# Parabolic Approximation Line Search for DNNs

**Maximus Mutschler and Andreas Zell**
University of Tübingen
Sand 1, D-72076 Tübingen, Germany
{maximus.mutschler, andreas.zell}@uni-tuebingen.de

## Abstract

A major challenge in current optimization research for deep learning is automatically finding optimal step sizes for each update step. The optimal step size is closely related to the shape of the loss in the update step direction. However, this shape has not yet been examined in detail. This work shows empirically that the mini-batch loss along lines in negative gradient direction is locally mostly convex and well suited for one-dimensional parabolic approximations. We introduce a simple and robust line search approach by exploiting this parabolic observation, which performs loss-shape-dependent update steps. Our approach combines well-known methods such as parabolic approximation, line search, and conjugate gradient to perform efficiently. It surpasses other step size estimating methods and competes with standard optimization methods on a large variety of experiments without the need for hand-designed step size schedules. Thus, it is of interest for objectives where step-size schedules are unknown or do not perform well. Our extensive evaluation includes multiple comprehensive hyperparameter grid searches on several datasets and architectures. Finally, we provide a general investigation of exact line searches in the context of batch losses and exact losses, including their relation to our line search approach.

## 1 Introduction

Automatic determination of optimal step sizes for each update step of stochastic gradient descent is a major challenge in current optimization research for deep learning [3, 5, 12, 29, 38, 45, 48, 52, 60]. One default approach to tackle this challenge is to apply line search methods. Several of these have been introduced for Deep Learning [12, 29, 38, 45, 60]. However, these approaches have not analyzed the shape of the loss functions in update step direction in detail, which is important since the optimal step size stands in strong relation to this shape. To shed light on this, our work empirically analyses the shape of the loss function along update step direction for deep learning scenarios often considered in optimization. We further elaborate on the properties found to define a simple and empirically justified optimizer.

Our contributions are as follows: **1:** Our empirical analysis suggests that the mini-batch loss in the negative gradient direction mostly shows locally convex shapes. Furthermore, we show that parabolic approximations are well suited to estimate the minima in these directions (Section 3). **2:** Exploiting this parabolic observation, we build a simple line search optimizer that constructs its loss function-dependent learning rate schedule. The performance of our optimization method is extensively analyzed, including a comprehensive comparison to other optimization methods (Sections 4,5). **3:** We provide a convergence analysis that backs our empirical results under strong assumptions (Section 4.4). **4:** We provide a general investigation of exact line searches on batch losses and their relation to line searches on the exact loss as well as their relation to our line search approach (Section 6) and, finally, analyze the relation of our approach to interpolation (Section 7).

The full-batch loss $\mathcal{L}$ is defined as the average over realizations of a mini-batch-wise loss function $\mathcal{L}_{\mathbb{B}}$: $\mathcal{L} : \mathbb{R}^m \to \mathbb{R}, \boldsymbol{\theta} \mapsto n^{-1} \sum_{i=1}^n \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta})$ with $\mathbb{B}_i$ being a batch of input samples, and $n$ the amount of batches in the dataset. $\boldsymbol{\theta} \in \mathbb{R}^m$ denotes the parameters to be optimized. Note, that we consider a sample as one mini-batch of multiple inputs. In this work, we consider $\mathcal{L}_{\mathbb{B}_t}$ at optimization step $t$ in negative gradient direction:

$$l_t(s) : \mathbb{R} \to \mathbb{R}, s \mapsto \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t + s \frac{-\boldsymbol{g}_t}{||\boldsymbol{g}_t||}), \tag{1}$$

where $\boldsymbol{g}_t$ is $\nabla_{\boldsymbol{\theta}_t} \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)$. We denote $l_t$ as *mini-batch loss along a line* and $s$ is the step on this line. The motivation of our work builds upon the following assumption:

**Assumption 1.** *(Informal) The position $\boldsymbol{\theta}_{min} = \boldsymbol{\theta}_t + s_{min} \frac{-\boldsymbol{g}_t}{||\boldsymbol{g}_t||}$ of a minimum of $l_t$ is a well enough estimator for the position of a minimum of the full-batch loss $\mathcal{L}$ on the same line to perform a successful optimization process.*

We empirically analyze Assumption 1 further in section 6.

## 2 Related work

Our optimization approach is based on well-known methods, such as line search, the non-linear conjugate gradient method, and quadratic approximation, which can be found in Numerical Optimization [28]. In addition, [28, §3.5] describes a similar line search routine for the deterministic setting. The concept of parabolic approximations is also exploited by the well-known line search of More and Thunte [40]. Our work contrasts common optimization approaches in deep learning by directly exploiting the local parabolic observation (see Section 3) of $l_t$. Similarly, *SGD-HD* [3] performs update steps towards a minimum of $l_t$, by performing gradient descent on the learning rate. Concurrently, [10] explored a similar direction as this work by analyzing possible line search approximations for DNN loss landscapes but does not exploit these for optimization.

The recently published *Stochastic Line-Search* (*SLS*) [60] is an optimized backtracking line search based on the Armijo condition, which samples, like our approach, additional mini-batch losses from the same batch and checks the Armijo condition on these. [60] assumes that the model interpolates the data. The latter states formally that the gradient with respect to each batch converges to zero at the optimum of the convex function. This implies that if $\mathcal{L}(\boldsymbol{\theta})$ is minimized at $\boldsymbol{\theta}^*$ and thus $\nabla_{\boldsymbol{\theta}^*} \mathcal{L}(\boldsymbol{\theta}^*) = 0$ holds, then for all $\mathcal{L}_{\mathbb{B}_i}$ we have that $\nabla_{\boldsymbol{\theta}^*} \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta}^*) = 0$. *SLS* exhibits competitive performance against multiple optimizers on several DNN tasks. [45] introduces a related idea but does not provide empirical results for DNNs. [12] also uses a backtracking Armijo line search, but with the aim to regulate the optimal batch size.

The methodically appealing but complex *Probabilistic Line Search* (*PLS*) [38] approximates $\mathcal{L}$ along lines with a Gaussian Process posterior based on realizations of $\mathcal{L}_{\mathbb{B}_i}$ for different batches $\mathbb{B}_i$. The mean of this Posterior is a cubic spline. A suitable update step size is found by exploiting a probabilistic formulation of the Wolf conditions.

*Gradient Only Line Search* (*GOLS1*) [29] GOLS1 searches for a minimum on lines by looking for a sign change of the first directional mini-batch derivative in search direction. For each measurement, a different $\mathbb{B}_i$ is chosen. Both approaches show that they can optimize successfully on several machine learning problems and can compete against plain *SGD*.

From the perspective of assumptions about the shape of the loss landscape, second order methods such as *oLBFGS* [55], *KFRA* [7], *L-SR1* [47], *QUICKPROP* [15], *S-LSR1* [4], and *KFAC* [39] generally assume that the loss function can be approximated locally by a parabola of the same dimension as the loss function. Adaptive methods such as *SGD* with momentum [51], *ADAM* [30], *ADAGRAD* [14], *ADABOUND* [37], *AMSGRAD* [49] or *RMSProp* [59] focus more on the handling of noise than on shape assumptions. In addition, methods exist that approximate the loss function in specific directions: The *L4* adaptation scheme [52] as well as *ALIG* [5] estimate step sizes by approximating the loss function linearly in negative gradient direction, whereas our approach approximates the loss function parabolically in negative gradient direction. Finally, *COCOB* [44] has to be mentioned, an alternative learning rate free approach, which automatically estimates step directions and sizes with a reward-based coin betting concept.
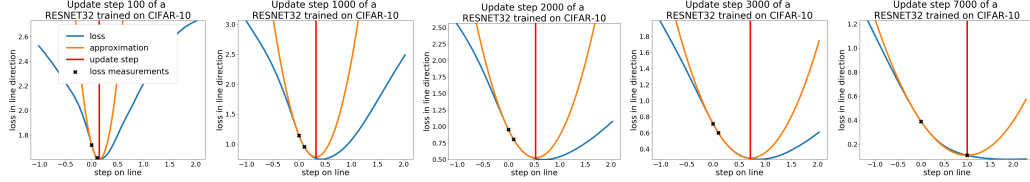
Figure 1: Representative mini-batch losses along lines in negative normalized gradient direction (blue), parabolic approximations (orange) and the position of the approximated minima (red). Further plots are provided in Appendix A.

## 3 Empirical analysis of the shape of mini-batch losses along lines

In this section, we analyze mini-batch losses along lines (see Eq. 1) during the training of multiple architectures and show that they locally exhibit mostly convex shapes, which are well suited for parabolic approximations. We focus on CIFAR-10, as it is extensively analyzed in optimization research for deep learning. However, we observed similar results on MNIST, CIFAR-100, and ImageNet, considering random samples. We analyzed mini-batch losses along the lines in the first 10000 SGD update step directions for four commonly used architectures (ResNet-32 [23], DenseNet-40 [25], EfficientNet [58], MobileNetV2 [54]). Along each line, we sampled 50 losses and performed a parabolic approximation (see Section 4). An unbiased selection of our results on a ResNet32 is shown in Figure 1. Further results are given in Appendix A. In accordance with [61], we conclude that the mini-batch losses along the analyzed lines tend to be locally convex. In addition, one-dimensional parabolic approximations of the form $f(s) = as^2 + bs + c$ with $a \neq 0$ are well suited to estimate the position of a minimum along such directions.

To substantiate the later observation, we analyzed the angle between the line direction and the gradient at the -by a parabola- estimated minimum during training. A position is a local extremum or saddle point of the loss along a line if and only if the angle between the line direction and the gradient at the position is $90°$, if measured on the same mini-batch. [1] Measuring step sizes and update step adaptations factors (see Sections 4.1 and 4.3) were chosen to fit the mini-batch loss along the line decently. As shown in Figures 2 and 3, the parabolic observation holds well for several architectures trained on MNIST, CIFAR-10, CIFAR-100 and ImageNet. The observation fits best for MNIST and gets worse for more complex tasks such as ImageNet. We can ensure that the extrema found are minima since we additionally plotted the mini-batch loss along the line for each update step.

In addition, we analyzed the mini-batch loss along lines in conjugate-like directions and random directions. Mini-batch losses along lines in conjugate like directions also tend to have convex shapes (see Appendix D.4 Figure 17 ). However, mini-batch losses along lines in random directions rarely exhibit convex shapes.
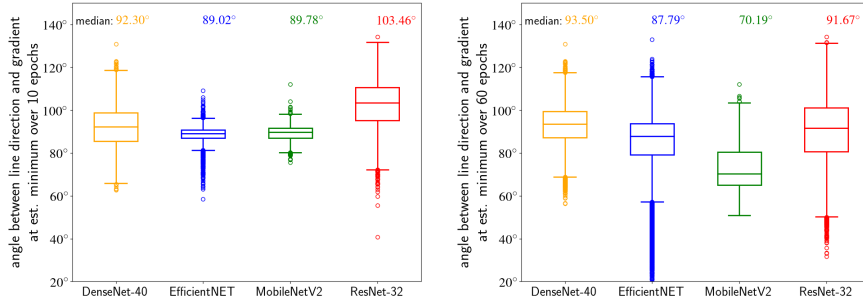


Figure 2: Angles between the line direction and the gradient at the estimated minimum measured on the same mini-batch. If the angle is $90°$, the estimated minimum is a real local minimum. We know from additional line plots that the found extrema or saddle points are minima. Left: measurement over the first 10 epochs. Right: measurement over the first 60 epochs. Update step adaptation (see Section 4.3) is applied.

---

[1]This holds because if the directional derivative of the measured gradient in line direction is 0, the current position is an extremum or saddle point along the line and the angle is $90°$. If the position is not an extremum or saddle point, the directional derivative is not 0 [28].
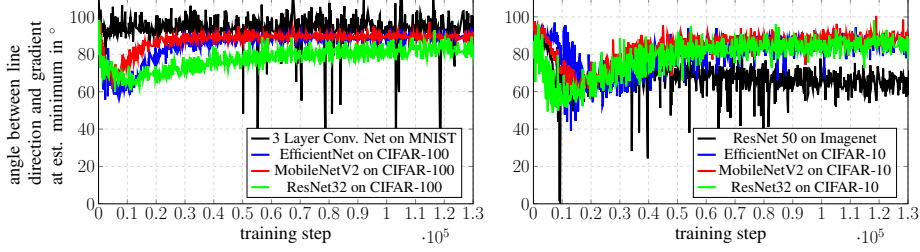
Figure 3: Angles between the line direction and the gradient at the estimated minimum measured on the same batch plotted over a whole training process on several networks and datasets. This figure clarifies that the parabolic observation holds also on further datasets and during the training process. It fits best for MNIST and becomes worse for ImageNet. Measuring step sizes and update step adaptations factors (see Sections 4.1,4.3) were used to fit the loss along the line decently.

