

SCoTTi: Save Computation at Training Time with an adaptive framework

Ziyu Li Enzo Tartaglione Van-Tam Nguyen

LTCI, Télécom Paris, Institut Polytechnique de Paris, France

{ziyu.li, enzo.tartaglione, van-tam.nguyen}@telecom-paris.fr

Abstract

On-device training is an emerging approach in machine learning where models are trained on edge devices, aiming to enhance privacy protection and real-time performance. However, edge devices typically possess restricted computational power and resources, making it challenging to perform computationally intensive model training tasks. Consequently, reducing resource consumption during training has become a pressing concern in this field. To this end, we propose SCoTTi (Save Computation at Training Time), an adaptive framework that addresses the aforementioned challenge. It leverages an optimizable threshold parameter to effectively reduce the number of neuron updates during training which corresponds to a decrease in memory and computation footprint. Our proposed approach demonstrates superior performance compared to the state-of-the-art methods regarding computational resource savings on various commonly employed benchmarks and popular architectures, including ResNets, MobileNet, and Swin-T.¹

1. Introduction

In recent years, there has been a growing interest in the development and deployment of machine-learning models due to their exceptional performance in various domains. However, the remarkable success of these models has come at a cost, as they require substantial computational resources for training and deployment, which makes training on resource-constrained edge devices very difficult [1]. To this end, two distinct research directions have emerged, each focusing on a different aspect of model optimization.

The first line of research centers around pruning algorithms, which aim to reduce the size of machine learning models without significantly compromising their performance [2]. Pruning algorithms selectively remove redundant or less influential components of a model, such as connections or parameters, while preserving its overall functionality [3, 4]. By doing so, these algorithms target

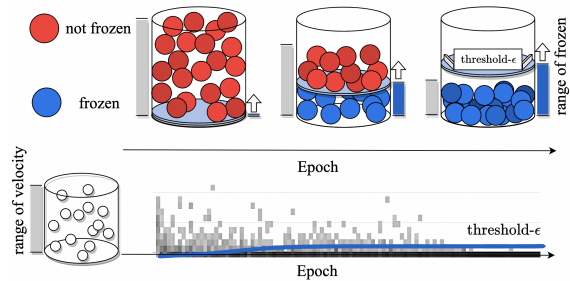


Figure 1. Our adaptive framework SCoTTi is based on an optimizable threshold ϵ . During the training, some neurons learn their own function (their change rate is measured through a velocity); all the neurons whose velocity is lower than some threshold ϵ will be stopped from updating.

the reduction of the memory footprint and computational requirements of the model without sacrificing its accuracy [5].

In parallel, researchers have also been exploring ways to save computational resources during the training process itself [6, 7, 8]. Training deep learning models typically involves numerous iterations over vast amounts of data, which can be computationally expensive and time-consuming, especially when moving to microcontrollers in the context of TinyML models [9, 10]. Therefore, finding methods to reduce the computational burden without compromising the training quality has become an important research objective [11], also considering that some steps are being moved to design hardware supporting dynamic pruning [12]. The lottery hypothesis [13, 14] highlights that only a small fraction of the tens of thousands of neurons in a neural network significantly impact the model’s performance. To save resources, we only need to update the portion of neurons that play a key role.

In our research, we updated only a subset of neurons during backpropagation and showed its feasibility through experiments. One approach is to optimize hyper-parameters like the learning rate through gradient descent [15], to ensure dynamic and faster convergence across the learning. In most instances, a trade-off between accuracy and saved Floating Point Operations (FLOPs) is inevitable, as we will

¹This article has been accepted for publication at the Resource Efficient Deep Learning for Computer Vision Workshop (ICCV 2023 workshop).

demonstrate in Tab. 1 (Sec. 4.3). To address this issue, we propose SCoTTi, a framework that allows dynamic selection of the neurons requiring an update at the current optimization step. More specifically, we leverage over the concept of *neuron velocity* (introduced in [16]) which is an estimator of learning convergence at the single neuron scale, to determine the sub-network to be updated. We propose a scheme where the velocity threshold ϵ , which determines whether neurons having a certain (low) velocity are in a frozen state or not, is learned at training time. Fig. 1 pictures the practical effect of SCoTTi: ϵ is naturally, through stochastic gradient descent, naturally increased as the average velocities (in the distribution at the bottom) in average drop to zero. SCoTTi has the big advantage of not requiring relevant hyper-parameters tuning (as both the learning rate and ϵ are learned), which ensures flexibility to the specific task and at the same time efficiency in terms of FLOPs computation at training time and to the best of our knowledge, is the first approach proposing a dynamic learned threshold to determine which part of the network to be learned.

The rest of the paper is organized as follows. In Sec. 2 we review the relevant literature concerning both pruning for deep neural networks and efficient hyper-parameters tuning. Next, in Sec. 3 we describe SCoTTi, first introducing the working mechanisms for the ultimate optimizer Sec. 3.1 and of neurons at equilibrium Sec. 3.2, and then describing the method proposing as well a rationale behind its working principle. Then, in Sec. 4 we experiment with our proposed training scheme over some deep ANNs on many different datasets, and finally, Sec. 5 concludes while providing further directions for future research.

2. Related Works

In this section, we will provide an overview of some of the most popular approaches to save computation when employing deep neural networks. To this end, we will divide the most-relevant literature into pruning-based approaches, distributed learning, and the most recent new directions to efficient parameters update.

