

# Introduction to Learning from Streaming Data

KEEPER Workshop Tutorial 2024

Nuwan Gunasekara<sup>\*1</sup>, Sepideh Pashami<sup>1</sup>,  
<https://nuwangunasekara.github.io/KEEPER2025/>

\* Corresponding author: [heitor.gomes@vuw.ac.nz](mailto:heitor.gomes@vuw.ac.nz)



# Our goals

- Introduce attendees to several machine-learning tasks for streaming data, such as:

Classification, regression, concept drifts, anomaly detection

- Enable attendees to **apply** and **extend** the concepts demonstrated using Python and ***capymoa***

# Outline

- **Machine Learning for Streaming Data (intro)**
  - Learning cycle
  - Evaluation
  - CapyMOA
  - [01\\_KEEPER2025\\_introduction.ipynb](#)
- **Supervised Learning**
  - Classification
  - Regression
  - [02\\_KEEPER2025\\_supervised.ipynb](#)
- **Unsupervised Learning**
  - Anomaly detection
  - [03\\_KEEPER2025\\_anomaly\\_detection.ipynb](#)



**Notebooks:**

<https://nuwangunasekara.github.io/KEEPER2025/>

# **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

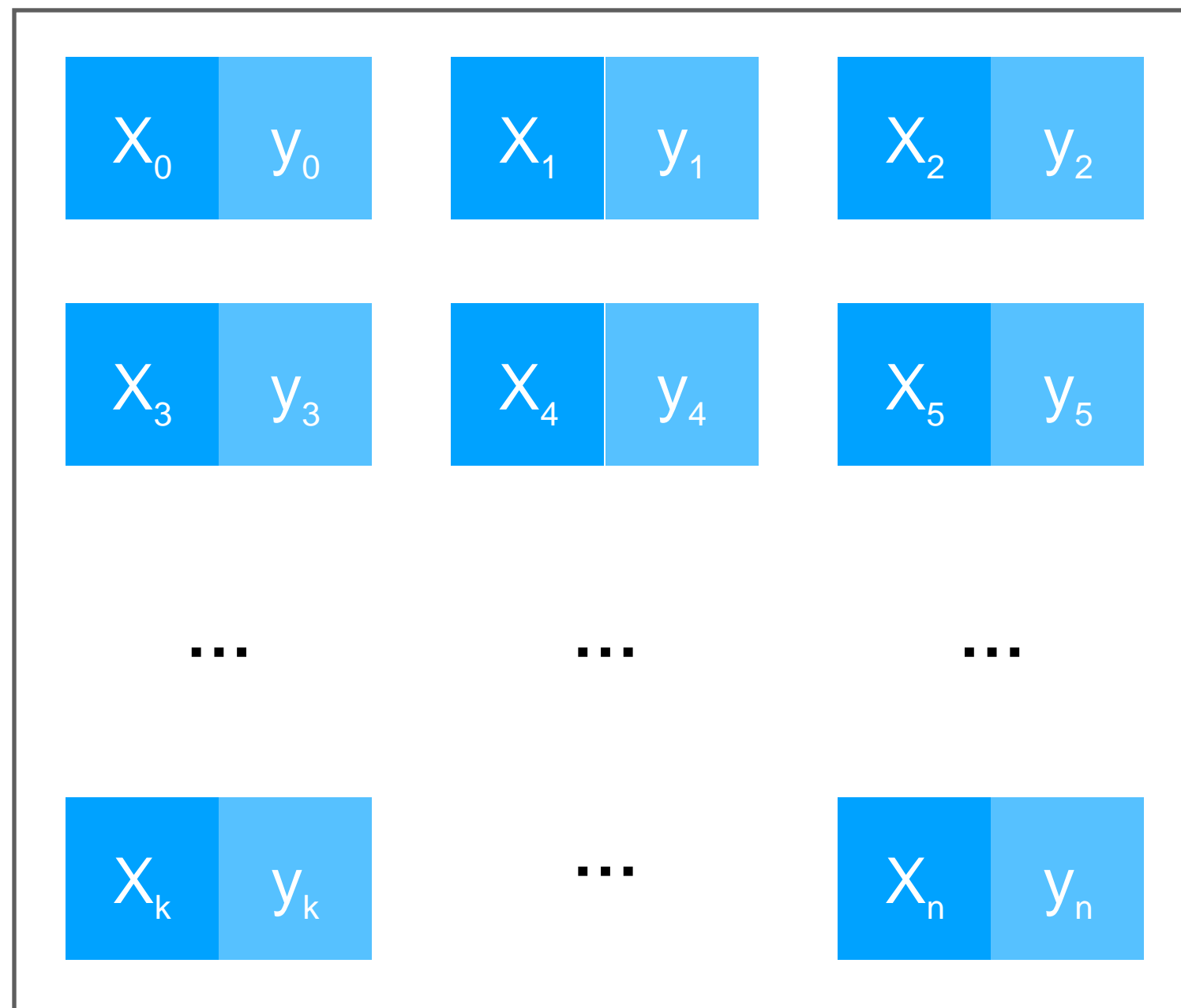
# Stream Learning

ML for **Batch** (“static”) data

**vs.**

ML for **Streaming** (“online”) data

# ML for Batch data



Fixed size dataset

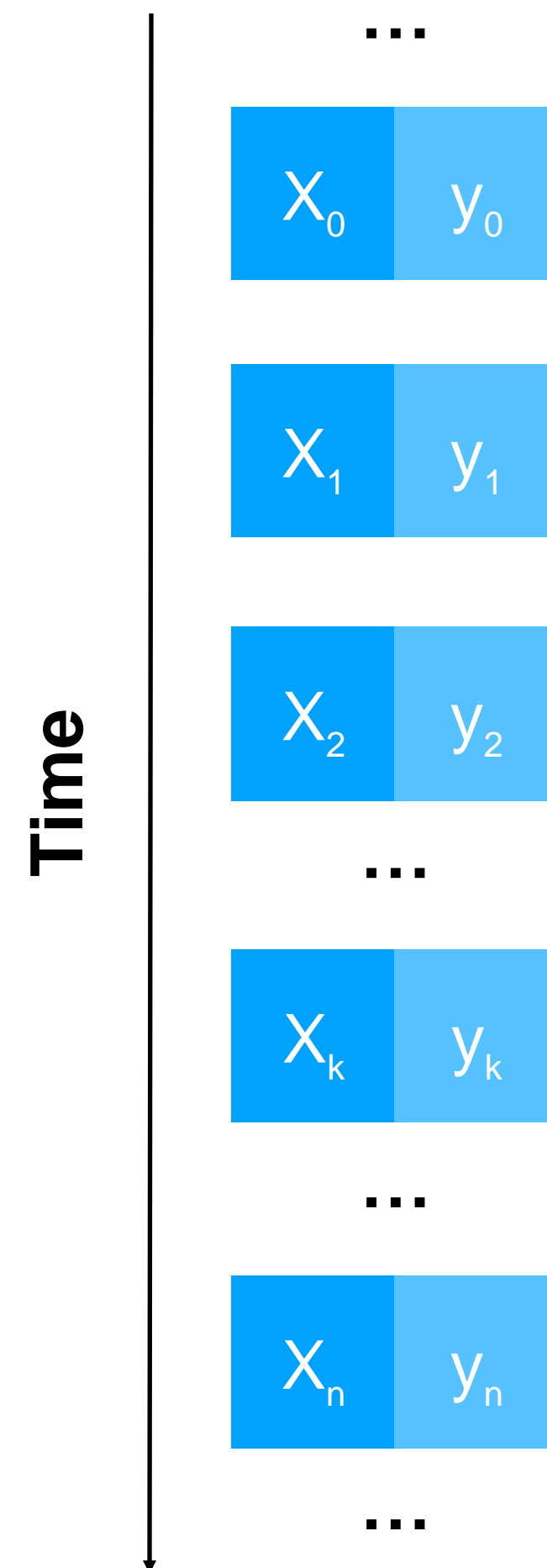
