

570 **A Convergence Proofs**

571 As an exercise to support our larger point about the irrelevance of this type of result when comparing
 572 LAMB and Adam, below we derive a convergence bound for Adam in a similar manner to the
 573 LAMB optimizer convergence bound in You et al. [2019]. We note that all of these bounds are loose
 574 upper bounds on the worst case behavior of the algorithms, so there is no reason that comparing
 575 them reflects the relative behaviors of optimizers in reality. For example in Equation 3 below, we
 576 follow similar operations as the LAMB bound derivation and simply switch a $-$ to a $+$ for algebraic
 577 convenience.

578 We define the following as our optimization objective

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{s \in \mathbb{P}}[\ell(x, s)] + \frac{\lambda}{2} \|x\|^2, \quad (1)$$

579 with an optimal solution(s) x^* . $x \in \mathbb{R}^d$ are the neural network parameters, ℓ a smooth and possibly
 580 nonconvex loss function, \mathbb{P} a data distribution, and λ the regularization strength.

581 Let T be the number of training steps, h the number of neural network layers, b the batch size, η the
 582 learning rate, n the mini-batch size, and $\phi(v) : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ a function that is layerwise multiplied by
 583 the learning rate in LARS and LAMB updates. Let L be a vector of the layerwise Lipschitz constants
 584 for the neural network, and L_{avg} the mean of L . Let s be a training step uniformly sampled from
 585 $\{1, 2, \dots, T\}$.

586 We define the stochastic minibatch estimate of the true gradient as $\mathbb{E}[g^{(i)}] = \nabla_i f(x)$ and assume that
 587 its variance is bounded by $\mathbb{E}[g^{(i)} - \nabla_i f(x)]^2 \leq \sigma_i^2$ layerwise for a vector of standard deviations
 588 $\sigma := [\sigma^{(1)}, \dots, \sigma^{(h)}]$ and elementwise for $\tilde{\sigma} := [\tilde{\sigma}^{(1)}, \dots, \tilde{\sigma}^{(h)}]$.

589 Next, let $\eta_t = \eta = \sqrt{\frac{2(f(x_1) - f(x^*))}{\alpha_u^2 \|L\|_1 T}} \forall t \in [T]$, $b = T$, $\alpha_l \leq \phi(v) \leq \alpha_u \forall v > 0$, $\alpha_l, \alpha_u > 0$.
 590 Crucially, additionally let $b = T$, $\beta_1 = 0$, $\lambda = 0$. Under these conditions You et al. [2019] show the
 591 convergence rate for LARS is

$$\left(\mathbb{E} \left[\frac{1}{\sqrt{h}} \sum_{i=1}^h \|\nabla_i f(x_s)\| \right] \right)^2 \leq \mathcal{O} \left(\frac{(f(x_1) - f(x^*)) L_{avg}}{T} + \frac{\|\sigma\|_1^2}{Th} \right).$$

592 They also derive the convergence rate of LAMB as

$$\mathbb{E}[\|\nabla f(x_a)\|^2] \leq \mathcal{O} \left(\sqrt{\frac{G^2 d}{h(1 - \beta_2)}} \times \left[\sqrt{\frac{2(f(x_1) - f(x^*)) \|L\|_1}{T}} + \frac{\|\tilde{\sigma}\|_1}{\sqrt{T}} \right] \right).$$

593 Additionally, for $\beta_2 = 0$, the convergence rate of LAMB can be derived as

$$\left(\mathbb{E} \left[\frac{1}{\sqrt{d}} \|\nabla f(x_a)\|_1 \right] \right)^2 \leq \mathcal{O} \left(\frac{(f(x_1) - f(x^*)) L_{avg}}{T} + \frac{\|\tilde{\sigma}\|_1^2}{Th} \right),$$

594 Below we derive a similar bound for the $\beta_2 > 0$ case for Adam updates. We note that the $\beta_2 = 0$
 595 case where the bound depends on L_{avg} instead of $\|L\|_1$ can be very similarly derived for Adam, but
 596 is also a very unrealistic condition in practice.

597 *Proof.* Under the assumption $\beta_1 = 0$, $\lambda = 0$, one could write the Adam update rule as follows:

$$x_{t+1}^{(i)} = x_t^{(i)} - \eta_t \sqrt{1 - \beta_2^t} \frac{g_t^{(i)}}{\sqrt{v_t^{(i)}}},$$

598 where $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ for all $i \in [h]$.