2.1. Pruning methods

Pruning and quantization have emerged as prominent model compression techniques and have witnessed extensive development and application in recent years. Three mainstream pruning strategies - weight pruning [2], filter pruning [17], and neurons pruning [18] - have rapidly become hotspots of research. Weight pruning achieves model sparsity by removing small-magnitude weights in deep neural networks, thereby reducing the model’s parameter count and computational complexity [19, 20]. Filter pruning focuses on entire convolutional filters (channels), eliminating unimportant filters to further compress the model [21]. Neurons pruning simplifies the model by removing entire neu-

rons that make minor contributions to the model’s output. Additionally, model quantization further reduces model size and computational requirements by representing model parameters with lower precision, which is especially beneficial for edge devices with limited memory and computation capabilities [22, 23]. The application of these pruning and quantization methods significantly enhances model efficiency on edge devices, enabling the deployment of more complex models and achieving outstanding performance in real-time applications. Despite the big interest of the community, pruning, and quantization methods also exhibit notable limitations. Firstly, pruning requires multiple iterations to remove model parameters, which entails significant computation and time consumption [24]. Quantization, on the other hand, may result in some loss of model accuracy due to reduced precision [25]. Additionally, while these techniques effectively optimize resource consumption during model inference on edge devices, they offer limited assistance in training new models and require specific techniques to address this issue [26]. Consequently, retraining models on edge devices becomes a challenging task, calling for innovative solutions to overcome these restrictions.

2.2. Distributed learning

So far, most methods for training on edge devices are based on the idea of distributed learning [27]. These methods allow model training on edge devices without uploading raw data to a central server. On-device training methods based on distributed learning or federated learning offer significant advantages, including enhanced privacy protection as training occurs locally on edge devices without requiring data uploads to central servers [28, 29]. This approach reduces data transmission costs and alleviates the burden on central servers, contributing to lower inference latency. Moreover, the near real-time model updates make these methods suitable for time-sensitive applications.

Distributed learning or federated learning-based method has, however, certain limitations. The limited computational and storage resources of edge devices can lead to slower training efficiency and a prolonged convergence time [30]. Communication overhead remains a concern, as model updates among devices might result in substantial network traffic, particularly in scenarios with a large number of devices [31]. The heterogeneity of edge devices poses challenges, potentially affecting training stability and convergence [32]. Furthermore, the complexity of federated learning algorithms requires careful consideration of privacy preservation, model aggregation, and model drift, making algorithm design and optimization demanding tasks [33].

2.3. Efficient parameters update schemes

The lottery ticket hypothesis [13] posits that within a neural network, only a small subset of neurons plays a crucial role in achieving high accuracy, while the majority of neurons are less significant. Based on this hypothesis, if we focus on updating only this critical subset of neurons during the training process, we can attain comparable accuracy to updating all the neurons. The advantage of this approach is that it leads to more efficient training and improved generalization performance. By updating only the essential neurons, the training process focuses on the most informative parts of the model, mitigating the negative influence of noise and less relevant parameters. The challenge to identify such sub-networks efficiently is, however, still ongoing [14]. Some first steps have been moved in such a direction, achieving similar accuracy to a fully trained model while requiring fewer iterations and computational resources [16]. SCoTTi puts itself in such a frame, building on top of recent advances in the field [16, 15] and providing a simple yet effective approach that automatically tunes the updated sub-graph at training time, leading to computation savings at training time. In the next section, we will ground and provide more details on the working principles of SCoTTi.

3. SCoTTi

In this section, we will describe our proposed approach, which is pictured in Fig. 2. It builds on top of two approaches: *Gradient descent the ultimate optimizer* [34] (reviewed in Sec. 3.1) and *To update or not to update? Neurons at equilibrium in deep models* [16] (reviewed in Sec. 3.2). The innovation of SCoTTi is to make these algorithms work synergetically and optimize over a (formerly) static hyper-parameter (ϵ) that determines a threshold above which a neuron should be kept updated. This hyper-parameter is optimized along the training process and hence does not require fine optimization (as described in Sec. 3.3).

3.1. Learning the learning rate

The optimization of Deep Neural Networks (DNN) is a critical task that requires finding suitable hyperparameters (eg. learning rate) for the optimizer. The optimization of hyperparameters has been introduced in a previous work [34]; however, this method suffers from three main limitations. Firstly, the manual differentiation of optimizer update rules is a laborious and error-prone process that needs to be repeated for each variant of the optimizer. Secondly, the method only focuses on tuning the learning rate and neglects other important hyperparameters such as the momentum coefficient. Lastly, the method introduces an additional hyperparameter, the hyper-learning rate, which also

requires tuning. These limitations highlight the need for an improved approach to address these challenges in hyperparameters optimization.

The ultimate optimizer [15] is a recently proposed optimization scheme that operates based on the principle of dynamically updating the original fixed learning rate and/or other hyperparameters in real-time. It is designed to complement and enhance the functionality of general optimizers such as Adam [35], SGD [36], and others. The primary objective of the ultimate optimizer is to adaptively adjust the learning rate and/or other hyperparameters throughout the training process to find the most appropriate and effective one/ones for the current state of the model. By dynamically updating the hyperparameters like the learning rate, the ultimate optimizer aims to improve the convergence speed, optimization efficiency, and overall performance of the neural network model:

$$\alpha_{t+1} = \alpha_t - \eta \cdot \frac{\partial f(\omega_t)}{\partial \alpha_t} \quad (1)$$

$$\omega_{t+1} = \omega_t - \alpha_{t+1} \cdot \frac{\partial f(\omega_t)}{\partial \omega_t}. \quad (2)$$

As shown in (2), the ultimate optimizer modified the traditional approach of updating weight ω by treating the learning rate α as a learnable parameter, as outlined in (1), the approach of updating learning rate α which is based on the gradient descent algorithm, an additional hyperparameter η is required to control the step size for updating the learning rate α through gradient descent.

