

---

# A Large Batch Optimizer Reality Check: Traditional, Generic Optimizers Suffice Across Batch Sizes

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Recently the LARS and LAMB optimizers have been proposed for training neural  
2 networks faster using large batch sizes. LARS and LAMB add layer-wise normal-  
3 ization to the update rules of Heavy-ball momentum and Adam, respectively, and  
4 have become popular in prominent benchmarks and deep learning libraries. How-  
5 ever, without fair comparisons to standard optimizers, it remains an open question  
6 whether LARS and LAMB have any benefit over traditional, generic algorithms. In  
7 this work we demonstrate that standard optimization algorithms such as Nesterov  
8 momentum and Adam can match or exceed the results of LARS and LAMB at large  
9 batch sizes. Our results establish new, stronger baselines for future comparisons  
10 at these batch sizes and shed light on the difficulties of comparing optimizers for  
11 neural network training more generally.

## 12 1 Introduction

13 In recent years, hardware systems employing GPUs and TPUs have enabled neural network training  
14 programs to process dramatically more data in parallel than ever before. The most popular way to  
15 exploit these systems is to increase the batch size in the optimization algorithm (i.e. the number  
16 of training examples processed per training step). On many workloads, modern systems can scale  
17 to larger batch sizes without significantly increasing the time per step [Jouppi et al., 2017, Wang  
18 et al., 2019], thus proportionally increasing the number of training examples processed per second.  
19 If researchers can use this increased throughput to reduce the time required to train each neural  
20 network, then they should achieve better results by training larger models, using larger datasets, and  
21 by exploring new ideas more rapidly.

22 As the capacity for data parallelism continues to increase, practitioners can take their existing,  
23 well-tuned training configurations and re-train with larger batch sizes, hoping to achieve the same  
24 performance in less training time [e.g. Ying et al., 2018]. On an idealized data-parallel system with  
25 negligible overhead from increasing the batch size, they might hope to achieve *perfect scaling*, a  
26 proportional reduction in training time as the batch size increases.

27 However, achieving perfect scaling is not always straightforward. Changing the batch size changes  
28 the training dynamics, requiring the training hyperparameters (e.g. learning rate) to be carefully  
29 re-tuned in order to maintain the same level of validation performance.<sup>1</sup> In addition, smaller batch  
30 sizes provide implicit regularization from gradient noise that may need to be replaced by other forms  
31 of regularization when the batch size is increased. Finally, even with perfect tuning, increasing

---

<sup>1</sup> Although there are heuristics for adjusting the learning rate as the batch size changes, these heuristics inevitably break down sufficiently far from the initial batch size and it is also not clear how to apply them to other training hyperparameters (e.g. momentum).

32 the batch size eventually produces diminishing returns. After a critical batch size, the number of  
33 training steps cannot be decreased in proportion to the batch size – the number of epochs must  
34 increase to match the validation performance of the smaller batch size. See Shallue et al. 2019 for a  
35 survey of the effects of data parallelism on neural network training. Once these effects are taken into  
36 account, there is no strong evidence that increasing the batch size degrades the maximum achievable  
37 performance on any workload. At the same time, the ever-increasing capacity for data parallelism  
38 presents opportunities for new regularization techniques that can replace the gradient noise of smaller  
39 batch sizes and new optimization algorithms that can extend perfect scaling to larger batch sizes by  
40 using more sophisticated gradient information [Zhang et al., 2019].

41 You et al. [2017] proposed the LARS optimization algorithm in the hope of speeding up neural  
42 network training by exploiting larger batch sizes. LARS is a variant of stochastic gradient descent  
43 (SGD) with momentum [Polyak, 1964] that applies layer-wise normalization before applying each  
44 gradient update. Although it is difficult to draw strong conclusions from the results presented in the  
45 LARS paper,<sup>2</sup> the MLPerf<sup>3</sup> Training benchmark<sup>4</sup> adopted LARS as one of two allowed algorithms  
46 in the closed division for ResNet-50 on ImageNet and it became the *de facto* standard algorithm for  
47 that benchmark task. With MLPerf entrants competing to find the fastest-training hyperparameters  
48 for LARS, the first place submissions in the two most recent MLPerf Training competitions used  
49 LARS to achieve record training speeds with batch sizes of 32,678 and 65,536, respectively. No  
50 publications or competitive submissions to MLPerf have attempted to match these results with a  
51 standard optimizer (e.g. Momentum or Adam). However, MLPerf entrants do not have a strong  
52 incentive (nor are necessarily permitted by the rules) to explore other algorithms because MLPerf  
53 Training is a systems benchmark that requires algorithmic equivalence between submissions to make  
54 fair comparisons. Moreover, since the main justification for LARS is its excellent performance on  
55 ResNet-50 at large batch sizes, more work is needed to quantify any benefit of LARS over standard  
56 algorithms at any batch size.

57 You et al. [2019] later proposed the LAMB optimizer to speed up pre-training for BERT [Devlin  
58 et al., 2018] using larger batch sizes after concluding that LARS was not effective across workloads.  
59 LAMB is a variant of Adam [Kingma and Ba, 2014] that adds a similar layer-wise normalization step  
60 to LARS. You et al. [2019] used LAMB for BERT pre-training with batch sizes up to 65,536 and  
61 claimed that Adam cannot match the performance of LAMB beyond batch size 16,384.

62 In this paper, we demonstrate that standard optimizers, without any layer-wise normalization tech-  
63 niques, can match or improve upon the large batch size results used to justify LARS and LAMB. In  
64 Section 2, we show that Nesterov momentum [Nesterov, 1983] matches the performance of LARS on  
65 the ResNet-50 benchmark with batch size 32,768. We are the first to match this result with a standard  
66 optimizer. In Section 3, contradicting the claims in You et al. [2019], we show that Adam obtains  
67 better BERT pre-training results than LAMB at the largest batch sizes, resulting in better downstream  
68 performance metrics after fine-tuning.

69 In addition, we establish a new state-of-the-art for BERT pretraining speed, reaching an F1 score of  
70 90.46 in 7,818 steps using Adam at batch size 65,536 (we report training speed in steps because our  
71 focus is algorithmic efficiency, but since we compare LARS and LAMB to simpler optimizers, fewer  
72 training steps corresponds to faster wall-time in an optimized implementation – our BERT result  
73 with Adam also improves upon the wall-time record of LAMB reported in You et al. 2019). Taken  
74 together, our results establish stronger training speed baselines for these tasks and batch sizes, which  
75 we hope will assist future work aiming to accelerate training using larger batch sizes.

