

Improving Optimizer Evaluation in Deep Learning

ICLR 2022 Workshop - ML Evaluation Standards

Frank Schneider

April 29, 2022

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



imprs-is

The State of Deep Learning Optimization

A crowded field of methods and hyperparameters

A huge number of optimization methods

AcceleGrad	AMSBound	K-BFGS/K-BFGS(L)	RMSProp
ACClip	AMSGrad	KF-QN-CNN	RMSTerov
AdaAlter	AngularGrad	KFAC	S-SGD
AdaBatch	ArmijoLS	KFLR/KFRA	SAdam
AdaBayes	ARSG	L4Adam/L4Momentum	Sadam/SAMSGrad
AdaBelief	ASAM	LAMB	SALR
AdaBlock	AutoLRS	LaProp	SAM
AdaBound	AvaGrad	LARS	SC-Adagrad-/RMSPProp
AdaComp	BAdam	LHOPT	SDProp
Adadelat	BGAdam	LookAhead	SGD
Adafactor	BPGrad	M-SVAG	SGD-BB
AdaFix	BRMSProp	MADGRAD	SGD-G2
AdaFom	BSGD	MAS	SGDEM
AdaFTRL	C-ADAM	MEKA	SGDHess
Adagrad	CADA	MTAdam	SGDM
ADAHESSIAN	Cool Momentum	MVRC-1/MVRC-2	SGDR
Adai	CProp	Nadam	SHAdagrad
AdaLoss	Curveball	NAMSB/NAMSG	Shampoo
Adam	Dadam	ND-Adam	SignAdam++
Adam ⁺	DeepMemory	Nero	SignSGD
AdamAL	DGNOpt	Nesterov	SKQN/S4QN
AdaMax	DiffGrad	Noisy Adam/Noisy K-FAC	SMS
AdamBS	EAdam	NosAdam	SMG
AdamNC	EKFAC	Novograd	SNGM
AdaMod	Eve	NT-SGD	SoftAdam
AdamP/SGDP	Expectigrad	Padam	SRSGD
AdamT	FastAdaBelief	PAGE	Step-Tuned SGD
AdamW	FRSGD	PAL	SWATS
AdamX	G-AdaGrad	PolyAdam	SWNTS
ADAS	GADAM	Polyak	TAdam
AdaS	Gadam	PowerSGD/PowerSGDM	TEKFAC
AdaScale	GOALS	Probabilistic Polyak	VAdam
AdaSGD	GOLS-1	ProbLS	VR-SGD
AdaShift	Grad-Avg	PStorm	vSGD-b/vSGD-g/vSGD-l
AdaSqrt	GRAPES	QHAdam/QHM	vSGD-fd
Adathm	Gravilon	RAdam	WNGrad
AdaX/AdaX-W	Gravity	Ranger	YellowFin
AEGD	HAdam	RangerLars	Yogi
ALI-G	HyperAdam		

The State of Deep Learning Optimization

A crowded field of methods and hyperparameters

A huge number of optimization methods

AcceleGrad	AMSBound	K-BFGS/K-BFGS(L)	RMSProp
ACClip	AMSGrad	KF-QN-CNN	RMSTerov
AdaAlter	AngularGrad	KFAC	S-SGD
AdaBatch	ArmijoLS	KFLR/KFRA	SAdam
AdaBayes	ARSG	L4Adam/L4Momentum	Sadam/SAMSGrad
AdaBelief	ASAM	LAMB	SALR
AdaBlock	AutoLRS	LaProp	SAM
AdaBound	AvaGrad	LARS	SC-Adagrad/RMSProp
AdaComp	BAdam	LHOPT	SDProp
Adadelat	BGAdam	LookAhead	SGD
Adafactor	BPGrad	M-SVAG	SGD-BB
AdaFix	BRMSProp	MADGRAD	SGD-G2
AdaFom	BSGD	MAS	SGDEM
AdaFTRL	C-ADAM	MEKA	SGDHess
Adagrad	CADA	MTAdam	SGDM
ADAHESSIAN	Cool Momentum	MVRC-1/MVRC-2	SGDR
Adai	CProp	Nadam	SHAdagrad
AdaLoss	Curveball	NAMSB/NAMSG	Shampoo
Adam	Dadam	ND-Adam	SignAdam++
Adam ⁺	DeepMemory	Nero	SignSGD
AdamAL	DGNOpt	Nesterov	SKQN/S4QN
AdaMax	DiffGrad	Noisy Adam/Noisy K-FAC	SMS
AdamBS	EAdam	NosAdam	SMG
AdamNC	EKFAC	Novograd	SNGM
AdaMod	Eve	NT-SGD	SoftAdam
AdamP/SGDP	Expectigrad	Padam	SRSGD
AdamT	FastAdaBelief	PAGE	Step-Tuned SGD
AdamW	FRSGD	PAL	SWATS
AdamX	G-AdaGrad	PolyAdam	SWNTS
ADAS	GADAM	Polyak	TAdam
AdaS	Gadam	PowerSGD/PowerSGDM	TEKFAC
AdaScale	GOALS	Probabilistic Polyak	VAdam
AdaSGD	GOLS-1	ProbLS	VR-SGD
AdaShift	Grad-Avg	PStorm	vSGD-b/vSGD-g/vSGD-l
AdaSqrt	GRAPES	QHAdam/QHM	vSGD-fd
Adathm	Gravilon	RAdam	WNGrad
AdaX/AdaX-W	Gravity	Ranger	YellowFin
AEGD	HAdam	RangerLars	Yogi
ALI-G	HyperAdam		

Tuning hyperparameters

SGD:

$$\theta = \theta - \eta g$$

The State of Deep Learning Optimization

A crowded field of methods and hyperparameters

A huge number of optimization methods

AcceleGrad	AMSBound	K-BFGS/K-BFGS(L)	RMSProp
ACClip	AMSGrad	KF-QN-CNN	RMSTerov
AdaAlter	AngularGrad	KFAC	S-SGD
AdaBatch	ArmijoLS	KFLR/KFRA	SAdam
AdaBayes	ARSG	L4Adam/L4Momentum	Sadam/SAMSGrad
AdaBelief	ASAM	LAMB	SALR
AdaBlock	AutoLRS	LaProp	SAM
AdaBound	AvaGrad	LARS	SC-Adagrad/RMSProp
AdaComp	BAdam	LHOPT	SDProp
Adadelata	BGAdam	LookAhead	SGD
Adafactor	BPGrad	M-SVAG	SGD-BB
AdaFix	BRMSProp	MADGRAD	SGD-G2
AdaFom	BSGD	MAS	SGDEM
AdaFTRL	C-ADAM	MEKA	SGDHess
Adagrad	CADA	MTAdam	SGDM
ADAHESSIAN	Cool Momentum	MVRC-1/MVRC-2	SGDR
Adai	CProp	Nadam	SHAdagrad
AdaLoss	Curveball	NAMSB/NAMSG	Shampoo
Adam	Dadam	ND-Adam	SignAdam++
Adam+	DeepMemory	Nero	SignSGD
AdamAL	DGNOpt	Nesterov	SKQN/S4QN
AdaMax	DiffGrad	Noisy Adam/Noisy K-FAC	SMS
AdamBS	EAdam	NosAdam	SMG
AdamNC	EKFAC	Novograd	SNGM
AdaMod	Eve	NT-SGD	SoftAdam
AdamP/SGDP	Expectigrad	Padam	SRSGD
AdamT	FastAdaBelief	PAGE	Step-Tuned SGD
AdamW	FRSGD	PAL	SWATS
AdamX	G-AdaGrad	PolyAdam	SWNTS
ADAS	GADAM	Polyak	TAdam
AdaS	Gadam	PowerSGD/PowerSGDM	TEKFAC
AdaScale	GOALS	Probabilistic Polyak	VAdam
AdaSGD	GOLS-1	ProbLS	VR-SGD
AdaShift	Grad-Avg	PStorm	vSGD-b/vSGD-g/vSGD-l
AdaSqrt	GRAPES	QHAdam/QHM	vSGD-fd
Adathm	Gravilon	RAdam	WNGrad
AdaX/AdaX-W	Gravity	Ranger	YellowFin
AEGD	HAdam	RangerLars	Yogi
ALI-G	HyperAdam		

