

AI benchmarking infrastructures: lessons from Benchopt

Thomas Moreau



MIND

inria
informatics mathematics

Benchopt

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
- ▶ Top-1/Top-5 accuracy reported on a fixed test set



The ImageNet competition

- ▶ Annual competition since 2010
- ▶ Evaluate image classification methods with 14M labeled images among 1k categories
- ▶ Top-1/Top-5 accuracy reported on a fixed test set
- ▶ **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
- ▶ Top-1/Top-5 accuracy reported on a fixed test set
- ▶ **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.

Also many dataset repos:

- ▶ OpenML, UCI, PhysioNet,
HuggingFace, etc.



Many benchmarks in AI

Many benchmarks followed ImageNet:

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

Also many dataset repos:

- ▶ OpenML, UCI, PhysioNet, HuggingFace, etc.

⇒ “Benchmarks” are now ubiquitous in AI research.



What makes a benchmark?

A benchmark is defined by three components:

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

⇒ A benchmark is not only a dataset with any type of evaluation protocol.

What makes 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

What makes 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.

What makes a benchmark?

A benchmark is defined by three components:

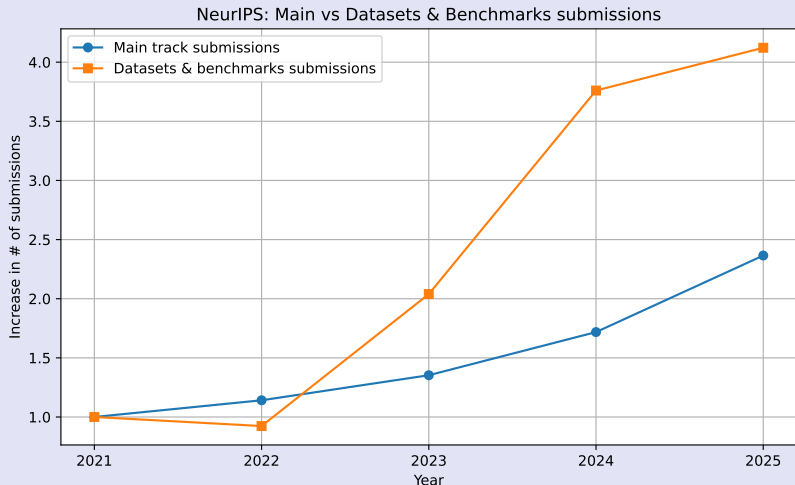
- ▶ **Objective:** what is being measured?
- ▶ **Dataset:** on what evidence?
- ▶ **Solvers/Methods:** What are we comparing?

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

Too many benchmarks in AI?

In the recent years, many benchmarks have been proposed:



⇒ Most of them don't have long-term maintenance plan.

Benchmark goals in AI

kaggle



Task-specific

Short-term

Challenge/Competition
→ push limits quickly

Long-term

SOTA tracking
→ measure progress

Generalizable

Research question
→ empirical study

Field Benchmark
→ extensible



Benchmark goals in AI

kaggle



Task-specific → **Short-term**
Challenge/Competition
→ push limits quickly

Long-term
SOTA tracking
→ measure progress

Generalizable → **Research question**
→ empirical study

Field Benchmark
→ extensible



Takeaway

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

Benchmark goals in AI

kaggle



Task-specific

Short-term

Long-term

Challenge/Competition
→ push limits quickly

SOTA tracking
→ measure progress

Generalizable

Research question
→ empirical study

Field Benchmark
→ extensible



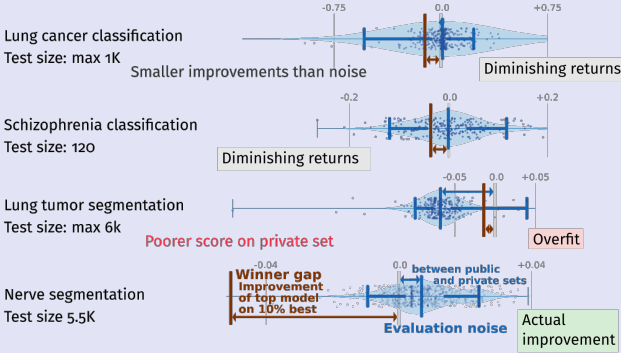
Takeaway

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

Challenges in AI benchmarking

► Futile benchmarks

[Varoquaux and Cheplygina 2022]



