

Should I try multiple optimizers when fine-tuning pre-trained Transformers for NLP tasks? Should I tune their hyperparameters?

Nefeli Gkouti^{1, 2, 3} Prodromos Malakasiotis^{1, 4}
 Stavros Toumpis¹ Ion Androutsopoulos^{1, 3}

¹ Department of Informatics, Athens University of Economics and Business, Greece

² Department of Informatics and Telecommunications,
 National and Kapodistrian University of Athens, Greece

³ Archimedes/Athena RC, Greece

⁴ Workable

nefeli.gkouti@athenarc.gr, {rulller, toumpis, ion}@aueb.gr

Abstract

NLP research has explored different neural model architectures and sizes, datasets, training objectives, and transfer learning techniques. However, the choice of optimizer during training has not been explored as extensively. Typically, some variant of Stochastic Gradient Descent (SGD) is employed, selected among numerous variants, using unclear criteria, often with minimal or no tuning of the optimizer’s hyperparameters. Experimenting with five GLUE datasets, two models (DistilBERT and DistilRoBERTa), and seven popular optimizers (SGD, SGD with Momentum, Adam, AdaMax, Nadam, AdamW, and AdaBound), we find that when the hyperparameters of the optimizers are tuned, there is no substantial difference in test performance across the five more elaborate (adaptive) optimizers, despite differences in training loss. Furthermore, tuning just the learning rate is in most cases as good as tuning all the hyperparameters. Hence, we recommend picking any of the best-behaved adaptive optimizers (e.g., Adam) and tuning only its learning rate. When no hyperparameter can be tuned, SGD with Momentum is the best choice.

and [Monro, 1951](#); [Kiefer and Wolfowitz, 1952](#)).¹ Optimizer selection seems to be based on unclear criteria and anecdotal information. Furthermore, most optimizers have several hyperparameters, often minimally tuned (e.g., only the learning rate is tuned) or left to their default values. Hence, when models need to be trained (e.g., pre-trained or fine-tuned), it is unclear if the available computing resources should be used to try multiple optimizers, tune their hyperparameters, both, or none.

Our work is inspired by [Schmidt et al. \(2021\)](#), who experimented with 15 optimizers and 8 tasks from DeepOBS ([Schneider et al., 2019](#)). Their most striking finding was that trying several optimizers with default hyperparameters was almost as beneficial as (and cheaper than) picking any single (competent) optimizer and tuning its hyperparameters. Hence, practitioners would be advised to try multiple optimizers with defaults, rather than selecting a single optimizer (e.g., based on anecdotal evidence) and tuning its hyperparameters, when they cannot tune both the choice of optimizer and hyperparameters (which is expensive). In fact, tuning the hyperparameters of a single optimizer was only slightly better than using its defaults, which also advocates against hyperparameter tuning. However, [Schmidt et al. \(2021\)](#) considered only one NLP task (character-level language modeling) with an RNN, acknowledging that their findings may not hold with more complicated models, such as Transformers. They also found indications that the best optimizer may depend on the model and task.

We complement the work of [Schmidt et al. \(2021\)](#) from an NLP perspective by investigating empirically if it is worth (a) trying multiple optimiz-

1 Introduction

NLP research has investigated how different neural model architectures (e.g., RNNs, CNNs, Transformers), model sizes, datasets, training (or pre-training) objectives, and transfer-learning techniques (e.g., pre-training and fine-tuning) affect performance and efficiency. However, the effects of using different optimizers to minimize the training loss have not been explored as extensively. Adam ([Kingma and Ba, 2014](#)) is a popular choice, but there are dozens of alternatives, mostly variants of Stochastic Gradient Descent (SGD) ([Robbins](#)

¹ [Schmidt et al. \(2021\)](#) list more than a hundred optimizers that have been used in deep learning.

ers and/or (b) tuning their hyperparameters (and which ones), when fine-tuning pre-trained Transformer encoders. We experiment with five tasks from GLUE (Wang et al., 2018), using seven popular optimizers, namely SGD, SGD with Momentum (SGDM) (Qian, 1999), Adam and AdaMax (Kingma and Ba, 2014), Nadam (Dozat, 2016), AdamW (Loshchilov and Hutter, 2019), and AdaBound (Luo et al., 2019). For each task and optimizer, we fine-tune DistilBERT (Sanh et al., 2019) and DistilRoBERTa,² two efficient pre-trained Transformers that allow us to complete the many experiments of this work with our limited budget. We consider three cases: using the default hyperparameters of the optimizers, tuning all their hyperparameters, or tuning only their learning rates.

With the exception of the two non-adaptive optimizers considered, i.e., plain SGD and SGDM, which are largely unaffected by hyperparameter tuning, the test performance of the other five (adaptive) optimizers improves substantially when their hyperparameters are tuned, unlike smaller overall improvements reported by Schmidt et al. (2021). Interestingly, in most cases tuning only the learning rate is as good as (and much cheaper than) tuning all the hyperparameters. Furthermore, when hyperparameters (or just the learning rate) are tuned, the adaptive optimizers have very similar test scores in most cases, unlike plain SGD and SGDM, which are clearly the worst and second worst, respectively. This parity of test performance of the adaptive optimizers is obtained despite occasional differences in the training loss they reach. When no hyperparameter can be tuned (e.g., due to limited budget), SGDM is the best choice and AdaBound occasionally works relatively well, but the other optimizers considered are much worse. Trying multiple optimizers with defaults (Schmidt et al., 2021) is reasonably good too, because of the good untuned performance of SGDM and AdaBound. However, our experiments suggest that picking just one of the best-behaved adaptive optimizers i.e., with consistent top-performance across datasets, e.g., Adam, and tuning only its learning rate is the best strategy.

2 Optimizers Considered

All the optimizers we consider aim at tuning the weights vector θ of a model, so that a loss function $f(\theta) = \frac{1}{K} \sum_{k=1}^K f_k(\theta)$ is minimized, where $f_k(\theta)$ is the loss of the k -th training example, and K the

Algorithm 1 Gradient Descent (GD)

Stochastic GD (SGD) SGD with Momentum (SGDM)

1: **Input:**

- initial time step $t \leftarrow 0$; initial weight vector θ_0
- learning rate $\epsilon > 0$
- momentum parameter $\alpha \in [0, 1)$
- initial velocity $v_0 \leftarrow 0$

2: **while** stopping criterion not met **do**

3: select all examples ($m \leftarrow K$)

 sample mini-batch of $m \ll K$ examples

 sample mini-batch of $m \ll K$ examples

4: $g_t \leftarrow \frac{1}{m} \sum_{k=1}^m \nabla f_k(\theta_t)$

5: $\theta_{t+1} \leftarrow \theta_t - \epsilon g_t + \alpha v_t$

6: $v_{t+1} \leftarrow \alpha v_t - \epsilon g_t$; $t \leftarrow t + 1$

size of the training set. The gradient of $f(\theta)$,

$$g(\theta) \triangleq \nabla f(\theta) = \frac{1}{K} \sum_{k=1}^K \nabla f_k(\theta), \quad (1)$$

is of special interest as it points to the direction along which $f(\theta)$ increases the fastest. The optimizers we consider are iterative, i.e., they create a sequence of points $\{\theta_t, t = 0, 1, \dots\}$, such that each θ_{t+1} is selected by taking a step away from the previous point θ_t in an attempt to decrease $f(\theta)$. The step could be towards the opposite direction of $g(\theta)$. However, exactly computing $g(\theta)$ at each step is too costly for large training set sizes K .

2.1 Non-adaptive Optimizers

The simplest optimizer is pure **Gradient Descent (GD)** (Algorithm 1). At each iteration the gradient $g(\theta)$ is computed exactly (line 4, as in Eq. 1), then the next point θ_{t+1} is selected to be towards the opposite direction. The **learning rate** $\epsilon > 0$ is a hyperparameter that affects the size of the steps. Computing the exact gradient at each step, however, is too costly for large training sets (i.e., large K).

Stochastic Gradient Descent (SGD) (Algorithm 1) *estimates* the gradient at each step, by using a sample (mini-batch) of $m \ll K$ examples (line 4); the next point is selected as in GD (line 5).

Stochastic Gradient Descent with Momentum (SGDM) (Polyak, 1964) (Algorithm 1) aims to accelerate learning by suppressing oscillations in the created sequence of points. It maintains an exponentially weighted moving average of past gradients, termed the **velocity**, v_t (line 6). The direction

²<https://huggingface.co/distilroberta-base>

and size of the next step (line 5) are now determined by a linear combination of the latest gradient estimate \mathbf{g}_t and the velocity (\mathbf{v}_t). Intuitively, the sequence of points $\{\boldsymbol{\theta}_t, t = 0, 1, \dots\}$ resembles the movement of a ball traveling (in the space of weight vectors) towards points of lower altitude (loss), but also subject to its own momentum.

The optimizers above are called non-adaptive, because their learning rate ϵ is fixed. The optimizers we discuss next modify ϵ while creating the sequence of points and are, hence, called adaptive.

2.2 Adaptive Optimizers

The loss function $f(\boldsymbol{\theta})$ of large neural models corresponds to a complicated hyper-surface. In some directions, the loss may change rapidly, in other directions slowly. This makes choosing the learning rate ϵ difficult: if ϵ is too large, the minimum along the chosen direction of a step may be overshoot; if, however, ϵ is too small, progress towards smaller values of the loss function will be slow. This difficulty of choosing a good, uniform, ϵ a priori is one of the reasons adaptive optimizers are so useful.

Adaptive Moments (Adam) (Kingma and Ba, 2014) follows Algorithm 2. Like SGDM, it uses the concept of momentum by maintaining a velocity-type vector \mathbf{s}_t (line 6), termed the **first moment**. It also adapts the learning rate by scaling it (line 10) roughly inversely proportionally to the square root of an exponentially weighed average \mathbf{r}_t (line 7) of the component-wise squared gradient estimate, \mathbf{g}_t^2 , termed the **second moment**. Lines 8 and 9 normalize the two moments to take into account biases due to their initial values $\mathbf{0}$ (Kingma and Ba, 2014).

Nadam (Dozat, 2016) is identical to Adam, except that it uses Nesterov momentum (Nesterov, 1983), which has been shown to be somewhat superior to plain momentum (Sutskever et al., 2013). The only difference is in lines 8, 9 in Algorithm 2, which ensure that the momentum step incorporates the estimated gradient at the location where it is used.

AdamW (Loshchilov and Hutter, 2019) (Algorithm 2) is based on the empirical observation that smaller weights tend to overfit less. Hence, it adds (in line 10) a term $-\lambda\boldsymbol{\theta}_{t-1}$, where $\lambda \in (0, 1)$, to decay the weight vector towards the origin.