Tuning hyperparameters

SGD:

$$\theta = \theta - \eta g$$

ADAM:

$$m = \beta_1 m + (1 - \beta_1)g$$

$$v = \beta_2 v + (1 - \beta_2)g \odot g$$

$$\hat{m} = \frac{m}{1 - \beta_1^t} \quad \hat{v} = \frac{v}{1 - \beta_2^t}$$

$$\theta = \theta - \frac{\eta}{\sqrt{\hat{v}} + \epsilon} \odot \hat{m}$$

The State of Deep Learning Optimization

A crowded field of methods and hyperparameters

A huge number of optimization methods

AcceleGrad	AMSBound	K-BFGS/K-BFGS(L)	RMSProp
ACClip	AMSGrad	KF-QN-CNN	RMSTerov
AdaAlter	AngularGrad	KFAC	S-SGD
AdaBatch	ArmijoLS	KFLR/KFRA	SAdam
AdaBayes	ARSG	L4Adam/L4Momentum	Sadam/SAMSGrad
AdaBelief	ASAM	LAMB	SALR
AdaBlock	AutoLRS	LaProp	SAM
AdaBound	AvaGrad	LARS	SC-Adagrad/RMSProp
AdaComp	BAdam	LHOPT	SDProp
Adadelata	BGAdam	LookAhead	SGD
Adafactor	BPGrad	M-SVAG	SGD-BB
AdaFix	BRMSProp	MADGRAD	SGD-G2
AdaFom	BSGD	MAS	SGDEM
AdaFTRL	C-ADAM	MEKA	SGDHess
Adagrad	CADA	MTAdam	SGDM
ADAHESSIAN	Cool Momentum	MVRC-1/MVRC-2	SGDR
Adai	CProp	Nadam	SHAdagrad
AdaLoss	Curveball	NAMSB/NAMSG	Shampoo
Adam	Dadam	ND-Adam	SignAdam++
Adam+	DeepMemory	Nero	SignSGD
AdamAL	DGNOpt	Nesterov	SKQN/S4QN
AdaMax	DiffGrad	Noisy Adam/Noisy K-FAC	SMS
AdamBS	EAdam	NosAdam	SMG
AdamNC	EKFAC	Novograd	SNGM
AdaMod	Eve	NT-SGD	SoftAdam
AdamP/SGDP	Expectigrad	Padam	SRSGD
AdamT	FastAdaBelief	PAGE	Step-Tuned SGD
AdamW	FRSGD	PAL	SWATS
AdamX	G-AdaGrad	PolyAdam	SWNTS
ADAS	GADAM	Polyak	TAdam
AdaS	Gadam	PowerSGD/PowerSGDM	TEKFAC
AdaScale	GOALS	Probabilistic Polyak	VAdam
AdaSGD	GOLS-1	ProbLS	VR-SGD
AdaShift	Grad-Avg	PStorm	vSGD-b/vSGD-g/vSGD-l
AdaSqrt	GRAPES	QHAdam/QHM	vSGD-fd
Adathm	Gravilon	RAdam	WNGrad
AdaX/AdaX-W	Gravity	Ranger	YellowFin
AEGD	HAdam	RangerLars	Yogi
ALI-G	HyperAdam		

Benchmarks

Tuning hyperparameters

SGD:

$$\theta = \theta - \eta g$$

ADAM:

$$m = \beta_1 m + (1 - \beta_1)g$$

$$v = \beta_2 v + (1 - \beta_2)g \odot g$$

$$\hat{m} = \frac{m}{1 - \beta_1^t} \quad \hat{v} = \frac{v}{1 - \beta_2^t}$$

$$\theta = \theta - \frac{\eta}{\sqrt{\hat{v}} + \epsilon} \odot \hat{m}$$

The State of Deep Learning Optimization

A crowded field of methods and hyperparameters

A huge number of optimization methods

AcceleGrad	AMSBound	K-BFGS/K-BFGS(L)	RMSProp
ACClip	AMSGrad	KF-QN-CNN	RMSTerov
AdaAlter	AngularGrad	KFAC	S-SGD
AdaBatch	ArmijoLS	KFLR/KFRA	SAdam
AdaBayes	ARSJ	L4Adam/L4Momentum	Sadam/SAMSGrad
AdaBelief	ASAM	LAMB	SALR
AdaBlock	AutoLRS	LaProp	SAM
AdaBound	AvaGrad	LARS	SC-Adagrad-/RMSPProp
AdaComp	BAdam	LHOPT	SDProp
Adadelat	BGAdam	LookAhead	SGD
Adafactor	BPGrad	M-SVAG	SGD-BB
AdaFix	BRMSProp	MADGRAD	SGD-G2
AdaFom	BSGD	MAS	SGDEM
AdaFTRL	C-ADAM	MEKA	SGDHess
Adagrad	CADA	MTAdam	SGDM
ADAHESSIAN	Cool Momentum	MVRC-1/MVRC-2	SGDR
Adai	CProp	Nadam	SHAdagrad
AdaLoss	Curveball	NAMSB/NAMSG	Shampoo
Adam	Dadam	ND-Adam	SignAdam++
Adam+	DeepMemory	Nero	SignSGD
AdamAL	DGNOpt	Nesterov	SKQN/S4QN
AdaMax	DiffGrad	Noisy Adam/Noisy K-FAC	SMS
AdamBS	EAdam	NosAdam	SMG
AdamNC	EKFAC	Novograd	SNGM
AdaMod	Eve	NT-SGD	SoftAdam
AdamP/SGDP	Expectigrad	Padam	SRSGD
AdamT	FastAdaBelief	PAGE	Step-Tuned SGD
AdamW	FRSGD	PAL	SWATS
AdamX	G-AdaGrad	PolyAdam	SWNTS
ADAS	GADAM	Polyak	TAdam
AdaS	Gadam	PowerSGD/PowerSGDM	TEKFAC
AdaScale	GOALS	Probabilistic Polyak	VAdam
AdaSGD	GOLS-1	ProbLS	VR-SGD
AdaShift	Grad-Avg	PStorm	vSGD-b/vSGD-g/vSGD-l
AdaSqrt	GRAPES	QHAdam/QHM	vSGD-fd
Adathm	Gravilon	RAdam	WNGrad
AdaX/AdaX-W	Gravity	Ranger	YellowFin
AEGD	HAdam	RangerLars	Yogi
ALI-G	HyperAdam		

Benchmarks

Tuning hyperparameters

SGD:

$$\theta = \theta - \eta g$$

ADAM:

$$m = \beta_1 m + (1 - \beta_1)g$$

$$v = \beta_2 v + (1 - \beta_2)g \odot g$$

$$\hat{m} = \frac{m}{1 - \beta_1^t} \quad \hat{v} = \frac{v}{1 - \beta_2^t}$$

$$\theta = \theta - \frac{\eta}{\sqrt{\hat{v}} + \epsilon} \odot \hat{m}$$

Debugging Tools

The How, What & Why of Deep Learning Optimization

Making neural network training more user-friendly in three steps



How can we fairly compare deep learning optimizers? DEEPOBS is a Python package for benchmarking optimizers.



What are the best deep learning optimization methods? An extensive empirical comparison of fifteen popular deep learning optimizers.




Why do/don't deep learning optimizers work? COCKPIT is a visual and statistical debugger specifically designed for deep learning.



The How, What & Why of Deep Learning Optimization

Making neural network training more user-friendly in three steps

The logo for DeepOBS, featuring the text "DeepOBS" in white on a dark grey background with a faint pattern of neural network connections.

DeepOBS

How can we fairly compare deep learning optimizers? DEEPOBS is a Python package for benchmarking optimizers.

The logo for Benchmark, featuring the text "Benchmark" in white on a light grey background with a faint neural network diagram.

Benchmark

What are the best deep learning optimization methods? An extensive empirical comparison of fifteen popular deep learning optimizers.

The logo for Cockpit, featuring the text "Cockpit" in white on a light grey background with a faint dashboard-like interface.

Cockpit

Why do/don't deep learning optimizers work? COCKPIT is a visual and statistical debugger specifically designed for deep learning.






The Status Quo of Benchmarking Deep Learning Optimizers



Why we need an independent & standardized optimizer benchmark protocol

DEEPOBS (Schneider, Balles, et al. 2019)

Deep learning benchmarks require time

-  Use of rather small test problems
-  Use only a few test problems
-  Repeated work & incomparable results

Misaligned incentives




-  Own optimizer gets more attention than the competition
-  Potential for cherry-picking

The Status Quo of Benchmarking Deep Learning Optimizers



Why we need an independent & standardized optimizer benchmark protocol

DEEPOBS (Schneider, Balles, et al. 2019)

Deep learning benchmarks require time

-  Use of rather small test problems
-  Use only a few test problems
-  Repeated work & incomparable results

Misaligned incentives

-  Own optimizer gets more attention than the competition
-  Potential for cherry-picking

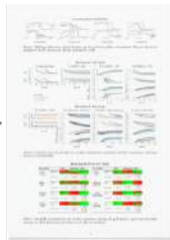
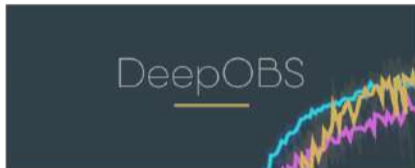
Solution: Independent & standardized benchmark protocol

DEEPOBS: The Deep Learning Optimizer Benchmark Suite

Fairer, faster & more comparable evaluations with less work



DEEPOBS (Schneider, Balles, et al. 2019)



1. Run the optimizer automatically on multiple meaningful test problems.
2. Compare to the provided state of the art without extra costs.
3. Plot the results & the standardized summary evaluation.

The How, What & Why of Deep Learning Optimization

Making neural network training more user-friendly in three steps



DeepOBS

How can we fairly compare deep learning optimizers? DEEPOBS is a Python package for benchmarking optimizers.



Benchmark

What are the best deep learning optimization methods? An extensive empirical comparison of fifteen popular deep learning optimizers.



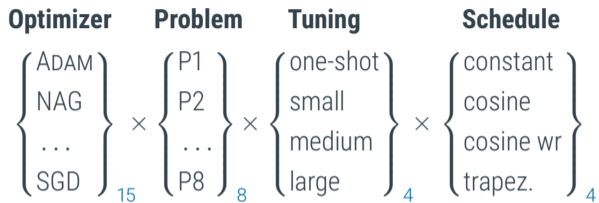
Cockpit

Why do/don't deep learning optimizers work? COCKPIT is a visual and statistical debugger specifically designed for deep learning.



