm to approximate the average gradient. We prove its approximation accuracy, validate the proof using our handcrafted metric to compare batch-loss minimization efficiency between methods, and test our method on various domains and models. Our algorithm in its current form primarily targets very deep models with many nonlinear layers.

Due to the tendency to increase model depth along with its width (Tan & Le, 2019) and the popularity of certain nonlinear activation functions, our approach may offer insights for future improvements in practical deep learning. Our primary target is to significantly improve sample efficiency, even at the cost of a moderate increase in computation time, which is essential for practical applications in fields like deep reinforcement learning or reinforcement learning from human feedback (Kirk et al., 2023). These methods have been used in popular chatbots, such as OpenAI’s ChatGPT (OpenAI, 2023) and Anthropic’s Claude (Kirk et al., 2023).

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107



(a) *Example 1.* The average gradient suggests a different direction for updating a particular parameter. (b) *Example 2.* If the average gradient decreases in the same direction as the gradient, it additionally provides more information about the loss landscape.

Figure 1: *Comparison of Gradient and Average Gradient.* The latter accurately reflects the influence of a parameter update on loss (as described by Equation 3, under the assumptions that  $f$  represents the visualized loss function, with  $x$  and  $x'$  denoting the parameter values before and after an update, respectively). The plots refer to a simple case with only one parameter of the model. However, they can also be understood to present the loss contribution of a parameter during a weight update involving multiple model parameters. Appendix F presents visualizations involving two parameters.

## 1.2 GRADIENT OPTIMIZATION AND AVERAGING

Gradient optimization dominates deep learning with optimizers like Stochastic Gradient Descent (Liu et al., 2020), RMSProp (Tieleman et al., 2012), Adam (Kingma & Ba, 2014) or Nadam (Dozat, 2016). The leading algorithms for training do not change frequently over the years. However, our algorithm or its variants may be used along with first-order optimizers.

Gradient averaging is commonly used in machine learning, but in a distinct scenario than in our approach. Momentum is the running average of gradients over subsequent batches (Liu et al., 2020; Kingma & Ba, 2014; Dozat, 2016). It prevents falling into local minimums and may accelerate learning. Similarly, averaging model parameters may improve convergence and learning speed (Ruppert, 1988; Polyak & Juditsky, 1992; Merity et al., 2017; Wei et al., 2023; Sun et al., 2010), though it requires a significant amount of memory. The technique can be described as averaging a function of gradients, as the averaged model parameters over subsequent updates depend on the gradient values.

Accumulating gradients over a batch is inherent in machine learning. In practice, it is equivalent to averaging gradients computed for multiple inputs. However, this approach alone does not take into account the information about a parameter update (Fig. 1), which remains unknown during its computation. Consequently, it does not guarantee the accuracy of computing the influence of the unknown parameter update on the loss. Nevertheless, batching remains fully compatible with our method and is employed in our implementation.

Our approach is more closely related to some second-order optimization methods (Tan & Lim, 2019) rather than momentum-based or parameter-averaging techniques. This is due to the utilization of information about the curvature of a loss function during each parameter update (Fig. 1). Recently, one of the most popular algorithms for second-order optimization of neural networks is L-BFGS (Berahas et al., 2016). However, the current methods in this field are impractical for training large models due to their computational inefficiency or substantial memory requirements.

The integrated gradient, closely related to the average gradient, is used in some neural-network explainability techniques (Sundararajan et al., 2017; Khorram et al., 2021; Sattarzadeh et al., 2021). However, the approximation algorithms for the integral of the gradient used in the literature are very inefficient to compute for every parameter update of a model due to the calculation of the Riemann sum (Hughes-Hallett et al., 2021).

## 2 METHODS

### 2.1 ALGORITHM

All of the best and most popular optimizers for training large neural networks rely on the gradient. Consequently, they explicitly ignore how loss function in terms of model parameters behaves in the range between before and after a potential weight update (Fig. 1). The definition of the gradient implies, that it reflects the accurate influence on loss only for learning rates approaching to zero, which does not hold in practice. Consequently, gradient-based optimizers do not calculate the accurate influence on loss of a potential weight update, which may significantly slow down the learning of very deep models with many nonlinear operators, as our experiments show. The average gradient solves the described problem. Our algorithm efficiently approximates the average gradient, providing more reliable information on the update direction that minimizes the loss. The average gradient (contrary to the gradient) is directly proportional to the loss delta (Fig. 1; Equation 14 in Appendix B), hence it accurately describes the influence on loss of a parameter delta.

In our algorithm, given a sequential model, the average gradient is approximated and propagated according to the equation proven in Appendix B:

$$\mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell \cong \mathcal{AVG}_{\theta_k} \frac{\partial \mathbf{x}_k}{\partial \theta_k} \cdot \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k} \cdot \dots \cdot \mathcal{AVG}_{\mathbf{x}_{n-1}} \frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_{n-1}} \cdot \mathcal{AVG}_{\mathbf{x}_n} \nabla_{\mathbf{x}_n} \ell \quad (1)$$

where  $\ell$  is a loss function,  $\theta_k$  are parameters of a layer no.  $i$  and  $(\mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$  are inputs and outputs of subsequent layers of a neural network. The notation  $\nabla_{\mathbf{x}} f$  refers to the gradient of some function  $f$  for an argument  $\mathbf{x}$ , and  $\frac{\partial f}{\partial \mathbf{x}}$  denotes the Jacobian. The average operator  $\mathcal{AVG}$  of gradients or Jacobians is defined in Appendix A and aligns with intuition. The averages are aggregated with respect to the parameters of a model ( $\theta_k$ ) or the outputs of subsequent layers  $(\mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$ . The average gradients are propagated in the same manner as the gradients in the standard backpropagation algorithm. The computation based on Equation 1 is fast and memory efficient because the procedure is similar to the standard backpropagation of gradients, which is done according to:

$$\nabla_{\theta_k} \ell = \frac{\partial \mathbf{x}_k}{\partial \theta_k} \cdot \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k} \cdot \dots \cdot \frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_{n-1}} \cdot \nabla_{\mathbf{x}_n} \ell \quad (2)$$

The version of our algorithm that consists of two iterations (Algorithm 1) first performs the standard backpropagation (Equation 2) through layer outputs  $\mathbf{x}$ , and model parameters  $\theta$  along with parameter update of an optimizer (in the experiments it is RMSProp) to new weight values  $\theta'$ . Then it is assumed that the absolute value of the parameter delta  $|\theta - \theta'|$  of the RMSProp optimizer is good enough to retain it. The second backpropagation is performed for eventual negations of update directions only, where, conversely, the *average* gradient is propagated (Algorithm 2). Importantly, the range on which the gradient is averaged equals  $[\theta, \theta']$  (between parameters before and after the estimated potential update; Algorithm 3). The average derivatives of each nonlinear activation are calculated as follows:

$$\mathcal{AVG}_{t \in [x, x']} f'(t) = \frac{\int_x^{x'} f'(t) dt}{x' - x} = \frac{f(x') - f(x)}{x' - x} \quad (3)$$

where  $f$  means an activation function (in the experiments it is either ELU or Tanh activation),  $x$  means an input scalar assuming forward propagation using the  $\theta$  weights, and  $x'$  means the corresponding input number assuming forward pass for the  $\theta'$ . The equation is the one-dimensional analogy of the average gradient and the Jacobian, both of which are defined in Appendix A.

In the case of applying an activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , or  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , to a layer output  $\mathbf{x} = \langle x_1, x_2, \dots, x_n \rangle$  (assuming parameters  $\theta$ ), which changes to  $\mathbf{x}' = \langle x'_1, x'_2, \dots, x'_n \rangle$  during the forward pass with updated parameters  $\theta'$ :

$$\mathcal{AVG}_{\mathbf{x}, \mathbf{x}'} \frac{\partial \mathbf{f}}{\partial \mathbf{t}} = \text{diag}(\langle \mathcal{AVG}_{t_1 \in [x_1, x'_1]} f'(t_1), \mathcal{AVG}_{t_2 \in [x_2, x'_2]} f'(t_2), \dots, \mathcal{AVG}_{t_n \in [x_n, x'_n]} f'(t_n) \rangle) \quad (4)$$

where each term  $\mathcal{AVG}_{(\cdot)} f'(\cdot)$  is defined in Equation 3.

Let us define a typical layer, denoted as  $k$ , which is parameterized by  $\theta_k$ . This layer could be a convolutional layer, a fully-connected layer, or another operator that is linear over all or most of

its domain. Let us assume that the layer no.  $k$  outputs  $\mathbf{y}_k$ , which is then passed to an activation  $f_k$ . Consequently each part of Equation 1 can be approximated as:

$$\begin{aligned} \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k} &= \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial f_k}{\partial \mathbf{x}_k} \cong \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{y}_k}{\partial \mathbf{x}_k} \cdot \mathcal{AVG}_{t \in [\mathbf{y}_k, \mathbf{y}'_k]} \frac{\partial f}{\partial t} \\ \mathcal{AVG}_{\theta_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \theta_k} &= \mathcal{AVG}_{\theta_k} \frac{\partial f_k}{\partial \theta_k} \cong \mathcal{AVG}_{\theta_k} \frac{\partial \mathbf{y}_k}{\partial \theta_k} \cdot \mathcal{AVG}_{t \in [\mathbf{y}_k, \mathbf{y}'_k]} \frac{\partial f}{\partial t} \end{aligned} \quad (5)$$

where the approximation, instead of equality, is the consequence of chaining averages of Jacobians, which can be proven analogously to Equation 1 (see Appendix B). The average operator  $\mathcal{AVG}$  of Jacobians is defined in Appendix A.  $\mathcal{AVG}_{t \in [\mathbf{y}_k, \mathbf{y}'_k]} \frac{\partial f}{\partial t}$  is defined in Equation 4. Generally, the vast majority of applied neural network operators are either nonlinear activations or linear functions in by far most of their domains (e.g., max pooling, convolution, fully connected, or ReLU). In the case of the nonlinear activations, equations no. 3 and 4 are used to compute the average Jacobians. For linear transformations, such as  $\mathbf{y}_k(\mathbf{x}_k)$  and  $\mathbf{y}_k(\theta_k)$ , the average gradients and Jacobians are easy and fast to compute. However, for implementation simplicity and a slight speedup of computations, broader estimates of the average Jacobians from Equation 5 are applied:

$$\begin{aligned} \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k} &= \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial f_k}{\partial \mathbf{x}_k} \cong \frac{\partial \mathbf{y}_k}{\partial \mathbf{x}_k} \cdot \mathcal{AVG}_{t \in [\mathbf{y}_k, \mathbf{y}'_k]} \frac{\partial f}{\partial t} \\ \mathcal{AVG}_{\theta_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \theta_k} &= \mathcal{AVG}_{\theta_k} \frac{\partial f_k}{\partial \theta_k} \cong \frac{\partial \mathbf{y}_k}{\partial \theta_k} \cdot \mathcal{AVG}_{t \in [\mathbf{y}_k, \mathbf{y}'_k]} \frac{\partial f}{\partial t} \end{aligned} \quad (6)$$

which use the non-averaged Jacobian  $\frac{\partial \mathbf{y}_k}{\partial \mathbf{x}_k}$ . Therefore, intuitively, the broad estimation of  $\mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k}$  is approximately between  $\frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k}$  and  $\mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k}$ , and analogously for  $\mathcal{AVG}_{\theta_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \theta_k}$ .

---

**Algorithm 1** *Simplified algorithm version for 2 iterations.* Back and forward propagation would be called two times in optimal implementation, where memory requirement would be the same as for Adam optimizer. Over the whole paper, we refer to the optimal implementation as the one that minimizes recomputations, avoids costly statistics during training, and is machine-code optimized to the same extent as optimizers from mainstream libraries. The lines marked as redundant within comments in curly brackets are unnecessary for the optimal operation of the below pseudocode.

