

# Reproducibility and frugality in AI benchmarking: lessons from Benchopt

Thomas Moreau



MIND



## A word about me

- ▶ Researcher at Inria – MIND
- ▶ <https://tommoral.github.io>
- ▶ [thomas.moreau@inria.fr](mailto:thomas.moreau@inria.fr)
- ▶ **Research topics:** Time-series, Physiological signals, Inverse Problems, Bilevel Optimization, Unrolling, Pattern Learning, Point Processes, .
- ▶ OSS maintainer/contributor



## The Era of Benchmarks: AI as an empirical science

# The ImageNet competition

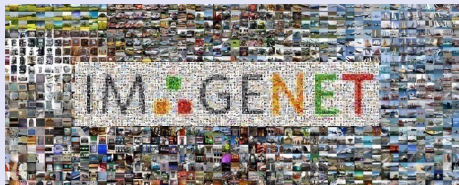
- ▶ Annual competition since 2010
- ▶ Evaluate image classification methods with 14M labeled images among 1k categories





# The ImageNet competition

- ▶ Annual competition since 2010
- ▶ Evaluate image classification methods with 14M labeled images among 1k categories
- ▶ Boosted AI and Deep Learning research when Alex Krizhevsky won in 2012.



# The ImageNet competition

- ▶ Annual competition since 2010
- ▶ Evaluate image classification methods with 14M labeled images among 1k categories
- ▶ Boosted AI and Deep Learning research when Alex Krizhevsky won in 2012.



⇒ Demonstrates the importance of benchmarks to drive research in AI.

# Many benchmarks in AI

## Many benchmarks followed ImageNet:

- ▶ Natural Language Processing: GLUE, SuperGLUE
- ▶ Reinforcement Learning: Atari, MuJoCo, OpenAI Gym
- ▶ Others: fastMRI, DAWNBench, MLPerf, etc.



# Many benchmarks in AI

Many benchmarks followed  
ImageNet:

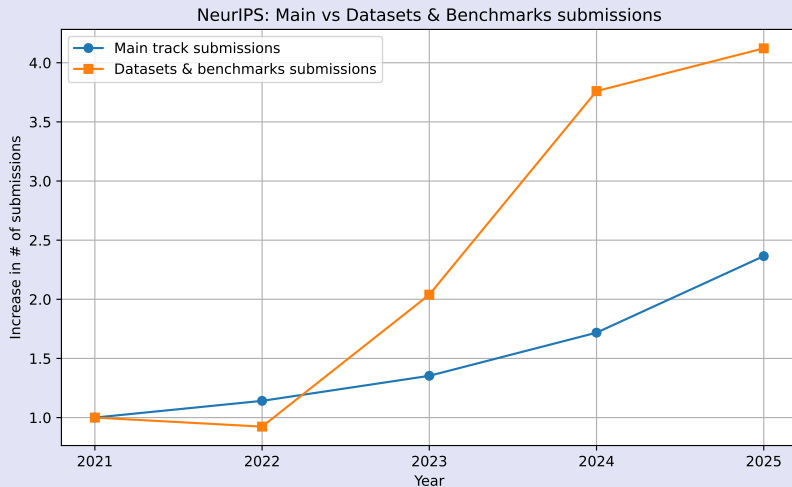
- ▶ Natural Language Processing: GLUE, SuperGLUE
- ▶ Reinforcement Learning: Atari, MuJoCo, OpenAI Gym
- ▶ Others: fastMRI, DAWNBench, MLPerf, etc.



⇒ Benchmarks are now ubiquitous in AI research.

# Too many benchmarks in AI?

In the recent years, many benchmarks have been proposed:



# Benchmark goals in AI

	Short-term progress	Long-term evaluation
<b>Task-specific</b>	Challenge/Competition → push limits quickly	SOTA tracking → measure progress
<b>Generalizable</b>	Research question → empirical study	Benchmark framework → stable & extensible

# Benchmark goals in AI

	Short-term progress	Long-term evaluation
<b>Task-specific</b>	Challenge/Competition → push limits quickly	SOTA tracking → measure progress
<b>Generalizable</b>	Research question → empirical study	Benchmark framework → stable & extensible

## Takeaway

Most attention goes to the top-left quadrant for fast progress, but solid science requires the bottom-right.

# The three components of a benchmark

A benchmark is defined by three components:

- ▶ **Objective:** what is being measured?
- ▶ **Dataset:** on what evidence?
- ▶ **Solvers/Methods:** What are we comparing?



# The three components of a benchmark

A benchmark is defined by three components:

- ▶ **Objective:** what is being measured?
- ▶ **Dataset:** on what evidence?
- ▶ **Solvers/Methods:** What are we comparing?

Different benchmark goals emphasize different components.

- ▶ **Challenge benchmarks:** Fixed dataset + single metric → stresses solver comparison.  
Focus on performance

# The three components of a benchmark

A benchmark is defined by three components:

- ▶ **Objective:** what is being measured?
- ▶ **Dataset:** on what evidence?
- ▶ **Solvers/Methods:** What are we comparing?

Different benchmark goals emphasize different components.

- ▶ **Challenge benchmarks:** Fixed dataset + single metric → stresses solver comparison. Focus on performance
- ▶ **SOTA tracking benchmarks:** Multiple metrics and datasets, compare solvers to measure progress over time.  
Risk: test-set overfitting / cherry-picked metrics.

# The three components of a benchmark

A benchmark is defined by three components:

- ▶ **Objective:** what is being measured?
- ▶ **Dataset:** on what evidence?
- ▶ **Solvers/Methods:** What are we comparing?

Different benchmark goals emphasize different components.