599 Since the function f is L -smooth, we have the following:

$$\begin{aligned}
f(x_{t+1}) &\leq f(x_t) + \langle \nabla_i f(x_t), x_{t+1}^{(i)} - x_t^{(i)} \rangle + \sum_{i=1}^h \frac{L_i}{2} \|x_{t+1}^{(i)} - x_t^{(i)}\|^2 \\
&\leq f(x_t) - \underbrace{\eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} [\nabla_i f(x_t)]_j \sqrt{1 - \beta_2^t} \frac{g_{t,j}^{(i)}}{\sqrt{v_{t,j}^{(i)}}}}_{T_1} + \sum_{i=1}^h \frac{L_i \eta_t^2 d}{2(1 - \beta_2)} \quad (2)
\end{aligned}$$

600 Where the last term comes from the fact that $1 - \beta_2^t \leq 1$. We bound term T_1 in the following manner,
601 in line with [You et al., 2019]:

$$\begin{aligned}
T_1 &= -\eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} [\nabla_i f(x_t)]_j \sqrt{1 - \beta_2^t} \frac{g_{t,j}^{(i)}}{\sqrt{v_{t,j}^{(i)}}} \\
&\leq -\eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} \frac{\sqrt{1 - \beta_2^t}}{G} [\nabla_i f(x_t)]_j g_{t,j}^{(i)} \\
&\quad - \eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} \left([\nabla_i f(x_t)]_j \sqrt{1 - \beta_2^t} \frac{g_{t,j}^{(i)}}{\sqrt{v_{t,j}^{(i)}}} \right) \mathbf{1}(\text{sign}([\nabla_i f(x_t)]_j) \neq \text{sign}(g_{t,j}^{(i)}))
\end{aligned}$$

602 Relying on the following inequalities: $\sqrt{v_t} \leq G$ and $1 - \beta_2^t > 1 - \beta_2$.
603 Taking expectation, we have the following:

$$\begin{aligned}
\mathbb{E}[T_1] &\leq -\eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} \frac{\sqrt{1 - \beta_2^t}}{G} \mathbb{E} \left[[\nabla_i f(x_t)]_j g_{t,j}^{(i)} \right] \\
&\quad - \eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} \frac{\sqrt{1 - \beta_2^t}}{G} \mathbb{E} \left[\left([\nabla_i f(x_t)]_j g_{t,j}^{(i)} \right) \mathbf{1}(\text{sign}([\nabla_i f(x_t)]_j) \neq \text{sign}(g_{t,j}^{(i)})) \right] \\
\mathbb{E}[T_1] &\leq -\eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} \frac{\sqrt{1 - \beta_2^t}}{G} \mathbb{E} \left[[\nabla_i f(x_t)]_j g_{t,j}^{(i)} \right] \\
&\quad + \eta_t \sum_{i=1}^h \sum_{j=1}^{d_i} \sqrt{1 - \beta_2^t} \mathbb{E} \left[\left([\nabla_i f(x_t)]_j \right) \mathbb{P}(\text{sign}([\nabla_i f(x_t)]_j) \neq \text{sign}(g_{t,j}^{(i)})) \right] \quad (3)
\end{aligned}$$

604 similarly what is shown in signsgd, we bound the probability by first relaxing the condition, then
605 applying Markov's and then Jensen's inequality:

$$\begin{aligned}
\mathbb{P}(\text{sign}([\nabla_i f(x_t)]_j) \neq \text{sign}(g_{t,j}^{(i)})) &\leq \mathbb{P} \left(|[\nabla_i f(x_t)]_j - g_{t,j}^{(i)}| \geq |g_{t,j}^{(i)}| \right) \\
&\leq \frac{\mathbb{E} \left[|[\nabla_i f(x_t)]_j - g_{t,j}^{(i)}| \right]}{|[\nabla_i f(x_t)]_j|} \\
&\leq \frac{\sqrt{\mathbb{E} \left[([\nabla_i f(x_t)]_j - g_{t,j}^{(i)})^2 \right]}}{|[\nabla_i f(x_t)]_j|} \\
&= \frac{\tilde{\sigma}_{t,i}}{|[\nabla_i f(x_t)]_j|} \\
&\leq \frac{\tilde{\sigma}_i}{\sqrt{n} |[\nabla_i f(x_t)]_j|}
\end{aligned}$$

606 where the last inequality is from the fact that $\tilde{\sigma}_{t,j}$ is the minibatch variance at time t with batch size
 607 n . Substituting this into our derivation of T_1

$$\mathbb{E}[T_1] \leq -\eta_t \frac{\sqrt{1-\beta_2}}{G} \|\nabla f(x_t)\|^2 + \eta_t \sqrt{1-\beta_2} \sum_{i=1}^h \sum_{j=1}^{d_i} \frac{\tilde{\sigma}_i}{\sqrt{n}}$$

608 and replacing this with our definition of T_1 in Eq. (2) we get

$$\mathbb{E}[f(x_{t+1})] \leq f(x_t) - \eta_t \frac{\sqrt{1-\beta_2}}{G} \|\nabla f(x_t)\|^2 + \eta_t \sqrt{1-\beta_2} \frac{\|\tilde{\sigma}\|_1}{\sqrt{n}} + \frac{\|L\|_1 \eta_t^2 d}{2(1-\beta_2)}. \quad (4)$$

