# Neural Network Training through the Lens of Benchmarking and Debugging

## SPP 2353 - Summer School

**Frank Schneider**

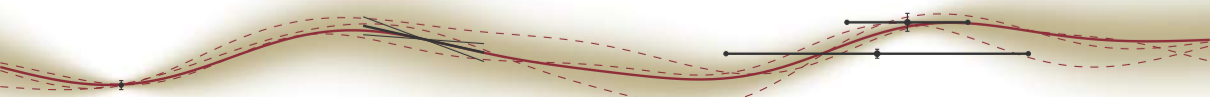April 4, 2023

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

erc

Federal Ministry
of Education
and Research

imprs-is

**LLMs**
ChatGPT/GPT-4



**Generative AI**
Stable Diffusion/DALL-E/Midjourney



**AI4Science**
AlphaFold



**SPP 2353**

# How do we train neural networks?

# 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 OPT Training Logbook
A real-world case study of how hard it is to train a neural network

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Zhang et al. 2022)

**OPT: Open Pre-trained Transformer Language Models**

Susan Zhang, Stephen Roller, Naman Goyal,
Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li,
Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig,
Punit Singh Koura, Anjali Sridhar, Tianlu Wang, Luke Zettlemoyer

Meta AI

{susanz, roller, naman}@fb.com

► Meta AI trained an open-source GPT-3-like large language model (LLM) called OPT.

# The OPT Training Logbook

A real-world case study of how hard it is to train a neural network

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Zhang et al. 2022)

- ▶ Meta AI trained an open-source GPT-3-like large language model (LLM) called OPT.
- ▶ Training clearly requires more than just "ADAM with a default learning rate".

We use an AdamW optimizer (Loschilov and Hutter, 2017) with $(\beta_1, \beta_2)$ set to $(0.9, 0.95)$, and weight decay of 0.1. We follow a linear learning rate schedule, warming up from 0 to the maximum learning rate over the first 2000 steps in OPT-175B, or over 375M tokens in our smaller baselines, and decaying down to 10% of the maximum LR over 300B tokens. A number of mid-flight changes to LR were also required (see Section 2.5). Our batch sizes range from 0.5M to 4M depending on the model size (see Table 1) and is kept constant throughout the course of training.

We use a dropout of 0.1 throughout, but we do not apply any dropout to embeddings. We clip gradient norms at 1.0, except for some mid-flight changes that reduce this threshold down from 1.0 to 0.3 (see Section 2.5). We also in-

# The OPT Training Logbook

A real-world case study of how hard it is to train a neural network

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Zhang et al. 2022)

Figure 1: **Empirical LR schedule.** We found that lowering learning rate was helpful for avoiding instabilities.

- ▶ Meta AI trained an open-source GPT-3-like large language model (LLM) called OPT.
- ▶ Training clearly requires more than just "ADAM with a default learning rate".
- ▶ The result is a highly empirical and complicated learning rate schedule.

# The OPT Training Logbook
A real-world case study of how hard it is to train a neural network

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Zhang et al. 2022)

**Run is stuck in loop of "lower loss scale"**

Uh oh. The loss exploded. We don't know why.

**Remember to document your actions in the logbook.**

Actions:
1. **Don't panic.**
2. Ping the group chat to discuss potential options. If no one responds make a decision by yourself. Letting nodes idle costs $2500/hour so

▶ Meta AI trained an open-source GPT-3-like large language model (LLM) called OPT.

▶ Training clearly requires more than just "ADAM with a default learning rate".

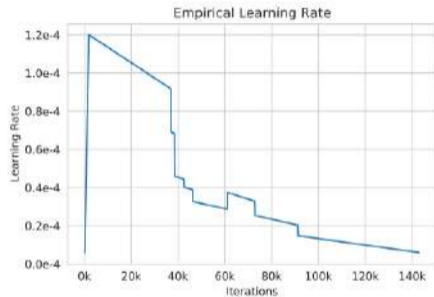▶ The result is a highly empirical and complicated learning rate schedule.

▶ 114-paged logbook detailing the struggles.