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. Saqib Sarfraz^{1,2}, Me-Yen Chen¹, Lukas Layer¹, Kunyu Peng², Marius Koutakis²

PNAS

RESEARCH ARTICLE | COMPUTER SCIENCES

OPEN ACCESS

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

Elfrat Shimron^{a,1}, Jonathan I. Tamir^{b,c,d}, Ke Wang^e, and Michael Lustig^f

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

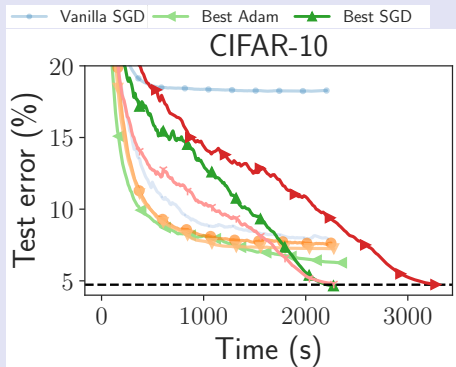
Robin M. Schmidt^{*1}, Frank Schneider^{*1}, Philipp Hennig^{1,2}

Unclear improvement!

Challenges in AI benchmarking

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

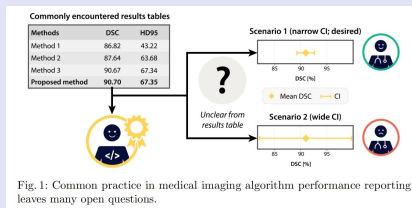
[Moreau et al. 2022]



Challenges in AI benchmarking

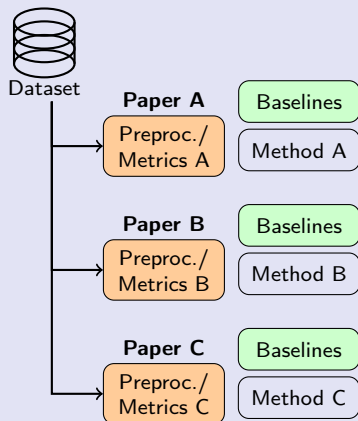
- ▶ Futile benchmarks
- ▶ Lack of proper baselines hinders scientific progress.
- ▶ Reproducing benchmarks is hard.
- ▶ Statistical validity is often missing.

[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.

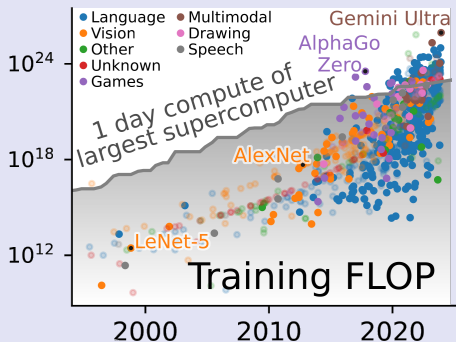


Each paper independently rebuilds preprocessing, baselines, and evaluation → duplicated cost.

Challenges in AI benchmarking

- ▶ Futile benchmarks
- ▶ Lack of proper baselines hinders scientific progress.
- ▶ Reproducing benchmarks is hard.
- ▶ Statistical validity is often missing.
- ▶ Benchmarking cost is duplicated across groups.
- ▶ Using AI infrastructure is complicated for many researchers.

[Varoquaux et al. 2025]



What AI research needs: shared evaluation infrastructure

Scientific fields mature when they build **shared infrastructure**:

- ▶ **Compute** – HPC clusters, cloud, SLURM, ...
- ▶ **Data** – HuggingFace, OpenML, UCI
- ▶ **Evaluation** – ???

What AI research needs: shared evaluation infrastructure

Scientific fields mature when they build **shared infrastructure**:

- ▶ **Compute** – HPC clusters, cloud, SLURM, ...
- ▶ **Data** – HuggingFace, OpenML, UCI
- ▶ **Evaluation** – ???

Evaluation is currently **rebuilt from scratch** for every paper:
challenges are tackled individually

What AI research needs: shared evaluation infrastructure

Scientific fields mature when they build **shared infrastructure**:

- ▶ **Compute** – HPC clusters, cloud, SLURM, ...
- ▶ **Data** – HuggingFace, OpenML, UCI
- ▶ **Evaluation** – ???

