

Practical Challenges and Insights in Concept Drift Research

Discovering Drift Phenomena in Evolving Landscape (DELTA 2024)
ACM SIGKDD 2024 workshop

Heitor Murilo Gomes
<https://heitorgomes.com/>

Victoria University of Wellington, New Zealand



<https://capymoia.org/>



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.

Research interests: ML for data streams, ensemble learning, semi-supervised learning, concept drift detection/adaptation, ...

PI of a few research projects ranging from applied to fundamental research, i.e. ML for energy prediction, novel methodologies for stream learning, ...

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

<https://heitorgomes.com/>



About this talk

Brief introduction

Practical challenges

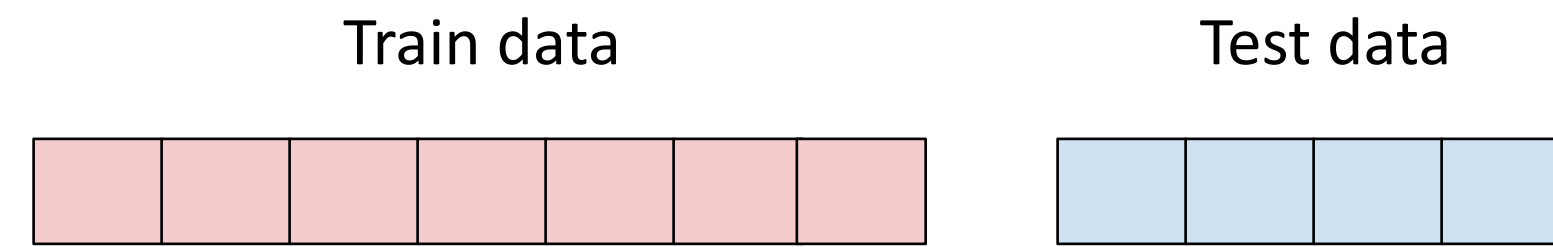
- Simulation
- Evaluation
- Using drift detectors

Other challenges

- Recurrent concept drifts
- Delayed & Unsupervised drift detection

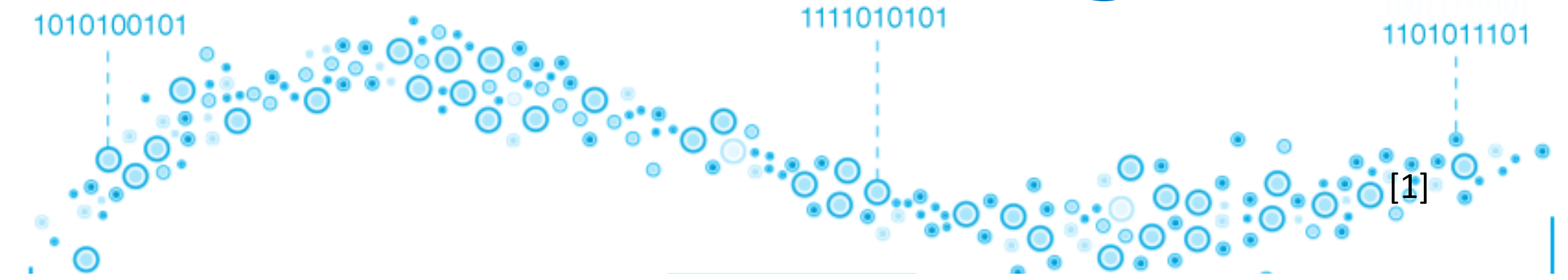
Brief Introduction

Batch Learning



vs

Stream Learning (SL)



Learns from **static** data

Uses a **large** amount of **computing**

Can only **predict** after (**extensive**) **training**.

Should **re-train** after a **concept drift**

Learns from a **stream of data**

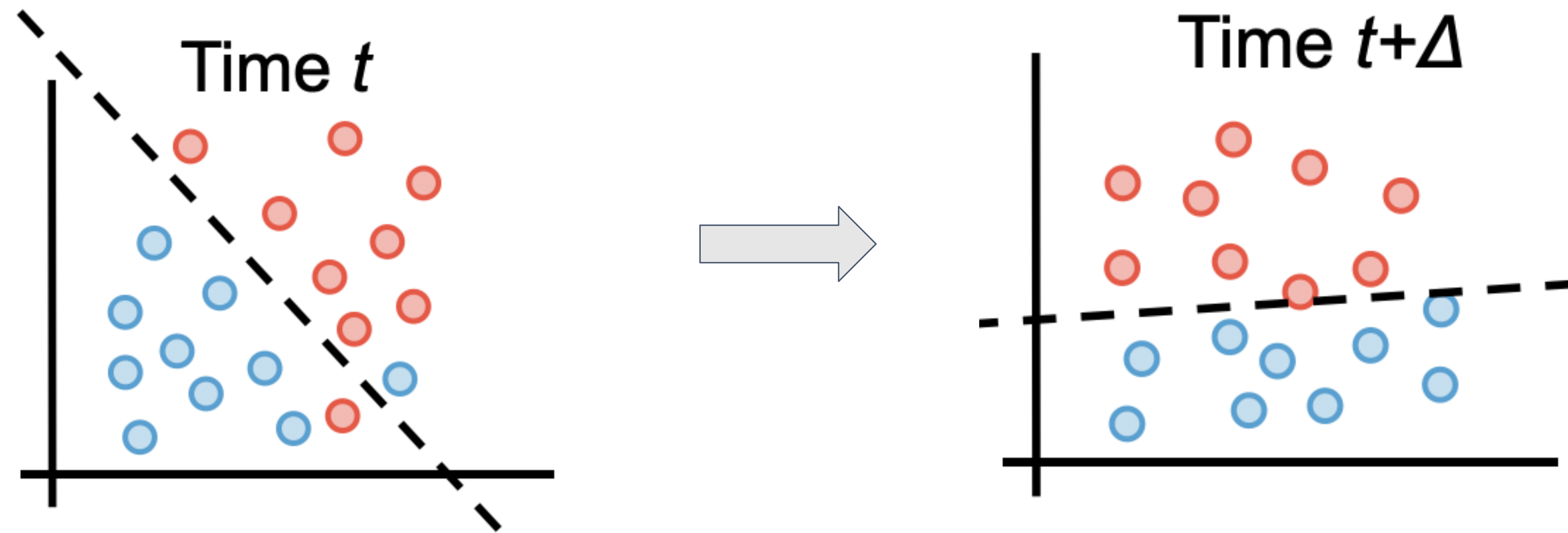
Incrementally online learn form instance/mini-batch.