# 4 The line search algorithm

We introduce **P**arabolic **A**pproximation **L**ine Search (*PAL*), which exploits the observation that parabolic approximations are suited to estimate the minimizer of the mini-batch loss along lines. This simple approach combines well-known methods from basic optimization such as parabolic approximation and line search [28] to perform an efficient line search. We note that the general idea of this method can be applied to any optimizer that provides an update step direction.

## 4.1 Parameter update rule

An intuitive explanation of *PAL*'s parameter update rule based on a parabolic approximation is given in Figure 4. Since $l_t$ (see Eq.1) is assumed to exhibit a convex and almost parabolic shape, we approximate it with $\hat{l}_t(s) = as^2 + bs + c$ with $a \in \mathbb{R}^+$ and $b, c \in \mathbb{R}$. Consequently, we need three measurements of $l_t$ to define $a, b$ and $c$. Those are given by the current loss $l_t(0)$, the derivative in gradient direction $l'_t(0) = -||\boldsymbol{g}_t||$ (see Eq. 4) and an additional loss $l_t(\mu)$ with measuring distance $\mu \in \mathbb{R}^+$. It is simple to show that $a = \frac{l_t(\mu) - l_t(0) - l'_t(0)\mu}{\mu^2}$, $b = l'_t(0)$, and $c = l_t(0)$. The update step $s_{\text{upd}}$ to the minimum of the parabolic approximation $\hat{l}_t(s)$ is thus given by:

$$s_{\text{upd}_t} = -\frac{\hat{l}'_t(0)}{\hat{l}''_t(0)} = -\frac{b}{2a} = \frac{-l'_t(0)}{2\frac{l_t(\mu) - l_t(0) - l'_t(0)\mu}{\mu^2}} \quad (2)$$

Note, that $\hat{l}''_t(0)$ is the second derivative of the approximated parabola and is only identical to the exact directional derivative $\frac{-\boldsymbol{g}_t}{||\boldsymbol{g}_t||} H(\mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)) \frac{-\boldsymbol{g}_t^T}{||\boldsymbol{g}_t||}$ if the parabolic approximation fits. The normalization of the gradient to unit length (Eq.1) was chosen to have the measuring distance $\mu$ independent of the gradient size and of weight scaling. Note that two network inferences are required to determine $l_t(0)$ and $l_t(\mu)$. Consequently, *PAL* needs two forward passes and one backward pass through a model. Further on, $\mathcal{L}_{\mathbb{B}_t}$ may include random components, but to ensure continuity during one line search, drawn random numbers have to be reused for each value determination of $\mathcal{L}_{\mathbb{B}_t}$ at $t$ (e.g., for Dropout [57]). The memory required by *PAL* is similar to *SGD* with momentum, since only the last update direction has to be saved. A basic version of *PAL* is given in Algorithm 1.
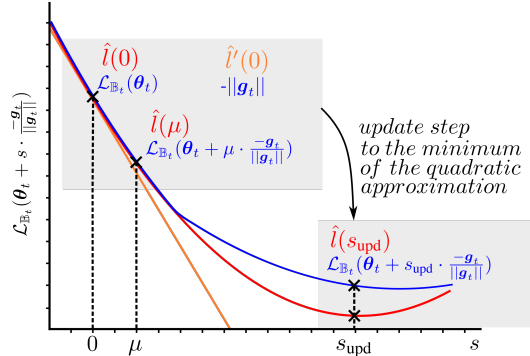


Figure 4: Basic idea of *PAL*'s parameter update rule. The blue curve is the mini-batch loss along the negative gradient starting at $\mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)$. It is defined as $l_t(s) = \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t + s\frac{-\boldsymbol{g}_t}{||\boldsymbol{g}_t||})$ where $\boldsymbol{g}_t$ is $\nabla_{\boldsymbol{\theta}_t}\mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)$. The red curve is its parabolic approximation $\hat{l}(s)$. $l_t(0)$, $l_t(\mu)$ and $\boldsymbol{g}_t$ (orange) are the three parameters needed to determine the update step $s_{upd}$ to the minimum of the parabolic approximation.

4

**Algorithm 1** The basic version of our proposed line search algorithm. See Section 4 for details.

**Input:** $\mu$: measuring step size
**Input:** $\mathcal{L}_{\mathbb{B}_t}$: mini-batch loss function
**Input:** $\boldsymbol{\theta}_0$: initial parameter vector
1: $t \leftarrow 0$
2: **while** $\boldsymbol{\theta}_t$ not converged **do**
3:     $\mathbb{B}_t \leftarrow \text{sampleBatch}()$
4:     $l_0 \leftarrow \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)$       # $l_0 = l_t(0)$ see Eq. 1
5:     $\boldsymbol{g}_t \leftarrow -\nabla_{\boldsymbol{\theta}_t} \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)$
6:     $l_\mu \leftarrow \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t + \mu \frac{\boldsymbol{g}_t}{||\boldsymbol{g}_t||})$
7:     $b \leftarrow -||\boldsymbol{g}_t||$

8:     $a \leftarrow \frac{l_\mu - l_0 - b\mu}{\mu^2}$
9:     **if** proper curvature **then**
10:       $s_{\text{upd}} \leftarrow -\frac{b}{2a}$
11:     **else**
12:       # set $s_{\text{upd}}$ according to section 4.2
13:     **end if**
14:     $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + s_{\text{upd}} \frac{\boldsymbol{g}_t}{||\boldsymbol{g}_t||}$
15:     $t \leftarrow t + 1$
16: **end while**
17: **return** $\boldsymbol{\theta}_t$

## 4.2 Case discrimination of parabolic approximations

Since not all parabolic approximations are suitable for parameter update steps, the following cases are considered separately. Note that $b = \hat{l}'(0) = l'(0)$ and $a = 0.5\hat{l}''(0) \approx 0.5 l''(0)$.. **1:** $a > 0$ and $b < 0$: the parabolic approximation has a minimum in line direction, thus, the parameter update is done as described in Section 4.1. **2:** $a \leq 0$ and $b < 0$: the parabolic approximation has a maximum in negative line direction or is a line with a negative slope. In those cases, a parabolic approximation is inappropriate. $s_{upd}$ is set to $\mu$, since the second measured point has a lower loss than the first. **3:** Since $b = -||g_t||$ cannot be greater than 0, the only case left is an extremum at the current position ($l'(0) = 0$). In this case, no weight update is performed. However, the mini-batch loss function changes with the next batch, which likely does not have an extremum at exactly the same point. In accordance with Section 3, cases two and three appeared very rarely in our experiments.

## 4.3 Additions

Here multiple additions for Algorithm 1 to fine-tune the performance and handle degenerate cases are introduced. Our hyperparameter sensitivity analysis (Appendix D.6) suggests that the influence of the introduced hyperparameters on the optimizer's performance is low. Thus, they only need to be adapted to fine-tune the results. The full version of *PAL* including all additions is given in Appendix B Algorithm 2.

**Direction adaptation:** Instead of following the direction of the negative gradient we follow an adapted conjugate-like direction $\boldsymbol{d}_t$:

$$\boldsymbol{d}_t = -\nabla_{\boldsymbol{\theta}_t} \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t) + \beta \boldsymbol{d}_{t-1} \quad \boldsymbol{d}_0 = -\nabla_{\boldsymbol{\theta}_0} \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_0) \tag{3}$$

with $\beta \in [0, 1]$. Since now an adapted direction is used, $l_t'(0)$ changes to:

$$l_t'(0) = \nabla_{\boldsymbol{\theta}_t} \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t) \frac{\boldsymbol{d}_t}{||\boldsymbol{d}_t||} \tag{4}$$

This approach aims to find a more optimal search direction than the negative gradient. We implemented and tested the formulas of Fletcher-Reeves [16], Polak-Ribière [50], Hestenes-Stiefel [24] and Dai-Yuan [11] to determine conjugate directions under the assumption that the loss function is a quadratic. However, choosing a constant $\beta$ of value $0.2$ or $0.4$ performs equally. The influence of $\beta$ and dynamic update steps on *PAL*'s performance is discussed in Appendix D.5. In the analyzed scenario, $\beta$ can both increase and decrease the performance, whereas dynamic update steps mostly increase the performance. A combination of both is needed to achieve optimal results.

**Update step adaptation:** Our preliminary experiments revealed a systematic error caused by constantly approximating with slightly too narrow parabolas. Therefore, $s_{\text{upd}}$ is multiplied by a parameter $\alpha \geq 1$ (compare to Eq. 2). This is useful to estimate the minimum's position along a line more exactly but has minor effects on training performance.

**Maximum step size:** To hinder the algorithm from failing due to inaccurate parabolic approximations, we use a maximum step size $s_{\text{max}}$. The new update step is given by $\min(s_{\text{upd}}, s_{\text{max}})$. However, most of our experiments with $s_{\text{max}} = 10^{0.5} \approx 3.16$ never reached this step size and still performed well.

## 4.4 Theoretical considerations

Usually, convergence in deep learning is shown for convex stochastic functions with an L-Lipschitz continuous gradient. However, since our approach originates from empirical results, it is not given that a profound theoretical analysis is possible. In order to show any convergence guarantees for parabolic approximations, we have to fall back to uncommonly strong assumptions, which lead to quadratic models. Since convergence proofs on quadratics are of minor importance for most readers, our derivations can be found in Appendix C.

# 5 Evaluation

## 5.1 Experimental design

We performed a comprehensive evaluation to analyze the performance of *PAL* on a variety of deep learning optimization tasks. Therefore, we tested *PAL* on commonly used architectures on CIFAR-10 [31], CIFAR-100 [31] and ImageNet [13]. For CIFAR-10 and CIFAR-100, we evaluated on DenseNet40 [25], EfficientNetB0 [58], ResNet32 [23] and MobileNetV2 [54]. On ImageNet we evaluated on DenseNet121 and ResNet50. In addition, we considered an RNN trained on the Tolstoi war and peace text prediction task. We compare *PAL* to *SLS* [60], whose Armijo variant is state-of-the-art in the line search field for DNNs at the time of writing. In addition, we compare against the following well studied and widely used first order optimizers: *SGD* with momentum [51], *ADAM* [30], and *RMSProp* [59] as well as against *SGDHD* [3] and *ALIG* [5], which automatically estimate learning rates in negative gradient direction and, finally, against the coin betting approach *COCOB* [44]. To perform a fair comparison, we compared various hyperparameter combinations of commonly used hyperparameters for each optimizer. In addition, we utilize those combinations to analyze the hyperparameter sensitivity for each optimizer. Since a grid search on Imagenet was too expensive, the best hyperparameter configuration of CIFAR-100's evaluation was used to test hyperparameter transferability. A detailed explanation of the experiments including hyperparameters and data augmentations used, are given in Appendix D.8. All in all, we trained over 4500 networks with Tensorflow 1.15 [1] on Nvidia Geforce GTX 1080 TI graphic cards. Since *PAL* is a line search approach, the predefined learning rate schedules of SGD and the generated schedules of *SLS*, *ALIG*, *SGDHD* and *PAL* were compared. Due to normalization, *PAL*'s learning rate is given by $s_{\mathrm{upd}_t}/||\boldsymbol{d}_t||$.

## 5.2 Results

A selection of our results is given in Figure 5. The results of other architectures trained on CIFAR-10, CIFAR-100, Imagenet and Tolstoi are found in Appendix D Figures 13,14,15. A table with exact numerical results of all experiments is provided in Appendix D.9.

