

Practical Machine Learning for Streaming Data

ACM SIGKDD Tutorial (hands-on)
2024

Heitor Murilo Gomes¹, Albert Bifet^{2,3}

[https://adaptive-machine-learning.github.io/
kdd2024_ml_for_streams/](https://adaptive-machine-learning.github.io/kdd2024_ml_for_streams/)



[1] Victoria University of Wellington, New Zealand, [2] University of Waikato, New Zealand,
[3] TELECOM Paris, LCTI, France.

Heitor Murilo Gomes

Senior lecturer at the Victoria University of Wellington (VuW) in New Zealand. Before joining VuW, Heitor was co-director of the AI Institute at the University of Waikato.

PI for a few research projects ranging from applied to fundamental research (i.e. ML for energy distribution, novel SSL approaches for DS, ...).

Leads the **capymoa** open source library for data stream learning, and provide support for **MOA** (Massive On-line Analysis).

<https://heitorgomes.com/>



Albert Bifet

Professor of AI and the Director of the AI Institute at University of Waikato, and Professor of Big Data at Data, Intelligence and Graphs (DIG) LTCl, Télécom Paris, IP Paris. Co-chair of the NZ AI Researchers Association. Leads the TAI AO Environmental Data Science project and co-author of the book “Machine Learning from Data Streams” published at MIT Press.

Co-leads the open source project **MOA** Massive On-line Analysis, and provide advice/support for several other projects such as **capymoa**

<https://albertbifet.com/>



About this tutorial

- **Our goal:** Introduce attendees to diverse machine-learning tasks for streaming data applications

About this tutorial

- **Our goal:** Introduce attendees to diverse machine-learning tasks for streaming data applications
 - Classification, regression, ensemble learning, prediction intervals, concept drifts, partially and delayed streams, clustering, anomaly detection, ...

About this tutorial

- **Our goal:** Introduce attendees to diverse machine-learning tasks for streaming data applications
 - Classification, regression, ensemble learning, prediction intervals, concept drifts, partially and delayed streams, clustering, anomaly detection, ...
- **Beyond the introduction:** Enable attendees to apply and extend the concepts introduced using python notebooks and *capymoa*

Outline

- **Machine Learning for Streaming Data (intro)**

[KDD_2024_introduction.ipynb](#)

- Learning cycle
- Evaluation procedures
- Introduction to *capymoa*

- **Concept drifts**

[KDD_2024_drift.ipynb](#)

- Simulation, Detection & Evaluation

- **Supervised Learning**

[KDD_2024_supervised.ipynb](#)

- Classification
- Ensemble learning

- **Supervised Learning (cont.)**

- Regression
- Prediction Intervals

- **Advanced Topics**

[KDD_2024_advanced.ipynb](#)

- Partially and delayed labeled streams
- Clustering
- Anomaly detection

Notebooks: https://adaptive-machine-learning.github.io/kdd2024_ml_for_streams/



Tutorial Logistics

Schedule

- 10:00 to 10:45 - Slides + Code
- **10:45 to 11:00** - Break
- 11:00 to 11:45 - Slides + Code
- **11:45 to 12:00** - Work on examples/ exercises
- 11:00 to 11:45 - Slides + Code
- **12:45 to 13:00** - Solve exercises + final considerations

Installation: <https://capymoa.org/installation.html>



- Interleave slides with code
- Not a lot of time to work on the exercises, but we can use the time on the end as well

Notebooks: https://adaptive-machine-learning.github.io/kdd2024_ml_for_streams/



Machine Learning for Streaming Data

Stream Learning

What are data streams?

Sequences of items, possibly infinite, each item having a timestamp, and so a temporal order

Stream Learning

What are data streams?

Sequences of items, possibly infinite, each item having a timestamp, and so a temporal order

Machine learning for streaming data (or Stream learning)

Data items arrive one by one, and we would like to **build and maintain models**, such as patterns or predictors, of these items in real time (or near real time)

Stream Learning: Examples

Sensor data (IoT): energy demand prediction, environmental monitoring, traffic flow

Marketing and e-commerce: product recommendation, click stream analysis, sentiment analysis (social networks)

Cybersecurity: malware detection, spam detection, intrusion detection

And many more!*

Stream Learning

When should we abstract the data as a continuous stream?

Stream Learning

When should we abstract the data as a continuous stream?

can't store all the data; or

shouldn't store all the data

Stream Learning: **can't store**

Storing all the data may exceed the
available storage capacity or cause
practical limitations

The **volume or velocity** of incoming
data may be too high to store and
process in its entirety

Stream Learning: **shouldn't store**

Storing all the data may not be desirable due to privacy concerns, compliance requirements, or the nature of the problem

For example, if we are only interested in **real-time analysis** or **immediate decision-making**

Stream Learning

Using a stream abstraction, we can process the data incrementally, **focusing** on the **most recent** or **relevant** data points, and **discard** or **aggregate** the older data as needed