Should use **limited computing** resources.

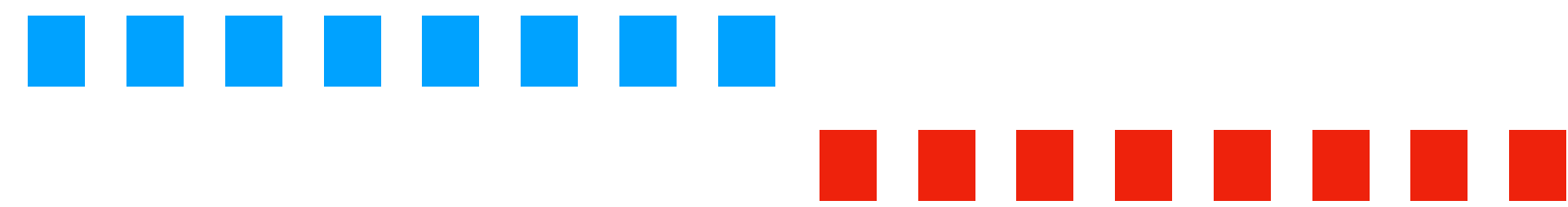
Able to **predict** at **any given moment**.

Should **adapt** to **concept drifts online**.

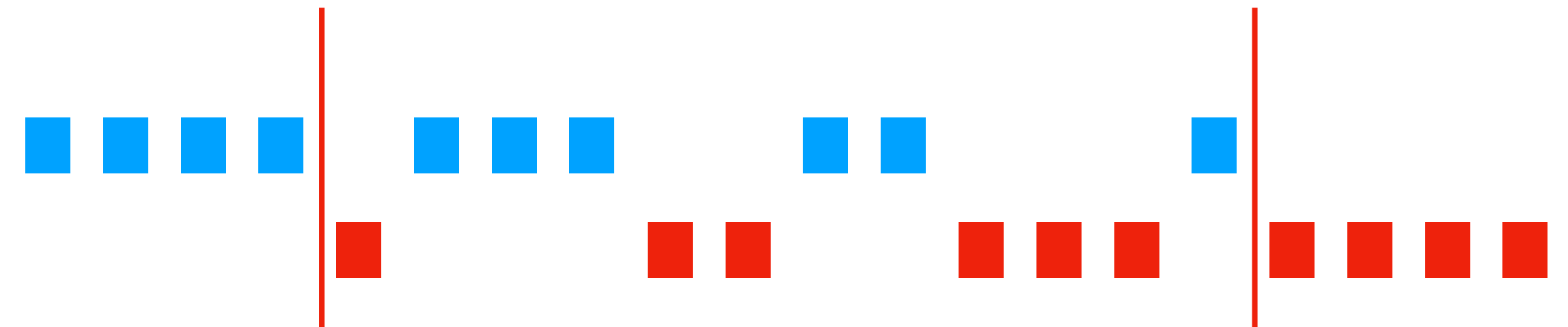
Concept Drift



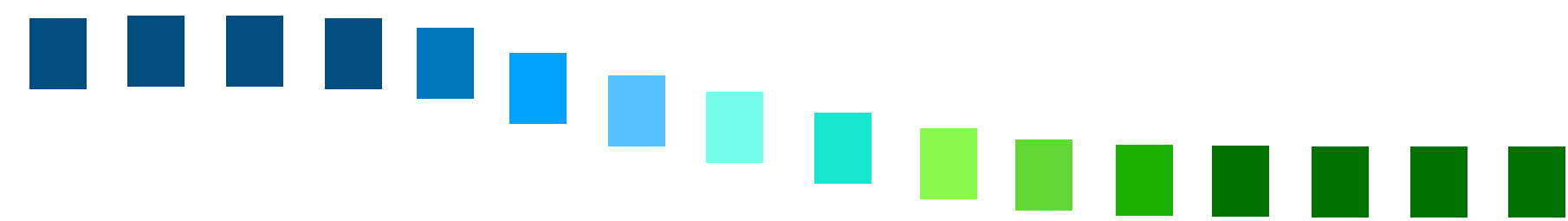
Concept Drift (categorisation)



