# Reproducibility in AI benchmarking: lessons from Benchopt

Thomas Moreau

# A word about me

- ▶ Researcher at Inria – MIND

- ▶ `https://tommoral.github.io`

- ▶ 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.**



$\Rightarrow$ Demonstrates the importance of benchmarks to drive research in AI.
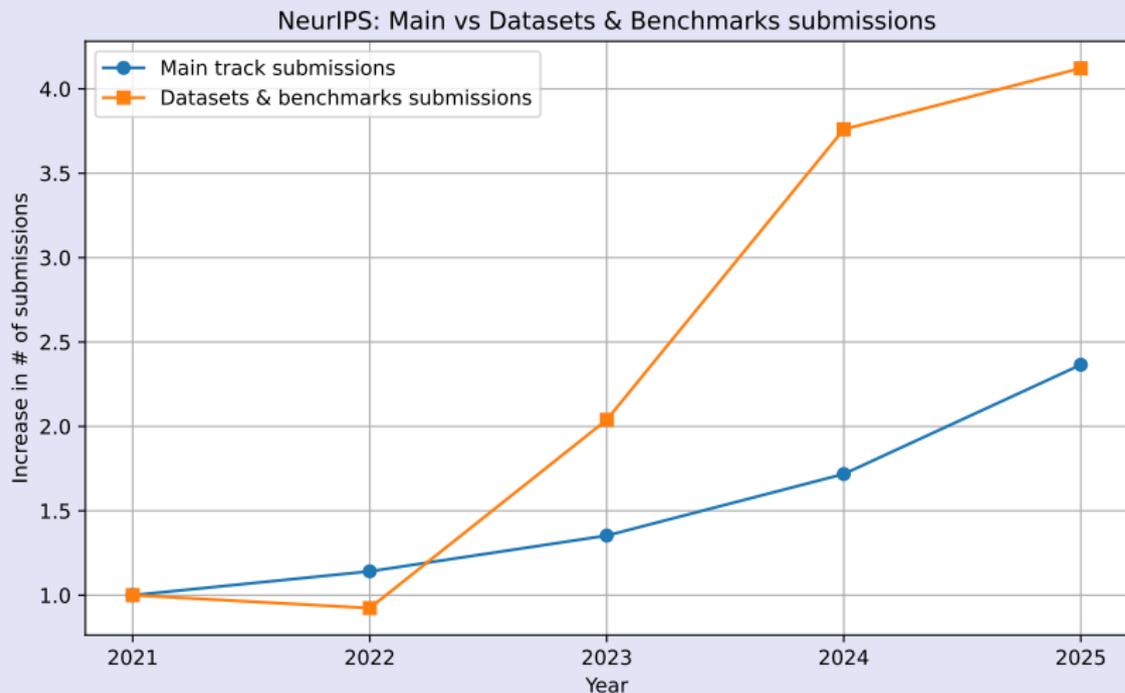
**Many benchmarks followed ImageNet:**

- ▶ Natural Language Processing: GLUE, SuperGLUE

- ▶ Reinforcement Learning: Atari, MuJoCo, OpenAI Gym

- ▶ Others: fastMRI, DAWNBench, MLPerf, etc.

**Many benchmarks followed ImageNet:**

► Natural Language Processing: GLUE, SuperGLUE

► Reinforcement Learning: Atari, MuJoCo, OpenAI Gym

► Others: fastMRI, DAWNBench, MLPerf, etc.



$\Rightarrow$ Benchmarks are now ubiquitous in AI research.

# Too many benchmarks in AI?

In the recent years, many benchmarks have been proposed:



NeurIPS: Main vs Datasets & Benchmarks submissions

$\Rightarrow$ Most of them don't have long-term maintainance plan, and are quickly abandoned

## Benchmark goals in AI

|  | **Short-term progress** | **Long-term evaluation** |
|---|---|---|
| **Task-specific** | Challenge/Competition $\rightarrow$ push limits quickly | SOTA tracking $\rightarrow$ measure progress |
| **Generalizable** | Research question $\rightarrow$ empirical study | Benchmark framework $\rightarrow$ stable & extensible |