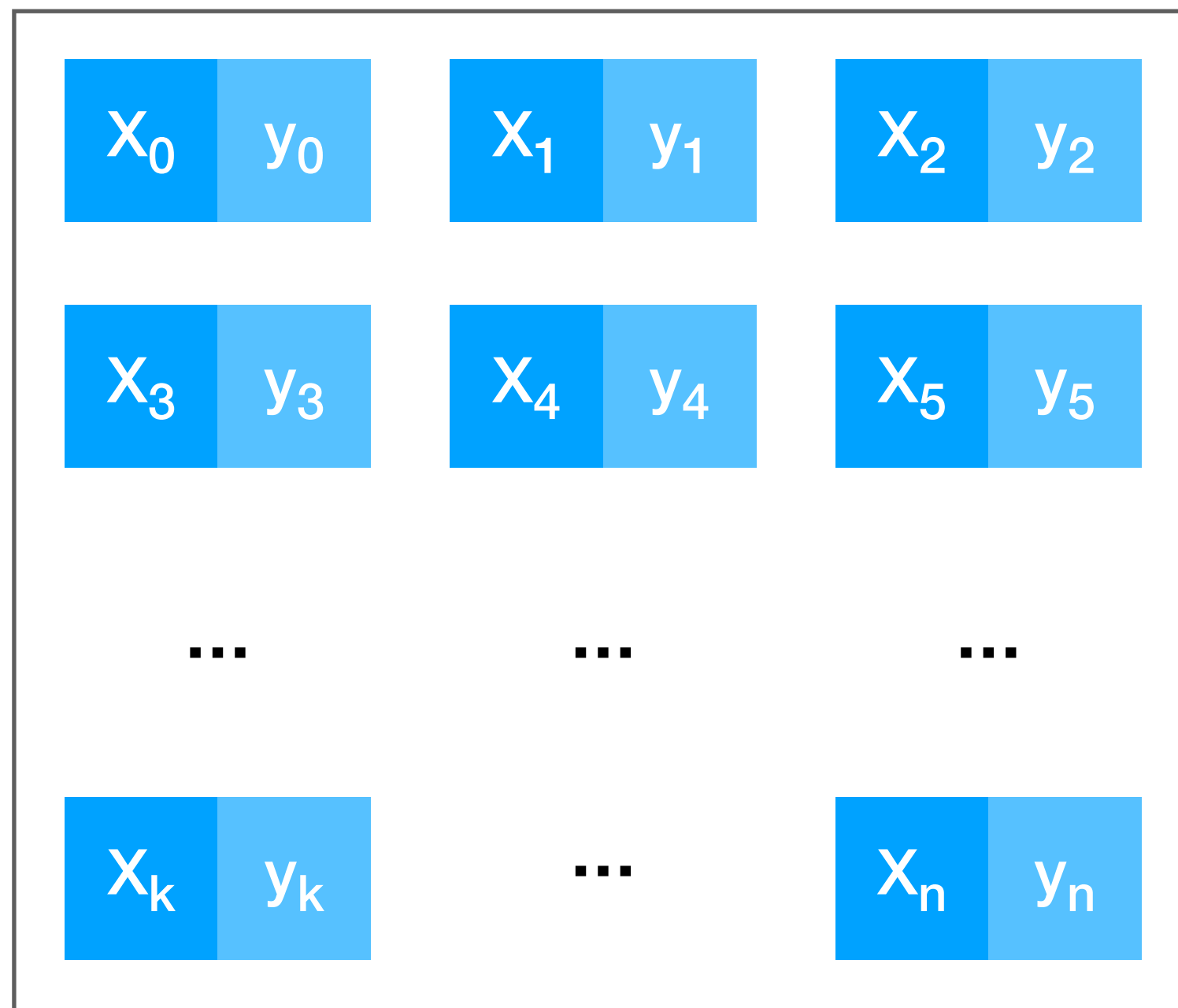
Stream Learning

ML for **Batch** (“static”) data

vs.

ML for **Streaming** (“online”) data

ML for Batch data



Fixed size dataset

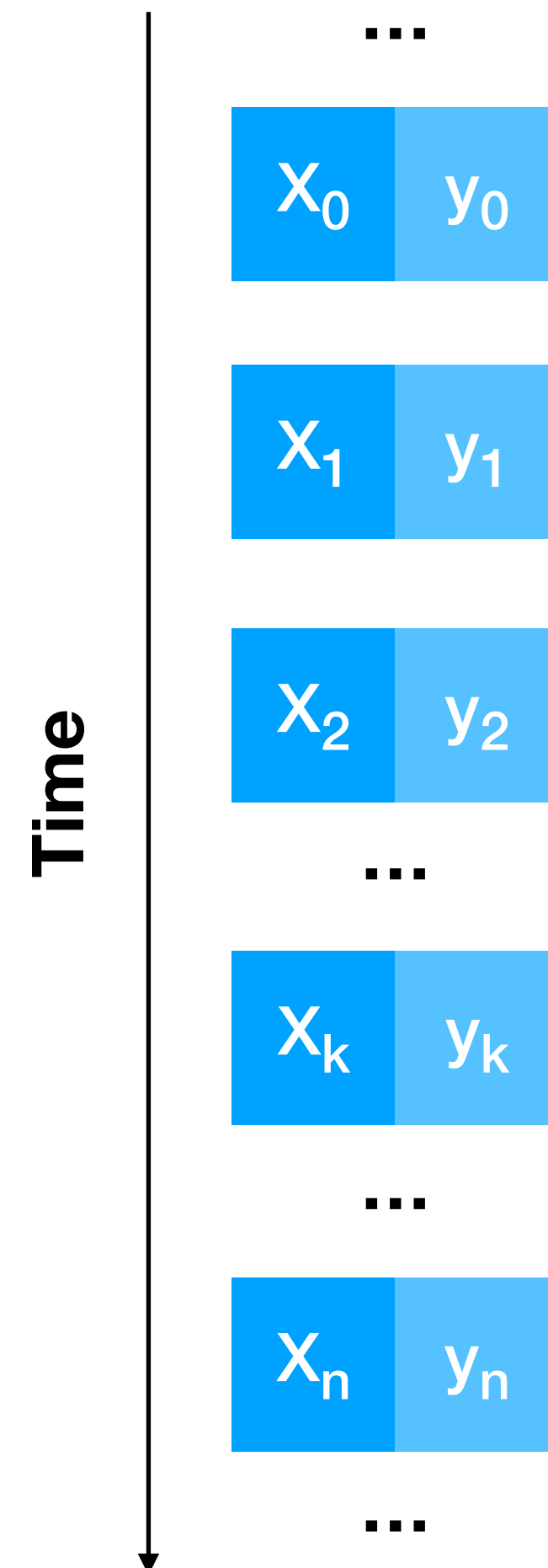
Random access to any instance

Well-defined phases
(Train, Validation, Test)

Challenges

noise, missing data,
imbalance, high
dimensionality, ...