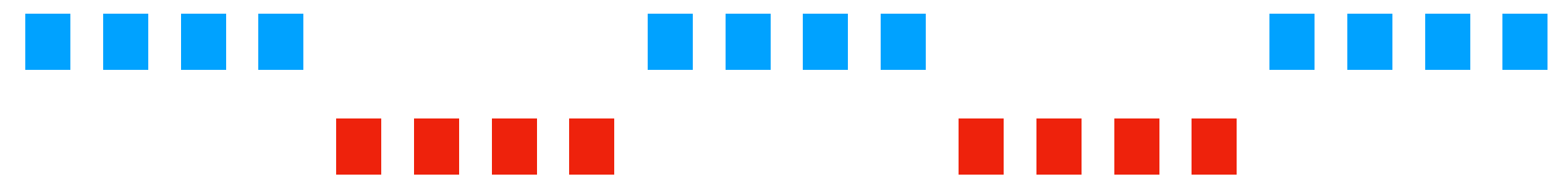
abrupt drift



gradual drift



incremental drift



recurrent concept drift

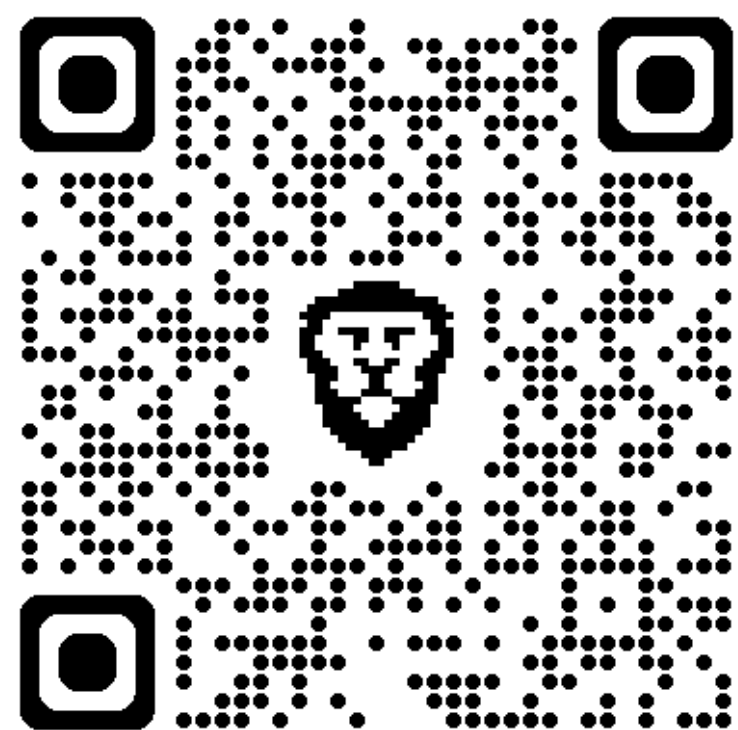
Practical challenges

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

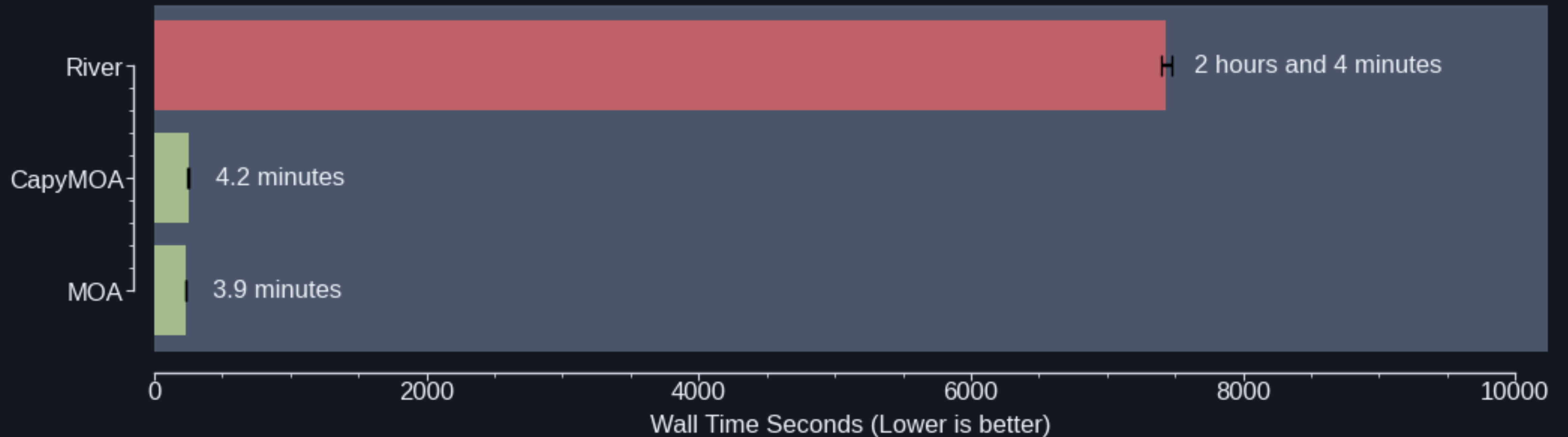
[1] Bifet, A., Holmes, G., Pfahringer, B., Kranen, P., Kremer, H., Jansen, T., & Seidl, T. (2010). Moa: Massive online analysis, a framework for stream classification and clustering. In *Workshop on applications of pattern analysis* (pp. 44-50). PMLR.