The Benchmarking Process

Many dimensions to explore at once



The Benchmarking Process

Many dimensions to explore at once



Benchmark (Schmidt et al. 2021)



● AMSBOUND

● AMSGRAD

● ADABELIEF

● ADABOUND

● ADADELTA

● ADAGRAD

● ADAM

● LA(MOM.)

● LA(RADAM)

● MOMENTUM

● NAG

● NADAM

● RADAM

● RMSPROP

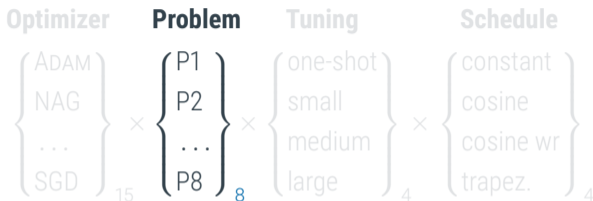
● SGD

The Benchmarking Process

Many dimensions to explore at once



Benchmark (Schmidt et al. 2021)



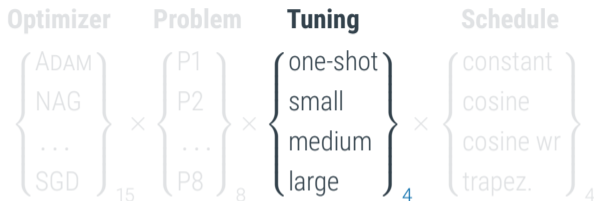
	Data set	Model	Task	Metric	Batch size	Budget <i>in epochs</i>	Approx. run time
P1	Artificial	Noisy quadratic	Minimization	Loss	128	100	< 1 min
P2	MNIST	VAE	Generative	Loss	64	50	10 min
P3	Fashion-MNIST	Simple CNN: 2c2d	Classification	Accuracy	128	100	20 min
P4	CIFAR-10	Simple CNN: 3c3d	Classification	Accuracy	128	100	35 min
P5	Fashion-MNIST	VAE	Generative	Loss	64	100	20 min
P6	CIFAR-100	All-CNN-C	Classification	Accuracy	256	350	4 h 00 min
P7	SVHN	Wide ResNet 16-4	Classification	Accuracy	128	160	3 h 30 min
P8	War and Peace	RNN	Character Prediction	Accuracy	50	200	5 h 30 min

The Benchmarking Process

Many dimensions to explore at once



Benchmark (Schmidt et al. 2021)



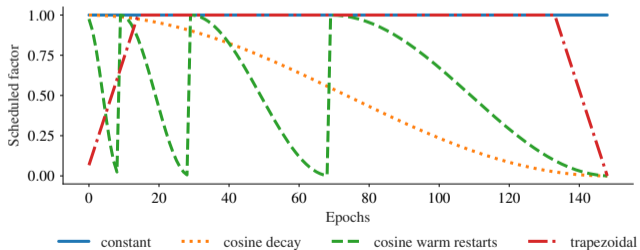
- ▶ **One-Shot** - 1 Run
No tuning, uses default hyperparameters
- ▶ **Small** - 25 Runs
Tuned via random search
- ▶ **Medium** - 50 Runs
Tuned via random search, superset of *small budget*
- ▶ **Large** - 75 Runs
Tuned via random search, refined search spaces

The Benchmarking Process

Many dimensions to explore at once

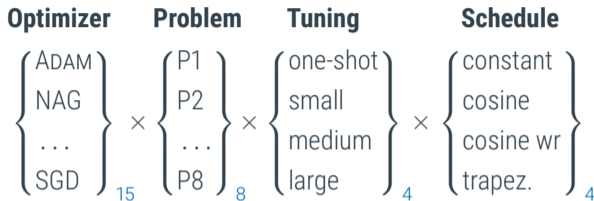


Benchmark (Schmidt et al. 2021)



The Benchmarking Process

Many dimensions to explore at once



1,920 **Configurations**

fifteen optimizers, eight problems, four budgets, & four schedules

>50,000 **Individual Runs**

includes tuning runs & runs stochastic fidelity

128 **Settings** each optimizer is tested in

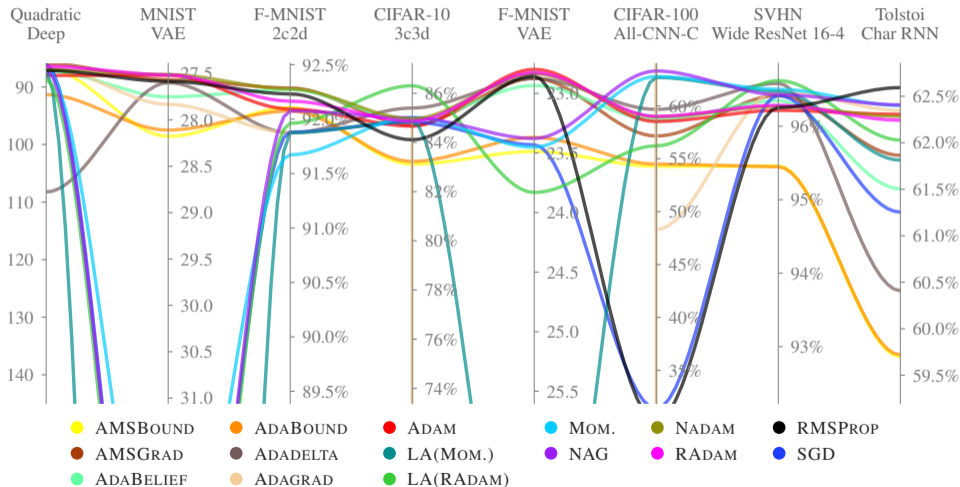
eight problems, four budgets, & four schedules

The Benchmarking Results



Why ADAM is still a good choice | *large budget with a trapezoidal schedule*

Benchmark (Schmidt et al. 2021)

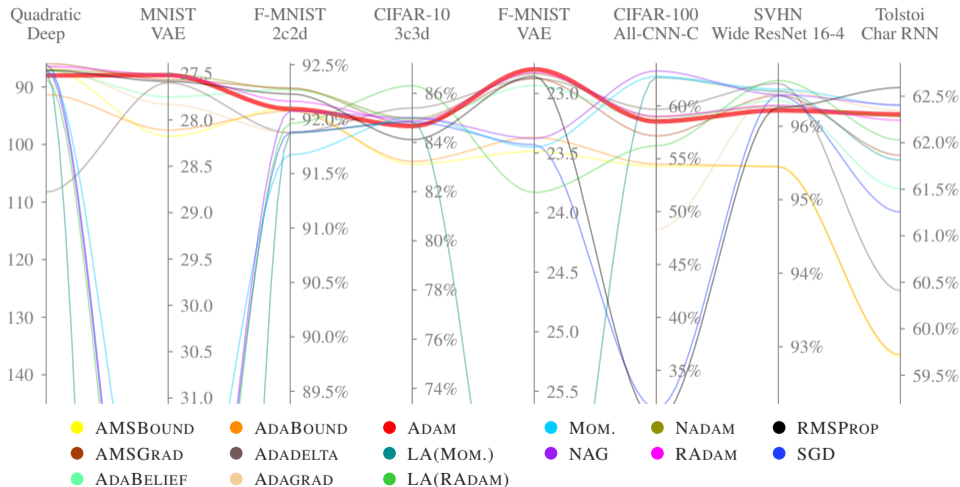


The Benchmarking Results



Why ADAM is still a good choice | *large budget with a trapezoidal schedule*

Benchmark (Schmidt et al. 2021)

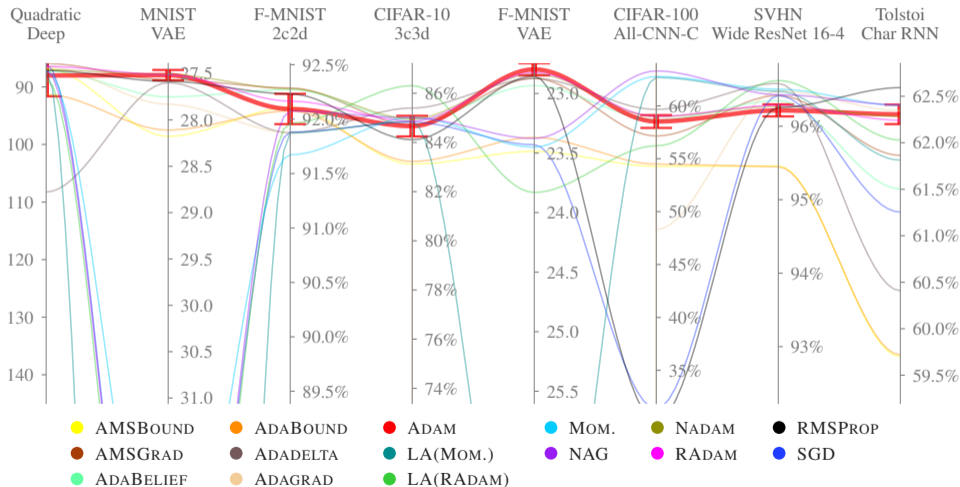


The Benchmarking Results



Why ADAM is still a good choice | *large budget with a trapezoidal schedule*

Benchmark (Schmidt et al. 2021)