ML for Streaming data



Continuous flow of data

Limited time to inspect
data points

Interleaved phases
(Train, Validation, Test)

Challenges

Concept drifts, concept
evolution, strict memory/
processing requirements,
may more and...
inherit all those from batch

Batch vs. Streaming

Batch data



The output is a **trained model**

Streaming data



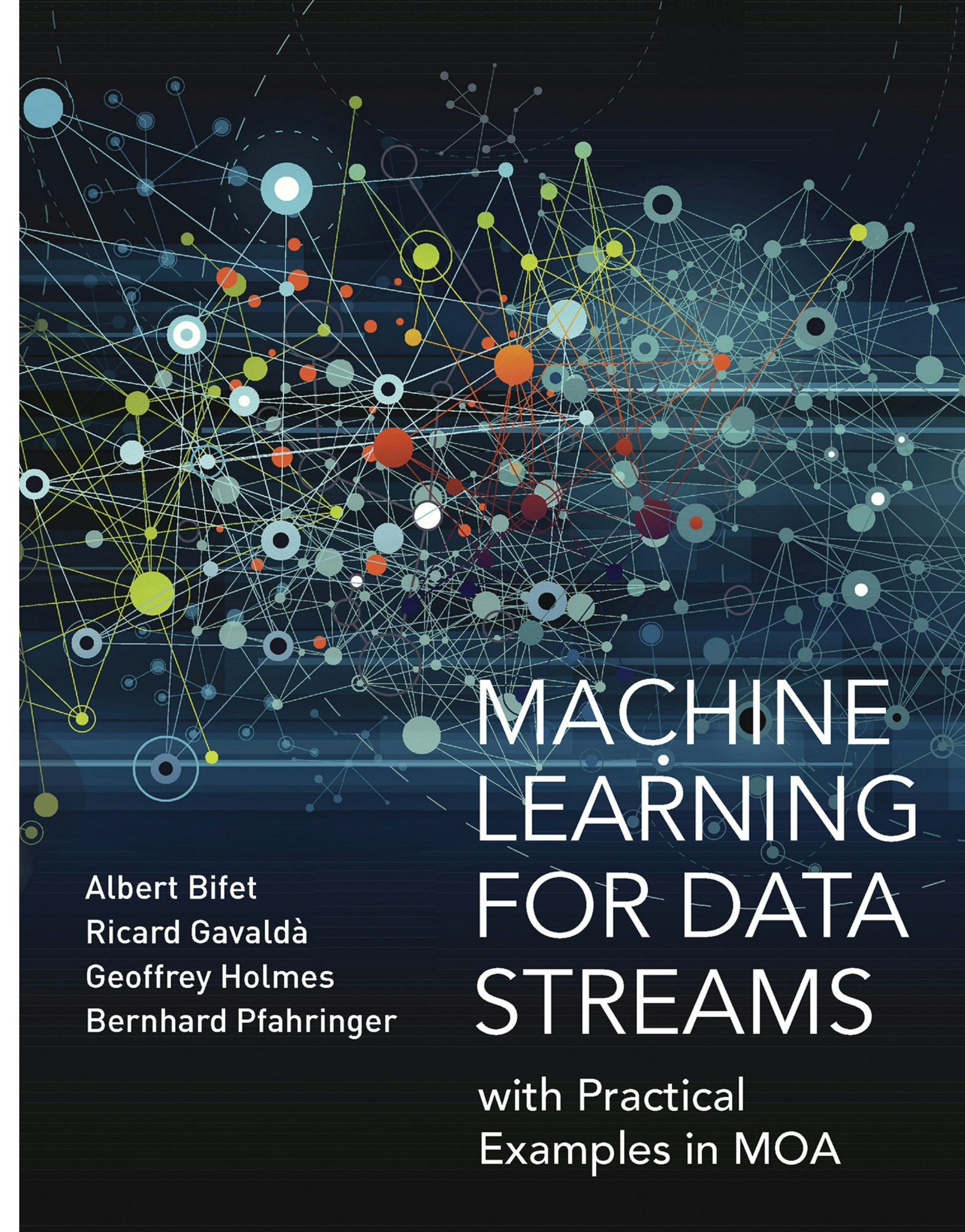
The output is a trainable model

Other resources

The “Machine Learning for Data Streams with Practical Examples in MOA” textbook is a resource intended to help students and practitioners enter the field of machine learning and data mining for data streams. **The online version of the book will remain available online for free.**

This textbook can also be ordered on Amazon.