76 In addition to the contributions mentioned above, we demonstrate several key effects that are often  
77 overlooked by studies aiming to establish the superiority of new optimization algorithms. We show  
78 that future work must carefully disentangle regularization and optimization effects when comparing a  
79 new optimizer to baselines. We also report several under-documented details used to generate the  
80 best LARS and LAMB results, a reminder that future comparisons should document any novel tricks  
81 and include them in baselines. Finally, our results add to existing evidence in the literature on the  
82 difficulty of performing independently rigorous hyperparameter tuning for optimizers and baselines.

---

<sup>2</sup> The modified AlexNet on ImageNet benchmark did not have well-established accuracy targets from prior work and LARS used a more general learning rate schedule than the momentum baseline. For ResNet-50 on ImageNet, LARS achieved sub-par accuracy numbers and was not compared to any other optimizer at the same batch size, leaving open the possibility that a generic optimizer would scale just as well as LARS. <sup>3</sup> MLPerf is a trademark of MLCommons.org. <sup>4</sup> <https://mlperf.org/training-overview>

83 In particular, we show that the optimal shape of the learning rate schedule is optimizer-dependent (in  
84 addition to the scale), and that differences in the schedule can dominate optimizer comparisons at  
85 smaller step budgets and become less important at larger step budgets.

## 86 1.1 Related work

87 Shallue et al. [2019] and Zhang et al. [2019] explored the effects of data parallelism on neural network  
88 training for different optimizers, finding no evidence that larger batch sizes degrade performance  
89 and demonstrating that different optimizers can achieve perfect scaling up to different critical batch  
90 sizes. You et al. [2017, 2019] developed the LARS and LAMB optimizers in the hope of speeding up  
91 training by achieving perfect scaling beyond standard optimizers. Many other recent papers have  
92 proposed new optimization algorithms for generic batch sizes or larger batch sizes [see Schmidt  
93 et al., 2020]. Choi et al. [2019] and Schmidt et al. [2020] demonstrated the difficulties with fairly  
94 comparing optimizers, showing that the hyperparameter tuning protocol is a key determinant of  
95 optimizer rankings. The MLPerf Training benchmark [Mattson et al., 2019] provides a competitive  
96 ranking of neural network training systems, but does not shed much light on the relative performance  
97 of optimizers because entrants are limited in the algorithms they can use and the hyperparameters  
98 they can tune.

## 99 2 Matching LARS on ImageNet

100 The MLPerf training benchmark for ResNet-50 v1.5 on ImageNet [Mattson et al., 2019] aims to  
101 reach 75.9% validation accuracy in the shortest possible wall-clock time. In the closed division of  
102 the competition, entrants must choose between two optimizers, SGD with momentum or LARS, and  
103 are only allowed to tune a specified subset of the optimization hyperparameters, with the remaining  
104 hyperparameter values set by the competition rules.<sup>5</sup> The winning entries in the two most recent  
105 competitions used LARS with batch size 32,768 for 72 training epochs<sup>6</sup> and LARS with batch size  
106 65,536 for 88 training epochs,<sup>7</sup> respectively. Kumar et al. [2019] later improved the training time  
107 for batch size 32,768 by reaching the target accuracy in 64 epochs. These are currently the fastest  
108 published results on the ResNet-50 benchmark. However, it has been unclear whether LARS was  
109 necessary to achieve these training speeds since no recent published results or competitive MLPerf  
110 submissions have used another optimizer. In this section, we describe how we matched the 64 epoch,  
111 32,768 batch size result of LARS using standard Nesterov momentum.<sup>8</sup>

112 A fair benchmark of training algorithms or hardware systems must account for stochasticity in  
113 individual training runs. In the MLPerf competition, the benchmark metric is the mean wall-clock  
114 time of 5 trials after the fastest and slowest trials are excluded. Only 4 out of the 5 trials need to reach  
115 the target accuracy and there is no explicit limit on the number of times an entrant can try a different  
116 set of 5 trials. Since our goal is to compare algorithms, rather than systems, we aim to match the  
117 LARS result in terms of training steps instead (but since Nesterov momentum is computationally  
118 simpler than LARS, this would also correspond to faster wall-clock time on an optimized system).  
119 Specifically, we measure the median validation accuracy over 50 training runs with a fixed budget of  
120 2,512 training steps<sup>9</sup> at a batch size of 32,768. When we ran the published LARS training pipeline,<sup>10</sup>  
121 LARS achieved a median accuracy of 75.97% and reached the target in 35 out of 50 trials. We  
122 consider the LARS result to be matched by another optimizer if the median over 50 trials exceeds the  
123 target of 75.9%.

### 124 2.1 Nesterov momentum at batch size 32k

125 This section describes how we used the standard Nesterov momentum optimizer to train the ResNet-  
126 50 v1.5 on ImageNet to 75.9% validation accuracy in 2,512 update steps at a batch size of 32,768,  
127 matching the best published LARS result at this batch size. Although we implemented our own  
128 training program, the only logical changes we made to the published LARS pipeline were to the  
129 optimizer and the optimization hyperparameters. Our model implementation and data pre-processing  
130 pipeline were identical to those required under the MLPerf closed division rules (see Appendix B).

<sup>5</sup> <https://git.io/JtknD>

<sup>6</sup> <https://mlperf.org/training-results-0-6>

<sup>7</sup> <https://mlperf.org/training-results-0-7> <sup>8</sup> The 88 epoch, 65,536 batch size result is faster in terms of wall-clock time but requires more training epochs, indicating that it is beyond LARS's perfect scaling regime. Although LARS obtains diminishing returns when increasing the batch size from 32,768 to 65,536, future work could investigate whether Nesterov momentum drops off more or less rapidly than LARS.

<sup>9</sup> Corresponding to 64 training epochs in Kumar et al. [2019]. <sup>10</sup> <https://git.io/JtsLQ>

131 We present two Nesterov momentum hyperparameter configurations that achieve comparable per-  
 132 formance to LARS. Configuration A achieved a median accuracy of 75.97% (the same as LARS)  
 133 and reached the target accuracy in 34 out of 50 trials. Configuration B is a modified version of  
 134 Configuration A designed to make as few changes as possible to the LARS hyperparameters; it  
 135 achieved a median accuracy of 75.92% and reached the target in 29 out of 50 trials. See Appendix D.1  
 136 for the complete hyperparameter configurations.