Part I
**Fundamentals**

Why training a neural network is hard

# Supervised Learning with Neural Networks

A simplified overview of what sounds like an optimization problem

# Supervised Learning with Neural Networks

A simplified overview of what sounds like an optimization problem

**Input**

**Images**
- Dog
- Cat
- Human

**Audio**

**Text**
- "Hallo?"
- "Please?"
- "I'm sorry"

**Model**

$f_\theta$

**Output**

**Classes**
- Dog
- Cat
- Human

**Numbers**
- 12
- 0.5
- 42

**Text**
- "Hello?"
- "Yes!"
- "No thanks"

1. **Data**
   $\mathbb{D}_{\text{train}} = \{(\boldsymbol{x}^{(i)}, y^{(i)})\}$

2. **Model**
   $f_{\boldsymbol{\theta}}$

3. **Loss**
   $L$

4. **Training Algorithm**
   $\theta_{t+1} = \dots$

**Setting**

We want to learn a function which is able to correctly predict the output $y \in \mathbb{Y}$ given some features $\boldsymbol{x} \in \mathbb{X}$. A **loss function** $\ell(\hat{y}, y) : \mathbb{Y} \times \mathbb{Y} \to \mathbb{R}^+$ quantifies how different the predictions of a model $\hat{y}$ are from the true targets $y$.

**True Risk/Expected Loss**

Given the **true underlying data distribution** $P_{\text{true}}(\boldsymbol{x}, y)$, we want to minimize

$$L_{P_{\text{true}}}(f_{\boldsymbol{\theta}}) = \mathbb{E}_{(\boldsymbol{x},y) \sim P_{\text{true}}(\boldsymbol{x},y)} \left[ \ell\left( f(\boldsymbol{x}; \boldsymbol{\theta}), y \right) \right]$$

**Empirical Risk/Empirical Loss**

However, we only have access to a finite **training set** $\mathbb{D}_{\text{train}} = \{ (\boldsymbol{x}^{(1)}, y^{(1)}), \dots, (\boldsymbol{x}^{(N)}, y^{(N)}) \}$ and instead minimize

$$L_{\mathbb{D}_{\text{train}}}(f_{\boldsymbol{\theta}}) = \frac{1}{N} \sum_{i=1}^{N} \underbrace{\ell(f(\boldsymbol{x}^{(i)}, \boldsymbol{\theta}), y^{(i)})}_{=: \ell^{(i)}}$$

**Stochastic Gradient Descent**

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \boldsymbol{g}_{\mathbb{B}^{(t)}}$$

**(Heavy Ball) MOMENTUM**

$$\boldsymbol{v}^{(t)} = \rho\boldsymbol{v}^{(t-1)} + \eta\boldsymbol{g}_{\mathbb{B}^{(t)}}$$
$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \boldsymbol{v}^{(t)}$$

**NESTEROV ACCELERATED GRADIENT**

$$\boldsymbol{v}^{(t)} = \rho\boldsymbol{v}^{(t-1)} + \eta\frac{1}{B}\sum_{i=1}^{B}\nabla_{\boldsymbol{\theta}}\ell(f(\boldsymbol{x}^{(i)}, \boldsymbol{\theta} - \rho\boldsymbol{v}^{(t-1)}), y^{(i)})$$

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \boldsymbol{v}^{(t)}$$

**RMSPROP**

$$\boldsymbol{s}^{(t)} = \rho\boldsymbol{s}^{(t-1)} + (1 - \rho)\boldsymbol{g}_{\mathbb{B}^{(t)}} \odot \boldsymbol{g}_{\mathbb{B}^{(t)}}$$
$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \frac{\eta}{\sqrt{\boldsymbol{s}^{(t)}} + \varepsilon} \odot \boldsymbol{g}_{\mathbb{B}^{(t)}}$$

**ADAM**

$$\boldsymbol{m}^{(t)} = \beta_1 \boldsymbol{m}^{(t-1)} + (1 - \beta_1)\boldsymbol{g}_{\mathbb{B}^{(t)}}$$

$$\boldsymbol{v}^{(t)} = \beta_2 \boldsymbol{v}^{(t-1)} + (1 - \beta_2)\boldsymbol{g}_{\mathbb{B}^{(t)}} \odot \boldsymbol{g}_{\mathbb{B}^{(t)}}$$

$$\hat{\boldsymbol{m}}^{(t)} = \frac{\boldsymbol{m}^{(t)}}{1 - \beta_1^t}$$

$$\hat{\boldsymbol{v}}^{(t)} = \frac{\boldsymbol{v}^{(t)}}{1 - \beta_2^t}$$

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \frac{\eta}{\sqrt{\hat{\boldsymbol{v}}^{(t)}} + \varepsilon} \odot \hat{\boldsymbol{m}}^{(t)}$$

## 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 |
| Adadelta | 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 | SM3 |
| 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-I | 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 | | |

8

# The State of Neural Network Training

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 |
| Adadelta | 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 | SM3 |
| 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-I | 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 Neural Network Training