Random access to any instance

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

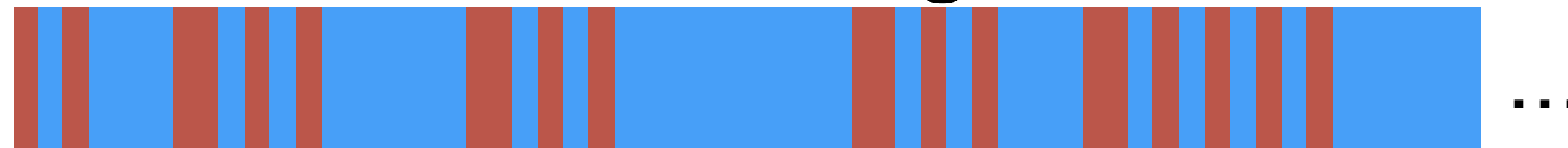
# Batch vs. Streaming

## Batch data



The output is a **trained model**

## Streaming data



The output is a trainable model

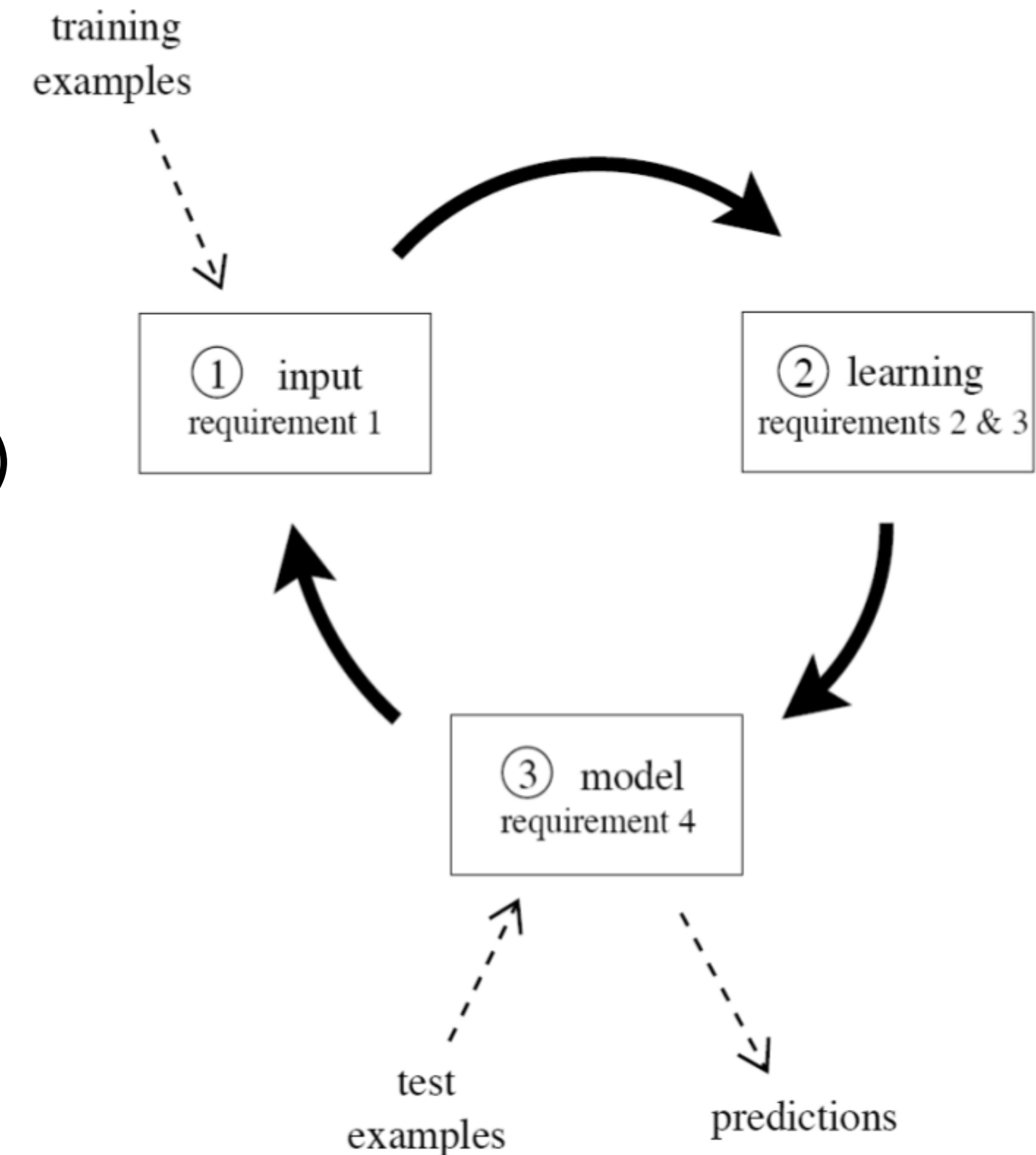


# The Learning Cycle

# The Learning Cycle

## Requirements

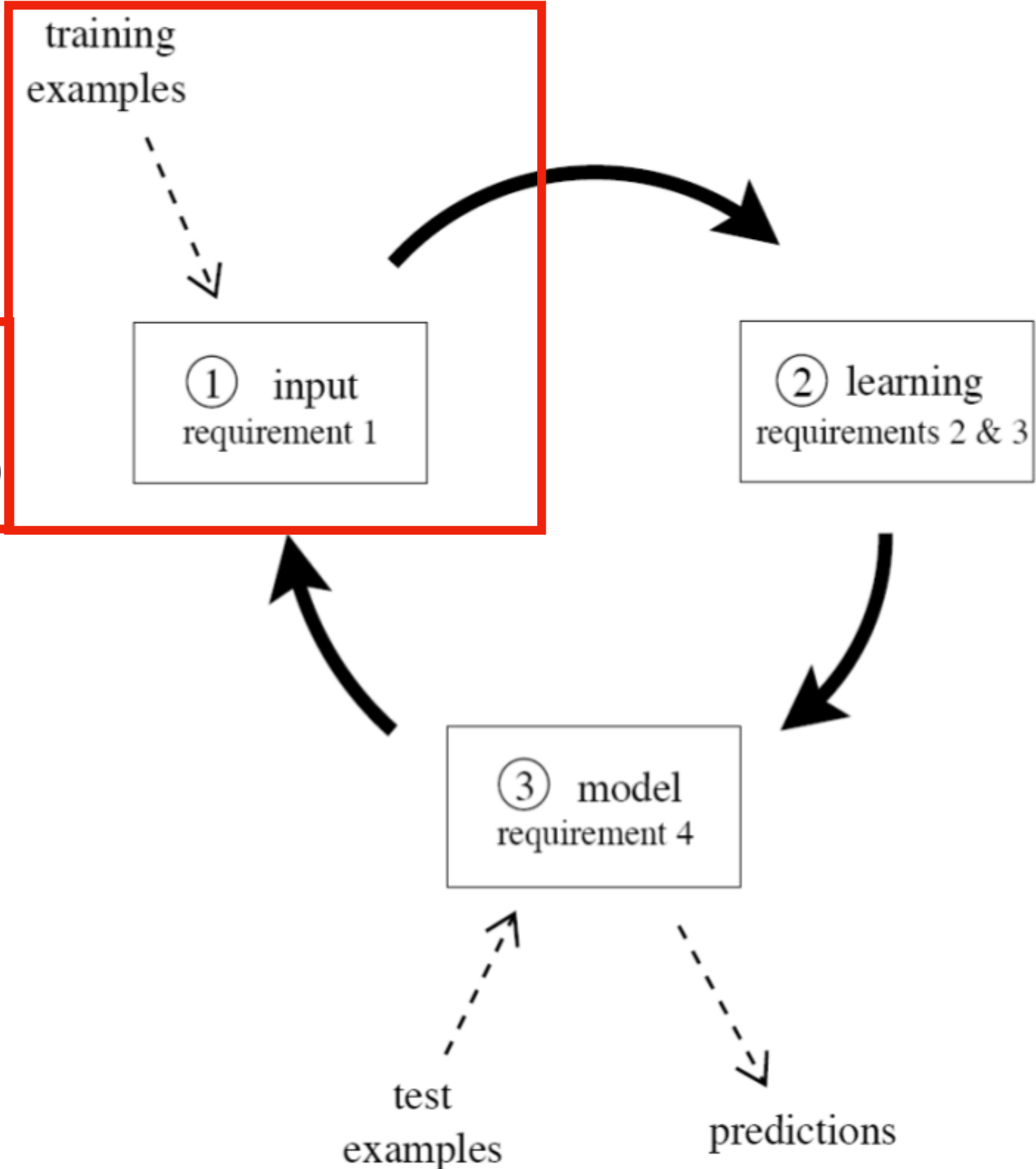
1. Process an example at a time, and **inspect it only once** (at most)
2. Use a **limited** amount of **memory**
3. Work in a **limited** amount of **time**
4. Be **ready to predict at any point**