3.2. Neurons at Equilibrium

It is widely recognized in the literature that Deep Learning models commonly employed in state-of-the-art scenarios tend to be over-parametrized [37]. This observation has spurred two lines of research: reducing the size of these models through pruning algorithms, or finding ways to save computational resources during the training process. The former approach has received significant attention and has been extensively investigated. However, the latter approach had been at an impasse until a recent study [13] suggested its viability.

As shown in the lottery ticket hypothesis, only a small fraction of the total number of neurons in the model can be considered effective or influential in determining the model's performance (ie. The model's performance heavily relies on a subset of neurons that play a decisive role in capturing and representing the relevant patterns and information in the given data.). As inferred by the authors of Neurons at Equilibrium (NEq) [16], the existence of redundancy or inefficiency in the training process, where a large number of neurons may get updates despite not playing a decisive role in capturing the underlying patterns or

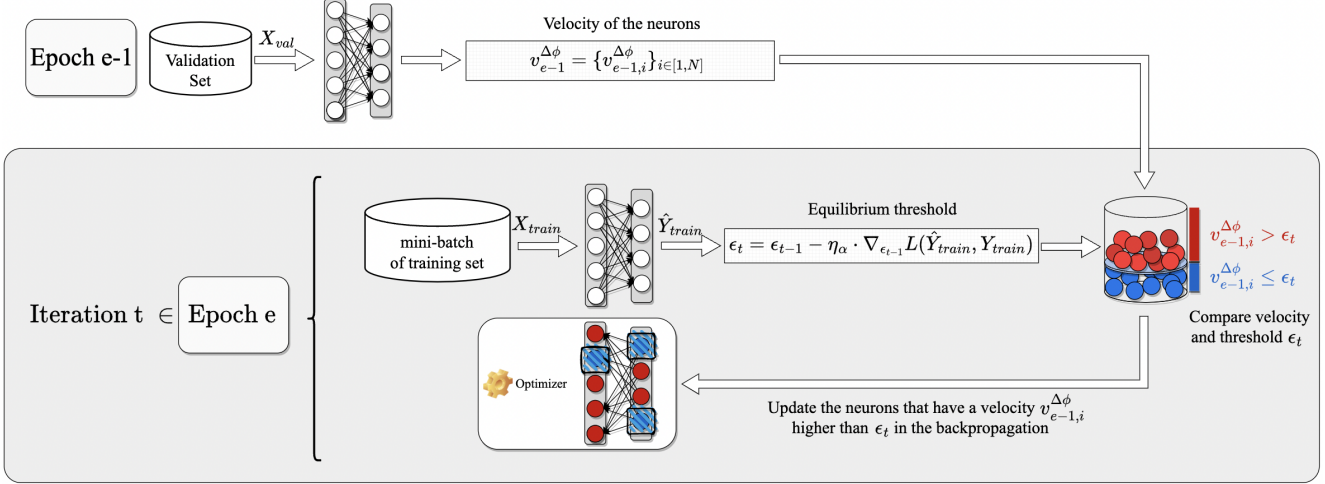


Figure 2. In the given figure, X_{train} and X_{val} correspond to the input of the training set Ξ_{train} and the validation set Ξ_{val} , respectively. Y_{train} corresponds to the labels of the training set. \hat{Y}_{train} represents the outputs obtained from the model when the corresponding inputs from the training set are fed into it. We update the threshold ϵ by the model's loss on the training set and compare it to the velocity of each neuron, suppressing updates to neurons with a velocity less than ϵ , where velocity is a value that indicates the similarity between the outputs of the same neuron in two consecutive epochs.

information in the data. By identifying and mitigating these unnecessary updates, it is possible to optimize the training process and improve the overall efficiency of the model. To this end, we want to identify the crucial neurons that require updating in each epoch based on a mask that is generated by employing a fixed threshold- ϵ .

The method for determining whether a neuron requires a further update or not follows. First, the cosine similarity between the outputs of the same neuron in two consecutive epochs is calculated below and also shown in Fig. 3:

$$\phi_e = \cos(\theta) = \hat{y}_{e-1} \cdot \hat{y}_e = \sum_{\xi \in \Xi_{val}} \sum_{n=1}^N \hat{y}_e^{n,\xi} \cdot \hat{y}_{e-1}^{n,\xi}. \quad (3)$$

We further use the obtained similarity to calculate the corresponding velocity as shown below:

$$\Delta\phi_e = \phi_e - \phi_{e-1} \quad (4)$$

$$v_e = \Delta\phi_e - \mu_{eq} \cdot v_{e-1}. \quad (5)$$

This is then compared to some threshold ϵ to determine if the neuron should be updated. The condition determining whether a neuron is in a frozen state is:

$$|v_e| \leq \epsilon. \quad (6)$$

When the inequality (6) is satisfied, it indicates that the neuron has reached an equilibrium state at epoch e .²

²The state of a neuron at epoch $e+1$ depends on two factors: its parameters and the working region (domain). When frozen, while the parameters do not change, the working region might (unless all the neurons of the previous layers are frozen). Hence, the state of a neuron might unfreeze.