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 |
| Adadelta | 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 | SM3 |
| 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-I | 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:**

$$\theta = \theta - \boxed{\eta}\, \boldsymbol{g}$$

**ADAM:**

$$\boldsymbol{m} = \boxed{\beta_1}\,\boldsymbol{m} + (1 - \beta_1)\boldsymbol{g}$$

$$\boldsymbol{v} = \boxed{\beta_2}\,\boldsymbol{v} + (1 - \beta_2)\boldsymbol{g} \odot \boldsymbol{g}$$

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

$$\theta = \theta - \frac{\boxed{\eta}}{\sqrt{\hat{\boldsymbol{v}}} + \boxed{\varepsilon}} \odot \hat{\boldsymbol{m}}$$

## 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 |
| Adadelta | 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 | SM3 |
| 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-I | 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:**

$$\theta = \theta - \boxed{\eta}\, \boldsymbol{g}$$

**ADAM:**

$$\boldsymbol{m} = \boxed{\beta_1}\, \boldsymbol{m} + (1 - \beta_1)\boldsymbol{g}$$

$$\boldsymbol{v} = \boxed{\beta_2}\, \boldsymbol{v} + (1 - \beta_2)\boldsymbol{g} \odot \boldsymbol{g}$$

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

$$\theta = \theta - \frac{\boxed{\eta}}{\sqrt{\hat{\boldsymbol{v}}} + \boxed{\varepsilon}} \odot \hat{\boldsymbol{m}}$$

**Benchmarks**

8

## 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 |
| Adadelta | 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 |
| AdamBS | DiffGrad | Noisy Adam/Noisy K-FAC | SM3 |
| AdamMC | EAdam | NosAdam | SMG |
| AdaMod | EKFAC | Novograd | SNGM |
| AdamP/SGDP | Eve | NT-SGD | SoftAdam |
| AdamT | Expectigrad | Padam | SRSGD |
| AdamW | FastAdaBelief | PAGE | Step-Tuned SGD |
| AdamX | FRSGD | PAL | SWATS |
| ADAS | G-AdaGrad | PolyAdam | SWNTS |
| AdaS | GADAM | Polyak | TAdam |
| AdaScale | Gadam | PowerSGD/PowerSGDM | TEKFAC |
| AdaSGD | GOALS | Probabilistic Polyak | VAdam |
| AdaShift | GOLS-I | ProbLS | VR-SGD |
| AdaSqrt | Grad-Avg | PStorm | vSGD-b/vSGD-g/vSGD-l |
| Adathm | GRAPES | QHAdam/QHM | vSGD-fd |
| AdaX/AdaX-W | Gravilon | RAdam | WNGrad |
| AEGD | Gravity | Ranger | YellowFin |
| ALI-G | HAdam | RangerLars | Yogi |
| | HyperAdam | | |

## Tuning hyperparameters

**SGD:**

$$\theta = \theta - \boxed{\eta}\, \boldsymbol{g}$$

**ADAM:**

$$\boldsymbol{m} = \boxed{\beta_1}\,\boldsymbol{m} + (1 - \beta_1)\boldsymbol{g}$$

$$\boldsymbol{v} = \boxed{\beta_2}\,\boldsymbol{v} + (1 - \beta_2)\boldsymbol{g} \odot \boldsymbol{g}$$

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

$$\theta = \theta - \frac{\boxed{\eta}}{\sqrt{\hat{\boldsymbol{v}}} + \boxed{\varepsilon}} \odot \hat{\boldsymbol{m}}$$

**Benchmarks**

**Debugging Tools**

UNIVERSITÄT TÜBINGEN
EBERHARD KARLS

▶ **Training $\neq$ optimization**
A fundamental and crucial difference.

▶ **Stochasticity is a core property**
A primary source of the challenge.

▶ **Training requires more than an optimizer**
Update rule + hyperparameters + tuning methods + schedules + ....
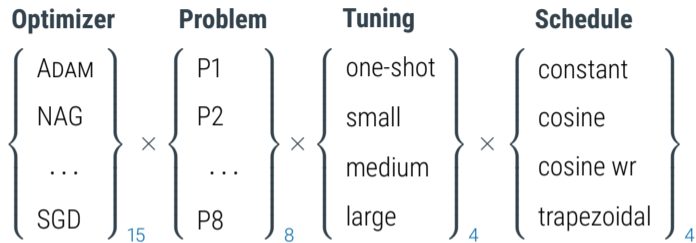
Part II
**Benchmarking Training Algorithms**