Evaluation is currently **rebuilt from scratch** for every paper: challenges are tackled individually

The missing layer

A community-maintained, extensible, validated framework for

- ▶ defining evaluation protocols
- ▶ contributing datasets & baselines
- ▶ reproducing & extending results

⇒ This is what Benchopt provides.

Reproducible method comparison with Benchopt



References



Making runnable benchmarks with benchopt



benchopt provides a framework to organize and run benchmarks

Examples of existing benchmarks:

- ▶ **Image Classification (resnet)**
- ▶ **NanoGPT Optimization**
- ▶ **TSFMs evaluation**
- ▶ **Inverse problem resolution**
- ▶ **Unsup. Domain Adaptation**
- ▶ **Bilevel Optimization**
- ▶ **Brain Computer Interface**
- ▶ **...**

Benchopt scope and contributors

	Short-term	Long-term
Task-specific	Challenge	SOTA tracking
General.	Research Q.	Field Benchmark

⇒ Community-maintained,
multi-institution, extensible
benchmarks

Contributors from...



Structure of a benchmark

3 components: Objective, Datasets, Solvers

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

Modular & extendable

- ▶ New metric? modify objective

Structure of a benchmark

3 components: Objective, Datasets, Solvers

```
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

3 components: Objective, Datasets, Solvers

```
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

3 components: Objective, Datasets, Solvers

```
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

3 components: Objective, Datasets, Solvers

```
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:

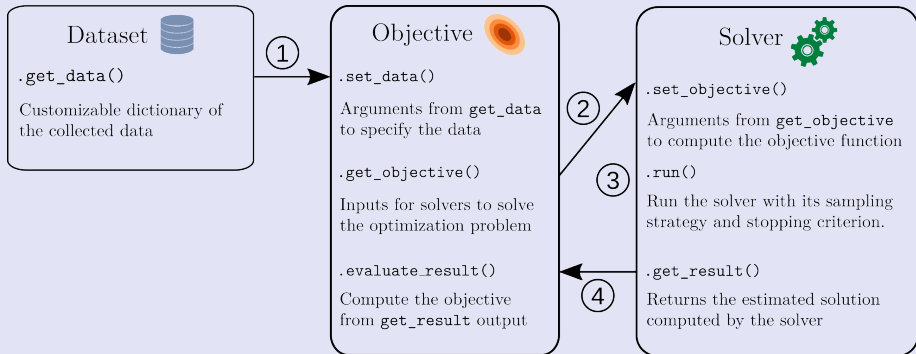
https://github.com/benchopt/template_benchmark

https://github.com/benchopt/template_benchmark_ml

A modular framework to create benchmarks

3 components: Objective, Dataset, Solver

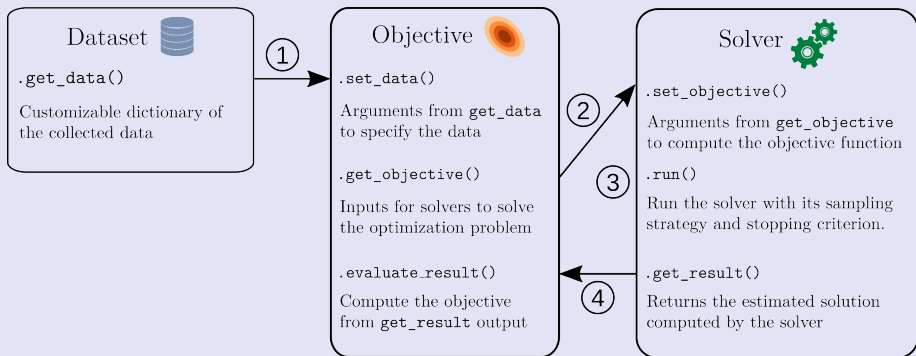
Dependency relation between Dataset - Objective - Solver



A modular framework to create benchmarks

3 components: Objective, Dataset, Solver

Dependency relation between Dataset - Objective - Solver



⇒ Benchopt defines the interface between components.