In most cases *PAL* decreases the training loss faster and to a lower value than the other optimizers (row 1 of Figures 5,13,14,15). Considering validation and test accuracy, *PAL* surpasses *ALIG*, *SGDHD* and *COCOB*, competes with *RMSProp* and *ADAM* but gets surpassed by *SGD* (rows 2,3 of Figures 5,13,14,15). However, *RMSProp*, *ADAM* and *SGD* were tuned with a step size schedule. If we compare *PAL* to their basic implementations without a schedule, which roughly corresponds to the first plateau reached in row 2 of Figures 5,13,14,15, *PAL* would surpass the other optimizers and shows that it can find a well performing step size schedule. This is especially interesting for problems for which default schedules might not work.

*SLS* decreases the training loss further than the other optimizers on a few problems but shows weak performance and poor generalization on most. This contrasts to the results of [60], where *SLS* behaves robustly and excels. To exclude the possibility of errors on our side, we reimplemented *SLS* experiment on ResNet34 and could reproduce a similar well performance as in [60] (Appendix D.3). Our results suggest that the interpolation assumption on which *SLS* is based is not always valid for the considered tasks.

Considering the box plots of Figures 5 and 14, which represent the sensitivity to hyperparameter combinations, one would likely try on a new unknown objective, we can see that *PAL* has a strong tendency to exhibit low sensitivity in combination with good performance. To emphasize this statement, a sensitivity analysis of *PAL*'s hyperparameters (Appendix Figure 19) shows that *PAL* performs well on a wide range for each hyperparameter on a ResNet32.
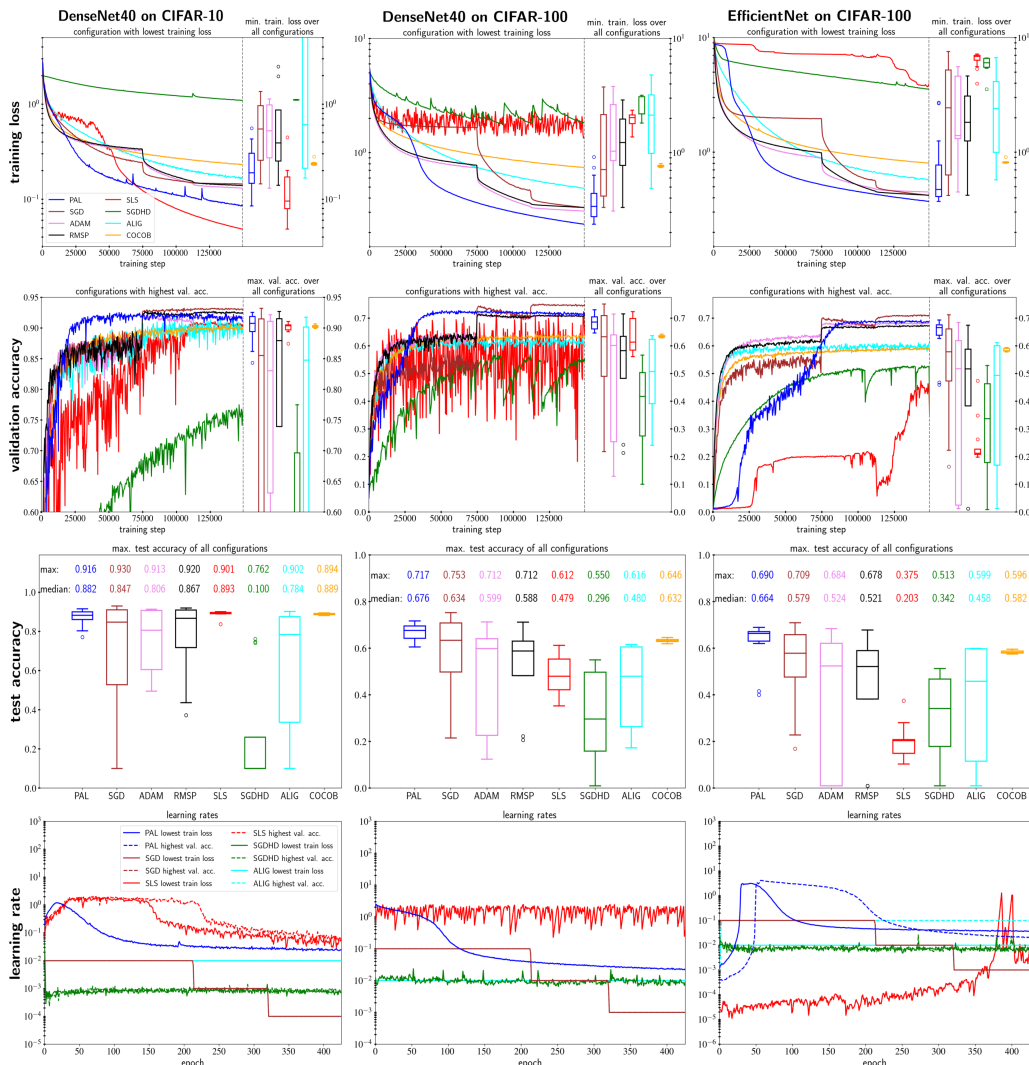
Figure 5: Comparison of *PAL* against *SLS*, *SGD*, *ADAM*, *RMSProp*, *ALIG*, *SGDHD* and *COCOB* on train. loss (row 1), val. acc. (row 2), test. acc. (row 3) and *SLS*, *SGD*, *ALIG*, *SGDHD* and *PAL* on learning rates (row 4). Comparison is done across several datasets and models. Further results are found in Appendix D.1 Figure (13,14,15). Results are averaged over 3 runs. Box plots result from comprehensive hyperparameter grid searches in plausible intervals. Learning rates are averaged over epochs. *PAL* surpasses, *ALIG*, *SGDHD*, and *COCOB* and competes against all other optimizers except against SGD.

On wall-clock-time *PAL* performs as fast as *SLS* but slower than the other optimizers, which achieve similar speeds (Appendix D.2). However, an automatic, well-performing learning rate schedule might compensate for the slower speed depending on the scenario.

Considering the learning rate schedules of *PAL* (row 4 of Figures 5,13,14,15) we achieved unexpected results. *PAL*, which estimates the learning rate directly from approximated local shape information, does not follow a schedule that is similar to the one of *SLS*, *ALIG*, *SGDHD* or any of the commonly used handcrafted schedules such as piece-wise constant or cosine decay. However, it achieves similar results. An interesting side result is that *ALIG* and *SGDHD* tend to perform best if hyperparameters are chosen in a way that the learning rate is only changed slightly and, therefore, virtually an SGD training with a fixed learning rate is performed.
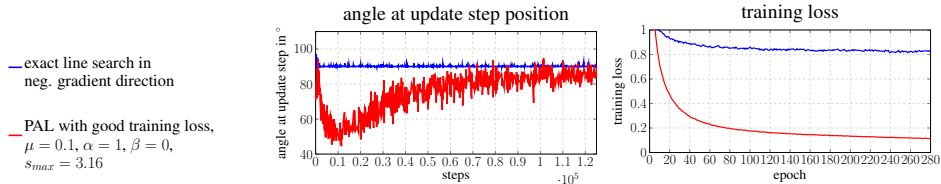
Figure 6: Comparison of *PAL* against an exact line search. The first plot shows the angle between the direction and gradient vector at the update step position. A ResNet32 was trained on CIFAR-10. One can observe that an exact line search exhibits poor performance.

# 6 On the exactness of line searches on mini-batch losses

In this section, we investigate the general question of whether line searches that estimate the position of a minimum of mini-batch losses exactly are beneficial. In Figure 2 we showed that *PAL* can perform an almost exact line search on batch losses if we use a fixed update step adaptation factor (Section 4.3). However, *PAL*'s best hyperparameter configuration does not perform an exact line search (see Figure 6). Consequently, we analyzed how an exact line search, which exactly estimates a minimum along the line, behaves. We implemented an inefficient binary line search (see Appendix E), which measured up to 20 values on each line to estimate the position of a minimum. The results, given in Figure 6, show that an optimal line search does not optimize well. Thus, the reason why *PAL* performs well is not the exactness of its update steps. In fact, slightly inexact update steps seem to be beneficial.

These results query Assumption 1, which assumes that the position of a minimum on a line in negative gradient direction of $\mathcal{L}_{\mathbb{B}_t}$ is a suitable estimator for the minimum of the full-batch loss $\mathcal{L}$ on this line to perform a successful optimization process. To investigate this further, we measured $\mathcal{L}$ and the distribution of mini-batch losses for multiple SGD update directions on a ResNet32. Our results suggest, as exemplarily shown in Figure 7 that the position of $\mathcal{L}_{\mathbb{B}_i}$'s minimum along a line in negative gradient direction is not always a good estimator for the position of $\mathcal{L}$'s corresponding minimum. This explains why exact line searches on the batch loss perform weakly.

Corollaries are that the empirical loss on the investigated lines also tends to be locally convex and that the optimal step size tends to be smaller than the step size given by the batch loss on such lines. This is a possible explanation why the slightly too narrow parabolic approximations of *PAL* without update step adaptation perform well.
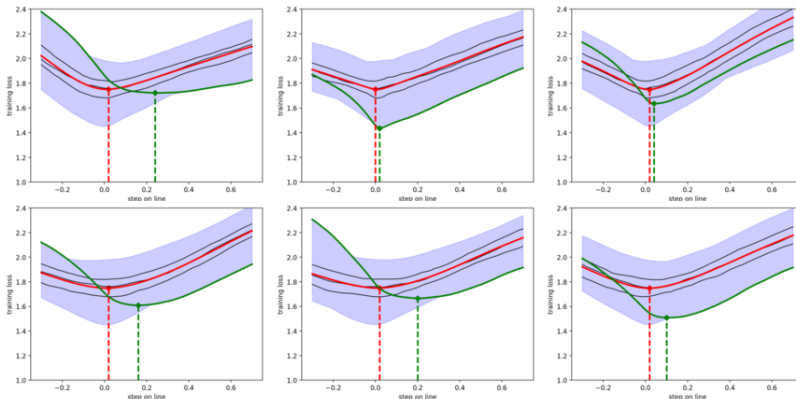


Figure 7: Distributions (blue) overall mini-batch losses along representative lines during a training of a ResNet32 on CIFAR-10. The full-batch loss, which is the mean value of the distribution, is given in red. Quartiles are given in black. The mini-batch loss, whose negative gradient defines the search direction, is given in green. It can be observed that the minimum of the green mini-batch loss is not always an adequate estimator of the minimum of the full-batch loss along the line.

8

# 7 PAL and Interpolation

This section analyzes whether the reason why *PAL* performs well is related to the interpolation condition. Formally, interpolation requires that the gradient with respect to each sample converges to zero at the optimum. We repeated the experiments of the *SLS* paper (see [60] Section 7.2 and 7.3), which analyze the performance on problems for which interpolation hold or does not hold.
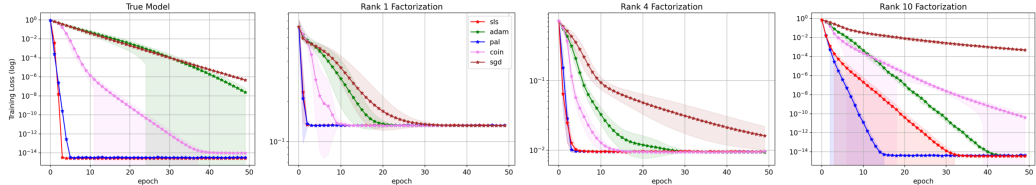


Figure 8: The matrix factorization problem of [60, §7.2] . For $k = 1$ and $k = 4$ interpolation does not hold. Rank 1 factorization is under-parameterized, whereas rank 4 and rank 10 factorizations are over-parameterized.
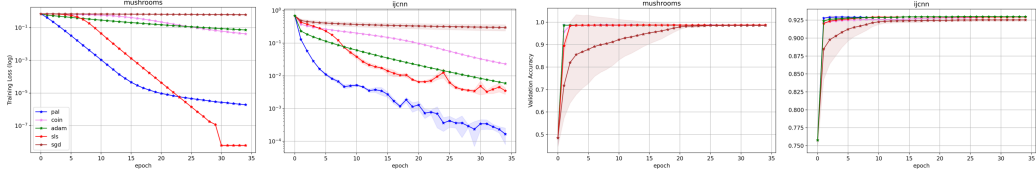


Figure 9: Binary classification task of [60, §7.3] using a softmax loss and RBF kernels for mushrooms and ijcnn datasets. With RBF kernels, the mushrooms dataset is linear separable in kernel-space with the selected kernel bandwidths, while the ijcnn dataset is not.

