

The Results of the Inaugural **ALGOPERF** Competition

How to train a neural net 30% faster

Frank Schneider

February 11, 2025

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



Federal Ministry
of Education
and Research

imprs-is

The ALGO PERF competition

A (very) short summary

ML ● Commons

The goal of the ALGO PERF: TRAINING ALGORITHMS benchmark & competition is to measure speed-ups in neural network training due to algorithmic improvements.

What are the best algorithms to train neural networks?

The ALGOPERF competition

A (very) short summary

ML ● Commons

The goal of the ALGOPERF: TRAINING ALGORITHMS benchmark & competition is to measure speed-ups in neural network training due to algorithmic improvements.

What are the best algorithms to train neural networks?

Why?

- ▶ There was no established protocol to benchmark deep learning training methods.
- ▶ There are lots of subtle pitfalls when evaluating training methods, such as hyperparameter tuning, training horizons, isolating the algorithm, etc. (see Dahl et al. (2023))
- ▶ Unreasonably hard task for researchers to perform a convincing, fair, and practically relevant comparison with strong baselines.

The Key Features of ALGOPERF

Training real-world deep learning workloads as fast as possible

Task	Dataset	Model	Loss	Metric	Validation Target	Maximum Runtime
Clickthrough rate prediction	CRITEO 1TB	DLRMSMALL	Cross Entropy	Cross Entropy	0.123 735	7703
MRI reconstruction	FASTMRI	U-NET	L1	SSIM	0.7344	8859
Image classification	IMAGENET	RESNET-50 ViT	Cross Entropy	Error Rate	0.225 69	63 008
			Cross Entropy	Error Rate	0.226 91	77 520
Speech recognition	LIBRISPEECH	CONFORMER DEEPSPEECH	CTC	Word Error Rate	0.085 884	61 068
			CTC	Word Error Rate	0.119 936	55 506
Molecular property prediction	OGBG	GNN	Cross Entropy	mAP	0.280 98	18 477
Translation	WMT	TRANSFORMER	Cross Entropy	BLEU	30.8491	48 151

The Key Features of ALGOPERF

Isolating algorithmic improvements

Submissions can only modify the training **algorithm** and must leave all other aspects untouched

- ▶ `update_params`: Typically involves optimizers such as SGD, ADAM, or custom methods.
- ▶ `init_optimizer_state`: Define a method's internal state, e.g. to define learning rate schedules.
- ▶ `data_selection`: How to construct batches of data.
- ▶ `get_batch_size`: Batch sizes for each workload, e.g. the largest batch size fitting in memory.
- ▶ (In the external tuning ruleset) `hyperparameter_search_space`: A *workload-agnostic* tuning space for a method's hyperparameters.

The Key Features of ALGOPERF

Two distinct rulesets simulating different use cases

External Tuning Ruleset

Parallel tuning across 5 tuning trials

Fastest trial counts for scoring

Submissions must define a workload-agnostic search space

Simulates training with parallel resources, e.g. multiple devices

Examples: Learning rate tuning using a log grid or a list of five hyperparameter configurations

Self-Tuning Ruleset

No additional tuning, i.e. a single trial

All computations are “on-the-clock”

Any required workload-adaptation must be part of the method

Simulates (sequential) training using a single device

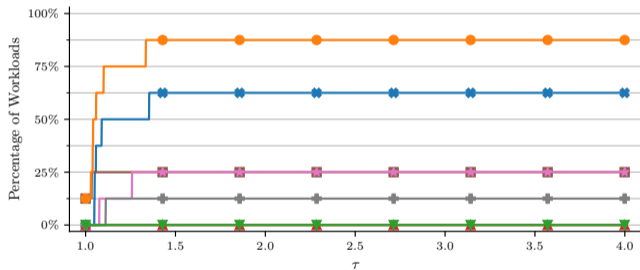
Examples: ADAM with default hyperparameters or inner-loop tuning during the run

Repeat this process five times across different *studies* (with different random seeds) and take the median for a more robust final score.