<https://moa.cms.waikato.ac.nz/book-html/>



Albert Bifet
Ricard Gavaldà
Geoffrey Holmes
Bernhard Pfahringer

MACHINE LEARNING FOR DATA STREAMS

with Practical
Examples in MOA

The Learning Cycle

Batch vs. Stream

Batch data

The model is updated through several passes over the training data

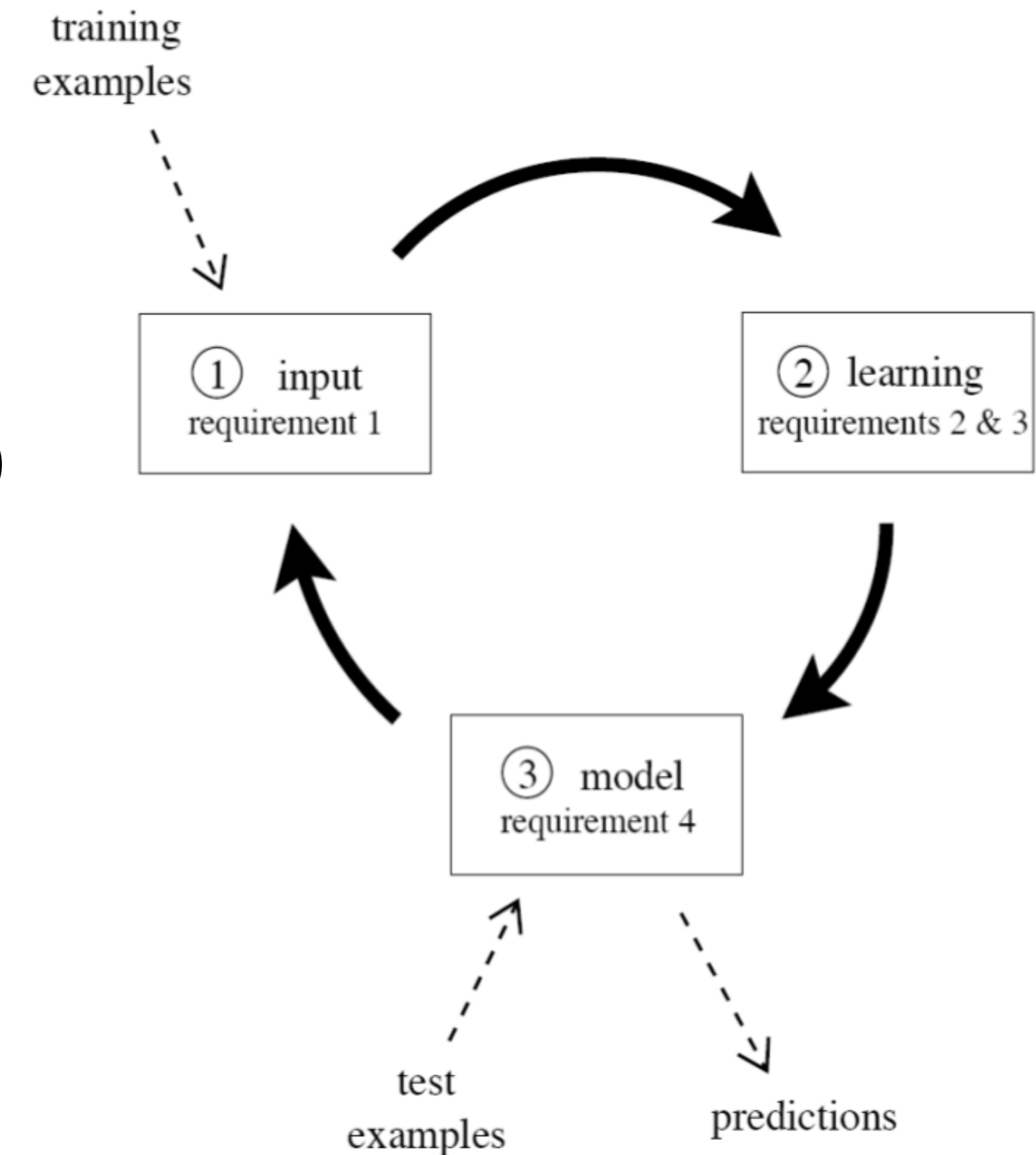
Streaming data

The model is updated after observing every new data point

The Learning Cycle

Requirements

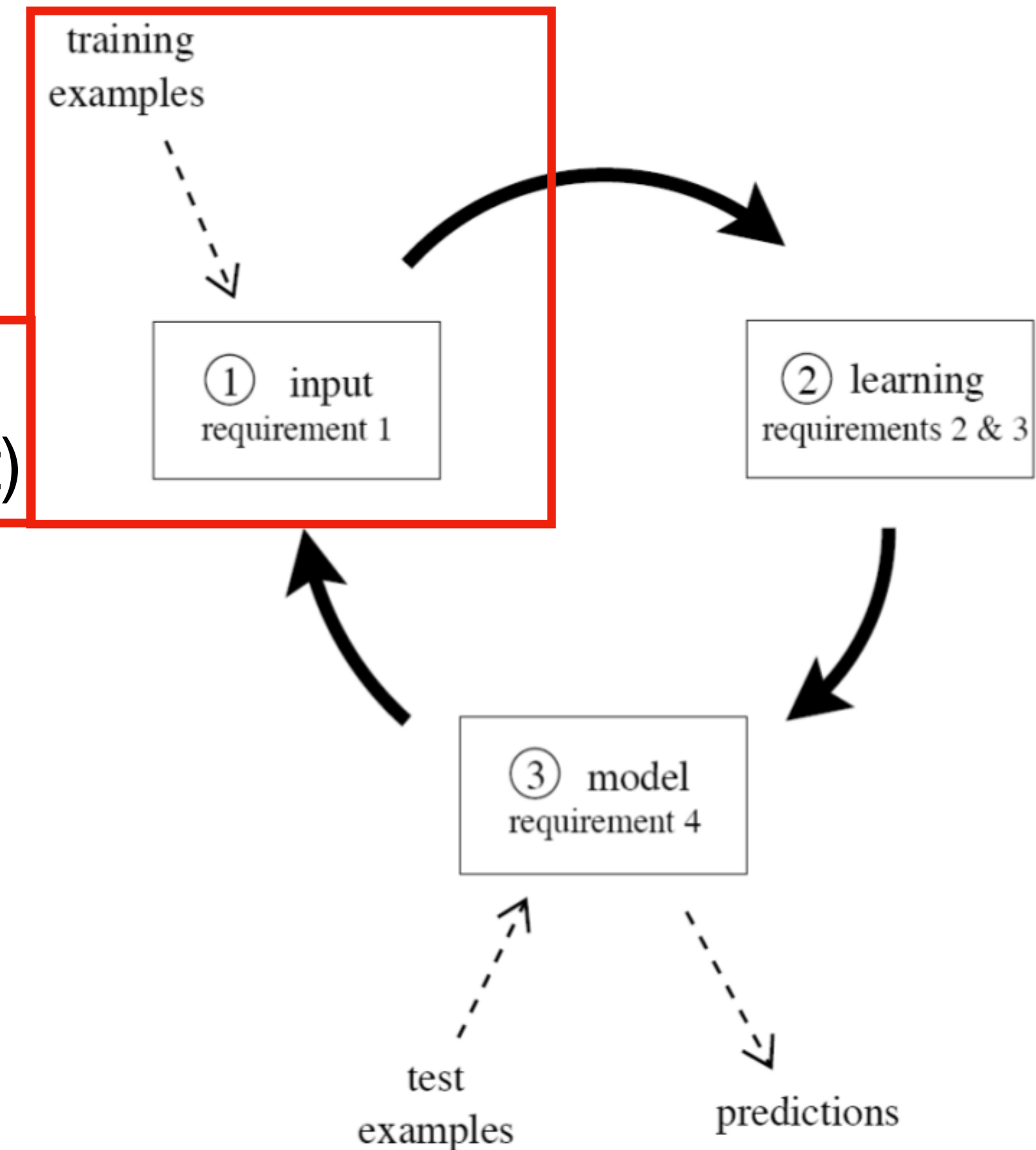
1. Process an example at a time, and **inspect it only once** (at most)
2. Use a **limited** amount of **memory**
3. Work in a **limited** amount of **time**
4. Be **ready to predict at any point**