609 We then arrive at the final bound by summing Eq. (4) to step T and cancelling consecutive terms via
 610 the telescoping sum, followed by rearranging and then multiplying through by $\frac{G}{T\eta_t\sqrt{1-\beta_2}}$

$$\begin{aligned} \mathbb{E}[f(x_{T+1})] &\leq f(x_1) - \frac{\eta_t \sqrt{1-\beta_2}}{G} \sum_{t=1}^T \|\nabla f(x_t)\|^2 + T\eta_t \sqrt{1-\beta_2} \frac{\|\tilde{\sigma}\|_1}{\sqrt{n}} + \frac{T\|L\|_1 \eta_t^2 d}{2(1-\beta_2)}. \\ \frac{\eta_t \sqrt{1-\beta_2}}{G} \sum_{t=1}^T \|\nabla f(x_t)\|^2 &\leq f(x_1) - f(x^*) + T\eta_t \sqrt{1-\beta_2} \frac{\|\tilde{\sigma}\|_1}{\sqrt{n}} + \frac{T\|L\|_1 \eta_t^2 d}{2(1-\beta_2)} \\ \frac{1}{T} \sum_{t=1}^T \|\nabla f(x_t)\|^2 &\leq G \left(\frac{f(x_1) - f(x^*)}{T\eta_t \sqrt{1-\beta_2}} + \frac{\|\tilde{\sigma}\|_1}{\sqrt{n}} + \frac{\|L\|_1 \eta_t d}{2(1-\beta_2)^{\frac{3}{2}}} \right) \end{aligned}$$

611 Taking $\eta_t = \eta = \sqrt{\frac{2(f(x_1) - f(x^*))}{T\|L\|_1(1-\beta_2)d}}$ and letting $n = T$ as is similarly done in [You et al., 2019], we
 612 can recover a bound that, up to some constants, is similar to the bound for LAMB:

$$\begin{aligned} \mathbb{E}[\|\nabla f(x_t)\|^2] &\leq \mathcal{O} \left(G \left(\frac{f(x_1) - f(x^*)}{T\sqrt{\frac{2(f(x_1) - f(x^*))}{T\|L\|_1(1-\beta_2)d}} \sqrt{1-\beta_2}} + \frac{\|\tilde{\sigma}\|_1}{\sqrt{n}} + \frac{\|L\|_1 \sqrt{\frac{2(f(x_1) - f(x^*))}{T\|L\|_1(1-\beta_2)d}} d}{2(1-\beta_2)^{\frac{3}{2}}} \right) \right) \\ &= \mathcal{O} \left(G \left(\frac{1}{2} \sqrt{\frac{2(f(x_1) - f(x^*))\|L\|_1 d}{T}} + \frac{\|\tilde{\sigma}\|_1}{\sqrt{n}} + \frac{1}{2(1-\beta_2)^2} \sqrt{\frac{2(f(x_1) - f(x^*))\|L\|_1 d}{T}} \right) \right) \\ &= \mathcal{O} \left(G \left(1 + \frac{1}{(1-\beta_2)^2} \right) \sqrt{\frac{2(f(x_1) - f(x^*))\|L\|_1 d}{T}} + \frac{G\|\tilde{\sigma}\|_1}{\sqrt{T}} \right) \end{aligned}$$

613 □

614 B Additional experiment details

615 B.1 ResNet-50 training benchmark

616 All experiments were run on Google TPUs [Jouppi et al., 2017]. We typically trained on TPUv2-256
 617 or TPUv3-128 in order to accommodate the 32,768 batch size. The ResNet-50 experiments used Jax
 618 [Bradbury et al., 2018] using the Flax library, with code released here. The BERT experiments were
 619 run using TensorFlow [Abadi et al., 2015] version 1.15. We used the standard train/validation split
 620 from the previous literature and MLPerf competition.

621 For ImageNet, we used the following sequence of TensorFlow functions for pre-processing:¹⁸

```
622 tf.image.sample_distorted_bounding_box
623 tf.image.decode_and_crop_jpeg
624 tf.image.resize
625 tf.image.random_flip_left_right
626 tf.image.convert_image_dtype
```

¹⁸ Full code available at <https://git.io/JtgtE>

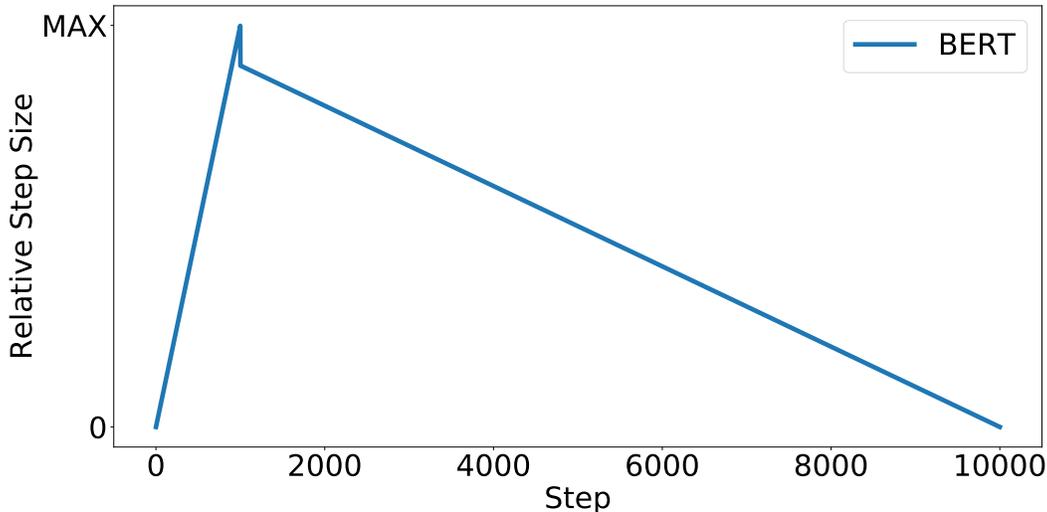


Figure 2: An illustration of the sudden drop in the BERT learning rate schedule in the official codebase.

