

# Benchmarking Neural Network Training Algorithms

ELLIS Workshop

**Frank Schneider**

July 19, 2023

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



MAX PLANCK INSTITUTE  
FOR INTELLIGENT SYSTEMS



imprs-is

# Neural Networks are extremely useful models...

...but!

In practice, neural networks are...

- ▶ ...**slow** to train (training can easily take days or weeks)
- ▶ ...**tedious** (demanding human intervention and fiddling)
- ▶ ...**expensive** (requiring dozens of trial runs)

# Neural Networks are extremely useful models...

...but!

In practice, neural networks are...

- ▶ ...**slow** to train (training can easily take days or weeks)
- ▶ ...**tedious** (demanding human intervention and fiddling)
- ▶ ...**expensive** (requiring dozens of trial runs)

**We all want neural network training to be faster, more automatic, and more efficient!**

# The Three Pillars of Neural Network Training

Our focus is on the (training) algorithms

## Hardware

- ▶ Leveraging accelerators (e.g. GPUs, TPUs, etc.)
- ▶ Maximizing accelerator utilization throughout training
- ▶ Reducing hardware bottlenecks

## Software

- ▶ Convenient deep learning frameworks (e.g. PYTORCH, JAX, etc.)
- ▶ Efficient software implementations
- ▶ Putting ML models into production

## Algorithms

- ▶ Efficient **training algorithms** (e.g. ADAM, SHAMPOO, etc.)
- ▶ Powerful deep learning models (e.g. TRANSFORMERS, etc.)
- ▶ Faster tuning methods (e.g. BayesOpt, etc.)





# The Current Benchmarks for Training Algorithms are Insufficient

and this holds back the entire field

- ▶ **Unclear SOTA**

There is no established protocol to train neural networks.

- ▶ **No reliable way to detect algorithmic improvements**

Let alone understand what novel methods are most promising.

- ▶ **Training algorithms are assumed to not be useful until widely adopted**

This chicken-and-egg problem is troubling for practitioners training neural networks and developers of new training algorithms.

# The Current Benchmarks for Training Algorithms are Insufficient

and this holds back the entire field

- ▶ **Unclear SOTA**

There is no established protocol to train neural networks.

- ▶ **No reliable way to detect algorithmic improvements**

Let alone understand what novel methods are most promising.

- ▶ **Training algorithms are assumed to not be useful until widely adopted**

This chicken-and-egg problem is troubling for practitioners training neural networks and developers of new training algorithms.

**We desperately need new benchmarks for neural network training algorithms!**



# MLCommons Introduces the AlgoPef: Training Algorithm Benchmark

An unprecedented effort to find faster deep learning training algorithms

## ML ● Commons Algorithms Working Group



**George Dahl**

Google



**Frank Schneider**

University of Tübingen

### **AlgoPerf: Training Algorithm Benchmark**

**A competition to measure neural network training speedups due to algorithmic changes.**

- ▶ A competitive benchmark with open submissions
- ▶ Compete on time-to-result over multiple deep learning workloads
- ▶ A huge large-scale effort by 25+ researchers from Google, University of Tübingen, University of Toronto, Meta AI, etc.

# Challenge I: What is the target?

Which algorithm trains the fastest depends on what it means for training to be complete

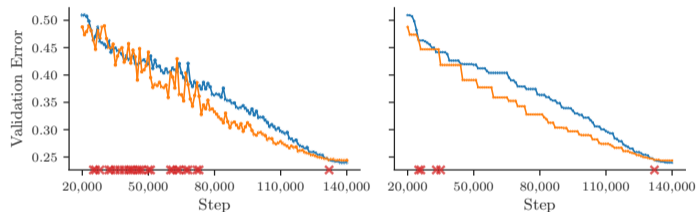


Figure 1: *Left*: Validation error for two different runs (—, —) of ADAM on RESNET-50 on IMAGENET. *Right*: The *best* validation error obtained so far. The runs intersect multiple times (✕).