The Key Features of ALGOPERF

Aggregate scoring using performance profiles

- ▶ **Workload scores** = median wall-clock runtimes to reach the target $t_{s,w}$.
- ▶ **Performance ratio** = workload score relative to the fastest workload score, i.e. $r_{s,w} = \frac{t_{s,w}}{\min_{s \in S} t_{s,w}}$.
- ▶ **Performance profile** = plot the fraction of workloads where a submission is less than τ away from the fastest submission, i.e. workloads where $r_{s,w} \leq \tau$.
- ▶ **Benchmark score** = integrate the performance profile, i.e. $B_s \in [0, 1]$.



Key Takeaways I

The key features of ALGOPERF



- ▶ **A competitive time-to-results benchmark**
→ Strong baselines.
- ▶ **Fixed hardware, workloads, and evaluation protocol**
→ Submissions need to innovate on the *training algorithms*.
- ▶ **Fully-specified algorithms that need to perform well across multiple workloads**
→ No cherry-picking, general-purpose methods with complete training recipes & properly account for hyperparameter tuning.

The Key Results of ALGOPERF

New SOTA in neural network training methods

External Tuning Ruleset

Submission	Authors	Institutions
DISTRIBUTED SHAMPOO	Shi, Lee, et al.	Meta Platforms
SCHEDULE FREE ADAMW	Defazio, Yang, Mishchenko	Meta AI, Samsung AI
GENERALIZED ADAM	Dahl, Medapati, et al.	Google
CYCLIC LR	Ajroldi, Orvieto, Geiping	MPI-IS, ELLIS Tübingen
NADAMP	Dahl, Medapati, et al.	Google
BASELINE		
AMOS	Tian	Google
CASPR ADAPTIVE	Duvvuri, Dhillon, Hsieh	UT Austin, UCLA, Google
LAWA QUEUE	Ajroldi, Orvieto, Geiping	MPI-IS, ELLIS Tübingen
LAWA EMA	Ajroldi, Orvieto, Geiping	MPI-IS, ELLIS Tübingen
SCHEDULE FREE PRODIGY	Defazio, Yang, Mishchenko	Meta AI, Samsung AI

Self-Tuning Ruleset

Submission	Authors	Institutions
SCHEDULE FREE ADAMW	Defazio, Yang, Mishchenko	Meta AI, Samsung AI
BASELINE		
NADAMW SEQUENTIAL	Dahl, Medapati, et al.	Google
SINV6 75	Moudgil	Mila, Concordia University
SINV6	Moudgil	Mila, Concordia University
ADAMG	Pang	Michigan State University

The Key Results of ALGOPERF

New SOTA in neural network training methods

External Tuning Ruleset

Submission	Authors	Institutions	Score
DISTRIBUTED SHAMPOO	Shi, Lee, et al.	Meta Platforms	0.7784
SCHEDULE FREE ADAMW	Defazio, Yang, Mishchenko	Meta AI, Samsung AI	0.7077
GENERALIZED ADAM	Dahl, Medapati, et al.	Google	0.6383
CYCLIC LR	Ajroldi, Orvieto, Geiping	MPI-IS, ELLIS Tübingen	0.6301
NADAMP	Dahl, Medapati, et al.	Google	0.5909
BASELINE			0.5707
AMOS	Tian	Google	0.4918
CASPR ADAPTIVE	Duvvuri, Dhillon, Hsieh	UT Austin, UCLA, Google	0.4722
LAWA QUEUE	Ajroldi, Orvieto, Geiping	MPI-IS, ELLIS Tübingen	0.3699
LAWA EMA	Ajroldi, Orvieto, Geiping	MPI-IS, ELLIS Tübingen	0.3384
SCHEDULE FREE PRODIGY	Defazio, Yang, Mishchenko	Meta AI, Samsung AI	0