627 B.2 BERT pre-training

628 We used the same experimental setup as the official BERT codebase¹⁹ and the standard train/test
 629 split from the previous literature. This matches the experimental setup of You et al. [2019]. We
 630 trained on Google TPUs, using TPUv3-256 or TPUv3-512 for the 32,768 batch size experiments, and
 631 TPUv3-1024 for the 65,536 batch size experiments.

632 We trained the two pretraining objectives on the combined Wikipedia and Books corpus [Zhu et al.,
 633 2015] datasets (2.5B and 800M words, respectively). We used sequence lengths of 128 and 512,
 634 respectively, for the pretraining tasks. We ran the fine-tuning phase on the SQuAD v1.1 question
 635 answering task. In order to match You et al. [2019], we report the F1-score on the dev set as the target
 636 metric. We followed the fine-tuning protocol described in the LAMB optimizer setup and did not
 637 perform any additional tuning for fine-tuning.

638 We tuned Adam hyperparameters using quasi-random search [Bousquet et al., 2017] in a simple
 639 search space. Hyperparameters included learning rate η , β_1 , β_2 , the polynomial power for the
 640 learning rate warmup p_{warmup} , and weight decay λ . We fixed the ϵ in Adam to 10^{-11} for all BERT
 641 experiments. See Appendix D.2 for the search spaces. We selected the best trial using the masked
 642 language model accuracy over 10k examples from the training set. The number of training steps for
 643 each of the phases, as well as the warmup steps are identical to You et al. [2019] and are listed in
 644 Appendix D.2. Each phase of pretraining used completely independent Adam hyperparameters. We
 645 found the final hyperparameters within 30 trials of random search for each of the phases, except for
 646 the second phase of 65,536 batch size which used 130 trials.

647 C Nesterov ablations

648 To explore the sensitivity of our best Nesterov momentum configuration (Configuration A), we ablated
 649 several elements of the experiment pipeline, one at a time, and tested their impact on performance.
 650 Figure 4 shows the results of these experiments. “Base” refers to Nesterov momentum Configuration
 651 A (Table 5). “ResNet version” is the same point as “Base” but with ResNet version 1.0 instead of
 652 version 1.5. “BN init” is the same point as “Base” but with $\gamma_0 = 1.0$ instead of 0.4138. “Virtual BN”
 653 is the same point as “Base” but with a virtual batch size of 256 instead of 64, which is the largest that
 654 fits in a single TPUv3 core. “BN & LR tuning” is Configuration B (Table 5), the same point as “Base”
 655 but with p_{decay} , t_{warmup} , η_0 , ρ , ϵ set to their values in the LARS pipeline. Finally, “L2 variables” is
 656 the same point as “Base” but where the L2 regularization is applied to all variables. The only ablation

¹⁹ <https://github.com/google-research/bert>

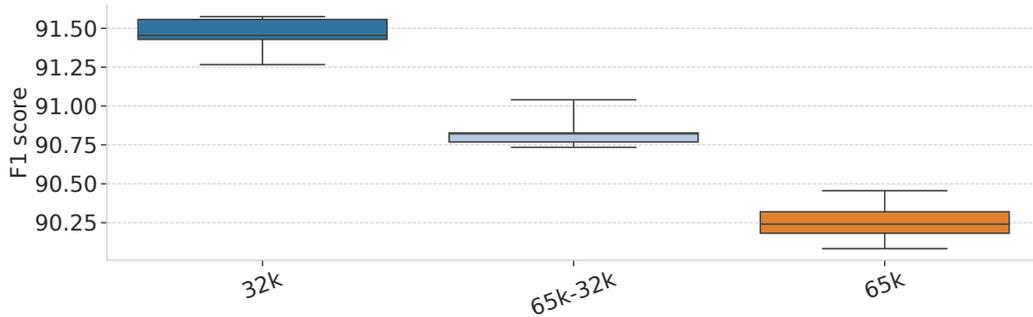


Figure 3: 6 finetuning runs starting from the same pretraining checkpoint to show the stability of our results, at each of the 32,768, mixed 65,536-32,768, and 65,536 batch size settings.