The Learning Cycle

Requirements

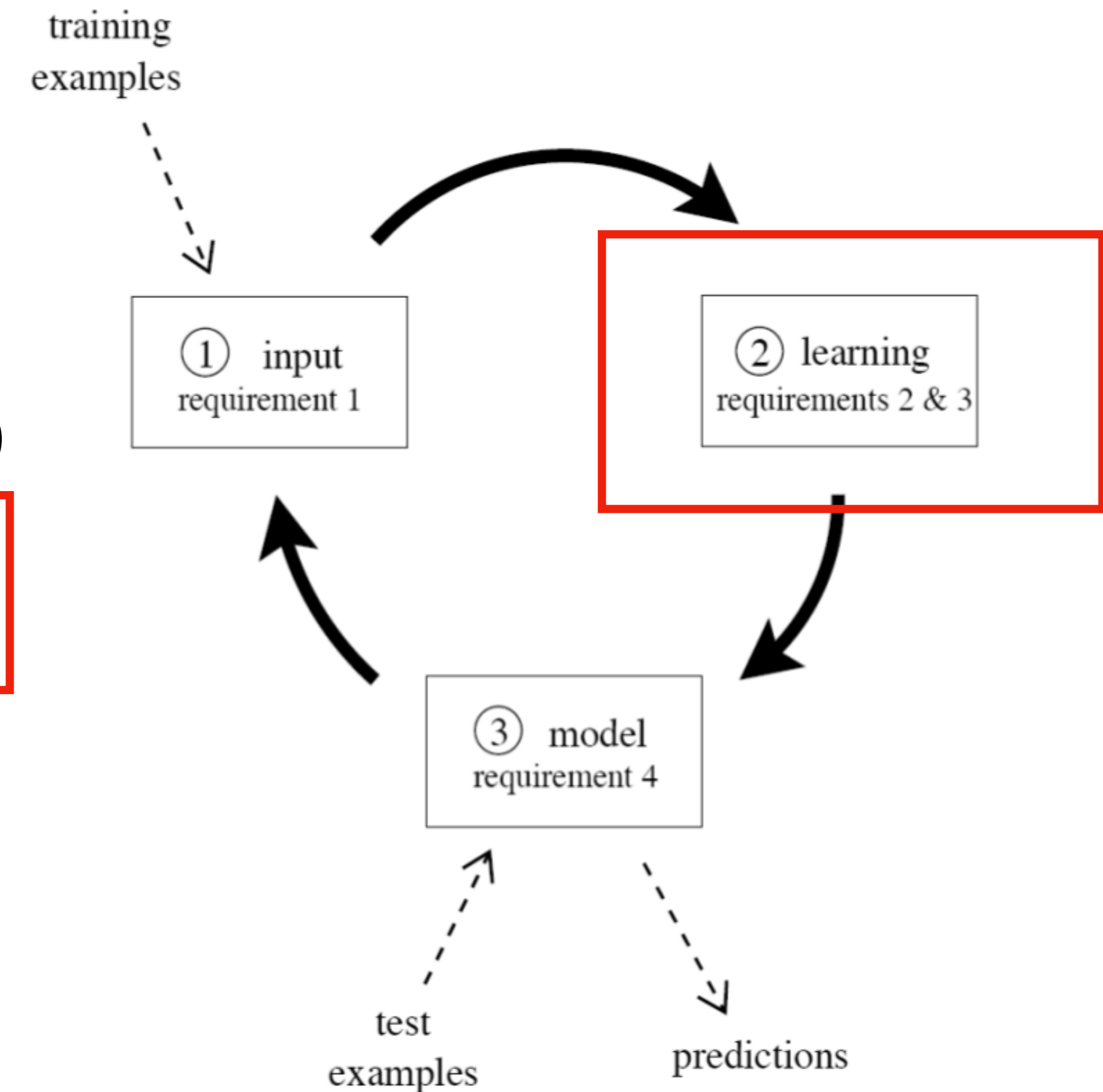
1. Process an example at a time, and **inspect it only once** (at most)
2. Use a **limited** amount of **memory**
3. Work in a **limited** amount of **time**
4. Be **ready to predict at any point**



The Learning Cycle

Requirements

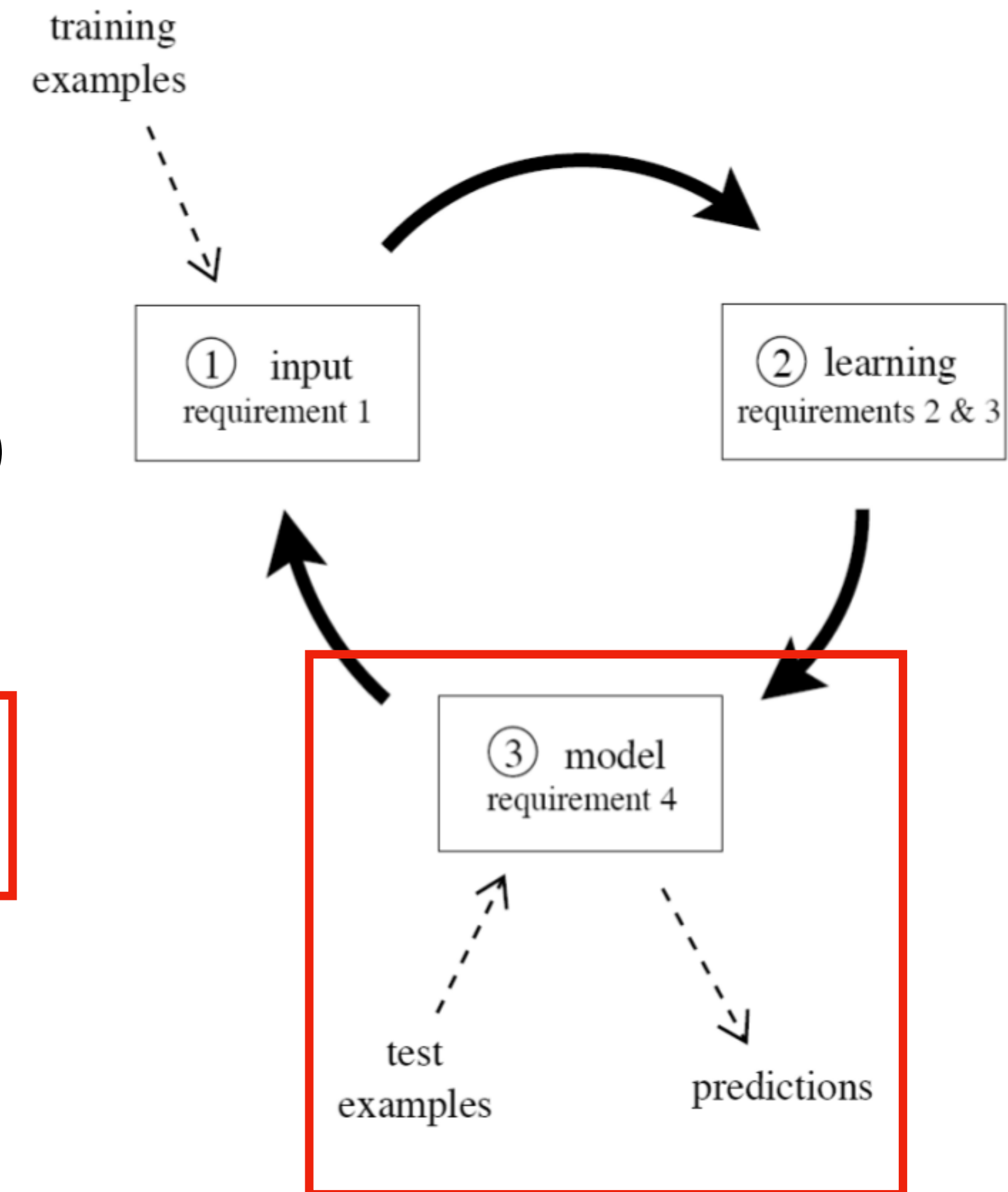
1. Process an example at a time, and **inspect it only once** (at most)
2. Use a **limited** amount of **memory**
3. Work in a **limited** amount of **time**
4. Be **ready to predict at any point**