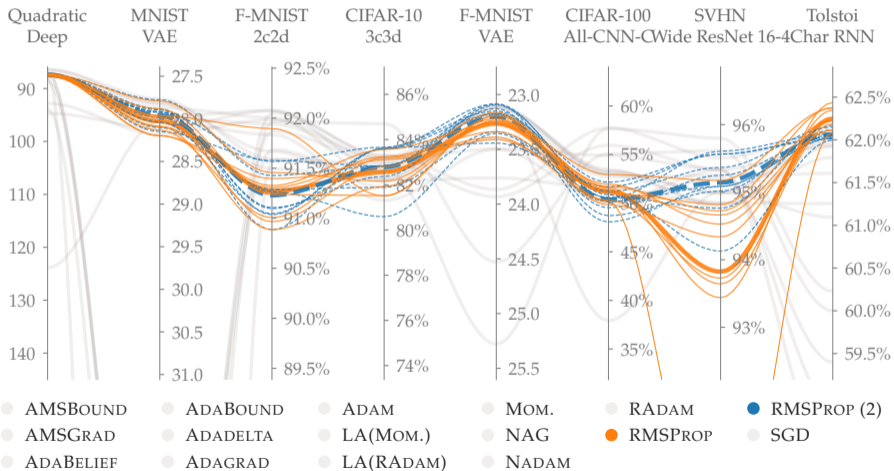


The Benchmarking Results



Individual results can change after re-tuning but overall trends remain

Benchmark (Schmidt et al. 2021)



MLCOMMONS

A natural extension of our benchmarking efforts

ML ● Commons Algorithms Working Group

Chairs:



George Dahl

Google



Frank Schneider

University of Tübingen

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

MLCOMMONS

Benchmarking more than just optimizers

Training Algorithm Track

- ▶ **Tests training algorithms** (update rules, data selection, hyperparameter spaces, etc.)
- ▶ **Scored by “time-to-result”** on multiple tasks
- ▶ **No manual workload-specific adaptation** is allowed

Model Track

- ▶ **Tests models** (architecture, initialization, data augmentation, etc.)
- ▶ **Scored by “time-to-result”** on a single task
- ▶ **Uses standard optimizer** and tuning procedures

MLCOMMONS

Three ways to interact and contribute!

- ▶ **Join the Working Group!**

Consider contributing to the working group to improve the benchmark.

- ▶ **Submit to the Competition!**

Consider submitting your novel algorithm to our competition.

- ▶ **Check the Results!**

Consider reading our (future) summary of the competition.

The How, What & Why of Deep Learning Optimization

Making neural network training more user-friendly in three steps



DeepOBS

How can we fairly compare deep learning optimizers? DEEPOBS is a Python package for benchmarking optimizers.



Benchmark

What are the best deep learning optimization methods? An extensive empirical comparison of fifteen popular deep learning optimizers.



Cockpit

Why do/don't deep learning optimizers work? COCKPIT is a visual and statistical debugger specifically designed for deep learning.



Deep Learning Requires New Debuggers

As coding is replaced by learning, we need to replace our debuggers



Classic
Programming

Functions

```
def f(x):  
    return x + 1  
  
def g(x):  
    return x * 2  
  
def h(x):  
    return g(f(x))  
  
def i(x):  
    return f(g(x))  
  
def j(x):  
    return h(g(x))  
  
def k(x):  
    return g(h(x))  
  
def l(x):  
    return f(h(x))  
  
def m(x):  
    return h(f(x))  
  
def n(x):  
    return g(f(g(x)))  
  
def o(x):  
    return f(g(f(x)))  
  
def p(x):  
    return h(g(f(x)))  
  
def q(x):  
    return f(h(g(x)))  
  
def r(x):  
    return g(h(f(x)))  
  
def s(x):  
    return h(f(g(x)))  
  
def t(x):  
    return g(f(h(x)))  
  
def u(x):  
    return f(h(g(x)))  
  
def v(x):  
    return h(g(f(x)))  
  
def w(x):  
    return g(h(f(x)))  
  
def x(x):  
    return f(h(g(x)))  
  
def y(x):  
    return h(g(f(x)))  
  
def z(x):  
    return g(h(f(x)))
```

Level of Abstraction

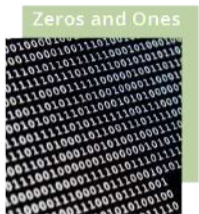
Deep Learning Requires New Debuggers

As coding is replaced by learning, we need to replace our debuggers



COCKPIT (Schneider, Dangel, et al. 2021)

Classic
Programming



Deep Learning Requires New Debuggers

As coding is replaced by learning, we need to replace our debuggers



COCKPIT (Schneider, Dangel, et al. 2021)

Classic
Programming

Zeros and Ones



Debugger



Functions



Level of Abstraction

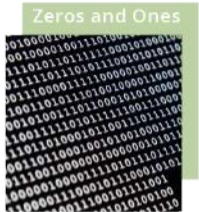
Deep Learning Requires New Debuggers

As coding is replaced by learning, we need to replace our debuggers

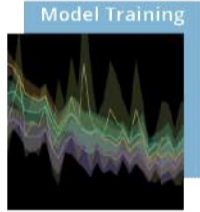


COCKPIT (Schneider, Dangel, et al. 2021)

Classic
Programming



Deep
Learning



Level of Abstraction

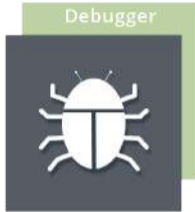
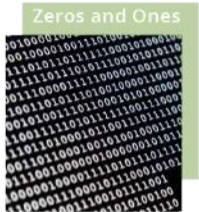
Deep Learning Requires New Debuggers

As coding is replaced by learning, we need to replace our debuggers

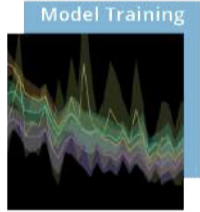


COCKPIT (Schneider, Dangel, et al. 2021)

Classic
Programming



Deep
Learning

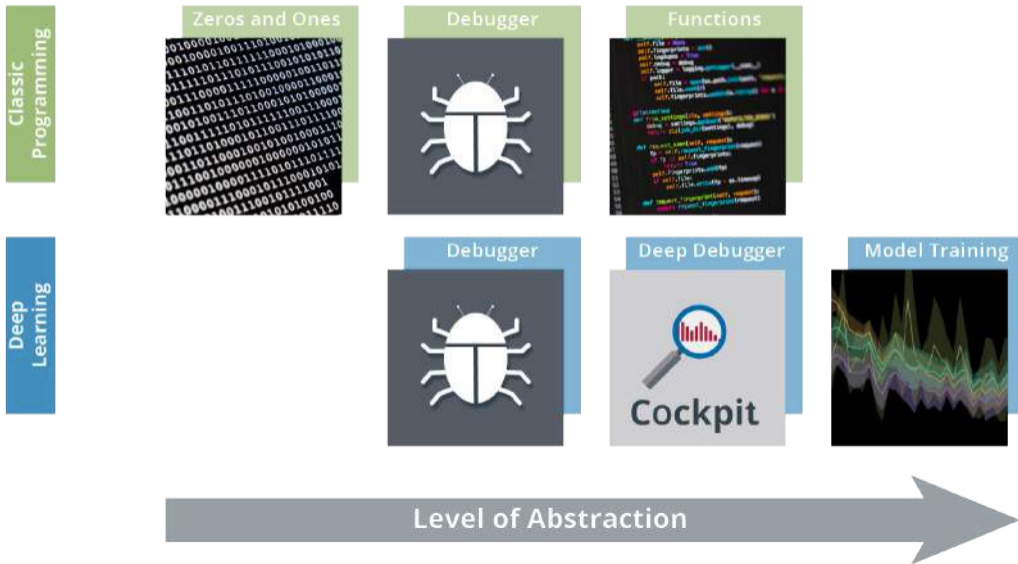


Deep Learning Requires New Debuggers



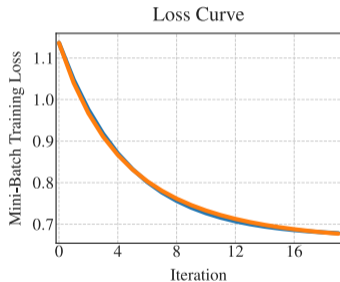
As coding is replaced by learning, we need to replace our debuggers

COCKPIT (Schneider, Dangel, et al. 2021)



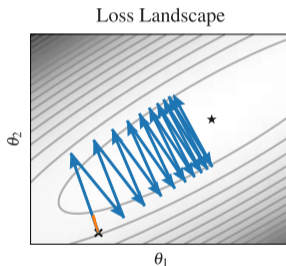
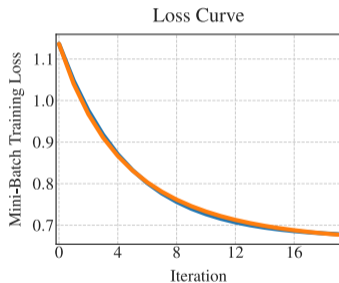
Loss Curves Do Not Tell the Full Story

Why we need better observables in neural network training



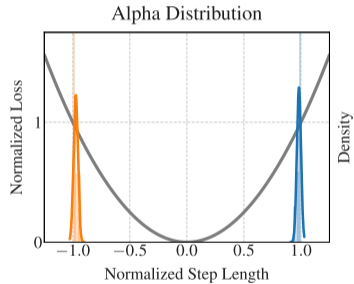
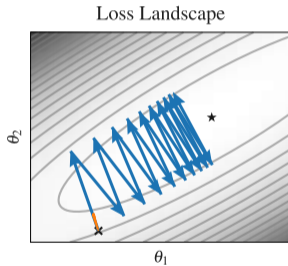
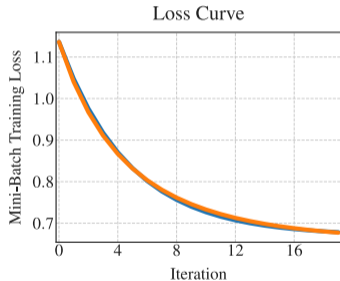
Loss Curves Do Not Tell the Full Story

Why we need better observables in neural network training



Loss Curves Do Not Tell the Full Story

Why we need better observables in neural network training



COCKPIT in Action!