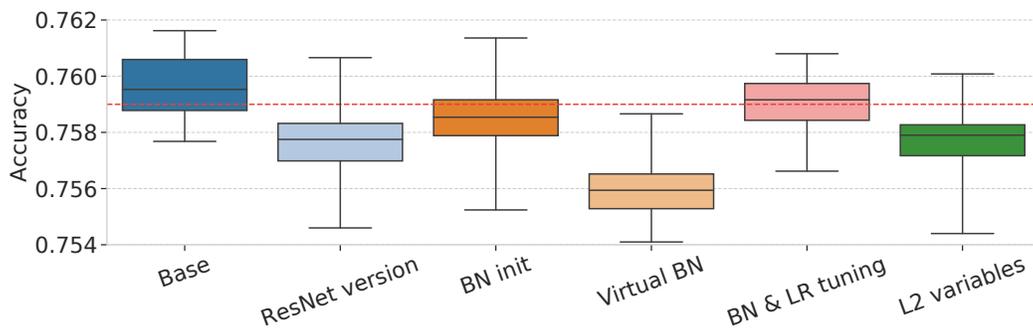


Figure 4: Distributions over 50 training runs for each ablation study around our best Nesterov momentum configuration (Configuration A). The dotted red line is at the target accuracy of 75.9%, and the boxes show the min, max, and quartiles of the distribution of accuracies over the 50 training runs.

657 whose median over 50 seeds continues to beat the target 75.9% accuracy (noted by the dotted red
 658 line) is “BN & LR tuning”, with the rest having between 0.1%-0.3% drops in median accuracy.

659 D Hyperparameter tuning

660 D.1 Nesterov momentum training speed on ResNet-50

661 We considered two configurations of Nesterov hyperparameters: Configuration A, where we tuned a
 662 wide set of hyperparameters in the experiment pipeline, and Configuration B, where we reverted the
 663 less impactful hyperparameters to the same values as the LARS baseline (or in the case of p_{warmup} ,
 664 a simpler value). We included Configuration B in order to demonstrate the minimal set of changes
 665 to the baseline necessary to still reach the target accuracy. The hyperparameter values for these
 666 configurations can be found in Table 5.

667 D.2 Adam on BERT

668 The search space used to tune Adam on BERT for all phases of the pipeline can be found in Table 6,
 669 which yielded our best Adam results on BERT in Table 7.

	Configuration A	Configuration B	LARS
t_{warmup}	638	706	706
p_{warmup}	2.497	2.0	1.0
p_{decay}	1.955	2.0	2.0
ρ	0.94	0.9	0.9
ϵ	4×10^{-6}	10^{-5}	10^{-5}
η_{peak}	7.05	7.05	29.0
η_{final}	6×10^{-6}	6×10^{-6}	10^{-4}
$1 - \mu$	0.02397	0.02397	0.071
λ	5.8×10^{-5}	5.8×10^{-5}	10^{-4}
τ	0.15	0.15	0.10
γ_0	0.4138	0.4138	0.0

Table 5: Nesterov momentum Configurations A and B.

Hyperparameter	Range	Scaling
p	{1, 2}	Discrete
η	$[10^{-5}, 1.0]$	Log
$1 - \beta_1$	$[10^{-2}, 0.5]$	Log
$1 - \beta_2$	$[10^{-2}, 0.5]$	Log
λ	$[10^{-3}, 10]$	Log

Table 6: The search space used to tune Adam on BERT for all phases of the pipeline. λ refers to weight decay and p refers to the polynomial power in the learning rate schedule for both the warmup and decay phases.

Batch size	Phase	Seq len	Warmup steps	Train steps	Learning rate	β_1	β_2	λ	p
32,768	1	128	3,125	14,063	5.9415×10^{-4}	0.934271	0.989295	0.31466	1
32,768	2	512	781	1,562	2.8464×10^{-4}	0.963567	0.952647	0.31466	1
65,536	1	128	2,000	7,037	1.3653×10^{-3}	0.952378	0.86471	0.19891	2
32,768	2	512	781	1,562	2.8464×10^{-4}	0.952647	0.963567	0.19891	2
65,536	2	512	390	781	6.1951×10^{-5}	0.65322	0.82451	0.19891	2

Table 7: Best hyperparameters from tuning Adam on BERT-Large pretraining. λ refers to weight decay and p refers to the polynomial power in the learning rate schedule for both the warmup and decay phases. All trials used $\epsilon = 10^{-11}$.

670 D.3 Less stringent step budget on ResNet-50

671 All trials used a cosine decay learning rate schedule and tuned the initial learning rate η and L2
672 regularization or weight decay parameter²⁰ λ according to Table 8. We used 50 or more trials to
673 search in the ‘‘Initial Range’’ and then 25 trials to search in the refined ‘‘Final Range.’’ Finally, we
674 ran the best point from the latter for 5 random seeds. When LARS or LAMB were used alongside
675 a different optimizer for the batch normalization and ResNet-50 bias parameters, we set $\lambda = 0$ on
676 the batch normalization and ResNet-50 bias parameters. When LAMB was used all parameters, the
677 majority of trials diverged during training – it took **67 trials** to get 25 trials that did not NaN during
678 training. Our trial budgets refer to the number of feasible trials, i.e. trials that do not diverge during
679 training.

²⁰ As suggested in You et al. [2019], we used L2 regularization for LARS and weight decay for LAMB. For consistency, we used L2 regularization for Nesterov momentum (which is more analogous to LARS) and weight decay for Adam (which is more analogous to LAMB).