[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? Efficiency

Adaptive Random Forest Ensemble (ARF100) on RTG2Abrupt



Simulating Concept Drifts

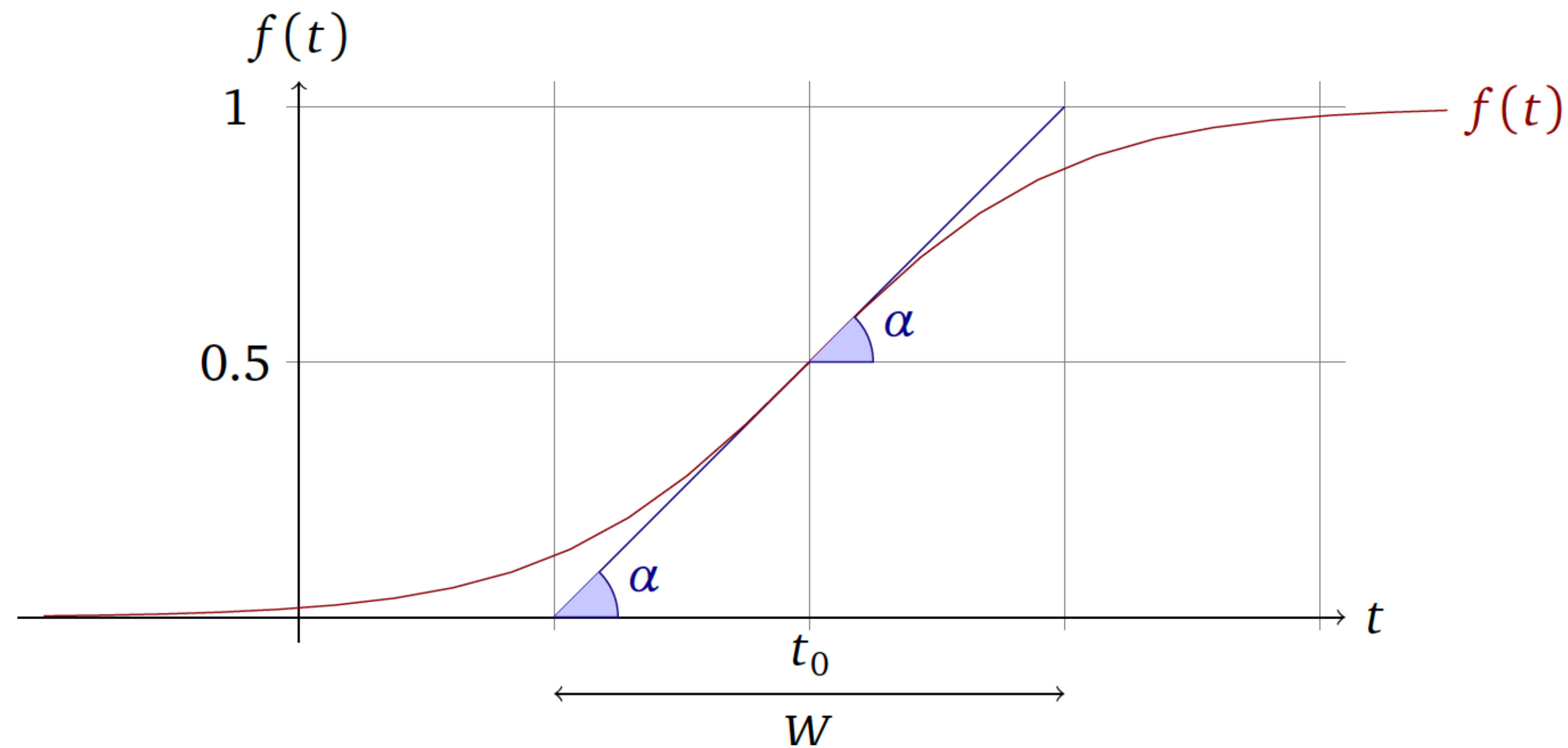
Concept drift is hard to define in a real data stream

Thus, studying it using real data can be challenging

One approach is to use synthetic data for studying and benchmarking algorithms

Concept Drift Framework

“Model a concept drift event as a **weighted combination of two pure distribution** that characterizes the target concepts before and after the drift.” [Bifet et al, 2011]



Recursive definition

- Most tools (MOA[1], river[2], scikit-multiflow[3], ...) uses a recursive approach to specify concept drift locations like:

CDS(CDS(SEA(1), SEA(2), 1000), SEA(3), 2000)

- Where we specify the drift **position** and the **width** of a drift (if it is a Gradual Drift) recursively.
- This can lead to some confusion depending on where the recursion is placed

[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.

Explicit list

In CapyMOA [1], the concepts and drifts are clearly outlined on a list format through a **DriftStream**

Drift position + drift width:

- The **start** and **end** of a **concept** is determined by the presence of an **AbruptDrift** or **GradualDrift** object

```
DriftStream([SEA(1), AbruptDrift(position=1000), SEA(2), GradualDrift(position=2000, width=500), SEA(3)])
```

The GradualDrift can also be specified in terms of start and end:

```
DriftStream([SEA(1), AbruptDrift(position=1000), SEA(2), GradualDrift(start=1750, end=2250), SEA(3)])
```

- This can make the drift locations more explicit and easy for new comers

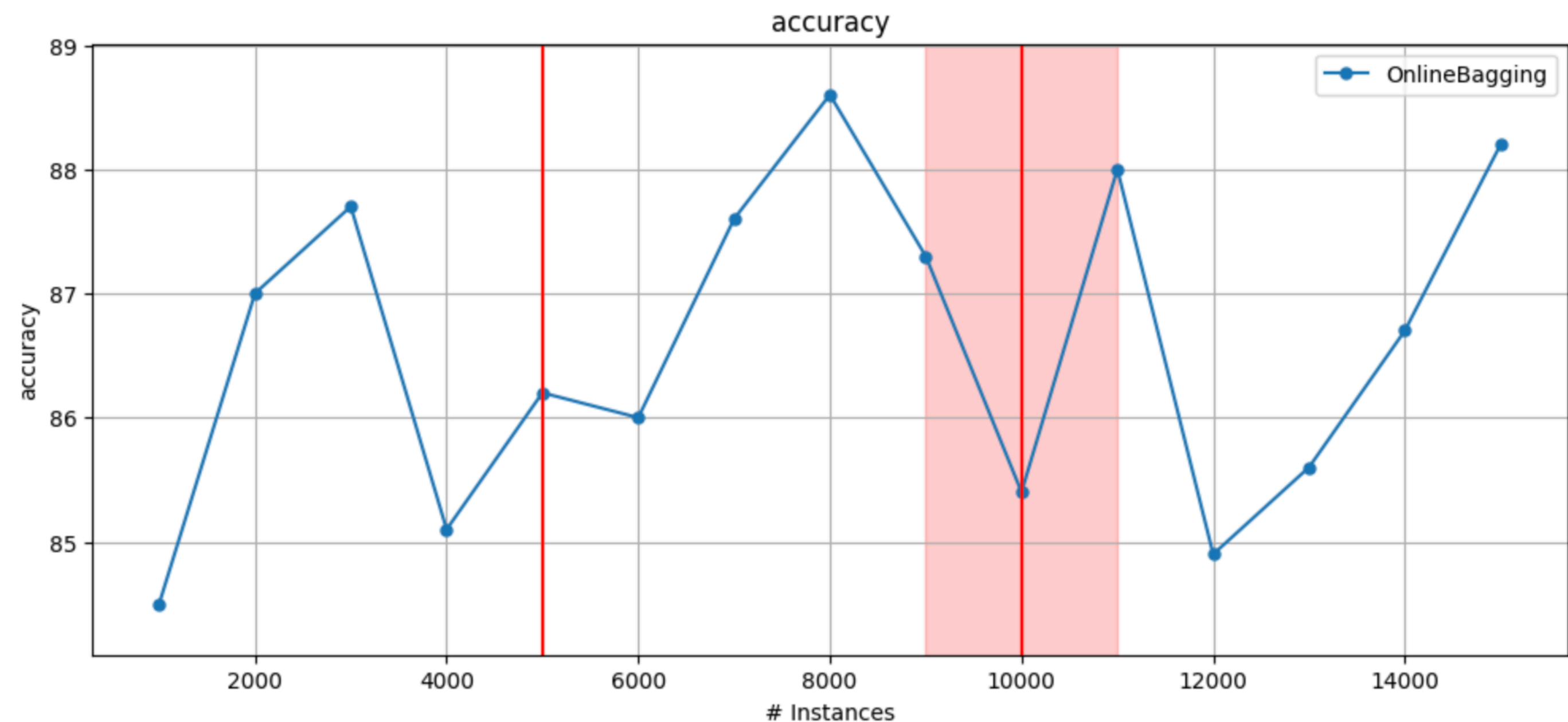
Code example

```
DriftStream(stream=[SEA(function=1),  
AbruptDrift(position=5000),  
SEA(function=3),  
GradualDrift(position=10000, width=2000),  
SEA(function=1)])
```