- ▶ **Challenge benchmarks:** Fixed dataset + single metric → stresses solver comparison. Focus on performance
- ▶ **SOTA tracking benchmarks:** Multiple metrics and datasets, compare solvers to measure progress over time. Risk: test-set overfitting / cherry-picked metrics.
- ▶ **Research benchmarks:** Fixed set of methods evaluated, with broad range of metrics. Risk: incomplete / quickly outdated

# Challenges in AI benchmarking

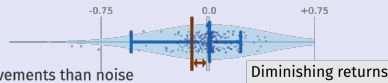
[Varoquaux and Cheplygina 2022]

## ► Futile benchmarks

Lung cancer classification

Test size: max 1K

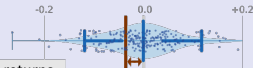
Smaller improvements than noise



Schizophrenia classification

Test size: 120

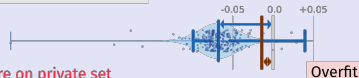
Diminishing returns



Lung tumor segmentation

Test size: max 6k

Poorer score on private set



Nerve segmentation

Test size 5.5K



# Challenges in AI benchmarking

- ▶ Futile benchmarks
- ▶ Lack of proper baselines hinders scientific progress.

## Do we really need Foundation Models for multi-step-ahead Epidemic Forecasting?

Position: Quo Vadis, Unsupervised Time Series Anomaly Detection?

M. Saquib Sarfraz<sup>1,2</sup>, Mei-Yun Chen<sup>1</sup>, Lukas Layer<sup>1</sup>, Kunyu Peng<sup>2</sup>, Marius Koutakis<sup>2</sup>

PNAS

RESEARCH ARTICLE

COMPUTER SCIENCES

OPEN ACCESS

## Implicit data crimes: Machine learning bias arising from misuse of public data

Efrat Shimron<sup>1</sup>, Jonathan I. Tamir<sup>2</sup>, Ke Wang<sup>3</sup>, and Michael Lustig<sup>4</sup>

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

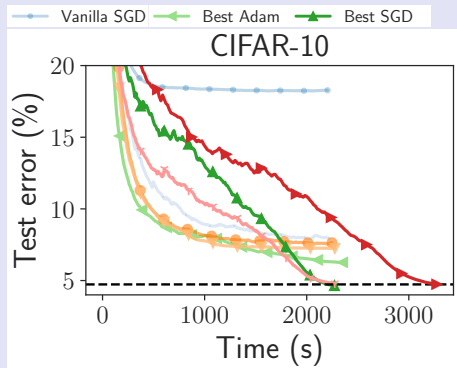
Robin M. Schmidt<sup>1</sup>, Frank Schneider<sup>1</sup>, Philipp Hennig<sup>1,2</sup>

# Unclear improvement!

# Challenges in AI benchmarking

[Moreau et al. 2022]

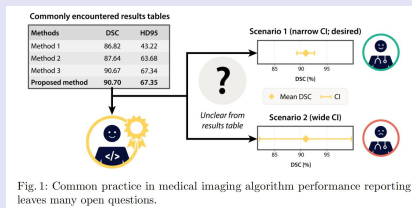
- ▶ Futile benchmarks
- ▶ Lack of proper baselines hinders scientific progress.
- ▶ Reproducing benchmarks is hard.



# Challenges in AI benchmarking

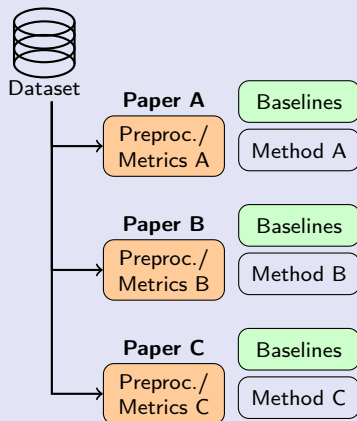
- ▶ Futile benchmarks
- ▶ Lack of proper baselines hinders scientific progress.
- ▶ Reproducing benchmarks is hard.
- ▶ Statistical validity is often missing.

[Christodoulou et al. 2024]



# Challenges in AI benchmarking

- ▶ Futile benchmarks
- ▶ Lack of proper baselines hinders scientific progress.
- ▶ Reproducing benchmarks is hard.
- ▶ Statistical validity is often missing.
- ▶ Benchmarking cost is duplicated across groups.



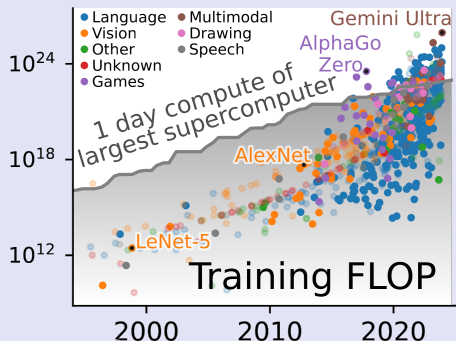
Each paper independently rebuilds preprocessing, baselines, and evaluation → duplicated cost.



# Challenges in AI benchmarking

[Varoquaux et al. 2025]

- ▶ Futile benchmarks
- ▶ Lack of proper baselines hinders scientific progress.
- ▶ Reproducing benchmarks is hard.
- ▶ Statistical validity is often missing.
- ▶ Benchmarking cost is duplicated across groups.