---

**Input:** *model*: Neural Network Model  
*dataset*: Training Dataset  
*lossFn*: Loss Function  
*optimizer*: Optimizer

**for all** *batch*  $\in$  *dataset* **do**

*modelOutput*  $\leftarrow$  *model*(*batch.x*) {It is assumed that *model*'s layers' results are kept inside *model*}

*modelLoss*  $\leftarrow$  *LossFn*(*modelOutput*, *batch.y*)

*modelCopy*  $\leftarrow$  *model* {Copy *model*}

*modelCopyOutput*  $\leftarrow$  *modelCopy*(*batch.x*) {This inference is redundant if *modelCopy* gets also intermediate layers' results copied}

*modelCopyLoss*  $\leftarrow$  *LossFn*(*modelCopyOutput*, *batch.y*) {This computation is also redundant, since it is the same as *modelLoss*}

*Backpropagate*(*modelCopyLoss*) {Compute the gradients using the standard backpropagation procedure. Assume that the gradients are stored inside *modelCopy*}

*optimizer.Step*(*modelCopy*) {Perform weight update on *modelCopy* (using the gradients stored inside *modelCopy*)}

*modelCopy*(*batch.x*) {Execute inference to store new layer-wise results in *modelCopy*}

**AveragedBackpropagation**(*model*, *modelCopy*, *modelLoss*) {The procedure is described as Algorithm 2. The parameters of the *model* are modified within}

**end for**

---

An algorithm version with  $n$  backpropagation iterations computes  $(n - 1)$  times the approximated gradient average, each time based on the previous. The intuition behind this is that a better estimate of the averaged derivatives of nonlinear activations is computed after every iteration (Equation 3; backpropagated according to equations no. 4, 6 and 1). Consequently, each time a more precise estimate of the optimal parameter update  $(\Delta\theta)^*$  is obtained (where optimality means that the average gradient is accurately estimated and the parameter update is compliant with it). Therefore, once again

---

**Algorithm 2** *Averaged Backpropagation Algorithm* calculates the approximation of the average gradient.

---

```

216 procedure AVERAGEDBACKPROPAGATION(
217   model: Neural Network Model
218   modelAfterUpdate: model After Candidate Update of
219     Parameters
220   modelLoss: model's Loss)
221   Backpropagate(modelLoss,
222     model.Layers.Last().Output) {Compute the gradient of modelLoss in terms of the last
223     layer's output. Let us assume that the gradient is assigned to the grad property of the output
224     variable (model.Layers.Last().Output)}
225   for index  $\leftarrow$  (Count(model.Layers) - 1) to 0 do
226      $\theta = \text{model.Layers}[\textit{index}].\theta$  {To simplify notation}
227     if IsNonlinear(model.Layers[index]) then {Calculation of either the gradient or its aver-
228       age, which corresponds to the terms in Equation 6}
229       BackpropagateThroughNonlinearLayer(model.Layers[index].Output,
230         model.Layers[index].Input, modelAfterUpdate.Layers[index].Output,
231         modelAfterUpdate.Layers[index].Input) {Procedure described as Algorithm 3. Let
232         us assume that activations are separate layers (like in equations no. 4 and 6)}
233     else
234       Backpropagate(model.Layers[index].Output, model.Layers[index].Input) {The typ-
235       ical backpropagation procedure. It propagates the gradient through a layer. Let us assume
236       that the gradient is assigned to the grad property of the input variable}
237      $\theta.\textit{averagedGrad} \leftarrow \theta.\textit{grad}$  {In this case, for a linear layer, the gradient is treated as its
238     average (compare equations no. 6 and 5)}
239     end if
240      $\theta' = \text{modelAfterUpdate.Layers}[\textit{index}].\theta$  {Notation simplification}
241      $\theta \leftarrow \theta + |\theta' - \theta| \cdot \textit{sgn}(\theta.\textit{averagedGrad})$  {Update by the absolute value of optimizer's
242     update from Algorithm 1:  $|\theta' - \theta|$ , but in the direction of the approximated gradient average
243      $\textit{sgn}(\theta.\textit{averagedGrad})$ }
244   end for
245 end procedure

```

---

the average gradient can be refined to more precisely match the parameter update  $(\Delta\theta)^*$ , and so on (in a loop).

The  $n$ -iteration version of the method is labeled as Algorithm 4 in Appendix G (for two iterations it is a little slower than Algorithm 1 due to additional model-state copies, moreover it is more complex). The optimal implementation of the  $n$ -iteration algorithm variant would be slightly more than  $n$  times slower than the gradient-based RMSProp training, where  $n$  is the number of iterations. There are exactly  $n$  backward passes, optimally  $n$  inferences, and some additional copy operations  $C$  of a model (optimally  $|C| \leq n + 1$ , but Algorithm 4 in Appendix G is suboptimal in this respect). However, the copies are generally significantly faster than forward or backward passes, because it is just needed to copy blocks of data, that are not bigger than the memory used during forward or backward pass. Furthermore, creating the copies by saving results of weight updates directly into different memory addresses only slightly increases the execution time.

An interesting way of comparing the gradient-based RMSProp optimization with our algorithm is to examine the average loss deltas for all weight updates of both algorithms. The purpose of the evaluation approach is to validate the proof in Appendix B.3. However, such a comparison focuses on loss measurements for a *single batch* each time, making it an imperfect predictor of performance on the *whole dataset*. The first iteration of our method is the gradient-based RMSProp procedure, hence the change in loss for RMSProp  $\Delta_{RMSProp}$  is known for both the same model parameters and data as in the case of the loss delta of our method. Therefore, the sum of loss differences after the updates of both approaches can be easily and measurably compared relatively to the sum of loss

**Algorithm 3** *Backpropagation Through Nonlinear Layer*. It is assumed that each input number influences a corresponding single output scalar. This is because, in the experiments, the only operators assumed to be nonlinear during backpropagation of the gradient average are certain activation functions:  $\mathbb{R} \rightarrow \mathbb{R}$ ).

---

```

274 procedure BACKPROPAGATETHROUGHNONLINEARLAYER(LayerOutput: Tensor
275 LayerInput: Tensor
276 LayerOutputAfterUpdate: Tensor
277 LayerInputAfterUpdate: Tensor)
278   f ← Layer function
279   for all (outputNum, inputNum, outputNumUpdated, inputNumUpdated) ∈ Zip(
280     LayerOutput, LayerInput, LayerOutputAfterUpdate, LayerInputAfterUpdate) do
281     {The commonly used Zip function illustrates iterating through multiple tensors at once}
282     if |inputNumUpdated − inputNum| >  $\epsilon$  then {Check if the difference in the inputs is
283       higher than a tiny constant  $\epsilon$ . The condition prevents division by zero. In the experiments
284        $\epsilon \approx 1.19\text{e-}7$ }
285        $\text{AVG}_{x \in [\text{inputNum}, \text{inputNumUpdated}]} f'(x) = \frac{\text{outputNumUpdated} - \text{outputNum}}{\text{inputNumUpdated} - \text{inputNum}}$  {Eq. 3 and
286       4}
287       inputNum.averagedGrad ←  $\text{AVG}_x f'(x) \cdot \text{outputNum.averagedGrad}$ 
288       {Propagate the average gradient backward using the chain rule. Equations 3 and 4 define
289       the term  $\text{AVG}_{t \in [y_k, y'_k]} \frac{\partial f}{\partial t}$  in Equation 6, which is part of Equation 1}
290     else
291       inputNum.averagedGrad ←  $f'(\text{inputNum}) \cdot \text{outputNum.averagedGrad}$  {In this
292       case  $\text{inputNum} \approx \text{inputNumUpdated}$ , thus  $\text{AVG}_{x \in [\dots]} f'(x) \approx f'(\text{inputNum})$  for
293       activation functions. The backpropagation towards input complies with the chain rule
294       (equations no. 6 and 1)}
295     end if
296   end for
297 end procedure

```

---

deltas of RMSProp:

$$\begin{aligned}
 \mathcal{RD}_{AG, RMSProp} &= \frac{\text{AVG}_{b \in B} \Delta_{AG} - \text{AVG}_{b \in B} \Delta_{RMSProp}}{|\text{AVG}_{b \in B} \Delta_{RMSProp}|} \\
 &= \frac{\sum_{b \in B} (\ell_b(\theta'_{AG,b}) - \ell_b(\theta_b))}{|\sum_{b \in B} (\ell_b(\theta'_{RMSProp,b}) - \ell_b(\theta_b))|} - \text{sgn}(\sum_{b \in B} (\ell_b(\theta'_{RMSProp,b}) - \ell_b(\theta_b)))
 \end{aligned} \tag{7}$$

$\mathcal{RD}_{AG, RMSProp}$  is the relative difference in avg. loss deltas of *RMSProp* and the method based on the average gradient (*AG*). The  $\text{AVG}$  operator denotes the arithmetic average.  $B$  is the set of all batches.  $\Delta_{AG,b}$  is the loss delta assuming a batch  $b$  after our algorithm’s update of model parameters  $\theta_b$  to new values  $\theta'_{AG,b}$ . Notation for RMSProp is analogous.  $\ell_b$  is the loss, assuming data of a batch  $b$ .  $\text{sgn}$  is the sign function.  $\mathcal{RD}$  would not be as useful when using momentum because the metric compares the aggregated loss of a single batch per parameter update, whereas momentum contributes to a decrease in loss over many batches per a single parameter update. Without the momentum,  $\mathcal{RD}$  significantly increases the statistical confidence in comparing training algorithms because, for *the same* model weights, the losses are compared for each weight update. Keeping the same parameter values for each loss delta reduces the variance of  $\mathcal{RD}$ , resulting in a decrease in errors when comparing methods.

## 2.2 MODELS AND TRAINING

Our algorithm was tested on three different models with nonlinear ELU (Clevert et al., 2015) and Tanh activations. Model A has a small number of layers (Table 1), and the second one, Model B, is much deeper, with 30 nonlinear layers (Table 2; not counting max pooling as nonlinear). The third model is a convolutional neural network with 46 nonlinear layers (see Appendix H.1 for the training details of this model). It was assumed that Model A is trained for 15 epochs, while Model B – 500 in the case of the gradient-based RMSProp training, and 300 for our method. Grid search was used to find the optimal learning rates for the standard RMSProp training over the course of all 500 epochs,

while our method was optimized only for 200 epochs (out of 300 during testing). The objective of the hyperparameter search was to minimize the loss that is the smallest over a training. The results of the search for optimal learning rates are shown in Table 3. The epoch counts are tailored to ensure that the training achieves minimal or near-minimal test loss values before the final epoch of the gradient-based RMSProp training. The only loss function used in this research is cross-entropy loss, and the batch size is set to 128 in all experiments.

Table 1: *Model A*

Layers	Output Shape	Parameter Count
Convolution 2D ( $3 \times 3$ )		
+ ELU	(8, 26, 26)	80
Convolution 2D ( $3 \times 3$ )		
+ ELU	(8, 24, 24)	584
Convolution 2D ( $5 \times 5$ )		
+ ELU		
stride = 2 padding = 2	(16, 12, 12)	3216
Convolution 2D ( $3 \times 3$ )		
+ ELU	(16, 10, 10)	2320
Convolution 2D ( $3 \times 3$ )		
+ ELU	(16, 8, 8)	2320
Convolution 2D ( $5 \times 5$ )		
+ ELU		
stride = 2 padding = 2	(16, 4, 4)	6416
Flatten	256	
Linear + Softmax	10	2570
		17506

Table 2: *Model B* is designed to test the performance of our algorithm on deep neural networks to achieve a reasonable time of many trainings for statistical significance of the results. The practicality of the architecture is not prioritized.

Layers	Output Shape	Parameter Count
Convolution 2D ( $3 \times 3$ )		
+ ELU	(8, 26, 26)	80
Max Pooling 2D ( $2 \times 2$ )	(8, 13, 13)	
Convolution 2D ( $3 \times 3$ )		
+ ELU	(16, 11, 11)	1168
Max Pooling 2D ( $2 \times 2$ )	(16, 5, 5)	
Flatten	400	
Linear + Tanh	10	4010
<b>26</b> $\times$		
Linear + Tanh	10	<b>26</b> $\times$ 110
Linear + Softmax	10	330
		8228

Models A and B were trained on two popular image datasets: MNIST (LeCun & Cortes, 2010) and Fashion MNIST (Xiao et al., 2017). Both datasets have the same input size ( $28 \times 28 \times 1$ ), but their image characteristics are *significantly* different. Moreover, since the method does not have any hyperparameters apart from the learning rate, it is less likely to overfit to a specific experimental setup (model, dataset, and learning rate) and show good results on it while experiencing deficient performance on other setups. We further validated the performance of our algorithm on a deep sequential convolutional model using an NLP benchmark, specifically the IMDB dataset. All details are presented in Appendix H.