# The Learning Cycle

# Requirements

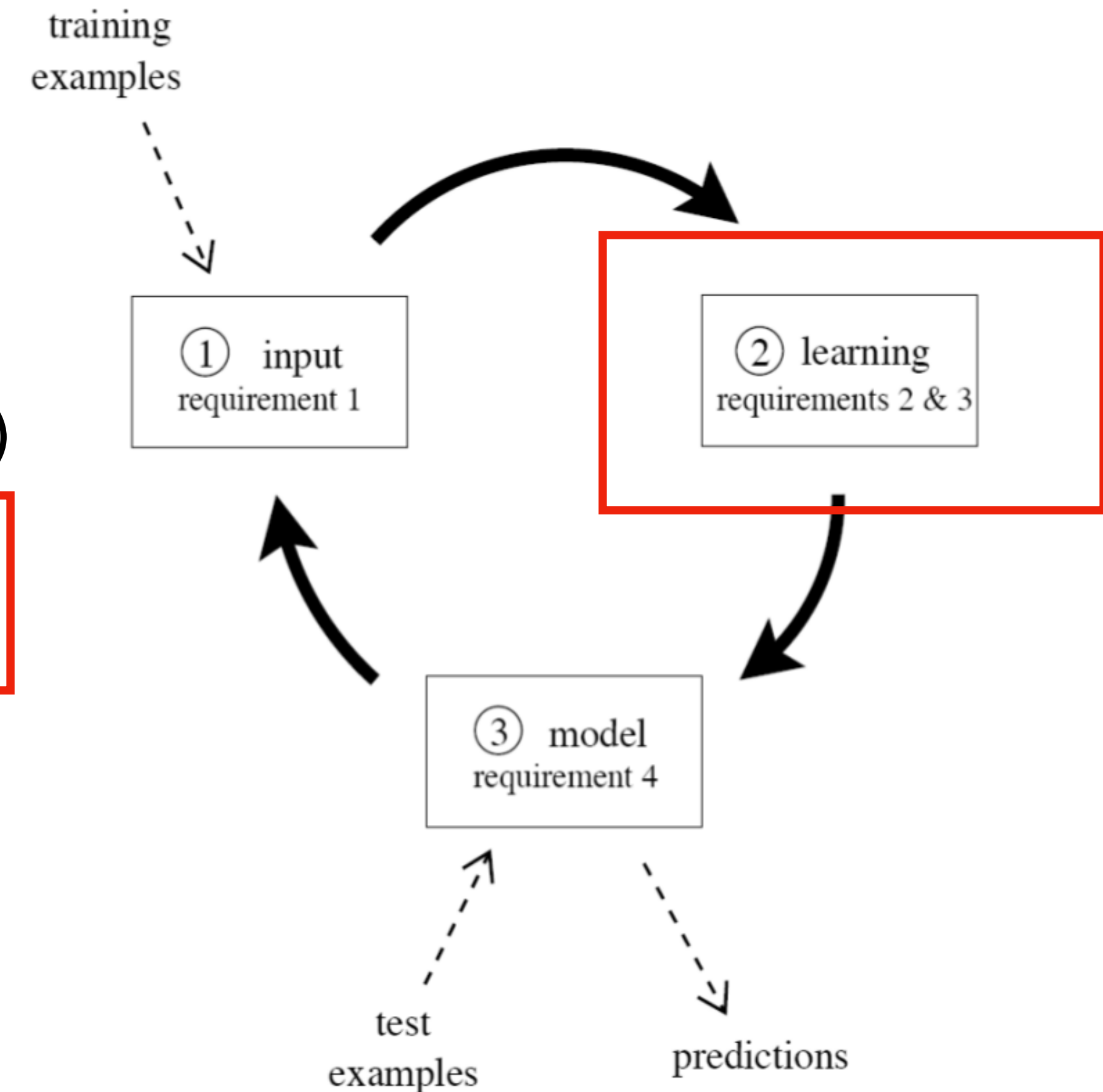
1. Process an example at a time, and **inspect it only once** (at most)
2. Use a **limited** amount of **memory**
3. Work in a **limited** amount of **time**
4. Be **ready to predict at any point**