## Benchmark goals in AI

|  | **Short-term progress** | **Long-term evaluation** |
|---|---|---|
| **Task-specific** | Challenge/Competition $\rightarrow$ push limits quickly | SOTA tracking $\rightarrow$ measure progress |
| **Generalizable** | Research question $\rightarrow$ empirical study | Benchmark framework $\rightarrow$ 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.
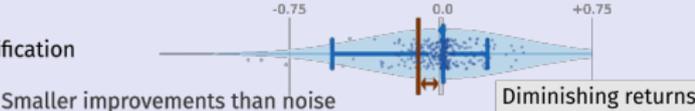
- ▶ **Challenge benchmarks:** Fixed dataset + single metric → stresses solver comparison.                    Focus on performance
- ▶ **SOTA tracking benchmarks:** Multiple metrics and datasets, compare solvers to measure progress over time.

  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
Diminishing returns

Schizophrenia classification
Test size: 120
Diminishing returns

Lung tumor segmentation
Test size: max 6k
Poorer score on private set
Overfit

Nerve segmentation
Test size 5.5K
Winner gap Improvement of top model on 10% best
between public and private sets
Evaluation noise
Actual improvement

- ▶ 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[1,2]  Mei-Yen Chen[1]  Lukas Layer[1]  Kunyu Peng[2]  Marios Koulakis[2]



PNAS    RESEARCH ARTICLE | COMPUTER SCIENCES    OPEN ACCESS

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

Efrat Shimron [a,1], Jonathan I. Tamir [b,c,d], Ke Wang[a], and Michael Lustig[a]



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

[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.

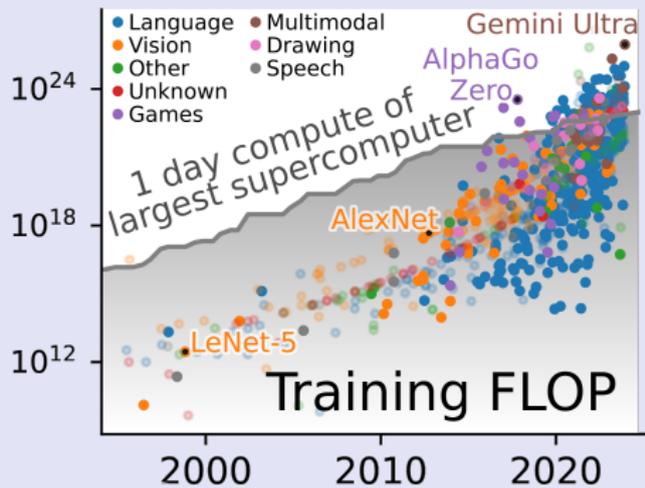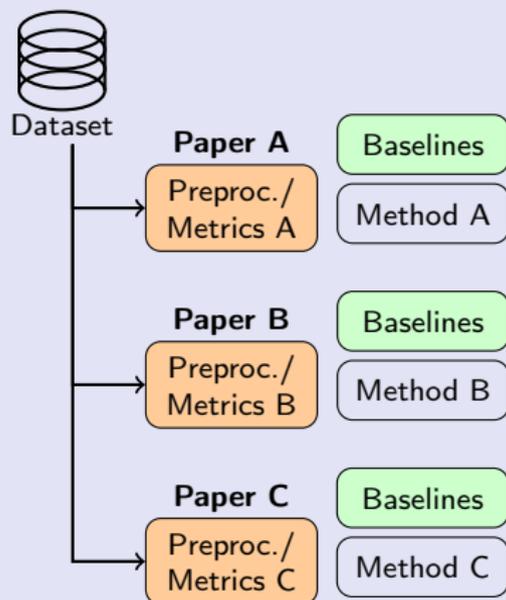▶ Reproducing benchmarks is hard.

▶ Statistical validity is often missing.

[Christodoulou et al. 2024]



Fig. 1: Common practice in medical imaging algorithm performance reporting leaves many open questions.

# 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.

## 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 $\rightarrow$ duplicated cost.

# Reproducible method comparison with Benchopt

 . . .

