

# Descending through a Crowded Valley – Benchmarking Deep Learning Optimizers

Robin M. Schmidt, Frank Schneider, Philipp Hennig  
ICML 2021

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



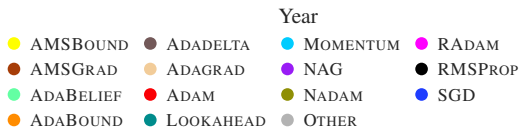
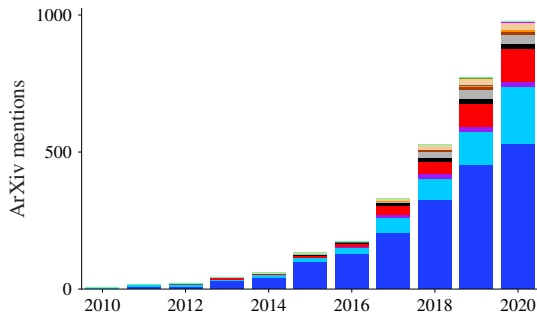
Max Planck Institute for  
**Intelligent Systems**  
imprs-is



some of the presented work is supported  
by the European Research Council.

# Optimization for Deep Learning

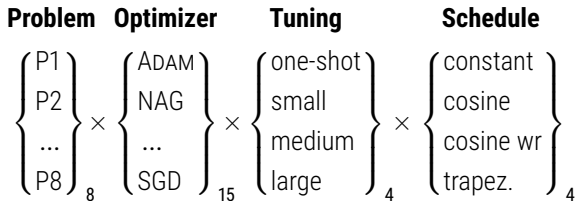
A crowded valley of methods



Name	Ref.	Name	Ref.
AccleGrad	(Levy et al., 2018)	HyperAdam	(Wang et al., 2019b)
ACClip	(Zhang et al., 2020)	K-BFGS/K-BFGS(L)	(Goldfarb et al., 2020)
AdaAlter	(Xie et al., 2019)	KF-QN-CNN	(Ren & Goldfarb, 2021)
AdaBatch	(Devarakonda et al., 2017)	KFAC	(Martens & Grosse, 2015)
AdaBayes/AdaBayes-SS	(Aitchison, 2020)	KFLR/KFRA	(Botev et al., 2017)
AdaBelief	(Zhang et al., 2020)	LAdams/L4Momentum	(Rolnick & Martinus, 2018)
AdaBlock	(Yun et al., 2019)	LAMB	(You et al., 2020)
AdaBound	(Luo et al., 2019)	LaProp	(Ziyin et al., 2020)
AdaComp	(Chen et al., 2018)	LARS	(You et al., 2017)
AdaDelta	(Zeiler, 2012)	LHOPT	(Almeida et al., 2021)
Adafactor	(Shazeer & Stern, 2018)	LookAhead	(Zhang et al., 2019)
AdaFix	(Bae et al., 2019)	M-SVAG	(Balles & Hennig, 2018)
AdaFom	(Chen et al., 2019a)	MADGRAD	(Defazio & Jelastic, 2021)
AdaFTRL	(Orabona & Pili, 2015)	MAS	(Landro et al., 2020)
Adagrad	(Duchi et al., 2011)	MEKA	(Chen et al., 2020b)
ADAHESIAN	(Yao et al., 2020)	MTadam	(Malkiel & Wolf, 2020)
Adai	(Xie et al., 2020)	MVRC-l/MVRC-2	(Chen & Zhou, 2020)
AdaLoss	(Teixeira et al., 2019)	Nadam	(Dozat, 2016)
Adam	(Kingma & Ba, 2015)	NAMSB/NAMSG	(Chen et al., 2019b)
Adam+	(Liu et al., 2020b)	ND-Adam	(Zhang et al., 2017a)
AdamAL	(Tao et al., 2019)	Nero	(Liu et al., 2021b)
AdaMax	(Kingma & Ba, 2015)	Nesterov	(Nesterov, 1983)
AdamBS	(Liu et al., 2020c)	Noisy Adam/Noisy K-FAC	(Zhang et al., 2018)
AdamNC	(Reddi et al., 2018)	NosAdam	(Huang et al., 2019)
AdaMod	(Ding et al., 2019)	Novograd	(Ginsburg et al., 2019)
AdamP/SGDP	(Heo et al., 2021)	NT-SGD	(Zhou et al., 2021b)
AdamT	(Zhou et al., 2020)	Padam	(Chen et al., 2020a)
AdamW	(Loshchilov & Hutter, 2019)	PAGE	(Li et al., 2020b)
AdamX	(Tran & Phong, 2019)	PAL	(Mutschler & Zell, 2020)
ADAS	(Eliyahu, 2020)	PolyAdam	(Orvieto et al., 2019)
AdaS	(Hosseini & Piatanotis, 2020)	Polyak	(Polyak, 1964)
AdaScale	(Johnson et al., 2020)	PowerSGD/PowerSGDM	(Vogels et al., 2019)
AdaSGD	(Wang & Wiens, 2020)	Probabilistic Polyak	(de Roos et al., 2021)
AdaShift	(Zhou et al., 2019)	ProBS	(Mahseeci & Hennig, 2017)
AdaSqrt	(Hu et al., 2019)	PSIorm	(Xu, 2020)
Adahm	(Sun et al., 2019)	QHAdam/QHM	(Ma & Yarats, 2019)
AdaX/AdaX-W	(Li et al., 2020a)	RADAM	(Liu et al., 2020a)
AEGL	(Liu & Tian, 2020)	Ranger	(Wright, 2020b)
ALI-G	(Berrada et al., 2020)	RangerLars	(Franklin, 2020)
AMSBound	(Luo et al., 2019)	RMSProp	(Tejerman & Hinton, 2012)
AMSGrad	(Reddi et al., 2018)	RMSIerov	(Choi et al., 2019)