# The Learning Cycle

## Requirements

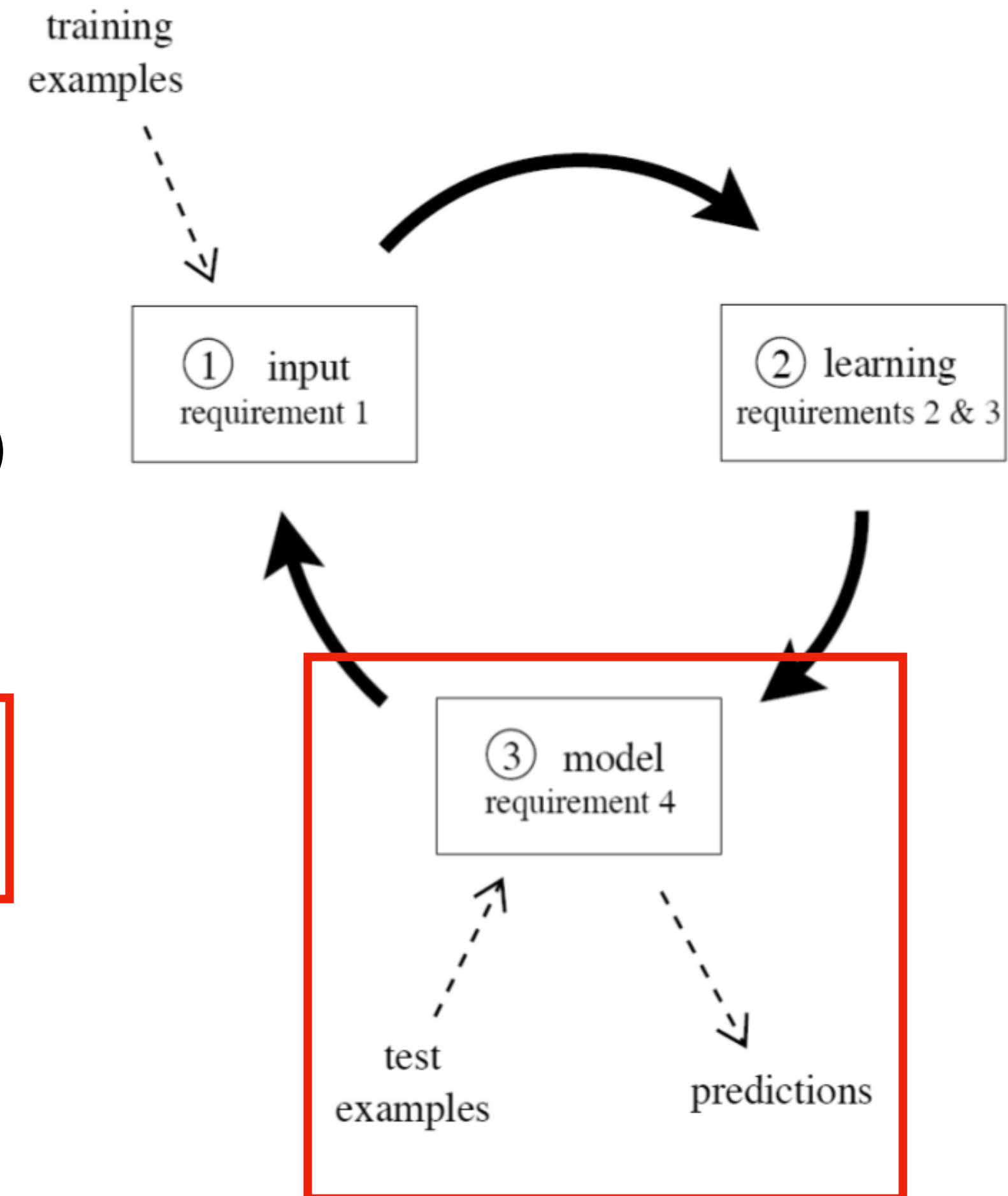
1. Process an example at a time, and **inspect it only once** (at most)
2. Use a **limited** amount of **memory**
3. Work in a **limited** amount of **time**
4. Be **ready to predict at any point**