137 To achieve these results, we tuned the hyperparameters of the training pipeline from scratch using  
 138 Nesterov momentum. We ran a series of experiments, each of which searched over a hand-designed  
 139 hyperparameter search space using quasi-random search [Bousquet et al., 2017]. Between each  
 140 experiment, we modified the previous search space and/or tweaked the training program to include  
 141 optimization tricks and non-default hyperparameter values we discovered in the state-of-the-art LARS  
 142 pipeline. The full sequence of experiments we ran, including the number of trials, hyperparameters  
 143 tuned, and search space ranges, are provided in Appendix D.4. Once we had matched the LARS  
 144 result with Configuration A, we tried setting each hyperparameter to its value in the LARS pipeline in  
 145 order to find the minimal set of changes that still achieved the target result, producing Configuration  
 146 B. The remainder of this section describes the hyperparameters we tuned and the techniques we  
 147 applied on the journey to these results.

### 148 2.1.1 Nesterov Momentum Optimizer

149 Nesterov momentum is a variant of classical or “heavy-ball” momentum defined by the update rule

$$v_{t+1} = \mu v_t + \nabla \ell(\theta_t),$$

$$\theta_{t+1} = \theta_t - \eta_t (\mu v_{t+1} + \nabla \ell(\theta_t)),$$

150 where  $v_0 = 0$ ,  $\theta_t$  is the vector of model parameters after  $t$  steps,  $\nabla \ell(\theta_t)$  is the gradient of the loss  
 151 function  $\ell(\theta)$  averaged over a batch of training examples,  $\mu$  is the momentum, and  $\eta_t$  is the learning  
 152 rate for step  $t$ . We prefer Nesterov momentum over classical momentum because it tolerates larger  
 153 values of its momentum parameter [Sutskever et al., 2013] and sometimes outperforms classical  
 154 momentum, although the two algorithms perform similarly on many tasks [Shallue et al., 2019, Choi  
 155 et al., 2019]. We tuned the Nesterov momentum  $\mu$  in Configurations A and B. We discuss the learning  
 156 rate schedule  $\{\eta_t\}$  separately in Section 2.1.4.

### 157 2.1.2 Batch normalization

158 The ResNet-50 v1.5 model uses batch normalization [Ioffe and Szegedy, 2015], defined as

$$\text{BN}(x^{(l)}) = \left( \frac{x^{(l)} - \text{mean}(x^{(l)})}{\sqrt{\text{var}(x^{(l)}) + \epsilon}} \right) \times \gamma^{(l)} + \beta^{(l)},$$

159 where  $x^{(l)}$  is a vector of pre-normalization outputs from layer  $l$ ,  $\text{mean}(\cdot)$  and  $\text{var}(\cdot)$  denote the  
 160 element-wise sample mean and variance across the batch of training examples,<sup>11</sup> and  $\gamma^{(l)}$  and  $\beta^{(l)}$   
 161 are trainable model parameters.

162 Batch normalization introduces the following tuneable hyperparameters:  $\epsilon$ , the small constant added  
 163 to the sample variance; the initial values of  $\gamma^{(l)}$  and  $\beta^{(l)}$ ; and  $\rho$ , which governs the exponential  
 164 moving averages of the scaling factors used in evaluation. The LARS pipeline uses  $\epsilon = 10^{-5}$  and  
 165  $\rho = 0.9$ . It sets the initial value of  $\beta^{(l)}$  to 0.0 everywhere, but the initial value of  $\gamma^{(l)}$  depends on  
 166 the layer: it sets  $\gamma^{(l)}$  to 0.0 in the final batch normalization layer of each residual block, and to 1.0  
 167 everywhere else. In Configuration A, we tuned  $\epsilon$ ,  $\rho$ , and  $\gamma_0$ , the initial value of  $\gamma^{(l)}$  in the final batch  
 168 normalization layer of each residual block. In Configuration B, we used the same values as LARS for  
 169  $\epsilon$  and  $\rho$ , but we found that choosing  $\gamma_0$  between 0.0 and 1.0 was important for matching the LARS  
 170 result with Nesterov momentum.

### 171 2.1.3 Regularization

172 In Configuration A, we tuned both the L2 regularization coefficient  $\lambda$  and label smoothing  
 173 coefficient  $\tau$  [Szegedy et al., 2016]. The LARS pipeline uses  $\lambda = 10^{-4}$  and  $\tau = 0.1$ .

<sup>11</sup> In a distributed training environment the mean and variance are commonly computed over a subset of the full batch. The LARS pipeline uses a “virtual batch size” of 64, which we also use to avoid changing the training objective [Hoffer et al., 2017].

174 Crucially, the LARS pipeline does not apply L2 regularization to the bias variables of the  
 175 ResNet model nor the batch normalization parameters  $\gamma^{(l)}$  and  $\beta^{(l)}$  (indeed, the published  
 176 LARS pipeline does not even apply LARS to these parameters – it uses Heavy-ball momen-  
 177 tum). This detail is extremely important for both LARS and Nesterov momentum to achieve  
 178 the fastest training speed. Configuration B used the same  $\lambda$  and  $\tau$  as Configuration A.  
 179

### 180 2.1.4 Learning rate schedule

181 The LARS pipeline uses a piecewise polynomial schedule

$$\eta_t = \begin{cases} \eta_{\text{init}} + (\eta_{\text{peak}} - \eta_{\text{init}}) \left(\frac{t}{t_{\text{warmup}}}\right)^{p_{\text{warmup}}}, & t \leq t_{\text{warmup}} \\ \eta_{\text{final}} + (\eta_{\text{peak}} - \eta_{\text{final}}) \left(\frac{T-t}{T-t_{\text{warmup}}}\right)^{p_{\text{decay}}}, & t > t_{\text{warmup}}, \end{cases}$$

182 with  $\eta_{\text{init}} = 0.0$ ,  $\eta_{\text{peak}} = 29.0$ ,  $\eta_{\text{final}} = 10^{-4}$ ,  $p_{\text{warmup}} = 1$ ,  
 183  $p_{\text{decay}} = 2$ , and  $t_{\text{warmup}} = 706$  steps. In Configuration A, we re-  
 184 tuned all of these hyperparameters with Nesterov momentum.  
 185 In Configuration B, we set  $\eta_{\text{init}}$ ,  $p_{\text{decay}}$ , and  $t_{\text{warmup}}$  to the same  
 186 values as LARS, changing only  $p_{\text{warmup}}$  from 1 to 2 and re-  
 187 scaling  $\eta_{\text{peak}}$  and  $\eta_{\text{final}}$ .