The Learning Cycle

Requirements

1. Process an example at a time, and **inspect it only once** (at most)
2. Use a **limited** amount of **memory**
3. Work in a **limited** amount of **time**
4. Be **ready to predict at any point**



Evaluation

Evaluation overview

Aspects concerning predictive performance evaluation:

- **Evaluation metrics.** How errors are considered?
- **Evaluation framework.** How past predictions influence the current metric?

Other measurements (e.g. wall-clock time, CPU time, ...)

Evaluation Framework

(Periodic) Holdout

- Interleave test and train subsets
- After testing in one window (or chunk), we observe the average over that window. We use the following window for training.

Evaluation Framework

Interleaved test-then-train (or cumulative)

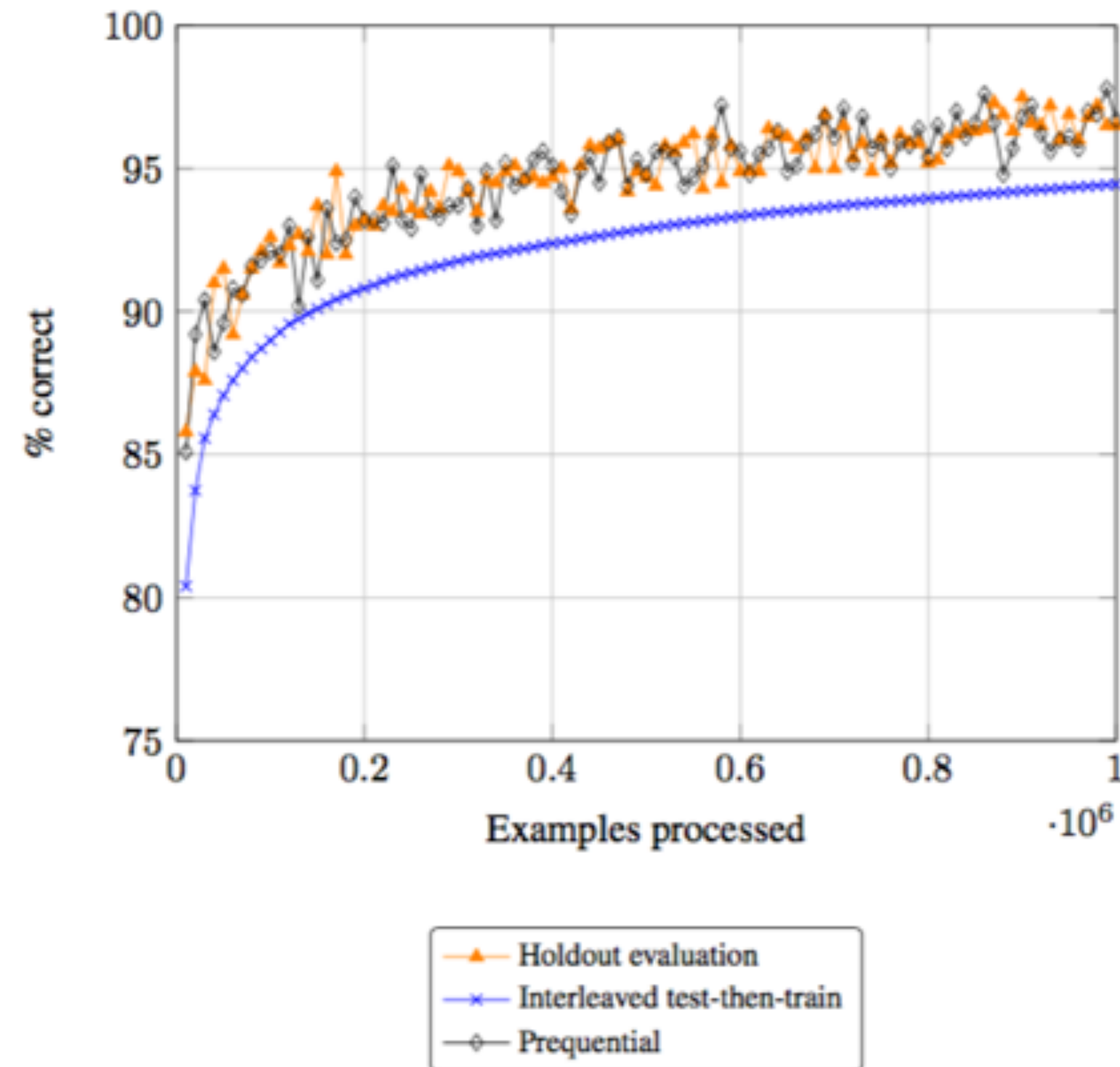
- Every instance is used for testing first, then for training
- At any point during execution, we observe the average over all instances seen so far

Evaluation Framework

Prequential evaluation

- Similar to interleaved test-then-train, but we observe the metrics over a window of the latest instances (optionally, we can use a fading factor)