Self-Tuning Ruleset

Submission	Authors	Institutions	Score
SCHEDULE FREE ADAMW	Defazio, Yang, Mishchenko	Meta AI, Samsung AI	0.8542
BASELINE			0.8194
NADAMW SEQUENTIAL	Dahl, Medapati, et al.	Google	0.3308
SINV6 75	Moudgil	Mila, Concordia University	0.1420
SINV6	Moudgil	Mila, Concordia University	0.0903
ADAMG	Pang	Michigan State University	0

The Key Results of ALGOPERF

More intuitive speedup numbers

- ▶ DISTRIBUTED SHAMPOO is on average 28 % faster than the (external tuning) BASELINE.
- ▶ SCHEDULE FREE ADAMW is on average 8 % faster than the (self-tuning) BASELINE.
- ▶ Comparisons across rulesets:
 - ▶ DISTRIBUTED SHAMPOO is on average 24 % faster than (self-tuning) SCHEDULE FREE ADAMW.
 - ▶ (self-tuning) SCHEDULE FREE ADAMW 10 % faster than the (external tuning) BASELINE.¹

¹Across the seven workloads both methods trained successfully.

The Key Results of ALGO PERF

Robustness is a major aspect of training methods

	CRITEO 1TB	FASTMRI	RESNET	VIT	CONFORMER	DEEPSPEECH	OGBG	WMT
DISTRIBUTED SHAMPOO	0.65	0.15	<i>inf</i>	0.43	0.78	0.62	0.18	0.80
SCHEDULE FREE ADAMW	0.67	0.13	<i>inf</i>	0.57	0.92	0.78	0.29 [‡]	0.33
GENERALIZED ADAM	0.83	0.18	0.97	0.84	<i>inf</i>	0.68	0.31 [‡]	0.63
CYCLIC LR	0.67	0.25	<i>inf</i>	0.81	0.94	0.70	0.38 [‡]	0.49
NADAMP	0.80	0.22	<i>inf</i>	0.88	0.94	0.60	0.43 [‡]	0.80
BASELINE	0.94	0.23	<i>inf</i>	0.91	0.90	0.65	0.42 [‡]	0.86
AMOS	<i>inf</i>	0.33	<i>inf</i>	0.65	0.71	0.57	0.60 [*]	0.68
CASPR ADAPTIVE	NaN	0.13	<i>inf</i>	0.58	<i>inf</i>	0.75	0.12	0.67 [‡]
LAWA QUEUE	<i>inf</i>	0.22	<i>inf</i>	0.66	<i>inf</i>	<i>inf</i>	0.25	0.56
LAWA EMA	0.69	0.29	<i>inf</i>	0.80	<i>inf</i>	<i>inf</i>	0.57 [*]	0.89
SCHEDULE FREE PRODIGY	NaN	0.21 [‡]	<i>inf</i>	<i>inf</i>	<i>inf</i>	<i>inf</i>	0.61 [*]	<i>inf</i>

Key Takeaways II

The key results of ALGOPERF

- ▶ **Significant improvements in neural net training**
 - The winners provide 28 % and 10 % faster training vs. their baseline.
 - Novel methods can improve over ADAM (e.g. non-diagonal, second-order, ...).
- ▶ **Despite these improvements, there is ample potential left to explore**
 - No submission dominates across workloads.
 - For external tuning, 5 submissions were the fastest across the 8 workloads.
- ▶ **ALGOPERF can meaningfully evaluate training algorithms and identify practically useful methods**
 - The results are robust to many of our benchmarking decisions.
 - The benchmark must evolve and improve alongside the methods.

Summary & What This Means for You

Results of the Inaugural ALGOPERF Competition

- ▶ With SHAMPOO & SCHEDULE-FREE, we have two new exciting SOTA training algorithms.
 - ▶ **Try them out and let us know your results!**
- ▶ There is even more potential for future improvement thanks to the clear signal provided by ALGOPERF.
 - ▶ **Help us try out SOAP, MUON, AdEMAMix, ...!**
- ▶ The benchmark needs to evolve and improve with the submissions.
 - ▶ **Help us shape the next iteration of ALGOPERF!**



...and so many more!