AdaMax (Kingma and Ba, 2014) is identical to Adam, except that the two moments \mathbf{s}_t and \mathbf{r}_t are not normalized, and the second moment (\mathbf{r}_t) is computed using line 7 in Algorithm 3, thus it is

Algorithm 2 Adaptive Moment Optimization (Adam)

Nesterov-Accelerated Adam (Nadam)

Adam with decoupled weight decay (AdamW)

(all vector operations are elementwise)

1: **Input:**

- initial time step $t \leftarrow 0$; initial weight vector $\boldsymbol{\theta}_0$
- learning rate $\epsilon > 0$; decay rates $\rho_1, \rho_2 \in [0, 1)$
- small constant $\delta > 0$
- first moment $\mathbf{s}_0 \leftarrow \mathbf{0}$; second moment $\mathbf{r}_0 \leftarrow \mathbf{0}$
- regularization factor $\lambda \in (0, 1)$

2: **while** stopping criterion not met **do**

3: $t \leftarrow t + 1$

4: sample mini-batch of m examples

5: $\mathbf{g}_t \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla f_i(\boldsymbol{\theta}_{t-1})$

6: $\mathbf{s}_t \leftarrow \rho_1 \mathbf{s}_{t-1} + (1 - \rho_1) \mathbf{g}_t$

7: $\mathbf{r}_t \leftarrow \rho_2 \mathbf{r}_{t-1} + (1 - \rho_2) \mathbf{g}_t^2$

8: $\hat{\mathbf{s}}_t \leftarrow \frac{\mathbf{s}_t}{1 - \rho_1^t} \quad \frac{\mathbf{s}_t}{1 - \rho_1^t} \quad \frac{\rho_1 \mathbf{s}_t}{1 - \rho_1^{t+1}} + \frac{(1 - \rho_1) \mathbf{g}_t}{1 - \rho_1^t}$

9: $\hat{\mathbf{r}}_t \leftarrow \frac{\rho_2}{1 - \rho_2^t} \frac{\mathbf{r}_t}{1 - \rho_2^t}$

10: $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \epsilon \frac{\hat{\mathbf{s}}_t}{\delta + \sqrt{\hat{\mathbf{r}}_t}} \quad -\lambda \boldsymbol{\theta}_{t-1}$

no longer an exponentially weighed average. This variation of Adam was proposed in the same paper that introduced Adam, as a more stable variant.

AdaBound (Luo et al., 2019) (Algorithm 3) ensures that extreme values for the learning rate are avoided, by bounding it by both a dynamic upper bound η_t^u and a dynamic lower bound η_t^l that start from infinity and zero, respectively, and then converge to a finite common value ϵ as the time step t increases. The operation $[x]_l^u$ in line 8 clips the vector x elementwise so that the output lies in the interval $[l, u]$. AdaBound initially behaves like Adam, and gradually transforms to SGD.

The optimizers we consider include simple non-adaptive baselines (SGD, SGDM), the most commonly used adaptive optimizer (Adam), its sibling AdaMax, as well as adaptive optimizers that incorporate influential ideas, in particular Nesterov momentum (Nadam), weight decay (AdamW), and dynamic bounds (AdaBound). AdamW and AdaBound are also two of the most recent optimizers.

3 Experiments

3.1 Datasets and Evaluation Measures

We experiment with five GLUE tasks (Wang et al., 2018).³ The datasets of all tasks are in English. Each experiment is repeated with five random training/development/test splits, and we report average

³To speed up our experiments, we do not use the other four GLUE tasks (QQP, QNLI, RTE, WNLI), which are all textual inference/paraphrasing tasks, like MRPC and MNLI.

Algorithm 3 AdaMax AdaBound
 (all vector operations are elementwise)

1: **Input:**

- initial time step $t \leftarrow 0$; initial weight vector θ_0
- learning rate $\epsilon > 0$; decay rates $\rho_1, \rho_2 \in [0, 1)$
- small constant $\delta > 0$
- first moment $s_0 \leftarrow \mathbf{0}$; second moment $r_0 \leftarrow \mathbf{0}$
- lower bound function η_t^l
- upper bound function η_t^u

2: **while** stopping criterion not met **do**

3: $t \leftarrow t + 1$

4: sample mini-batch of m examples

5: $\mathbf{g}_t \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla f_i(\theta_{t-1})$

6: $\mathbf{s}_t \leftarrow \rho_1 \mathbf{s}_{t-1} + (1 - \rho_1) \mathbf{g}_t$

7: $\mathbf{r}_t \leftarrow \max\{\rho_2 \mathbf{r}_{t-1}, |\mathbf{g}_t|\} \cdot \rho_2 \mathbf{r}_{t-1} + (1 - \rho_2) \mathbf{g}_t^2$

8: $\eta_t \leftarrow \left[\frac{\epsilon}{\sqrt{r}} \right] \eta_t^u$

9: $\theta_t \leftarrow \theta_{t-1} - \frac{\epsilon}{1 - \rho_1^t} \cdot \frac{\mathbf{s}_t}{\mathbf{r}_t} \cdot \theta_{t-1} - \eta_t \cdot \mathbf{s}_t$

scores and standard deviations over the repetitions.

SST-2 (Socher et al., 2013) is a binary sentiment classification dataset with 68.8k sentences from movie reviews (one label per sentence). To speed up our experiments, we sample 18k (from the 68.8k) sentences anew in each of the five repetitions and split them into training (15k), development (1.5k) and test (1.5k) subsets. The class distribution of the 68.8k sentences (55% positive sentiment) is preserved in all subsets of every split.

MRPC (Dolan and Brockett, 2005) contains 5.8k sentence pairs from online news. Each pair is classified as containing paraphrases (sentences with the same meaning) or not. We use 80% of the 5.8k pairs for training, 10% for development, 10% for testing, preserving the class distribution (67% paraphrases).

CoLA (Warstadt et al., 2019) contains 9.6k word sequences labeled to indicate if each sequence is a grammatically correct sentence or not. We use 80% of the sequences for training, 10% for development and 10% for testing, preserving the class distribution (70% acceptable).

STS-B (Cer et al., 2017) contains 7.2k sentence pairs from news headlines, video/image captions, and natural language inference data. Each pair is annotated with a similarity score from 1 to 5. In each repetition, we sample 80% of the 7.2k pairs for training, 10% for development, 10% for testing.

MNLI (Williams et al., 2018) contains 393k premise-hypothesis pairs for training, and 19.6k

pairs for development. The task is to predict if the premise entails, contradicts, or is neutral to the hypothesis. To speed up the experiments, in each repetition we sample (anew) 50k from the 393k pairs for training, 9.8k from the 19.6k for development, and the remaining 9.8k for testing, preserving the original class distribution (balanced).⁴

Evaluation measures: We use the measures adopted by GLUE (Wang et al., 2018), i.e., Accuracy for SST-2 and MNLI, Macro-F1 for MRPC, Matthews correlation (Matthews, 1975) for CoLA, and Pearson correlation (Kirch, 2008) for STS-B.

3.2 Experimental Setup

Transformer models: Given the volume of the experiments and our limited resources, we fine-tune: (i) DistilBERT (Sanh et al., 2019), a distilled BERT-base (Devlin et al., 2019) with 40% fewer parameters that runs 60% faster, but retains 95% of BERT-base’s performance on GLUE, according to its creators; and (ii) DistilRoBERTa, a similarly distilled version of RoBERTa-base (Liu et al., 2019).

Hyperparameter tuning: For each optimizer, model, task, and data split (Section 3.1), we try 30 different combinations (trials) of hyperparameter values, as selected by Optuna (Akiba et al., 2019), seeking to maximize the evaluation measure of the task on the development subset.⁵ In each trial, we retain the weights from the epoch with the best development score. The hyperparameter search space of each optimizer includes the default values proposed by its creators, with the exception of the learning rate ϵ of adaptive optimizers, since it is standard practice when fine-tuning Transformers with adaptive optimizers to use much smaller ϵ .⁶ We repeat the experiments, tuning only ϵ . We also report results with default hyperparameter values.

Loss functions: We minimize cross-entropy for the classification tasks (SST-2, MRPC, CoLA, and MNLI), and mean squared error for STS-B.

3.3 Experimental Results

We include in the main paper only results using DistilBERT. DistilRoBERTa results are reported in Appendix B, and lead to the same conclusions.

⁴We also ensure that development and test sets contain 50% ‘in-domain’ (matched) pairs and 50% ‘out-of-domain’.

⁵In Optuna, we use Tree-Structured Parzen Estimation (Bergstra et al., 2011, 2013) and median-based pruning.

⁶More details on the hyperparameter search space and the selected values are provided in Appendix A.