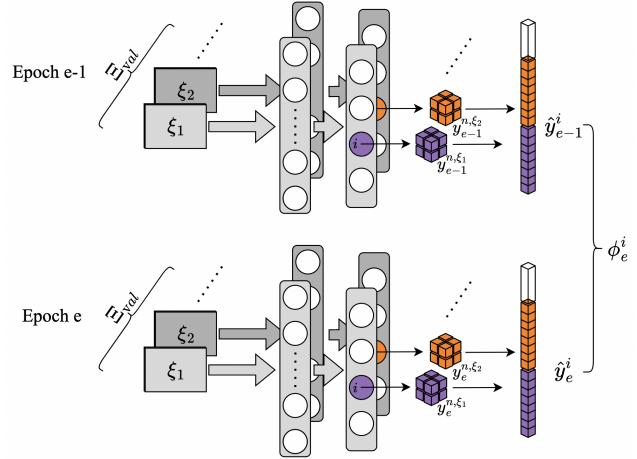


Figure 3. For a given epoch e , the model receives the ξ -th sample from the validation set Ξ_{val} . The output of the i -th neuron depends on both the model's parameters and the specific sample on the validation set. These outputs are squeezed, concatenated and the obtained vector is then normalized, obtaining \hat{y}_e^i . And we denote that ϕ_e^i is the cosine similarity of the feature vector \hat{y}_e^i and the feature vector \hat{y}_{e-1}^i which is obtained at epoch $e-1$.

Consequently, the neuron will be detached from the back-propagation graph, indicating that it should not be updated at epoch e . The effectiveness of this approach depends on two main factors: the optimality of the learning schedule, and the optimal, dynamic tuning of ϵ . While the first is addressed in Sec. 3.1, in the next section we will propose an approach to dynamically learn ϵ .

3.3. Learning the velocity threshold

In our study, we employ the NEq method in conjunction with the ultimate optimizer to update the weights. Additionally, we introduce the threshold parameter ϵ as an additional optimizable parameter to the update paradigm. After rearranging the terms of (6), the previous inequality can be expressed as follows:

$$|v_e| - \epsilon \leq 0. \quad (7)$$

We define the i -th neuron to be in the equilibrium state at epoch e if (7) is satisfied. Our approach is currently based on the SGD optimizer (but extendible to any other optimizer), therefore to incorporate ϵ as a learnable parameter that can be optimized by the hyper optimizer as shown in Fig. 2, we integrate the inequality into the weight update paradigm, as illustrated below. We utilize step function $\Theta(\cdot)$ to serve as the activation function for this component:

$$\omega_t = \omega_{t-1} - \alpha \cdot \frac{\partial \mathcal{L}}{\partial \omega_t} \cdot \Theta(|v_e| - \epsilon). \quad (8)$$

As shown in (8), α represents the learned learning rate, ϵ presents the learned threshold, ω_t represents the weight update at epoch t , $\frac{\partial \mathcal{L}}{\partial \omega_t}$ refers to the gradient of the loss function, and $\Theta(\cdot)$ is the step function that ensures the inequality condition is met. The hyperparameters α and ϵ are iteratively updated by another hyper optimizer, following the gradient of the loss function for each parameter, denoted as $\frac{\partial \mathcal{L}}{\partial \alpha}$ and $\frac{\partial \mathcal{L}}{\partial \epsilon}$, respectively. This update scheme ensures that the parameters are adjusted in a direction that minimizes the loss function, thereby leading to a reduction in overall loss. By incorporating the $\Theta(\cdot)$ term, compared to the original update paradigm, our approach presents two major advantages.

We introduce ϵ as a learnable parameter The inclusion of $\Theta(\cdot)$ enables us to incorporate ϵ into the update paradigm as an optimizable parameter. Specifically, ϵ is optimized along the direction of loss reduction, ensuring that updates are made in a manner that progressively minimizes the overall loss. This allows us to dynamically adjust the threshold for determining neuron equilibrium during the optimization process.

SCoTTi is consistent with “freezing neurons at equilibrium” The use of the $\Theta(\cdot)$ function ensures that if a neuron is in equilibrium, $\Theta(\cdot)$ will filter these ones. Conversely, if the neuron is not in equilibrium, $\Theta(\cdot)$ will allow updates. This aligns with the concept of NEq, where only non-equilibrium neurons undergo updates while equilibrium neurons remain unchanged.

In the following, we will show how to update α and ϵ by a parameter η via gradient descent per iteration and derive an iterative update paradigm, being for ϵ .

$$\begin{aligned} \epsilon_{t+1} &= \epsilon_t - \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial \epsilon_t} = \epsilon_t - \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial w_t} \cdot \frac{\partial w_t}{\partial \alpha_t} \\ &= \epsilon_t - \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial w_t} \\ &\quad \cdot \frac{\partial \left[w_{t-1} - \alpha_t \cdot \frac{\partial \mathcal{L}(w_{t-1})}{\partial w_{t-1}} \cdot \Theta(|v_e| - \epsilon_t) \right]}{\partial \epsilon_t} \\ &= \epsilon_t - \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial w_t} \cdot \frac{\partial \mathcal{L}(w_{t-1})}{\partial w_{t-1}} \cdot \alpha_t, \end{aligned} \quad (9)$$

while for α

$$\begin{aligned} \alpha_{t+1} &= \alpha_t - \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial \alpha_t} = \alpha_t - \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial w_t} \cdot \frac{\partial w_t}{\partial \alpha_t} \\ &= \alpha_t - \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial w_t} \\ &\quad \cdot \frac{\partial \left[w_{t-1} - \alpha_t \cdot \frac{\partial \mathcal{L}(w_{t-1})}{\partial w_{t-1}} \cdot \Theta(|v_e| - \epsilon_t) \right]}{\partial \alpha_t} \\ &= \alpha_t + \eta \cdot \frac{\partial \mathcal{L}(w_t)}{\partial w_t} \cdot \frac{\partial \mathcal{L}(w_{t-1})}{\partial w_{t-1}} \cdot \Theta(|v_e| - \epsilon_t). \end{aligned} \quad (10)$$

Indeed, the iterative update paradigm enables us to optimize both α and ϵ simultaneously. The iterative process continues until the values of α and ϵ converge to their optimal values, indicating that the model has reached an optimal state.