### 3 RESULTS

For the shallow model A, all of the training algorithms are approximately equal (Fig. 2a, Fig. 2b). The relative difference in summed loss deltas (Equation 7) revealed that the algorithm based on the average gradient is only marginally better than the standard RMSprop according to  $\mathcal{RD} = 1.20e-3 \pm 2.7e-4$  (0.12% faster minimization of loss with 0.027% of SEM error) on MNIST and  $\mathcal{RD} = 5.86e-3 \pm 2.79e-3$  on Fashion MNIST in the case of two iterations. For five iterations,  $\mathcal{RD} = 6.47e-4 \pm 9.8e-5$  on MNIST and  $\mathcal{RD} = 2.37e-3 \pm 4.5e-4$  on Fashion MNIST.

The results of Model B are much more interesting. The version of the algorithm with two iterations is about three times faster at minimizing the median of training losses on both datasets (Fig. 2c; Fig. 2d). Moreover, the mean losses tend to be considerably lower than those for standard RMSProp training, even when repeating the experiments using different weight initialization (see Appendix I). Despite the minority of epochs with high oscillations, the method utilizing the average gradient is approximately two to three times faster in minimizing the mean loss, although this is not clearly visible in the plots. Furthermore, for both versions of our algorithm on both datasets, during from 49.3% to 70% of epochs, the average training loss was lower with statistical significance (SEM) than for the gradient-based RMSProp. Conversely, our algorithm was worse in that respect during from 0.667% to 2.33% of epochs with statistical significance. The average of minimal training losses on MNIST for the five iterations is  $0.0393 \pm 0.0058$ , which is significantly lower than  $0.0883 \pm 0.0117$

Table 3: *Learning rates*. All hyperparameter searches of Model A consist of five trainings for each learning rate (LR), while in the case of Model B, it is one training, unless stated otherwise. For Model B, the losses do not directly predict the performance of the methods, because different epoch counts are used between the methods. The standard error of the mean is used as the confidence range for the losses, while for the LRs, the maximum distance to the next best LRs on both sides represents the errors. The LRs used in the experiments are listed in the "Learning Rate" column.

Dataset	Model	Method	Learning Rate	The Most Important Hyperparameter Search Results [Learning Rate: Avg. of Min. Training Loss]	
MNIST	Model A	RMSProp	$8e-4$	$6e-4 : 8.04e-3; 7e-4 : 6.48e-3; 8e-4 : 5.69e-3$	
		2 Iterations	$8e-4$	$9e-4 : 5.94e-3; 10e-4 : 7.70e-3; 11e-4 : 7.39e-3$	
		5 Iterations	$8e-4$	$8e-4 : 0.00555; 9e-4 : 0.00692; 1e-3 : 0.00832$	
				$8e-4 : 0.00514; 9e-4 : 0.00580; 1e-3 : 0.00678$	
				$1.5e-4 : 0.194; 2e-4 : 0.0979; 2.5e-4 : 0.0651$	
				$3e-4 : 0.0683; 3.5e-4 : 0.191; 4e-4 : 0.0759$	
				The best learning rate of the search <i>after</i> the experiments (10 trainings per LR in $\{1.5e-4, 2e-4, \dots, 5.5e-4\}$ ): $(3.5e-4 \pm 1.5e-4) : (0.0856 \pm 0.0139)$ , (matches the performance in our experiments in Section 3)	
				The loss for a high learning rate (10 trainings):	
		Model B	RMSProp	$2.5e-4$	$1.5e-3 : (2.09 \pm 0.05)$
			2 Iterations	$7.5e-4$	The learning rate is guessed
		5 Iterations	$7.5e-4$	The learning rate is guessed	
Fashion MNIST	Model A	RMSProp	$1.5e-3$	$1e-3 : 0.201; 1.25e-3 : 0.186; 1.5e-3 : 0.183$	
		2 Iterations	$1.9e-3$	$1.75e-3 : 0.189; 2e-3 : 0.183; 2.25e-3 : 0.193$	
		5 Iterations	$1.5e-3$	$1.8e-3 : 0.186; 1.9e-3 : 0.179; 2e-3 : 0.180$	
				$1.5e-3 : 0.178; 1.6e-3 : 0.179; 1.7e-3 : 0.200$	
				$2e-4 : 0.356; 2.5e-4 : 0.331; 3e-4 : 0.285$	
				$3.5e-4 : 0.349; 4e-4 : 0.487; 4.5e-4 : 0.459$	
				The best learning rate of the search <i>after</i> the experiments (10 trainings per LR in $\{2e-4, 2.5e-4, \dots, 6e-4\}$ ): $(4e-4 \pm 1.5e-4) : (0.318 \pm 0.016)$ , (matches the performance in our experiments in Section 3)	
				The loss for a high learning rate (10 trainings):	
		Model B	RMSProp	$3e-4$	$9e-4 : (0.641 \pm 0.168)$
			2 Iterations	$9e-4$	$6e-4 : 0.330; 9e-4 : 0.242; 1.2e-3 : 0.355$
		5 Iterations	$9e-4$	$9e-4 : 0.243; 1.5e-3 : 0.276$	

for the standard RMSProp. Meanwhile, the two iterations are also perform better than the gradient-based RMSProp, but without statistical significance, achieving  $0.0747 \pm 0.0188$ . Even better averages of minimal training losses were obtained on Fashion MNIST, with the five-iteration and two-iteration versions achieving  $0.254 \pm 0.017$  and  $0.257 \pm 0.014$  respectively, compared to  $0.314 \pm 0.008$  by the gradient-based training.

Plots of the test losses of Model B look very similar to the training losses (Appendix C), showing significant improvements in generalization, which correspond to the lower training losses. On MNIST, the average of best accuracies over training for five iterations is equal to  $(97.87 \pm 0.09)\%$ , which is significantly higher than  $(96.80 \pm 0.78)\%$  and  $(96.75 \pm 0.55)\%$  for the two-iteration version and gradient-based algorithm, respectively. On Fashion MNIST, the analogous results are  $(88.09 \pm 0.35)\%$ ,  $(87.54 \pm 0.55)\%$  and  $(86.57 \pm 0.29)\%$ , respectively. Appendix D presents the accuracy plots.

For Model B, the  $\mathcal{RD}$  metric (Equation 7) provides a very high confidence of superiority of the average gradient for the high learning rates used for the trainings based on our algorithm (Table 3). On MNIST for two and five iterations, it equals  $10.41 \pm 1.94$  and  $1.43 \pm 0.29$ , respectively. On Fashion MNIST, it is  $0.58 \pm 0.14$  and  $0.24 \pm 0.04$  for both variants, respectively.

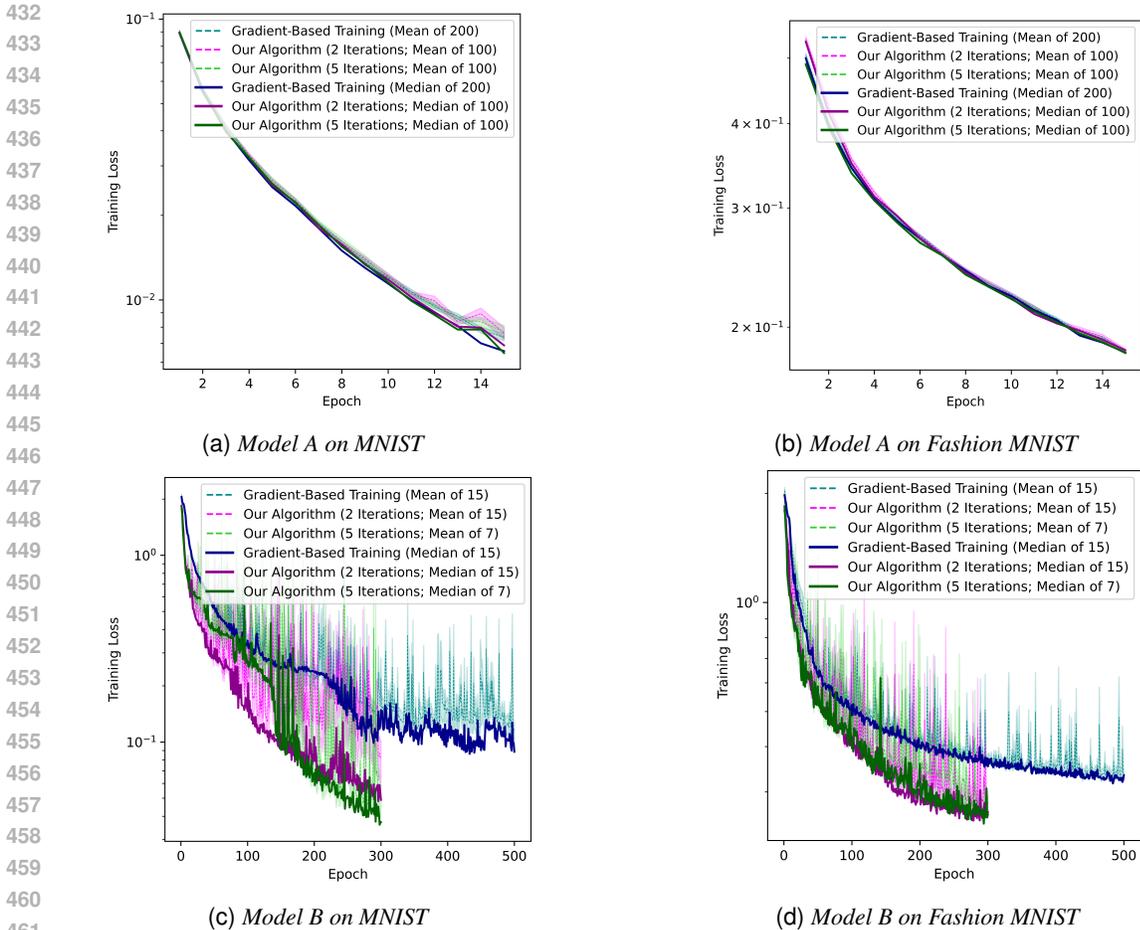


Figure 2: Training losses. Only mean curves contain confidence ranges (SEM).

Importantly, for Model B, the average-gradient algorithm dominated also for the learning rates that are optimal for the standard RMSProp training. Multiple metrics favored our algorithm with statistical significance, i.e.,  $\mathcal{RD} \in [0.0611 \pm 0.0004, 1.07 \pm 0.31]$ , despite training counts equal to only two or three (for each of the four experiments).

In the case of Model B, our implementation of the two-iteration variant of the algorithm based on the average gradient (Alg. 1) is nearly three times slower per epoch than the training based on the gradient, while the five iterations (Alg. 4 in Appendix G) are almost eight times slower per epoch. The estimated runtime of optimal implementation is slightly more than two times longer for the two iterations per epoch when compared to the gradient-based RMSProp, and around six to seven times longer for the five iterations.

For the deep convolutional model on the IMDB dataset, sample efficiency of the two iteration variant of our algorithm achieved about 55% gain in sample efficiency compared to the gradient baseline. The analogous result using four iterations falls between 25% and 30%. See Appendix H.2 for details.

## 4 CONCLUSIONS

Surprisingly, modifying the gradient on nonlinear activations in very deep models can significantly increase sample efficiency for some deep models, which is a direct conclusion of our experiments. For the MNIST and Fashion MNIST benchmarks, the algorithm based on the average gradient offers significant benefits compared to the standard RMSProp training for the deep model with many stacked fully-connected layers and nonlinear activations: (a) About a threefold increase in sample efficiency

486 in terms of median loss, and about two to three times faster mean loss reduction. This is reached  
 487 by only two iterations, which optimally require a little more than double the time of computation  
 488 per epoch in comparison with the gradient-based RMSProp training. Meanwhile, our suboptimal  
 489 implementation of the two-iteration version of the algorithm needs nearly three times more runtime  
 490 per epoch than the training based on the gradient. Therefore, the presented method is not only more  
 491 sample-efficient, but it is also faster and saves energy. (b) Outstanding performance on higher  
 492 learning rates, which may offer significant benefits in terms of both electricity and time spent on  
 493 hyperparameter searches. (c) Considerably better generalization, at least in a reasonable epoch  
 494 count. The increase in sample efficiency and good performance across a wider range of learning  
 495 rates is confirmed by experiments using different weight initialization (see Appendix I).