**References**

▶ **TM**, Massias, M., Gramfort, A., Ablin, P., Bannier, P.-A., Charlier, B.,
Dagréou, M., la Tour, T. D., Durif, G., Dantas, C. F., Klopfenstein, Q.,
Larsson, J., Lai, E., Lefort, T., Malézieux, B., Moufad, B., Nguyen, B. T.,
Rakotomamonjy, A., Ramzi, Z., Salmon, J., and Vaiter, S. (2022). Benchopt:
Reproducible, efficient and collaborative optimization benchmarks. In *NeurIPS*

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**
- **. . .**

## 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

**3 components**: Objective, Dataset, Solver

Dependency relation between Dataset - Objective - Solver



**Dataset**

`.get_data()`

Customizable dictionary of the collected data

**Objective**

`.set_data()`

Arguments from `get_data` to specify the data

`.get_objective()`

Inputs for solvers to solve the optimization problem

`.evaluate_result()`

Compute the objective from `get_result` output

**Solver**

`.set_objective()`

Arguments from `get_objective` to compute the objective function

`.run()`

Run the solver with its sampling strategy and stopping criterion.

`.get_result()`

Returns the estimated solution computed by the solver

① ② ③ ④

**3 components**: Objective, Dataset, Solver

Dependency relation between Dataset - Objective - Solver



⇒ Benchopt defines the interface between components.

## Explicit requirements and parameters

```python
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

```python
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

```python
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)
```

$\Rightarrow$ Reproducible benchmark by design!

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

**Steps**: Install, 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.
- ▶ **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 visualization.



**Ali Rahimi** @alirahimi0 · Oct 22
Replying to @mathusmassias
first, thank you for taking the time to massage the code into a benchopt module. second, benchopt looks like a great tool! varying n_iter then timing is what i wanted to do, but didn't take the time to code it up. glad benchopt does it. i'll poke around and report in a few days.

# 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.

$\Rightarrow$ Many novel methods but unclear improvements.

$\Rightarrow$ 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/

# Example: Large scale-optimization for Deep learning

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



https://github.com/tommoral/benchmark_nanogpt/

# 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

$\Rightarrow$ Limited incentives in AI to maintain a codebase

## The benchopt roadmap

**Grow the benchmark collection**

- ▶ Domain adaptation (SKADA)

  **⑤ skada-bench** (Public)
  generated from benchopt/template_benchmark

- ▶ Brain-computer interfaces (MOABB)

  **⑥ benchmark_bci** (Public)
  generated from benchopt/template_benchmark

- ▶ ... 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?

  → this talk, Section 4

- ▶ **Foundation model evaluation:** few-shot, prompted, fine-tuned – incomparable by standard metrics.

- ▶ **Frugality:** benchmark under compute constraints, not just final accuracy.

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

**References**

▶ Eve, C., Varoquaux, G., and **TM** (2026). Crossing the Validation crisis: Cross-validation reduces benchmarking variance surprisingly well. Preprint

## Reproducible research

Three levels of reproducibility:

▶ **Repeatability:** reproduce the exact same results.
                Bitwise reproducibility – control seeds, hardware, versions.

▶ **Reusability:** provide tools others can apply to their use case.
                        Clean code, documentation, proper packaging.

▶ **Replicability:** extend results with new methods or datasets.
                The hard one in AI – stochasticity makes rankings unstable.

## Reproducible research

Three levels of reproducibility:

▶ **Repeatability:** reproduce the exact same results.
Bitwise reproducibility – control seeds, hardware, versions.

▶ **Reusability:** provide tools others can apply to their use case.
Clean code, documentation, proper packaging.

▶ **Replicability:** extend results with new methods or datasets.
The hard one in AI – stochasticity makes rankings unstable.

$\Rightarrow$ Replicability requires statistically valid evaluation protocols.

## Evaluating decision functions & learning algorithms

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

$$R_{\mathcal{D}}(g) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \sim \mathcal{D}} \ell(g(x), y)$$

The oracle score is:

$$\mathcal{R}^*(g) = \mathbb{E}_{(X,Y) \sim p_{X,Y}} [\ell(g(X), Y)]$$

## Evaluating decision functions & learning algorithms

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

$$R_{\mathcal{D}}(g) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \sim \mathcal{D}} \ell(g(x), y)$$

The oracle score is:

$$\mathcal{R}^*(g) = \mathbb{E}_{(X,Y) \sim p_{X,Y}} [\ell(g(X), Y)]$$

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

$$\begin{aligned}
\mathcal{R}^*_{n_{\text{tr}}}(F_\lambda) &= \mathbb{E}_{\mathcal{D}^{\text{tr}}} \mathbb{E}_\xi [\mathcal{R}^*(F_\lambda(\mathcal{D}^{\text{tr}}, \xi))] \\
&= \mathbb{E}_{\mathcal{D}^{\text{te}}} \mathbb{E}_{\mathcal{D}^{\text{tr}}} \mathbb{E}_\xi [R_{\mathcal{D}^{\text{te}}}(F_\lambda(\mathcal{D}^{\text{tr}}, \xi))]
\end{aligned}$$

