# The **AlgoPerf**: Training Algorithms Benchmark

## Faster neural network training through better training algorithms

**Frank Schneider**

September 14, 2023

# 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 fiddeling)
- ► ...**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 fiddeling)
- ► ...**expensive** (requiring dozens of trial runs)

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

# The state of deep learning training methods

A confusingly crowded field of methods & hyperparameters

## A huge number of training methods...

| Name | Ref. | Name | Ref. | Name | Ref. |
|------|------|------|------|------|------|

*from "Descending through a Crowded Valley"*
*(Schmidt, Schneider, Hennig; 2021)*

# The state of deep learning training methods

A confusingly crowded field of methods & hyperparameters

## A huge number of training methods...



*from "Descending through a Crowded Valley"*
*(Schmidt, Schneider, Hennig; 2021)*

## ...and training tricks

- ► OneCycle scheduler, gradient checkpointing
- ► Genetic Algorithm for Hyperparameters
- ► Avoid batches that lead to NaN/inf losses
- ► One cycle, low fidelity training, SGD with restarts
- ► Proximal optimization for regularizers
- ► Line searches for the maximum learning rate
- ► Normalized updates
- ► Distributed Shampoo, Normformer, GLU
- ► Weight averaging

- ► FreezeOut
- ► A different epsilon value!
- ► Check hyperparameter performance over multiple seeds
- ► Lowering the learning rate!
- ► Normalizing data works better than batch or layer norm
- ► Mixed precision training
- ► Train with a small subset
- ► Cyclic and one cycle LR
- ► Label smoothing
- ► ...

*from the NeurIPS "HITY Workshop"*
*(Schneider et al.; 2022)*

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

# The state of "benchmarking" in current deep learning optimizer papers
No standardized procedure to follow

- ► Each paper "invents" their own evaluation protocol.
- ► Unreasonably hard to perform a convincing, informative, and practically relevant comparison with strong baselines.
- ► Lot's of subtle pitfalls with tuning, problem-specification, etc.

# Example of pitfalls when comparing training algorithms

Which algorithms 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 (✖).

# Example of pitfalls when comparing training algorithms

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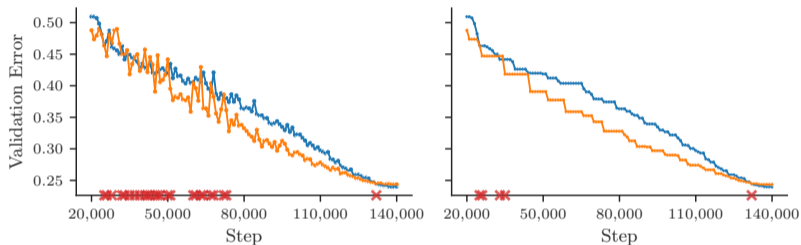


Figure 1: *Left:* Validation error for two different runs (—, —) of ADAM on RESNET-50 on IMAGENET. *Right:* The *best* validation error obtained so far. The runs intersect multiple times (✖).

- ► Directly comparing training curves is ill-posed.
- ► Without defining the target in advance, we can champion any method (moving the goal post after the experiment).
- ► Tuning goals must align.

# MLCommons introduces the **AlgoPerf**: Training Algorithm Benchmark

An unprecedented effort to find faster deep learning training algorithms

ML
●Commons

**Algorithms Working Group**

**AlgoPerf: Training Algorithm Benchmark**
**A standardized benchmark competition to measure neural network training speedups due to algorithmic changes.**

An open large-scale effort by 25+ researchers from Google, University of Tübingen, University of Toronto, Meta AI, etc.

Chaired by

**George Dahl**      Google

**Frank Schneider**      University of Tübingen

# The key features of AlgoPerf

A standardized competitive time-to-results benchmark

► **A competitive time-to-results benchmark.**
  → Everyone's submissions are everyone's **strong baselines**.

# The key features of AlgoPerf

A standardized competitive time-to-results benchmark

- ► **A competitive time-to-results benchmark.**
  → Everyone's submissions are everyone's **strong baselines**.
- ► **Fixed hardware, workloads, and process.**
  → Submissions need to provide **training algorithm improvements**.

# The key features of AlgoPerf

A standardized competitive time-to-results benchmark

- ▶ **A competitive time-to-results benchmark.**
  → Everyone's submissions are everyone's **strong baselines**.
- ▶ **Fixed hardware, workloads, and process.**
  → Submissions need to provide **training algorithm improvements**.
- ▶ **Aggregate across a variety of realistic workloads using performance profiles.**
  → No specialized solutions but **general-purpose methods**.

# The key features of AlgoPerf

A standardized competitive time-to-results benchmark

- ▶ **A competitive time-to-results benchmark.**
  → Everyone's submissions are everyone's **strong baselines**.
- ▶ **Fixed hardware, workloads, and process.**
  → Submissions need to provide **training algorithm improvements**.
- ▶ **Aggregate across a variety of realistic workloads using performance profiles.**
  → No specialized solutions but **general-purpose methods**.
- ▶ **Explicitly accounts for hyperparameter tuning by providing search spaces.**
  → **Runable training algorithms**, not algorithm templates.

## The **AlgoPerf** Training Algorithms Benchmark

We need you!

| **Read the Rules** | github.com/mlcommons/algorithmic-efficiency/blob/main/RULES.md |
| --- | --- |
| **Read the Paper** | arxiv.org/abs/2306.07179 |
| **Submit!** | Call for Submission coming soon |

► **Benchmark experts** Join the effort and tell us how to improve!

► **ML community** Help us spread the word!

► **Algorithm researchers** Submit! AlgoPerf is the easiest way to convincingly demonstrate the capabilities of your training algorithm.