Weights Optimizer	Bias/BN Optimizer	Name	Initial Range	Final Range	Best
Nesterov	Nesterov	η	np.logspace(-.5, .5, 10)	[0.8, 3]	1.173
Nesterov	Nesterov	λ	np.logspace(-4, -3, 10)	$[3 \times 10^{-4}, 10^{-3}]$	3.026×10^{-4}
LARS	Heavy-ball momentum	η	np.logspace(0, 2, 10)	[10, 40]	14.49
LARS	Heavy-ball momentum	λ	np.logspace(-5, -2, 10)	$[5 \times 10^{-5}, 2 \times 10^{-4}]$	1.708×10^{-4}
LARS	LARS	η	[1, 30]	[10, 30]	14.18
LARS	LARS	λ	$[10^{-4}, 10^{-1}]$	$[5 \times 10^{-5}, 5 \times 10^{-4}]$	5.278×10^{-5}
Adam ($\epsilon = 10^{-8}$)	Adam ($\epsilon = 10^{-8}$)	η	$[10^{-3}, 1]$	$[4 \times 10^{-3}, 2 \times 10^{-2}]$	0.004596
Adam ($\epsilon = 10^{-8}$)	Adam ($\epsilon = 10^{-8}$)	λ	$[10^{-2}, 4]$	$[2 \times 10^{-1}, 1]$	0.6182
Adam ($\epsilon = 10^{-6}$)	Adam ($\epsilon = 10^{-6}$)	η	np.logspace(-3, 0, 10)	$[3 \times 10^{-3}, 10^{-2}]$	3.332×10^{-3}
Adam ($\epsilon = 10^{-6}$)	Adam ($\epsilon = 10^{-6}$)	λ	np.logspace(-2, 0.5, 6)	[0.5, 2]	1.055
LAMB	LAMB	η	np.logspace(-4, 0, 30)	$[4 \times 10^{-3}, 5 \times 10^{-2}]$	0.01134
LAMB	LAMB	λ	np.logspace(-5, -2, 4)	$[1 \times 10^{-2}, 0.1]$	0.02657
LAMB	Adam ($\epsilon = 10^{-8}$)	η	$[10^{-3}, 1]$	$[10^{-2}, 8 \times 10^{-2}]$	0.02569
LAMB	Adam ($\epsilon = 10^{-8}$)	λ	$[10^{-2}, 4]$	[1, 8]	2.500
LAMB	Adam ($\epsilon = 10^{-6}$)	η	np.logspace(-3, 0, 10)	$[10^{-2}, 8 \times 10^{-2}]$	0.03378
LAMB	Adam ($\epsilon = 10^{-6}$)	λ	np.logspace(-2, 0.5, 6)	[1, 8]	4.197

Table 8: Search spaces used for the 6,000 step, cosine learning rate schedule experiments. All hyperparameters were tuned on a logarithmic scale, except for those which define a discrete sequence of points to evaluate such as “np.logspace”.

	Range	Scaling
η_0	$[10^{-3}, 50.0]$	Log
η_{decay_factor}	$\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$	Discrete
$1 - \mu$	$[10^{-3}, 1.0]$	Log
λ	$[10^{-5}, 10^{-1}]$	Log
τ	$[10^{-2}, 2 \times 10^{-1}]$	Linear

Table 9: First search space of the Nesterov tuning journey. The search spaces were mostly by informed guesses by the authors. λ refers to weight decay, which is applied to all variables. Tuned for 251 trials. Trained for 2,815 steps (“72 epochs” as defined by MLPerf epoch calculations). We used a linear learning rate decay schedule that decays for all training steps, starting from η_0 and ending at $\eta_0 \times \eta_{decay_factor}$. Virtual batch size 128.

680 D.4 Nesterov ResNet50 search space chronology

681 Below we list the sequence of search spaces we used to arrive at our final values in Table 5. Given
682 that the final results reported in papers are rarely found in a single iteration of experiments, we believe
683 that it is important to document the full journey to arriving at our results.

684 Note that although we tuned a wide range of hyperparameters to match the LARS result with Nesterov
685 momentum, we later realized that many of these hyperparameters could be reverted to the values
686 from the LARS pipeline (see Table 5). We started tuning with a training budget of 2,815 steps, which
687 is the number of steps in the MLPerf 0.6 submission. We sometimes would decrease this to 2,658
688 steps to test how decreasing the training budget would affect tuning performance, before eventually
689 moving to the 2,512 steps used to generate the results in the main text.

	Range	Scaling
η_0	$[10^{-3}, 50.0]$	Log
$\eta_{\text{decay_factor}}$	$\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$	Discrete
$1 - \mu$	$[10^{-3}, 1.0]$	Log
λ	$[10^{-5}, 10^{-1}]$	Log
τ	$[10^{-2}, 2 \times 10^{-1}]$	Linear

Table 10: Same as Table 9 but trained for 2,658 steps (“68 epochs” as defined by MLPerf epoch calculations) for 50 trials.