3.4. Overview on the method

As presented in Algorithm 1, we incorporate a mask in the ultimate optimizer to constrain the weight updates for a part of neurons (in the form of a dictionary, \mathcal{E}). Additionally, we introduce an additional parameter ϵ , into the standard iterative update paradigm for the weights. This enables us to learn the threshold ϵ for neuron equilibrium and control the updates accordingly in real time. Focusing on the hyper-optimization step, we extend the hyper-optimizer for the learning rate (line 1) to include updates for ϵ (line 2). By optimizing both of them simultaneously, we can fine-tune the threshold for determining neuron equilibrium and further enhance the model’s performance. When performing the update in line 2, to maintain the gradient flow without disruption, we employ a Straight-Through Estimator (STE) [38]. STE allows for backpropagation through the weight update operation without modifying the gradient. This technique enables us to incorporate non-differentiable

Algorithm 1 Learned threshold ϵ approach

Require:

- \mathcal{L}_{t-1} : loss calculated on the current minibatch,
- \mathcal{E} : a dictionary used to store the neurons that are at equilibrium state,
- v_e : a dictionary used to record the velocity between epochs $e - 1$ and e ,
- ω_{t-1} : parameters of the model at iteration $t - 1$,
- α_{t-1} : learnable learning rate, at iteration $t - 1$,
- ϵ_{t-1} : learnable threshold, at iteration $t - 1$,
- η_α : hyper-learning rate for α ,
- η_ϵ : hyper-learning rate for ϵ .

- 1: $\alpha_t \leftarrow \alpha_{t-1} - \eta_\alpha \cdot \nabla_{\alpha_{t-1}} \mathcal{L}_{t-1}$
 - 2: $\epsilon_t \leftarrow \epsilon_{t-1} - \eta_\epsilon \cdot \nabla_{\epsilon_{t-1}} \mathcal{L}_{t-1}$
 - 3: **for** all the i -th neurons in the model **do**
 - 4: **if** i is in \mathcal{E} **then**
 - 5: $G_{t-1}^i = 0$
 - 6: **else**
 - 7: $G_{t-1}^i = \nabla_{\omega_{t-1}^i} \mathcal{L}_{t-1}$
 - 8: **end if**
 - 9: $\omega_t^i \leftarrow \omega_{t-1}^i - \alpha_t \cdot G_{t-1}^i \cdot \Theta(|v_e^i| - \epsilon_t)$
 - 10: **end for**
 - 11: **return** $\alpha_t, \epsilon_t, \omega_t$
-

operations, such as binarization which we used here (neurons in a frozen state), into the training process while preserving the gradient information. We can successfully skip the gradient computation for the neurons marked in a frozen state (line 4) and successively we perform a standard update for the model’s parameters (line 9).

4. Experiment

In this section, we present our experiments describing datasets, architectures, and learning policies. We first propose an ablation study (Sec. 4.2) and then we describe the main experimental results (Sec. 4.3). We have performed our experiments on an NVIDIA Tesla V100 equipped with 32GB and developed the code using PyTorch 1.13.1.³

4.1. Datasets

CIFAR-10 [39]. The CIFAR-10 dataset is a widely recognized benchmark in image classification. It consists of 60,000 color images divided into 10 classes representing specific object categories, such as airplanes, automobiles, birds, and cats. With 6,000 images per class, the dataset is evenly balanced. The images have a resolution of 32x32

pixels, making them computationally efficient and suitable for scenarios with limited computational resources.

CIFAR-100 [39]. The CIFAR-100 dataset builds upon the CIFAR-10 dataset by introducing 100 fine-grained object categories. It contains 60,000 color images, with 600 images per class. CIFAR-100 offers a higher level of label granularity compared to CIFAR-10, enabling more challenging classification tasks. The images in this dataset also have a resolution of 32x32 pixels.

Tiny ImageNet [43]. The Tiny ImageNet dataset is more extensive than CIFAR-10 and CIFAR-100, consisting of 200 diverse object categories. It includes a total of 100,000 color images, with 500 images per class. The images in Tiny ImageNet have a resolution of 64x64 pixels, providing a higher level of detail compared to the previous 2 datasets.

Clipart [45]. The Clipart dataset contains 73,810 images representing various everyday objects and scenes in a clipart-style visual format. It consists of 345 object categories, making it a rich resource for studying domain adaptation to clipart-like graphics.

Painting [45]. The Painting dataset consists of 76,174 high-quality images inspired by diverse artistic painting styles. It includes 345 object categories, providing researchers with ample opportunities to study domain adaptation in the context of artistic representations.

4.2. Ablation study

We performed our ablation study training a ResNet-56 [42] on the CIFAR-100 dataset. The baseline models were trained with Stochastic Gradient Descent (SGD) as the optimizer, initialized with a learning rate α of 0.1, momentum μ of 0.9, and weight decay of 5×10^{-4} for a total of 250 epochs. The learning rate α was reduced by a factor of 10 after 100 and 150 epochs. For the experiments based on the Ultimate Optimizer (including SCoTTi), we uniformly set the hyperlearning rate for learning rate to 1.5×10^{-5} . For all other parameters unrelated to the learning rate, we remain consistent with the experiments that use SGD as the optimizer.

In general, a trade-off between FLOPs and accuracy can be observed in Fig. 4: for less than the model’s accuracy experiences a significant drop when the average FLOPs value falls below a certain value. In these MFLOPs, we observe a degradation in the performance of SCoTTi. This can be observed by testing several representative ϵ for each of the two frameworks using a fixed threshold ϵ (the blue points approach and green points approach are) to determine the relationship between FLOPs of backpropagation and accuracy. For our proposed framework with optimizable threshold ϵ (in red), we initialize ϵ to zero, and since the ϵ updating process is based on the gradient descent method, a hyperparameter η_ϵ is needed to control the step size when updating

³The code is publicly available at <https://github.com/liziyu403/SCoTTi-Save-Computation-at-Training-Time-with-an-adaptive-framework>.