Figure 8 shows that *PAL* such as *SLS* converge faster to an artificial optimization floor on non-over-parameterized models ($k = 4$) of the matrix factorization problem of [60, §7.2]. In the interpolation case *PAL* and *SLS* converge linearly to machine precision. On the binary classification problem of [60] Section 7.3, which uses a softmax loss and RBF kernels on the mushrooms and ijcnn datasets, we observe that *PAL* and *SLS* converge fast on the mushrooms task, for which the interpolation condition holds (Figure 9). However, *PAL* converges faster on the ijcnn task, for which the interpolation condition does not hold.

The results indicate that the interpolation condition is beneficial for *PAL*, but, *PAL* performs also robust when it is likely not satisfied (see Figure 5,13,14,15. In those experiments *PAL* mostly performs competitive but *SLS* does not. However, the relation of the parabolic observation to interpolation needs to be investigated more closely in future.

# 8 Conclusions

This work tackles a major challenge in current optimization research for deep learning: to automatically find optimal step sizes for each update step. In detail, we focus on line search approaches to deal with this challenge. We introduced a robust and straightforward line search approach based on one-dimensional parabolic approximations of mini-batch losses. The introduced algorithm is an alternative to *SGD* for objectives where default decays are unknown or do not work.

Loss functions of DNNs are commonly perceived as being highly non-convex. Our analysis suggests that this intuition does not hold locally since lines of loss landscapes across models and datasets can be approximated parabolically to high accuracy. This new knowledge might further help to explain why update steps of specific optimizers perform well.

To gain deeper insights into line searches in general, we analyzed how an expensive but exact line search on batch losses behaves. Intriguingly, its performance is weak, which lets us conclude that the minor inaccuracies of the parabolic approximations are beneficial for training.

## Retrospective from 10.2021

After the publication of this paper, we conducted further research in this field and observed the following: *PAL*'s performance is highly dependent on the weight initialization and batch size. In general, it does not work well with Pytorch's default weight initialization. Note that the experiments here were done with Tensorflow 1.5. In addition, it is hard to find hyperparameters that compete with SGD. Our further works exploring the empirics-based line search direction are [42, 43]. [42] performs a comprehensive empirical study of the shape of the full-batch loss in line direction and the performance of several optimizers along such lines. [43] introduces a line search approach that fits the full-batch loss with a parabola and can also handle smaller batch sizes.

## Ethics Statement

Since we understand our work as basic research, it is highly error-prone to estimate its *specific* ethical aspects and future positive or negative social consequences. As optimization research influences the whole field of deep learning, we refer to the following works, which discuss the ethical aspects and social consequences of AI and Deep Learning in a comprehensive and general way: [6, 41, 63].

## References

[1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.

[2] L. Balles. Probabilistic line search tensorflow implementation, 2017.

[3] A. G. Baydin, R. Cornish, D. M. Rubio, M. Schmidt, and F. Wood. Online learning rate adaptation with hypergradient descent. *arXiv preprint arXiv:1703.04782*, 2017.

[4] A. S. Berahas, M. Jahani, and M. Takác. Quasi-newton methods for deep learning: Forget the past, just sample. *CoRR*, abs/1901.09997, 2019.

[5] L. Berrada, A. Zisserman, and M. P. Kumar. Training neural networks for and by interpolation. *arXiv preprint arXiv:1906.05661*, 2019.

[6] N. Bostrom and E. Yudkowsky. The ethics of artificial intelligence. *The Cambridge handbook of artificial intelligence*, 1:316–334, 2014.

[7] A. Botev, H. Ritter, and D. Barber. Practical gauss-newton optimisation for deep learning. *arXiv preprint arXiv:1706.03662*, 2017.

[8] A. Botev, H. Ritter, and D. Barber. Practical gauss-newton optimisation for deep learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 557–565. JMLR. org, 2017.

[9] C.-A. Brust, S. Sickert, M. Simon, E. Rodner, and J. Denzler. Neither quick nor proper - evaluation of quickprop for learning deep neural networks. *CoRR*, abs/1606.04333, 2016.

[10] Y. Chae and D. N. Wilke. Empirical study towards understanding line search approximations for training neural networks. *arXiv preprint arXiv:1909.06893*, 2019.

[11] Y.-H. Dai and Y. Yuan. A nonlinear conjugate gradient method with a strong global convergence property. *SIAM Journal on optimization*, 10(1):177–182, 1999.

[12] S. De, A. Yadav, D. Jacobs, and T. Goldstein. Big batch sgd: Automated inference using adaptive batch sizes. *arXiv preprint arXiv:1610.05792*, 2016.

[13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.

[14] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for on learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.

[15] S. E. Fahlman et al. An empirical study of learning speed in back-propagation networks. 1988.

[16] R. Fletcher and C. M. Reeves. Function minimization by conjugate gradients. *The computer journal*, 7(2):149–154, 1964.

[17] S. Fort and S. Jastrzebski. Large scale structure of neural network loss landscapes. In *Advances in Neural Information Processing Systems*, pages 6706–6714, 2019.

[18] X. Gastaldi. Shake-shake regularization. *CoRR*, abs/1705.07485, 2017.

[19] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010.

[20] I. J. Goodfellow, O. Vinyals, and A. M. Saxe. Qualitatively characterizing neural network optimization problems. *arXiv preprint arXiv:1412.6544*, 2014.

[21] Google. Tensorflow adam optimizer documentation. `https://www.tensorflow.org/api_docs/python/tf/train/AdamOptimizer`. Accessed: 2019-11-12.

[22] H. He, G. Huang, and Y. Yuan. Asymmetric valleys: Beyond sharp and flat local minima. In *Advances in Neural Information Processing Systems*, pages 2549–2560, 2019.

[23] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[24] M. R. Hestenes and E. Stiefel. *Methods of conjugate gradients for solving linear systems*, volume 49. NBS Washington, DC, 1952.

[25] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, volume 1, page 3, 2017.

[26] P. Izmailov, D. Podoprikhin, T. Garipov, D. Vetrov, and A. G. Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.

[27] P. Izmailov, D. Podoprikhin, T. Garipov, D. Vetrov, and A. G. Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.

[28] S. W. Jorge Nocedal. *Numerical Optimization*. Springer series in operations research. Springer, 2nd ed edition, 2006.

[29] D. Kafka and D. Wilke. Gradient-only line searches: An alternative to probabilistic line searches. *arXiv preprint arXiv:1903.09383*, 2019.

[30] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.

[31] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.

[32] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.

[33] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–48. Springer, 2012.

[34] H. Li, Z. Xu, G. Taylor, and T. Goldstein. Visualizing the loss landscape of neural nets. *CoRR*, abs/1712.09913, 2017.

[35] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han. On the variance of the adaptive learning rate and beyond, 2019.

[36] D. G. Luenberger, Y. Ye, et al. *Linear and nonlinear programming*, volume 2. Springer, 1984.

[37] L. Luo, Y. Xiong, and Y. Liu. Adaptive gradient methods with dynamic bound of learning rate. In *International Conference on Learning Representations*, 2019.

[38] M. Mahsereci and P. Hennig. Probabilistic line searches for stochastic optimization. *Journal of Machine Learning Research*, 18, 2017.

[39] J. Martens and R. Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417, 2015.

[40] J. J. Moré and D. J. Thuente. Line search algorithms with sufficient decrease. *ACM Transactions on Mathematical Software (TOMS)*, 20(3):286–307, 1994.

[41] L. Muehlhauser and L. Helm. The singularity and machine ethics. In *Singularity Hypotheses*, pages 101–126. Springer, 2012.

[42] M. Mutschler and A. Zell. Empirically explaining sgd from a line search perspective. *ICANN*, 2021.

[43] M. Mutschler and A. Zell. Using a one dimensional parabolic model of the full-batch loss to estimate learning rates during training. *arxiv*, 2021.

[44] F. Orabona and T. Tommasi. Training deep networks without learning rates through coin betting. In *Advances in Neural Information Processing Systems*, pages 2160–2170, 2017.

[45] C. Paquette and K. Scheinberg. A stochastic line search method with convergence rate analysis. *arXiv preprint arXiv:1807.07994*, 2018.

[46] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017.

[47] V. Ramamurthy and N. Duffy. L-sr1: a second order optimization method. 2017.

[48] N. M. P. L. Rates. Automatic and simultaneous adjustment of learning rate and momentum for stochastic gradient descent.

[49] S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.

[50] G. Ribière and E. Polak. Note sur la convergence de directions conjugées. *Rev. Francaise Informat Recherche Opertionelle*, 16:35–43, 1969.

[51] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.

[52] M. Rolinek and G. Martius. L4: Practical loss-based stepsize adaptation for deep learning. In *Advances in Neural Information Processing Systems*, pages 6434–6444, 2018.

[53] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533, 1986.

[54] M. Sandler, A. G. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018.

[55] N. N. Schraudolph, J. Yu, and S. Günter. A stochastic quasi-newton method for online convex optimization. In M. Meila and X. Shen, editors, *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, volume 2 of *Proceedings of Machine Learning Research*, pages 436–443, San Juan, Puerto Rico, 21–24 Mar 2007. PMLR.

[56] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.

[57] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.

[58] M. Tan and Q. V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.

[59] T. Tieleman and G. Hinton. Lecture 6.5-rmsprop, coursera: Neural networks for machine learning. *University of Toronto, Technical Report*, 2012.

[60] S. Vaswani, A. Mishkin, I. Laradji, M. Schmidt, G. Gidel, and S. Lacoste-Julien. Painless stochastic gradient: Interpolation, line-search, and convergence rates. pages 3727–3740, 2019.

[61] C. Xing, D. Arpit, C. Tsirigotis, and Y. Bengio. A walk with sgd. *arXiv preprint arXiv:1802.08770*, 2018.

[62] Y. Yamada, M. Iwamura, T. Akiba, and K. Kise. Shakedrop regularization for deep residual learning. 2018.

[63] E. Yudkowsky et al. Artificial intelligence as a positive and negative factor in global risk. *Global catastrophic risks*, 1(303):184, 2008.

[64] M. D. Zeiler. Adadelta: An adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.

# A  Further line plots



Figure 10: **DenseNet40** mini-batch losses along in negative gradient direction (blue) combined with our parabolic approximation (orange) and the position of the minimum (red). The unit of the horizontal axis is the change of $\theta$.



Figure 11: Mini-batch losses along lines of **MobileNetV2**. For explanations see Figure 10. During training the parabolic approximation fits less accurately on the right hand side and during the first 150 steps it does not fit at all, however, the minimum of the parabola is still a good estimator for a low loss value on the line.



Figure 12: Mini-batch losses along lines of **EfficientNet**. For explanations see Figure 10. The parabolic approximation fits only on the left hand side and during the first 150 steps it does not fit at all, however, the minimum of the parabola is still a good estimator for a low loss value on the line.