## Benchmarking goal: compare and rank methods

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

$$\mathcal{R}^*_{n_{\mathrm{tr}}}(F_\lambda) \overset{?}{<} \mathcal{R}^*_{n_{\mathrm{tr}}}(F_{\lambda'})$$

## Benchmarking goal: compare and rank methods

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

$$\mathcal{R}^*_{n_{\text{tr}}}(F_\lambda) \overset{?}{<} \mathcal{R}^*_{n_{\text{tr}}}(F_{\lambda'})$$

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

$$R_{\mathcal{D}^{\text{te}}}(F_\lambda(\mathcal{D}^{\text{tr}}, \xi)) \overset{?}{<} R_{\mathcal{D}^{\text{te}}}(F_{\lambda'}(\mathcal{D}^{\text{tr}}, \xi'))$$

$\Rightarrow$ How to ensure that our benchmark results reflect the true oracle ranking?

# DE L'Importance
## DE
# la validation croisée

OÙ 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.

*ERECTO DECUS*

A Grenoble,
Chez BARRILLOT & FILS.

- **Study set**: simulates realistic conditions (100–10k samples); split into train/test via Cross-Validation (CV).
- **Outer set**: independent test evaluations with growing size $\alpha \cdot n_{\text{te}}$.
- **Benchmarking set**: very large – acts as quasi-oracle $\widehat{R}^{\star}(g)$.

## How CV reduces variance – theory

For $K$ i.i.d. splits (*ShuffleSplit / MCCV*):

$$\mathrm{Var}[\widehat{R}_K] = \frac{1}{K}\,\sigma_{\mathsf{HO}}^2 + \frac{K-1}{K}\,\tau$$

- $\sigma_{\mathsf{HO}}^2$: variance of a single hold-out evaluation.
- $\tau = \mathrm{Cov}[E_1, E_2] > 0$: covariance between two splits (shared samples).

## How CV reduces variance – theory

For $K$ i.i.d. splits (*ShuffleSplit / MCCV*):

$$\text{Var}[\widehat{R}_K] = \frac{1}{K}\sigma_{\text{HO}}^2 + \frac{K-1}{K}\tau$$

- $\sigma_{\text{HO}}^2$: variance of a single hold-out evaluation.
- $\tau = \text{Cov}[E_1, E_2] > 0$: covariance between two splits (shared samples).

### Key insight
As $K \to \infty$, variance floors at $\tau$ – **irreducible** due to finite data.
But in practice, the gains for moderate $K$ are *much larger than expected*.

## How CV reduces variance – theory

For $K$ i.i.d. splits (*ShuffleSplit / MCCV*):

$$\text{Var}[\widehat{R}_K] = \frac{1}{K}\,\sigma_{\text{HO}}^2 + \frac{K-1}{K}\,\tau$$

- $\sigma_{\text{HO}}^2$: variance of a single hold-out evaluation.
- $\tau = \text{Cov}[E_1, E_2] > 0$: covariance between two splits (shared samples).

### Key insight
As $K \to \infty$, variance floors at $\tau$ – **irreducible** due to finite data.
But in practice, the gains for moderate $K$ are *much larger than expected*.