Both approaches will generate a similar output.

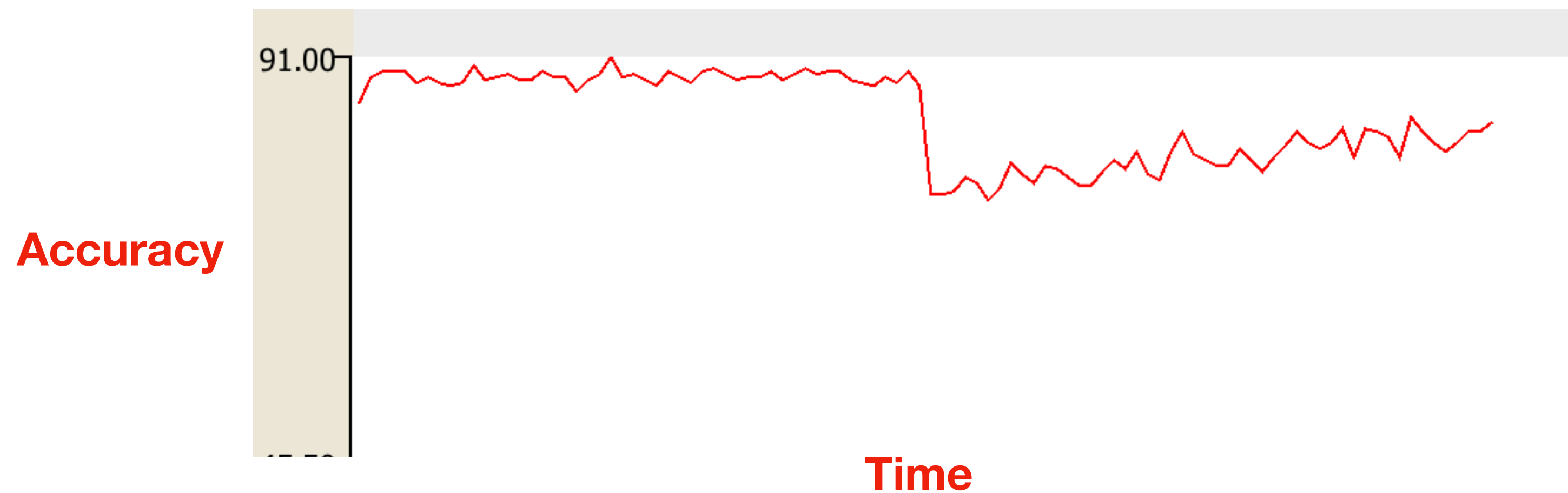
The second one use CapyMOA generic API to invoke the MOA CLI

```
Stream(moa_stream=ConceptDriftStream(),  
      CLI='-s  
(ConceptDriftStream -s  
generators.SEAGenerator -d  
(generators.SEAGenerator -f 3) -p 5000 -w  
1) -d generators.SEAGenerator -w 2000 -p  
10000 -r 1')
```



Evaluation (detectors)

Common approach (proxy): “**Attach the method to a classifier, if the accuracy goes up, then the detector works**”



Not necessarily the detector is successful in detecting changes, maybe it is just randomly resetting the classifier!

We must use **specific metrics** to **evaluate** a detector

Evaluation (detectors)