# 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

**Cumulative (test-then-train):** At any point during execution, we observe the average over all instances seen so far

**Windowed (prequential):** Similar to cumulative, but we observe the metrics over a window of the latest instances



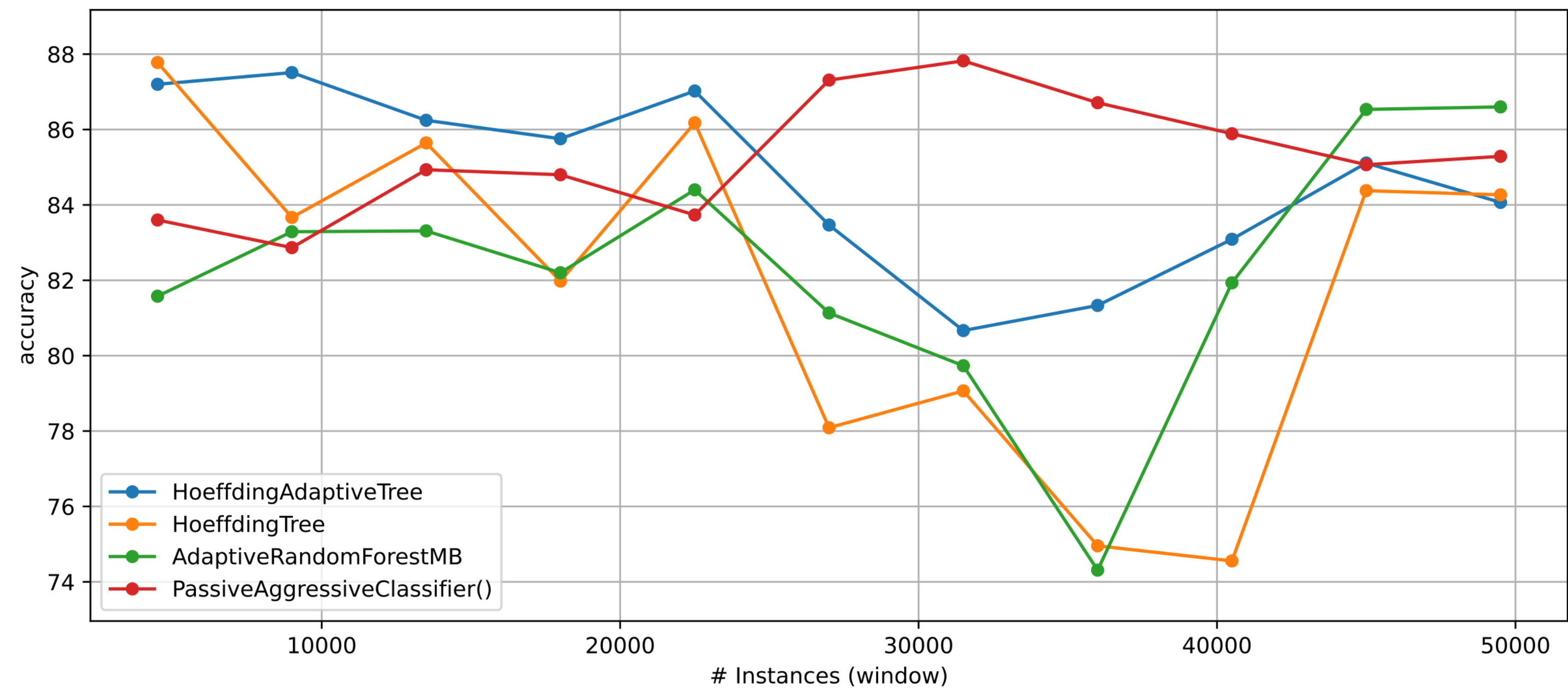
# Evaluation Framework (example)

In *capymoa* `prequential_evaluation(...)` will return both results

## Cumulative

Algorithm	Accuracy (cumulative)
HoeffdingAdapt.	84.6861
HoeffdingTree	81.6604
AdaptiveRand.	81.9076
PassiveAggr.	85.2445