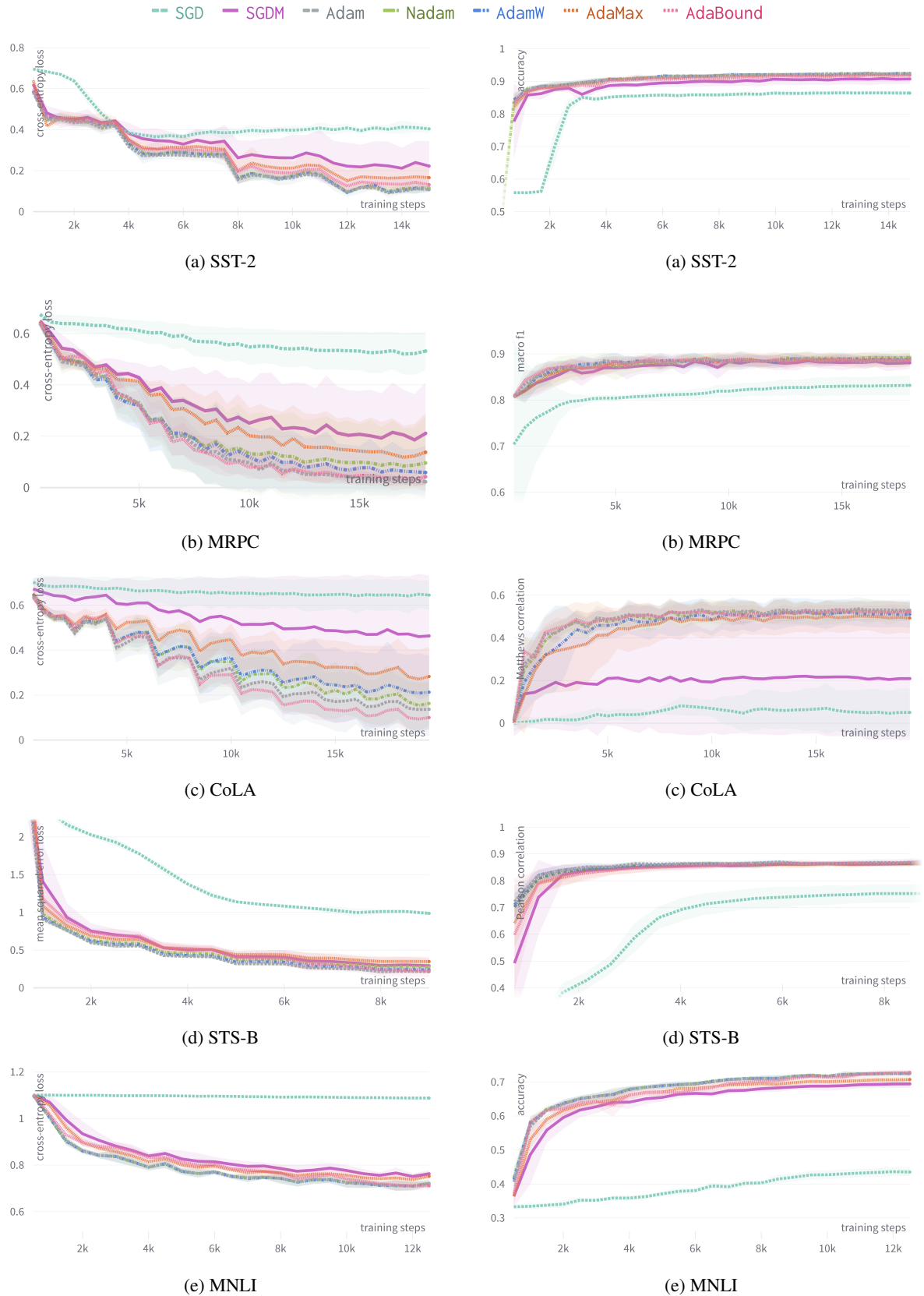


Figure 1: **Training loss (left) and evaluation score on development data (right) with all hyperparameters of the optimizers tuned, as a function of training steps, using DistilBERT.** For each dataset, we use **five random data splits**, and plot the **average and standard deviation (shadow)** over the five splits. **Plain SGD** is clearly the **worst**, but adding **Momentum (SGDM)** turns it to a **competent** optimizer in terms of development scores, except for CoLA. The five **adaptive optimizers** (Adam, Nadam, AdamW, AdaMax, AdaBound) have **almost identical** development score curves across the tasks, **despite occasional differences in the training losses** they reach.

Optimizer	SST-2 Accuracy	MRPC Macro-F1	CoLA Matthews	STS-B Pearson	MNLI Accuracy
AdaBound	91.61 (1.05)	81.34 (1.17)	0.53 (0.03)	0.87 (0.01)	72.33 (0.35)
AdamW	91.93 (0.73)	81.01 (1.09)	0.51 (0.05)	0.87 (0.01)	72.17 (0.38)
AdaMax	91.73 (1.10)	80.88 (0.68)	0.51 (0.07)	0.86 (0.01)	70.44 (0.70)
Nadam	91.76 (0.97)	81.85 (3.75)	0.53 (0.03)	0.86 (0.01)	72.04 (0.36)
Adam	91.63 (0.75)	80.81 (1.83)	0.52 (0.03)	0.87 (0.01)	72.08 (0.37)
SGDM	90.53 (1.78)	79.51 (1.58)	0.22 (0.26)	0.87 (0.01)	68.84 (1.62)
SGD	86.17 (1.16)	65.45 (9.94)	0.09 (0.14)	0.74 (0.03)	43.64 (0.65)

Table 1: Evaluation scores on **test data** with **all hyperparameters** of the optimizers **tuned**, using **DistilBERT**. For each dataset, we use **five random splits** of training/development/test data (the same for all optimizers) and report the **average** test score and the **standard deviation** over the five splits. **Plain SGD** is clearly the **worst**, but **SGDM** is **competitive**, except for CoLA. The **five adaptive optimizers** (top five) **all perform very similarly**.

Optimizer	SST-2 Accuracy	MRPC Macro-F1	CoLA Matthews	STS-B Pearson	MNLI Accuracy
AdaBound	91.40 (1.16)	81.59 (1.98)	0.49 (0.06)	0.87 (0.01)	72.74 (0.85)
AdamW	91.87 (0.88)	81.32 (1.75)	0.54 (0.03)	0.87 (0.01)	72.20 (0.28)
AdaMax	89.52 (1.16)	81.12 (0.59)	0.48 (0.07)	0.84 (0.01)	66.80 (0.28)
Nadam	92.79 (0.37)	80.91 (1.69)	0.55 (0.04)	0.87 (0.01)	72.36 (0.37)
Adam	91.75 (0.91)	80.81 (1.83)	0.52 (0.02)	0.86 (0.01)	72.19 (0.27)
SGDM	89.79 (1.23)	80.79 (0.79)	0.46 (0.10)	0.85 (0.01)	67.31 (0.61)
SGD	86.17 (1.16)	65.45 (9.94)	0.09 (0.14)	0.74 (0.03)	43.64 (0.65)

Table 2: Evaluation scores on **test data**, having **tuned only the learning rate**, with all other hyperparameters of the optimizers set to their defaults, using **DistilBERT**. Again, we use five random splits and report the average and standard deviation. **In most cases, tuning only the learning rate leads to very similar results as tuning all the hyperparameters** (cf. Table 1). Exceptions include AdaMax, which now lags behind on CoLA and (more noticeably) on MNLI, as well as AdaBound, whose performance deteriorates on CoLA (comparing to Table 1). Interestingly, SGDM is now competitive to the five adaptive optimizers on CoLA (where it lagged behind in Table 1).

Figure 1 shows the training loss for each task and optimizer (left column) and the corresponding evaluation score on development data (right column), as a function of training steps, using DistilBERT, when all the hyperparameters of the optimizers are tuned. For each curve, we plot the average and standard deviation (shadow) over the five data splits (Section 3.1). SGD clearly struggles to learn the training data in all five tasks (left), which is also reflected in its development scores (right). However, adding Momentum turns it to a competent optimizer (SGDM) in terms of development scores. An exception is CoLA, where SGDM is clearly worse in development score than the five more elaborate (adaptive) optimizers (Adam, Nadam, AdamW, AdaMax, and AdaBound), in accordance with its poor training loss, though it still outperforms SGD. The adaptive optimizers have almost identical development score curves across the tasks, with some minor differences in CoLA where AdaMax and AdamW are slightly worse, again reflecting their inferior training losses. Otherwise, differences in the training losses reached by the five adaptive optimizers (when there are any) do not lead to visible

differences in development scores.⁷

Table 1 tells a similar story, now evaluating on test data, again using DistilBERT. Again, SGD is clearly the worst, but SGDM is competitive, except for CoLA. The five adaptive optimizers all perform very similarly. Interestingly, in most cases tuning only the learning rate (Table 2) leads to very similar results as (and is much cheaper than) tuning all the hyperparameters (Table 1). Some exceptions are reported in the caption of Table 2.

Table 3 shows test results with default hyperparameter values, again using DistilBERT. SGDM is not affected by the lack of hyperparameter tuning (cf. Table 2) and is now the best overall. Plain SGD is also not particularly affected, and actually performs much better on MRPC and CoLA untuned (cf. Table 2). Overall, it seems that the defaults of the two non-adaptive optimizers are good global choices; tuning their hyperparameters occasionally overfits the development data. By contrast, the adaptive optimizers are negatively affected by the lack of tuning. AdaBound is the least affected and is now (Table 3) overall the second best and clearly

⁷Curves for the cases where we tune only the learning rate or use the defaults can be found in Appendix B.

Optimizer	SST-2 Accuracy	MRPC Macro-F1	CoLA Matthews	STS-B Pearson	MNLI Accuracy
AdaBound	90.88 (1.41)	76.14 (2.63)	0.19 (0.26)	0.86 (0.01)	65.58 (8.48)
AdamW	55.83 (0.06)	62.57 (0.00)	0.00 (0.00)	0.23 (0.22)	35.34 (0.01)
AdaMax	59.49 (7.35)	62.57 (0.00)	0.00 (0.00)	0.50 (0.15)	35.34 (0.01)
Nadam	55.80 (0.00)	62.57 (0.00)	0.00 (0.00)	0.14 (0.11)	35.34 (0.01)
Adam	55.83 (0.06)	62.57 (0.00)	0.00 (0.00)	0.27 (0.17)	35.34 (0.01)
SGDM	89.89 (1.13)	80.00 (0.51)	0.49 (0.03)	0.86 (0.01)	67.56 (0.24)
SGD	86.35 (0.96)	70.90 (1.59)	0.33 (0.04)	0.74 (0.03)	44.65 (0.20)

Table 3: Evaluation scores on **test data**, with **all hyperparameters** of the optimizers set to **defaults**, using **DistilBERT**. Again, we use five random splits and report the average and standard deviation. **SGDM is not affected** by the lack of hyperparameter tuning (cf. Table 2) and is **now the best** overall. Plain SGD is also not particularly affected, and actually performs much better on MRPC and CoLA untuned (cf. Table 2). Although negatively affected by the use of defaults, **AdaBound** is now the best optimizer on SST-2 and STS-B, and much better than the other four adaptive optimizers (but worse than SGDM) on the other datasets.

better than the other adaptive optimizers. This may be related to the fact that AdaBound behaves similarly to SGD at the end of training (Section 2). All seven optimizers have very similar default learning rates (Appendix A), hence the superior out-of-the-box performance of SGDM and AdaBound is not due to different default learning rates.

Trying multiple optimizers with defaults (Schmidt et al., 2021) is competitive too. Based on development scores (Fig. 3), one would select the same optimizers (per task) whose (test) scores are shown in bold in Table 3. The resulting test scores are competitive, but worse than the best scores of Tables 1–2. Also, the competitive scores of trying multiple optimizers with defaults are due only to the good out-of-the-box performance of SGDM and (to a lesser extent) AdaBound; the other optimizers perform much worse with defaults, hence trying them would have been a waste of resources.

Therefore, based on our experiments, we recommend picking just one adaptive optimizer and tuning only its learning rate. Among the adaptive optimizers we considered, we recommend picking AdamW, Nadam, or Adam, since AdaBound and AdaMax were not top performers across all datasets when tuning only the learning rate (Tables 2, 14).

4 Related Work

DeepOBS (Schneider et al., 2019) is an optimizer benchmarking suite that includes several classical datasets, models (e.g., CNNs, RNNs), optimizer implementations (currently SGD, SGDM, Adam), and facilities to compare optimizers. However, it includes only one NLP task (character-level language modeling with an RNN) and no Transformer models. A similar observation can be made for the more recent AlgoPerf benchmark (Dahl et al.,

2023), which includes Transformers, but only one NLP task (machine translation). MultiTask (Metz et al., 2020) also considers only three NLP tasks (language modelling with characters or words/subwords, text classification), all with RNNs.