Important: we need the ground-truth of drift location for some of these

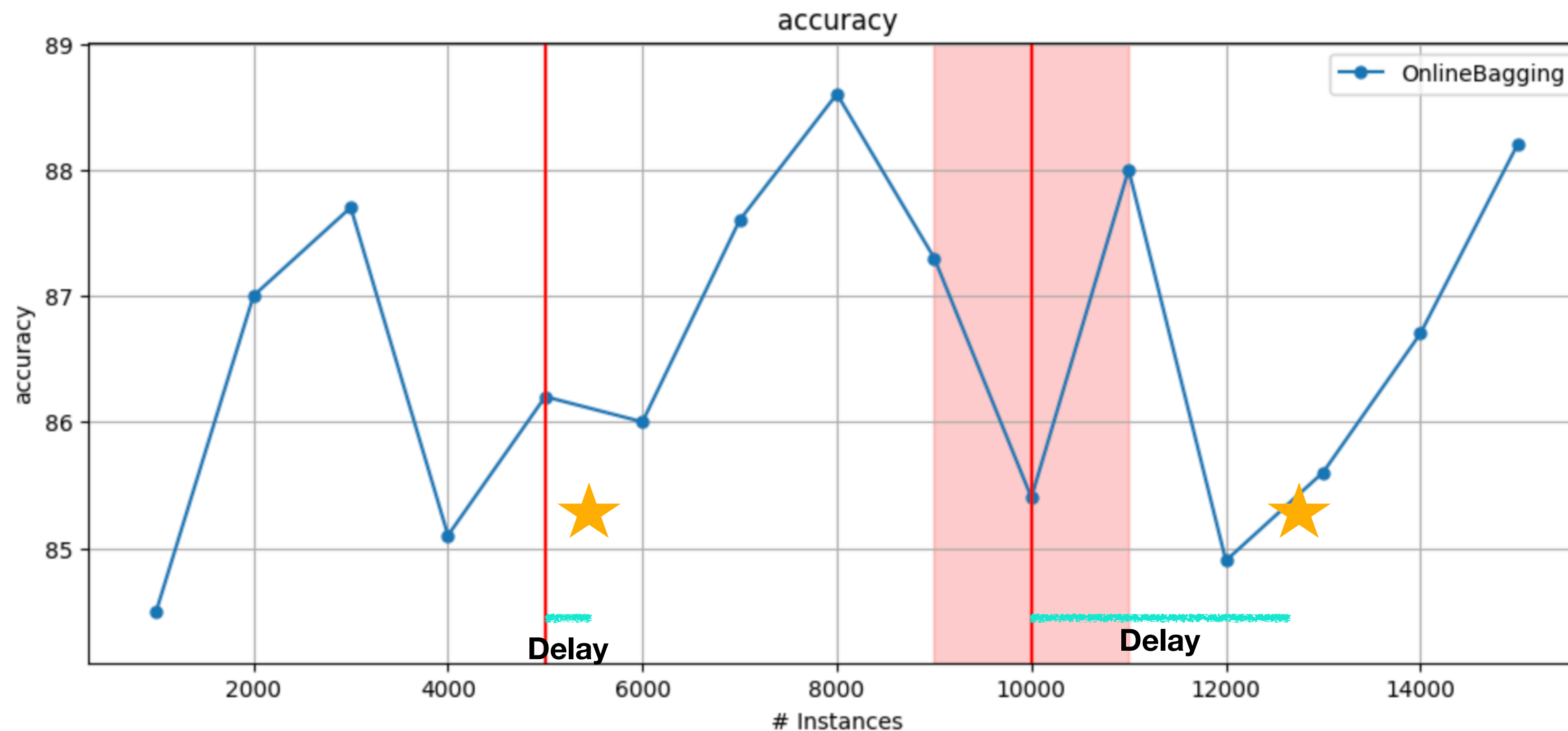
Some Metrics:

- Recall, Precision, ...
- Number of detections
- Mean Time between False Alarms (MTFA)
- Mean Time to Detection (MTD)
- And others: MDR, ARL, MTR, ...

Evaluation (detectors)

We might want to only account for detections if they are within a **max_delay**

Let's assume in the example below that **yellow stars** are detections, we can observe some delay between the drifts (red vertical lines or rectangles (gradual)) and detections.

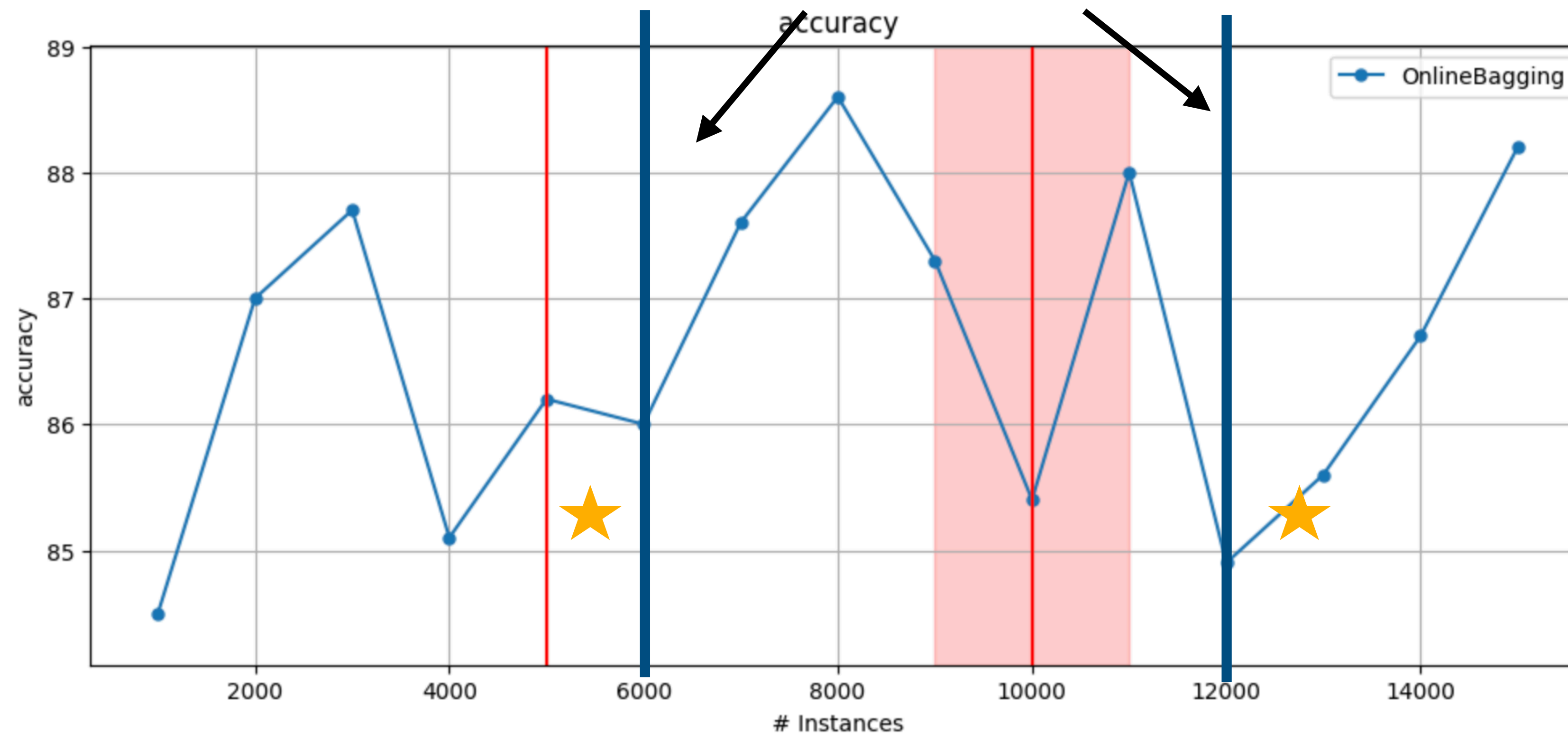


Evaluation (detectors)

We might want to only account for detections if they are within a **max_delay**

Let's assume in the example below that **yellow stars** are detections, we can observe some delay between the drifts (red vertical lines or rectangles (gradual)) and detections.