496 On the other hand, for a deep sequential convolutional model trained on the IMDB dataset, sample  
 497 efficiency is improved by about 55% when using only two iterations of our algorithm (Appendix H.2).  
 498 This is the only significant benefit of our algorithm in this experiment, as the variant using more  
 499 iterations achieved efficiency between that of the vanilla RMSProp and the two-iteration variant.

500 The  $\mathcal{RD}$  (Equation 7) confirms the outstanding results of the other measures. The score of  $\mathcal{RD} =$   
 501  $10.41 \pm 1.94$ , achieved by the two iterations on MNIST, corresponds to the average speed of batch-loss  
 502 minimization that is  $(1141 \pm 194)\%$  of the speed of the gradient-based RMSProp while using the  
 503 same absolute values of weight updates. In the other cases of deep models, the average speed of  
 504 batch-loss minimization ranges from  $(2.10 \pm 0.18)\%$  to  $(243 \pm 29)\%$ . Therefore, even a relatively  
 505 slight speedup in batch-loss minimization (such as 2.1% on the IMDB dataset) can contribute to  
 506 a significantly higher gain in sample efficiency. Moreover, it is crucial to note that the highest of  
 507 the mentioned gains occur at learning rates that are three times higher than the optimal rates for  
 508 gradient-based training. Generally, high learning rate values may enable rapid learning because model  
 509 parameters are adjusted faster. Nevertheless, the average gradient is also superior in terms of the  
 510 average speed of batch-loss minimization when using the optimal learning rates for gradient-based  
 511 training across all tested models, with statistical significance. This validates the proof in Appendix B,  
 512 as both the metric and the proof focus on the efficiency of batch-loss minimization. On the other hand,  
 513 refer to Appendix F for the limitations of our algorithm in estimating the locally optimal update.

514 Surprisingly, the algorithm version with five iterations is worse than the two iterations according to  
 515  $\mathcal{RD}$  with higher statistical confidence than for other measures. Across all experiments, the variant is  
 516 computationally inefficient in terms of the resources required to reduce the loss to a certain level.

517 In the case of the shallow model with nonlinear ELU activations, the method is only marginally better  
 518 than the standard gradient-based RMSProp training. This behavior is expected due to the scaling  
 519 properties of the algorithm (Appendix E).

## 521 5 DISCUSSION

522  
 523 The successful evaluation using different weight initialization techniques on the NLP and computer  
 524 vision benchmarks, using both deep convolutional architecture and the model based on fully-connected  
 525 layers with nonlinear activations, provides insight into significant improvements in sample efficiency,  
 526 at least for some models. Furthermore, the computational cost associated with these improvements is  
 527 modest. These results are especially important in the field of online learning, where sample efficiency  
 528 is crucial.

529 For very deep models without residual connections, gradient-based training tends to be inefficient  
 530 (Balduzzi et al., 2017), which we demonstrate how to mitigate. In general, the very deep structure  
 531 of human brains enables the learning of universal and complex patterns. Therefore, accurately  
 532 mimicking human brain model could potentially lead to satisfactory results. Our algorithm aims to  
 533 improve learning in scenarios involving neural structures that are very deep, a feature of provably  
 534 efficient biological brains that distinguishes them from current AI models. Therefore, the method  
 535 may contribute to the training of large models in the future, where sample efficiency is needed to  
 536 learn new tasks on the fly, akin to how people or some animals do.

537 However, at present, the potential modifications to the algorithm are even more intriguing. Not only  
 538 is it possible to efficiently calculate the average gradient for linear layers using Eq. 5 instead of Eq. 6,  
 539 but Eq. 1 can also be utilized to compute the average gradients over a much larger range than that of  
 a parameter update to capture the global trend of the loss landscape (see Appendix J for future work).

## 6 REPRODUCIBILITY STATEMENT

We put emphasis on providing detailed descriptions of all experiments. The algorithms (Alg. 4 in Appendix G and Alg. 1 in Section 2, with subprocedures labeled as Alg. 2 and Alg. 3) are described in detail in Section 2.1. The models (Tables 1, 2 and 4 in Appendix H), the learning rates (Tables 3 and 5 in Appendix H), and all other important experiment settings are described in Section 2.2 and Appendix H.1. The code, along with environment settings, is available under [...]. Appendix B contains one of our most important theoretical results: the proof of Equation 1 and its superiority over the gradient in minimizing the batch loss by accurately indicating how each model parameter individually contributes to the change in the batch loss (Equation 14). The proven potential for batch-loss minimization is verified not only by the  $\mathcal{RD}$  metric with high statistical significance but also by comparisons of training losses and other metrics (Section 3).

## REFERENCES

- David Balduzzi, Marcus Frean, Lennox Leary, JP Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? In *International conference on machine learning*, pp. 342–350. PMLR, 2017.
- Albert S Berahas, Jorge Nocedal, and Martin Takáč. A multi-batch l-bfgs method for machine learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- Timothy Dozat. Incorporating nesterov momentum into adam. 2016.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- Deborah Hughes-Hallett, Andrew M Gleason, Patti Frazer Lock, and Daniel E Flath. *Applied calculus*. John Wiley & Sons, 2021.
- Nikhil Ketkar. Stochastic gradient descent. *Deep learning with Python: A hands-on introduction*, pp. 113–132, 2017.
- Saeed Khorram, Tyler Lawson, and Li Fuxin. igos++ integrated gradient optimized saliency by bilateral perturbations. In *Proceedings of the Conference on Health, Inference, and Learning*, pp. 174–182, 2021.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 12 2014.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and diversity. *arXiv preprint arXiv:2310.06452*, 2023.
- Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- Yanli Liu, Yuan Gao, and Wotao Yin. An improved analysis of stochastic gradient descent with momentum. *Advances in Neural Information Processing Systems*, 33:18261–18271, 2020.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing lstm language models. *arXiv preprint arXiv:1708.02182*, 2017.
- R OpenAI. Gpt-4 technical report. *ArXiv*, 2303, 2023.

594 Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word  
595 representation. In *Proceedings of the 2014 conference on empirical methods in natural language*  
596 *processing (EMNLP)*, pp. 1532–1543, 2014.

597  
598 Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM*  
599 *journal on control and optimization*, 30(4):838–855, 1992.

600  
601 David Ruppert. Efficient estimations from a slowly convergent robbins-monro process. Technical  
602 report, Cornell University Operations Research and Industrial Engineering, 1988.

603  
604 Sam Sattarzadeh, Mahesh Sudhakar, Konstantinos N Plataniotis, Jongseong Jang, Yeonjeong Jeong,  
605 and Hyunwoo Kim. Integrated grad-cam: Sensitivity-aware visual explanation of deep convolu-  
606 tional networks via integrated gradient-based scoring. In *ICASSP 2021-2021 IEEE International*  
607 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1775–1779. IEEE, 2021.

608  
609 Xu Sun, Hisashi Kashima, Takuya Matsuzaki, and Naonori Ueda. Averaged stochastic gradient  
610 descent with feedback: An accurate, robust, and fast training method. In *2010 IEEE international*  
611 *conference on data mining*, pp. 1067–1072. IEEE, 2010.

612  
613 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
614 *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.

615  
616 Hong Hui Tan and King Hann Lim. Review of second-order optimization techniques in artificial  
617 neural networks backpropagation. In *IOP conference series: materials science and engineering*,  
618 volume 495, pp. 012003. IOP Publishing, 2019.

619  
620 Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks.  
621 In *International conference on machine learning*, pp. 6105–6114. PMLR, 2019.

622  
623 Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running  
624 average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31,  
625 2012.

626  
627 Andreas Veit, Michael J Wilber, and Serge Belongie. Residual networks behave like ensembles of  
628 relatively shallow networks. *Advances in neural information processing systems*, 29, 2016.

629  
630 Ziyang Wei, Wanrong Zhu, and Wei Biao Wu. Weighted averaged stochastic gradient descent:  
631 Asymptotic normality and optimality. *arXiv preprint arXiv:2307.06915*, 2023.

632  
633 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking  
634 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

## 635 A DEFINITION OF AVERAGE GRADIENT/JACOBIAN

636  
637 Let us define the average gradient of a function  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  for some row vector  $\mathbf{x} \in [\mathbf{a}, \mathbf{b}]$  (the  
638 formula is analogous to the one-dimensional case in Equation 3):

$$639 \text{AVG}_{\mathbf{x} \in [\mathbf{a}, \mathbf{b}]} \nabla_{\mathbf{x}} f = (\mathbf{b} - \mathbf{a})^{\circ -1} \circ \int_{\mathbf{a}}^{\mathbf{b}} \nabla_{\mathbf{x}} f \, d\mathbf{x} = (\mathbf{b} - \mathbf{a})^{\circ -1} \circ \int_0^1 \nabla_{\mathbf{a} + t \cdot (\mathbf{b} - \mathbf{a})} f \, dt \quad (8)$$

642  
643 where  $\circ$  denotes the elementwise operation of either multiplication or inversion  $((\cdot)^{\circ -1})$ . However,  
644 the cases of vector elements where division by zero occurs are handled differently, using the partial  
645 derivative  $\frac{\partial f}{\partial x_i}$ :

$$646 \forall i : b_i - a_i = 0 \implies \text{AVG}_{x_i \in [a_i, b_i]} \frac{\partial f}{\partial x_i} = \frac{\partial f}{\partial a_i} \quad (9)$$

If  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , then using to Equation 8:

$$\begin{aligned} \mathcal{AVG}_{\mathbf{x} \in [a, b]} \frac{\partial \mathbf{f}}{\partial \mathbf{x}} &= \begin{bmatrix} \mathcal{AVG}_{\mathbf{x} \in [a, b]} \nabla_{\mathbf{x}} f_1 \\ \mathcal{AVG}_{\mathbf{x} \in [a, b]} \nabla_{\mathbf{x}} f_2 \\ \dots \\ \mathcal{AVG}_{\mathbf{x} \in [a, b]} \nabla_{\mathbf{x}} f_m \end{bmatrix} = \begin{bmatrix} (\mathbf{b} - \mathbf{a})^{\circ-1} \circ \int_a^b \nabla_{\mathbf{x}} f_1 d\mathbf{x} \\ (\mathbf{b} - \mathbf{a})^{\circ-1} \circ \int_a^b \nabla_{\mathbf{x}} f_2 d\mathbf{x} \\ \dots \\ (\mathbf{b} - \mathbf{a})^{\circ-1} \circ \int_a^b \nabla_{\mathbf{x}} f_m d\mathbf{x} \end{bmatrix} \\ &= \begin{bmatrix} (\mathbf{b} - \mathbf{a})^{\circ-1} \\ (\mathbf{b} - \mathbf{a})^{\circ-1} \\ \dots \\ (\mathbf{b} - \mathbf{a})^{\circ-1} \end{bmatrix} \circ \begin{bmatrix} \int_0^1 \nabla_{\mathbf{a}+t(\mathbf{b}-\mathbf{a})} f_1 dt \\ \int_0^1 \nabla_{\mathbf{a}+t(\mathbf{b}-\mathbf{a})} f_2 dt \\ \dots \\ \int_0^1 \nabla_{\mathbf{a}+t(\mathbf{b}-\mathbf{a})} f_m dt \end{bmatrix} \end{aligned} \quad (10)$$

Again, the cases of vector elements where division by zero occurs are handled as follows:

$$\forall i : b_i - a_i = 0 \implies \mathcal{AVG}_{x_i \in [a_i, b_i]} \frac{\partial \mathbf{f}}{\partial x_i} = \frac{\partial \mathbf{f}}{\partial a_i} \quad (11)$$