As already noted, we were inspired by Schmidt et al. (2021), who experimented with 15 optimizers and 8 tasks (from DeepOBS), but only one NLP task (the only one of DeepOBS), without considering Transformers. They found that Adam remained a strong contender, with more recent variants failing to consistently outperform it. Tuning the hyperparameters of a single optimizer was overall only slightly better than using its default hyperparameter values (median improvement 3.4% for a tuning budget of 50 trials and diminishing returns for larger budgets). Trying several optimizers with defaults was almost as beneficial as (and cheaper than) picking any single (competent) optimizer and tuning it. Schmidt et al. acknowledged, however, that their findings may not hold with more complicated models, such as Transformers. They also found indications that the best optimizer may depend on the model and task. They employed random search for hyperparameter tuning, whereas we used Optuna (Section 3.2). They also experimented with four update schedules for the learning rate (constant, cosine decay, cosine with warm restarts, and trapezoidal) on top of the tuned rate, with results indicating that non-constant schedules add small gains; we considered only a constant schedule.

Wilson et al. (2017) reported that adaptive optimizers may lead to worse development or test performance than SGD, even in cases where the adaptive optimizers reach lower training losses. They considered artificial datasets, an image classification task (using a CNN), character-level language

modelling, and two parsing tasks, the latter three tasks using LSTM-based models. However, they only tuned the learning rate and the learning rate decay scheme, as pointed out by Choi et al. (2019).

Choi et al. (2019) introduced the notion of *inclusion* between optimizers. For example, SGD is a special case of (is included by) SGDM for $\alpha = 0$ (Algorithm 1). With an exhaustive hyperparameter search, an optimizer should never perform worse than an optimizer it includes. Indeed, Choi et al. show that with extensive hyperparameter tuning, inclusion relationships reflect end-task performance, and they criticize previous work by Wilson et al. (2017) and Schneider et al. (2019) for not having tuned all hyperparameters. We tuned all hyperparameters, but found that tuning only the learning rate was equally good. We also note that extensive tuning of the kind recommended by Choi et al. (e.g., with coarser to finer swaps to explore and define the search space anew per task) is computationally very expensive. Finally, we note that Choi et al. (2019) considered only image classification and language modelling, the latter with an LSTM and a (not pre-trained) Transformer.

Sivaprasad et al. (2020) pointed out that when comparing optimizers, it is important to consider how easy it is to reach reasonable performance with a limited budget (number of trials), rather than comparing performance scores obtained with very extensive (and costly) hyperparameter tuning, unlike the setting of Choi et al. (2019). They experimented with SGDM, Adagrad (Duchi et al., 2011), Adam, and AdamW, with a range of budgets, in nine tasks, of which only two were NLP tasks (sentiment analysis, news classification), without considering Transformers. They recommended using Adam and tuning only its learning rate, especially with low budgets, which agrees with our conclusions. Although we experimented with a single, relatively small budget (30 trials), we used the same budget for all optimizers, like Sivaprasad et al. (2020). Their evaluation protocol, which efficiently simulates multiple budgets, could be used in future extensions of our work, though it is incompatible with our use of Optuna, as it requires random search.

A particularly interesting research direction is to semi-automatically discover new optimization algorithms via evolutionary program searches and manual intervention. Chen et al. (2023) recently used this approach to produce Lion, an optimizer that, among other experiments, was reported to

be overall slightly better than AdamW on GLUE, when using the T5 model (Raffel et al., 2020), but tuning only the learning rates and decoupled weight decay hyperparameters of the two optimizers.

The optimizers we considered are first-order, i.e., they compute only the gradient of the loss function and not its Hessian (second-order partial derivatives), which would be prohibitively costly (M^2 second-order derivatives at each step, for a model with M weights). Interestingly, Liu et al. (2023) investigate an *approximate* second-order optimizer for use in language model pretraining.

Furthermore, the optimizers we considered use a single value of the learning rate at each step. *Backtracking* methods consider multiple learning rate values at each step, computing the loss for each one. Although backtracking is very common in traditional optimization, apparently it has not received sufficient attention in machine learning; an exception is the work of Truong and Nguyen (2021).

5 Conclusions

We investigated if it is worth (a) trying multiple optimizers and/or (b) tuning their hyperparameters (and which ones), when fine-tuning a pre-trained Transformer. We experimented with five GLUE datasets, two efficient pre-trained Transformer encoders (DistilBERT, DistilRoBERTa), and seven popular optimizers (SGD, SGDM, Adam, AdaMax, Nadam, AdamW, and AdaBound). With the exception of the two non-adaptive optimizers (SGD, SGDM), which were largely unaffected by hyperparameter tuning, the test performance of the other five (adaptive) optimizers improved substantially when they were tuned, unlike previously reported smaller overall gains (Schmidt et al., 2021). In most cases, tuning only the learning rate was as good as (and cheaper than) tuning all the hyperparameters. Furthermore, when hyperparameters (or just the learning rate) were tuned, all the adaptive optimizers had very similar test scores, unlike SGD and SGDM, which were clearly the worst and second worst, respectively. This parity of test performance of the adaptive optimizers was obtained despite occasional differences in the training loss they reached. When no hyperparameter was tuned (which might be the case with a low budget), SGDM was the best choice and AdaBound the second best; the other optimizers were much worse. Trying multiple optimizers with defaults (Schmidt et al., 2021) worked reasonably well too,

but only because of the good untuned performance of SGDM and (to a lesser extent) AdaBound; trying the other optimizers untuned would have been a waste of resources. Hence, we suggest picking just one adaptive optimizer and tuning only its learning rate; we recommend AdamW, Nadam, or Adam, which were consistently top performers when tuning only the learning rate.

Our work can help save substantial effort, computational resources, and energy, by reducing the number of hyperparameter tuning experiments practitioners perform. However, our findings need to be complemented by future additional experiments with more models, more NLP tasks (including pre-training tasks), and different tuning budgets (i.e., different maximum number of trials); see also Section 6. This is a computationally very expensive exploration that might be best handled by multiple groups performing and reporting similar studies. Towards this direction, we make all the code, data, and results of our experiments publicly available.⁸

6 Limitations

We examined different optimizers when *fine-tuning* pre-trained Transformers for NLP downstream tasks. Given our limited resources, we did not consider *pre-training* and we experimented only with two lightweight encoder-only models, i.e., DistilBERT and DistilRoBERTa. Thus, we limited ourselves to five *classification* datasets for single word sequences or pairs of sequences (Section 3.1), and did not consider sequence-to-sequence or sequence generation tasks, which would require encoder-decoder or decoder-only models. Also, given our limited resources, we considered seven popular optimizers (Section 2), among many more available.

During hyperparameter tuning, we restricted ourselves to the direct hyperparameters of the optimizers. Thus, we did not consider tuning the batch size, although it often affects the choice of learning rate.⁹ Also, given that the effect of a non-constant learning rate may vary substantially with respect to the optimizer and task (Schmidt et al., 2021), we experimented with a constant (but tuned) learning rate only, leaving the investigation of other update schedules (e.g., cosine decay) for future work.

Finally, we did not measure how the training speed is affected by the choice of optimizer and hyperparameter tuning (Schneider et al., 2019). In

our work, the training speed can only be indirectly inferred from the learning curves (Fig. 1–6) and depends on the size of each dataset. We also did not vary the tuning budget; we used a fixed budget of 30 trials in all experiments, which is close to the ‘small’ budget (25 trials) of Schmidt et al. (2021).

7 Ethical considerations

Selecting the best optimization scheme, which includes the choice of optimizer and tuning its hyperparameters, is much more expensive than training the final model with tuned hyperparameters (Metz et al., 2020). Thus, selecting the best optimization scheme has a large economic and environmental impact, which could be greatly reduced with the help of proper guidelines to narrow the search space without compromising performance. Of course, the effort to develop such guidelines (which includes our work) also requires substantial resources. However, by making code, data, and results publicly available (as we do), we believe that this effort will help save significant resources in the long run.

Acknowledgements

This work was partially supported by project MIS 5154714 of the National Recovery and Resilience Plan Greece 2.0 funded by the European Union under the NextGenerationEU Program. It was also supported by the TPU Research Cloud (TRC) program of Google.¹⁰

References

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, New York, NY, USA. Association for Computing Machinery.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. 2011. *Algorithms for hyper-parameter optimization*. In *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc.
- James Bergstra, Daniel Yamins, and David D. Cox. 2013. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *ICML’13*, page I–115–I–123, Atlanta, GA, USA.

⁸<https://github.com/nlp-aueb/nlp-optimizers>

⁹We used a fixed batch size of 4 in all experiments.

¹⁰<https://sites.research.google/trc/about/>

- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. 2023. [Symbolic discovery of optimization algorithms](#).
- Dami Choi, Christopher J. Shallue, Zachary Nado, Jaehoon Lee, Chris J. Maddison, and George E. Dahl. 2019. On empirical comparisons of optimizers for deep learning. *Computing Research Repository*, abs/1910.05446.
- George E. Dahl, Frank Schneider, Zachary Nado, Naman Agarwal, Chandramouli Shama Sastry, Philipp Hennig, Sourabh Medapati, Runa Eschenhagen, Priya Kasimbeg, Daniel Suo, Juhan Bae, Justin Gilmer, Abel L. Peirson, Bilal Muhammad Khan, Rohan Anil, Michael G. Rabbat, Shankar Krishnan, Daniel Snider, Ehsan Amid, Kongtao Chen, Chris J. Maddison, R. Aravind Vasudev, Michal Badura, Ankush Garg, and Peter Mattson. 2023. Benchmarking neural network training algorithms. *ArXiv*, abs/2306.07179.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- William B. Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the 3d International Workshop on Paraphrasing (IWP2005)*.
- Timothy Dozat. 2016. Incorporating Nesterov momentum into Adam. In *Proceedings of the 4th International Conference on Learning Representations*, San Juan, Puerto Rico.
- John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159.
- Jack Kiefer and Jacob Wolfowitz. 1952. A stochastic approximation method. *The Annals of Mathematical Statistics*, 23(3):462–466.
- Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference for Learning Representations*, San Diego, CA, USA.
- Wilhelm Kirch. 2008. [Pearson’s correlation coefficient](#). In *Encyclopedia of Public Health*, pages 1090–1091, Dordrecht. Springer Netherlands.
- Hong Liu, Zhiyuan Li, David Hall, Percy Liang, and Tengyu Ma. 2023. [Sophia: A scalable stochastic second-order optimizer for language model pre-training](#).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#).
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *Proceedings of the 7th International Conference on Learning Representations*, New Orleans, LA, USA.
- Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. 2019. Adaptive gradient methods with dynamic bound of learning rate. In *Proceedings of the 7th International Conference on Learning Representations*, New Orleans, LA, USA.
- Brian W Matthews. 1975. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451.
- Luke Metz, Niru Maheswaranathan, Ruoxi Sun, C. Daniel Freeman, Ben Poole, and Jascha Sohl-Dickstein. 2020. Using a thousand optimization tasks to learn hyperparameter search strategies. In *NeurIPS MetaLearning Workshop (MetaLearn 2020)*.
- Yurii Nesterov. 1983. A method for unconstrained convex minimization problem with the rate of convergence $o(1/k^2)$. *Proceedings of the USSR Academy of Sciences*, 269:543–547.
- Boris T. Polyak. 1964. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17.
- Ning Qian. 1999. On the momentum term in gradient descent learning algorithms. *Neural networks : the official journal of the International Neural Network Society*, 12(1):145–151.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Herbert Robbins and Sutton Monro. 1951. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. In *5th*

Workshop on Energy Efficient Machine Learning and Cognitive Computing @ NeurIPS 2019.