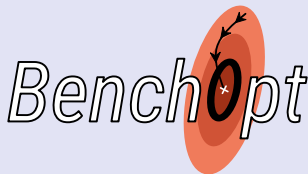
## Reproducible method comparison with Benchopt



## References

- ▶ **TM**, Massias, M., Gramfort, A., Ablin, P., Bannier, P.-A., Charlier, B., Dagr  ou, M., la Tour, T. D., Durif, G., Dantas, C. F., Klopfenstein, Q., Larsson, J., Lai, E., Lefort, T., Mal  zieux, B., Moufad, B., Nguyen, B. T., Rakotomamonjy, A., Ramzi, Z., Salmon, J., and Vaite  r, S. (2022). [Benchopt: Reproducible, efficient and collaborative optimization benchmarks](#). In *NeurIPS*

# Making runnable benchmarks with benchopt



benchopt provides a framework to organize and run benchmarks

Examples of existing benchmarks:

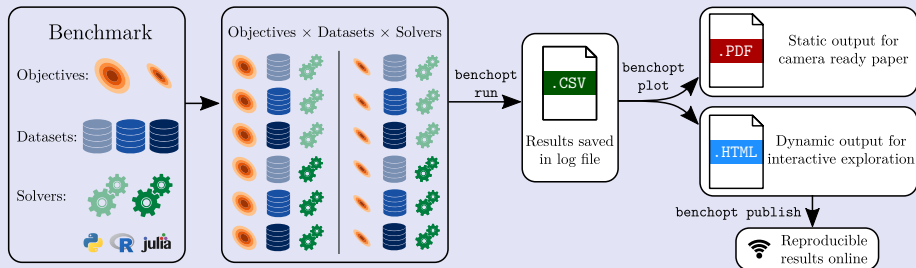
- ▶ **Image Classification (resnet)**
- ▶ **Logistic regression**
- ▶ **Lasso**
- ▶ **ICA**
- ▶ **Unsup. Domain Adaptation**
- ▶ **Bilevel Optimization**
- ▶ **Brain Computer Interface**
- ▶ **...**

## Contributors from...



# Components of a benchmark

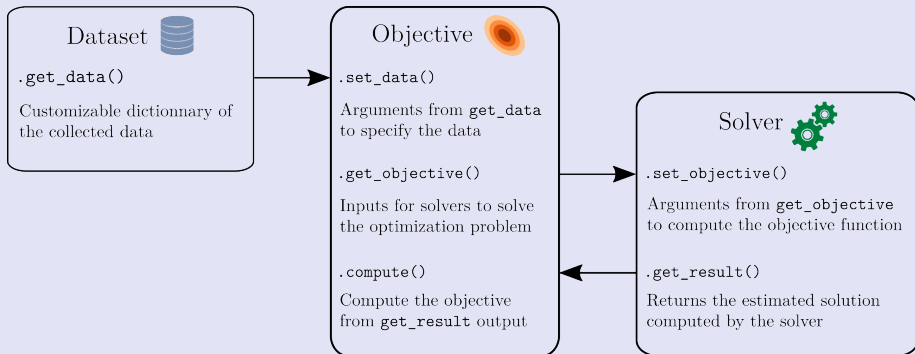
## 3 components: Objective, Dataset, Solver



# A modular framework to create benchmarks

## 3 components: Objective, Dataset, Solver

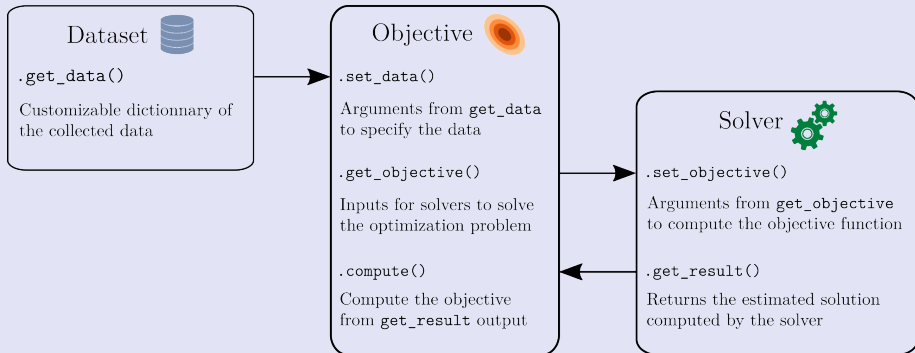
Dependency relation between Dataset - Objective - Solver



# A modular framework to create benchmarks

## 3 components: Objective, Dataset, Solver

Dependency relation between Dataset - Objective - Solver



⇒ Benchopt defines the interface between components.

# Explicit requirements and parameters

```
from benchopt import BaseSolver
from benchmark_utils import grad, init_func

class Solver(BaseSolver):
    name = "GD"
    requirements = ["numpy"]
    sampling_strategy = "callback"
    parameters = {"lr": [1, 1.9], "init": [0, 0.1]}

    def set_objective(self, X):
        self.X = X

    def run(self, cb):
        self.w = init_func(self.X, self.init)
        while cb():
            self.w -= self.lr * grad(self.X, self.w)

    def get_result(self):
        return dict(w=self.w)
```



# Eval multiple metrics at once

```
from benchopt import BaseObjective
from benchmark_utils import split, error

class Objective(BaseObjective):
    name = "Least Squares"
    url = "https://github.com/#ORG/#BENCHMARK_NAME"

    def set_data(self, X):
        self.X_train, self.X_test = split(X)

    def evaluate_result(self, w):
        return dict(
            train_error=error(w, self.X_train),
            test_error=error(w, self.X_test),
        )

    def get_objective(self):
        return dict(X=self.X_train)
```

# Explicit preprocessing

```
from benchopt import BaseDataset
from sklearn.datasets import load_digits

class Dataset(BaseDataset):
    name = "Digits"
    requirements = ["scikit-learn"]

    def get_data(self):
        X = load_digits(return_X_y=True)[0]
        X /= X.std()
        return dict(X=X)
```

# Explicit preprocessing

```
from benchopt import BaseDataset
from sklearn.datasets import load_digits

class Dataset(BaseDataset):
    name = "Digits"
    requirements = ["scikit-learn"]

    def get_data(self):
        X = load_digits(return_X_y=True)[0]
        X /= X.std()
        return dict(X=X)
```

⇒ Reproducible benchmark by design!