# B PAL with all additions:

Algorithm 2 provides the full description of *PAL* including the additions described in Section 4.3. After analyzing the best hyperparameter combinations of *PAL* over all experiments, we suggest to use values from the following parameter intervals: $\mu = [0.1, 1]$, $\alpha = [1.0, 1.6]$, $\beta = [0, 0.4]$, $s_{\max} = 3$. Where $\alpha, \beta$ and $s_{\max}$ usually have a low sensitivity. Thus, the basic implementation of *PAL* (Algorithm 1) performs already well and is always found in the upper qua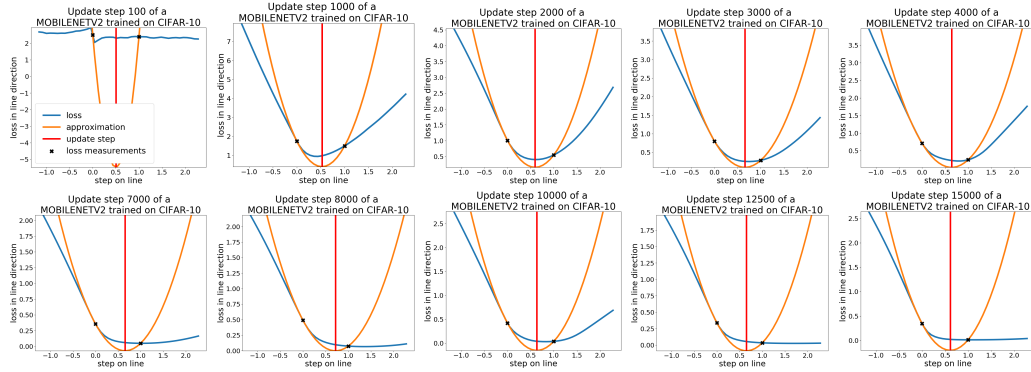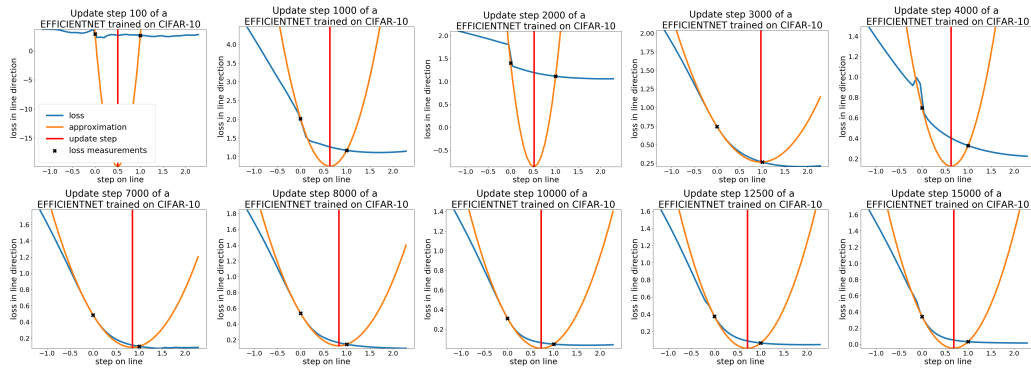rtile considering the performance of all analyzed hyperparameter configurations in our experiments. PyTorch and Tensorflow 1.5 implementations are provided at `https://github.com/cogsys-tuebingen/PAL`.

---

**Algorithm 2** *PAL*, our proposed line search algorithm for DNNs. See Section 4 for details.

---

**Input:** Hyperparameters: $\mu$: measuring step size, $\alpha$: update step adaptation, $\beta$: direction adaptation factor, $s_{\max}$: maximum step size.
**Input:** $\mathcal{L}_{\mathbb{B}_t}$: mini-batch loss function
**Input:** $\boldsymbol{\theta}_0$: initial parameter vector
1: $t \leftarrow 0$
2: $\boldsymbol{d}_t \leftarrow \boldsymbol{0}$
3: **while** $\boldsymbol{\theta}_t$ not converged **do**
4:     $\mathbb{B}_t \leftarrow \text{sampleBatch}()$
5:     $l_0 \leftarrow \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)$
6:     $\boldsymbol{d}_t \leftarrow -\nabla_{\boldsymbol{\theta}_t}\mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t) + \beta\boldsymbol{d}_{t-1}$
7:     $l_\mu \leftarrow \mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t + \mu\frac{\boldsymbol{d}_t}{||\boldsymbol{d}_t||})$
8:     $b \leftarrow \nabla_{\boldsymbol{\theta}_t}\mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)\frac{\boldsymbol{d}_t}{||\boldsymbol{d}_t||}$
9:     $a \leftarrow \frac{l_\mu - l_0 - b\mu}{\mu^2}$
10:    **if** $a > 0$ **and** $b < 0$ **then**
11:     $s_{\text{upd}} \leftarrow -\alpha\frac{b}{2a}$
12:    **else if** $a \leq 0$ **and** $b < 0$ **then**
13:     $s_{\text{upd}} \leftarrow \mu$
14:    **else**
15:     $s_{\text{upd}} \leftarrow 0$
16:    **end if**
17:    **if** $s_{\text{upd}} > s_{\max}$ **then**
18:     $s_{\text{upd}} \leftarrow s_{\max}$
19:    **end if**
20:    $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + s_{\text{upd}}\frac{\boldsymbol{d}_t}{||\boldsymbol{d}_t||}$
21:    $t \leftarrow t + 1$
22: **end while**
23: **return** $\boldsymbol{\theta}_t$

---

# C Further Theoretical Considerations

We have to note that the following derivation is based upon **strong assumptions** and if they are valid at all, they are likely more valid locally than globally. We assume that each slice of the loss function is a one-dimensional parabolic function:

**Assumption 2.** *Let $n \in \mathbb{N}$ be the number of parameters and let $\boldsymbol{l}, \boldsymbol{d} \in \mathbb{R}^n$ be vectors. Then for all $\boldsymbol{l}, \boldsymbol{d}$ there exists $a, b, c \in \mathbb{R}$ with $a > 0$, such that $\mathcal{L}_{\mathbb{B}}(\boldsymbol{l} + \boldsymbol{d}s) = as^2 + bs + c$ for all $s \in \mathbb{R}$.*

This **strong** assumption is a simplified adaptation to our empirical results that lines in negative gradient direction behave locally almost parabolic (see Section 3). For the following derivations we assume a basic *PAL* without the additions introduced in Section 4.3. Proofs are provided in Appendix C.1. At first we show that $\mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta})$ is a n-dimensional parabolic function:

**Lemma 1.** *Let $f : \mathbb{R}^n \to \mathbb{R}$ be a k-times continuously differentiable function. Furthermore, assume there exists $a, b, c \in \mathbb{R}$ with $a > 0$, such that $f(\boldsymbol{l} + \boldsymbol{d}s) = as^2 + bs + c$ for all $s \in \mathbb{R}$. Then there exist $z \in \mathbb{R}, \boldsymbol{r} \in \mathbb{R}^n$ and a positive definite Matrix $\boldsymbol{Q} \in \mathbb{R}^{n \times n}$ such that $f(\boldsymbol{x}) = c + \boldsymbol{r}^T\boldsymbol{x} + \boldsymbol{x}^T\boldsymbol{Q}\boldsymbol{x}$ for all $\boldsymbol{x} \in \mathbb{R}^n$.*

Now we show that *PAL* converges on $\mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta})$:

**Proposition 1.** *PAL converges on $f : \mathbb{R}^n \to \mathbb{R}, \boldsymbol{x} \mapsto c + \boldsymbol{r}^T\boldsymbol{x} + \boldsymbol{x}^T\boldsymbol{Q}\boldsymbol{x}$ with $\boldsymbol{Q} \in \mathbb{R}^{n \times n}$ hermitian and positive definite.*

We have to note that in this scenario, PAL is identical to the method of steepest descent for which the convergence, including convergence rates on quadratics, is already proven in [36, Page 235]. Nevertheless, we have attached our proof better adapted to this scenario.

For a noisy scenario where each batch defines a quadratic, *PAL* has no convergence guarantee. Given two shifted one-dimensional parabolas, $ax^2 + bx + c$ and $a(x + d)^2 + b(x + d) + c$, which are

presented to *PAL* alternately, *PAL* will always perform an update step to the minimum position of one of these but never to the minimum position of the average of both. By slightly changing the training procedure and assuming that each $\mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta})$ has the same $\boldsymbol{Q}$ this can be fixed:

**Proposition 2.** *If* $\mathcal{L}(\boldsymbol{\theta}) : \mathbb{R}^n \to \mathbb{R}$ $\boldsymbol{\theta} \mapsto \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} c_i + \boldsymbol{r}_i^T \boldsymbol{\theta} + \boldsymbol{\theta}^T \boldsymbol{Q}_i \boldsymbol{\theta}$ *and* $c_i + \boldsymbol{r}_i^T \boldsymbol{\theta} + \boldsymbol{\theta}^T \boldsymbol{Q}_i \boldsymbol{\theta} = \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta})$ *with* $m$ *being number of batches* $\mathbb{B}_i$. *(Each batch defines a parabola. The empirical loss* $\mathcal{L}(\boldsymbol{\theta})$ *is the mean of these parabolas). And for all* $i, j \in \mathbb{N}$ *it holds that* $\boldsymbol{Q}_i = \boldsymbol{Q}_j$ *and that* $\boldsymbol{Q}_i$ *is positive definite. Then* $\arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta})$ *holds.*

This implies that under Assumption 2 and a fixed $\boldsymbol{Q}$ the position of the minimum of the empirical loss is given by the average of the positions of the mini-batch losses' minima. The minimum position of full-batch loss is found by *PAL*, by slightly adapting *PAL* to search on one batch until it finds the minimum's position and then averaging the minimum of each batch. As a result, *PAL* converges in this noisy scenario. However, we have to emphasize at this point that our assumptions about $\boldsymbol{l}$ and $\boldsymbol{Q}$ are likely not valid for general deep learning scenarios. Nevertheless, if it is locally valid, this direction might be a further explanation, in addition to those of [17, 22], why averaging the weights of each update step in the last epoch (stochastic weight averaging) [27] performs well.

## C.1  Proofs

**Lemma 1.** *Let* $f : \mathbb{R}^n \to \mathbb{R}$ *be a k-times continuously differentiable function. Furthermore, assume there exists* $a, b, c \in \mathbb{R}$ *with* $a > 0$, *such that* $f(\boldsymbol{l} + \boldsymbol{d}s) = as^2 + bs + c$ *for all* $s \in \mathbb{R}$. *Then there exist* $z \in \mathbb{R}, \boldsymbol{r} \in \mathbb{R}^n$ *and a positive definite Matrix* $\boldsymbol{Q} \in \mathbb{R}^{n \times n}$ *such that* $f(\boldsymbol{x}) = c + \boldsymbol{r}^T \boldsymbol{x} + \boldsymbol{x}^T \boldsymbol{Q} \boldsymbol{x}$ *for all* $\boldsymbol{x} \in \mathbb{R}^n$.

*Proof.*

$$g(\boldsymbol{x}) = u + \boldsymbol{v}^T \boldsymbol{x} + \boldsymbol{x}^T \boldsymbol{W} \boldsymbol{x} \text{ for some } u \in \mathbb{R}, \boldsymbol{v} \in \mathbb{R}^n \text{ and } \boldsymbol{W} \in \mathbb{R}^{n \times n}$$

$$\Leftrightarrow \forall \boldsymbol{l}, \boldsymbol{d} \in \mathbb{R}^n \wedge ||\boldsymbol{d}|| = 1 : \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{l=1}^{n} \frac{\partial^3 g(\boldsymbol{l})}{\partial x_j, \partial x_k, \partial x_l} d_j d_k d_l = 0 \tag{5}$$

$\Rightarrow$ holds since we have a polynomial of degree 2 and its third derivative is always a $\boldsymbol{0}$ tensor.
$\Leftarrow$ holds since the reminder of the quadratic Taylor expansion is always 0.

In our case the right part is 0 since:

$$\sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{l=1}^{n} \frac{\partial^3 f(\boldsymbol{l})}{\partial x_j, \partial x_k, \partial x_l} d_j d_k d_l = \frac{\partial}{\partial s^3} f(\boldsymbol{l} + \boldsymbol{d}s) = 0 \tag{6}$$

In words: $f(\boldsymbol{x})$ is a parabolic function if and only if for each location $\boldsymbol{l}$ the third directional derivative of $f(\boldsymbol{l})$ in each direction $\boldsymbol{d}$ is 0. Which is the case, since the third derivative of each intersection is 0. $\boldsymbol{W}$ is positive definite since:

$$\forall \boldsymbol{d}, \boldsymbol{l} \in \mathbb{R}^n \wedge ||\boldsymbol{d}|| = 1 : \boldsymbol{d}^T \boldsymbol{W} \boldsymbol{d} = \frac{1}{2} \boldsymbol{d}^T \boldsymbol{H}(\boldsymbol{l}) \boldsymbol{d} = \frac{1}{2} \frac{\partial}{\partial s^2} f(\boldsymbol{l} + \boldsymbol{d}s) = a > 0 \tag{7}$$

where $\boldsymbol{H}$ is the Hessian. ∎

**Proposition 1.** *PAL converges on $f : \mathbb{R}^n \to \mathbb{R}, x \mapsto c + r^T x + x^T Q x$ with $Q \in \mathbb{R}^{n \times n}$ hermitian and positive definite.*

*Proof.*
For this prove we consider a basic *PAL* without the features introduced in Section 4.3. Note, that during the proof we will see, that $a > 0$ and $b < 0$. Thus, only the update step for this case has to be considered (see Section 4.2.