We can also specify **max_delay** to determine when a detection should be considered a TP



Using Drift Detectors

Ideally, using concept drift detectors should be straightforward

There should be some way to **update** it with new values, and some way of **detecting that a drift** has been detected

Using Drift Detectors

```
# Create a learner nb, create a detector, declare the stream, ...

while stream_sealdrift.has_more_instances() and i < max_instances:
    instance = stream_sealdrift.next_instance()
    pred = nb.predict(instance)
    evaluator.update(instance.y_index, pred)

    is_correct = int(pred == instance.y_index)

    detector.add_element(is_correct)

    if detector.detected_change():
        print('Change detected at instance: ' + str(i))

    nb.train(instance)

# ...
```

Recurrent Concept Drifts

Recurrent Concept Drifts

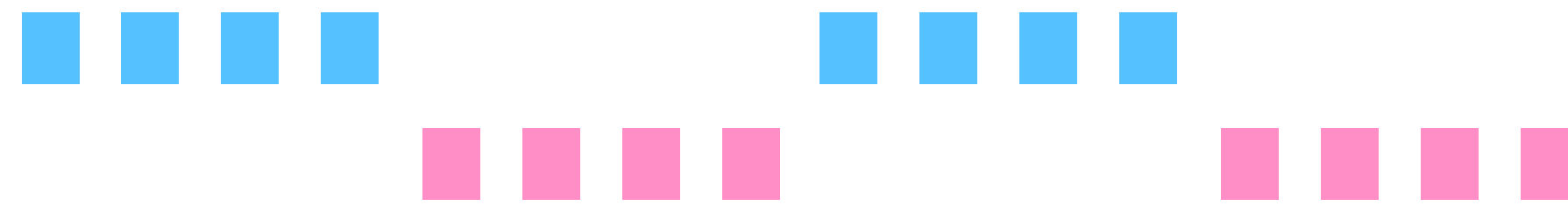
- by **concept transition**
- by **time of recurrence**