**Goal:** if you use the same setup, you don't need to re-run the baseline methods!

# Structure of a benchmark

```
benchmark/  
├── objective.py  
├── datasets/  
│   ├── dataset1.py  
│   └── dataset2.py  
└── solvers/  
    ├── solver1.py  
    └── solver2.py
```

## Modular & extendable

- ▶ New metric? modify objective

# Structure of a benchmark

```
benchmark/  
├── objective.py  
├── datasets/  
│   ├── dataset1.py  
│   └── dataset2.py  
└── solvers/  
    ├── solver1.py  
    └── solver2.py
```

## Modular & extendable

- ▶ New metric? modify objective
- ▶ New dataset? add a file

# Structure of a benchmark

```
benchmark/  
├── objective.py  
├── datasets/  
│   ├── dataset1.py  
│   └── dataset2.py  
└── solvers/  
    ├── solver1.py  
    └── solver2.py
```

## Modular & extendable

- ▶ New metric? modify objective
- ▶ New dataset? add a file
- ▶ New solver? add a file

# Structure of a benchmark

```
benchmark/  
├── objective.py  
├── datasets/  
│   ├── dataset1.py  
│   └── dataset2.py  
└── solvers/  
    ├── solver1.py  
    └── solver2.py
```

## Modular & extendable

- ▶ New metric? modify objective
- ▶ New dataset? add a file
- ▶ New solver? add a file

# Structure of a benchmark

```
benchmark/  
├── objective.py  
├── datasets/  
│   ├── dataset1.py  
│   └── dataset2.py  
└── solvers/  
    ├── solver1.py  
    └── solver2.py
```

## Modular & extendable

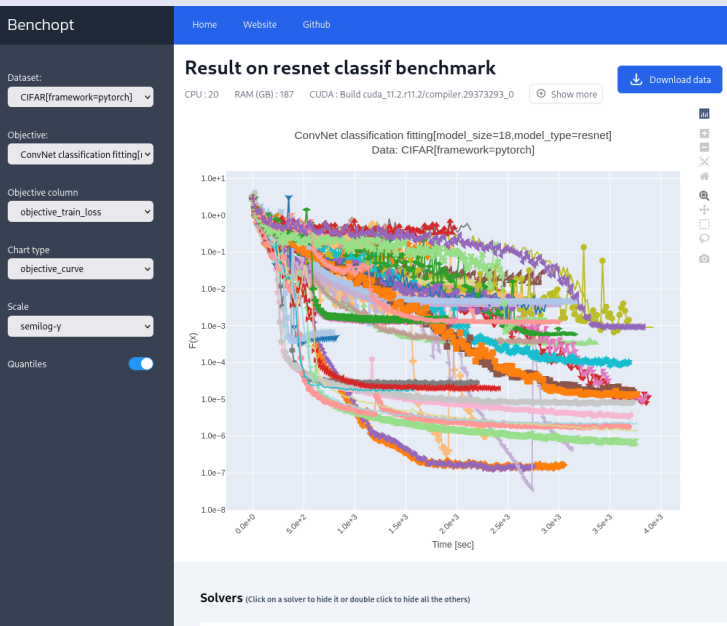
- ▶ New metric? modify objective
- ▶ New dataset? add a file
- ▶ New solver? add a file

Template to create a new benchmark:

[https://github.com/benchopt/template\\_benchmark](https://github.com/benchopt/template_benchmark)

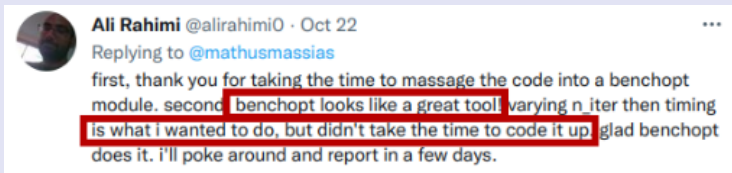


# Interactive results exploration



## Benchopt makes your life easy

- ▶ Build on previous benchmarks
- ▶ Use solvers in Python, R, Julia, binaries...
- ▶ Monitor any metric you want altogether (test/train loss, ...)
- ▶ Add parameters to solvers and use config files
- ▶ Share and publish HTML results
- ▶ Run all benchmarks in parallel, on HPC clusters...
- ▶ Cache runs and results
- ▶ and much more!



# Typical case: deep learning optimization

## A. List of optimizers and schedules considered

Table 2: List of optimizers considered for our benchmark. This is only a subset of all existing methods for deep learning.