$f(x)$ is convex since $Q$ is positive definite. Thus it has one minimum.
Without loss of generality we set $c = 0, r = 0, x_n \neq 0$

$$f(x) = x^T Q x \text{ and } \nabla_x f(x) = f'(x) = 2 Q x \tag{8}$$

The values of $f(x)$ along a line through $x$ in the direction of $-f'(x)$ are given by:

$$f(-f'(x)\hat{s} + x) \tag{9}$$

Now we expand the line function:

$$
\begin{aligned}
f(-f'(x)\hat{s} + x) &= f(-2Qx\hat{s} + x) \\
&= (-2Qx\hat{s} + x)^T Q(-2Qx\hat{s} + x) \\
&= \underbrace{4x^T Q^3 x}_{=:a} \hat{s}^2 + \underbrace{-4x^T Q^2 x}_{=:b} \hat{s} + \underbrace{x^T Q x}_{=:c}
\end{aligned} \tag{10}
$$

Here we see that $f(\hat{s})$ is indeed a parabolic function with $a > 0$, $b < 0$ and $c > 0$ since $Q^3, Q^2$ and $Q$ are positive definite.
The location of the minimum $s_{min}$ of $f(\hat{s})$ is given by:

$$\hat{s}_{min} = \arg \min_{\hat{s}} f(-f'(x)\hat{s} + x) = -\frac{b}{2a} \tag{11}$$

*PAL* determines $\hat{s}_{min}$ exactly with $\hat{s}_{min} = \frac{s_{upd}}{||f'(x)||}$ (see equation 1 and 2). $||f'(x)|| > 0$ since otherwise we are already in the minimum.

The value at the minimum is given by:

$$f(\hat{s}_{min}) = a(\frac{-b}{2a})^2 + b(\frac{-b}{2a}) + c = -\frac{b^2}{4a} + c = -\underbrace{\frac{(-x^T Q^2 x)^2}{x^T Q^3 x}}_{=:g(x)} + x^T Q x = -g(x) + f(x) \tag{12}$$

Since $Q^2$ and $Q^3$ are positive definite and $x \neq 0$:

$$g(x) > 0 \tag{13}$$

Now we consider the sequence $f(x_n)$, with $x_n$ defined by *PAL* (see Equation 1):

$$x_{n+1} = -\frac{f'(x_n)}{||f'(x_n)||} \hat{s}_{upd} + x_n = -f'(x_n)\hat{s}_{min} + x_n \tag{14}$$

It is easily seen by induction that:

$$0 < f(x_{n+1}) < f(x_n) = \sum_{i=0}^{n-1} -g(x_i) + f(x_0) < f(x_0). \tag{15}$$

$g(x_n)$ converges to 0. Since $\forall n : g(x_n) > 0$ and $\sum_{i=0}^{n-1} -g(x_i)$ is bounded.

Now we have to show that $x_n$ converges to 0.
We have:

$$g(x_n) = \frac{(x_n^T Q^2 x_n)^2}{x_n^T Q^3 x_n} = \frac{\langle x_n, Q^2 x_n \rangle^2}{\langle x_n, Q^3 x_n \rangle} \tag{16}$$

Now we use the theorem of Courant-Fischer:

$$\langle x, x \rangle \min\{\lambda_1, \ldots, \lambda_n\} \le \langle x, Ax \rangle \le \langle x, x \rangle \max\{\lambda_1, \ldots, \lambda_n\}$$

$$\text{for any symmetric } A \in \mathbb{R}^{n \times n} \text{ with } \lambda_1, \ldots, \lambda_n \tag{17}$$

And get:

$$g(\boldsymbol{x}_n) \ge \frac{\lambda^2_{\boldsymbol{Q}^2 \min} \langle \boldsymbol{x}_n, \boldsymbol{x}_n \rangle^2}{\lambda_{\boldsymbol{Q}^3 \max} \langle \boldsymbol{x}_n, \boldsymbol{x}_n \rangle} = C \frac{||\boldsymbol{x}_n||^4}{||\boldsymbol{x}_n||^2} = C ||\boldsymbol{x}_n||^2 \tag{18}$$

with

$$C = \frac{\lambda^2_{\boldsymbol{Q}^2 \min}}{\lambda_{\boldsymbol{Q}^3 \max}} > 0 \text{ since all } \lambda \text{ of the positive definite } \boldsymbol{Q} \text{ are positive} \tag{19}$$

Thus, we have:

$$g(\boldsymbol{x}_n) \ge C ||\boldsymbol{x}_n||^2 \ge 0 \tag{20}$$

Since $g(\boldsymbol{x}_n)$ converges to 0, $C ||\boldsymbol{x}_n||^2$ converges to 0.
This means, $\boldsymbol{x}_n$ converges to $\boldsymbol{0}$, which is the location of the minimum. $\blacksquare$

**Proposition 3.** *If $\mathcal{L}(\boldsymbol{\theta}) : \mathbb{R}^n \to \mathbb{R} \ \boldsymbol{\theta} \mapsto \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m c_i + \boldsymbol{r}_i^T \boldsymbol{\theta} + \boldsymbol{\theta}^T \boldsymbol{Q}_i \boldsymbol{\theta}$ and $c_i + \boldsymbol{r}_i^T \boldsymbol{\theta} + \boldsymbol{\theta}^T \boldsymbol{Q}_i \boldsymbol{\theta} = \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta})$ with $m$ being number the of batches $\mathbb{B}_i$. (Each batch defines a parabola. The empirical loss $\mathcal{L}(\boldsymbol{\theta})$ is the mean of these parabolas). And for all $i, j \in \mathbb{N}$ it holds that $\boldsymbol{Q}_i = \boldsymbol{Q}_j$ and that $\boldsymbol{Q}_i$ is positive definite. Then $\arg\min_{\theta} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \arg\min_{\theta} \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta})$ holds.*

*Proof.*
Since $\mathcal{L}(\boldsymbol{\theta})$ is a sum of convex functions, it is also convex and has one minimum.
At first we determine the derivative of $\mathcal{L}(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$:

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m (\boldsymbol{r}_i + 2\boldsymbol{Q}_i \boldsymbol{\theta}) = 2\boldsymbol{Q}\boldsymbol{\theta} + \frac{1}{m} \sum_{i=1}^m \boldsymbol{r}_i \tag{21}$$

Then we determine the minima:

$$\arg\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \Leftrightarrow \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{0} \Leftrightarrow \boldsymbol{\theta} = -\frac{1}{2} \left(\sum_{i=1}^m \boldsymbol{Q}_i\right)^{-1} \sum_{i=1}^m \boldsymbol{r}_i = -\frac{1}{2m} \boldsymbol{Q}^{-1} \sum_{i=1}^m \boldsymbol{r}_i \tag{22}$$

$$\arg\min_{\boldsymbol{\theta}} \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta}) = -\frac{1}{2} \boldsymbol{Q}^{-1} \boldsymbol{r}_i \tag{23}$$

Thus, we get:

$$\arg\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = -\frac{1}{2m} \boldsymbol{Q}^{-1} \sum_{i=1}^m \boldsymbol{r}_i = \frac{1}{m} \sum_{i=1}^m -\frac{1}{2} \boldsymbol{Q}^{-1} \boldsymbol{r}_i = \frac{1}{m} \sum_{i=1}^m \arg\min_{\boldsymbol{\theta}} \mathcal{L}_{\mathbb{B}_i}(\boldsymbol{\theta}) \tag{24}$$

$\blacksquare$

# D   Further experimental results

## D.1   Performance Comparison on ImageNet, CIFAR-10, CIFAR-100 and Tolstoi



Figure 13: Comparison of *PAL* on **CIFAR-10** against *SLS*, *SGD*, *ADAM*, *RMSProp*, *ALIG*, *SGDHD* and *COCOB* on train. loss (row 1), val. acc. (row 2), test. acc. (row 3) and *SLS*, *SGD*, *ALIG*, *SGDHD* and *PAL* on learning rates (row 4). Results are averaged over 3 runs. Box plots result from comprehensive hyperparameter grid searches in plausible intervals. Learning rates are averaged over epochs. *PAL* surpasses *SLS*, *ALIG*, *SGDHD* and competes against all other optimizers except against SGD. The learning rate schedule comparison shows that *PAL* performs competitive although elaborating significantly different schedules.

Figure 14: Comparison of *PAL* on **CIFAR-100** against *SLS*, *SGD*, *ADAM*, *RMSProp*, *ALIG*, *SGDHD* and *COCOB* on train. loss (row 1), val. acc. (row 2), test. acc. (row 3) and *SLS*, *SGD*, *ALIG*, *SGDHD* and *PAL* on learning rates (row 4)). Results are averaged over 3 runs. Box plots result from comprehensive hyperparameter grid searches in plausible intervals. Learning rates are averaged over epochs. *PAL* surpasses *SLS*, *ALIG*, *SGDHD* and competes against all other optimizers except against SGD. The learning rate schedule comparison shows that *PAL* performs competitive although elaborating significantly different schedules.

19

Figure 15: Comparison of *PAL* to *SGD*, *SLS*, *ADAM*, *RMSProp* on training loss, validation accuracy and learning rates on **Imagenet**, and a simple RNN, trained on the **Tolstoi** War and Peace dataset. Learning rates are averaged over epochs. For Imagenet the best hyperparameter configuration from the CIFAR-100 evaluation were used to test hyperparameter transferability.

20

## D.2 Wall-clock time comparison

Table 1: Required seconds per epoch of *PAL*, *SLS*, *ALIG*, *SGDHD*, *COCOB* and *SGD* on CIFAR-10. RMSP and ADAM reach a similar speed as SGD. The comparison was performed on a Nvidia Geforce GTX 1080 TI. *PAL* and *SLS* perform slower, since they have to measure additional losses, whereas the additional operations of *ALIG*, *SGDHD*, *COCOB* tend to be cheap.

| network | seconds / epoch *PAL* | *SLS* | *SGD* | *ALIG* | *SGDHD* | *COCOB* |
|---|---|---|---|---|---|---|
| ResNet32 | 20.9 | 21.7 | 10.7 | 11.0 | 11.1 | 16.4 |
| MobilenetV2 | 53.2 | 52.4 | 34.1 | 34.01 | 34.2 | 36.6 |
| EfficientNet | 55.5 | 52.2 | 30.7 | 31.2 | 32.2 | 37.5 |
| DenseNet40 | 88.8 | 87.5 | 59.7 | 61.3 | 64.6 | 61.4 |

## D.3 SLS ResNet34 test case re-implementation

In the shown experiments and in contrast to the evaluation of *SLS* in [60], we used Tensorflow default Xavier weight initialization [19] versus PyTorch default Lecun initialization [33]. In addition, we used L2 regularisation versus no regularization. Furthermore, default implementations of networks for both frameworks have small differences. All in all, those differences usually influence the optimizer performance only marginally, as given by the fact that all other investigated optimizers perform well. However, in this case of *SLS* we see significant differences.

To prove that our implementation of *SLS* is correct, we re-implemented [60]'s ResNet34 test case on CIFAR-10 in Tensorflow and achieved similar results as [60]. SLS shows well performance and is not significantly overfitting as it does in in Section 5.2.



Figure 16: On the re-implemented ResNet34 test case of [60] SLS shows well performance and is not significantly overfitting as it does in in Section 5.2

### D.4 Parabolic property in adapted directions:



Figure 17: Angles between the line direction and the gradient at the estimated minimum measured on the same batch are plotted over a whole training process on several networks on CIFAR-10. This figure clarifies, that parabolic observation is also valid if a **direction adaptation factor of 0.4** is applied. Measuring step sizes and update step adaptations factors (see Sections 4.1,4.3) were set to fit loss along the line decently.