### 188 2.1.5 Comparing Nesterov momentum and LARS

189 Table 1 shows the hyperparameter values for Configuration B that differ from the state-  
 190 of-the-art LARS pipeline. Aside from re-tuning the momentum, learning rate scale, and  
 191 regularization hyperparameters (whose optimal values are all expected to change with the  
 192 optimizer), the only changes are setting  $p_{\text{warmup}}$  to 2 instead of 1 and re-tuning  $\gamma_0$ .  
 193

194 Figure 1 shows the LARS learning rate schedule compared to the Nesterov momentum schedule. Even though  
 195 these schedules are similar, we found that each optimizer had a different optimal value of the warmup polynomial  
 196 power. As Table 2 shows, Nesterov momentum performs better with  $p_{\text{warmup}} = 2$  instead of 1, while the opposite  
 197 is true with LARS. As discussed in Agarwal et al. [2020], optimizers can induce implicit step size schedules that  
 198 strongly influence their training dynamics and solution quality, and it appears from Table 2 that the implicit step  
 199 sizes of Nesterov momentum and LARS may evolve differently, causing the shapes of their optimal learning rate  
 200 schedules to differ.  
 201  
 202  
 203  
 204  
 205  
 206

207 Although the main concern of a practitioner is validation performance, the primary task of an  
 208 optimization algorithm is to minimize training loss. Table 2 shows that Nesterov momentum achieves  
 209 higher training accuracy than LARS, despite similar validation performance. Thus, it may be more  
 210 appropriate to consider the layerwise normalization of LARS to be a regularization technique, rather  
 211 than an optimization technique.

212 Spending even more effort tuning LARS or Nesterov momentum would likely further improve the  
 213 current state-of-the-art for that optimizer. Meaningful optimizer comparisons are only possible  
 214 with independent and equally intensive tuning efforts, and we do not claim that either optimizer  
 215 outperforms the other on this benchmark. That said, if the main evidence for LARS’s utility as a  
 216 “large-batch optimizer” is its performance on this particular benchmark, then more evidence is needed  
 217 to quantify any benefit it has over traditional, generic optimizers like Nesterov momentum.

## 218 2.2 Lessons learned

219 In hindsight, it was only necessary to make a few changes to the LARS pipeline to match its  
 220 performance at batch size 32,768 with Nesterov momentum. However, Table 1 does not accurately  
 221 represent the effort required when attempting to match a highly tuned training-speed benchmark.

	Nesterov	LARS
$p_{\text{warmup}}$	2	1
$\eta_{\text{peak}}$	7.05	29.0
$\eta_{\text{final}}$	$6 \times 10^{-6}$	$10^{-4}$
$1 - \mu$	0.02397	0.071
$\lambda$	$5.8 \times 10^{-5}$	$10^{-4}$
$\tau$	0.15	0.10
$\gamma_0$	0.4138	0.0

Table 1: The hyperparameters of Configuration B that differ from state-of-the-art LARS at batch size 32,768 [Kumar et al., 2019].

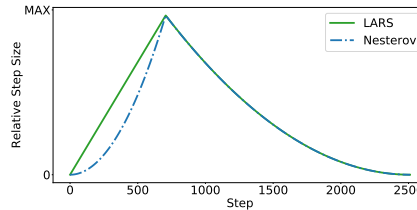


Figure 1: The learning rate schedules of LARS and Nesterov momentum Configuration B. Aside from re-scaling, the only difference is setting the warmup polynomial power to 2 instead of 1.

$p_{\text{warmup}}$	Nesterov	LARS	Optimizer	Train Acc	Test Acc
1	75.79%	75.97%	Nesterov	78.97%	75.93%
2	75.92%	75.69%	LARS	78.07%	75.97%

Table 2: **(Left)** The best warmup schedule differs for Nesterov momentum and LARS. Values are medians over 50 training runs after setting  $p_{\text{warmup}}$  without retuning other hyperparameters. **(Right)** Median train and test accuracies over 50 training runs for Nesterov momentum Configuration B and LARS.

222 Firstly, as described in Sections 2.1.2 and 2.1.3, the strong results of LARS depend partly on a few  
 223 subtle optimization tricks and non-default values of uncommonly-tuned hyperparameters. Fortunately,  
 224 in this case we could discover these tricks by examining the open-source code required for MLPerf  
 225 submissions, but machine learning research papers do not always report these important details.  
 226 Researchers can easily waste a lot of experiments and produce misleading results before getting all of  
 227 these details right. We demonstrate the importance of adding these tricks to our Nesterov momentum  
 228 pipeline in Appendix C; without these tricks (or some new tricks), we likely would not have been  
 229 able to match the LARS performance.

230 Secondly, the learning rate schedule really matters when trying to maximize performance with a  
 231 relatively small step budget. Both LARS and Nesterov momentum are sensitive to small deviations  
 232 from the optimized learning rate schedules in Figure 1, and neither schedule works as well for the  
 233 other optimizer. Although relatively minor changes were sufficient to match LARS with Nesterov  
 234 momentum, there is no way to know *a priori* how the optimal schedule will look for a new optimizer  
 235 Wu et al. [2018]. Even in toy settings where the optimal learning rate schedule can be derived, it  
 236 does not fit into commonly used schedule families and depends strongly on the optimizer Zhang  
 237 et al. [2019]. Indeed, this problem applies to the other optimization hyperparameters as well: it  
 238 is extremely difficult to know which are worth considering ahead of time. Finally, even when we  
 239 narrowed down our hyperparameter search spaces around the optimal point, the volume of our search  
 240 spaces corresponding to near-peak performance was small, likely due to the small step budget [Shallue  
 241 et al., 2019]. We investigate how these effects change with a less stringent step budget in Section 4.