Training the ALL-CNN-C network through the lens of COCKPIT



Summary

Improving Optimizer Evaluation in Deep Learning

- ▶ **DEEPOBS**, a benchmarking suite for deep learning optimizers.
- ▶ An empirical comparison of fifteen optimization methods for deep learning.
- ▶ **COCKPIT**, a visual debugging tool for deep learning.

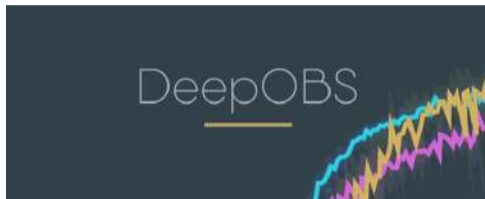


DEEPOBS: `pip install deepobs`

Benchmark: github.com/SirRob1997/Crowded-Valley---Results

COCKPIT: `pip install cockpit-for-pytorch`





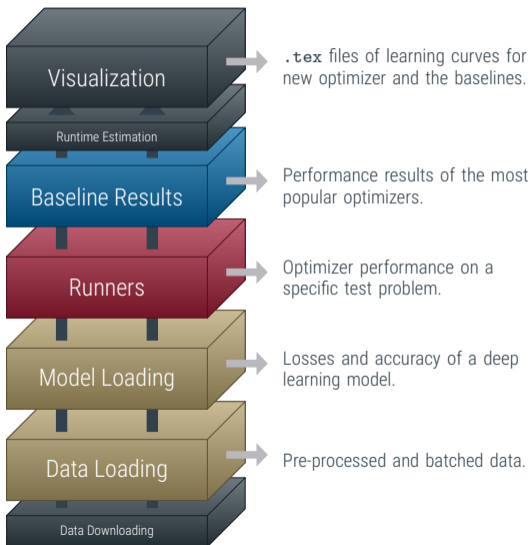
DeepOBS

The DEEPOBS Stack

A modular environment for every step of evaluating deep learning optimizers



DEEPOBS (Schneider, Balles, et al. 2019)



The DEEPOBS Test Problems

A wide selection of varied problems with room to grow.

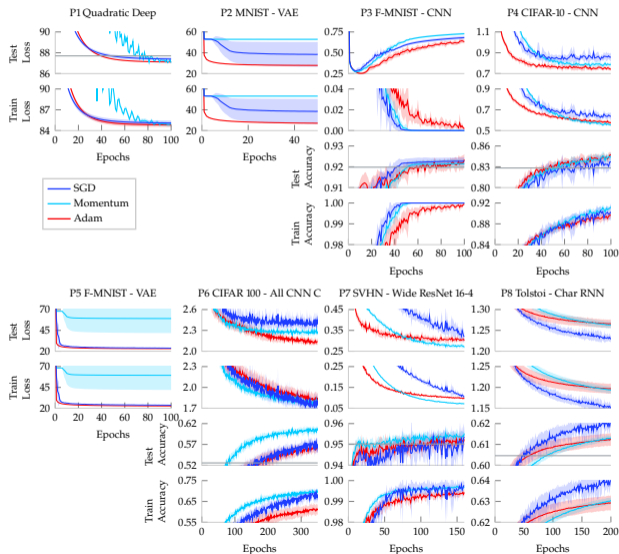
Data set	Model	Description	Conv	RNN	Drop	BN	L ²
● 2D	Noisy Beale	Noisy version of the Beale function (Beale 1958)					
	Noisy Branin	Noisy version of the Branin function (Branin 1972)					
	Noisy Rosenbrock	Noisy version of the Rosenbrock function (Rosenbrock 1960)					
● Quadratic	Deep	100-dimensional ill-conditioned noisy quadratic (Chaudhari et al. 2017)					
● MNIST (LeCun et al. 1998)	Log. Regr.	Logistic regression					
	MLP	Four layer fully connected network					
	2c2d	Two conv. and two fully connected layers	✓				
	VAE	Variational Autoencoder	✓		✓		
● FASHION ● MNIST (Xiao et al. 2017)	Log. Regr.	Logistic regression					
	MLP	Four layer fully connected network					
	2c2d	Two conv. and two fully connected layers	✓				
	VAE	Variational Autoencoder	✓		✓		
● CIFAR-10 (Krizhevsky et al. 2009)	3c3d	Three conv. and three fully connected layers	✓				✓
	VGG16	Adapted version of VGG16 (Simonyan et al. 2015)	✓		✓		✓
	VGG19	Adapted version of VGG19	✓		✓		✓
	3c3d	Three conv. and three fully connected layers	✓				✓
● CIFAR-100 (Krizhevsky et al. 2009)	VGG16	Adapted version of VGG16	✓		✓		✓
	VGG19	Adapted version of VGG19	✓		✓		✓
	All-CNN-C	The all convolutional net from Springenberg et al. (2015)	✓		✓		✓
	Wide ResNet-40-4	Wide Residual Network (Zagoruyko et al. 2016)	✓			✓	✓
● SVHN (Netzer et al. 2011)	3c3d	Three conv. and three fully connected layers	✓				✓
	Wide ResNet-16-4	Wide Residual Network	✓			✓	✓
● IMAGENET (Deng et al. 2009)	VGG16	VGG16	✓		✓		✓
	VGG19	VGG19	✓		✓		✓
	Inception-v3	Inception-v3 network as described by Szegedy et al. (2016)	✓		✓	✓	✓
● Tolstói	CharRNN	Recurrent Neural Network for character-level language modeling		✓	✓		

Results Output of DEEPOBS

A showcase of DEEPOBS' capabilities



DEEPOBS (Schneider, Balles, et al. 2019)

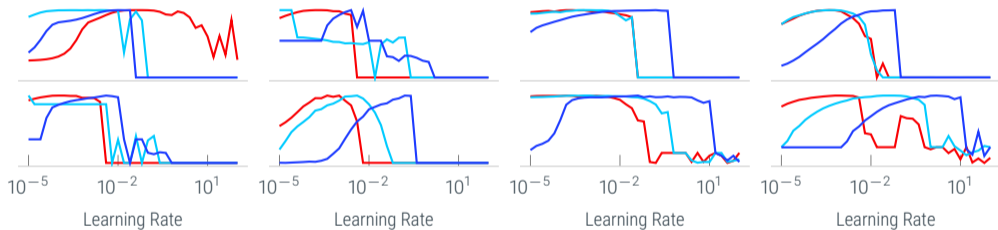


Results Output of DEEPOBS

A showcase of DEEPOBS' capabilities



DEEPOBS (Schneider, Balles, et al. 2019)

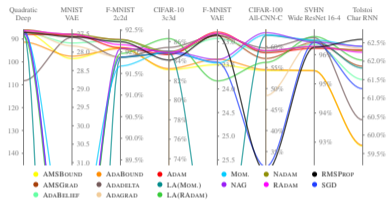


Results Output of DEEPOBS

A showcase of DEEPOBS' capabilities

Test Problem		SGD	Momentum	Adam
P1 Quadratic Deep	Performance	87.40	87.05	87.11
	Speed	51.1	70.5	39.9
	Tuneability	η : 1.58e-02	η : 2.51e-03 μ : 0.99	η : 3.98e-02 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08
P2 MNIST VAE	Performance	38.46	52.93	27.83
	Speed	1.0	1.0	1.0
	Tuneability	η : 3.98e-03	η : 2.51e-05 μ : 0.99	η : 1.58e-04 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08
P3 F-MNIST CNN	Performance	92.27 %	92.14 %	92.34 %
	Speed	40.6	59.1	40.1
	Tuneability	η : 1.58e-01	η : 2.51e-03 μ : 0.99	η : 2.51e-04 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08
P4 CIFAR-10 CNN	Performance	83.71 %	84.41 %	84.75 %
	Speed	42.5	40.7	36.0
	Tuneability	η : 6.31e-02	η : 3.98e-04 μ : 0.99	η : 3.98e-04 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08

Test Problem		SGD	Momentum	Adam
P5 F-MNIST VAE	Performance	23.80	59.23	23.07
	Speed	1.0	1.0	1.0
	Tuneability	η : 3.98e-03	η : 1.00e-05 μ : 0.99	η : 1.58e-04 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08
P6 CIFAR-100 All CNN C	Performance	57.06 %	60.33 %	56.15 %
	Speed	128.7	72.8	152.6
	Tuneability	η : 1.58e-01	η : 3.98e-03 μ : 0.99	η : 1.00e-03 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08
P7 SVHN Wide ResNet	Performance	95.37 %	95.53 %	95.25 %
	Speed	28.3	10.8	12.1
	Tuneability	η : 2.51e-02	η : 6.31e-04 μ : 0.99	η : 1.58e-04 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08
P8 Tolstoj Char RNN	Performance	62.07 %	61.30 %	61.23 %
	Speed	47.7	88.0	62.8
	Tuneability	η : 1.58e+00	η : 3.98e-02 μ : 0.99	η : 2.51e-03 β_1 : 0.9 β_2 : 0.999 ϵ : 1e-08