## Windowed



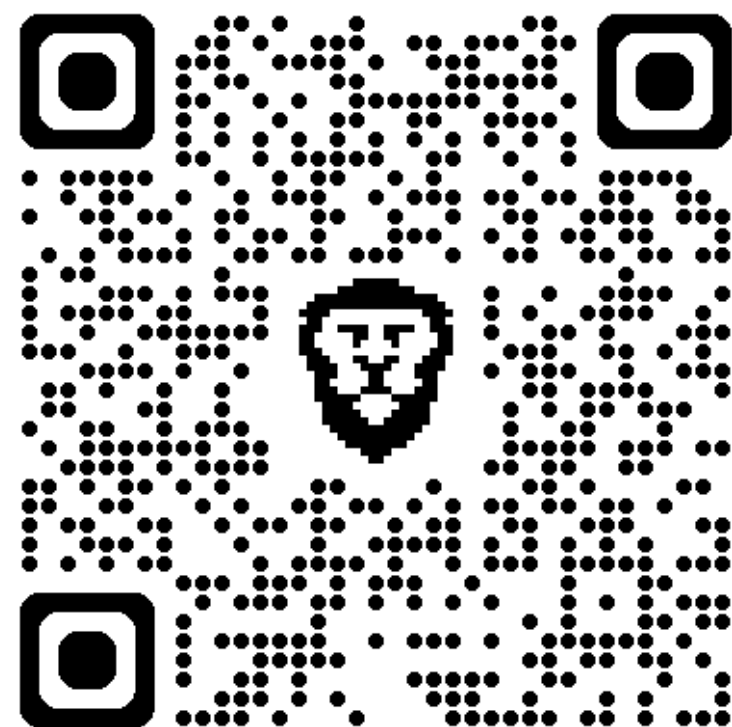


# CapyMOA

Machine learning  
for data streams

<https://capymoa.org/>

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





# CapyMOA

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

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

***capymoa*** is open-source and it was first publicly available on May 03, 2024

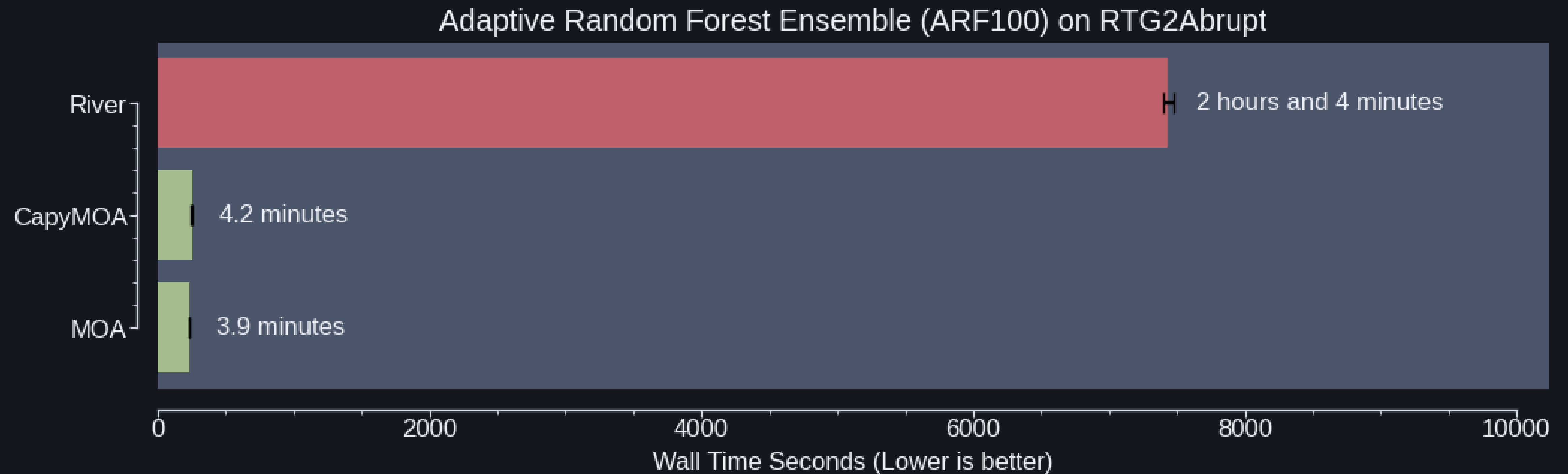
Other frameworks: **MOA** (java)<sup>1</sup>, **river** (python)<sup>2</sup> and **scikit-multiflow** (python)<sup>3</sup>

[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? Efficiency



# Why? Accessibility

Easy to configure and execute complex experiments

Code in Python, but take advantage of MOA (Java) objects

Allows access to existing and future MOA implementations

Integrate stream simulation with evaluation and visualisation

# Why? Accessibility

Simulate a data stream with 3 concepts drifts

```
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? Accessibility

Configure an ensemble with 100 learners and 4 jobs (multithreaded)

```
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? Accessibility

Calculate cumulative and windowed metrics

```
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? Accessibility