### 242 3 Stronger BERT pretraining speed baselines

243 You et al. [2019] developed the LAMB optimizer in the hope of speeding up training for BERT-Large  
 244 [Bidirectional Encoder Representations from Transformers, Devlin et al., 2018]. BERT training  
 245 consists of two phases. The “pretraining” phase has two objectives: (1) predicting masked tokens  
 246 based on the rest of the sequence (a masked language model), and (2) predicting whether two  
 247 given sentences follow one from another. Finally, the “fine-tuning” phase refines the model for a  
 248 downstream task of interest. BERT pretraining takes a considerable amount of time (up to 3 days on  
 249 16 Cloud TPU-v3 chips Jouppi et al. [2017]), whereas the fine-tuning phase is typically much faster.  
 250 Model quality is typically assessed on the downstream metrics, not on pretraining loss, making BERT  
 251 training a somewhat awkward benchmark for optimization research.

252 You et al. [2019] used LAMB for BERT pretraining with batch sizes up to 65,536 and claimed that  
 253 LAMB outperforms Adam batch size 16,384 and beyond. The LAMB optimizer has since appeared  
 254 in several NLP toolkits, including as Microsoft DeepSpeed and NVIDIA Multi-node BERT training,  
 255 and as a benchmark task in MLPerf v0.7.<sup>12</sup>

256 As shown in Table 3, we trained Adam (with decoupled weight decay) baselines that achieve better  
 257 results than both the LAMB and Adam results reported in You et al. [2019]. Our new Adam  
 258 baselines obtain better F1 scores on the development set of the SQuAD v1.1 task in the same number  
 259 of training steps as LAMB for both batch size 32,768 and the hybrid 65,536-then-32,768 batch  
 260 size training regime in You et al. [2019]. We also ran Adam at batch size 65,536 to reach nearly  
 261 the same F1 score as the hybrid batch size LAMB result, but in much fewer training steps. We  
 262 believe 7,818 steps is a new state-of-the-art for BERT pretraining speed [in our experiments, it  
 263 also improves upon the 76-minute record claimed in You et al., 2019]. Additionally, at batch  
 264 size 32,768 our Adam baseline got a better pretraining loss of 1.277 compared to LAMB’s 1.342.

<sup>12</sup> We do not consider the MLPerf task in this paper since it is a warm-start, partial training task.

265

266 We used the same experimental setup as You  
 267 et al. [2019], including two pretraining phases  
 268 with max sequence lengths of 128 and then 512.

269 In order to match You et al. [2019], we reported  
 270 the F1 score on the downstream SQuAD v1.1  
 271 task as the target metric, although this metric  
 272 introduces potential confounds: optimization  
 273 efficiency should be measured on the training

274 task using training and held-out data sets. Fortunately, in this case better pretraining performance  
 275 correlated a with higher F1 score after fine-tuning. See Appendix B.2 for additional experiment  
 276 details. We tuned Adam hyperparameters independently for each pretraining phase, specifically  
 277 learning rate  $\eta$ ,  $\beta_1$ ,  $\beta_2$ , the polynomial power for the learning rate warmup  $p_{warmup}$ , and weight  
 278 decay  $\lambda$ , using quasi-random search [Bousquet et al., 2017]. See Appendix D.2 for the search spaces.

279 In addition to hyperparameter tuning, our improved Adam results at these batch sizes are also likely  
 280 due to two implementation differences. First, the Adam implementation in You et al. [2019] comes  
 281 from the BERT open source code base, in which Adam is missing the standard bias correction.<sup>13</sup>  
 282 The Adam bias correction acts as an additional step size warm-up, thereby potentially improving the  
 283 stability in the initial steps of training. Second, the BERT learning rate schedule had a discontinuity  
 284 at the start of the decay phase due to the learning rate decay being incorrectly applied during warm-up  
 285 <sup>14</sup> (see Figure 2 in Appendix B). This peculiarity is part of the official BERT release and is present in  
 286 3000+ copies of the BERT Training code on GitHub.

## 287 4 Investigating a less stringent step budget

288 Part of what makes comparing optimizers so difficult is that the hyperparameter tuning tends to  
 289 dominate the comparisons [Choi et al., 2019]. Moreover, tuning becomes especially difficult when  
 290 we demand a fixed epoch budget even when dramatically increasing the batch size [Shallue et al.,  
 291 2019]. Fixing the epoch budget as the batch size increases is equivalent to demanding perfect scaling  
 292 (i.e. that the number of training steps decreases by the same factor that the batch size is increased).  
 293 We can view the role of hyperparameter tuning for large batch training as resisting the inevitable end  
 294 of perfect scaling. For example, it might be possible to extend perfect scaling using delicately tuned  
 295 learning rate schedules, but comparing optimizers under these conditions can make the learning rate  
 296 schedule dominate the comparison by favoring some algorithms over others. Therefore, in order to  
 297 better understand the behavior of LARS and LAMB compared to Nesterov Momentum and Adam, we  
 298 ran additional ResNet-50 experiments with a more generous 6,000 step budget (vs 2,512 in Section 2)  
 299 and a more simplistic cosine learning rate schedule. At batch size 32,768, this budget should let us  
 300 reach better validation accuracy than the MLPerf target of 75.9%.

301 Although not mentioned in You et al. [2017], the state-of-the-art MLPerf pipeline for “LARS” actually  
 302 uses both LARS and Heavy-ball Momentum, with Momentum applied to the batch normalization and  
 303 ResNet bias parameters and LARS applied to the other parameters. You et al. [2019] does not mention  
 304 whether LAMB was only applied to some parameters and not others. If layerwise normalization can  
 305 be harmful for some model parameters, this is critical information for practitioners using LARS or  
 306 LAMB, since it might not be obvious which optimizer to apply to which parameters. To investigate  
 307 this, we trained both pure LARS and LAMB configurations, as well as configurations that did not  
 308 apply layerwise normalization to the batch normalization and ResNet bias parameters. Moreover,  
 309 LAMB’s underlying Adam implementation defaults to  $\epsilon = 10^{-6}$ , rather than the typical  $10^{-7}$  or  
 310  $10^{-8}$ . In some cases,  $\epsilon$  can be a critical hyperparameter for Adam [Choi et al., 2019], so we included  
 311 Adam configurations with both  $\epsilon = 10^{-6}$  and  $\epsilon = 10^{-8}$ .