	Range	Scaling
η_0	$[10^{-1}, 20.0]$	Log
$\eta_{\text{decay_factor}}$	$\{10^{-5}, 10^{-4}, 10^{-3}\}$	Discrete
t_{decay}	$[2392, 2.658]$	Linear
$1 - \mu$	$[10^{-3}, 1.0]$	Log
λ	$[10^{-5}, 2 \times 10^{-1}]$	Log
τ	$[10^{-2}, 2 \times 10^{-1}]$	Linear

Table 11: λ refers to weight decay, which is now not applied to the bias and batch normalization variables. 50 trials. Trained for 2,658 steps. Linear learning rate decay schedule that decays for t_{decay} steps, starting from η_0 and ending at $\eta_0 \times \eta_{\text{decay_factor}}$. Virtual batch size 128.

	Range	Scaling
η_{peak}	$[10^{-1}, 32.0]$	Log
$\eta_{\text{decay_factor}}$	$\{10^{-5}, 10^{-4}, 10^{-3}\}$	Discrete
t_{decay}	$[2392, 2.658]$	Linear
$1 - \mu$	$[10^{-4}, 10^{-1}]$	Log
λ	$[10^{-4}, 10^{-1}]$	Log
τ	$[5 \times 10^{-2}, 0.15]$	Linear

Table 12: λ refers to weight decay, which is not applied to the bias and batch normalization variables. 50 trials. Trained for 2,658 steps. Linear warmup for 500 steps followed by a quadratic decay, which decays until step t_{decay} , and then is constant at the final learning rate $\eta_0 \times \eta_{\text{decay_factor}}$. Virtual batch size 128. We increased the max learning rate based off the larger learning rates used by LARS. **We also ran two additional studies which were the same except with 250 and 977 warmup steps.**

	Range	Scaling
η_{peak}	$[10^{-1}, 32.0]$	Log
$\eta_{\text{decay_factor}}$	$[3 \times 10^{-5}, 3 \times 10^{-4}]$	Log
t_{decay}	$[2533, 2.815]$	Linear
$1 - \mu$	$[10^{-4}, 10^{-1}]$	Log
λ	$[10^{-4}, 10^{-1}]$	Log
τ	$[5 \times 10^{-2}, 0.15]$	Linear

Table 13: λ refers to weight decay, which is not applied to the bias and batch normalization variables. 50 trials. Trained for 2,815 steps. Linear warmup for 500 steps followed by a quadratic decay, which decays until step t_{decay} , and then is constant at the final learning rate $\eta_0 \times \eta_{\text{decay_factor}}$. Virtual batch size 128.

	Range	Scaling
η_{peak}	$[10^{-1}, 32.0]$	Log
$\eta_{\text{decay_factor}}$	$[3 \times 10^{-5}, 3 \times 10^{-4}]$	Log
t_{decay}	$[2533, 2,815]$	Linear
$1 - \mu$	$[5 \times 10^{-3}, 10^{-1}]$	Log
λ	$[10^{-2}, 10^{-1}]$	Log
τ	$[5 \times 10^{-2}, 0.15]$	Linear

Table 14: λ refers to weight decay, which is not applied to the bias and batch normalization variables. 50 trials. Trained for 2,815 steps. Linear warmup for 500 steps followed by a quadratic decay, which decays until step t_{decay} , and then is constant at the final learning rate $\eta_0 \times \eta_{\text{decay_factor}}$. Virtual batch size 128.

	Range	Scaling
η_{peak}	$[10^{-1}, 32.0]$	Log
$\eta_{\text{decay_factor}}$	$[3 \times 10^{-5}, 3 \times 10^{-4}]$	Log
t_{decay}	$[2533, 2,815]$	Linear
$1 - \mu$	$[5 \times 10^{-3}, 10^{-1}]$	Log
λ	$[10^{-2}, 10^{-1}]$	Log
τ	$[5 \times 10^{-2}, 0.15]$	Linear

Table 15: The same as Table 14 except with virtual batch size 64.

	Range	Scaling
η_{peak}	$\{\{10^\alpha, 2 \times 10^\alpha, \dots, 9 \times 10^\alpha\} \forall \alpha \in \{-3, \dots, 2\}\} + \{100, \}$	Discrete
$\eta_{\text{decay_factor}}$	8.144×10^{-5}	–
t_{decay}	2250	–
$1 - \mu$	0.02397	–
λ	0.009992	–
τ	0.07786	–

Table 16: λ refers to weight decay, which is not applied to the bias and batch normalization variables. Trained for 2,815 steps. Virtual batch size 64. Using the best hyperparameters from Table 15, we swept over the peak learning rate in a discrete set of ten values per order of magnitude, **each for three random seeds**, to find the max stable learning rate.

	Range	Scaling
η_{peak}	4.118	–
$\eta_{\text{decay_factor}}$	8.144×10^{-5}	–
t_{decay}	2250	–
$1 - \mu$	0.02397	–
λ	$\{\{0.5 \times 10^\alpha, 10^\alpha, \dots\} \forall \alpha \in \{-3, \dots, 0\}\} + \{1.0, \}$	Discrete
τ	0.07786	–

Table 17: λ refers to weight decay, which is not applied to the bias and batch normalization variables. Trained for 2,815 steps. Virtual batch size 64. Using the best hyperparameters from Table 15, we swept over the weight decay in a discrete set of twenty values per order of magnitude, to test how high the regularization has to be in this region of hyperparameter space.