Creating a benchmark: the scientist's task

Define the 3 components for your problem:

Benchmark	Dataset	Solver	Objective
Lasso regression	(X, y)	sparse solver	suboptimality gap
DL optimization	dataloader + network	optimizer	val. accuracy
Brain-Computer Interface	EEG recordings + trial labels	decoder	classif. accuracy on 5-folds
Storage infra. evaluation	Large file dataset	dataloading strategy	throughput

Creating a benchmark: the scientist's task

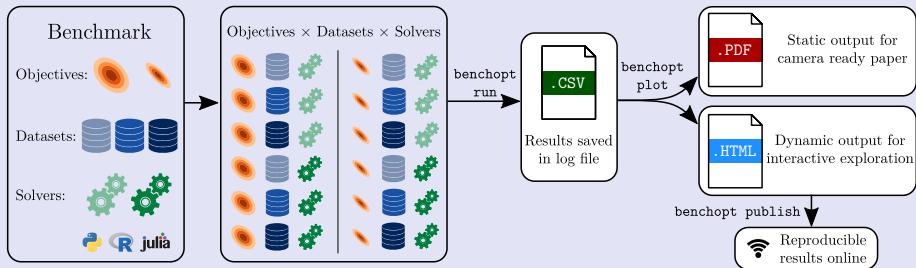
Define the 3 components for your problem:

Benchmark	Dataset	Solver	Objective
Lasso regression	(X, y)	sparse solver	suboptimality gap
DL optimization	dataloader + network	optimizer	val. accuracy
Brain-Computer Interface	EEG recordings + trial labels	decoder	classif. accuracy on 5-folds
Storage infra. evaluation	Large file dataset	dataloading strategy	throughput

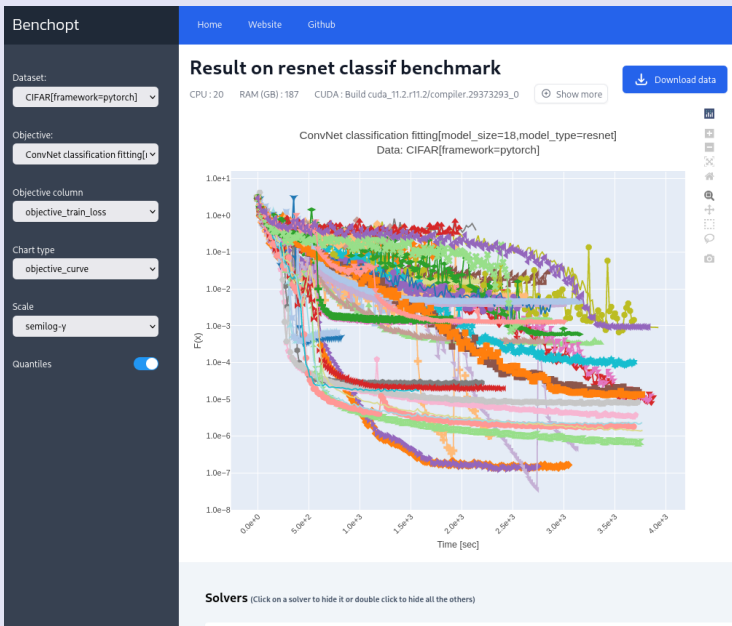
⇒ The benchmark structure is the same: only the components change.

Benchmark workflow

Steps: Install/Prepare, Test, Run, Explore, Publish

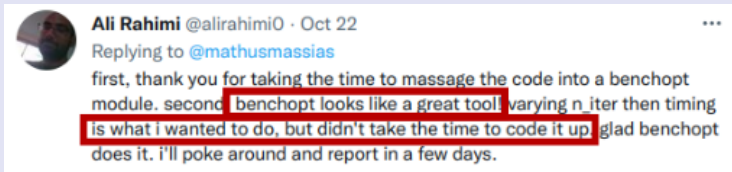


Interactive results exploration



Benchopt makes your life easy

- ▶ **Integrate broadly:** use implementations from Python, R, Julia, or binaries.
- ▶ **Scale experiments:** run in parallel locally or on HPC clusters, natively integrate with SLURM, ...
- ▶ **Save time:** cache results to avoid recomputing unchanged runs.
- ▶ **Trust comparisons:** control randomness with seeds and stable protocols.
- ▶ **Share outcomes:** merge and publish results from multiple runs, easy interactive visualization.