312 Table 4 shows the validation accuracy of these different configurations after training for 6,000  
 313 steps with batch size 32,768. In every case, we used a simple cosine decay learning rate sched-  
 314 ule and tuned the initial learning rate and weight decay using quasi-random search. We used  
 315 momentum parameters of 0.98 for Nesterov momentum and 0.929 for LARS, respectively, based  
 316 on the tuned values from Section 2. We used default hyperparameters for Adam and LAMB  
 317 except where specified. We set all other hyperparameters to the same values as the state-of-the-  
 318 art LARS pipeline, except we set  $\gamma_0 = 1.0$ . See Appendix D.3 for more details. As expected,

Batch size	Step budget	LAMB	Adam
32k	15,625	91.48	<b>91.58</b>
65k/32k	8,599	90.58	<b>91.04</b>
65k	7,818	–	<b>90.46</b>

Table 3: Using Adam for pretraining exceeds the reported performance of LAMB in You et al. [2019] in terms of F1 score on the downstream SQuAD v1.1 task.

<sup>13</sup> <https://git.io/JtY8d> <sup>14</sup> See <https://git.io/JtnQW> and <https://git.io/JtnQ8>.

319 highly tuned learning rate schedules and optimizer hyperparameters are no longer necessary with  
 320 a less stringent step budget. Multiple optimizer configurations in Table 4 exceed the MLPerf  
 321 target accuracy of 75.9% at batch size 32,768 with minimal tuning. Training with larger batch  
 322 sizes is *not* fundamentally unstable: stringent step budgets make hyperparameter tuning trickier.  
 323

Weights Optimizer	Bias/BN Optimizer	Top-1
Nesterov	Nesterov	76.7
LARS	Momentum	76.9
LARS	LARS	76.9
Adam ( $\epsilon = 10^{-8}$ )	Adam ( $\epsilon = 10^{-8}$ )	76.2
Adam ( $\epsilon = 10^{-6}$ )	Adam ( $\epsilon = 10^{-6}$ )	76.4
LAMB	LAMB	27.3
LAMB	Adam ( $\epsilon = 10^{-8}$ )	76.3
LAMB	Adam ( $\epsilon = 10^{-6}$ )	76.3

337 Table 4: Validation accuracy of ResNet-50 on ImageNet trained for 6,000 steps instead of 2,512. The  
 338 second column is the optimizer that was applied to the batch norm and ResNet bias variables. We  
 339 report the median top-1 accuracy over 5 seeds of the best hyperparameter setting in a refined search  
 340 space. See Appendix D.3 for details.  
 341  
 342  
 343  
 344

345 which optimizers are easiest to tune, experiments like these that operate outside the regime of highly  
 346 tuned learning rate schedules can serve as a starting point. In this experiment, LARS and LAMB do  
 347 not appear to have an advantage in how easy they are to tune even on a dataset and model that were  
 348 used in the development of both of those algorithms. LAMB is a variant of Adam and performs about  
 349 the same as Adam with the same value of  $\epsilon$ ; LARS is more analogous to Momentum and indeed  
 350 Nesterov momentum and LARS have similar performance.

## 351 5 Discussion

352 Our results show that standard, generic optimizers suffice for achieving strong results across batch  
 353 sizes. Therefore, any research program to create new optimizers for training at larger batch sizes  
 354 must start from the fact that Momentum, Adam, and likely other standard methods work fine at batch  
 355 sizes as large as those considered in this paper. The LARS and LAMB update rules have no more  
 356 to do with the batch size (or “large” batches) than the Momentum or Adam update rules. Although  
 357 You et al. [2019] presented convergence rate bounds for LARS and LAMB to support their claims  
 358 of superior performance, we show in Appendix A that Adam satisfies a similar bound to LAMB.  
 359 These bounds all rely on very unrealistic assumptions.<sup>15</sup> Most of all, they are loose upper bounds  
 360 on the worst case behavior of the algorithms, not accurate reflections of optimizer performance in  
 361 reality. Whether layer-wise normalization can be useful for optimization or regularization remains an  
 362 open question. However, if LARS and LAMB have any advantage over standard techniques, it is not  
 363 that they work dramatically better on the tasks and batch sizes in You et al. [2017, 2019]. This is  
 364 not to suggest that there is nothing interesting about studying neural network optimization at larger  
 365 batch sizes. For example, as gradient noise decreases, there may be opportunities to harness curvature  
 366 information and extend the region of perfect scaling [Zhang et al., 2019]. However, there is currently  
 367 no evidence that LARS and LAMB scale better than Momentum and Adam.

368 Our primary concern in this paper has been matching the state of the art—and establishing new  
 369 baselines—for *training speed* measurements of the sort used to justify new techniques and algorithms  
 370 for training with larger batch sizes. In contrast, many practitioners are more concerned with obtaining  
 371 the best possible validation error with a somewhat flexible training time budget. Part of the reason  
 372 why matching LARS at batch size 32,768 was non-trivial is because getting state of the art training

<sup>15</sup> All convergence bounds assume no momentum is used, and the  $L_{avg}$  bound for LAMB also assumes  $\beta_2 = 0$ , when it is typically 0.999. Additionally,  $L_{avg}$  could still be large if  $L_\infty$  is large, but we leave an empirical analysis of this to future work.

In Table 4, “pure LAMB” performs extremely poorly: LAMB only obtains reasonable results when it is *not* used on the batch normalization and ResNet bias parameters, suggesting that layerwise normalization can indeed be harmful on some parameters. “Pure LARS” and Nesterov momentum perform roughly the same at this step budget, but the MLPerf LARS pipeline, which is tuned for a more stringent step budget, does not use LARS on all parameters, at least suggesting that the optimal choice could be budget-dependent.

Many new neural net optimizers, including LAMB, are introduced alongside claims that the new optimizer does not require any—or at least minimal—tuning. Unfortunately, these claims require a lot of work to support, since they require trying the optimizer on new problems without using those problems during the development of the algorithm. Although our experiments here are not sufficient to determine



373 speed requires several tricks and implementation details that are not often discussed. It was not  
374 obvious to us *a priori* which ones would prove crucial. These details do not involve changes to the  
375 optimizer, but they interact with the optimizer in a regime where all hyperparameters need to be well  
376 tuned to stay competitive, making it necessary to re-tune everything for a new optimizer.