Evaluation Framework

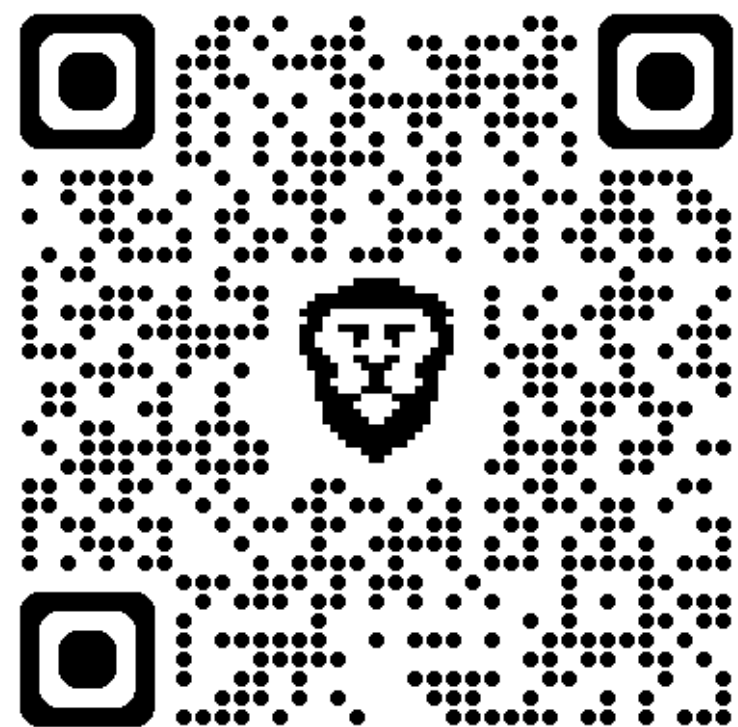


CapyMOA

Machine learning for data streams

<https://capymoa.org/>

<https://github.com/adaptive-machine-learning/CapyMOA>



v0.7.0



CapyMOA

A machine learning library for streaming data based on four pillars:

- **Efficiency**
- **Interoperability**
- **Accessibility**
- **Flexibility**

First released on May 03, 2024

Other frameworks: **MOA** (java)¹, **river** (python)² and **scikit-multiflow** (python)³

[1] Bifet, A., Holmes, G., Pfahringer, B., Kranen, P., Kremer, H., Jansen, T., & Seidl, T. (2010). Moa: Massive online analysis, a framework for stream classification and clustering. In *Workshop on applications of pattern analysis* (pp. 44-50). PMLR.

[2] Montiel, J., Halford, M., Mastelini, S.M., Bolmier, G., Sourty, R., Vaysse, R., Zouitine, A., Gomes, H.M., Read, J., Abdessalem, T. and Bifet, A., 2021. River: machine learning for streaming data in python. *Journal of Machine Learning Research*, 22(110), pp.1-8.

[3] Montiel, J., Read, J., Bifet, A., & Abdessalem, T. (2018). Scikit-multiflow: A multi-output streaming framework. *Journal of Machine Learning Research*, 19(72), 1-5.

Why another one?

Efficiency

A key aspect of stream learning are **efficient** implementations; they should learn from thousands of instances as quickly as possible (near real-time)

Interoperability

Use algorithms from MOA, scikit-learn and PyTorch through an unified streaming API

Accessibility

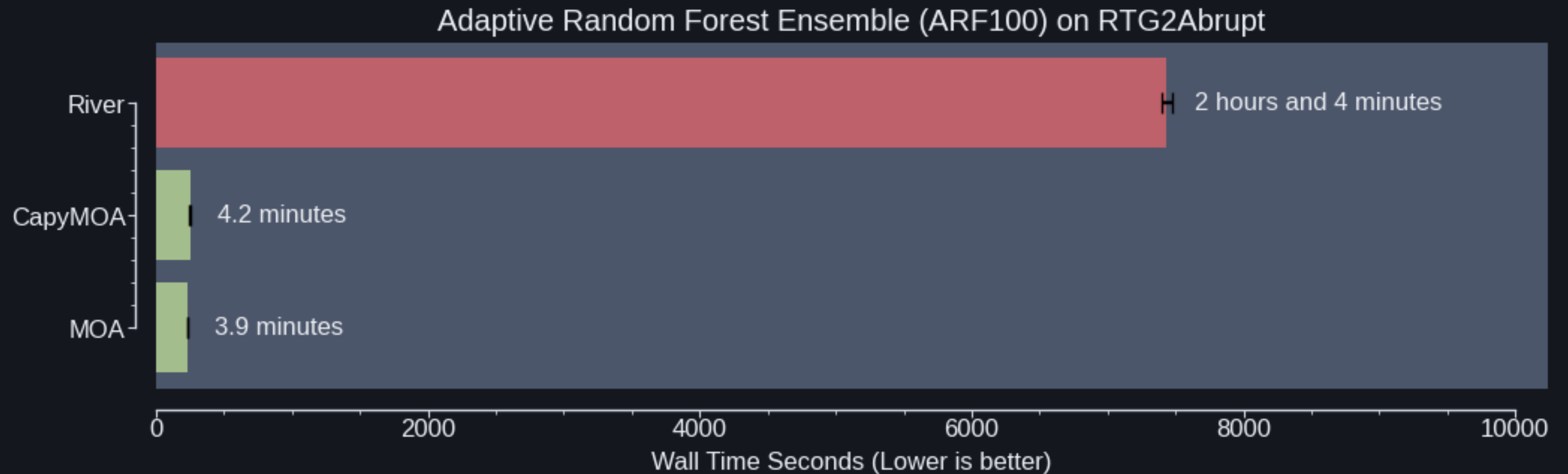
Must be easy to prototype experiments and extend functionality