Name	Ref.	Name	Ref.
AccuGrad	(Levy et al., 2018)	HypocoDam	(Wang et al., 2019b)
ACCGP	(Zhang et al., 2020)	K-RFGSG-BFGSL	(Goldfarb et al., 2020)
AdaAkr	(Ola et al., 2019)	KF-QN-CNN	(Ola & Goldfarb, 2021)
AdaBach	(Devoruchinda et al., 2017)	KMAC	(Sharma & Grosse, 2015)
AdaBayes/AdaBayes-S5	(Aichukon, 2020)	KFL/KFRA	(Bures et al., 2017)
AdaBelief	(Zhang et al., 2020)	LAdaDefLAdoement	(Runkel & Martins, 2019)
AdaBelief	(Vas et al., 2019)	LAMS	(Owens, 2020)
AdaBound	(Luo et al., 2019)	LaProp	(Ziyin et al., 2020)
AdaCmp	(Chen et al., 2019)	LARS	(Hoyer et al., 2017)
AdaDelta	(Zeiler, 2012)	LaPTT	(Amend et al., 2021)
AdaFactor	(Shawar & Stern, 2018)	LookAhead	(Zhang et al., 2019)
AdaFro	(Hao et al., 2019)	M-SVRG	(Bailin & Haining, 2019)
AdaFro	(Chen et al., 2019a)	MADGRAD	(Delella & Jafari, 2021)
AdaFTRL	(Chaboua & Pui, 2015)	MAS	(Lachry et al., 2020)
Adagrad	(Duchi et al., 2011)	MIRA	(Chen et al., 2020a)
ADAHESIAN	(Vas et al., 2020)	MTAdam	(Makris & Wolf, 2020)
Adam	(Ola et al., 2020)	MYRC-IMPVRC-2	(Chen & Zhou, 2020)
AdamAm	(Stevens et al., 2019)	Nadam	(Demp, 2016)
Adam	(Kingma & Ba, 2015)	NAMSHNMSG	(Chen et al., 2019b)
Adam <sup>+</sup>	(Li et al., 2020a)	NO-Adam	(Zhang et al., 2017a)
AdamAL	(Tao et al., 2019)	Noto	(Luo et al., 2021b)
AdaMax	(Kingma & Ba, 2015)	Noxton	(Nemey, 1983)
AdamNES	(Li et al., 2020b)	Noisy Adam/Noisy K-RAC	(Zhang et al., 2019)
AdamNC	(Rohdi et al., 2018)	NoxAdam	(Huang et al., 2019)
AdaMod	(Dang et al., 2019)	Novograd	(Zhang et al., 2019)
AdamPFGP	(Ola et al., 2021)	NTSGD	(Zhang et al., 2021b)
AdamT	(Zhou et al., 2020)	Padam	(Chen et al., 2020a)
AdamW	(Loshchilov & Hutter, 2019)	PaLR	(Li et al., 2020b)
AdamX	(Tao & Wang, 2019)	PdL	(Muehler & Zell, 2020)
ADAMS	(Hoyden, 2020)	PopAdam	(Dervoz et al., 2019)
AdaS	(Hosami & Papanicolaou, 2020)	PyAda	(Piyab, 1966)
AdaScale	(Johnson et al., 2020)	PowerSGD/PowerGEM	(Vogelstein, 2019)
AdaSGD	(Wang & Wern, 2020)	Probabilistic PyAda	(de Koo et al., 2021)
AdaShin	(Zhou et al., 2019)	PyALS	(Mishchenko & Haining, 2017)
AdaStep	(Hu et al., 2019)	Pyomo	(Ola, 2020)
AdaTanh	(Gao et al., 2019)	QD-Loss/GRM	(Ola & Tarras, 2019)
AdaUN/AdaW	(Li et al., 2020a)	RAdam	(Li et al., 2020a)
AFGD	(Luo & Tian, 2020)	Ranger	(Wright, 2020a)
ALG	(Berrada et al., 2020)	RangerAdam	(Wright, 2020b)
AMEGGrad	(Luo et al., 2019)	RMSProp	(Tejerman & Hinton, 2012)
AMEGGrad	(Rohdi et al., 2018)	RMSProp	(Chen et al., 2019)
AsynGrad	(Roy et al., 2021)	S-GD	(Chang et al., 2020)
AsynSGD	(Vasquez et al., 2019)	SAdam	(Wang et al., 2020b)
ASGD	(Ola et al., 2019b)	SAdamSAMSGrad	(Wang et al., 2019)
ASAM	(Khan et al., 2021)	SALR	(Chen et al., 2020)
AsynES	(Ola et al., 2021)	SAM	(Hest et al., 2021)
AsynES	(Gower et al., 2019)	SC-Adaptive/SC-RMSProp	(Mahmoud & Bera, 2017)
BAdam	(Jain et al., 2018)	SDPProp	(Mori et al., 2017)
BAdam	(Hao & Zhang, 2019)	SGD	(Robbins & Monro, 1951)
BRProp	(Zhang et al., 2017b)	SGD-BB	(Gao et al., 2016)
BRMSProp	(Aichukon, 2020)	SGD-G2	(Ayadi & Tarras, 2020)
BSGD	(Hu et al., 2020)	SGD-IRM	(Ramanathan-Kubera et al., 2021)
C-ADAM	(Stevens et al., 2020)	SGD-RMS	(Tao & Carlsberg, 2021)
CADA	(Chen et al., 2021)	SGD-IRM	(Luo & Luo, 2020)
Cine Momentum	(Berrada & Berrada, 2020)	SGD-IRM	(Loshchilov & Hutter, 2017)
CPop	(Pacheco & Kipriakou, 2019)	SHAdagrad	(Huang et al., 2020)
Cureball	(Korogonov et al., 2019)	Shampoo	(Aoi et al., 2020; Gupta et al., 2019)
Dadan	(Nasiri et al., 2019)	SuperMomentum	(Wang et al., 2019a)
DeepMemory	(Wright, 2020a)	SuperSGD	(Bennett et al., 2019)
DDNGP	(Li et al., 2021a)	SKONGS-AD	(Yang et al., 2020)
DPRFad	(Bures et al., 2020)	SNL	(Ola et al., 2019)
EAdam	(Vas & Gao, 2020)	SMG	(Tao et al., 2020)



Frank Schneider  
@frankstefansch1



#ICML 2021 Paper



Overwhelmed by the flood of optimizers for deep learning? We felt the same and performed an extensive benchmark. Joint work with @robinschmidt\_ & @PhilippHennig5.

Paper: [arxiv.org/abs/2007.01547](https://arxiv.org/abs/2007.01547)

Results: [github.com/SirRob1997/Cro...](https://github.com/SirRob1997/Cro...)

Video: [youtu.be/cz9RZlstFdE](https://youtu.be/cz9RZlstFdE)



Frank Schneider  
@frankstefansch1

Our results? There is no winner consistently outperforming the competition. Instead, Adam remains a strong contender for many problems.