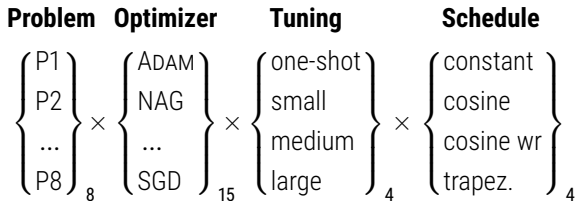
# Benchmark Setup

The many dimensions to explore



# Benchmark Setup

The many dimensions to explore

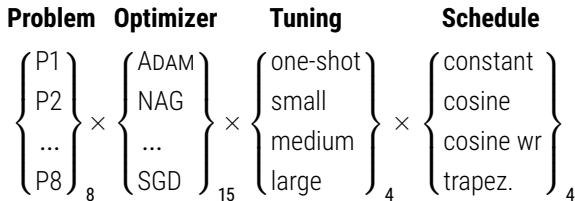


- |             |             |           |
|-------------|-------------|-----------|
| ● AMSBOUND  | ● ADAGRAD   | ● NAG     |
| ● AMSGRAD   | ● ADAM      | ● NADAM   |
| ● ADABELIEF | ● LA(MOM.)  | ● RADAM   |
| ● ADABOUND  | ● LA(RADAM) | ● RMSPROP |
| ● ADADELTA  | ● MOMENTUM  | ● SGD     |

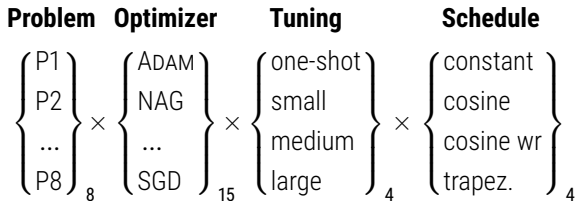


# Benchmark Setup

The many dimensions to explore



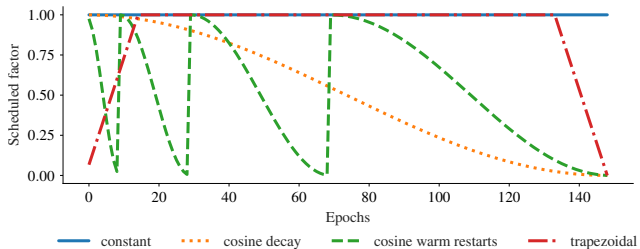
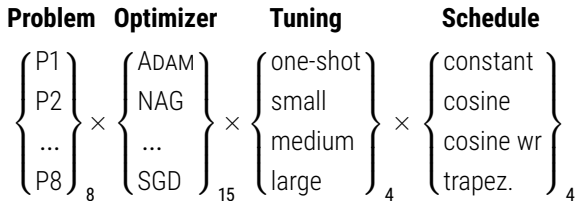
	<b>Data set</b>	<b>Model</b>	<b>Task</b>	<b>Approx. run time</b>
<b>P1</b>	Artificial	Noisy quadratic	Minimization	< 1 min
<b>P2</b>	MNIST	VAE	Generative	10 min
<b>P3</b>	Fashion-MNIST	Simple CNN: 2c2d	Classification	20 min
<b>P4</b>	CIFAR-10	Simple CNN: 3c3d	Classification	35 min
<b>P5</b>	Fashion-MNIST	VAE	Generative	20 min
<b>P6</b>	CIFAR-100	All-CNN-C	Classification	4 h 00 min
<b>P7</b>	SVHN	Wide ResNet 16-4	Classification	3 h 30 min
<b>P8</b>	War and Peace	RNN	Character Prediction	5 h 30 min