The quest for finding the state of the art in neural network training

# The Benchmarking Process
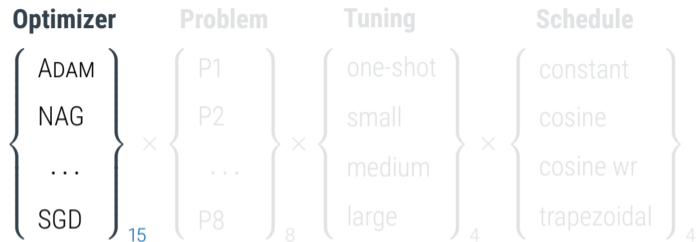
Many dimensions to explore at once

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schmidt et al. 2021)

**Optimizer**

$$\left\{ \begin{array}{c} \text{ADAM} \\ \text{NAG} \\ \ldots \\ \text{SGD} \end{array} \right\}_{15}$$

$\times$

**Problem**

$$\left\{ \begin{array}{c} \text{P1} \\ \text{P2} \\ \ldots \\ \text{P8} \end{array} \right\}_{8}$$

$\times$

**Tuning**

$$\left\{ \begin{array}{c} \text{one-shot} \\ \text{small} \\ \text{medium} \\ \text{large} \end{array} \right\}_{4}$$

$\times$

**Schedule**

$$\left\{ \begin{array}{c} \text{constant} \\ \text{cosine} \\ \text{cosine wr} \\ \text{trapezoidal} \end{array} \right\}_{4}$$

# The Benchmarking Process

Many dimensions to explore at once

| Optimizer | Problem | Tuning | Schedule |
|---|---|---|---|
| ADAM | P1 | one-shot | constant |
| NAG | P2 | small | cosine |
| ... | ... | medium | cosine wr |
| SGD | P8 | large | trapezoidal |
| 15 | 8 | 4 | 4 |

- ● AMSBOUND
- ● AMSGRAD
- ● ADABELIEF
- ● ADABOUND
- ● ADADELTA

- ● ADAGRAD
- ● ADAM
- ● LA(MOM.)
- ● LA(RADAM)
- ● MOMENTUM

- ● NAG
- ● NADAM
- ● RADAM
- ● RMSPROP
- ● SGD

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schmidt et al. 2021)

| | Optimizer | | Problem | | Tuning | | Schedule | |
|---|---|---|---|---|---|---|---|---|
| | ADAM | | P1 | | one-shot | | constant | |
| | NAG | × | P2 | × | small | × | cosine | × |
| | ... | | ... | | medium | | cosine wr | |
| | SGD | 15 | P8 | 8 | large | 4 | trapezoidal | 4 |

| | Data set | Model | Task | Metric | Batch size | Budget *in epochs* | 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

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schmidt et al. 2021)

| **Optimizer** | | **Problem** | | **Tuning** | | **Schedule** | |
|---|---|---|---|---|---|---|---|
| ADAM | | P1 | | one-shot | | constant | |
| NAG | × | P2 | × | small | × | cosine | × |
| . . . | | . . . | | medium | | cosine wr | |
| SGD | 15 | P8 | 8 | large | 4 | trapezoidal | 4 |

▶ **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

10

# The Benchmarking Process
Many dimensions to explore at once

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN
(Schmidt et al. 2021)

**Optimizer** **Problem** **Tuning** **Schedule**

$$
\left\{ \begin{matrix} \text{ADAM} \\ \text{NAG} \\ \cdots \\ \text{SGD} \end{matrix} \right\}_{15}
\times
\left\{ \begin{matrix} \text{P1} \\ \text{P2} \\ \cdots \\ \text{P8} \end{matrix} \right\}_{8}
\times
\left\{ \begin{matrix} \text{one-shot} \\ \text{small} \\ \text{medium} \\ \text{large} \end{matrix} \right\}_{4}
\times
\left\{ \begin{matrix} \text{constant} \\ \text{cosine} \\ \text{cosine wr} \\ \text{trapezoidal} \end{matrix} \right\}_{4}
$$

# The Benchmarking Process

Many dimensions to explore at once

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schmidt et al. 2021)

**Optimizer**    **Problem**    **Tuning**    **Schedule**

$$\left\{ \begin{array}{c} \text{ADAM} \\ \text{NAG} \\ \dots \\ \text{SGD} \end{array} \right\}_{15} \times \left\{ \begin{array}{c} \text{P1} \\ \text{P2} \\ \dots \\ \text{P8} \end{array} \right\}_{8} \times \left\{ \begin{array}{c} \text{one-shot} \\ \text{small} \\ \text{medium} \\ \text{large} \end{array} \right\}_{4} \times \left\{ \begin{array}{c} \text{constant} \\ \text{cosine} \\ \text{cosine wr} \\ \text{trapezoidal} \end{array} \right\}_{4}$$