Benchmark Code: github.com/mlcommons/algorithmic-efficiency

Blog Post: mlcommons.org/2024/08/mlc-algoperf-benchmark-competition

Results Paper: Accepted at ICLR 2025! openreview.net/forum?id=CtM5xjRSfm

Appendix

Performance profiles: External tuning ruleset

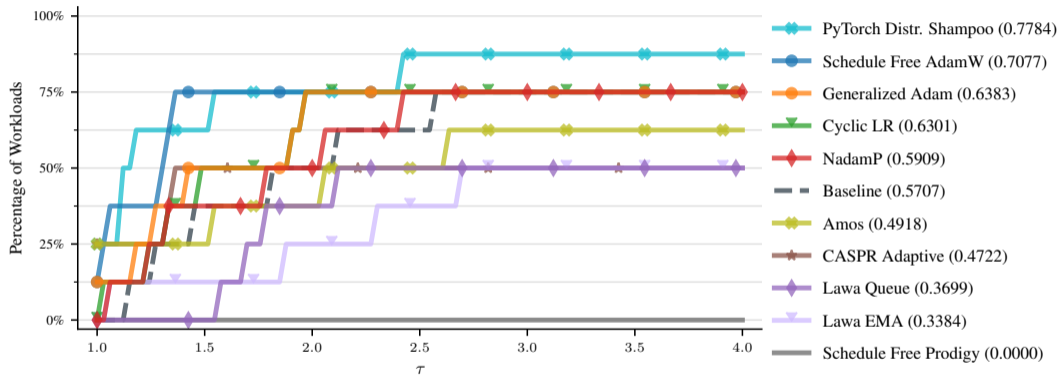


Figure 1: Performance profiles for the external tuning ruleset

Appendix

Performance profiles: Self-tuning ruleset

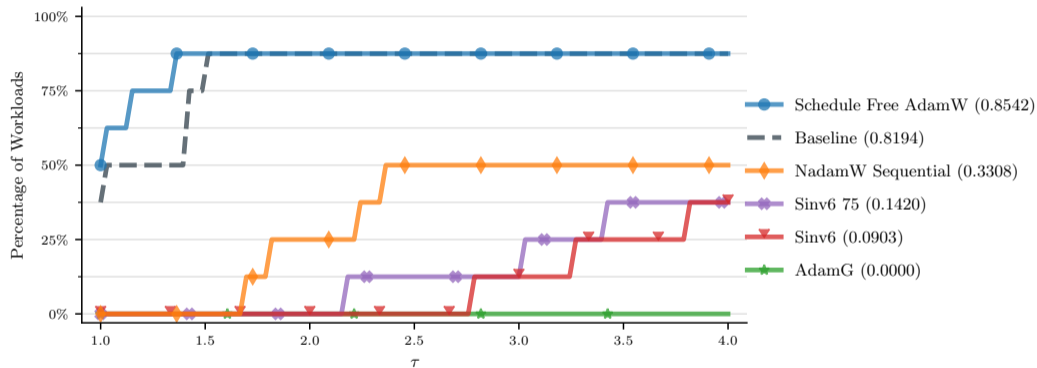


Figure 2: Performance profiles for the self-tuning ruleset

Appendix

Performance profiles: External tuning ruleset (without held-out workloads)