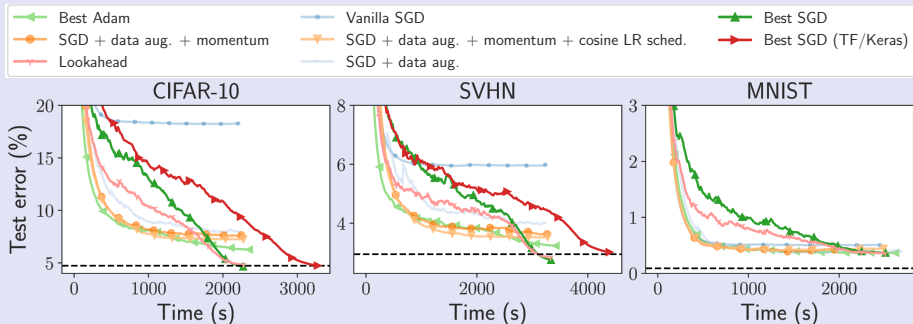
In some cases, just trying out a few optimizers with their default hyperparameters can work as well as tuning one specific method.

⇒ Many novel methods but unclear improvements.

⇒ But this benchmark cannot be easily reproduced!

# Example: Optimization for ResNet on image classification

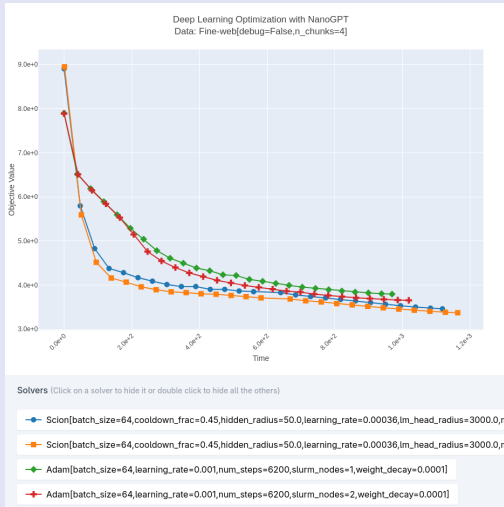
- ▶ Image classification with resnet18
- ▶ Various optimization strategies
- ▶ Compare pytorch and tensorflow
- ▶ Publish reproducible SOTA for baselines



[https://github.com/benchopt/benchmark\\_resnet\\_classif/](https://github.com/benchopt/benchmark_resnet_classif/)

# Example: Large scale-optimization for Deep learning (WIP)

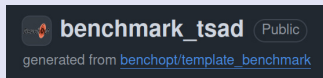
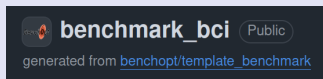
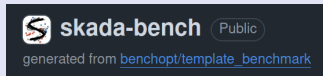
- ▶ Use modern large-scale datasets and models
- ▶ Classical optimization tricks
- ▶ Distributed training and mixed precision



[https://github.com/tommorale/benchmark\\_nanogpt/](https://github.com/tommorale/benchmark_nanogpt/)

# The benchopt roadmap

## Develop reference benchmarks



## Improve the benchmarking methodology



How to ensure that a method is better than another?

Investigating the ranking confidence depending on the number of samples and the CV strategy.

## Reproducibility in AI

# Reproducible research

Different goals:

- ▶ Reproduce the exact same results?

Repeatability | Bitwise reproducibility



# Reproducible research

Different goals:

- ▶ Reproduce the exact same results?

Repeatability | Bitwise reproducibility

- ▶ Provide tools for other to use on their use case?

Reusability

# Reproducible research

Different goals:

- ▶ Reproduce the exact same results?

Repeatability | Bitwise reproducibility

- ▶ Provide tools for other to use on their use case?

Reusability

- ▶ Run with new parameters to evaluate robust results?

# Reproducible research

Different goals:

- ▶ Reproduce the exact same results?

Repeatability | Bitwise reproducibility

- ▶ Provide tools for other to use on their use case?

Reusability

- ▶ Run with new parameters to evaluate robust results?

- ▶ Run with a new dataset to extend the validity?

# Reproducible research

Different goals:

- ▶ Reproduce the exact same results?

Repeatability | Bitwise reproducibility

- ▶ Provide tools for other to use on their use case?

Reusability

- ▶ Run with new parameters to evaluate robust results?

- ▶ Run with a new dataset to extend the validity?

- ▶ Extend the results with a new method, to show improvement?

Replicability | External reproducibility

# Reproducible research

Different goals:

- ▶ Reproduce the exact same results?

Repeatability | Bitwise reproducibility

- ▶ Provide tools for other to use on their use case?

Reusability

- ▶ Run with new parameters to evaluate robust results?

- ▶ Run with a new dataset to extend the validity?

- ▶ Extend the results with a new method, to show improvement?

Replicability | External reproducibility

# Reproducible research

Different goals:

- ▶ Reproduce the exact same results?

Repeatability | Bitwise reproducibility

- ▶ Provide tools for other to use on their use case?

Reusability

- ▶ Run with new parameters to evaluate robust results?

- ▶ Run with a new dataset to extend the validity?

- ▶ Extend the results with a new method, to show improvement?

Replicability | External reproducibility

The last point becomes complicated with AI

## Technical challenges for repeatability:

- ▶ Multiple frameworks: PyTorch, Tensorflow, jax,...
- ▶ Hardware dependence: CPU, GPU, TPUs,...
- ▶ Large scale datasets: hard to share, pre-process,...

## Technical challenges for repeatability:

- ▶ Multiple frameworks: PyTorch, Tensorflow, jax,...
- ▶ Hardware dependence: CPU, GPU, TPUs,...
- ▶ Large scale datasets: hard to share, pre-process,...

## Challenges for reuse:

- ▶ Usability: API, modularity, documentation, tests,...
- ▶ Maintenance: dependencies, code rot, long term support, ...



## Technical challenges for repeatability:

- ▶ Multiple frameworks: PyTorch, Tensorflow, jax,...
- ▶ Hardware dependence: CPU, GPU, TPUs,...
- ▶ Large scale datasets: hard to share, pre-process,...

## Challenges for reuse:

- ▶ Usability: API, modularity, documentation, tests,...
- ▶ Maintenance: dependencies, code rot, long term support, ...