## B PROOF OF EQUATION 1 AND ITS LOSS-MINIMIZATION POTENTIAL

### B.1 DEFINITION AND PROPERTIES OF AVERAGE GRADIENT OF LOSS

Using Equation 2, the average gradient  $\mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell$  can be defined without the approximation given in Equation 1:

$$\mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell = \mathcal{AVG}_{(\theta_k, \mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n)} \left( \frac{\partial \mathbf{x}_k}{\partial \theta_k} \cdot \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k} \cdot \dots \cdot \frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_{n-1}} \cdot \nabla_{\mathbf{x}_n} \ell \right) \quad (12)$$

where multiple variables are under the average operator  $(\theta_k, \mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$ . There are numerous ways to define how  $(\mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$  depend on the weights and biases  $\theta_k$ , as they all change together during a parameter update. To compute the average (Equation 12), it can be assumed that the parameters of the layer no.  $k$  and the outputs of the layers change linearly with respect to each other, as if they move from  $\theta_k$  to  $\theta'_k$  and from  $(\mathbf{x}_k, \dots, \mathbf{x}_n)$  to  $(\mathbf{x}'_k, \dots, \mathbf{x}'_n)$  after an update of the parameters of all layers. Under this assumption, the calculation is formulated as follows: while computing the average, the integral contains a function  $f_{\theta, k}(t) = \theta_k + t \cdot (\theta'_k - \theta_k)$  for the variable under integration  $t \in [0, 1]$  ( $\theta_k$  and  $\theta'_k$  denote model parameters before and after an update, respectively). Moreover, the integral involves each layer's output:  $\mathbf{f}_{\mathbf{x}, i}(t) = \mathbf{x}_i + t \cdot (\mathbf{x}'_i - \mathbf{x}_i)$ . Finally, the average gradient (Equation 12) is equal to:

$$\begin{aligned} \mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell &= \mathcal{AVG}_{f_{\theta, k}} \nabla_{f_{\theta, k}} \ell = \mathcal{AVG}_t \left( \frac{\partial \mathbf{f}_{\mathbf{x}, k}(t)}{\partial f_{\theta, k}(t)} \cdot \frac{\partial \mathbf{f}_{\mathbf{x}, k+1}(t)}{\partial \mathbf{f}_{\mathbf{x}, k}(t)} \cdot \dots \cdot \frac{\partial \mathbf{f}_{\mathbf{x}, n}(t)}{\partial \mathbf{f}_{\mathbf{x}, n-1}(t)} \cdot \nabla_{\mathbf{f}_{\mathbf{x}, n}(t)} \ell(t) \right) \\ &= \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x}, k}(t)}{\partial f_{\theta, k}(t)} \cdot \frac{\partial \mathbf{f}_{\mathbf{x}, k+1}(t)}{\partial \mathbf{f}_{\mathbf{x}, k}(t)} \cdot \dots \cdot \frac{\partial \mathbf{f}_{\mathbf{x}, n}(t)}{\partial \mathbf{f}_{\mathbf{x}, n-1}(t)} \cdot \nabla_{\mathbf{f}_{\mathbf{x}, n}(t)} \ell(t) dt \end{aligned} \quad (13)$$

which is more direct and easier to work with.

Importantly, unlike the gradient, the average gradient ( $\mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell$ ) is directly proportional to the loss-change impact of each model parameter separately  $\mathbf{l}_{\theta', k} - \mathbf{l}_{\theta, k}$  (of the shape of  $\theta_k$  and  $\theta'_k$ , unlike the scalar  $\ell$ ):

$$\begin{aligned} \mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell &= \mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \left( \sum_{j=0}^n \sum_{i=0}^{|\theta_j|} \ell_{\theta, j, i} \right) = \mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \left( \sum_{i=0}^{|\theta_k|} \ell_{\theta, k, i} \right) = \mathcal{AVG}_{\theta_k} \left( \text{diag} \left( \frac{\partial \mathbf{l}_{\theta, k}}{\partial \theta_k} \right) \right) = \\ &= \left\langle \mathcal{AVG}_{\theta_{k,1}} \ell'_{\theta, k, 1}, \dots, \mathcal{AVG}_{\theta_{k,n}} \ell'_{\theta, k, n} \right\rangle = \left\langle \frac{\int_{\theta_{k,1}}^{\theta'_{k,1}} \ell'_{\theta, k, 1} d\theta}{\theta'_{k,1} - \theta_{k,1}}, \dots, \frac{\int_{\theta_{k,n}}^{\theta'_{k,n}} \ell'_{\theta, k, n} d\theta}{\theta'_{k,n} - \theta_{k,n}} \right\rangle \\ &= \left\langle \frac{\ell_{\theta', k, 1} - \ell_{\theta, k, 1}}{\theta'_{k,1} - \theta_{k,1}}, \dots, \frac{\ell_{\theta', k, n} - \ell_{\theta, k, n}}{\theta'_{k,n} - \theta_{k,n}} \right\rangle = (\theta'_k - \theta_k)^{\circ-1} \circ (\mathbf{l}_{\theta', k} - \mathbf{l}_{\theta, k}) \propto \mathbf{l}_{\theta', k} - \mathbf{l}_{\theta, k} \end{aligned} \quad (14)$$

where  $\circ$  denotes the elementwise operation of either multiplication or inversion ( $(\cdot)^{\circ-1}$ ).  $\text{diag}(\frac{\partial \mathbf{l}_{\theta,k}}{\partial \theta_k})$  denotes diagonal elements of the Jacobian matrix.  $\ell_{\theta,k,i} \in \mathbf{l}_{\theta,k}$  represents the scalar loss contribution of a single model parameter  $(\theta_{k,i})$ , that can be defined as an integral of the gradient:  $\ell_{\theta,k,i} = \int_{C_1}^{\theta_{k,i}} \nabla_{\theta} \ell_{\theta} d\theta + C_2$ , for any constant scalars  $C_1$  and  $C_2$ . (Note that in this case,  $\ell_{\theta,k,i} \neq \ell_{\theta} + C_{\ell}$ , for any constant  $C_{\ell}$ , because the loss  $\ell$  also depends on other parameters than  $\theta_{k,i}$ .) Important properties: (a)  $\ell_{\theta} = C + \sum_{k=0}^n \sum_{i=0}^{|\theta_k|} \ell_{\theta,k,i}$  for a constant  $C$  that is invariant across updates of the model parameters  $\theta$ . (b) The elements of  $\mathbf{l}$  are related to the difference in loss during parameter update:  $\ell_{\theta'} - \ell_{\theta} = (\sum_{k=0}^n \sum_{i=0}^{|\theta'_k|} \ell_{\theta',k,i}) - (\sum_{k=0}^n \sum_{i=0}^{|\theta_k|} \ell_{\theta,k,i})$ . (c) The following equation is satisfied:  $\nabla_{\theta} \ell = \nabla_{\theta} (\sum_{k=0}^n \sum_{i=0}^{|\theta_k|} \ell_{\theta,k,i})$ . The simple one-dimensional visualization of the proportionality from Equation 14 ( $\mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell \propto \mathbf{l}_{\theta',k} - \mathbf{l}_{\theta,k}$ ) is shown in Fig. 1. Note that the property of proportionality does not hold for the gradient updates (which are utilized by Adam (Kingma & Ba, 2014), RMSProp (Tieleman et al., 2012), and SGD (Ketkar, 2017; Liu et al., 2020)). In the gradient case, during the update step of  $\theta$  weights,  $\theta'$  is not used in the calculation of itself. Therefore,  $\mathbf{l}_{\theta'} - \mathbf{l}_{\theta}$  cannot be computed yet, and the accurate influence on loss remains unknown, unlike for the average gradient (Equation 14). The cases of scalar parameters  $\theta_{k,i} \in \theta_k$  and  $\theta'_{k,i} \in \theta'_k$  where division by zero occurs are handled differently:

$$\forall i : \theta'_{k,i} - \theta_{k,i} = 0 \implies \mathcal{AVG}_{\theta_{k,i} \in [\theta_{k,i}, \theta'_{k,i}]} \frac{\partial \ell}{\partial \theta_{k,i}} = \frac{\partial \ell}{\partial \theta_{k,i}} \quad (15)$$

Assuming the functions  $f_{\theta,k}$  and  $f_{\mathbf{x},i}$  from Equation 13 are any functions (but differentiable with respect to each other), Equation 14 remains valid. Therefore, the crucial property of direct proportionality to the loss values does not depend on our previous assumptions about  $\theta_k$  and  $\mathbf{x}_i$ . The purpose of these assumptions is to provide a simple example, reduce reasoning abstraction, and simplify further proofs in Sections B.2 and B.3.

## B.2 PROOF OF OF EQUATION 1 WITHOUT SPECIFYING PRECISION OF APPROXIMATION

For some function  $f$  and some constants  $C_1, C_2, \dots, C_n$ :

$$\int C_1 \cdot C_2 \cdot \dots \cdot C_n \cdot f(x) dx = C_1 \cdot C_2 \cdot \dots \cdot C_n \cdot \int f(x) dx \quad (16)$$

Similarly, let us denote approximately constant functions as  $C'_1(x) \cong C_1, C'_2(x) \cong C_2, \dots, C'_n(x) \cong C_n$  for some  $x \in [a, b], a \neq b$ . The constant that precisely approximates each function  $C'_i(x)$ , is its average:  $C'_1(x) \cong \mathcal{AVG} C'_1(x) = C_1, C'_2(x) \cong \mathcal{AVG} C'_2(x) = C_2, \dots, C'_n(x) \cong \mathcal{AVG} C'_n(x) = C_n$ . Therefore, similarly to Equation 16:

$$\begin{aligned} \int_a^b C'_1(x) \cdot \dots \cdot C'_n(x) \cdot f(x) dx &\cong \mathcal{AVG}_{x \in [a,b]} C'_1(x) \cdot \dots \cdot \mathcal{AVG}_{x \in [a,b]} C'_n(x) \cdot \int_a^b f(x) dx \\ \int_a^b C'_1(x) \cdot \dots \cdot C'_n(x) \cdot f(x) dx &\cong \int_a^b \frac{C'_1(x)}{b-a} dx \cdot \dots \cdot \int_a^b \frac{C'_n(x)}{b-a} dx \cdot \int_a^b f(x) dx \end{aligned} \quad (17)$$

which is also approximately equal to both sides of Equation 16. In Equation 17, both approximations are equivalent, because  $\mathcal{AVG} C'_i(x) = \int_a^b C'_i(x) / (b-a) dx$ . For functions  $\mathbb{R}^n \rightarrow \mathbb{R}^m$ , equations no. 16 and 17 are analogous. Note that, in the general case, the different approximations of the terms  $C'_i(x) \cong C'_i(a)$  and  $C'_i(x) \cong C'_i(b)$  are worse than the average:  $C'_i(x) \cong \mathcal{AVG} C'_i(x) = C_i$  (which is used further in Section B.3).

Rapid changes in the gradient over the range of an update indicate that the update step is too large, leading to instability and reduced training effectiveness due to excessively large steps in the loss landscape. We assume effective learning, where gradients do not change significantly<sup>1</sup> between updates, ensuring the learning rate is appropriately sized. In this case, the gradient  $\nabla_{\theta_k} \ell$  does not change significantly<sup>2</sup> over the range of a weight update  $[\theta, \theta']$ . However, these assumptions are

<sup>1</sup>The magnitude of the gradient change need not be specified, as it suffices that it contributes to the approximations with unspecified bounds in Equations 18 and 19. The accuracy of these approximations is proven in Section B.3.