```
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, metric='accuracy')
```

Plot the windowed results



<https://capymoa.org/>

# CapyMOA team

- Heitor Murilo Gomes (project leader)<sup>1</sup>
- Anton Lee<sup>1</sup>
- Nuwan Gunasekara<sup>2</sup>
- Yibin Sun<sup>2</sup>
- Guilherme Cassales<sup>2</sup>
- Marco Heyden<sup>3</sup>
- Justin Liu<sup>2</sup>
- Jesse Read<sup>4</sup>
- Maroua Bahri<sup>5</sup>
- Marcus Botacin<sup>6</sup>
- Vitor Cerqueira<sup>7</sup>
- Albert Bifet<sup>2,9</sup>
- Bernhard Pfahringer<sup>2</sup>
- Yun Sing Koh<sup>8</sup>

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

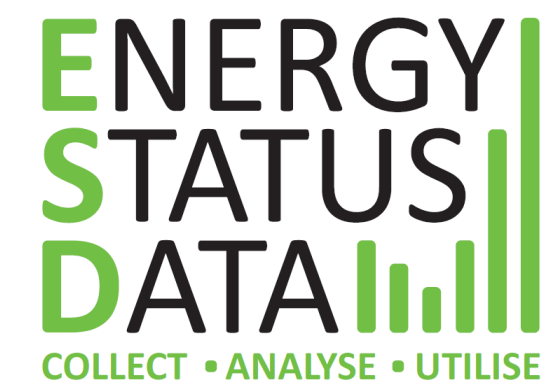
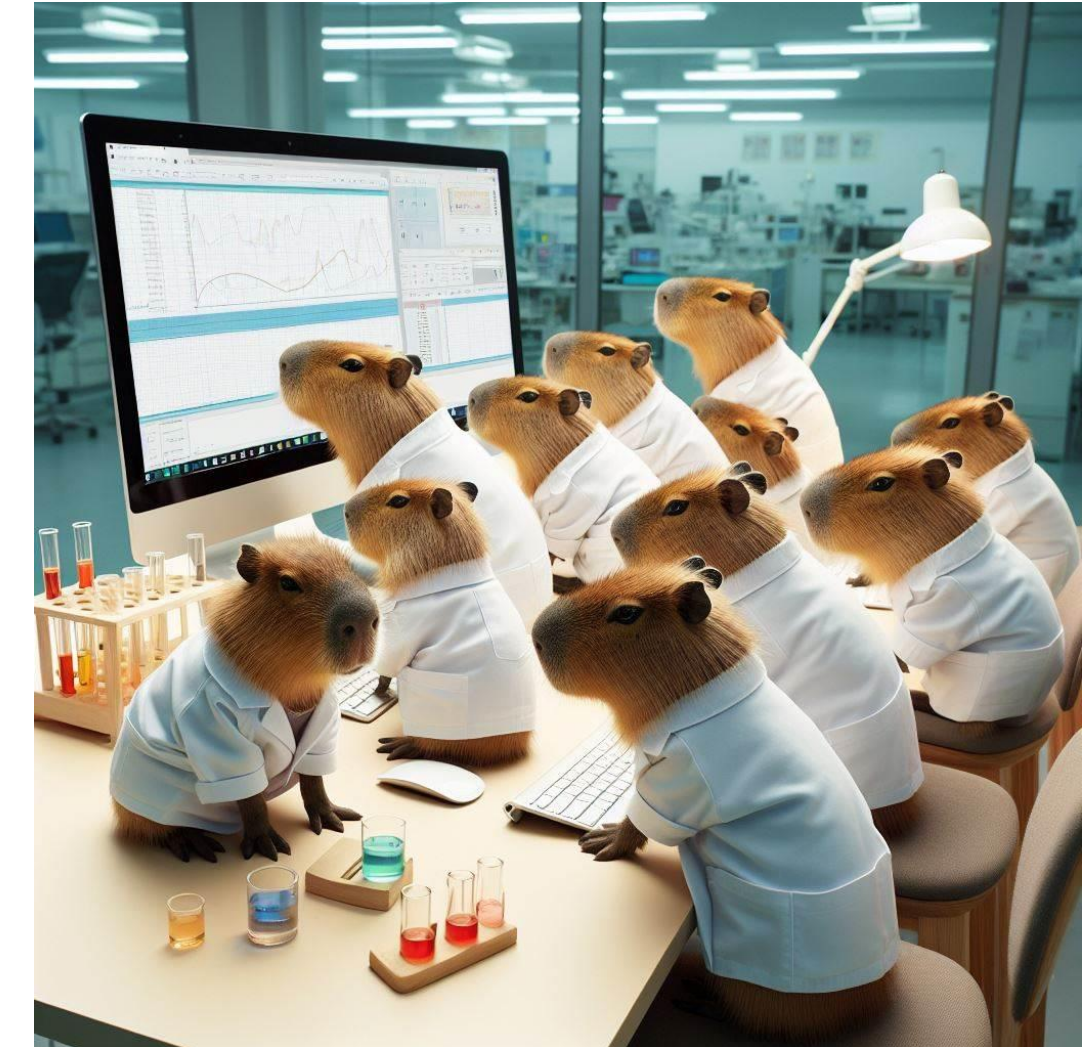
[5] INRIA Paris, France

[6] Texas A&M Engineering, USA

[7] Porto University, Portugal

[8] University of Auckland, New Zealand

[9] Télécom Paris, IP Paris, France





# CapyMOA summary

- Code in Python or Java, or both
- Integration with PyTorch and scikit-learn
- Streams, learners and evaluation are designed to interoperate with visualization
- Latest release (0.9.0): March, 2025
- 20 classifiers, 8 regressors, 11 drift detectors, 3 anomaly detectors, evaluation, data representation, ... as of 0.8.0



# Practical examples



**01\_KEEPER2025\_introduction.ipynb**

**Notebooks:**

<https://nuwangunasekara.github.io/KEEPER2025/>