## Statistical challenges for replicability:

- ▶ Stochastic algorithms: random initialization, data shuffling,...
- ▶ Data handling: splitting, data preprocessing, model selection,...

## Good practices to share code

Minimal requirement for a research project: (in Python)

- ▶ Write clean code:
  - ▶ Consistent style: *Use black or flake8.*
  - ▶ Use readable names.
  - ▶ No notebooks!
- ▶ Document your code: *docstring and README.md. optionally sphinx.*
- ▶ Determinist output: *Control the random seeds.*

Optional but advised

- ▶ Document dependencies.
- ▶ Proper package organization
- ▶ Add some tests: *Use pytest.*

## Going further: creating a package

If you really want to make your research *reproducible* by other in different contexts, you need to properly package it.

- ▶ Documentation: *Sphinx*.
- ▶ Test on multiple platforms: *Continuous Integration*.
- ▶ Release on pypi/conda-forge
- ▶ Talk about it ! :)

Example of package: [https://github.com/tomMoral/test\\_package](https://github.com/tomMoral/test_package)

## An overlooked aspect: longer term maintenance

**Many research results are not maintained after publication:**

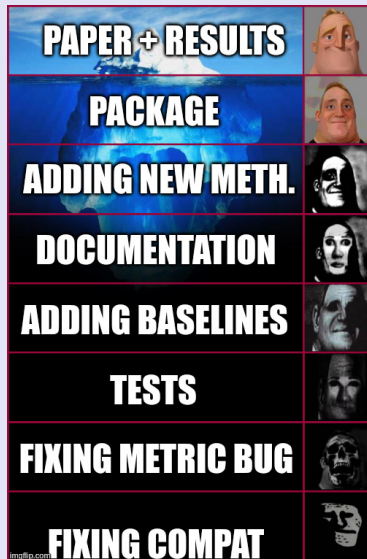
- ▶ Every PhD creates a package;
- ▶ Every post-doc abandons one;
- ▶ The ecosystem grows horizontally, not vertically

## An overlooked aspect: longer term maintenance

Many research results are not maintained after publication:

- ▶ Every PhD creates a package;
- ▶ Every post-doc abandons one;
- ▶ The ecosystem grows horizontally, not vertically

⇒ Limited incentives in AI to maintain a codebase



## Crossing the validation crisis in AI benchmarking: the challenge of statistical validity

## Evaluating decision functions & learning algorithms

In AI, we produce decision function  $g$ , that can be evaluated with:

$$S_{\mathcal{D}}(g) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \sim \mathcal{D}} m(g(X), y)$$

The oracle score is:

$$S^*(g) = \mathbb{E}_{\mathcal{D} \sim d} [S_{\mathcal{D}}(g)]$$

## Evaluating decision functions & learning algorithms

In AI, we produce decision function  $g$ , that can be evaluated with:

$$S_{\mathcal{D}}(g) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \sim \mathcal{D}} m(g(X), y)$$

The oracle score is:

$$S^*(g) = \mathbb{E}_{\mathcal{D} \sim d} [S_{\mathcal{D}}(g)]$$

For learning algorithms  $F_{\lambda}$  producing  $g = F_{\lambda}(\mathcal{D}_{train}, \xi)$ , the oracle score becomes:

$$\begin{aligned} S^*(F_{\lambda}) &= \mathbb{E}_{\mathcal{D}_{train} \sim d} \mathbb{E}_{\xi} [S^*(F_{\lambda}(\mathcal{D}_{train}, \xi))] \\ &= \mathbb{E}_{\mathcal{D}_{test} \sim d} \mathbb{E}_{\mathcal{D}_{train} \sim d} \mathbb{E}_{\xi} [S_{\mathcal{D}_{test}}(F_{\lambda}(\mathcal{D}_{train}, \xi))] \end{aligned}$$



## Benchmarking goal: compare and rank methods

Given two learning algorithms, we want to compare them according to their oracle score:

$$\mathcal{S}^*(F_\lambda) \stackrel{?}{>} \mathcal{S}^*(F_{\lambda'})$$

## Benchmarking goal: compare and rank methods

Given two learning algorithms, we want to compare them according to their oracle score:

$$\mathcal{S}^*(F_\lambda) \stackrel{?}{>} \mathcal{S}^*(F_{\lambda'})$$

In practice, we only have access to empirical estimates of the oracle score!

⇒ How to ensure that our benchmark results reflect the true oracle ranking?

DE L'Importance  
DE  
la validation croisée

OU DU RAPPORT QUE LES LOIX DOIVENT AVOIR AVEC LA CONS-  
TITUTION DE CHAQUE GOUVERNEMENT, LES MOEURS,  
LE CLIMAT, LA RELIGION, LE COMMERCE, &c.

*à quoi l'Auteur a ajouté*

Des recherches nouvelles sur les Loix Romaines touchant les  
Successions, sur les Loix Françoises, & sur les Loix Féodales.

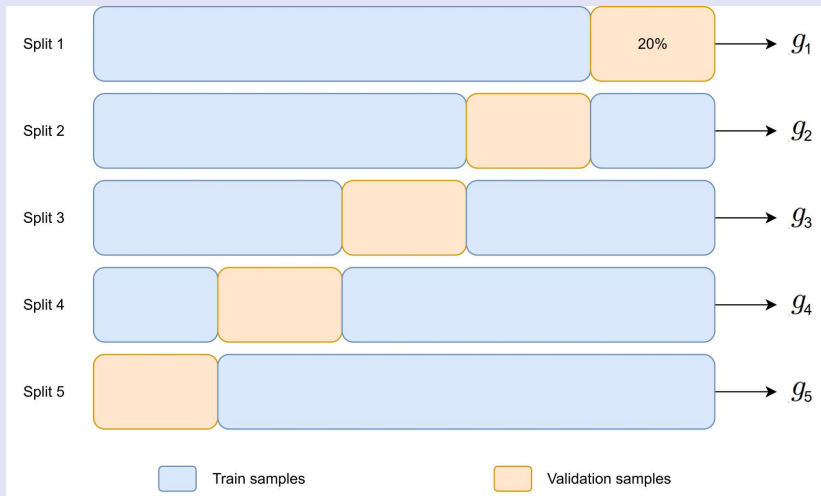
TOME PREMIER.