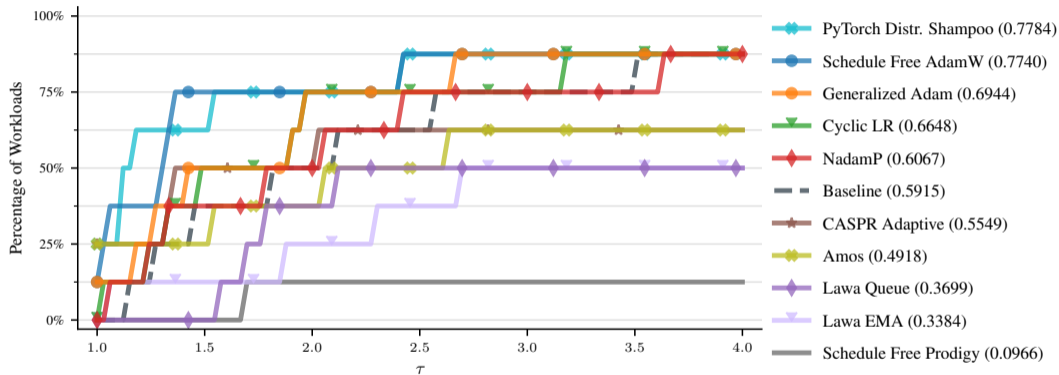


Figure 3: Performance profiles for the external tuning ruleset when ignoring held-out workloads.

Appendix

Performance profiles: Self-tuning ruleset (without held-out workloads)

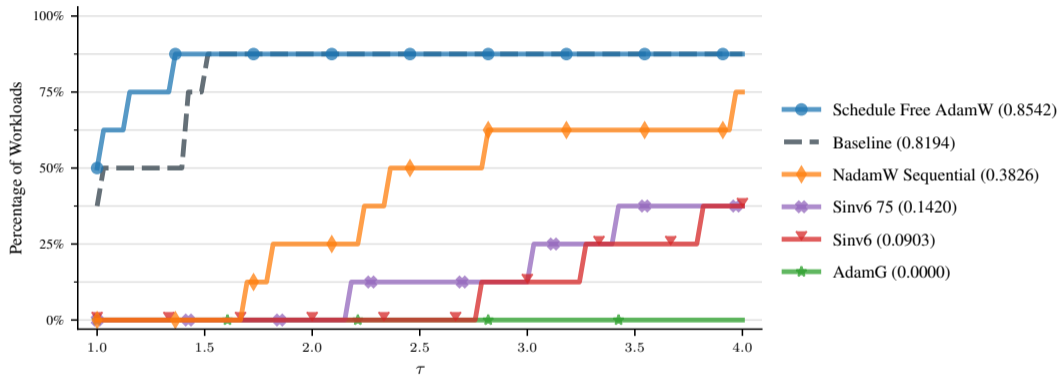


Figure 4: Performance profiles for the self-tuning ruleset when ignoring held-out workloads

Appendix

Performance profiles: External tuning ruleset (qualification set)

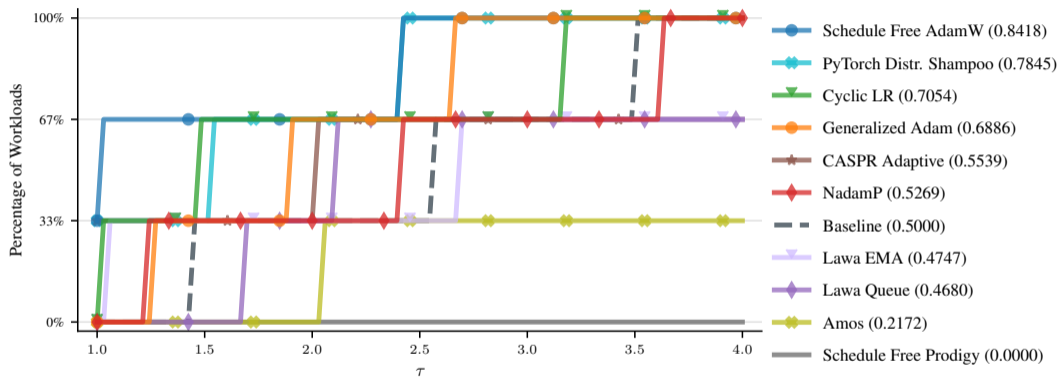


Figure 5: Performance profiles for the external tuning ruleset on the qualification set

Appendix

Performance profiles: Self-tuning ruleset (qualification set)