**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 Process
Many dimensions to explore at once

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN
(Schmidt et al. 2021)

| Optimizer | | Problem | | Tuning | | Schedule | |
|---|---|---|---|---|---|---|---|
| ADAM | | P1 | | one-shot | | constant | |
| NAG | × | P2 | × | small | × | cosine | |
| ... | | ... | | medium | | cosine wr | |
| SGD | 15 | P8 | 8 | large | 4 | trapezoidal | 4 |

**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*

10

# The Benchmarking Results

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

(Schmidt et al. 2021)

# The Benchmarking Results

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

(Schmidt et al. 2021)



Columns: Quadratic Deep | MNIST VAE | F-MNIST 2c2d | CIFAR-10 3c3d | F-MNIST VAE | CIFAR-100 All-CNN-C | SVHN Wide ResNet 16-4 | Tolstoi Char RNN

Legend:
- AMSBOUND
- ADABOUND
- ADAM
- Mom.
- NADAM
- RMSPROP
- AMSGRAD
- ADADELTA
- LA(MOM.)
- NAG
- RADAM
- SGD
- ADABELIEF
- ADAGRAD
- LA(RADAM)

# The Benchmarking Results

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

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schmidt et al. 2021)

Quadratic Deep · MNIST VAE · F-MNIST 2c2d · CIFAR-10 3c3d · F-MNIST VAE · CIFAR-100 All-CNN-C · SVHN Wide ResNet 16-4 · Tolstoi Char RNN

Legend: AMSBound, AdaBound, Adam, Mom., Nadam, RMSProp, AMSGrad, Adadelta, LA(Mom.), NAG, RAdam, SGD, AdaBelief, Adagrad, LA(RAdam)

ML
Commons

**Algorithms Working Group**

**George Dahl**      Google

**Frank Schneider**      University of Tübingen

**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 workloads

► A huge large-scale effort by 25+ researchers from Google, University of Tübingen, University of Toronto, Meta AI, etc.

UNIVERSITÄT
TÜBINGEN
EBERHARD KARLS

▶ **Unclear SOTA**
There is no established protocol to train neural networks.

▶ **Literature only provides families of algorithms**
The many confounding factors make it hard to benchmark training methods.

▶ **Babysitting and tuning**
The existing methods are clearly unsatisfying.

Part III
**Novel Debugging Tools**

Leveraging all available information for better deep learning debuggers

# Loss Curves Do Not Tell the Full Story

Why we need better observables in neural network training

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schneider et al. 2021)

Loss Curve

14

Loss Curve

Loss Landscape

# Loss Curves Do Not Tell the Full Story

Why we need better observables in neural network training

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schneider et al. 2021)

# COCKPIT in Action!

Training the ALL-CNN-C network through the lens of COCKPIT
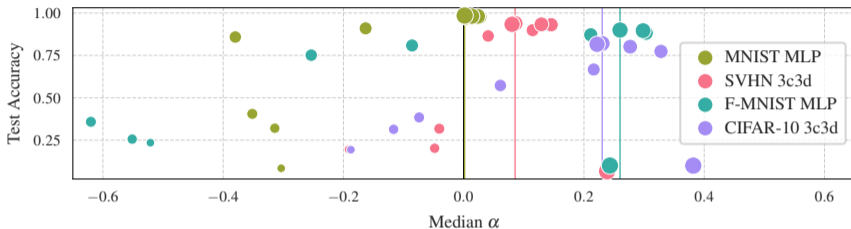
EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schneider et al. 2021)

# COCKPIT as a Tool for Finding Bugs and Inspiring Research

Neural network training requires systematically overstepping the local minima

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

(Schneider et al. 2021)

▶ **Additional available information**
Inexpensive but often hidden by existing software frameworks.

▶ **We can leverage it for better debugging tools**
Or as a path to better autonomous training methods.

▶ **Practical tips & lessons learned**
- ▶ Start with something easy that was shown to work (or ADAM).
- ▶ Tune your hyperparameters (learning rate + schedule, weight decay, $1 - \beta_1$, etc.).
- ▶ Great resource: "Deep Learning Tuning Playbook".

# Summary

► **The current training methods are unsatisfying**: There is no established protocol to train neural networks and all contenders require babysitting and tuning.

► **Benchmarking deep learning algorithms is crucial** but it is challenging and we need a better methodology to find the best training algorithms.

► **Using distributions and confidences** we can built better debugging tools and training methods for neural networks.

with the help of many more

| | | |
|---|---|---|
| **Benchmark** | github.com/SirRob1997/Crowded-Valley---Results | |
| **MLCommons** | github.com/mlcommons/algorithmic-efficiency | ML Commons |
| **Cockpit** | pip install cockpit-for-pytorch | Cockpit |