$$\Rightarrow \text{How to quantify these gains in interpretable units?}$$

# Sample gain: quantifying CV efficiency
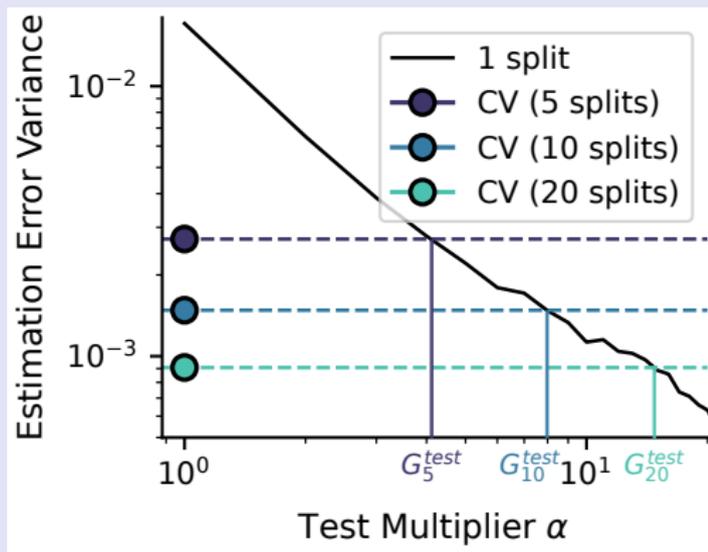
**Variance-Equivalent Test Sample Gain $G_K^{\text{test}}$:**

Number of times larger a *single* test set would need to be to match the variance of $K$-split CV:

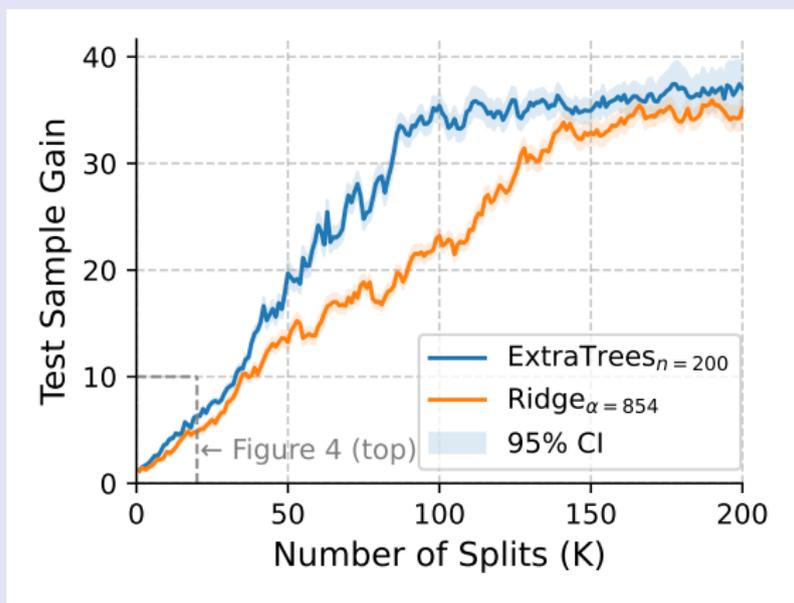$$G_K^{\text{test}} = \frac{N_K^{\text{equiv}}}{n_{\text{te}}}$$

A value of 5 means:
*K splits $\approx$ 5$\times$ more test data.*



Variance of estimation error vs. test-set multiplier $\alpha$ (1 split) compared to multi-split CV. $G_K^{\text{test}}$ is read off on the $x$-axis.
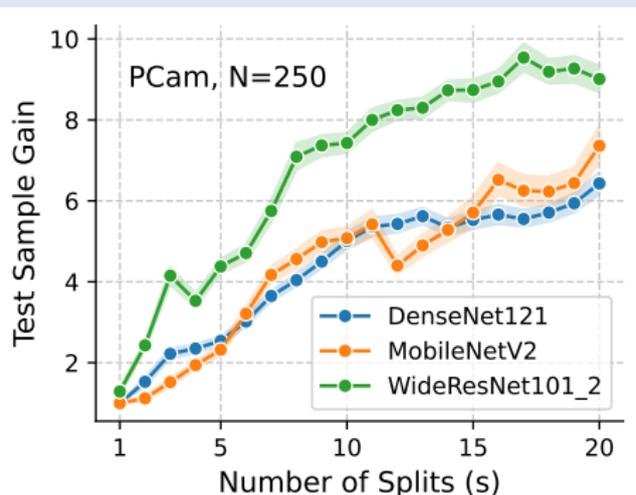
# Sample gains on synthetic data