<sup>2</sup>See footnote 1.

merely intended to build intuition and *are not necessary for this proof*. We do not yet assume any specific level of precision in how Equation 17 approximates Equation 13:

$$\begin{aligned}
\mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell &= \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x},k}(t)}{\partial f_{\theta,k}(t)} \cdot \frac{\partial \mathbf{f}_{\mathbf{x},k+1}(t)}{\partial \mathbf{f}_{\mathbf{x},k}(t)} \cdots \frac{\partial \mathbf{f}_{\mathbf{x},n}(t)}{\partial \mathbf{f}_{\mathbf{x},n-1}(t)} \cdot \nabla_{\mathbf{f}_{\mathbf{x},n}(t)} \ell(t) dt \\
&\cong \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x},k}(t)}{\partial f_{\theta,k}(t)} dt \cdot \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x},k+1}(t)}{\partial \mathbf{f}_{\mathbf{x},k}(t)} dt \cdots \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x},n}(t)}{\partial \mathbf{f}_{\mathbf{x},n-1}(t)} dt \cdot \int_0^1 \nabla_{\mathbf{f}_{\mathbf{x},n}(t)} \ell(t) dt \quad (18) \\
&= \mathcal{AVG}_{\theta_k} \frac{\partial \mathbf{x}_k}{\partial \theta_k} \cdot \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k} \cdots \mathcal{AVG}_{\mathbf{x}_{n-1}} \frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_{n-1}} \cdot \mathcal{AVG}_{\mathbf{x}_n} \nabla_{\mathbf{x}_n} \ell
\end{aligned}$$

Applying the notation of Equation 12 to Equation 18, we get:

$$\begin{aligned}
\mathcal{AVG}_{\theta_k} \nabla_{\theta_k} \ell &\cong \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x},k}(t)}{\partial f_{\theta,k}(t)} dt \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x},k+1}(t)}{\partial \mathbf{f}_{\mathbf{x},k}(t)} dt \cdots \int_0^1 \frac{\partial \mathbf{f}_{\mathbf{x},n}(t)}{\partial \mathbf{f}_{\mathbf{x},n-1}(t)} dt \int_0^1 \nabla_{\mathbf{f}_{\mathbf{x},n}(t)} \ell(t) dt \\
&= \int_{\theta_k}^{\theta'_k} \frac{\partial \mathbf{x}_k(\vartheta_k)}{\partial \vartheta_k} d\vartheta_k \int_{\mathbf{x}'_k}^{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}(\chi_k)}{\partial \chi_k} d\chi_k \cdots \int_{\mathbf{x}'_{n-1}}^{\mathbf{x}_{n-1}} \frac{\partial \mathbf{x}_n(\chi_{n-1})}{\partial \chi_{n-1}} d\chi_{n-1} \int_{\mathbf{x}_n}^{\mathbf{x}'_n} \nabla_{\mathbf{x}_n} \ell d\chi_n \\
&= \mathcal{AVG}_{\theta_k} \frac{\partial \mathbf{x}_k}{\partial \theta_k} \cdot \mathcal{AVG}_{\mathbf{x}_k} \frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_k} \cdots \mathcal{AVG}_{\mathbf{x}_{n-1}} \frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_{n-1}} \cdot \mathcal{AVG}_{\mathbf{x}_n} \nabla_{\mathbf{x}_n} \ell \quad (19)
\end{aligned}$$

where  $\theta_k, \mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n$  are all linear functions of  $t$  (previously denoted as  $f_{\theta,k}, f_{\mathbf{x},k}, f_{\mathbf{x},k+1}, \dots, f_{\mathbf{x},n}$ ). Therefore, the functions  $\mathbf{x}_k(\theta_k), \mathbf{x}_{k+1}(\mathbf{x}_k), \dots, \mathbf{x}_n(\mathbf{x}_{n-1})$  are known. The edge cases of those scalars within  $\theta_k, \mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n$  that do not depend on  $t$  are handled analogously to Equation 15, as in these cases the average gradient equals the gradient.

Despite the provided arguments on why the approximation is applied, the precision of the estimation is not specified, although it *is crucial*. Therefore, the accuracy of the approximation is described in Section B.3. Otherwise, if the precision of the estimation is not important, then Equation 19 ultimately proves Equation 1.  $\square$

The analogous reasoning can be applied to prove Equation 5.

In the algorithm, it is also assumed that the average gradient of the loss with respect to the output of the last layer, denoted as  $(\mathcal{AVG}_{\mathbf{x}_n} \nabla_{\mathbf{x}_n} \ell)$ , is replaced by the gradient  $(\nabla_{\mathbf{x}_n} \ell)$ . Moreover, in our implementation, the gradients replace the average gradients of layers that are approximately linear (using Equation 6 instead of Equation 5), resulting in a broader approximation in Equation 1. However, the presented reasoning still applies, including the proof of approximation accuracy in Section B.3. See Appendix F for comments on the limitations of our implementation of Equation 1.

### B.3 PROOF OF SUFFICIENT PRECISION OF APPROXIMATION

Referring to the content of the paragraphs just before and after Equation 17, the approximation in Equation 17 is more precise in the case of  $C'_i(x) \cong \mathcal{AVG} C'_i(x) = C'_i$  than in the case of approximating  $C'_i(x) \cong C'_i(a)$ . The average Jacobian of each term in Equation 1 can be denoted as  $\mathcal{AVG} C'_i(x)$ , while the Jacobian of each term in Equation 2 can be denoted as  $C'_i(a)$ . For the average Jacobian  $\mathcal{AVG} C'_i(x)$ , a better estimation in Equation 17 is obtained, as stated in the text near the equation. Consequently, applying Equation 17 to approximate Equation 13 results in a higher precision in estimating Equation 1 when averaging each Jacobian term separately, compared to utilizing the Jacobians without averaging. Therefore, a better approximation of the accurate average gradient is obtained compared to using the gradient.  $\square$  The average gradient is proportional to the change in loss after the corresponding parameter update (Equation 14). Therefore, approximating the average gradient more precisely than current gradient-based methods can lead to more efficient minimization of batch loss, for example, by using Eq. 1. Therefore, learning can be enhanced compared to the potential of gradient-based methods.

### C TEST LOSS CURVES OF MODEL B

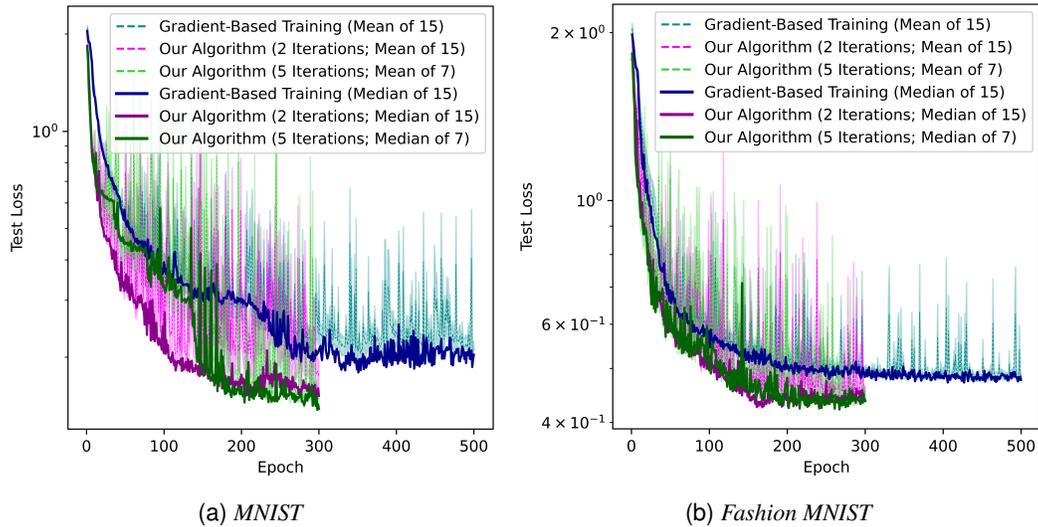


Figure 3: *Test losses of Model B. Only mean curves contain confidence ranges (SEM).*

### D TEST ACCURACY CURVES OF MODEL B

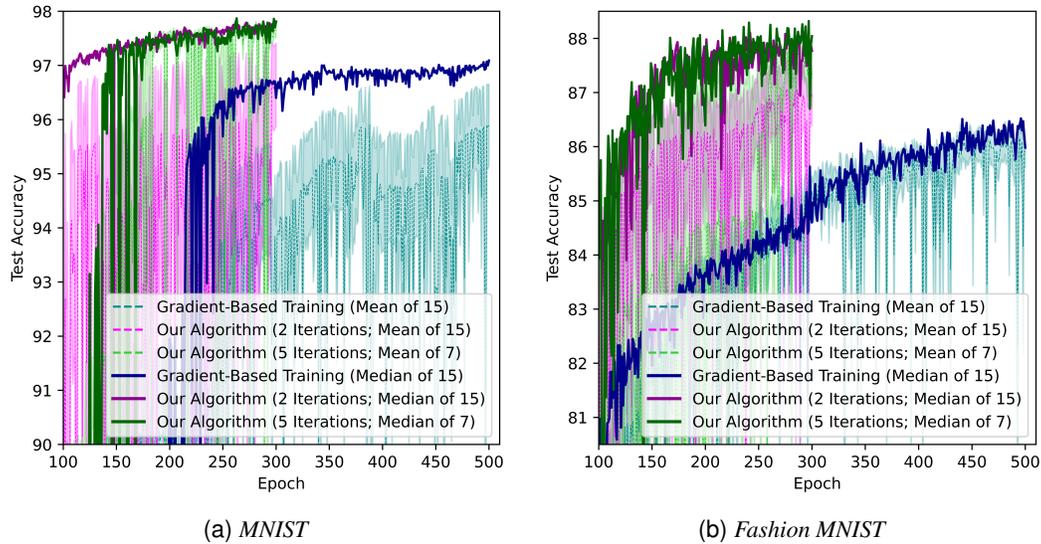


Figure 4: *Test accuracy of Model B. Only mean curves contain confidence ranges (SEM).*

### E SCALING IN TERMS OF MODEL DEPTH

The algorithm based on the average gradient aims to reduce errors of the predicted influence on loss of a parameter update. In the case of the gradient-based approach, the errors arise from the impaired prediction of how inputs to subsequent layers influence their outputs (Fig. 1). Let us model the errors as multiplicative, because each time a fraction of output may be influenced by the error (Balduzzi et al., 2017). Therefore, when compared to the gradient-based algorithm as a baseline, the multiplicative errors are reduced after backpropagation through each nonlinear layer

(by computing the average Jacobian of the layer). Consequently, the incorporation of the average gradient exponentially reduces the error in terms of a count of nonlinear layers (that are involved in the backpropagation process). This explains the huge performance-improvement gap between the models for the method based on the average gradient, which emerges from the difference in models' depths. However, the gap is also increased due to the linearity of the ELU activation function in most of its domain, where the gradient equals its average. In this case, our algorithm produces results similar to those of gradient-based optimization.

If the errors (of the predicted influence on loss of a parameter update) are enormous, then the learning is impossible. Therefore, the learning performance tends to decrease after the error reaches a certain value for a given model, learning rate, and other parameters. From that point onward, our algorithm more efficiently reduces the batch loss compared to the gradient-based approach by minimizing the error in the loss-influence prediction. Importantly, the improvements tend to increase with both the number of nonlinear layers in a model and the learning rate.

## F THREE-DIMENSIONAL COMPARISON OF THE GRADIENT AND THE AVERAGE GRADIENT

In our experiments, during a parameter update, in terms of the average reduction of loss for a batch, our algorithm lies between the gradient (red arrows in Figure 5) and the lowest average gradient (black arrows in Figure 5). Our algorithm does not always find a locally optimal solution (the best in the range of a single parameter update) because:

- a The average gradient is approximated (by using Equation 1 instead of Equation 3, Equation 6 as a substitute of Equation 5, and the non-averaged gradient of the loss with respect to the last layer output).
- b The optimal parameter update may be inaccurately estimated before the average gradient for this parameter update is calculated. Moreover, even after many iterations of Algorithm 4 (Appendix G), the update step may not converge to a locally optimal solution (black vectors in Figure 5).
- c After the first iteration of our algorithm, only the negations of the directions of changes in each parameter are possible. Thus, the search for locally optimal updates is bounded by  $2^{|\Theta|}$  combinations, where  $|\Theta|$  is the count of trainable parameters.

Nevertheless, the  $\mathcal{RD}$  metric (defined in Equation 7) indicates our algorithm minimizes the batch loss more efficiently on average compared to the gradient-based approach.