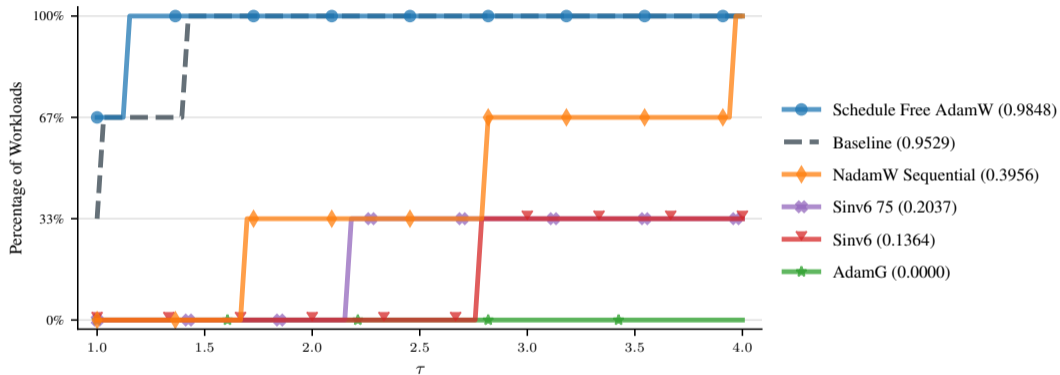


Figure 6: Performance profiles for the self-tuning ruleset on the qualification set

Appendix

Per-workload runtimes: Self-tuning ruleset

	CRITEO 1TB	FASTMRI	RESNET	ViT	CONFORMER	DEEPSPEECH	OGBG	WMT
ADAMG	<i>inf</i>	<i>inf</i>	<i>inf</i>	<i>inf</i>	<i>inf</i>	<i>inf</i>	<i>inf</i>	<i>inf</i>
BASELINE	0.75	0.22	<i>inf</i>	0.95	0.94	0.65	0.46	0.84
NADAMW SEQUENTIAL	2.96 [‡]	0.27	<i>inf</i>	1.58	<i>inf</i>	1.45	0.55	2.36 [‡]
SCHEDULE FREE ADAMW	0.75	0.15	<i>inf</i>	0.68	0.97	0.88	0.32	0.94
SINV6	NaN	0.49	<i>inf</i>	<i>inf</i>	<i>inf</i>	2.47	1.35 [*]	2.32
SINV6 75	NaN	0.45	<i>inf</i>	<i>inf</i>	<i>inf</i>	2.21	1.50 [*]	1.82