Robin M Schmidt, Frank Schneider, and Philipp Hennig. 2021. Descending through a crowded valley - benchmarking deep learning optimizers. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *ICML '21*, pages 9367–9376.

Frank Schneider, Lukas Balles, and Philipp Hennig. 2019. DeepOBS: A deep learning optimizer benchmark suite. In *Proceedings of the 7th International Conference on Learning Representations*, New Orleans, LA, USA.

Prabhu Teja Sivaprasad, Florian Mai, Thijs Vogels, Martin Jaggi, and François Fleuret. 2020. Optimizer benchmarking needs to account for hyperparameter tuning. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *ICML '20*, pages 9036–9045.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, WA, USA. Association for Computational Linguistics.

Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. 2013. On the importance of initialization and momentum in deep learning. In *Proceedings of the 30th International Conference on Machine Learning*, ICML '13, pages 1139–1147, Atlanta, GA, USA.

Tuyen Trung Truong and Hang-Tuan Nguyen. 2021. Backtracking gradient descent method and some applications in large scale optimisation. part 2: Algorithms and experiments. *Applied Mathematics & Optimization*, 84(3):2557–2586.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.

Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural network acceptability judgments. *Transactions of the Association for Computational Linguistics*, 7:625–641.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.

Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. 2017. The marginal value of adaptive gradient methods in machine learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, volume 30, Long Beach, CA, USA. Curran Associates, Inc.

Appendix

A Hyperparameter Tuning

Tables 4–12 show the *search space*, *default values*, and *tuned values* (averaged over the five random splits) of the hyperparameters across the optimizers and datasets, when using DistilBERT.¹¹ In most cases, the tuned hyperparameter values are different than the defaults. Regarding the learning rate (ϵ), as already noted in Section 3.2, the search space we use for adaptive optimizers does not include the default values (see Tables 4–6), because it is standard practice when fine-tuning Transformers with adaptive optimizers to use much smaller learning rates. Nevertheless, we observe (Tables 6, 12) that the tuned learning rates of SGDM and SGD are closer to the default, compared to the other optimizers, which may explain why SGDM is the best optimizer overall when using the defaults.

B Additional Results

Figure 2 shows the training loss per task and optimizer (left) and the corresponding evaluation score on development data (right), as a function of training steps, when *only the learning rate* of each optimizer is tuned, using DistilBERT. For each curve, we plot the average and standard deviation (shadow) over the five data splits (Section 3.1). As when tuning all hyperparameters (Fig. 1), SGD struggles to learn the training data in all five tasks (left), which is also reflected in its development scores (right). Again, adding Momentum to SGD turns it to a competent optimizer (SGDM) in terms of development scores, now even on CoLA (cf. Fig. 1). The five adaptive optimizers perform similarly overall in terms of development scores, except for AdaMax, which is visibly worse (along with SGDM) on CoLA and (to a larger extent) MNLI. Differences (when there are any) in the training loss of different optimizers do not always lead to substantial differences in development scores.

Figure 3 shows the corresponding curves when all hyperparameters are set to their *defaults*, again

¹¹The tuned hyperparameter values for DistilRoBERTa will be available in our code repo and lead to similar conclusions.

using DistilBERT. SGDM is not affected by the lack of tuning and is now the best overall, improving upon plain SGD, which is now overall a strong contender in terms of development scores (in agreement with the test scores of Table 3). Although negatively affected by the use of defaults, AdaBound is now overall the second best in terms of development scores (again, as in Table 3); it matches SGDM’s development scores on SST-2 and STS-B, but does not perform as well on the other datasets, although it is still better than the other adaptive optimizers. Again, differences in training loss are not always reflected to differences in development scores. In SST-2 and MRPC, for example, AdaBound reaches a much lower training loss than SGDM, but the development curve of AdaBound is almost identical (in SST-2) or worse (in MRPC) than the corresponding curve of SGDM.

Figure 4 shows the training losses and development scores when *all hyperparameters* are tuned, as in Fig. 1, but now using DistilRoBERTa instead of DistilBERT. The results are similar to those of Fig. 1, except that SGDM now performs better in terms of development scores on CoLA (where it lagged behind the other adaptive optimizers in Fig. 1) and it now performs poorly on STS-B and MNLI (where it was competent). Hence, these experiments confirm that SGDM is overall clearly better than SGD, but still worse than the adaptive optimizers, when all hyperparameters are tuned. Again, the five adaptive optimizers have very similar development scores, despite occasional larger differences in the training losses they reach.

Figure 5 shows the training losses and development scores when *only the learning rate* is tuned, as in Fig. 2, but now using DistilRoBERTa instead of DistilBERT. As in Fig. 2, AdaMax lags behind (but now only slightly) on CoLA and (more clearly) on MNLI in terms of development scores. The only important difference compared to Fig. 2 is that AdaBound is now also clearly worse than the other adaptive optimizers in development scores on CoLA and MNLI, where it is outperformed even by SGDM. Hence, these experiments confirm that tuning only the learning rate of adaptive optimizers is in most cases (but not always, AdaMax and AdaBound being exceptions on CoLA and MNLI) as good as tuning all their hyperparameters.

Figure 6 shows the training losses and development scores when using *defaults*, as in Fig. 3, but now using DistilRoBERTa instead of DistilBERT. As in Fig. 3, SGDM is now the best in terms of

development scores, and AdaBound is overall the second best. Again AdaBound eventually matches the development scores of SGDM on SST-2 and STSB, but not on the other datasets. The other adaptive optimizers perform overall poorly.

Tables 13–15 show results on *test data*, now using DistilRoBERTa. The best results are now slightly improved in most cases, as one might expect, compared to Tables 1–3 where DistilBERT was used. Otherwise the conclusions are very similar to those of Tables 1–3, they are summarized in the captions of Tables 13–15, and they are aligned with the findings of Fig. 4–6.

	AdamW	AdaMax	Nadam	AdaBound	Adam	SGDM	SGD
ϵ	[1e-7, 1e-5]	[1e-7, 1e-5]	[1e-7, 1e-5]	[1e-7, 1e-5]	[1e-7, 1e-5]	[1e-7, 1e-3]	[1e-7, 1e-3]
ρ_1	[0.8, 0.95]	[0.8, 0.95]	[0.8, 0.95]	[0.8, 0.95]	[0.8, 0.95]	—	—
ρ_2	[0.9, 0.99999]	[0.9, 0.99999]	[0.9, 0.99999]	[0.9, 0.99999]	[0.9, 0.99999]	—	—
δ	[1e-9, 1e-7]	[1e-9, 1e-7]	[1e-9, 1e-7]	[1e-9, 1e-7]	[1e-9, 1e-7]	—	—
α	—	—	[1e-4, 1e-2]	—	—	[0.7, 0.9999]	—
ϵ^*	—	—	—	[1e-2, 1e-1]	—	—	—
γ	—	—	—	[1e-4, 2e-3]	—	—	—

Table 4: **Hyperparameter search space** for all the optimization algorithms. ϵ^* and γ (not shown in Algorithm 3) are used by AdaBound’s lower and upper bound functions (η_t^l, η_t^u); γ controls the convergence speed of these functions, and ϵ^* is the learning rate used in the final training stages, where AdaBound transforms to SGD.

	AdamW	AdaMax	Nadam	AdaBound	Adam	SGDM	SGD
ϵ	1e-3	2e-3	2e-3	1e-3	1e-3	1e-3	1e-3
ρ_1	0.9	0.9	0.9	0.9	0.9	—	—
ρ_2	0.999	0.999	0.999	0.999	0.999	—	—
δ	1e-8	1e-8	1e-8	1e-8	1e-8	—	—
α	—	—	4e-3	—	—	0.9	—
ϵ^*	—	—	—	0.1	—	—	—
γ	—	—	—	1e-3	—	—	—

Table 5: **Default hyperparameter values** per optimizer. The search space for the learning rate (ϵ) does not include the defaults for adaptive optimizers, because it is standard practice when fine-tuning Transformers with adaptive optimizers to use much smaller learning rates.

	SST-2	MRPC	CoLA	MNLI	STS-B	Default	Search Space
AdaBound	8.17e-6	4.51e-6	8.27e-6	1.26e-6	2.06e-6	1e-3	[1e-7, 1e-5]
AdamW	9.80e-6	6.97e-6	5.59e-6	9.50e-6	7.99e-6	1e-3	[1e-7, 1e-5]
AdaMax	8.78e-6	8.44e-6	6.59e-6	9.19e-6	7.58e-6	2e-3	[1e-7, 1e-5]
Nadam	9.74e-6	6.46e-6	6.22e-6	9.58e-6	6.89e-6	2e-3	[1e-7, 1e-5]
Adam	9.61e-6	6.54e-6	7.09e-6	9.47e-6	9.13e-6	1e-3	[1e-7, 1e-5]
SGDM	8.03e-4	3.81e-4	3.41e-4	3.54e-4	6.49e-4	1e-3	[1e-7, 1e-3]
SGD	9.66e-4	7.69e-4	1.68e-4	9.58e-4	9.77e-4	1e-3	[1e-7, 1e-3]

Table 6: **Tuned learning rate (ϵ)** per optimizer and dataset, **averaged over the five random splits**, when **tuning all hyperparameters**, using **DistilBERT**. Default ϵ and search space also shown. All the tuned values are far from the defaults. The **tuned learning rate of SGDM (and SGD) is closer to the default**, comparing to the other optimizers, which may explain why SGDM is the best optimizer overall when using the optimizers with defaults.

	SST-2	MRPC	CoLA	MNLI	STS-B	Default	Search Space
AdaBound	0.87	0.90	0.86	0.86	0.88	0.90	[0.80, 0.95]
AdamW	0.92	0.85	0.88	0.88	0.88	0.90	[0.80, 0.95]
AdaMax	0.90	0.88	0.87	0.86	0.86	0.90	[0.80, 0.95]
Nadam	0.93	0.85	0.82	0.86	0.88	0.90	[0.80, 0.95]
Adam	0.89	0.86	0.88	0.88	0.87	0.90	[0.80, 0.95]

Table 7: **Tuned 1st momentum decay rate (ρ_1)** per optimizer (when applicable) and dataset, **averaged over the five random splits**, when **tuning all hyperparameters**, using **DistilBERT**. Default ρ_1 and search space also shown.