377 In neural network optimization research, training loss is rarely discussed in detail and evaluation  
378 centers on validation/test performance since that is what practitioners care most about. However,  
379 although we shouldn't *only* consider training loss, it is counter-intuitive and counter-productive to  
380 elide a careful investigation of the actual objective of the optimizer. If a new optimizer achieves better  
381 test performance, but shows no speedup on training loss, then perhaps it is *not* a better optimizer so  
382 much as an indirect regularizer.<sup>16</sup> Indeed, in our experiments we found that Nesterov momentum  
383 achieves noticeably better training accuracy on ResNet-50 than the LARS configuration we used,  
384 despite reaching roughly the same validation accuracy. Properly disentangling possible regularization  
385 benefits from optimization speed-ups is crucial if we are to understand neural network training,  
386 especially at larger batch sizes where we lose some of the regularization effect of gradient noise.  
387 Hypothetically, if the primary benefit of a training procedure is regularization, then it would be better  
388 to compare the method with other regularization baselines than other optimizers.

389 Ultimately, we only care about batch size to the extent that higher degrees of data parallelism lead  
390 to faster training. Training with a larger batch size is a means, not the end goal. New optimizers—  
391 whether designed for generic batch sizes or larger batch sizes—have the potential to dramatically  
392 improve algorithmic efficiency across multiple workloads, but our results show that standard opti-  
393 mizers can match the performance of newer alternatives on the workloads we considered. Indeed,  
394 despite the legion of new update rule variants being proposed in the literature, standard Adam and  
395 Momentum remain the workhorses of practitioners and researchers alike, while independent empirical  
396 comparisons consistently find no clear winner when optimizers are compared across a variety of  
397 workloads [Schmidt et al., 2020]. Meanwhile, as Choi et al. [2019] and our results underscore,  
398 comparisons between optimizers crucially depend on the effort spent tuning hyperparameters for each  
399 optimizer. Given these facts, we should regard with extreme caution studies claiming to show the  
400 superiority of one particular optimizer over others. Part of the issue stems from current incentives in  
401 the research community; we overvalue the novelty of new methods and undervalue establishing strong  
402 baselines to measure progress against. This is particularly problematic in the study of optimizers,  
403 where the learning rate schedule is arguably more important than the choice of the optimizer update  
404 rule itself! As our results show, the best learning rate schedule is tightly coupled with the optimizer,  
405 meaning that tuning the learning rate schedule for a new optimizer will generally favor the new  
406 optimizer over a baseline unless the schedule of the baseline is afforded the same tuning effort.

## 407 6 Conclusion

408 In this work, we demonstrated that standard optimizers, without any layer-wise normalization  
409 techniques, can match or exceed the large batch size results used to justify LARS and LAMB. Future  
410 work attempting to argue that a new algorithm is useful by comparing to baseline methods or results,  
411 including those established in this paper, faces a key challenge in showing that the gains are due to the  
412 new method and not merely due to better tuning or changes to the training pipeline (e.g. regularization  
413 tricks). Although gains from tuning will eventually saturate, we can, in principle, always invest more  
414 effort in tuning and potentially get better results for any optimizer. However, our goal should be  
415 developing optimizers that work better across many different workloads when taking into account the  
416 amount of additional tuning they require.

417 Moving forward, if we are to reliably make progress we need to rethink how we compare and evaluate  
418 new optimizers for neural network training. Given how sensitive optimizer performance is to the  
419 hyperparameter tuning protocol and how difficult it is to quantify hyperparameter tuning effort, we  
420 can't expect experiments with self-reported baselines to always lead to fair comparisons. Ideally, new  
421 training methods would be evaluated in a standardized competitive benchmark, where submitters of  
422 new optimizers do not have full knowledge of the evaluation workloads. Some efforts in this direction  
423 have started, for instance the MLCommons Algorithmic Efficiency Working Group<sup>17</sup>, but more work  
424 needs to be done to produce incentives for the community to publish well-tuned baselines and to  
425 reward researchers that conduct the most rigorous empirical comparisons.

<sup>16</sup> Deep learning folk wisdom is that “any method to make training less effective can serve as a regularizer,” whether it is a bug in gradients or a clever algorithm.

<sup>17</sup> <https://mlcommons.org/en/groups/research-algorithms/>

## 426 Checklist

- 427 1. For all authors...
- 428 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
429 contributions and scope? [Yes] See Sections 2, 3, 4
- 430 (b) Did you describe the limitations of your work? [Yes] We had a lengthy discussion of  
431 the limitations and scope of the work in Section 5
- 432 (c) Did you discuss any potential negative societal impacts of your work? [No] We did  
433 not discuss this in the main text. Our primary contribution is to improve experimental  
434 protocols for other methodological work, which is so removed from specific machine  
435 learning applications that it is hard to determine the net impact. That said, more  
436 effective experimental protocols should lead to more effective science which in turn  
437 should lead to more effective machine learning applications. Whether this development  
438 is positive or negative for society will depend on who stands to gain from the use of  
439 machine learning in future applied contexts. Additionally, although our work should, in  
440 the long run, save computational resources for individual researchers, in net across the  
441 community this may or may not produce an aggregate savings because more efficient  
442 machine learning training, by making larger scale projects more accessible, can lead  
443 to an increased demand for compute resources [York, 2006], which can have varying  
444 degrees of negative environmental impacts [Patterson et al., 2021].
- 445 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
446 them? [Yes]
- 447 2. If you are including theoretical results...
- 448 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Appendix A  
449 for a comprehensive description of the problem setting.
- 450 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.
- 451 3. If you ran experiments...
- 452 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
453 mental results (either in the supplemental material or as a URL)? [No] We will include  
454 a link to all code and all possible reproducibility instructions after the anonymized  
455 reviewing period is over.
- 456 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
457 were chosen)? [Yes] We are extremely detailed about our tuning procedures and dataset  
458 details, see Appendices B, D.
- 459 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
460 ments multiple times)? [Yes] While we do not report error bars in the tables in the main  
461 text, Appendices B.2, C contains box plots showing the quartiles of the distribution  
462 over random seeds.
- 463 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
464 of GPUs, internal cluster, or cloud provider)? [No] In Appendix B we state that we  
465 run on Google TPUs, however we do not tally up the total number of experiments run  
466 (although an interested reader could compute it from the information we provided in  
467 our detailed appendices given that we list all intermediate experiments, no matter how  
468 silly in hindsight).
- 469 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 470 (a) If your work uses existing assets, did you cite the creators? [Yes] We reference the  
471 relevant citations for all models, datasets, and techniques.
- 472 (b) Did you mention the license of the assets? [No]
- 473 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 474 (d) Did you discuss whether and how consent was obtained from people whose data you're  
475 using/curating? [N/A]
- 476 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
477 information or offensive content? [N/A]
- 478 5. If you used crowdsourcing or conducted research with human subjects...