Table 1. Table of experimental results. **FLOPs saved** refers to the total FLOPs saved at the end of model training as a percentage of the total FLOPs when the model is updated normally (no frozen); **Top-1** refers to the model’s Test Top1-Accuracy. Models marked with “*” indicate that they utilize a pre-trained model on ImageNet1k. The rows in gray are our proposed SCoTTi.

Dataset	Architecture	Optimization Approach			FLOPs saved	Top-1
		NEq [16]	Ultimate optim. [15]	Learned ϵ		
CIFAR-10 [39]	VGG-16 [40]		✓		00.00%	88.54%
					37.41%	89.86%
				✓	00.00%	92.76%
			✓	✓	30.98%	92.70%
			✓	✓	43.66%	92.58%
CIFAR-100 [39]	Swin-T* [41]		✓		00.00%	91.59%
					39.66%	90.96%
				✓	00.00%	91.65%
			✓	✓	48.84%	91.74%
			✓	✓	58.76%	91.77%
CIFAR-100 [39]	ResNet-32 [42]		✓		00.00%	68.42%
					38.80%	69.97%
				✓	00.00%	70.06%
			✓	✓	59.89%	69.24%
			✓	✓	60.59%	70.43%
CIFAR-100 [39]	ResNet-56 [42]		✓		00.00%	69.69%
					41.12%	71.40%
				✓	00.00%	70.15%
			✓	✓	56.58%	70.36%
			✓	✓	58.97%	71.55%
Clipart [43]	ResNet-18 [42]		✓		00.00%	73.21%
					38.06%	72.19%
				✓	00.00%	73.01%
			✓	✓	49.33%	72.60%
			✓	✓	53.86%	73.21%
Painting [43]	ResNet-18 [42]		✓		00.00%	64.51%
					27.94%	62.14%
				✓	00.00%	60.82%
			✓	✓	77.34%	63.46%
			✓	✓	76.92%	65.44%
Tiny ImageNet [43]	MobileNet-v2* [44]		✓		00.00%	55.69%
					53.28%	56.40%
				✓	00.00%	60.02%
			✓	✓	80.83%	60.53%
			✓	✓	86.44%	60.68%

the ϵ , and a few representative η_ϵ are selected for recording in our experiments, the learned curves for ϵ are shown as Fig. 5.

As Fig. 4 shows, in these two plots, the top plot represents the results obtained on the TRAIN set, and the bottom plot represents the results on the TEST set. The x-axis represents the average FLOPs required per epoch, while the y-axis represents the accuracy achieved. Our objective is to minimize the FLOPs while maximizing the accuracy.

Therefore, points closer to the upper left corner indicate better performance, as they represent lower FLOPs and higher accuracy simultaneously.

Our proposed optimizable ϵ approach demonstrates more pronounced and consistent performance compared to both the baseline and the NEq method based on the ultimate optimizer. We attribute this improvement to the initial small value of ϵ during training, which enables the network to effectively learn the features. In general, as the number of

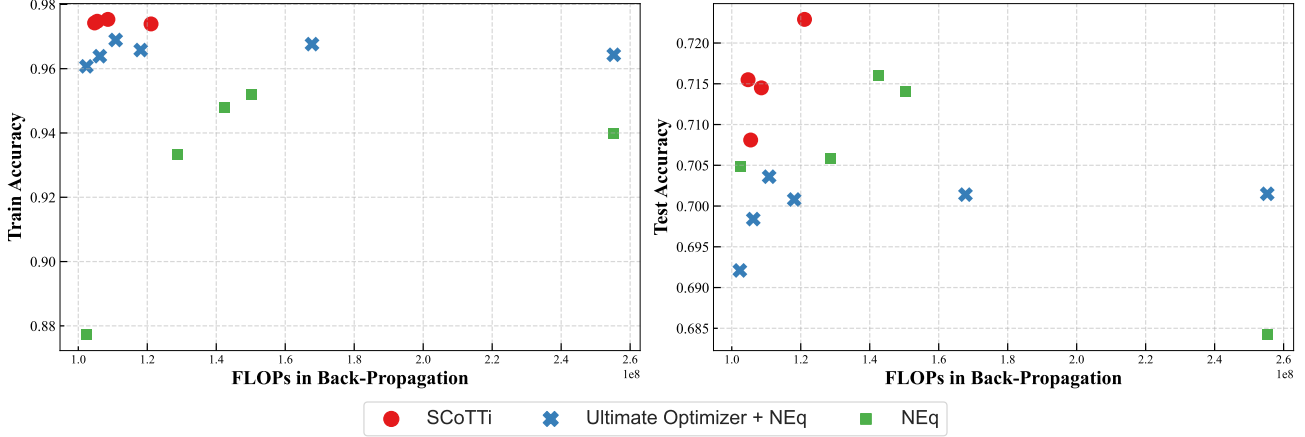


Figure 4. The x-axis represents the average of FLOPs required per epoch, while the y-axis represents the test accuracy achieved. Points closer to the upper left corner indicate better performance, as they represent lower FLOPs and higher accuracy simultaneously.