	SST-2	MRPC	CoLA	MNLI	STS-B	Default	Search Space
AdaBound	0.93	0.94	0.95	0.93	0.96	0.999	[0.9, 0.99999]
AdamW	0.92	0.96	0.95	0.98	0.94	0.999	[0.9, 0.99999]
AdaMax	0.92	0.97	0.94	0.93	0.91	0.999	[0.9, 0.99999]
Nadam	0.93	0.96	0.94	0.95	0.96	0.999	[0.9, 0.99999]
Adam	0.96	0.95	0.95	0.95	0.93	0.999	[0.9, 0.99999]

Table 8: **Tuned 2nd momentum decay rate (ρ_2)** per optimizer (when applicable) and dataset, **averaged over the five random splits**, when **tuning all hyperparameters**, using **DistilBERT**. Default ρ_2 and search space also shown.

	SST-2	MRPC	CoLA	MNLI	STSB	Default	Search Space
AdaBound	3.31e-08	2.09e-08	2.88e-08	4.18e-08	2.36e-08	1e-8	[1e-9, 1e-7]
AdamW	6.44e-08	3.37e-08	2.92e-09	3.53e-08	2.77e-08	1e-8	[1e-9, 1e-7]
AdaMax	2.18e-08	4.71e-08	5.59e-09	2.95e-08	2.97e-08	1e-8	[1e-9, 1e-7]
Nadam	1.22e-08	1.78e-08	6.24e-08	7.29e-09	3.46e-08	1e-8	[1e-9, 1e-7]
Adam	2.49e-08	2.67e-08	1.45e-08	4.23e-08	2.19e-08	1e-8	[1e-9, 1e-7]

Table 9: **Tuned small constant δ** for each optimizer (when applicable) and dataset, **averaged over the five random splits**, when **tuning all hyperparameters**, using **DistilBERT**. Default δ and search space also shown.

	SST-2	MRPC	CoLA	MNLI	STSB	Default	Search Space
Nadam	1.69e-3	3.21e-3	3.73e-3	1.01e-4	4.35e-3	4e-3	[1e-4, 1e-2]
SGDM	0.93	0.97	0.86	0.99	0.98	0.9	[0.7, 0.9999]

Table 10: **Tuned momentum strength (α)** per optimizer (when applicable) and dataset, **averaged over the five random splits**, when **tuning all hyperparameters**, using **DistilBERT**. Default α and search space also shown.

AdaBound	SST-2	MRPC	CoLA	MNLI	STSB	Default	Search Space
ϵ^*	7e-2	7.19e-2	6.35e-2	1.22e-1	1.1e-1	4e-3	[1e-2, 1e-1]
γ	3.03e-4	7.21e-4	2.04e-4	9.92e-4	8.46e-4	0.9	[1e-4, 2e-3]

Table 11: **Tuned AdaBound-specific hyperparameters (ϵ^* and γ)** per dataset, **averaged over the five random splits**, when **tuning all hyperparameters**, using **DistilBERT**. ϵ^* and γ are used by AdaBound’s lower and upper bound functions (η_t^l, η_t^u), γ controls the convergence speed of these functions, and ϵ^* is the learning rate used in the final training stages, where AdaBound transforms to SGD. Default values and search space are also shown.

	SST-2	MRPC	CoLA	MNLI	STSB	Default	Search Space
AdaBound	2.43e-6	1.87e-6	1.58e-6	5.42e-6	1.72e-7	1e-3	[1e-7, 1e-5]
AdamW	9.20e-6	7.57e-6	7.63e-6	9.66e-6	8.68e-6	1e-3	[1e-7, 1e-5]
AdaMax	9.67e-6	8.46e-6	8.04e-6	9.74e-6	8.80e-6	2e-3	[1e-7, 1e-5]
Nadam	9.63e-6	8.16e-6	7.83e-6	9.68e-6	8.53e-6	2e-3	[1e-7, 1e-5]
Adam	9.11e-6	7.99e-6	5.49e-6	9.58e-6	7.42e-6	1e-3	[1e-7, 1e-5]
SGDM	9.56e-4	8.35e-4	7.43e-4	9.35e-4	9.42e-4	1e-3	[1e-7, 1e-3]
SGD	9.66e-4	7.69e-4	1.68e-4	9.58e-4	9.77e-4	1e-3	[1e-7, 1e-3]

Table 12: **Tuned learning rate (ϵ)** per optimizer and dataset, **averaged over the five random splits**, when **tuning only the learning rate**, using **DistilBERT**. Default ϵ and search space also shown. As in Table 6, the **tuned learning rate of SGDM** (and SGD) is **closer to the default**, comparing to the other optimizers, which may explain why SGDM is the best optimizer overall when using the optimizers with defaults.

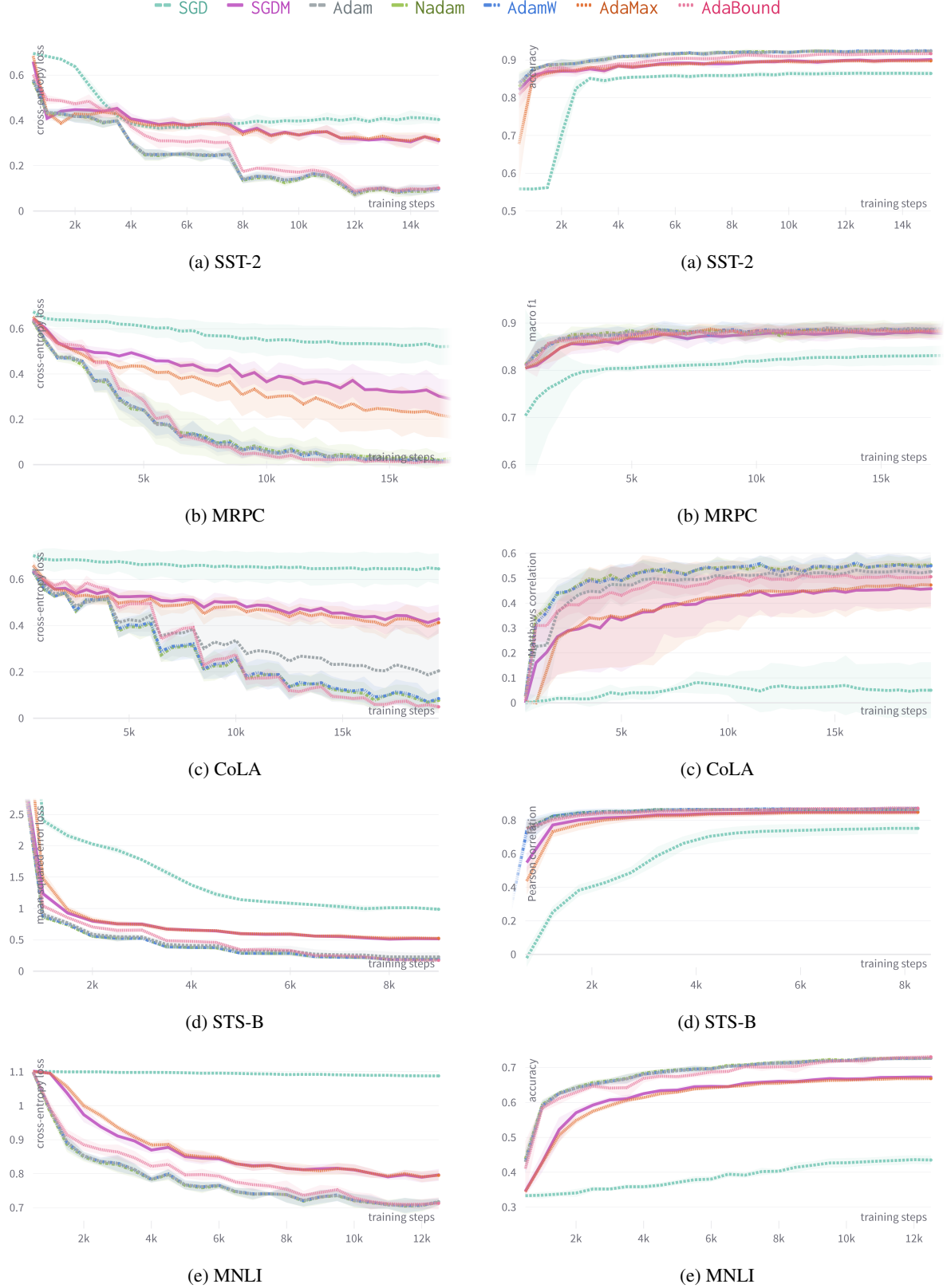


Figure 2: **Training loss (left) and evaluation score on development data (right) having tuned only the learning rate of the optimizers**, as a function of training steps, using **DistilBERT**. For each dataset, we use **five random data splits**, and plot the **average** and **standard deviation** (shadow). As when tuning all the hyperparameters (Fig. 1), SGD is clearly the **worst**, but adding **Momentum (SGDM)** turns it to a **competent** optimizer in terms of development scores. The five **adaptive optimizers** (Adam, Nadam, AdamW, AdaMax, AdaBound) **perform similarly overall** in terms of **development scores**, except for AdaMax which lags behind on CoLA and (more) on MNLI. Differences in training loss do not necessarily give rise to substantial differences in development scores.