### D.5 Influence of dynamic step sizes and the direction adaptation

This section analyses, whether *PAL*'s performance originates from dynamically chosen step sizes or from the the non-linear conjugate gradient like update step adaptation. We consider EfficientNets trained on CIFAR-10, since for those the update step adaptation factor $\beta$ is needed to achieve optimal results. We consider the following 6 scenarios: **1,2)** *PAL* without update step adaptation ($\beta = 0$) and with and without dynamic step sizes (Figure 18 left). **3,4)** *PAL* with a update step adaptation of 0.2 and with and without dynamic step sizes (Figure 18 middle). **5,6)** *PAL* with a update step adaptation of 0.4 with and without fixed step sizes (Figure 18 right). The case with fixed step sizes result in in normalized SGD (NSGD) with a momentum factor $\beta$. As fixed update step size we use the measuring step size $\mu$.

The results show that dynamic step sizes increase the performance always if direction adaptation is not applied and if it is applied in 6 out of 8 cases. Direction adaptation can increase or decrease the performance in both, the dynamic and the fixed step size cases. The best performance is achieved with a direction adaptation factor of 0.2 and a measuring step size of $10^{-1.5}$, which shows that both factors influence the best results in this scenario.



Figure 18: Analysis of the influences of dynamic step sizes and the direction adaptation factor $\beta$.

## D.6 Sensitivity analysis:

All in all *PAL* tends to have a low hyperparameter sensitivity as shown in Figure 19. Since $\mu$ is the most sensitive hyperparameter we analyzed its sensitivity over several further models trained on CIFAR-10 (see Figure 20).



Figure 19: Sensitivity analysis for PAL on a ResNet32 trained on CIFAR-10. The baseline parameters are: $\mu = 0.1, \beta = 0.2, \alpha = 1.0, s_{max} = 10$. It shows that $\beta$ should be chosen $\leq 0.6$. $\alpha$ has a low sensitivity, but with a value of $1.4$ it reaches best performance. $s_{max}$ has a low sensitivity and all investigated values perform similarly. $\mu$ should be chosen between $10^{-2}$ and $10^{-0.5}$.



Figure 20: Sensitivity of the measuring step size $\mu$ of PAL for several models on CIFAR-10. *PAL* shows low sensitivity.

## D.7 Comparison to Probabilistic Line-Search (PLS):

We used a empirically improved and only existing implementation of *PLS* [38] for Tensorflow 1 [2]. However, the sum of squared gradients has to be derived manually for each layer, which is a considerable amount of work for modern architectures. Consequently, we limit our comparison to a ResNet-32 trained on CIFAR-10. Figure 21 shows that *PAL* and *PLS* perform similarly in this scenario. Default hyperparameters were used.



Figure 21: Comparison of PAL to Probabilistic Line Search [38]

## D.8 Further experimental design details

### D.8.1 Training Procedure

On CIFAR-10 and CIFAR-100 we trained 150k steps. On Imagenet each network was trained for 500k steps. We performed a piecewise constant learning rate decay by dividing the learning rate by 10 at 50% and 75% of the steps.
The training set to evaluation set split was 45k to 15k for CIFAR-10 and CIFAR-100. At the time of writing, the default Tensorflow classes do not support the reuse of the same randomly sampled numbers for multiple inferences, therefore, we implemented and used our own Dropout [57] layer.
To get a fair comparison of the optimizers capabilities, we compare on the training loss, the validation accuracy and the test accuracy metrics. For all metrics we provide the median and the quartiles to analyze the hyperparameter sensitivity. For each hyperparameter combination we averaged our results over 3 runs using the seeds 1, 2 and 3 for reproducibility. All in all, we trained over 4500 networks with Tensorflow 1.15 [1] on Nvidia Geforce GTX 1080 TI graphic cards.

### D.8.2 Data Augmentation

On CIFAR-10 we performed the following augmentations [23]:
4 pixel padding and cropping, horizontal image flipping with probability 0.5.
On Imagenet we applied an initial random crop to 224x224 pixels. In addition, we applied lighting as described in [32]. For CIFAR-10, CIFAR-100 all images were normalized by channel-wise mean and variance. For the Tolstoi War and Peace dataset we omitted augmentation.

### D.8.3 Hyperparameter grid search

For our evaluation we used all combinations out of the following common used hyperparameters. The batch size is always 128 except for DenseNets trained with *ALIG*, *SGDHD* and *COCOB* for which we encountered memory overflows and had to reduce the batch size to 100. Weight decay is always $10^{-4}$.

On Imagenet, such a comprehensive grid search was not possible. In this case we compared with the best hyperparameter combinations found on Cifar-100.

*ADAM*:

| hyperparameter | symbol | values |
|---|---|---|
| learning rate | $\lambda$ | $\{1, 0.1, 0.01, 0.001, 0.0001\}$ |
| first momentum | $\beta_1$ | $\{0.9, 0.95\}$ |
| second momentum | $\beta_2$ | $\{0.999\}$ |
| epsilon | $\epsilon$ | $\{1e-8\}$ |

We did not vary the first or second momentum much, since [30] states that the values chosen are already good defaults.

*SGD*:

| hyperparameter | symbol | values |
|---|---|---|
| learning rate | $\lambda$ | $\{0.1, 0.01, 0.001, 0.0001\}$ |
| momentum | $\alpha$ | $\{0.85, 0.9, 0.95\}$ |

*RMSProp*:

| hyperparameter | symbol | values |
|---|---|---|
| learning rate | $\lambda$ | $\{0.1, 0.01, 0.001, 0.0001\}$ |
| discounting factor | $f$ | $\{0.9, 0.95\}$ |
| epsilon | $\epsilon$ | $\{1e-8\}$ |

*PAL*:

| hyperparameter | symbol | values |
|---|---|---|
| measuring step size | $\mu$ | $\{10^0, 10^{-0.5}, 10^{-0.1}, 10^{-0.15}\}$ |
| direction adaptation factor | $\beta$ | $\{0, 0.4\}$ |
| update step adaptation | $\alpha$ | $\{1, \frac{1}{0.8}\}$ |
| maximum step size | $s_{max}$ | $\{10^{0.5}(\approx 3.16)\}$ |

In our implementation we worked with a inverse update step adaptation $\gamma = \frac{1}{\alpha}$.

*SLS*:

| hyperparameter | symbol | values |
|---|---|---|
| initial step size | $\mu$ | $\{0.1, 1\}$ |
| step size decay | $\beta$ | $\{0.9, 0.99\}$ |
| step size reset | $\gamma$ | $\{2.0, 2.5\}$ |
| Armijo constant | $c$ | $\{0.1, 0.01\}$ |
| maximum step size | $\mu_{max}$ | $\{10.0\}$ |

*ALIG*:

| hyperparameter | symbol | values |
|---|---|---|
| maximal learning rate | $\lambda$ | $\{10, 1.0, 0.1, 0.01\}$ |
| momentum | $\beta$ | $\{0.85, 0.9, 0.95\}$ |

*COCOB*:

| hyperparameter | symbol | values |
|---|---|---|
| restriction factor | $\alpha$ | $\{25, 50, 75, 100, 125, 150, 175, 200\}$ |

*SGDHD*:

| hyperparameter | symbol | values |
|---|---|---|
| learning rate | $\lambda$ | $\{0.1, 0.01, 0.001\}$ |
| hyper gradient learning rate | $\beta$ | $\{0.1, 0.01, 0.001, 0.0001\}$ |

## D.9 Detailed numerical results

Table 2: Performance comparison of *PAL*, *RMSProp*, *ADAM*, *COCOB*, *SGDHD*, *ALIG* and *SGD*. All hyperparameter combinations given in Appendix D.8 were evaluated for each architecture. Results are averaged over 3 runs starting from different random seeds, except for training on ImageNet, for which results were not averaged. Note that tests on Imagenet were performed with the best hyperparameters found on CIFAR-100 to test the transferability of hyperparameters. Medians an Quartiles describe the distribution of results over reasonable hyper-parameter ranges.