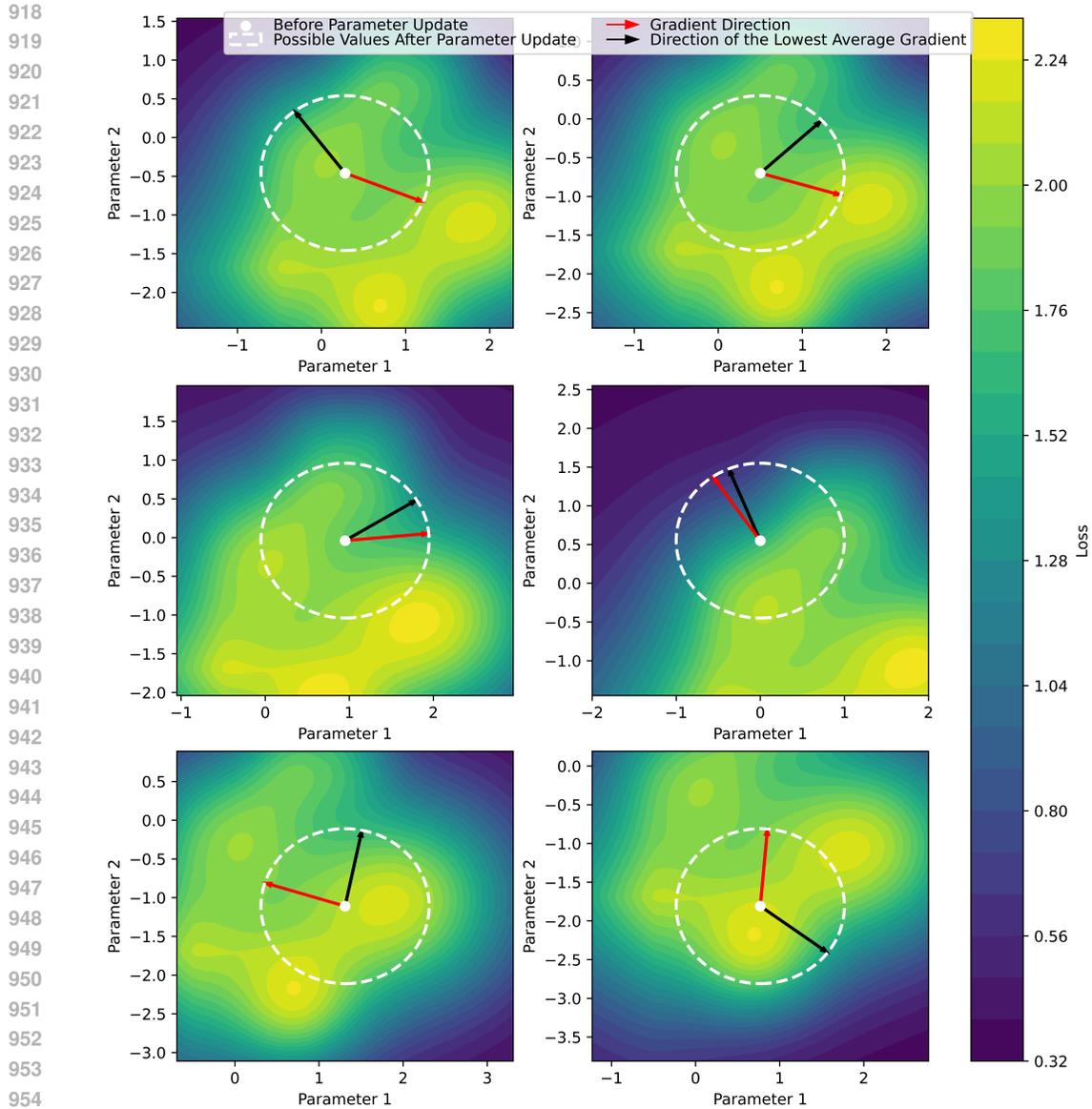


Figure 5: *Three-dimensional comparison of the gradient and the lowest average gradient in a few example scenarios.* The latter accurately reflects the influence on the loss of a parameter update. Furthermore, it accurately shows how each model parameter individually contributes to the change in the batch loss (Equation 14), which is utilized by our algorithm. Each plot illustrates the loss in terms of two example model parameters, assuming a specific magnitude for each parameter update (represented by the radius of each white circle). The arrows point to the loss values after an update based on the gradient and the average gradient. The average gradient is calculated for the update that minimizes it. Therefore, it points to the minimum loss on each white circle, although this minimum is not always achieved by the approximated average gradient computed by our algorithm.

---

## 972 G ALGORITHM VERSION WITH PARAMETERIZED NUMBER OF ITERATIONS

973  
974  
975 **Algorithm 4** *Algorithm Version with Parameterized Number of Iterations* (two or more). The number  
976 of iterations is equal to the number of backpropagation calls and inferences in optimal implementation.  
977 The memory requirement of the ideal implementation would be higher than that of Adam by only an  
978 additional scalar size per parameter of the model.

---

979 **Input:** *model*: Neural Network Model to Train  
980 *dataset*: Training Dataset  
981 *lossFn*: Loss Function  
982 *optimizer*: Optimizer  
983 *iterCount*: Number of Backpropagation Iterations  
984 **for all** *batch*  $\in$  *dataset* **do**  
985     *modelInitial*  $\leftarrow$  *model*  
986     *modelCopy*  $\leftarrow$  *model*  
987     *initialOutput*  $\leftarrow$  *modelCopy*(*batch.x*) {It is assumed that *modelCopy*'s layers' results are  
988     kept inside *modelCopy*}  
989     *initialLoss*  $\leftarrow$  *LossFn*(*initialOutput*, *batch.y*)  
990     *Backpropagate*(*initialLoss*) {Compute the gradients using the standard backpropagation  
991     procedure. Assume that the gradients are stored inside *modelCopy*}  
992     *optimizer.Step*(*modelCopy*) {Parameter update}  
993     *modelOutputAfterUpdate*  $\leftarrow$  *modelCopy*(*batch.x*)  
994     *modelLossAfterUpdate*  $\leftarrow$  *LossFn*(*modelOutputAfterUpdate*, *batch.y*)  
995     **for** *iter* = 1, ..., *iterCount* - 1 **do** {Loop (*iterCount* - 1) times, because one backward  
996     propagation is done}  
997         **if** *iter*  $\neq$  1 **then**  
998             *modelCopy*  $\leftarrow$  *model*  
999             *modelCopy*(*batch.x*) {For each layer, compute its output, and store it inside *modelCopy*}  
1000             *model*  $\leftarrow$  *modelInitial*  
1001             **end if**  
1002             *initialOutput*  $\leftarrow$  *model*(*batch.x*) {This computation is redundant if layer outputs are  
1003             copied from *modelInitial*}  
1004             *initialLoss*  $\leftarrow$  *LossFn*(*initialOutput*, *batch.y*) {Analogously, this computation is also  
1005             redundant}  
1006             *AveragedBackpropagation*(*model*, *modelCopy*, *initialLoss*) {The procedure is de-  
1007             scribed as Algorithm 2. The parameters of the *model* are modified within}  
1008             **end for**  
1009     **end for**

---

## 1009 H CONVOLUTIONAL NEURAL NETWORK ON IMDB

### 1010 H.1 METHODS

1011  
1012  
1013 We refer to Model C (Tab. 4) as our very deep convolutional model, which we tested on the IMDB  
1014 dataset. This model, primarily composed of convolutional layers, is designed to evaluate the perfor-  
1015 mance of our learning algorithm on a deep convolutional neural network without skip connections.  
1016 Skip connections simplify the learning task by enabling the network to leverage features that can be  
1017 extracted by shallower networks (Veit et al., 2016). Our primary goal is to assess the algorithm's  
1018 capabilities, rather than achieving state-of-the-art results.

1019 IMDB preprocessing includes: (a) Equal split for test and training sets. (b) Duplicate removal.  
1020 (c) Punctuation removal. (d) Tokenization. (e) Padding to the length of 122 (mean training-exam-  
1021 ple length), and keeping the final part of each review. (f) Lemmatization. (g) Vectorization using  
1022 GloVe embeddings (Pennington et al., 2014). Finally, the input has the shape (1, 122, 50), where  
1023 each input word is converted into its corresponding GloVe embedding with a length of 50. Model C  
1024 (Tab. 4) utilizes multiple convolutional layers of shape (1  $\times$  1), which are used to change the data  
1025 shape and, for each "pixel", to extract features from the outputs of different filters. The neighboring  
dimensions of each GloVe embedding do not have any special relationship compared to the distinct

Table 4: *Model C*.

Layers	Output Shape	Parameter Count
Convolution 2D ( $1 \times 50$ ), Tanh	(50, 122, 1)	2550
Convolution 2D ( $1 \times 1$ ), Tanh	(40, 122, 1)	2040
Convolution 2D ( $1 \times 1$ ), Tanh	(35, 122, 1)	1435
Convolution 2D ( $1 \times 1$ ), Tanh	(30, 122, 1)	1080
Convolution 2D ( $1 \times 1$ ), Tanh	(27, 122, 1)	837
Convolution 2D ( $1 \times 1$ ), Tanh	(24, 122, 1)	672
Convolution 2D ( $1 \times 1$ ), Tanh	(21, 122, 1)	525
Convolution 2D ( $1 \times 1$ ), Tanh	(18, 122, 1)	396
Convolution 2D ( $1 \times 1$ ), Tanh	(16, 122, 1)	304
Convolution 2D ( $1 \times 1$ ), Tanh	(14, 122, 1)	238
Convolution 2D ( $1 \times 1$ ), Tanh	(12, 122, 1)	180
Convolution 2D ( $1 \times 1$ ), Tanh	(10, 122, 1)	130
Convolution 2D ( $1 \times 1$ ), Tanh	(8, 122, 1)	88
Convolution 2D ( $1 \times 1$ ), Tanh	(6, 122, 1)	54
Convolution 2D ( $1 \times 1$ ), Tanh	(5, 122, 1)	35
Convolution 2D ( $3 \times 1$ ) with stride = 2, Tanh	(5, 60, 1)	80
<b>25×Convolution 2D (<math>3 \times 1</math>), Tanh</b>	<b>(5, 10, 1)</b>	<b>25×80</b>
Flatten	50	
Linear, Tanh	25	1275
Linear, Tanh	13	338
Linear, Tanh	7	98
Linear, Tanh	4	32
Linear, Softmax	2	10
		14397

ones. Therefore, we used convolutions with a filter-size dimension equal to either one or all features in the GloVe embeddings.

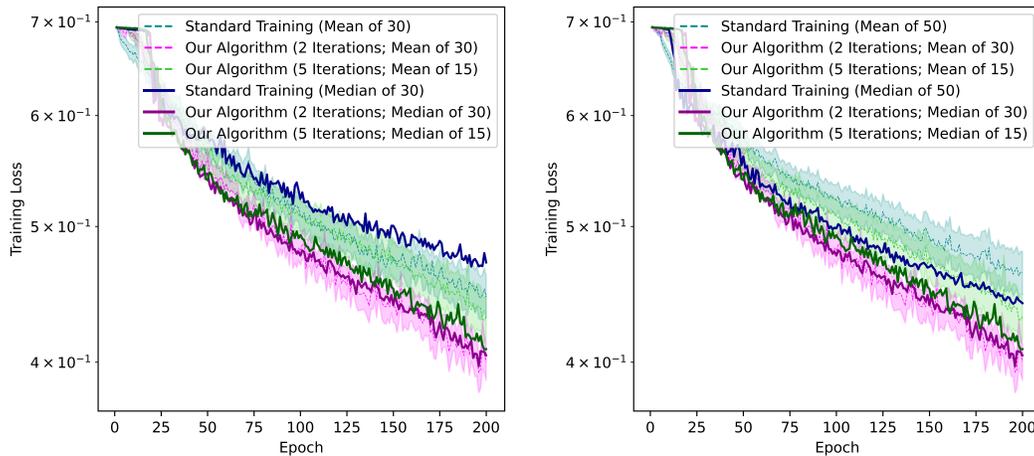
## H.2 RESULTS

Both versions of our algorithm were more sample-efficient than the gradient-based RMSProp, as indicated by the training loss (Fig. 6). In the case of gradient-based training, the trade-off between the mean and median of the training loss is visible in both Figs. 6a and 6b. The tendency for instability in training with a higher learning rate leads to the occurrence of outliers, also in terms of whole worse trainings, which lower the mean. However, the median is resistant to these outliers. This can also be observed in Figures 7b and 8b. To evaluate both the mean and median, considering the trade-off between them, we compared the methods by averaging the median and mean losses. This approach provides a consistent comparison result across both learning rates of the gradient-based RMSProp (Fig. 6). Using this evaluation method, the performance of gradient-based RMSProp at epoch 200 is approximately equal to the results of the two iterations of our method at epochs 125 and 130, in Figs. 6a and 6b, respectively. This translates to a sample efficiency between 53% and 60% higher in favor of the two iterations of our algorithm compared to the vanilla RMSProp. Surprisingly, the performance of the four iterations falls between the other methods, with the sample efficiency gain ranging from 25% to 30%. The  $\mathcal{RD} = 0.0394 \pm 0.0053$  metric also favors the two-iteration variant, outperforming the  $\mathcal{RD} = 0.0210 \pm 0.0018$  achieved by the four iterations.

The test-accuracy (Fig. 8) and test-loss curves (Fig. 7) should be interpreted in the context that the objective of the hyperparameter search is dependent solely on the training loss. In addition, considering the trade-off between the mean and median losses, which occurs between the lower and higher learning rates of the vanilla RMSProp, slightly better mean generalization in the gradient training for the low learning rate (Fig. 7a) does not imply generally better test performance. However, the comparison using the same learning rate (Fig. 7b) indicates that the two-iteration variant achieves the most stable test-loss performance.