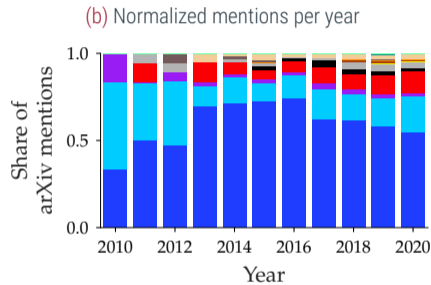
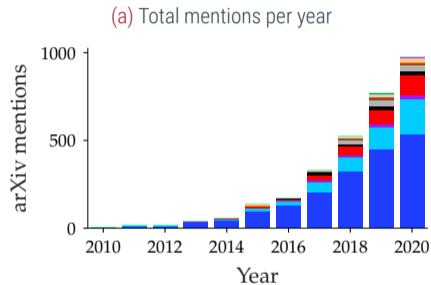
Benchmark

arXiv Mentions of Optimizers per Year



Our selected optimizers cover the most popular choices of the increasing number of methods

Benchmark (Schmidt et al. 2021)

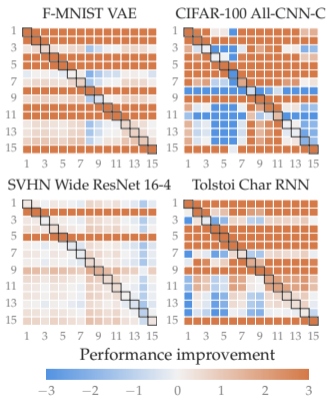
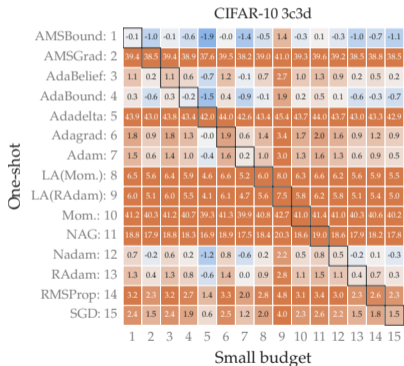


Results: Out-of-the-box Performance

Why trying out multiple optimizers can be better than tuning them



Benchmark (Schmidt et al. 2021)



Results: Out-of-the-box Performance



Why trying out multiple optimizers can be better than tuning them

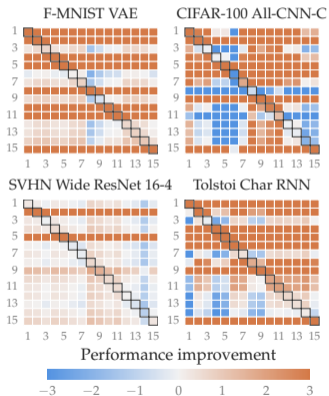
Benchmark (Schmidt et al. 2021)

- **Orange rows**
bad default hyperparameters
SGD, NAG, MOMENTUM,
AMSGRAD, ADADELTA

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
Adadelta: 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: Out-of-the-box Performance

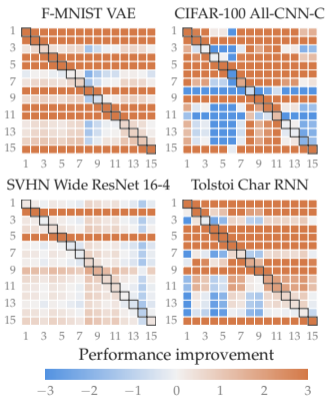
Why trying out multiple optimizers can be better than tuning them

Benchmark (Schmidt et al. 2021)

- ▶ **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
Adadelta: 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
RAadam: 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

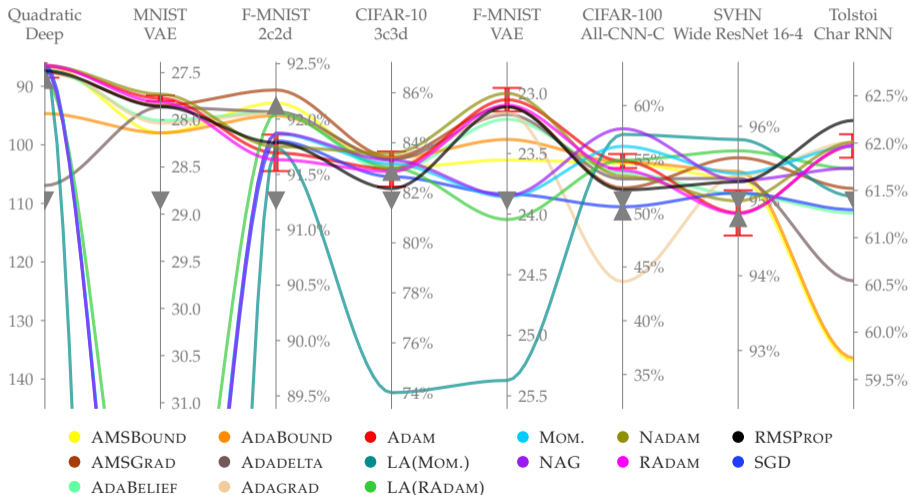


Results: Which Optimizer Should You Pick?



ADAM is still a reasonable choice | large budget without a learning rate schedule

Benchmark (Schmidt et al. 2021)

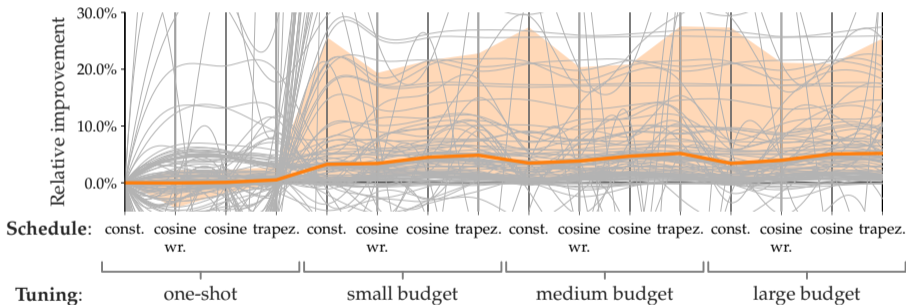


Results: Effect of Tuning and Schedules



Increasing the budget improves the performance on average with diminishing returns

Benchmark (Schmidt et al. 2021)

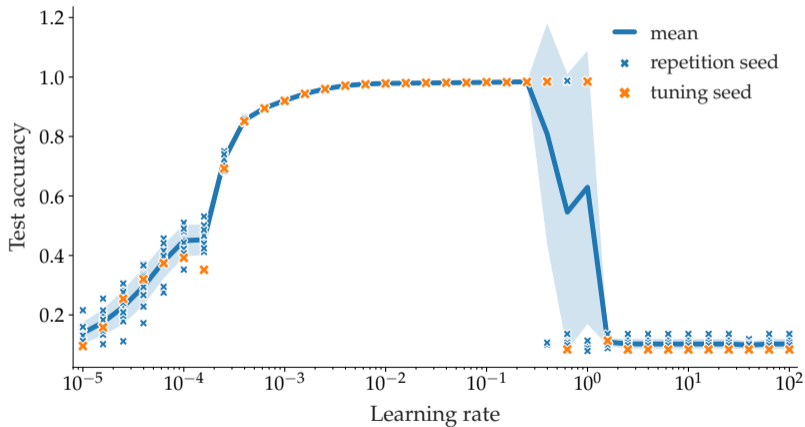


Results: Robustness to Random Seeds



Single seed tuning can result in hyperparameters in an unstable “danger zone” (<0.5% of cases)

Benchmark (Schmidt et al. 2021)

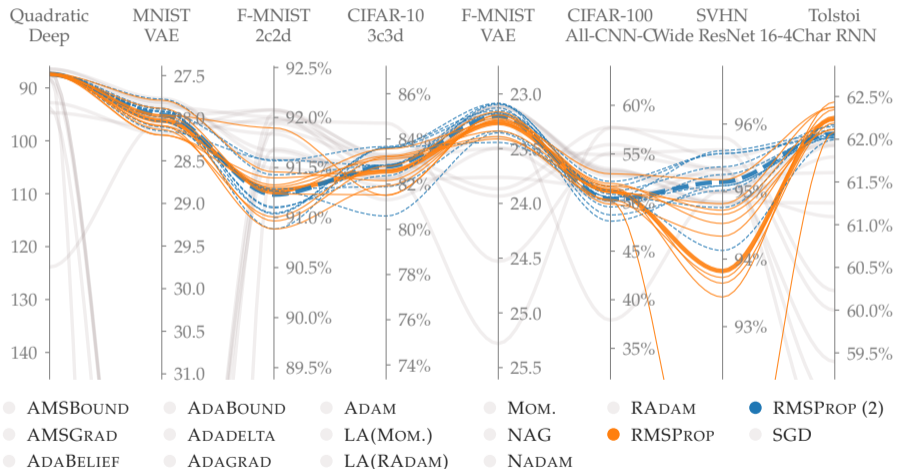


Results: Re-tuning Experiment



Individual results can change after re-tuning but overall trends remain

Benchmark (Schmidt et al. 2021)