# Challenge I: What is the target?

Which algorithm trains the fastest depends on what it means for training to be complete

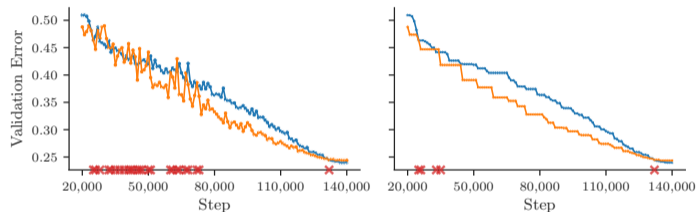


Figure 1: *Left*: Validation error for two different runs (—, —) of ADAM on RESNET-50 on IMAGENET. *Right*: The *best* validation error obtained so far. The runs intersect multiple times (✕).

- ▶ Directly comparing training curves is ill-posed
- ▶ Without defining the target in advance, we can champion any method
- ▶ We are effectively moving the goal post after the experiment



- ▶ **Defined Target Performances**

We define competitive validation and test targets for each workload that can be reliably achieved with currently popular methods.

- ▶ **A Time-To-Results Competition**

Measure the wall-clock runtime until the algorithm first hits the targets.

- ▶ **Fixed Hardware**

Submissions need to innovate on the training algorithm.

# Challenge II: Dependence on the Workload



Seemingly minor changes to the model can have a large effect on the performance of different training algorithms

Training Algorithm	RESNET-200 <i>(standard)</i>
NESTEROV	<b>0.2090</b>
ADAMW	0.2626

**Table 1:** Performance of training methods on different workloads. Shown is the best error rate (lower is better) achieved.

# Challenge II: Dependence on the Workload



Seemingly minor changes to the model can have a large effect on the performance of different training algorithms

Training Algorithm	RESNET-200 <i>(standard)</i>	RESNET-200 <i>Extra-BN</i>
NESTEROV	<b>0.2090</b>	No Feasible Trials
ADAMW	0.2626	0.2722

**Table 1:** Performance of training methods on different workloads. Shown is the best error rate (lower is better) achieved.

# Challenge II: Dependence on the Workload

Seemingly minor changes to the model can have a large effect on the performance of different training algorithms

Training Algorithm	RESNET-200 <i>(standard)</i>	RESNET-200 <i>Extra-BN</i>
NESTEROV	<b>0.2090</b>	No Feasible Trials
ADAMW	0.2626	0.2722

**Table 1:** Performance of training methods on different workloads. Shown is the best error rate (lower is better) achieved.

- ▶ How general is a performance improvement?
- ▶ Comparison with published results is dangerous
- ▶ Important to choose relevant workloads

# Solution II: Evaluate on Multiple Fixed Workloads

Finding general-purpose algorithmic improvements



- ▶ **Fixed Workloads**

We don't allow pipeline changes that are not part of the training algorithm.

- ▶ **Aggregate across Multiple Workloads**

Instead of specialized solutions, we want to find the best general-purpose method.

- ▶ **Use Randomized Workload Variants**

To test the robustness of the methods.



# Solution II: Evaluate on Multiple Fixed Workloads



Finding general-purpose algorithmic improvements

<b>Task</b>	<b>Dataset</b>	<b>Model</b>	<b>Loss</b>	<b>Metric</b>	<b>Validation Target</b>	<b>Test Target</b>	<b>Maximum Runtime</b>
Clickthrough rate prediction	<b>CRITEO 1TB</b>	<b>DLRMSMALL</b>	CE	CE	0.123 649	0.126 060	7703
MRI reconstruction	<b>FASTMRI</b>	<b>U-NET</b>	L1	SSIM	0.7344	0.741 652	8859
Image classification	<b>IMAGENET</b>	<b>RESNET-50</b>	CE	ER	0.225 69	0.3440	63 008
		<b>ViT</b>	CE	ER	0.226 91	0.3481	77 520
Speech recognition	<b>LIBRISPEECH</b>	<b>CONFORMER</b>	CTC	WER	0.078 477	0.046 973	101 780
		<b>DEEPSPEECH</b>	CTC	WER	0.1162	0.068 093	92 509
Molecular property prediction	<b>OGBG</b>	<b>GNN</b>	CE	mAP	0.280 98	0.268 729	18 477
Translation	<b>WMT 2016</b>	<b>TRANSFORMER</b>	CE	BLEU	30.8491	30.7219	48 151

# Challenge III: Hyperparameters

Often training algorithms are only templates, not runnable procedures

Search Space	Learning Rate	Weight Decay	$1 - \beta_1$	$\beta_2$
ADAMW NARROW	[2e-4, 5e-3]	[2e-2, 0.5]	0.1	0.999
ADAMW BROAD	[5e-6, 2e-2]	[5e-6, 2.0]	[1e-3, 1.0]	0.999

Table 2: Hyperparameter search spaces for two algorithms using ADAMW.