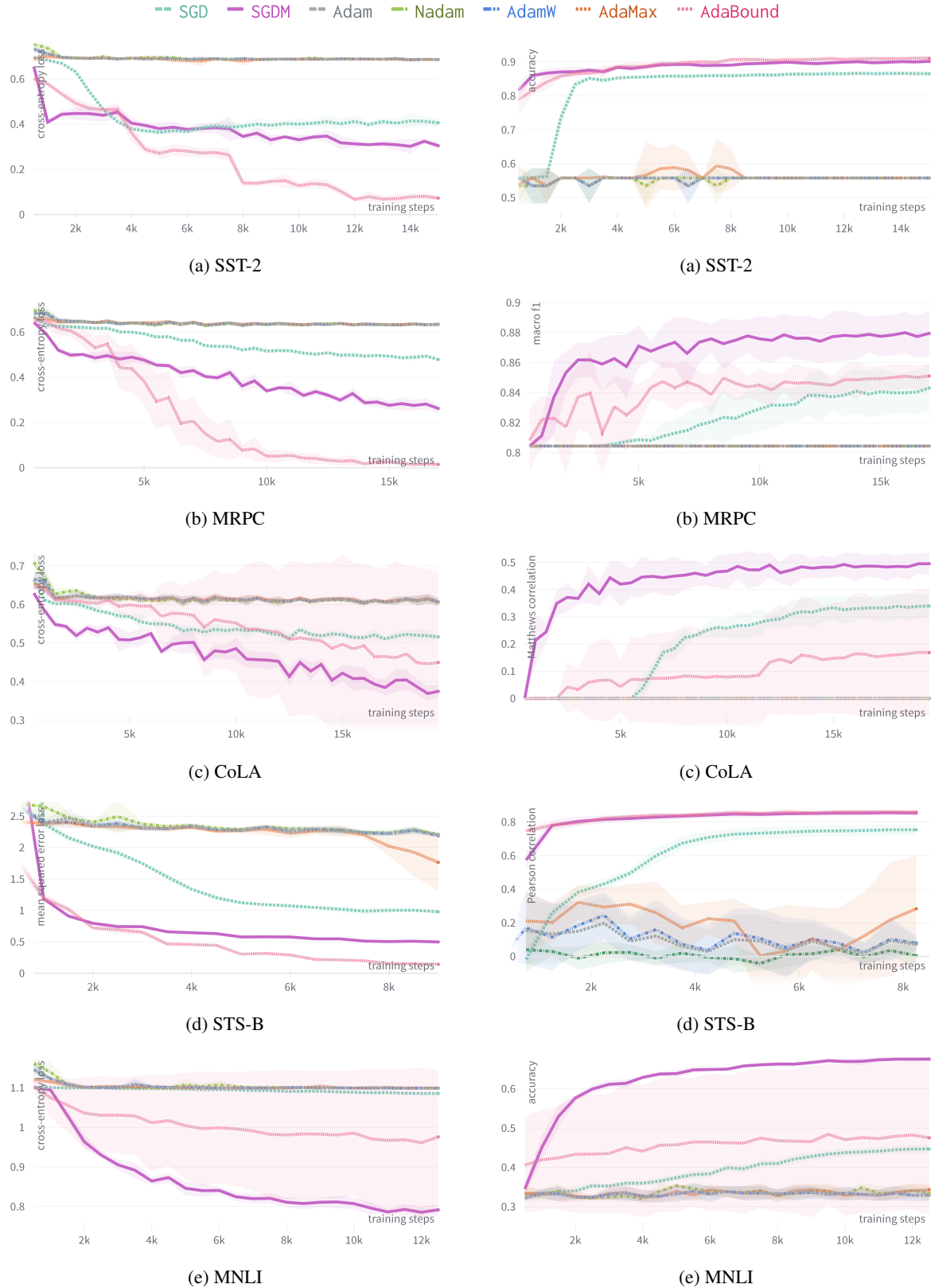


Figure 3: **Training loss (left) and evaluation score on development data (right) with all hyperparameters of the optimizers set to their defaults, as a function of training steps, using DistilBERT. Again, we use five random data splits, and plot the average and standard deviation (shadow). SGDM is not affected by the lack of hyperparameter tuning and is now the best overall, improving upon plain SGD. AdaBound matches SGDM in performance on SST-2 and STS-B, but does not perform as well on the other datasets, although it is still better than the other adaptive optimizers. Differences in training loss are not always reflected to differences in development scores.**

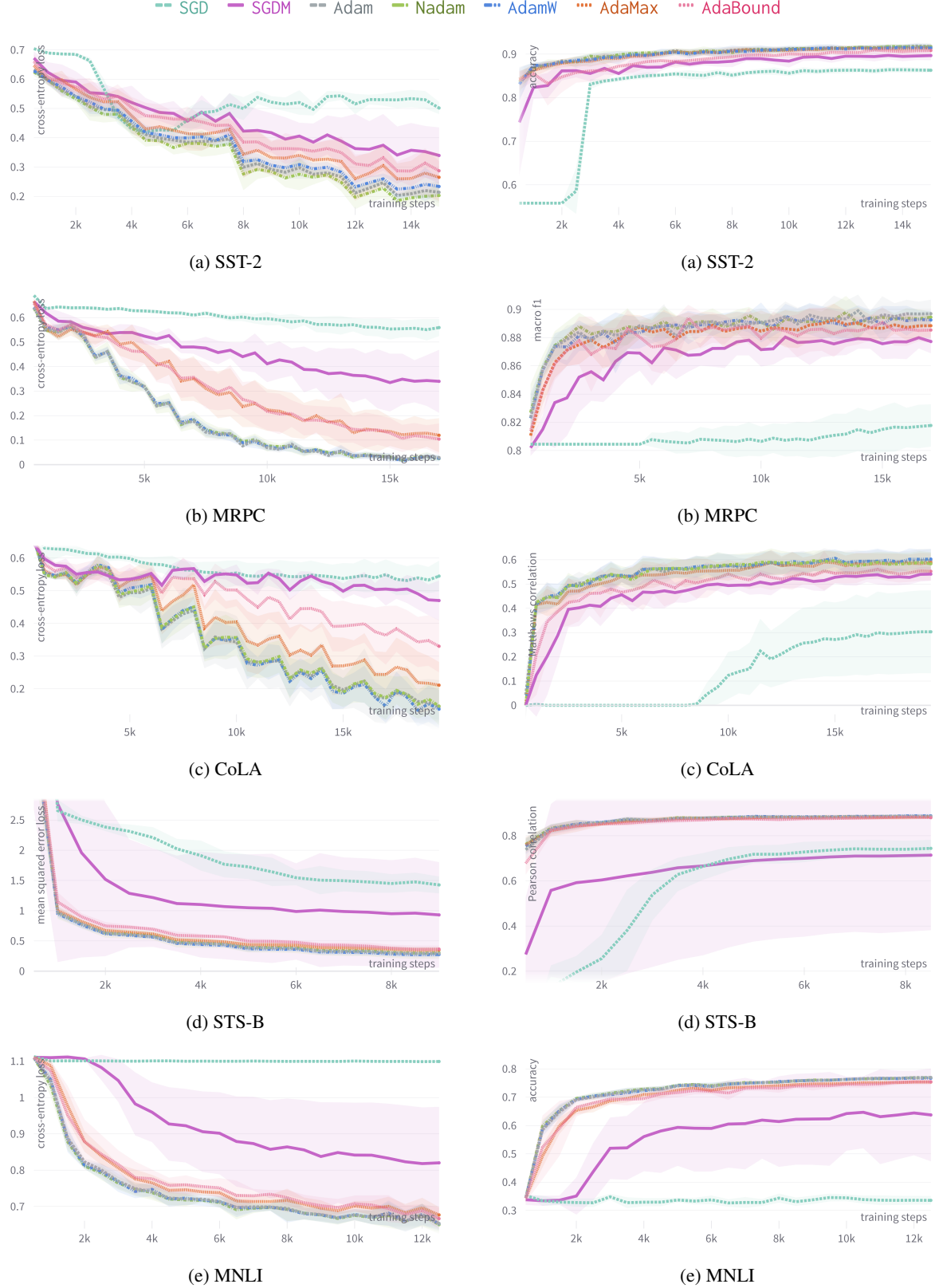


Figure 4: **Training loss (left) and evaluation score on development data (right) with all hyperparameters of the optimizers tuned, using DistilRoBERTa.** For each dataset, we use **five random data splits**, and plot the **average and standard deviation (shadow)**. The **results are similar** to those of the the experiments with DistilBERT (cf. Fig. 1), except that SGDM now performs better in development score on CoLA (where it lagged behind the other adaptive optimizers in Fig. 1) and it now performs poorly on STS-B and MNLI (where it was competent). Hence, these experiments confirm that **SGDM is overall clearly better than SGD**, but **still worse than the adaptive optimizers**, when all hyperparameters are tuned. Again, **the five adaptive optimizers perform very similarly**.

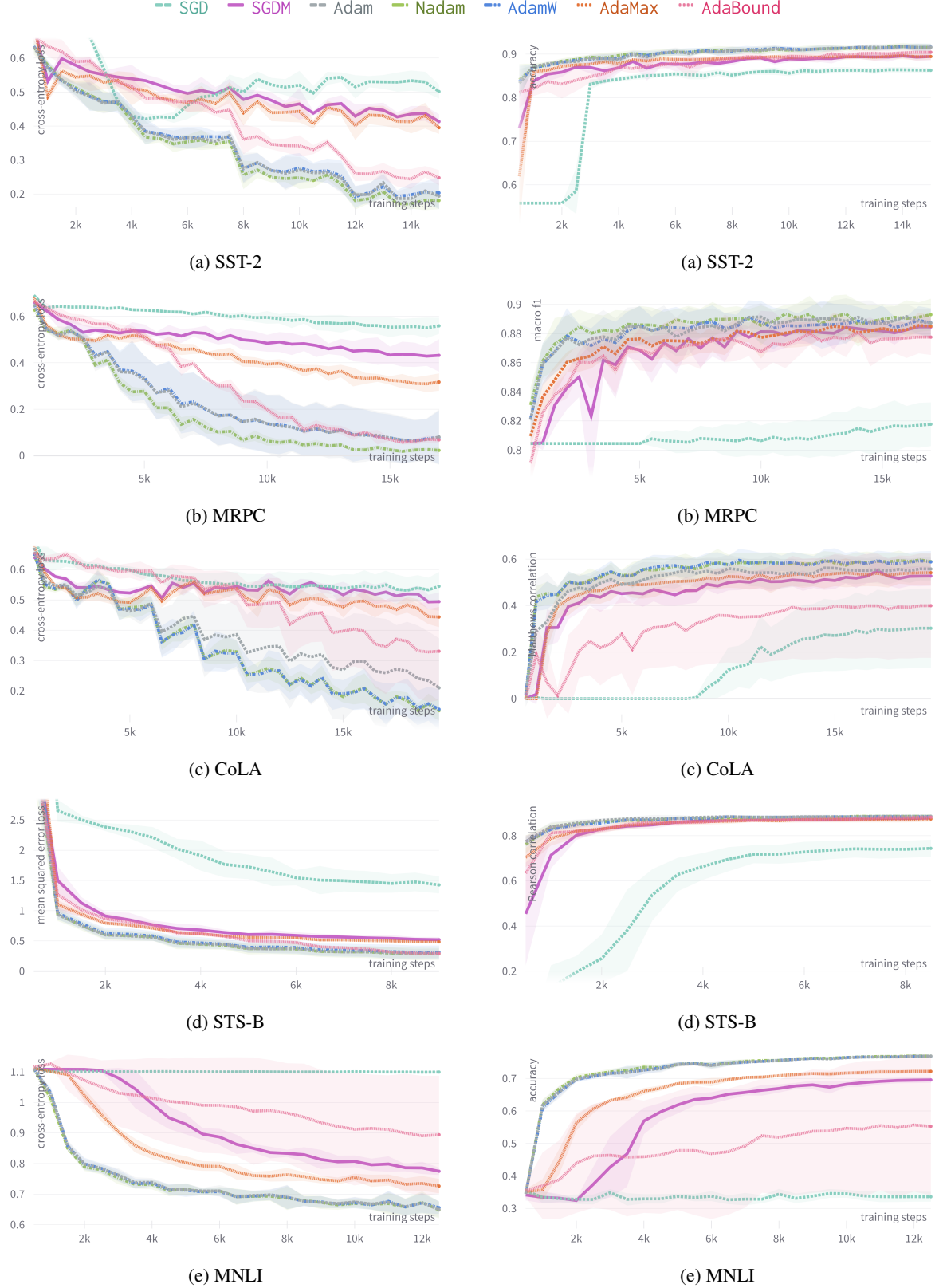


Figure 5: **Training loss (left) and evaluation score on development data (right) having tuned only the learning rate, using DistilRoBERTa.** For each dataset, we use **five random data splits**, and plot the **average and standard deviation** (shadow). As in the corresponding DistilBERT experiments (Fig. 2), **AdaMax** lags (now slightly) behind on CoLA and (more clearly) on MNLI in terms of development scores. The only important difference from Fig. 2 is that **AdaBound** is now also clearly worse than the other adaptive optimizers in development scores on CoLA and MNLI, where it is outperformed even by SGDM. Hence, these experiment confirm that **tuning only the learning rate** of adaptive optimizers is **in most cases (not always) as good as tuning all their hyperparameters**.