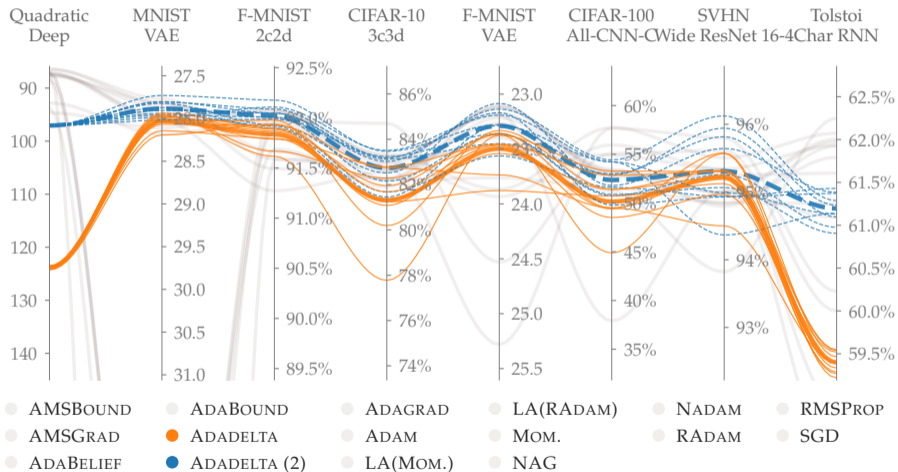


Results: Re-tuning Experiment



Individual results can change after re-tuning but overall trends remain

Benchmark (Schmidt et al. 2021)





Cockpit

List of Quantities in COCKPIT

A wide selection to choose from and easy to extend

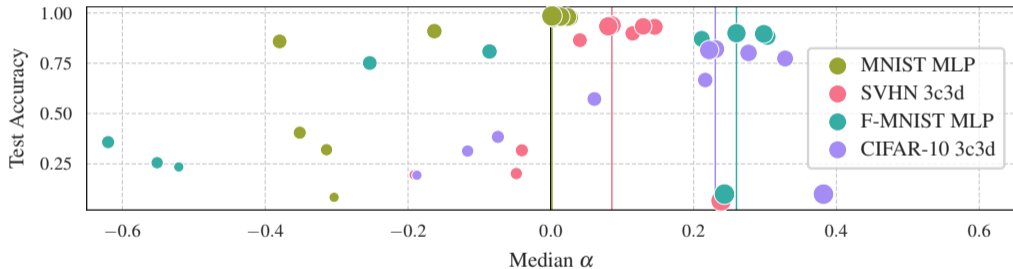
Name	Description	Configuration
Loss	Mini-batch training loss at current iteration, $L_{\mathbb{B}}$	<i>economy</i>
Distance	L^2 distance from initialization $\ \boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^{(0)}\ _2$	<i>economy</i>
UpdateSize	Update size of the current iteration $\ \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}\ _2$	<i>economy</i>
GradNorm	Mini-batch gradient norm $\ \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\ _2$	<i>economy</i>
Alpha	Normalized step on a noise-informed quadratic interpolation between $\boldsymbol{\theta}^{(t)}$, $\boldsymbol{\theta}^{(t+1)}$	<i>economy</i>
GradHist1d	Histogram of individual gradient elements, $\{\mathbf{g}_{\mathbb{B}}^{(i)}(\boldsymbol{\theta}_j)\}_{j=1, \dots, D}$	<i>economy</i>
NormTest	Normalized fluctuations of the residual norms $\ \mathbf{g}_{\mathbb{B}} - \mathbf{g}_{\mathbb{B}}^{(i)}\ $, proposed in Byrd et al. 2012	<i>economy</i>
InnerTest	Normalized fluctuations of $\mathbf{g}_{\mathbb{B}}^{(i)}$'s parallel components along $\mathbf{g}_{\mathbb{B}}$, proposed in Bollapragada et al. 2017	<i>economy</i>
OrthoTest	Normalized fluctuations of $\mathbf{g}_{\mathbb{B}}^{(i)}$'s orthogonal components along $\mathbf{g}_{\mathbb{B}}$, proposed in Bollapragada et al. 2017	<i>economy</i>
HessTrace	Exact or approximate Hessian trace, $\text{Tr}(\mathbf{H}_{\mathcal{B}}(\boldsymbol{\theta}))$, inspired by Yao et al. 2020	<i>business</i>
TICDiag	Relation between (diagonal) curvature and gradient noise, inspired by Thomas et al. 2020	<i>business</i>
HessMaxEV	Maximum Hessian eigenvalue, $\lambda_{\max}(\mathbf{H}_{\mathcal{B}}(\boldsymbol{\theta}))$, inspired by Yao et al. 2020	<i>full</i>
GradHist2d	Histogram of weights and individual gradient elements, $\{(\boldsymbol{\theta}_j, \mathbf{g}_{\mathbb{B}}^{(i)}(\boldsymbol{\theta}_j))\}_{j=1, \dots, D}$	<i>full</i>
Parameters	Parameter values $\boldsymbol{\theta}^{(t)}$ at the current iteration	-
Time	Time of the current iteration	-
CABS	Adaptive batch size for SGD, optimizes expected objective gain per cost, adapted from Balles et al. 2017	-
EarlyStopping	Evidence-based early stopping criterion for SGD, proposed in Mahsereci et al. 2017	-
TICTrace	Relation between curvature and gradient noise trace, inspired by Thomas et al. 2020	-
MeanGSNR	Average gradient signal-to-noise-ratio (GSNR), inspired by Liu et al. 2020	-

COCKPIT as a Tool for Inspiring Research

Neural network training requires systematically overstepping the local minima

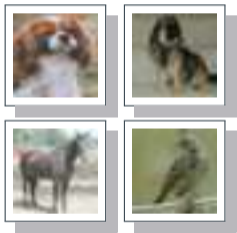


COCKPIT (Schneider, Dangel, et al. 2021)

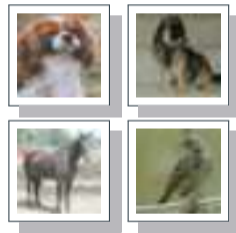


COCKPIT as a Tool for Finding Data Bugs

Normalized and raw data can *look* the same but *behave* differently



(a) Normalized Data



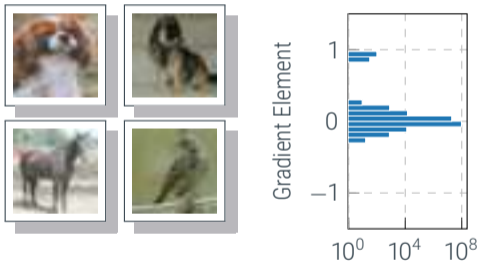
(b) Raw Data

COCKPIT as a Tool for Finding Data Bugs

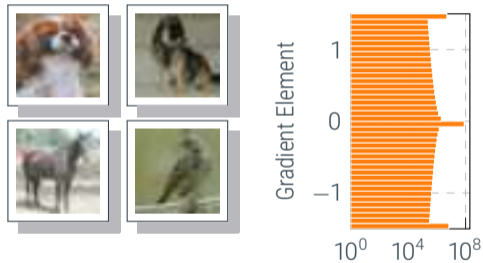


Normalized and raw data can *look* the same but *behave* differently

COCKPIT (Schneider, Dangel, et al. 2021)



(a) Normalized Data



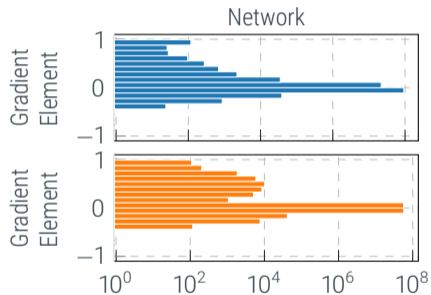
(b) Raw Data

COCKPIT as a Tool for Finding Model Bugs

Looking at the gradients layerwise can reveal model inefficiencies



COCKPIT (Schneider, Dangel, et al. 2021)



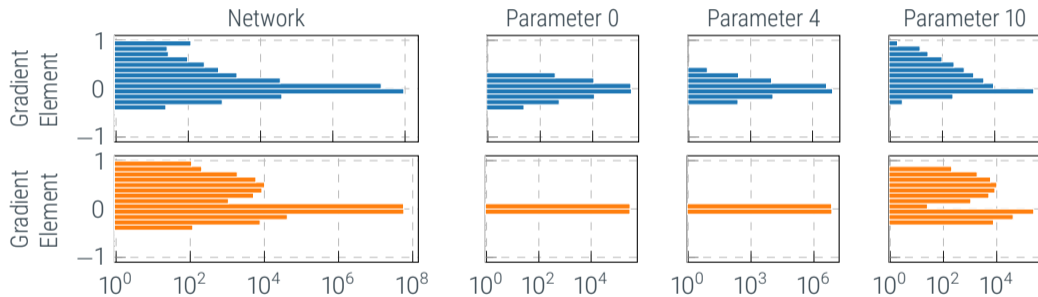
(a) Network Histogram

COCKPIT as a Tool for Finding Model Bugs

Looking at the gradients layerwise can reveal model inefficiencies



COCKPIT (Schneider, Dangel, et al. 2021)