- ✦ **One-Shot** - 1 Run  
No tuning, uses default hyperparameters
- ✦ **Small** - 25 Runs  
Tuned via random search
- ✦ **Medium** - 50 Runs  
Tuned via random search, superset of *small budget*
- ✦ **Large** - 75 Runs  
Tuned via random search, refined search spaces

# Benchmark Setup

The many dimensions to explore



# Results: Out-of-the-box performance

Why trying out optimizers can be better than tuning them

✦ **Orange rows**

*bad default hyperparameters*  
SGD, NAG, MOMENTUM,  
AMSGRAD, ADADELTA

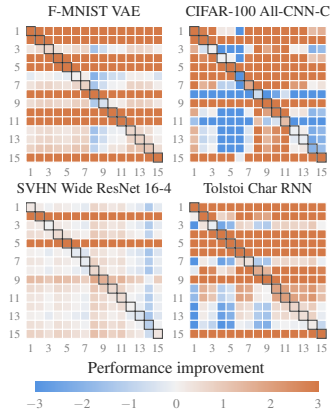
✦ **White & blue rows**

*good default hyperparameters*  
ADAM, NADAM, RADAM,  
AMSBOUND, ADABOUND

CIFAR-10 3c3d

One-shot	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
AMSBound: 1	-0.1	-1.0	-0.1	-0.6	-1.9	-0.0	-1.4	-0.5	1.4	-0.3	0.1	-0.3	-1.0	-0.7	-1.1
AMSGrad: 2	39.4	38.5	39.4	38.9	37.6	39.5	38.2	39.0	41.0	39.3	39.6	39.2	38.5	38.8	38.5
AdaBelief: 3	1.1	0.2	1.1	0.6	-0.7	1.2	-0.1	0.7	2.7	1.0	1.3	0.9	0.2	0.5	0.2
AdaBound: 4	0.3	-0.6	0.3	-0.2	-1.5	0.4	-0.9	-0.1	1.9	0.2	0.5	0.1	-0.6	-0.3	-0.7
Adadelata: 5	43.9	43.0	43.8	43.4	42.0	44.0	42.6	43.4	45.4	43.7	44.0	43.7	43.0	43.3	42.9
Adagrad: 6	1.8	0.9	1.8	1.3	-0.0	1.9	0.6	1.4	3.4	1.7	2.0	1.6	0.9	1.2	0.9
Adam: 7	1.5	0.6	1.4	1.0	-0.4	1.6	0.2	1.0	3.0	1.3	1.6	1.3	0.6	0.9	0.5
LA(Mom.): 8	6.5	5.6	6.4	5.9	4.6	6.6	5.2	6.0	8.0	6.3	6.6	6.2	5.6	5.9	5.5
LA(RAdam): 9	6.0	5.1	6.0	5.5	4.1	6.1	4.7	5.6	7.5	5.8	6.2	5.8	5.1	5.4	5.0
Mom.: 10	41.2	40.3	41.2	40.7	39.3	41.3	39.9	40.8	42.7	41.0	41.4	41.0	40.3	40.6	40.2
NAG: 11	18.8	17.9	18.8	18.3	16.9	18.9	17.5	18.4	20.3	18.6	19.0	18.6	17.9	18.2	17.8
Nadam: 12	0.7	-0.2	0.6	0.2	-1.2	0.8	-0.6	0.2	2.2	0.5	0.8	0.5	-0.2	0.1	-0.3
RAdam: 13	1.3	0.4	1.3	0.8	-0.6	1.4	0.0	0.9	2.8	1.1	1.5	1.1	0.4	0.7	0.3
RMSProp: 14	3.2	2.3	3.2	2.7	1.4	3.3	2.0	2.8	4.8	3.1	3.4	3.0	2.3	2.6	2.3
SGD: 15	2.4	1.5	2.4	1.9	0.6	2.5	1.2	2.0	4.0	2.3	2.6	2.2	1.5	1.8	1.5