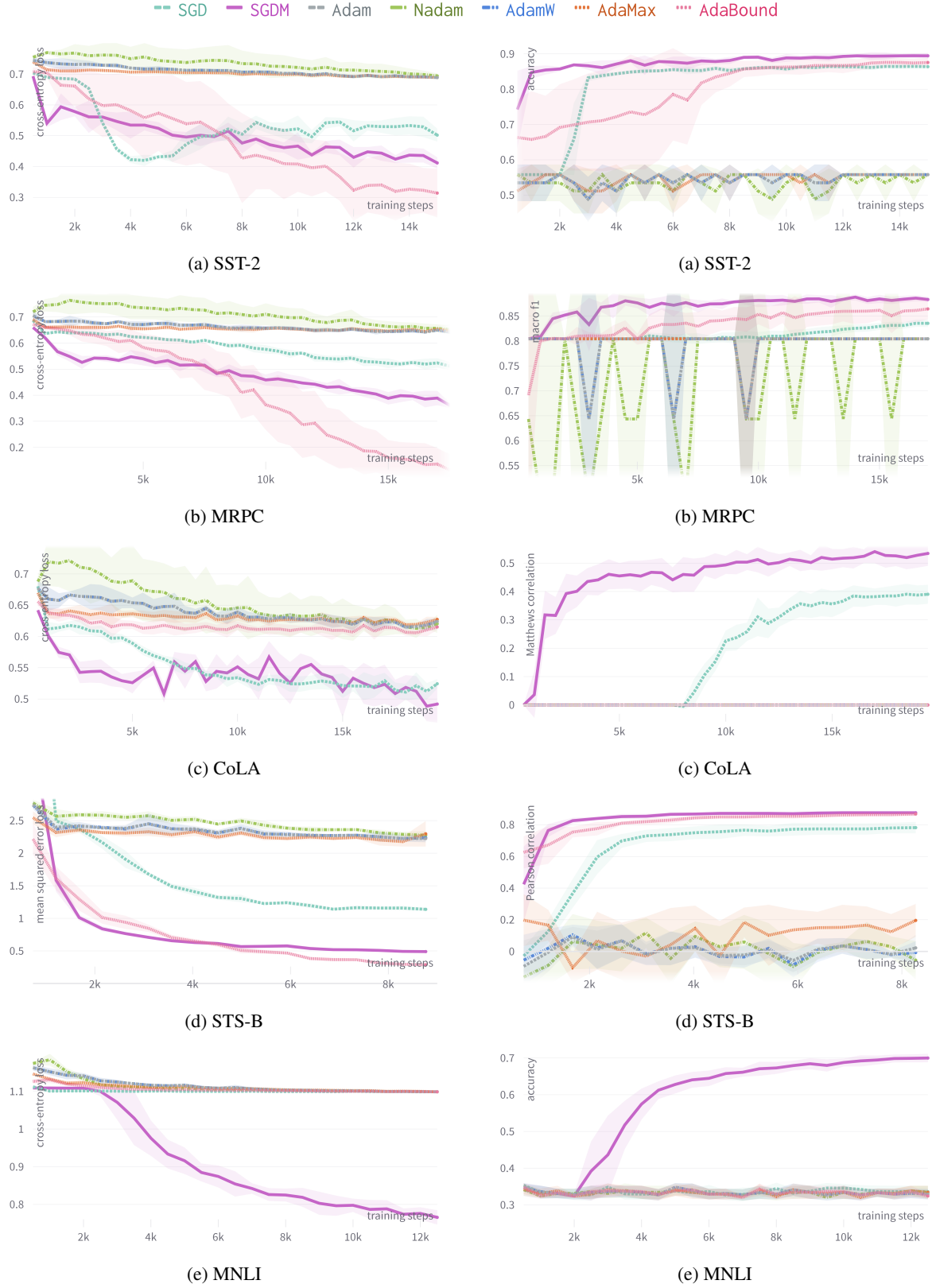


Figure 6: **Training loss (left) and evaluation score on development data (right) with all hyperparameters of the optimizers set to their defaults, using DistilRoBERTa.** Again, we use **five random data splits**, and plot the **average and standard deviation** (shadow). As in the corresponding experiments with DistilBERT (Fig. 3), **SGDM is now the best** in terms of development scores; again **AdaBound** eventually reaches almost the same performance as SGDM on SST-2 and STS-B, but not on the other datasets; the other adaptive optimizers perform overall poorly.

Optimizer	SST-2 Accuracy	MRPC Macro-F1	CoLA Matthews	STS-B Pearson	MNLI Accuracy
AdaBound	89.91 (0.47)	81.47 (2.21)	0.57 (0.03)	0.88 (0.01)	75.52 (1.08)
AdamW	91.39 (0.70)	81.88 (1.56)	0.58 (0.05)	0.88 (0.01)	76.74 (0.44)
AdaMax	91.94 (2.84)	82.33 (1.82)	0.56 (0.03)	0.88 (0.01)	75.40 (0.38)
Nadam	91.35 (0.87)	82.27 (1.83)	0.57 (0.02)	0.88 (0.00)	76.88 (0.49)
Adam	91.17 (0.65)	83.08 (1.34)	0.58 (0.03)	0.88 (0.00)	76.89 (0.47)
SGDM	89.45 (0.79)	80.70 (1.51)	0.52 (0.04)	0.71 (0.34)	65.37 (13.46)
SGD	86.08 (0.32)	65.30 (5.62)	0.27 (0.17)	0.74 (0.03)	35.86 (0.52)

Table 13: Evaluation scores on **test data** with **all hyperparameters tuned**, using **DistilRoBERTa**. For each dataset, we use **five random splits** and report the **average** test score and the **standard deviation**. As one would expect, the best scores (bold) are now slightly **better than those of DistilBERT** (cf. Table 1). Otherwise, the **findings are similar** to those of the experiments with DistilBERT (Table 1), except that SGDM now performs better on CoLA (where it lagged behind the other adaptive optimizers in Table 1) and it now performs poorly on STS-B and MNLI (where it was competent). Hence, these experiments confirm that **SGDM is overall clearly better than SGD**, but still **worse than the adaptive optimizers**, when all hyperparameters are tuned. Again, **the five adaptive optimizers perform very similarly**. These findings are also aligned with those of Fig. 4.

Optimizer	SST-2 Accuracy	MRPC Macro-F1	CoLA Matthews	STS-B Pearson	MNLI Accuracy
AdaBound	90.17 (0.93)	80.84 (1.72)	0.43 (0.15)	0.88 (0.00)	56.92 (19.86)
AdamW	91.08 (1.09)	81.15 (1.76)	0.59 (0.04)	0.88 (0.00)	76.86 (0.48)
AdaMax	89.17 (0.26)	80.96 (1.74)	0.53 (0.03)	0.87 (0.01)	72.34 (0.39)
Nadam	91.08 (0.91)	81.81 (1.85)	0.57 (0.03)	0.88 (0.00)	76.97 (0.37)
Adam	91.45 (0.95)	81.47 (2.05)	0.56 (0.05)	0.88 (0.00)	76.73 (0.59)
SGDM	89.11 (0.33)	81.68 (1.39)	0.53 (0.03)	0.87 (0.01)	69.67 (0.88)
SGD	86.08 (0.32)	65.30 (5.62)	0.27 (0.17)	0.74 (0.03)	35.86 (0.52)

Table 14: Evaluation scores on **test data**, having **tuned only the learning rate**, using **DistilRoBERTa**. Again, we use five random splits and report the average and standard deviation. As in the corresponding DistilBERT experiments (Table 2), **AdaMax** lags behind on CoLA and (more) on MNLI. The only important difference compared to Table 2 is that **AdaBound is now also clearly worse than the other adaptive optimizers** on CoLA and MNLI, where it is outperformed even by SGDM (but not plain SGD). Hence, these experiments confirm the conclusion that **tuning only the learning rate** of adaptive optimizers is **in most cases (but not always) as good as tuning all their hyperparameters**. These findings are aligned with those of Fig. 5. Again, the best scores (bold) are better than those of DistilBERT (Table 2), except for SST-2.

Optimizer	SST-2 Accuracy	MRPC Macro-F1	CoLA Matthews	STS-B Pearson	MNLI Accuracy
AdaBound	88.36 (0.62)	76.48 (2.54)	0.00 (0.00)	0.86 (0.00)	35.34 (0.01)
AdamW	55.79 (0.01)	62.57 (0.00)	0.00 (0.00)	0.14 (0.11)	35.33 (0.01)
AdaMax	55.80 (0.00)	62.57 (0.00)	0.00 (0.00)	0.30 (0.10)	35.33 (0.01)
Nadam	55.80 (0.01)	62.57 (0.00)	0.00 (0.00)	0.17 (0.03)	35.34 (0.01)
Adam	55.80 (0.00)	62.57 (0.00)	0.00 (0.00)	0.14 (0.11)	35.34 (0.01)
SGDM	89.20 (0.62)	82.28 (1.66)	0.54 (0.04)	0.87 (0.01)	70.14 (0.81)
SGD	86.29 (0.34)	69.14 (3.30)	0.36 (0.03)	0.78 (0.02)	35.90 (0.49)

Table 15: Evaluation scores on **test data**, with **all hyperparameters** of the optimizers set to their **defaults**, using **DistilRoBERTa**. Again, we use five random splits and report the average and standard deviation. As in the corresponding experiments with DistilBERT (Table 3), **SGDM is now the best** optimizer; again **AdaBound** performs well on SST-2 and STSB, now also relatively well on MRPC, but not on the other two datasets; the other adaptive optimizers perform much worse. These findings are aligned with those of Fig. 6. Again, the best scores (bold) are better than those of DistilBERT (Table 2), except for SST-2.