| dataset | network | optimizer | training loss min | median; p25; p75 | validation accuracy max | median; p25; p75 | test accuracy max | median; p25; p75 |
|---|---|---|---|---|---|---|---|---|
| CIFAR-10 | EfficientNet | COCOB | 0.659 | 0.824; 0.739; 0.855 | 0.857 | 0.837; 0.832; 0.845 | 0.843 | 0.824; 0.818; 0.832 |
| | | ALIG | 0.279 | 0.89; 0.464; 1.911 | 0.906 | 0.805; 0.451; 0.895 | 0.893 | 0.757; 0.297; 0.878 |
| | | SGDHD | 2.002 | 6.239; 4.357; 7.803 | 0.834 | 0.657; 0.18; 0.74 | 0.828 | 0.647; 0.179; 0.731 |
| | | SLS | 2.837 | 5.596; 4.681; 6.292 | 0.653 | 0.357; 0.211; 0.443 | 0.643 | 0.357; 0.216; 0.442 |
| | | RMSP | 0.154 | 0.637; 0.333; 1.261 | **0.93** | 0.864; 0.658; 0.902 | 0.919 | 0.854; 0.648; 0.889 |
| | | ADAM | 0.155 | 0.818; 0.292; 2.275 | 0.926 | 0.841; 0.211; 0.907 | 0.919 | 0.83; 0.1; 0.896 |
| | | SGD | 0.165 | 2.287; 0.343; 4.221 | **0.93** | 0.872; 0.794; 0.915 | **0.921** | 0.862; 0.784; 0.906 |
| | | PAL | **0.137** | **0.244**; 0.186; 0.388 | 0.927 | **0.912**; 0.906; 0.921 | 0.916 | **0.902**; 0.889; 0.908 |
| CIFAR-10 | MobileNetV2 | COCOB | 0.232 | **0.282**; 0.257; 0.295 | 0.879 | 0.87; 0.866; 0.876 | 0.865 | 0.852; 0.848; 0.865 |
| | | ALIG | 0.183 | 0.938; 0.347; 1.926 | 0.914 | 0.695; 0.233; 0.888 | 0.897 | 0.528; 0.1; 0.851 |
| | | SGDHD | 0.698 | 2.234; 1.835; 4.366 | 0.886 | 0.75; 0.298; 0.807 | 0.877 | 0.737; 0.295; 0.791 |
| | | SLS | 1.387 | 2.462; 2.011; 2.584 | 0.667 | 0.443; 0.407; 0.504 | 0.595 | 0.4; 0.343; 0.437 |
| | | RMSP | **0.085** | 0.493; 0.337; 0.918 | 0.938 | 0.872; 0.675; 0.895 | 0.929 | 0.865; 0.664; 0.882 |
| | | ADAM | 0.095 | 0.477; 0.314; 1.861 | 0.939 | 0.874; 0.309; 0.896 | 0.93 | 0.864; 0.289; 0.886 |
| | | SGD | 0.149 | 0.878; 0.204; 1.552 | **0.947** | **0.907**; 0.87; 0.933 | **0.94** | **0.899**; 0.859; 0.925 |
| | | PAL | 0.15 | 0.377; 0.205; 0.531 | 0.92 | 0.905; 0.896; 0.91 | 0.905 | 0.886; 0.877; 0.896 |
| CIFAR-10 | DenseNet40 | COCOB | 0.228 | 0.234; 0.23; 0.24 | 0.907 | 0.903; 0.901; 0.904 | 0.894 | 0.889; 0.885; 0.892 |
| | | ALIG | 0.188 | 0.604; 0.227; 2.903 | 0.918 | 0.848; 0.438; 0.902 | 0.902 | 0.784; 0.336; 0.875 |
| | | SGDHD | 1.094 | 2.279; 1.349; 2.908 | 0.775 | 0.341; 0.099; 0.696 | 0.762 | 0.1; 0.1; 0.26 |
| | | SLS | **0.065** | **0.115**; 0.104; 0.189 | 0.91 | 0.904; 0.897; 0.905 | 0.901 | **0.893**; 0.89; 0.897 |
| | | RMSP | 0.147 | 0.398; 0.256; 0.915 | 0.927 | 0.879; 0.737; 0.915 | 0.92 | 0.867; 0.717; 0.909 |
| | | ADAM | 0.138 | 0.749; 0.274; 1.028 | 0.922 | 0.777; 0.611; 0.91 | 0.913 | 0.806; 0.605; 0.907 |
| | | SGD | 0.147 | 0.794; 0.396; 1.746 | **0.932** | 0.855; 0.537; 0.914 | **0.93** | 0.847; 0.528; 0.914 |
| | | PAL | 0.099 | 0.217; 0.165; 0.343 | 0.925 | **0.907**; 0.894; 0.919 | 0.916 | 0.882; 0.861; 0.9 |
| CIFAR-10 | ResNet32 | COCOB | 0.125 | 0.128; 0.127; 0.129 | 0.888 | 0.886; 0.885; 0.887 | 0.878 | 0.872; 0.871; 0.874 |
| | | ALIG | 0.122 | 0.658; 0.279; 1.485 | 0.892 | 0.815; 0.47; 0.881 | 0.866 | 0.71; 0.367; 0.852 |
| | | SGDHD | 0.35 | 0.464; 0.413; 0.701 | 0.864 | 0.835; 0.791; 0.843 | 0.837 | 0.796; 0.766; 0.827 |
| | | SLS | **0.005** | **0.006**; 0.005; 0.827 | 0.871 | 0.856; 0.758; 0.869 | 0.846 | 0.824; 0.657; 0.839 |
| | | RMSP | 0.105 | 0.199; 0.129; 0.498 | 0.922 | 0.884; 0.804; 0.904 | 0.915 | 0.877; 0.792; 0.896 |
| | | ADAM | 0.105 | 0.332; 0.133; 1.004 | 0.917 | 0.875; 0.677; 0.881 | 0.914 | 0.868; 0.654; 0.873 |
| | | SGD | 0.098 | 0.131; 0.118; 0.322 | **0.939** | **0.899**; 0.85; 0.924 | **0.933** | **0.893**; 0.838; 0.92 |
| | | PAL | 0.05 | 0.105; 0.075; 0.195 | 0.921 | 0.893; 0.887; 0.906 | 0.903 | 0.88; 0.849; 0.888 |
| CIFAR-100 | DenseNet40 | COCOB | 0.739 | 0.761; 0.75; 0.772 | 0.642 | 0.633; 0.631; 0.637 | 0.646 | 0.632; 0.629; 0.637 |
| | | ALIG | 0.488 | 2.125; 0.988; 3.128 | 0.637 | 0.508; 0.391; 0.623 | 0.616 | 0.48; 0.264; 0.605 |
| | | SGDHD | 1.78 | 2.6; 2.179; 3.465 | 0.566 | 0.418; 0.274; 0.504 | 0.55 | 0.296; 0.159; 0.497 |
| | | SLS | 1.367 | 1.908; 1.446; 1.96 | 0.719 | 0.593; 0.572; 0.698 | 0.612 | 0.479; 0.422; 0.554 |
| | | RMSP | 0.348 | 1.238; 0.78; 1.972 | 0.716 | 0.583; 0.481; 0.634 | 0.712 | 0.588; 0.482; 0.631 |
| | | ADAM | 0.326 | 1.114; 0.859; 3.53 | 0.715 | 0.601; 0.165; 0.637 | 0.712 | 0.599; 0.226; 0.641 |
| | | SGD | 0.376 | 0.713; 0.431; 2.154 | **0.75** | 0.633; 0.489; 0.709 | **0.753** | 0.634; 0.498; 0.708 |
| | | PAL | **0.275** | **0.376**; 0.312; 0.459 | 0.73 | **0.686**; 0.66; 0.705 | 0.717 | **0.676**; 0.642; 0.695 |
| CIFAR-100 | EfficientNet | COCOB | 0.802 | 0.817; 0.807; 0.822 | 0.594 | 0.583; 0.581; 0.59 | 0.596 | 0.582; 0.58; 0.588 |
| | | ALIG | 0.57 | 2.4; 0.995; 4.085 | 0.612 | 0.494; 0.169; 0.6 | 0.599 | 0.458; 0.115; 0.597 |
| | | SGDHD | 3.545 | 6.528; 5.519; 8.917 | 0.529 | 0.337; 0.178; 0.463 | 0.513 | 0.342; 0.179; 0.468 |
| | | SLS | 3.731 | 6.713; 6.348; 6.857 | 0.474 | 0.212; 0.208; 0.227 | 0.375 | 0.203; 0.149; 0.208 |
| | | RMSP | 0.422 | 1.823; 1.253; 2.968 | 0.675 | 0.517; 0.383; 0.588 | 0.678 | 0.521; 0.382; 0.59 |
| | | ADAM | 0.45 | 1.394; 1.312; 4.606 | 0.684 | 0.518; 0.025; 0.619 | 0.684 | 0.524; 0.01; 0.621 |
| | | SGD | 0.42 | 2.44; 0.633; 5.214 | **0.712** | 0.579; 0.473; 0.661 | **0.709** | 0.579; 0.476; 0.658 |
| | | PAL | **0.372** | **0.471**; 0.409; 0.772 | 0.693 | **0.666**; 0.638; 0.676 | 0.69 | **0.664**; 0.63; 0.671 |
| CIFAR-100 | MobileNetV2 | COCOB | 0.486 | **0.513**; 0.492; 0.536 | 0.644 | 0.63; 0.626; 0.637 | 0.644 | 0.63; 0.623; 0.638 |
| | | ALIG | 0.323 | 2.396; 0.817; 4.247 | 0.661 | 0.41; 0.034; 0.623 | 0.652 | 0.229; 0.01; 0.602 |
| | | SGDHD | 1.485 | 3.307; 2.425; 7.002 | 0.593 | 0.476; 0.39; 0.545 | 0.589 | 0.456; 0.385; 0.525 |
| | | SLS | 3.857 | 5.086; 5.031; 5.64 | 0.332 | 0.2; 0.099; 0.203 | 0.197 | 0.081; 0.052; 0.126 |
| | | RMSP | 0.198 | 1.518; 0.718; 3.368 | 0.728 | 0.593; 0.43; 0.635 | 0.727 | 0.593; 0.431; 0.634 |
| | | ADAM | 0.218 | 1.873; 0.776; 4.524 | 0.729 | 0.528; 0.025; 0.593 | 0.729 | 0.533; 0.02; 0.595 |
| | | SGD | 0.4 | 0.974; 0.473; 2.151 | **0.733** | 0.657; 0.57; 0.7 | **0.736** | 0.659; 0.573; 0.701 |
| | | PAL | **0.181** | 0.602; 0.314; 1.571 | 0.726 | **0.666**; 0.574; 0.689 | 0.722 | **0.664**; 0.509; 0.681 |
| CIFAR-100 | ResNet32 | COCOB | 0.498 | 0.569; 0.524; 0.673 | 0.609 | 0.608; 0.607; 0.608 | 0.605 | 0.602; 0.599; 0.604 |
| | | ALIG | 0.537 | 1.932; 0.995; 3.572 | 0.597 | 0.491; 0.19; 0.58 | 0.587 | 0.414; 0.144; 0.549 |
| | | SGDHD | 0.881 | 1.359; 1.06; 1.772 | 0.601 | 0.539; 0.472; 0.586 | 0.599 | 0.517; 0.431; 0.571 |
| | | SLS | 2.62 | 2.808; 2.78; 2.82 | 0.399 | 0.388; 0.384; 0.392 | 0.363 | 0.305; 0.274; 0.33 |
| | | RMSP | 0.519 | 1.019; 0.807; 2.083 | 0.661 | 0.599; 0.455; 0.651 | 0.656 | 0.603; 0.455; 0.65 |
| | | ADAM | 0.402 | 1.772; 0.768; 3.038 | 0.659 | 0.513; 0.262; 0.564 | 0.658 | 0.519; 0.255; 0.567 |
| | | SGD | 0.375 | **0.474**; 0.4; 1.522 | **0.697** | 0.614; 0.494; 0.672 | **0.694** | 0.616; 0.502; 0.667 |
| | | PAL | **0.339** | 0.485; 0.369; 1.424 | 0.662 | **0.636**; 0.546; 0.652 | 0.663 | **0.621**; 0.512; 0.647 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| TOLSTOI | RNN | COCOB | 1.506 | 1.56; 1.533; 1.593 | 0.589 | 0.58; 0.573; 0.584 | 0.582 | 0.572; 0.566; 0.577 |
| | | ALIG | 1.501 | 1.562; 1.528; 1.766 | 0.591 | 0.579; 0.523; 0.586 | 0.584 | 0.571; 0.513; 0.577 |
| | | SGDHD | 2.282 | 2.433; 2.379; 2.445 | 0.375 | 0.338; 0.336; 0.348 | 0.369 | 0.334; 0.332; 0.344 |
| | | SLS | 3.128 | 3.149; 3.136; 3.156 | 0.169 | 0.159; 0.158; 0.165 | 0.168 | 0.158; 0.157; 0.164 |
| | | RMSP | **1.475** | **1.509**; 1.492; 1.556 | **0.599** | **0.591**; 0.579; 0.595 | **0.592** | **0.583**; 0.572; 0.587 |
| | | ADAM | 1.516 | 1.655; 1.596; 1.681 | 0.588 | 0.567; 0.55; 0.578 | 0.581 | 0.561; 0.543; 0.571 |
| | | SGD | 1.496 | 1.872; 1.56; 2.675 | 0.594 | 0.483; 0.278; 0.573 | 0.587 | 0.476; 0.275; 0.566 |
| | | PAL | 1.528 | 1.569; 1.547; 1.588 | 0.587 | 0.581; 0.577; 0.586 | 0.579 | 0.571; 0.556; 0.575 |
| ImageNet | ResNet50 | COCOB | 2.0 | – | 0.51 | – | 0.518 | – |
| | | ALIG | 1.854 | – | 0.539 | – | 0.512 | – |
| | | SGDHD | 2.742 | – | 0.498 | – | 0.495 | – |
| | | RMSP | 9.485 | – | 0.286 | – | 0.28 | – |
| | | ADAM | 1.863 | – | 0.562 | – | 0.559 | – |
| | | SLS | 3.808 | – | 0.286 | – | 0.069 | – |
| | | SGD | 1.123 | – | **0.656** | – | **0.65** | – |
| | | PAL | **0.773** | – | 0.608 | – | 0.608 | – |
| ImageNet | DenseNet121 | COCOB | 6.9 | – | 0.006 | – | 0.006 | – |
| | | ALIG | 2.142 | – | 0.533 | – | 0.512 | – |
| | | SGDHD | 2.939 | – | 0.343 | – | 0.362 | – |
| | | RMSP | 6.901 | – | 0.0 | – | 0.0 | – |
| | | ADAM | 6.901 | – | 0.001 | – | 0.0 | – |
| | | SLS | 7.768 | – | 0.001 | – | 0.001 | – |
| | | SGD | 3.308 | – | 0.458 | – | 0.452 | – |
| | | PAL | **1.228** | – | **0.617** | – | **0.611** | – |

# E  Binary Line Search

The optimal binary line search we compared *PAL* against. Since the line decreases in negative gradient direction, at first a extrapolation phase performs as many steps forward as the loss does not increase. Afterwards a binary search is performed. This approach is valid if the underlying line is convex. For simple readability we chose Python 3.6 syntax.

```
Input: max_num_of_search_steps                                          1
def binary_line_search(last_loss, step, counter, is_extrapolate):       2
    if counter == max_num_of_search_steps:                              3
        return last_loss                                                4
    counter += 1                                                        5
    if is_extrapolate:                                                  6
        current_loss = do_step_on_line(step)                            7
        if current_loss < last_loss:                                    8
            return binary_line_search(current_loss, step,counter,       9
                is_extrapolate)
        else:                                                           10
            is_extrapolate = False                                      11
            do_step_on_line(-step,get_loss=False)                       12
    if not is_extrapolate:                                              13
        loss_right = do_step_on_line(0.5*step, True)                    14
        if loss_right < last_loss:                                      15
            return binary_line_search(loss_right,                       16
                0.5*step,counter, is_extrapolate)
        loss_left = do_step_on_line(-1*step, True)                      17
        if loss_left < last_loss:                                       18
            return binary_line_search(loss_left,                        19
                0.5*step,counter, is_extrapolate)
        do_step_on_line(0.5*step,get_loss=False)                        20
        if loss_right >= last_loss and loss_left >= last_loss:          21
            return binary_line_search(loss_left, 0.5*step,              22
                counter, is_extrapolate)
        else:                                                           23
                # this state is not possible                            24
```