Appendix

ResNet near misses: DISTRIBUTED SHAMPOO

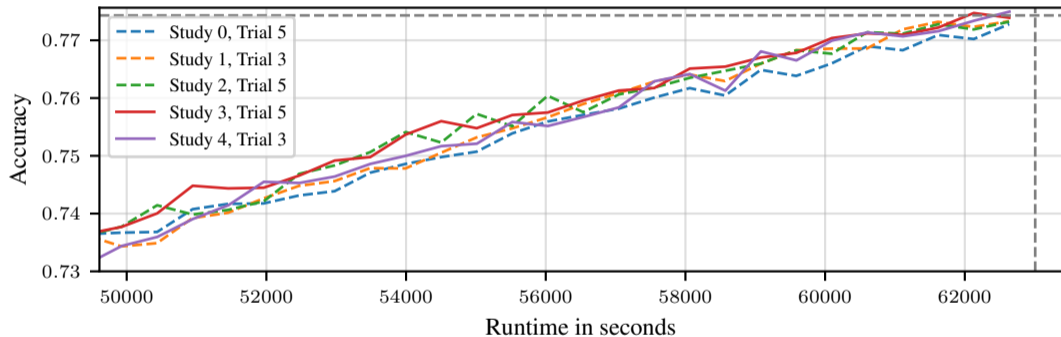


Figure 7: PYTORCH DISTRIBUTED SHAMPOO

Appendix

ResNet near misses: NADAMP

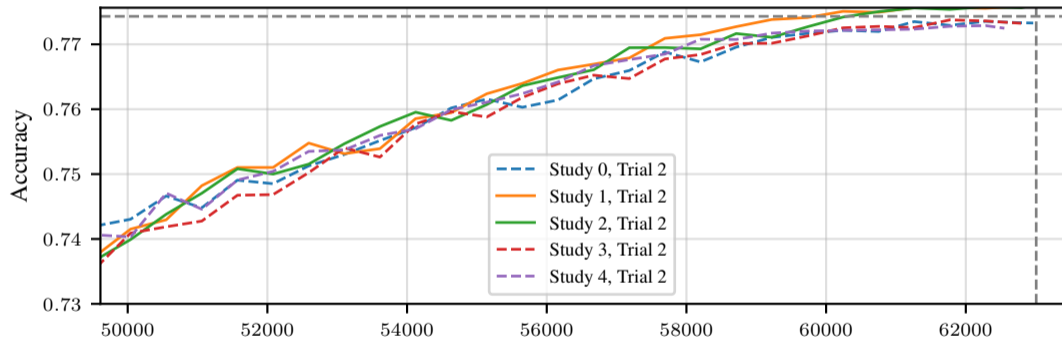


Figure 8: NADAMP

Appendix

ResNet near misses: BASELINE

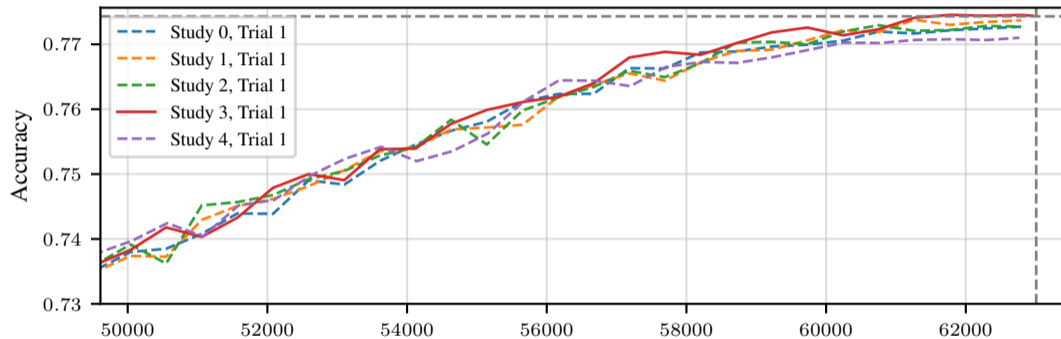


Figure 9: BASELINE

Appendix

Benchmark scores as a function of τ_{\max} : External tuning

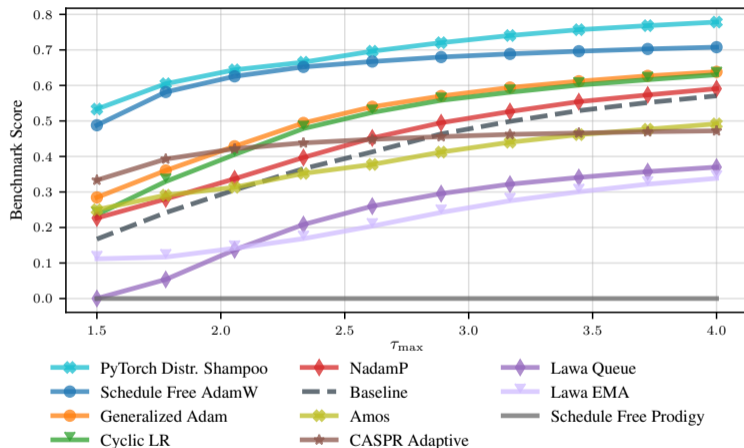


Figure 10: Benchmark scores as a function of τ_{\max} (external tuning).