Table 5: *Learning rates* for the IMDB dataset. The "Learning Rate" column presents the final chosen learning rates for the experiments. A repeated loss of 0.6931 is equivalent to the lack of training. The table includes the results of the final experiments; however, the results are clipped to 150 epochs for the variants of our algorithm. The best results and the learning rates chosen for the experiments are marked in bold.

Method	Learning Rate	Most Important Hyperparameter Search Results [Learning Rate: Avg. of Min. Training Loss]
		9.545e-5 : 0.6719; 1.193e-4 : 0.5605
		1.491e-4 : 0.4121; 1.864e-4 : 0.5106
		2.330e-4 : 0.3567; 2.912e-4 : 0.4026
	<b>3.641e-4</b> (stable trainings, small number of outliers, low average losses) and	<b>3.641e-4 : 0.3370</b> ; 4.551e-4 : 0.3407
	<b>6.906e-4</b> (slightly higher average training losses, but lower median losses)	5.689e-4 : 0.5243; <b>6.906e-4 : 0.2801</b>
		8.384e-4 : 0.4414; 1.018e-4 : 0.4565
		1.236e-3 : 0.6931; 1.500e-3 : 0.6931
		Repeated trainings:
<b>RMSProp</b> (200 epochs)		<b>3.641e-4 : (0.4307 ± 0.0143)</b> ; (30 trainings)
		<b>6.906e-4 : (0.4435 ± 0.0154)</b> ; (50 trainings)
		3.641e-4 : 0.4948; <b>4.733e-4 : 0.3373</b>
		6.153e-4 : 0.4385; 7.999e-4 : 0.4316
		1.040e-3 : 0.4096; 1.352e-3 : 0.6931
		Repeated trainings:
<b>2 Iterations</b> (150 epochs)	<b>6.906e-4</b> (low training losses, easy to compare with <b>RMSProp</b> due to matching learning rate)	<b>4.733e-4 : (0.4574 ± 0.0127)</b> ; (15 trainings)
	<b>6.906e-4</b> (easy to compare with <b>RMSProp</b> and <b>2 Iterations</b> due to matching learning rates)	<b>6.906e-4 : (0.4225 ± 0.0125)</b> ; (30 trainings)
		<b>4.734e-4 : 0.3766</b> ; 6.153e-4 : 0.4410
		Repeated trainings:
<b>5 Iterations</b> (150 epochs)		<b>6.906e-4 : (0.4567 ± 0.0129)</b> ; (15 trainings)

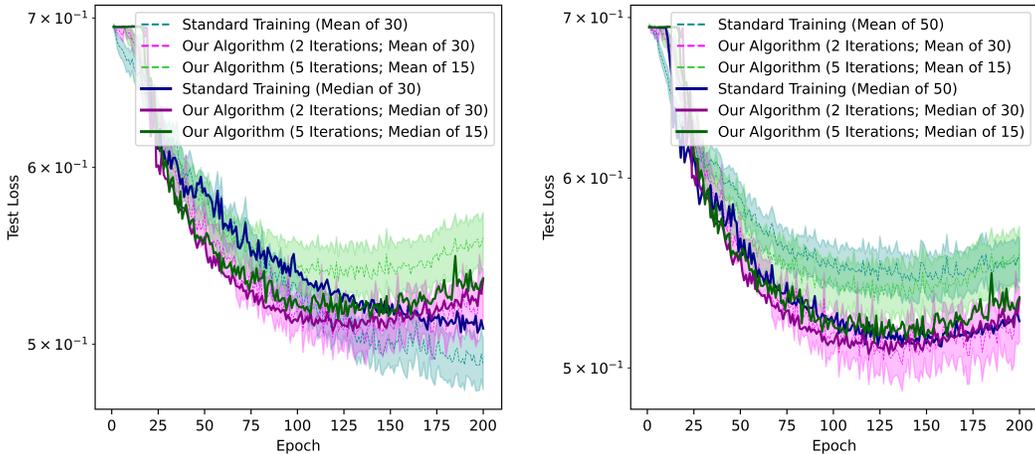


(a) Gradient-based RMSProp with the lower learning rate of **3.641e-4**. (b) Gradient-based RMSProp with the lower learning rate of **6.906e-4**.

Figure 6: *Training loss of Model C*. Only mean curves contain confidence ranges (SEM).

Due to suboptimal backpropagation of the average gradient through activations in our implementation, it has a bigger computational overhead for models applying activations to large feature maps. Therefore, our implementation is computationally slower relatively to the gradient-based training for Model C than in the case of Model B.

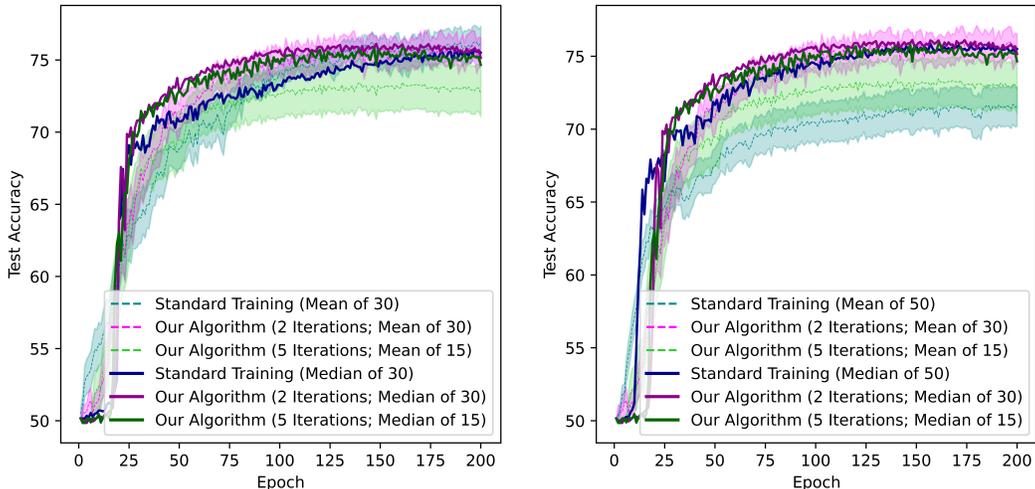
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148



(a) Gradient-based RMSProp with the lower learning rate of  $3.641e-4$ . (b) Gradient-based RMSProp with the lower learning rate of  $6.906e-4$ .

Figure 7: *Test loss of Model C*. Only mean curves contain confidence ranges (SEM).

1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169



(a) Gradient-based RMSProp with the lower learning rate of  $3.641e-4$ . (b) Gradient-based RMSProp with the lower learning rate of  $6.906e-4$ .

Figure 8: *Test accuracy of Model C*. Only mean curves contain confidence ranges (SEM).

## I EXPERIMENTS WITH ALTERNATIVE WEIGHT INITIALIZATION FOR MODEL B

1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

We repeated the experiments for Model B on the MNIST and Fashion MNIST datasets due to suboptimal parameter initialization, which resulted in vanishing gradients at the start of the training. During the repeated experiments, we initialized the weights using the Glorot uniform method (Glorot & Bengio, 2010), which is specifically designed to initialize layers with nonlinear activations such as Tanh or Sigmoid. A gradient-magnitude gain of  $\frac{5}{3}$  for Tanh activations was used, as recommended by the PyTorch library. Biases were initialized to zero. The gradient magnitudes were examined to ensure they fell within a satisfactory range after initialization. Training length was reduced to 125 epochs for the gradient-based RMSProp and 50 epochs for two iterations of our method to test a 2.5x learning speedup.

As expected, similar magnitudes of learning rates performed well in the trainings using different weight initializations (compare Tables 3 and 6). The losses using our method are significantly lower

Table 6: Results for different learning rates.

Dataset	Method	The Most Important Results [Learning Rate: Avg. of Min. Training Loss]	
MNIST	<b>RMSProp</b> (125 epochs)	8 trainings per each learning rate: 1e-4 : 0.0945 ± 0.0062; 1.5e-4 : 0.0821 ± 0.0121; 1.75e-4 : 0.0750 ± 0.0102; 2e-4 : 0.0519 ± 0.0057; 2.25e-4 : 0.0862 ± 0.0171; <b>2.5e-4 : 0.0417 ± 0.0035</b> ; 2.75e-4 : 0.0515 ± 0.0094; 3e-4 : 0.0749 ± 0.0187; 3.5e-4 : 0.0536 ± 0.0095; 4e-4 : 0.0604 ± 0.0188;	
		50 trainings: <b>2.5e-4 : 0.0478 ± 0.0032</b> ;	
		3 trainings per each learning rate: 1.54e-4 : 0.128 ± 0.0068; 4.61e-4 : 0.0351 ± 0.0029;	
		8 trainings per each learning rate: 6e-4 : 0.0320 ± 0.0013; 8e-4 : 0.0256 ± 0.0019; <b>9e-4 : 0.0245 ± 0.0005</b> ; 1e-3 : 0.0255 ± 0.0021;	
		<b>2 Iterations</b> (50 epochs) 1.1e-3 : 0.0253 ± 0.0016; 1.2e-3 : 0.0246 ± 0.0005; 1.4e-3 : 0.0255 ± 0.0002;	
	Fashion MNIST	<b>RMSProp</b> (125 epochs)	8 trainings per each learning rate: 2.5e-4 : 0.329 ± 0.022; 3e-4 : 0.352 ± 0.024; <b>3.5e-4 : 0.315 ± 0.025</b> ; 4e-4 : 0.347 ± 0.024; 4.5e-4 : 0.374 ± 0.029; 5e-4 : 0.396 ± 0.025; <b>5.5e-4 : 0.310 ± 0.024</b> ; 6e-4 : 0.371 ± 0.025; 6.5e-4 : 0.368 ± 0.034; 7e-4 : 0.407 ± 0.021;
			50 trainings: <b>3.5e-4 : 0.344 ± 0.010</b> ;
			3 trainings per each learning rate: 1.33e-4 : 0.467 ± 0.002; 1.75e-4 : 0.463 ± 0.001;
			8 trainings per each learning rate: 6e-4 : 0.269 ± 0.003; 8e-4 : 0.255 ± 0.002; 9e-4 : 0.255 ± 0.002; 1e-3 : 0.254 ± 0.003; 1.1e-3 : 0.253 ± 0.003; <b>1.2e-3 : 0.245 ± 0.001</b> ; 1.4e-3 : 0.252 ± 0.003; 1.5e-3 : 0.261 ± 0.004;
			<b>2 Iterations</b> (50 epochs) 1.6e-3 : 0.253 ± 0.003; 1.8e-3 : 0.273 ± 0.008; 2e-3 : 0.262 ± 0.000;

after 2.5 times fewer epochs. Therefore, the two iterations of our method increase sample efficiency by more than 2.5 times. Good performance of the method across different parameter-initialization distributions is essential for its practical application as it contributes to robustness. Importantly, our algorithm maintains its performance gain compared to the gradient-based RMSProp in scenarios involving vanishing gradients, as demonstrated in the main experiments.

## J FUTURE WORK

Interesting directions for further experiments include: (a) Computing the average gradients over a much larger range than that of a parameter update to capture the global trend of the loss landscape. (b) More accurate approximation of the average Jacobians using Equation 5 instead of Equation 6. This would enable computing the average Jacobians of linear operators. Therefore, the algorithm based on the average gradient may enhance trainings of deep models without nonlinear activations. Moreover, the usage of Equation 5 may further improve the performance in the case of many nonlinear activations because of the increased precision in approximating the average gradient. (c) Incorporation of the momentum into our algorithm. Preferably Nesterov momentum (Dozat, 2016) should be used. If not, the average gradient would also be calculated for the momentum part of the update step. This could often reverse the direction of the momentum for a model parameter, thereby impairing the effectiveness of the entire momentum procedure. (d) Development of similar algorithms, but with update steps, that, for a given model parameter, vary in size over the iterations

1242 of the average-gradient computation. By adjusting the step size of each model parameter to the  
1243 absolute value of the average gradient, the learning process may be enhanced. (e) Tests of the  
1244 method on large and very deep architectures, that are used in practice and contain many nonlinear  
1245 layers. (f) More research on how the method scales up (Appendix E), also in relation to the number  
1246 of neurons in layers of neural networks. (g) Experiments with learning without forgetting (Li &  
1247 Hoiem, 2017) and online learning. Sample efficiency may be very beneficial there.

1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295