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.
AccGrad	(Lyu et al., 2018)	HyperAdam	(Wang et al., 2019b)
ACCP	(Zhang et al., 2020)	K-RFGSG (RFGSL)	(Gaidfad et al., 2020)
AdaAkr	(Oke et al., 2019)	KF-QN-CNN	(Oke & Gokhraf, 2021)
AdaBach	(Devadas et al., 2017)	KRAC	(Sharma & Gao, 2015)
AdaBayes/AdaBayes5	(Aichison, 2020)	KFL/KFRA	(Oke et al., 2017)
AdaBelief	(Zhang et al., 2020)	Lakshminf.Adamovans	(Ratnak & Maras, 2019)
AdaBound	(Luo et al., 2019)	LARS	(Chen et al., 2020)
AdaBound	(Luo et al., 2019)	LaProp	(Zou et al., 2020)
AdaCmp	(Chen et al., 2018)	LARS	(Liu et al., 2017)
AdaDelta	(Gohar, 2012)	LRFPT	(Amirgholizadeh, 2021)
AdaDelc	(Shauer & Stern, 2011)	LookAhead	(Zhang et al., 2019)
AdaF1	(Hao et al., 2019)	MS-ADAM	(Bhat & Huang, 2019)
AdaFem	(Chen et al., 2019a)	MADGRAD	(Defazio & Johnson, 2021)
AdaFTRL	(Rubanu & Pui, 2017)	MAS	(Gholvafard et al., 2020)
Adagrad	(Duchi et al., 2011)	MIRA	(Chen et al., 2020a)
ADAHESSEAN	(Yao et al., 2020)	MfAdam	(Makni & Wolf, 2020)
AdaL	(Oke et al., 2020)	MFAC-IMVBC-2	(Chen & Zhou, 2020)
AdaLms	(Sutton et al., 2019)	NAdam	(Dong, 2016)
Adam	(Kingma & Ba, 2015)	NMSBINAMSG	(Chang et al., 2019b)
Adam*	(Liu et al., 2020b)	NO-Adam	(Zhang et al., 2017a)
AdamAL	(Tan et al., 2019)	Noto	(Liu et al., 2020a)
Adamax	(Kingma & Ba, 2015)	NoWarm	(Nemayev, 2015)
AdamES	(Liu et al., 2020a)	Noisy Adam/Noisy K-RAC	(Zhang et al., 2019)
AdamNC	(Kohli et al., 2018)	NoxAdam	(Huang et al., 2019)
AdaMod	(Ding et al., 2019)	Noxgrad	(Zhang et al., 2019)
AdaP/PGSP	(Oke et al., 2021)	NTSGD	(Zhang et al., 2019a)
AdamT	(Zhou et al., 2020)	Palam	(Chen et al., 2020a)
AdamW	(Loshchilov & Hutter, 2019)	PIRG	(Liu et al., 2020b)
AdamX	(Tan & Peng, 2019)	PNL	(Moshchuk & Zettl, 2020)
ADAS	(Elyasbi, 2020)	PopAdam	(Owens et al., 2019)
AdaS	(Hossein & Panayiotou, 2020)	PyAd	(Petrov, 1960)
AdaScale	(Johnson et al., 2020)	PowerSGD/PowerGEM	(Vergos et al., 2019)
AdaSGD	(Wang & Wiers, 2020)	PyRobust/PyRob	(de Kroon et al., 2021)
AdaShw	(Zhou et al., 2019)	PyRob	(Mishchenko & Huang, 2017)
AdaSpn	(Hao et al., 2019)	PSwarm	(Oke, 2020)
AdaSpn	(Oke et al., 2019)	QRAdam/GEM	(Ma & Taylor, 2019)
AdaSUN/AdaW	(Li et al., 2020a)	RAdam	(Liu et al., 2020)
ADAF	(Liu & Tian, 2020)	Ranger	(Wright, 2020a)
ADAF	(Berrada et al., 2020)	RangerAdam	(Wright, 2020b)
AMCBound	(Luo et al., 2019)	RMSProp	(Fedorak & Hinton, 2012)
AMSGrad	(Kohli et al., 2018)	RMSizes	(Chen et al., 2019)
AsquaterGrad	(Roy et al., 2021)	S-GRD	(Chang et al., 2020)
AspajLS	(Vucelja et al., 2019)	SAVan	(Wang et al., 2020a)
ASGD	(Oke et al., 2019b)	SAdamSAMSGrad	(Chen et al., 2019)
ASAM	(Kwon et al., 2021)	SALR	(Chen et al., 2020)
Asud RES	(Oke et al., 2021)	SAM	(Poon et al., 2021)
AsudGrad	(Goussier et al., 2019)	SC-Adaptive/SC-RMSProp	(Makridakis & Hinton, 2017)
BAkers	(Oke et al., 2018)	SDProp	(Mao et al., 2017)
BCAdam	(Hao & Zhang, 2019)	SGD	(Robbins & Monro, 1951)
BBFwd	(Zhang et al., 2017a)	SGD-BB	(Gao et al., 2016)
BRMSProp	(Aichison, 2020)	SGD-G2	(Ayadi & Fattouh, 2020)
BSGD	(Hao et al., 2020)	SGEM	(Ramanam-Krishnan et al., 2021)
C-ADAM	(Sutton et al., 2020)	SGERms	(Tan & Cardo, 2021)
CADA	(Chen et al., 2021)	SGERM	(Liu & Luo, 2020)
Class Momentum	(Keriven & Robik, 2020)	SGLR	(Loshchilov & Hutter, 2017)
CPW	(Phan-Chai & Kijakiatkai, 2019)	SHAdagrad	(Huang et al., 2020)
CurvBall	(Krogh et al., 2019)	Shampoo	(Oishi et al., 2020; Gupta et al., 2019)
Dalton	(Nouri et al., 2020)	SignAdam++	(Wang et al., 2020)
DeepMemory	(Wright, 2020a)	SignSGD	(Bottani et al., 2019)
DDNGy	(Liu et al., 2021a)	SKNGS+GN	(Yang et al., 2020)
DFPwd	(Huby et al., 2020)	SM2	(Oishi et al., 2019)
EAdam	(Yuan & Gao, 2020)	SM3	(Tan 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