Variance-Equivalent Test Sample Gain $G_K^{\text{test}}$ on simulated data ($n_{\text{tr}} = 1000$).
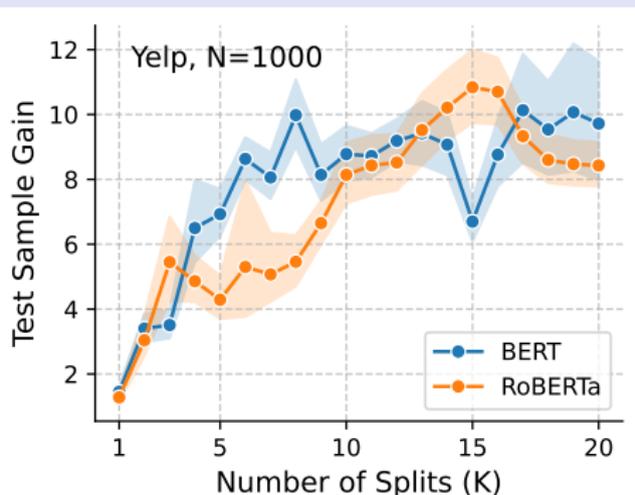
$\Rightarrow$ Gains of 30–40$\times$ with $K = 200$ – diminishing returns set in later than expected.

# Sample gains on real datasets (PCam & Yelp)



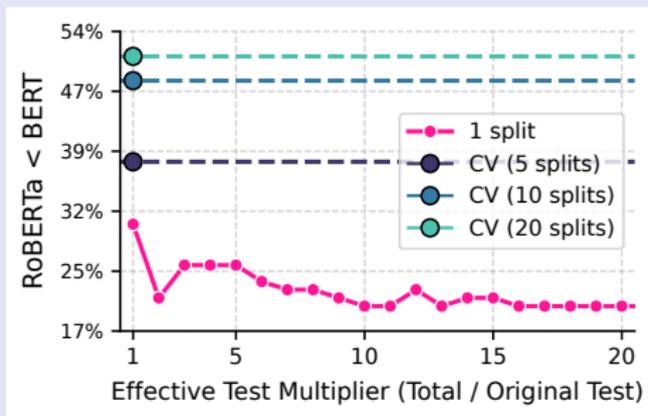**PatchCamelyon** ($N = 250$)
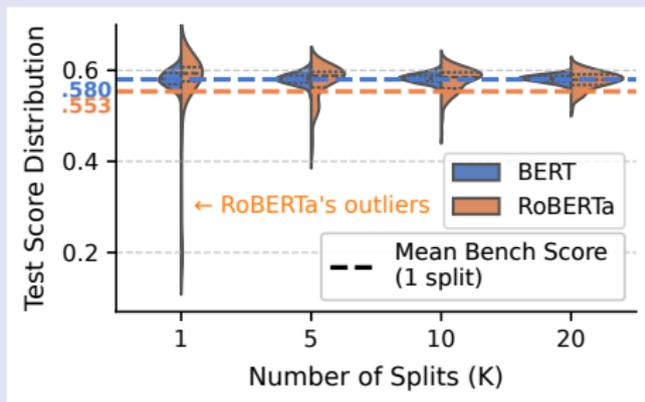
**Yelp Review Full** ($N = 1000$)

Consistent gains across medical imaging and NLP. BERT at $N = 1000$: $G_{20}^{\text{test}} \approx 10$ (20 splits $\approx 10\times$ more test data).

# The RoBERTa case: ranking of means ≠ mean of rankings



Proportion of runs where `RoBERTa` < `BERT` (oracle ranking) is retrieved, vs. number of splits. Single splits favour `RoBERTa` 70–80% of the time.

Score distributions at $N = 3750$: `RoBERTa` has severe low outliers that only multi-split CV averages out.

⇒ CV exposes pathological behaviours that single splits mask.

## Conclusion

Reproducible research needs more than just releasing code:

- ▶ Reusable $\rightarrow$ Clean and Documented.
- ▶ Extendable $\rightarrow$ Proper packaging and maintenance.
- ▶ Statistically valid $\rightarrow$ Proper evaluation protocols.

## Conclusion

Reproducible research needs more than just releasing code:

- ▶ Reusable $\rightarrow$ Clean and Documented.
- ▶ Extendable $\rightarrow$ Proper packaging and maintenance.
- ▶ Statistically valid $\rightarrow$ Proper evaluation protocols.

### Takeaway on statistical validity

- ▶ Stochasticity is intrinsic to AI – single splits can be misleading.
- ▶ CV substantially reduces benchmarking variance (sample gain up to $10\times$ on real datasets).
- ▶ **When resources allow, maximize CV splits.**

Don't hesitate to star the benchopt repo!