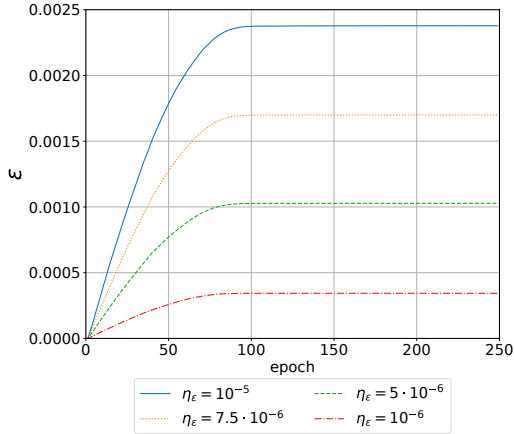


Figure 5. The hyperparameter ϵ is updated using gradient descent: $\epsilon_t = \epsilon_{t-1} - \eta_\epsilon \cdot \nabla_{\epsilon_{t-1}} \mathcal{L}_{t-1}$, with different curves corresponding to different values of η_ϵ .

training iterations increases, overfitting of certain features often occurs. However, the gradual increment of ϵ in our method helps mitigate the overfitting issue to some extent. This contributes to the enhanced and stable performance observed in our proposed approach.

4.3. Experimental results

We have chosen a selection of the most representative models currently available: VGG16, ResNet-18/32/56, MobileNet-v2, and Swin-T. Among them, MobileNet-v2 and Swin-T use a pre-trained model, on ImageNet1K.

In all the experiments based on the SGD scheduler, we set the initial learning rate α to 0.1 (initialized to 0.01 for MobileNet-v2 and Swin-T) and reduced it by a factor of 10 after 100 and 150 epochs (after 30 and 60 for MobileNet-

v2 and Swin-T). The momentum μ was set to 0.9, weight decay to 5×10^{-4} , and the threshold ϵ to 0.001, with a total of 250 epochs (90 epochs for MobileNet-v2 and Swin-T). The parameter settings for the ultimate optimizer-based experiments were largely the same as those for the SGD-based experiments, with the key difference being the use of a hyper-optimizer to optimize the initial learning rate instead of the scheduler. η_α played a crucial role in the ultimate optimizer as it was used to optimize the learning rate. After conducting tests, we identified a suitable value of η_α for each architecture, ensuring optimal performance. For SCoTTi, we primarily adopted the hyperparameters of the ultimate optimizer, except that ϵ was initialized to 0, and η_ϵ was set to half of η_α .

The experimental results are presented in Tab. 1. Our proposed approach is versatile, allowing for substantial FLOPs reduction across various datasets and architectures, while also slightly improving accuracy.

5. Conclusion

This study has presented the SCoTTi, an adaptive training framework for on-device training. The primary objective of SCoTTi was to reduce the number of FLOPs without compromising the model’s accuracy. Remarkably, the results showed that in the majority of experimental cases, SCoTTi not only achieved state-of-the-art performance but also demonstrated the potential to enhance model accuracy in some of the most popular architectures on downstream tasks. Future works for the SCoTTi framework include investigating its adaptability to diverse domains such as natural language processing and audio processing, tailoring it to various hardware platforms for optimal on-device training, and conducting large-scale deployment experiments to assess its viability for real-world applications.

Acknowledgements

The research leading to these results has received funding from the project titled “PC2-NF-NAI” in the frame of the program “PEPR 5G et Réseaux du futur”, and by Hi!PARIS Center on Data Analytics and Artificial Intelligence.

References

- [1] Neil C. Thompson, Kristjan H. Greenewald, Keeheon Lee, and Gabriel F. Manso. The computational limits of deep learning. *CoRR*, abs/2007.05558, 2020. [1](#)
- [2] Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In D. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 2. Morgan-Kaufmann, 1989. [1](#), [2](#)
- [3] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. [1](#)
- [4] Yawei Li, Kamil Adamczewski, Wen Li, Shuhang Gu, Radu Timofte, and Luc Van Gool. Revisiting random channel pruning for neural network compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 191–201, June 2022. [1](#)
- [5] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training, 2018. [1](#)
- [6] Xun Jiao, Vahideh Akhlaghi, Yu Jiang, and Rajesh K Gupta. Energy-efficient neural networks using approximate computation reuse. In *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 1223–1228. IEEE, 2018. [1](#)
- [7] Jason Servais and Ehsan Atoofian. Adaptive computation reuse for energy-efficient training of deep neural networks. *ACM Transactions on Embedded Computing Systems (TECS)*, 20(6):1–24, 2021. [1](#)
- [8] Yuedong Yang, Guihong Li, and Radu Marculescu. Efficient on-device training via gradient filtering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3811–3820, 2023. [1](#)
- [9] Danilo Pietro Pau and Fabrizio Maria Aymone. Suitability of forward-forward and pepita learning to mlcommons-tiny benchmarks. In *2023 IEEE International Conference on Omni-layer Intelligent Systems (COINS)*, pages 1–6. IEEE, 2023. [1](#)
- [10] Yinghao Wang, Rémi Nahon, Enzo Tartaglione, Pavlo Mozharovskiy, and Van-Tam Nguyen. Optimized preprocessing and tiny ml for attention state classification. In *2023 IEEE Statistical Signal Processing Workshop (SSP)*, pages 695–699, 2023. [1](#)
- [11] Li Huang, Yifeng Yin, Zeng Fu, Shifa Zhang, Hao Deng, and Dianbo Liu. Loadaboost: Loss-based adaboost federated machine learning with reduced computational complexity on iid and non-iid intensive care data. *Plos one*, 15(4):e0230706, 2020. [1](#)
- [12] Yang Wang, Yubin Qin, Dazheng Deng, Jingchuan Wei, Tianbao Chen, Xinhan Lin, Leibo Liu, Shaojun Wei, and Shouyi Yin. Trainer: An energy-efficient edge-device training processor supporting dynamic weight pruning. *IEEE Journal of Solid-State Circuits*, 57(10):3164–3178, 2022. [1](#)
- [13] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019. [1](#), [3](#)
- [14] Enzo Tartaglione. The rise of the lottery heroes: why zero-shot pruning is hard. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 2361–2365. IEEE, 2022. [1](#), [3](#)
- [15] Kartik Chandra, Audrey Xie, Jonathan Ragan-Kelley, and Erik Meijer. Gradient descent: The ultimate optimizer. *Advances in Neural Information Processing Systems*, 35:8214–8225, 2022. [1](#), [3](#), [7](#)
- [16] Andrea Bragagnolo, Enzo Tartaglione, and Marco Grangetto. To update or not to update? neurons at equilibrium in deep models. *Advances in Neural Information Processing Systems*, 35:22149–22160, 2022. [2](#), [3](#), [7](#)
- [17] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *International Conference on Learning Representations*, 2017. [2](#)
- [18] Abdolghani Ebrahimi and Diego Klabjan. Neuron-based pruning of deep neural networks with better generalization using kronecker factored curvature approximation. In *2023 International Joint Conference on Neural Networks (IJCNN)*, pages 1–9. IEEE, 2023. [2](#)
- [19] George Retsinas, Athena Elafrou, Georgios I. Goumas, and Petros Maragos. Weight pruning via adaptive sparsity loss. *CoRR*, abs/2006.02768, 2020. [2](#)
- [20] Tianyun Zhang, Shaokai Ye, Kaiqi Zhang, Jian Tang, Wujie Wen, Makan Fardad, and Yanzhi Wang. A systematic dnn weight pruning framework using alternating direction method of multipliers. In *Proceedings of the European conference on computer vision (ECCV)*, pages 184–199, 2018. [2](#)
- [21] Lucas Liebenwein, Cenk Baykal, Harry Lang, Dan Feldman, and Daniela Rus. Provable filter pruning for efficient neural networks. In *International Conference on Learning Representations*, 2020. [2](#)
- [22] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tijmen Blankevoort. A white paper on neural network quantization. *CoRR*, abs/2106.08295, 2021. [2](#)