Appendix

Benchmark scores as a function of τ_{\max} : Self-tuning

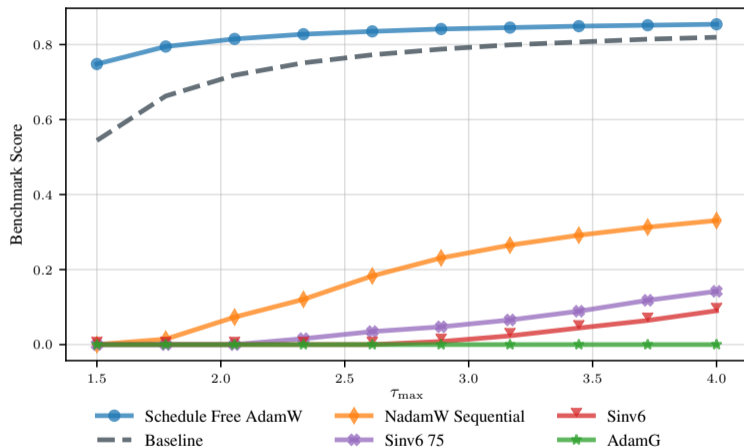


Figure 11: Benchmark scores as a function of τ_{\max} (self-tuning).

Appendix

(a) External tuning ruleset

	Full		CRITEO 1TB		FASTMRI		RESNET		ViT		CON-FORMER		DEEP SPEECH		OGBG		WMT	
	Score	Rank	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.
PYTORCH DISTR. SHAMPOO	0.78	1	0.75	1	0.75	1	0.89	1	0.75	1	0.75	1	0.75	1	0.77	2	0.81	1
SCHEDULE FREE ADAMW	0.71	2	0.67	2	0.67	2	0.81	2	0.68	2	0.68	3	0.68	2	0.81	1	0.67	2
GENERALIZED ADAM	0.64	3	0.60	3	0.61	4	0.59	6	0.63	3	0.73	2	0.59	3	0.73	3	0.63	3
CYCLIC LR	0.63	4	0.58	4	0.62	3	0.72	3	0.62	4	0.59	4	0.59	4	0.72	4	0.60	4
NADAMP	0.59	5	0.54	6	0.57	5	0.68	4	0.58	5	0.55	5	0.53	5	0.68	5	0.60	5
BASELINE	0.57	6	0.53	8	0.55	6	0.65	5	0.56	6	0.52	7	0.52	6	0.65	6	0.58	6
AMOS	0.49	7	0.56	5	0.50	7	0.56	7	0.44	7	0.42	9	0.42	8	0.56	7	0.47	8
CASPR ADAPTIVE	0.47	8	0.54	7	0.40	8	0.54	8	0.41	8	0.54	6	0.41	9	0.40	8	0.54	7
LAWA QUEUE	0.37	9	0.42	9	0.32	9	0.42	9	0.31	9	0.42	8	0.42	7	0.33	10	0.31	10
LAWA EMA	0.34	10	0.25	10	0.31	10	0.39	10	0.28	10	0.39	10	0.39	10	0.39	9	0.32	9
SCHEDULE FREE PRODIGY	0.00	11	0.00	11	0.00	11	0.00	11	0.00	11	0.00	11	0.00	11	0.00	11	0.00	11

Appendix

(a) Self-tuning ruleset

		Full		CRITEO 1TB		FASTMRI		RESNET		ViT		CON- FORMER		DEEP SPEECH		OGBG		WMT	
		Score	Rank	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.	S.	R.
SCHEDULE ADAMW	FREE	0.85	1	0.83	1	0.83	1	0.98	1	0.83	1	0.83	1	0.85	1	0.83	1	0.84	1
BASILINE		0.82	2	0.79	2	0.82	2	0.94	2	0.81	2	0.79	2	0.79	2	0.81	2	0.79	2
NADAMW SEQUENTIAL		0.33	3	0.38	3	0.27	3	0.38	3	0.30	3	0.38	3	0.29	3	0.27	3	0.38	3
SINV6 75		0.14	4	0.16	4	0.12	4	0.16	4	0.16	4	0.16	4	0.13	4	0.16	4	0.08	4
SINV6		0.09	5	0.10	5	0.07	5	0.10	5	0.10	5	0.10	5	0.09	5	0.10	5	0.04	5
ADAMG		0.00	6	0.00	6	0.00	6	0.00	6	0.00	6	0.00	6	0.00	6	0.00	6	0.00	6