(a) Network Histogram

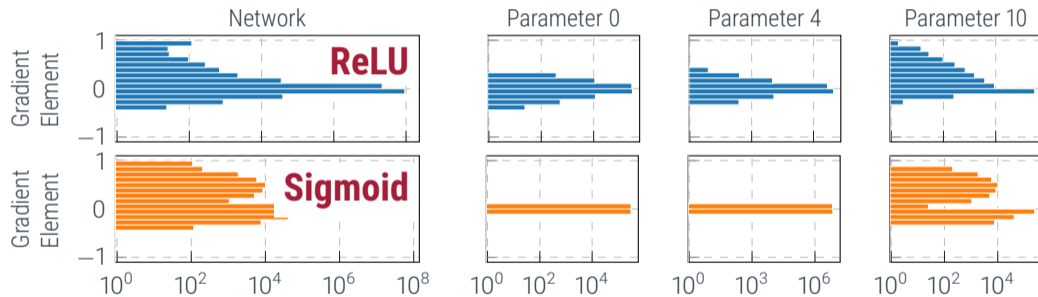
(b) Layer-wise Histograms

COCKPIT as a Tool for Finding Model Bugs

Looking at the gradients layerwise can reveal model inefficiencies



COCKPIT (Schneider, Dangel, et al. 2021)



(a) Network Histogram

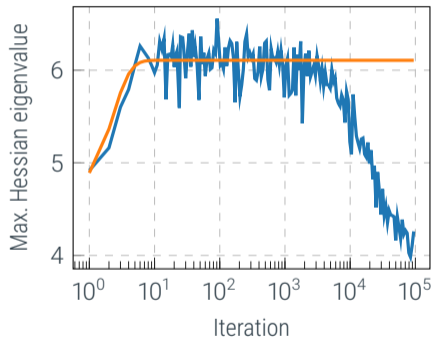
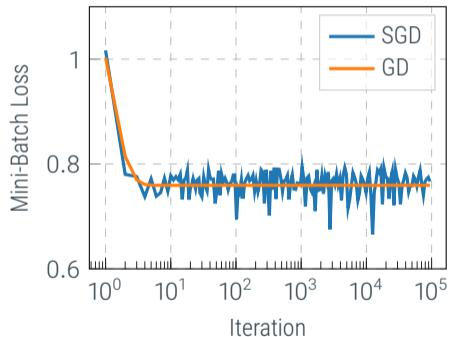
(b) Layer-wise Histograms

Observing Implicit Regularization with COCKPIT

An effect that cannot be seen by monitoring the only the loss



COCKPIT (Schneider, Dangel, et al. 2021)

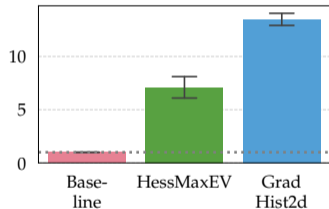
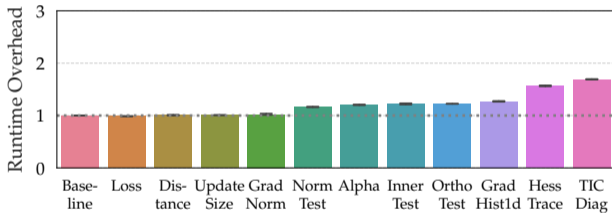


Efficient Implementation for Practical Overhead



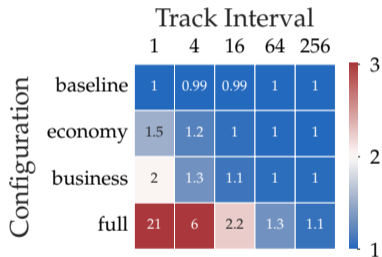
CockPIT's quantities are affordable in practice: CIFAR-10 and 3c3d

CockPIT (Schneider, Dangel, et al. 2021)



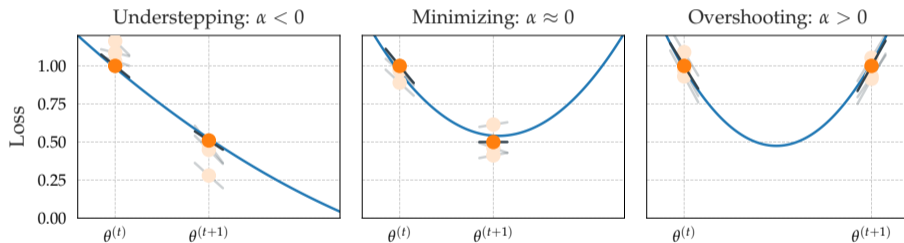
Efficient Implementation for Practical Overhead

Combining quantities and reducing the tracking interval reduces the overhead



Illustrating the Alpha Quantity

Are we understepping or overshooting?



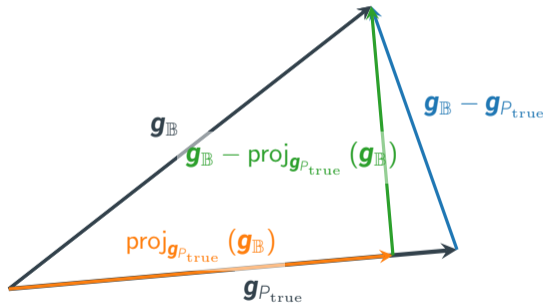
Illustrating the Gradient Tests



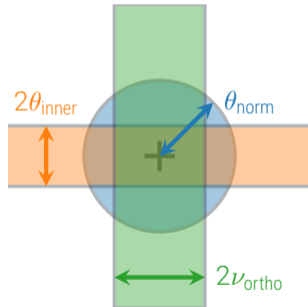
Defining geometric constraints between the mini-batch and the true gradient

COCKPIT (Schneider, Dangel, et al. 2021)

(a) Gradient test vectors



(b) Visualization in COCKPIT





- ▶ Balles, Lukas, Javier Romero, and Philipp Hennig (2017). "Coupling Adaptive Batch Sizes with Learning Rates". In: *Uncertainty in Artificial Intelligence - Proceedings of the 33rd Conference, UAI*.
- ▶ Beale, Evelyn M. L. (1958). *On an iterative method for finding a local minimum of a function of more than one variable*. Statistical Techniques Research Group, Section of Mathematical Statistics, Department of Mathematics, Princeton University.
- ▶ Bollapragada, Raghu, Richard Byrd, and Jorge Nocedal (2017). "Adaptive Sampling Strategies for Stochastic Optimization". In: *SIAM Journal on Optimization*.
- ▶ Branin, Franklin H. (1972). "Widely convergent method for finding multiple solutions of simultaneous nonlinear equations". In: *IBM Journal of Research and Development*.
- ▶ Byrd, Richard H. et al. (2012). "Sample Size Selection in Optimization Methods for Machine Learning". In: *Math. Program.*
- ▶ Chaudhari, Pratik et al. (2017). "Entropy-SGD: Biasing gradient descent into wide valleys". In: *5th International Conference on Learning Representations, ICLR*.
- ▶ Deng, Jia et al. (2009). "ImageNet: A Large-Scale Hierarchical Image Database". In: *The IEEE Conference on Computer Vision and Pattern Recognition, CVPR*.



- ▶ Krizhevsky, Alex and Geoffrey Hinton (2009). *Learning multiple layers of features from tiny images*. Technical Report.
- ▶ LeCun, Yann et al. (1998). "Gradient-Based Learning Applied to Document Recognition". In: *Proceedings of the IEEE*.
- ▶ Liu, Jinlong et al. (2020). "Understanding Why Neural Networks Generalize Well Through GSNR of Parameters". In: *8th International Conference on Learning Representations, ICLR*.
- ▶ Mahsereci, Maren et al. (2017). *Early Stopping without a Validation Set*. arXiv: 1703.09580.
- ▶ Netzer, Yuval et al. (2011). "Reading digits in natural images with unsupervised feature learning". In: *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*.
- ▶ Rosenbrock, Howard H. (1960). "An automatic method for finding the greatest or least value of a function". In: *The Computer Journal*.
- ▶ Schmidt, Robin M., Frank Schneider, and Philipp Hennig (2021). "Descending through a Crowded Valley - Benchmarking Deep Learning Optimizers". In: *38th International Conference on Machine Learning, ICML*.
- ▶ Schneider, Frank, Lukas Balles, and Philipp Hennig (2019). "DeepOBS: A Deep Learning Optimizer Benchmark Suite". In: *7th International Conference on Learning Representations, ICLR*.



- ▶ Schneider, Frank, Felix Dangel, and Philipp Hennig (2021). "Cockpit: A Practical Debugging Tool for the Training of Deep Neural Networks". In: *Advances in Neural Information Processing Systems 34, NeurIPS*.
- ▶ Simonyan, Karen and Andrew Zisserman (2015). "Very Deep Convolutional Networks for Large-Scale Image Recognition". In: *3rd International Conference on Learning Representations, ICLR*.
- ▶ Springenberg, Jost T. et al. (2015). "Striving for simplicity: The all convolutional net". In: *3rd International Conference on Learning Representations, ICLR (workshop track)*.
- ▶ Szegedy, Christian et al. (2016). "Rethinking the Inception Architecture for Computer Vision". In: *The IEEE Conference on Computer Vision and Pattern Recognition, CVPR*.
- ▶ Thomas, Valentin et al. (2020). "On the interplay between noise and curvature and its effect on optimization and generalization". In: *23rd International Conference on Artificial Intelligence and Statistics, AISTATS*.
- ▶ Xiao, Han, Kashif Rasul, and Roland Vollgraf (2017). *Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms*. arXiv: 1708.07747.
- ▶ Yao, Zhewei et al. (2020). *ADAHESIAN: An Adaptive Second Order Optimizer for Machine Learning*. arXiv: 2006.00719.



- ▶ Zagoruyko, Sergey and Nikos Komodakis (2016). "Wide Residual Networks". In: *Proceedings of the British Machine Vision Conference*.