Results: github.com/SirRob1997/Cro...

Video: 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.

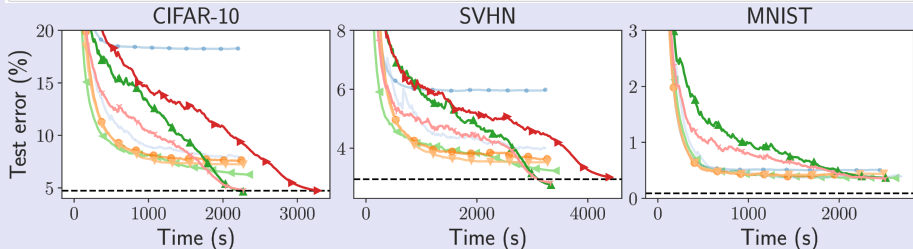
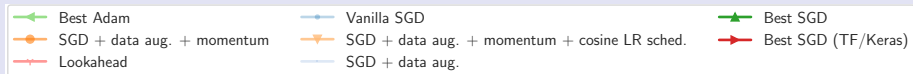
⇒ But this benchmark cannot be easily reproduced!

Example: Optimization for ResNet on image classification

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



Z. Ramzi



https://github.com/benchopt/benchmark_resnet_classif/

Example: Large scale-optimization for Deep learning

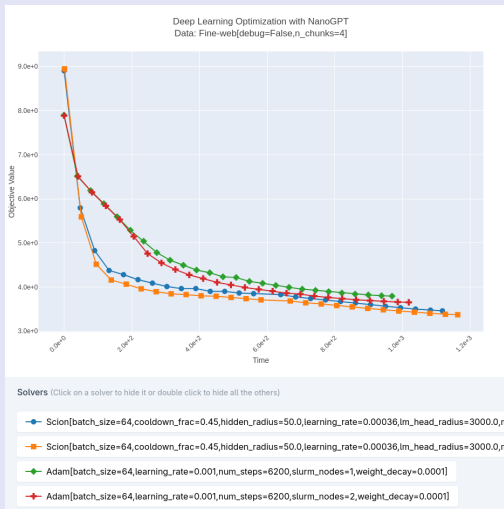
- ▶ Use modern large-scale datasets and models
- ▶ Classical optimization (schedulers, momentum, ...)
- ▶ Distributed training and mixed precision