- [23] Jiwei Yang, Xu Shen, Jun Xing, Xinmei Tian, Houqiang Li, Bing Deng, Jianqiang Huang, and Xian-sheng Hua. Quantization networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7308–7316, 2019. 2
- [24] Enzo Tartaglione, Andrea Bragagnolo, and Marco Grangetto. Pruning artificial neural networks: A way to find well-generalizing, high-entropy sharp minima. In *Artificial Neural Networks and Machine Learning–ICANN 2020: 29th International Conference on Artificial Neural Networks, Bratislava, Slovakia, September 15–18, 2020, Proceedings, Part II 29*, pages 67–78. Springer, 2020. 2
- [25] Aojun Zhou, Anbang Yao, Kuan Wang, and Yurong Chen. Explicit loss-error-aware quantization for low-bit deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9426–9435, 2018. 2
- [26] Bingyan Liu, Yifeng Cai, Yao Guo, and Xiangqun Chen. Transtailor: Pruning the pre-trained model for improved transfer learning. In *Proceedings of the AAAI conference on artificial intelligence*, pages 8627–8634, 2021. 2
- [27] Joost Verbraeken, Matthijs Wolting, Jonathan Katzy, Jeroen Kloppenburg, Tim Verbelen, and Jan S Rellermeyer. A survey on distributed machine learning. *Acm computing surveys (csur)*, 53(2):1–33, 2020. 2
- [28] Yang Lu, Zhengxin Yu, and Neeraj Suri. Privacy-preserving decentralized federated learning over time-varying communication graph. *ACM Transactions on Privacy and Security*, 26(3):1–39, 2023. 2
- [29] Maria Florina Balcan, Avrim Blum, Shai Fine, and Yishay Mansour. Distributed learning, communication complexity and privacy. In *Conference on Learning Theory*, pages 26–1. JMLR Workshop and Conference Proceedings, 2012. 2
- [30] Emre Ozfatura, Sennur Ulukus, and Deniz Gündüz. Straggler-aware distributed learning: Communication–computation latency trade-off. *Entropy*, 22(5), 2020. 2
- [31] Luping WANG, Wei Wang, and Bo LI. Cmf1: Mitigating communication overhead for federated learning. pages 954–964, 07 2019. 2
- [32] Alexander Brecko, Erik Kajati, Jiri Koziorek, and Iveta Zolotova. Federated learning for edge computing: A survey. *Applied Sciences*, 12(18), 2022. 2
- [33] Tien-Dung Cao, Tram Truong-Huu, Hien Tran, and Khanh Tran. A federated deep learning framework for privacy preservation and communication efficiency. *Journal of Systems Architecture*, 124:102413, 2022. 2
- [34] Atilim Gunes Baydin, Robert Cornish, David Martinez Rubio, Mark Schmidt, and Frank Wood. Online learning rate adaptation with hypergradient descent. In *International Conference on Learning Representations*, 2018. 3
- [35] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 3
- [36] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *International Conference on Computational Statistics*, 2010. 3
- [37] Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. SGD learns over-parameterized networks that provably generalize on linearly separable data. In *International Conference on Learning Representations*, 2018. 3
- [38] Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013. 5
- [39] A Krizhevsky. Learning multiple layers of features from tiny images. *Master’s thesis, University of Tront*, 2009. 6, 7
- [40] Shuying Liu and Weihong Deng. Very deep convolutional neural network based image classification using small training sample size. In *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, pages 730–734, 2015. 7
- [41] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 7
- [42] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 6, 7
- [43] 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*, pages 248–255, 2009. 6, 7
- [44] Andrew Howard, Andrey Zhmoginov, Liang-Chieh Chen, Mark Sandler, and Menglong Zhu. Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation. In *CVPR*, 2018. 7
- [45] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1406–1415, 2019. 6