Flexibility

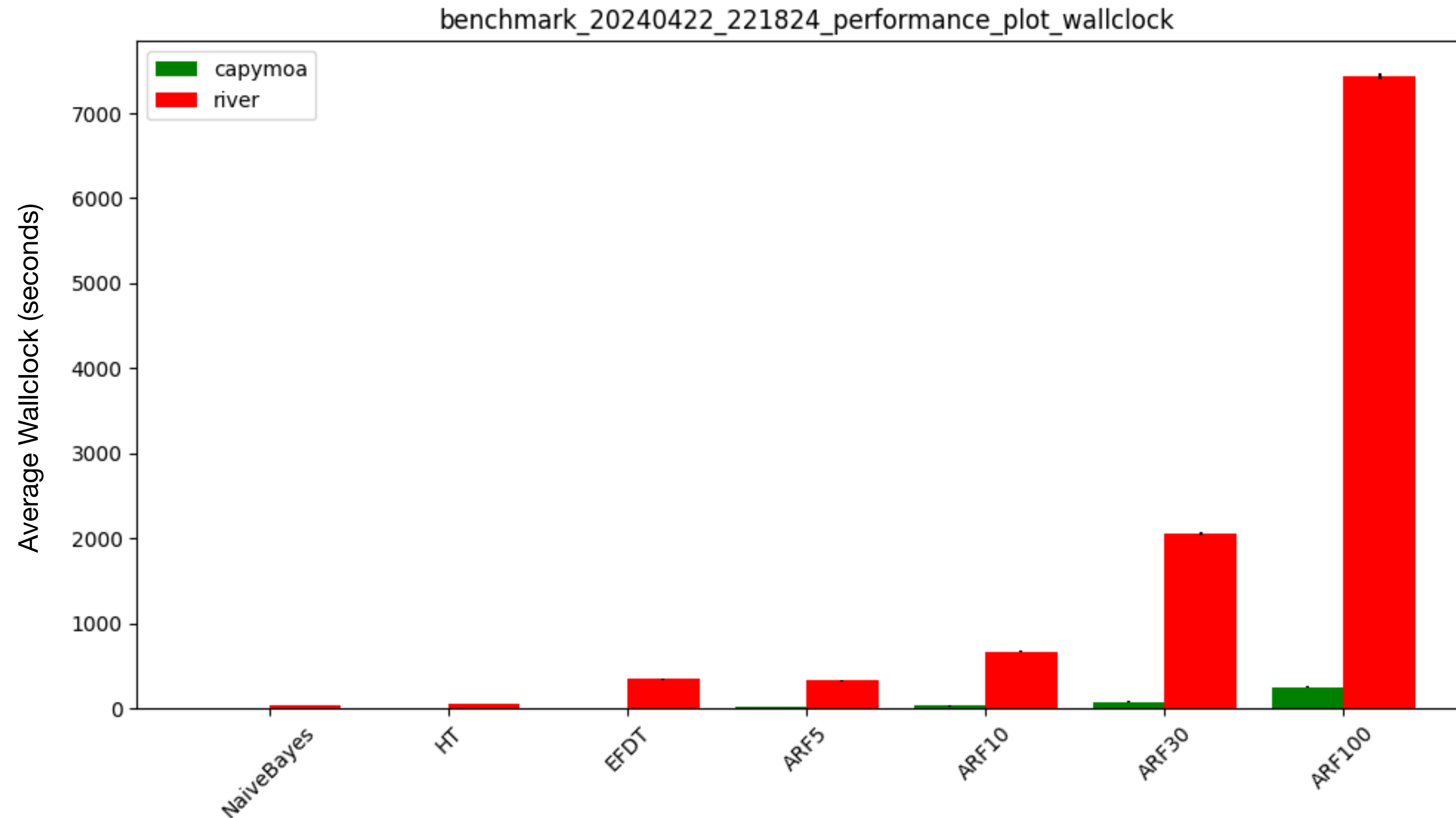
Advanced APIs for stream learning (concept drift, evaluation, visualisation, ...)

Code in either Python or Java*

Why? Efficiency



Why? Efficiency



Dataset: RandomTreeGenerator with 2 abrupt drifts, 100k instances, 30 attributes, and 5 classes.

Experiment: 5 repetitions

Algorithms: NaiveBayes, HoeffdingTree (HT), Extremely Fast Decision Tree (EFDT), and AdaptiveRandomForest (the suffix indicates the number of base learners)

Reproducibility: <https://github.com/adaptive-machine-learning/CapyMOA/blob/main/notebooks/benchmarking.py>

Why? Accessibility

Make it easier to configure and execute complex experiments

MOA's learning curve tends to be steeper as it is implemented in Java

CapyMOA allow users to code in Python, while allowing advanced users to take advantage from MOA objects directly from Python

```
from capymoa.stream.generator import SEA
from capymoa.stream.drift import DriftStream, AbruptDrift,
GradualDrift
from capymoa.classifier import AdaptiveRandomForestClassifier
from capymoa.evaluation import prequential_evaluation
from capymoa.evaluation.visualization import plot_windowed_results

SEA3drifts = DriftStream(stream=[SEA(1),
                                AbruptDrift(10000),
                                SEA(2),
                                GradualDrift(start=20000,
                                              end=25000),
                                SEA(3),
                                AbruptDrift(45000),
                                SEA(1)])

arf =
AdaptiveRandomForestClassifier(schema=SEA3drifts.get_schema(),
                               ensemble_size=100,
                               number_of_jobs=4)

results = prequential_evaluation(stream=SEA3drifts,
                                learner=arf,
                                window_size=1000,
                                max_instances=50000)

print(f"Cumulative accuracy = {results['cumulative'].accuracy()}")
print(f"wallclock = {results['wallclock']} seconds")
display(results['windowed'].metrics_per_window())
plot_windowed_results(results, ylabel='Accuracy')
```


Why? Interoperability

Easy access to **MOA**, **PyTorch** and **Scikit-learn** learners

Wraps and expose MOA's API through a modern and simple API

Standardise evaluation and benchmarking without losing **flexibility**



Why? Flexibility

CapyMOA facilitates preparing and executing experiments by incorporating:

- High-level evaluation functions
- Build-in visualisation tools
- Advanced Concept Drift API; and more
- We will show several examples throughout this tutorial

CapyMOA team

- Heitor Murilo Gomes (project leader)¹
- Anton Lee¹
- Nuwan Gunasekara²
- Yibin Sun²
- Guilherme Cassales²
- Marco Heyden³
- Justin Liu²
- Jesse Read⁴
- Maroua Bahri⁵
- Marcus Botacin⁶
- Vitor Cerqueira⁷
- Albert Bifet^{2,9}
- Bernhard Pfahringer²
- Yun Sing Koh⁸

And many other individual contributors

[1] Victoria University of Wellington, New Zealand

[2] University of Waikato, New Zealand

[3] KIT, Germany

[4] École polytechnique, IP Paris, France

[5] INRIA Paris, France

[6] Texas A&M Engineering, USA

[7] Porto University, Portugal

[8] University of Auckland, New Zealand

[9] Télécom Paris, IP Paris, France



CapyMOA summary

- Allows access to existing and future MOA implementations
- Code in Python or Java, or both
- Integration with PyTorch and scikit-learn
- Streams, learners and evaluation are designed to interoperate with visualisation
- Latest release (0.7.0): August 03, 2024
- 20 classifiers, 8 regressors, 11 drift detectors, 3 anomaly detectors, evaluation, data representation, ... as of 0.7.0



Hands-on example

KDD_2024_introduction.ipynb