S. Vaiter



T. Silveti-Falls



https://github.com/benchopt/benchmark_nanogpt/

Example: Inverse problem resolution

- ▶ Provide a series of tasks: one dataset + one physics
- ▶ An easy interface to evaluate reconstruction models
- ▶ Target to auto-run and upload on huggingface hub

```
from dbench import run_benchmark
import deepinv as dinv

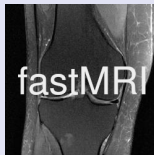
# Create reconstruction method
my_solver = dinv.models.RAM()
run_benchmark(
    my_solver, "benchmark_name"
)
```



J. Tachella



Dataset example



HF hub leaderboard

A screenshot of the Hugging Face Hub leaderboard for the fastMRI dataset. The table shows various models and their performance metrics. The top row is highlighted in blue.

Model	Score	Rank	Model	Score	Rank
deepinv/...	0.9999	1	deepinv/...	0.9999	1
deepinv/...	0.9998	2	deepinv/...	0.9998	2
deepinv/...	0.9997	3	deepinv/...	0.9997	3
deepinv/...	0.9996	4	deepinv/...	0.9996	4
deepinv/...	0.9995	5	deepinv/...	0.9995	5
deepinv/...	0.9994	6	deepinv/...	0.9994	6
deepinv/...	0.9993	7	deepinv/...	0.9993	7
deepinv/...	0.9992	8	deepinv/...	0.9992	8
deepinv/...	0.9991	9	deepinv/...	0.9991	9
deepinv/...	0.9990	10	deepinv/...	0.9990	10

<https://github.com/deepinv/benchmarks>

Example: Large-scale Inverse problem resolution

- ▶ A few very large-scale tasks: radio-astronomy, CT, ...
- ▶ Inference, Training and workflows
- ▶ Target measuring distributed performances

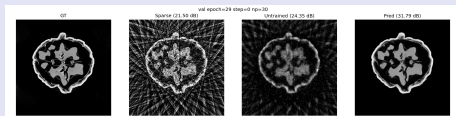
Inference

Training

Workflow



B.
Malezieux



https://github.com/bmalezieux/benchmark_invprob_largescale

Many research results are not maintained after publication:

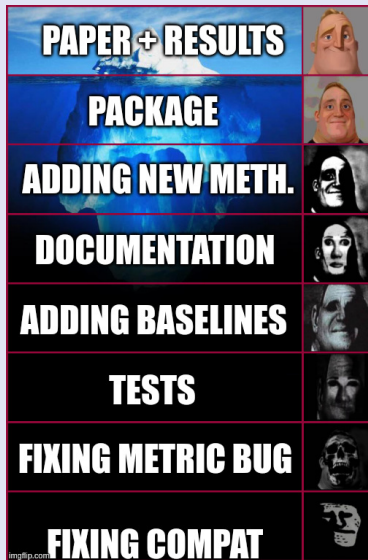
- ▶ Every PhD creates a package;
- ▶ Every post-doc abandons one;
- ▶ The ecosystem grows horizontally, not vertically

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



The benchopt roadmap

Grow the benchmark collection

- ▶ Zero-th order optimization
- ▶ Benchmarking dataloading strategies on HPC infrastructures
- ▶ ... and more community contributions

⇒ Open question: how to benchmark foundation models?

Improve benchmarking methodology

- ▶ **Statistical validity:** how many samples / CV splits to trust a ranking?
→ work w/ C. Eve and G. Varoquaux
- ▶ **Foundation model evaluation:** few-shot, prompted, fine-tuned – incomparable by standard metrics.
- ▶ **Frugality:** benchmark under compute constraints, not just final accuracy.

Thank you for your attention!

Three asks:

1. **Contribute** a solver or dataset to an existing benchmark.
2. **Consider Benchopt** as infrastructure for your next benchmark, in a paper or with broader scope.
3. **Star the repo** – it matters for visibility!



BenchOpt



Want to run a benchmark sprint?
Come talk to me!