- 479 (a) Did you include the full text of instructions given to participants and screenshots, if  
480 applicable? [N/A]
- 481 (b) Did you describe any potential participant risks, with links to Institutional Review  
482 Board (IRB) approvals, if applicable? [N/A]
- 483 (c) Did you include the estimated hourly wage paid to participants and the total amount  
484 spent on participant compensation? [N/A]

## 485 References

- 486 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S.  
487 Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew  
488 Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath  
489 Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah,  
490 Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent  
491 Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg,  
492 Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on  
493 heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from  
494 tensorflow.org.
- 495 Naman Agarwal, Rohan Anil, Elad Hazan, Tomer Koren, and Cyril Zhang. Disentangling adaptive  
496 gradient methods from learning rates. *arXiv preprint arXiv:2002.11803*, 2020.
- 497 Olivier Bousquet, Sylvain Gelly, Karol Kurach, Olivier Teytaud, and Damien Vincent. Critical hyper-  
498 parameters: No random, no cry. *arXiv*, 2017. URL <https://arxiv.org/abs/1706.03200>.
- 499 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal  
500 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and  
501 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL  
502 <http://github.com/google/jax>.
- 503 Dami Choi, Christopher J Shallue, Zachary Nado, Jaehoon Lee, Chris J Maddison, and George E  
504 Dahl. On empirical comparisons of optimizers for deep learning. *arXiv preprint arXiv:1910.05446*,  
505 2019.
- 506 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
507 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 508 Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the gener-  
509 alization gap in large batch training of neural networks. *arXiv preprint arXiv:1705.08741*,  
510 2017.
- 511 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by  
512 reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- 513 Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa,  
514 Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. In-datacenter performance analysis of  
515 a tensor processing unit. In *Proceedings of the 44th Annual International Symposium on Computer*  
516 *Architecture*, pages 1–12, 2017.
- 517 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
518 *arXiv:1412.6980*, 2014.
- 519 Sameer Kumar, Victor Bitorff, Dehao Chen, Chiachen Chou, Blake Hechtman, HyoukJoong Lee,  
520 Naveen Kumar, Peter Mattson, Shibo Wang, Tao Wang, et al. Scale mlperf-0.6 models on google  
521 tpu-v3 pods. *arXiv preprint arXiv:1909.09756*, 2019.
- 522 Peter Mattson, Christine Cheng, Cody Coleman, Greg Diamos, Paulius Micikevicius, David Patterson,  
523 Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, David Brooks, Dehao Chen, Debojy-  
524 oti Dutta, Udit Gupta, Kim Hazelwood, Andrew Hock, Xinyuan Huang, Atsushi Ike, Bill Jia,  
525 Daniel Kang, David Kanter, Naveen Kumar, Jeffery Liao, Guokai Ma, Deepak Narayanan, Tayo  
526 Oguntebi, Gennady Pekhimenko, Lillian Pentecost, Vijay Janapa Reddi, Taylor Robie, Tom St.  
527 John, Tsuguchika Tabaru, Carole-Jean Wu, Lingjie Xu, Masafumi Yamazaki, Cliff Young, and

- 528 Matei Zaharia. MLPerf training benchmark. *arXiv preprint arXiv:1910.01500*, 2019. URL  
529 <https://arxiv.org/abs/1910.01500>.
- 530 Yurii E Nesterov. A method for solving the convex programming problem with convergence rate  
531  $O(1/k^2)$ . In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983.
- 532 David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild,  
533 David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training. *arXiv*  
534 *preprint arXiv:2104.10350*, 2021.
- 535 Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR*  
536 *Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- 537 Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley–  
538 benchmarking deep learning optimizers. *arXiv preprint arXiv:2007.01547*, 2020.
- 539 Christopher J Shallue, Jaehoon Lee, Joseph Antognini, Jascha Sohl-Dickstein, Roy Frostig, and  
540 George E Dahl. Measuring the effects of data parallelism on neural network training. *Journal of*  
541 *Machine Learning Research*, 20(112):1–49, 2019.
- 542 Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization  
543 and momentum in deep learning. In *ICML*, 2013.
- 544 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking  
545 the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer*  
546 *vision and pattern recognition*, pages 2818–2826, 2016.
- 547 Yu Emma Wang, Gu-Yeon Wei, and David Brooks. Benchmarking tpu, gpu, and cpu platforms for  
548 deep learning. *arXiv preprint arXiv:1907.10701*, 2019.
- 549 Yuhuai Wu, Mengye Ren, Renjie Liao, and Roger Grosse. Understanding short-horizon bias in  
550 stochastic meta-optimization. *arXiv preprint arXiv:1803.02021*, 2018.
- 551 Chris Ying, Sameer Kumar, Dehao Chen, Tao Wang, and Youlong Cheng. Image classification at  
552 supercomputer scale. *arXiv preprint arXiv:1811.06992*, 2018.
- 553 Richard York. Ecological paradoxes: William stanley jevons and the paperless office. *Human Ecology*  
554 *Review*, pages 143–147, 2006.
- 555 Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv*  
556 *preprint arXiv:1708.03888*, 2017.
- 557 Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan  
558 Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep  
559 learning: Training bert in 76 minutes. In *International Conference on Learning Representations*,  
560 2019.
- 561 Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George Dahl, Chris  
562 Shallue, and Roger B Grosse. Which algorithmic choices matter at which batch sizes? insights  
563 from a noisy quadratic model. In *Advances in Neural Information Processing Systems*, pages  
564 8196–8207, 2019.
- 565 Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and  
566 Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching  
567 movies and reading books. In *Proceedings of the 2015 IEEE International Conference on Computer*  
568 *Vision (ICCV)*, ICCV '15, page 19–27, USA, 2015. IEEE Computer Society. ISBN 9781467383912.  
569 doi: 10.1109/ICCV.2015.11. URL <https://doi.org/10.1109/ICCV.2015.11>.