# Challenge III: Hyperparameters

Often training algorithms are only templates, not runnable procedures

Search Space	Learning Rate	Weight Decay	$1 - \beta_1$	$\beta_2$
ADAMW NARROW	[2e-4, 5e-3]	[2e-2, 0.5]	0.1	0.999
ADAMW BROAD	[5e-6, 2e-2]	[5e-6, 2.0]	[1e-3, 1.0]	0.999

Table 2: Hyperparameter search spaces for two algorithms using ADAMW.

Workload		ADAMW NARROW	ADAMW BROAD
CRITEO 1TB	DLRMSMALL	<b>0.124 01</b>	0.124 087
FASTMRI	U-NET	<b>0.734 746</b>	0.734 311
IMAGENET	RESNET-50	<b>0.232 56</b>	0.243 34
	ViT	<b>0.219 92</b>	0.236 16
LIBRISPEECH	CONFORMER	<b>0.075 989</b>	0.080 673
	DEEPSPEECH	<b>0.112 706</b>	0.120 674
OGBG	GNN	<b>0.282 14</b>	0.276 307
WMT 2016	TRANSFORMER	<b>31.3523</b>	30.9950

Table 3: Performance across multiple workloads for two ADAMW training methods. Bolded number highlights the better performance.

# Challenge III: Hyperparameters

Often training algorithms are only templates, not runnable procedures

Search Space	Learning Rate	Weight Decay	$1 - \beta_1$	$\beta_2$
ADAMW NARROW	[2e-4, 5e-3]	[2e-2, 0.5]	0.1	0.999
ADAMW BROAD	[5e-6, 2e-2]	[5e-6, 2.0]	[1e-3, 1.0]	0.999

Table 2: Hyperparameter search spaces for two algorithms using ADAMW.

Workload		ADAMW NARROW	ADAMW BROAD
CRITEO 1TB	DLRMSMALL	<b>0.124 01</b>	0.124 087
FASTMRI	U-NET	<b>0.734 746</b>	0.734 311
IMAGENET	RESNET-50	<b>0.232 56</b>	0.243 34
	ViT	<b>0.219 92</b>	0.236 16
LIBRISPEECH	CONFORMER	<b>0.075 989</b>	0.080 673
	DEEPSPEECH	<b>0.112 706</b>	0.120 674
OGBG	GNN	<b>0.282 14</b>	0.276 307
WMT 2016	TRANSFORMER	<b>31.3523</b>	30.9950

Table 3: Performance across multiple workloads for two ADAMW training methods. Bolded number highlights the better performance.

- ▶ Methods have hyperparameters that need to be set/tuned
- ▶ We can't just ignore them as they are workload-dependent
- ▶ Same update rule with different hyperparameters is effectively a different algorithm



- ▶ **Hyperparameter (Search Spaces) are Part of the Submission**

Submitters have to provide search spaces or default values.

- ▶ **A Competitive Benchmark**

Everyone submits their algorithm and how it should be tuned, therefore generating strong baselines.

- ▶ **Two Tuning Rulesets**

- ▶ In the **self-tuning** ruleset, everything is done on the clock (e.g. line-searches, freeze-thaw, default hyperparameters, etc.).
- ▶ The **external** tuning ruleset allows parallel resources and only the fastest trial counts.

# Summary

- ▶ Neural networks are useful but expensive models
- ▶ We are currently unable to identify which training algorithms are “best”
- ▶ We created the **AlgoPerf: Training Algorithms Benchmark** to find faster deep learning training algorithms
  - ▶ A competitive, time-to-result benchmark
  - ▶ Running on fixed hardware and workloads
  - ▶ Computing a joint score across multiple realistic workloads

**Read the Rules**

[github.com/mlcommons/algorithmic-efficiency/blob/main/RULES.md](https://github.com/mlcommons/algorithmic-efficiency/blob/main/RULES.md)

**Read the Paper**

[arxiv.org/abs/2306.07179](https://arxiv.org/abs/2306.07179)

**Submit!**

Call for Submission coming soon