	Range	Scaling
η_{peak}	4.118	–
$\eta_{\text{decay_factor}}$	8.144×10^{-5}	–
t_{decay}	2250	–
$1 - \mu$	0.02397	–
λ	0.009992	–
τ	0.07786	–
ρ	{0.0, 0.1, 0.3, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.995, 0.999}	Discrete
ϵ	{ $10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}$ }	Discrete

Table 18: λ refers to weight decay, which is not applied to the bias and batch normalization variables. Trained for 2,815 steps. Virtual batch size 64. Using the best hyperparameters from Table 15, we swept over batch normalization hyperparameters.

	Range	Scaling
η_{peak}	[2.0, 8.0]	Log
$\eta_{\text{decay_factor}}$	$[4 \times 10^{-5}, 1.6 \times 10^{-4}]$	Linear
t_{decay}	[2100, 2400]	Linear
$1 - \mu$	[0.012, 0.04]	Log
λ	$[7 \times 10^{-3}, 7 \times 10^{-2}]$	Log
τ	[0.04, 0.1]	Linear
ρ	[0.45, 0.55]	Linear
ϵ	$[5 \times 10^{-6}, 5 \times 10^{-5}]$	Linear

Table 19: λ refers to weight decay, which is not applied to the bias and batch normalization variables. 50 trials. Trained for 2,815 steps. Linear warmup for 500 steps followed by a quadratic decay, which decays until step t_{decay} , and then is constant at the final learning rate $\eta_0 \times \eta_{\text{decay_factor}}$. Virtual batch size 64. Peak learning rate range was consolidated based off the results of Table 16. The weight decay range was consolidated based off the results of Table 17.

	Range	Scaling
t_{warmup}	[300, 800]	Linear
p_{warmup}	[0.7, 2.0]	Linear
p_{decay}	1.8	–
η_0	[0.1, 1.0]	Log
η_{peak}	[5.0, 9.0]	Log
η_{final}	$[10^{-5}, 5 \times 10^{-5}]$	Log
$1 - \mu$	0.02397	–
λ	5×10^{-5}	–
τ	0.15	–
γ_0	[0.0, 0.6]	Linear
ρ	0.94	–
ϵ	4×10^{-6}	–

Table 20: Here we switched λ to refer to L2 regularization. We also began training for 2,512 steps, which is the final “64 epochs” used in the Nesterov results reported in the main text. Because of this more stringent step budget, we focused on the learning rate schedule. t_{decay} was set to all remaining steps after the warmup was finished. Tuned for 229 trials. Virtual batch size 64.

	Range	Scaling
t_{warmup}	638	–
p_{warmup}	[1.5, 3.0]	Linear
p_{decay}	[1.5, 2.5]	Linear
η_0	0.12	–
η_{peak}	7.05	–
η_{final}	$[10^{-6}, 5 \times 10^{-4}]$	Log
$1 - \mu$	0.02397	–
λ	$[5 \times 10^{-5}, 1 \times 10^{-3}]$	Log
τ	0.15	–
γ_0	[0.4, 1.0]	Linear
ρ	0.94	–
ϵ	4×10^{-6}	–

Table 21: Here we began focusing more on the shape of the learning rate schedule, as well as retuning the L2 regularization. λ refers to L2. Several values were picked from the best trial of Table 20. Trained for 2,512 steps. Tuned for 15 trials. Virtual batch size 64.

	Range	Scaling
t_{warmup}	638	–
p_{warmup}	[1.5, 3.0]	Linear
p_{decay}	[1.5, 2.5]	Linear
η_0	0.12	–
η_{peak}	7.05	–
η_{final}	$[10^{-6}, 5 \times 10^{-4}]$	Log
$1 - \mu$	0.02397	–
λ	$[1 \times 10^{-5}, 1 \times 10^{-4}]$	Log
τ	0.15	–
γ_0	[0.4, 1.0]	Linear
ρ	0.94	–
ϵ	4×10^{-6}	–

Table 22: Here we focus in more on tuning the L2 regularization. λ refers to L2. Trained for 2,512 steps steps. Tuned for 37 trials. Virtual batch size 64.

	Range	Scaling
t_{warmup}	638	–
p_{warmup}	[1.5, 3.0]	Linear
p_{decay}	[1.5, 2.5]	Linear
η_0	0.12	–
η_{peak}	7.05	–
η_{final}	$[10^{-6}, 5 \times 10^{-4}]$	Log
$1 - \mu$	0.02397	–
λ	$[5 \times 10^{-5}, 6 \times 10^{-5}]$	Linear
τ	0.15	–
γ_0	[0.4, 1.0]	Linear
ρ	0.94	–
ϵ	4×10^{-6}	–

Table 23: Again we dial in more on a tighter tuning range for the L2 regularization. λ refers to L2. Trained for 2,512 steps steps. Tuned for 37 trials. Virtual batch size 64.