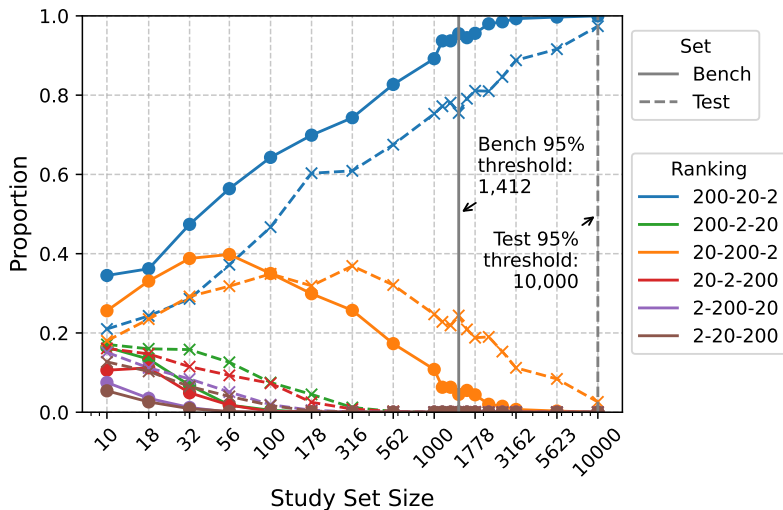
A Grenoble,  
Chez BARRILLOT & FILS.

---

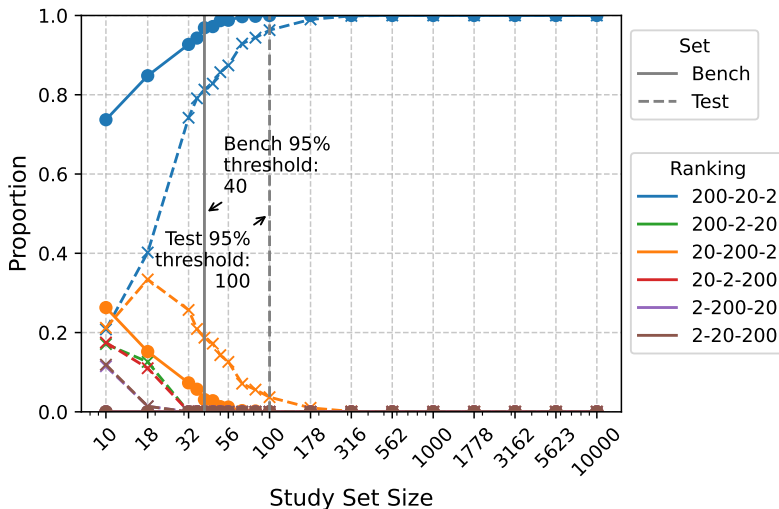
# The cross-validation setting



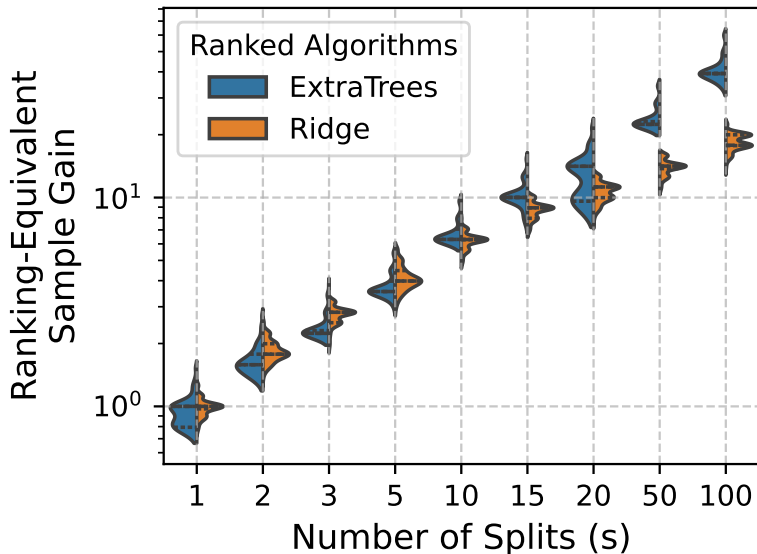
# The benefit of CV over a single train/test split



# The cross-validation setting



## Sample gain of CV over a single split



Stochasticity is intrinsic to AI applications.

⇒ How to ensure that benchmark results are statistically reproducible?



Stochasticity is intrinsic to AI applications.

⇒ How to ensure that benchmark results are statistically reproducible?

One answer is to use proper statistical tools to assess significance of the results.

As this is costly (need many repetitions or large datasets), we need to have a way to easily share and reuse past results!

Reproducible research needs more than just releasing code:

- ▶ Reusable → Clean and Documented.
- ▶ Extendable → Proper packaging and maintenance.
- ▶ Statistically valid → Proper evaluation protocols.

Reproducible research needs more than just releasing code:

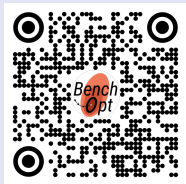
- ▶ Reusable → Clean and Documented.
- ▶ Extendable → Proper packaging and maintenance.
- ▶ Statistically valid → Proper evaluation protocols.

Need proper tools to make it possible!

Also need to build a community around these tools, to share the maintenance and the fun!

# Thank you for your attention!

*Bench*Opt



Don't hesitate to star the benchopt repo!