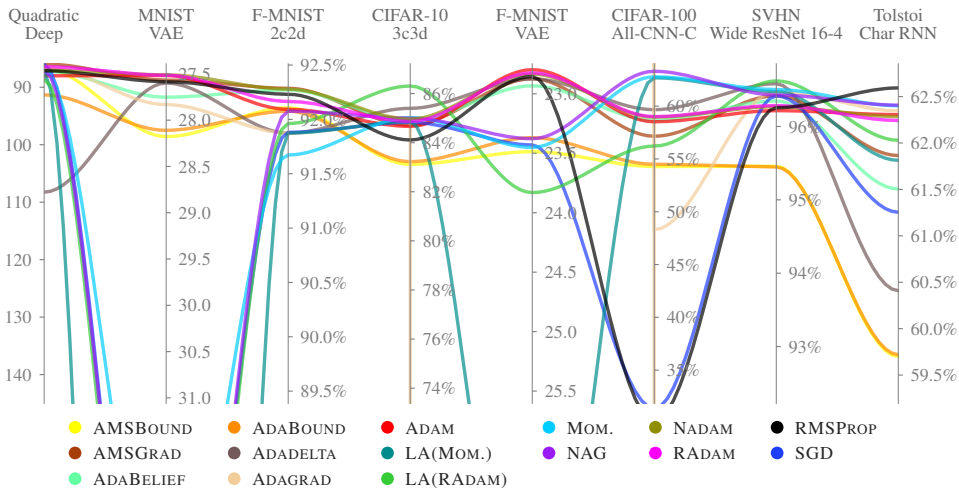
Small budget



# Results: Which optimizer to pick?



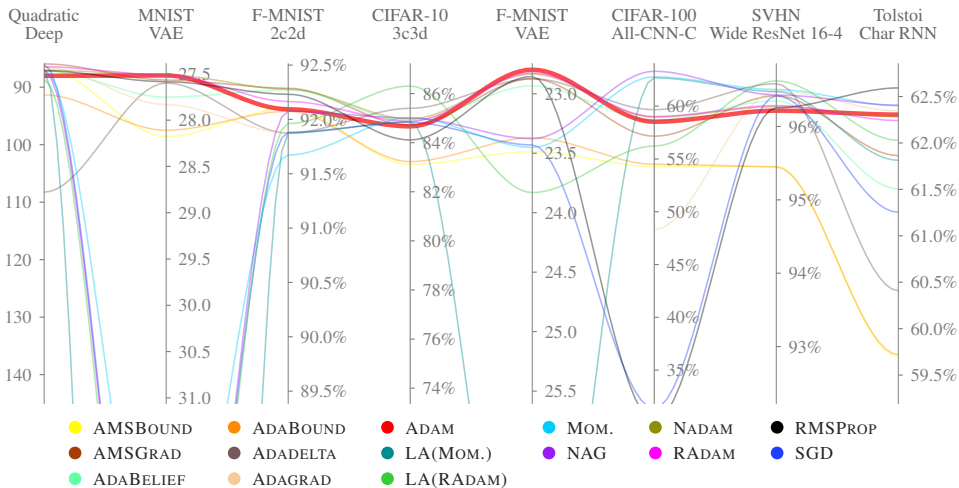
Why ADAM is still a good choice





# Results: Which optimizer to pick?

Why ADAM is still a good choice





---

## Descending through a Crowded Valley – Benchmarking Deep Learning Optimizers

Robin M. Schmidt, Frank Schneider, Philipp Hennig

---

- ✦ No method significantly and consistently outperforms the competition.
- ✦ ADAM remains a viable choice that often ranks near the top.
- ✦ Trying out different optimizers helps about as much as tuning the parameters of one specific method.



**Paper**

arXiv

2007.01547

**Results**



<https://github.com/SirRob1997/Crowded-Valley---Results>

**Framework**

DEEPOBS

<https://github.com/fsschneider/DeepOBS>