Recurrent Concept Drifts

- **by concept transition**

Recurrent Concept Drifts

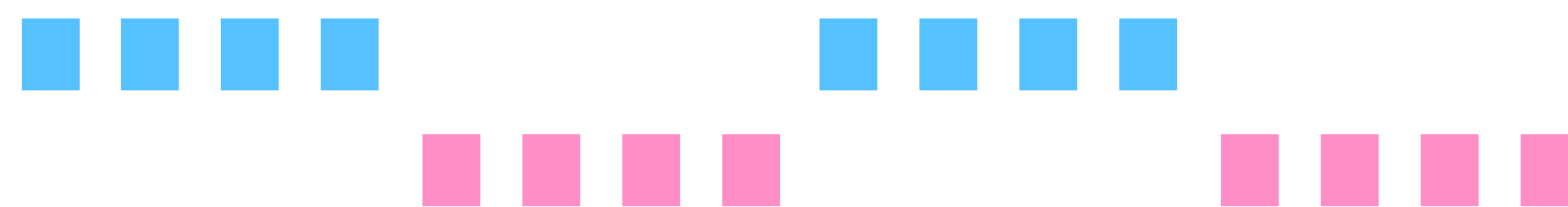
- by **concept transition**



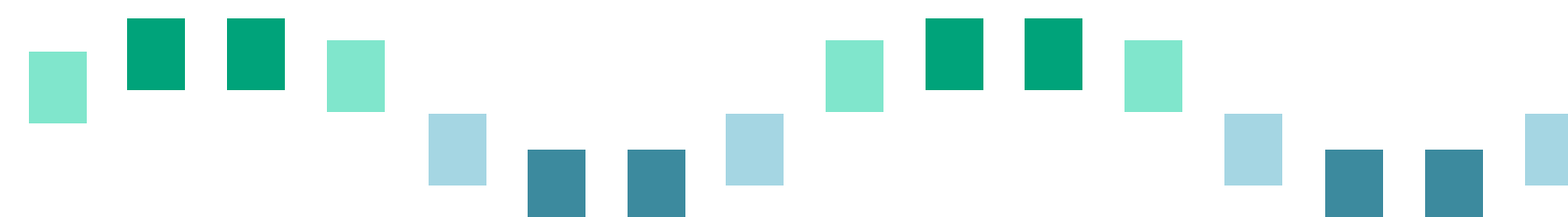
a) abrupt recurrent drift

Recurrent Concept Drifts

- by **concept transition**



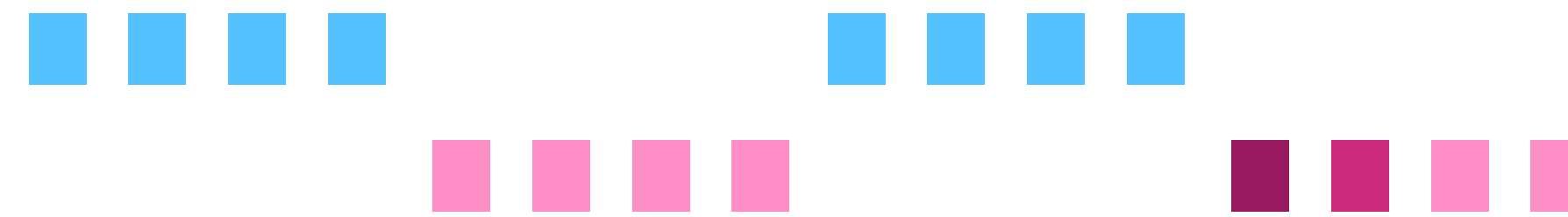
a) abrupt recurrent drift



b) incremental recurrent drift

Recurrent Concept Drifts

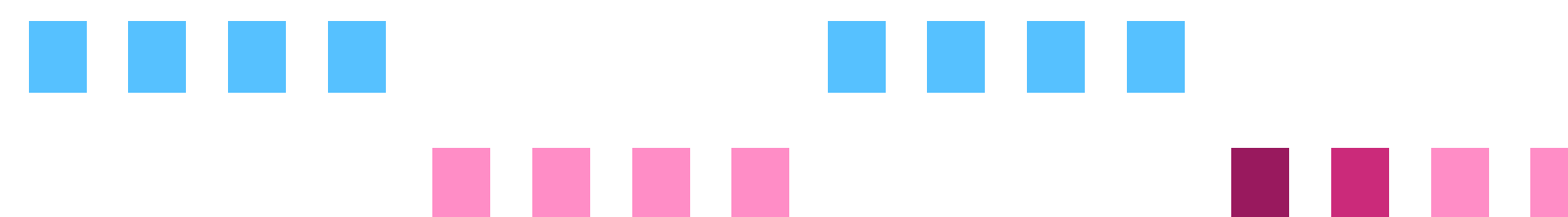
- by **concept transition** cont..



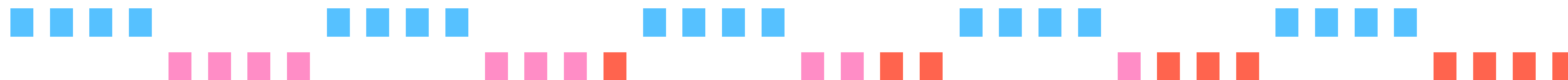
d) partial recurrent drift

Recurrent Concept Drifts

- by **concept transition** cont..



d) partial recurrent drift



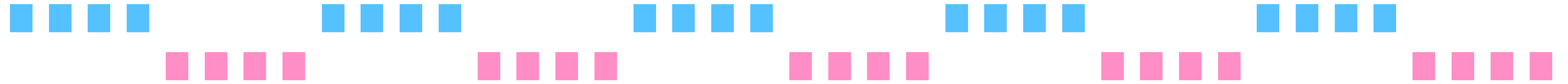
e) evolving recurrent drift

Recurrent Concept Drifts

- **by time of recurrence**

Recurrent Concept Drifts

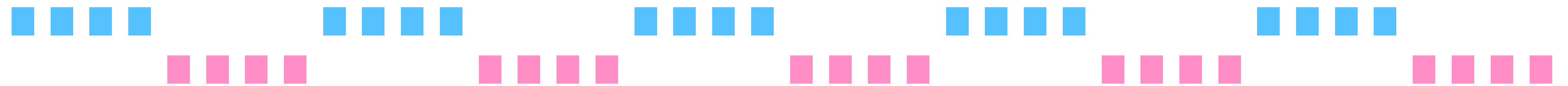
- by time of recurrence



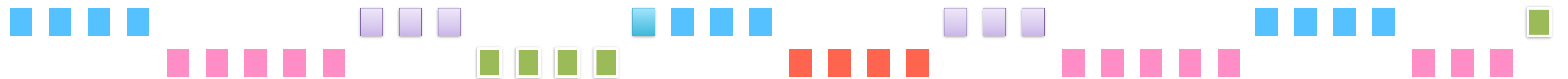
a) periodic recurrent drifts

Recurrent Concept Drifts

- by time of recurrence



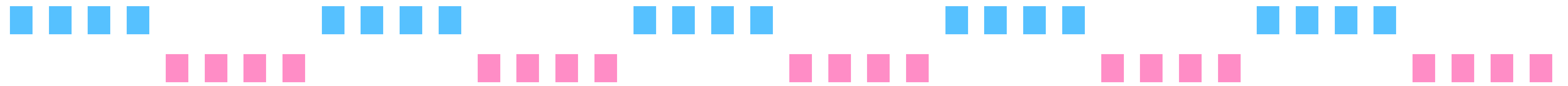
a) periodic recurrent drifts



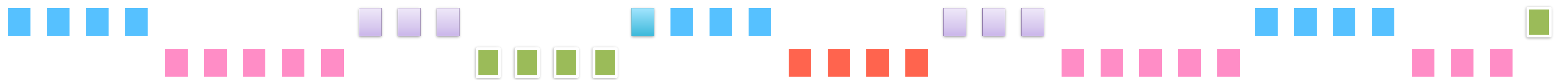
b) semi-periodic recurrent drifts

Recurrent Concept Drifts

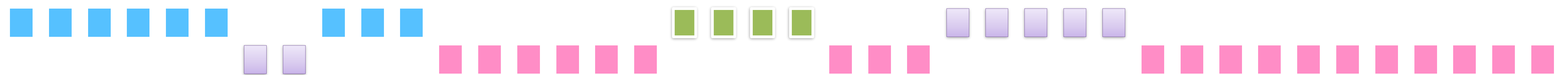
- by time of recurrence



a) periodic recurrent drifts



b) semi-periodic recurrent drifts



c) random recurrent drifts

Drift Re-Cap

- Many types of **drifts**
 - abrupt, incremental, gradual and recurrent
- Many types of **recurrent concept drifts**
 - by transition
 - by recurrence

Learners that cope with recurrent CD

- Ideally, these methods attempt to retain the knowledge acquired on a given concept and reuse it whenever that concept reappears.
- Simple example: Maintain a pool of learners, somehow identify if a concept is related to a given learner from the pool.
- Challenges:
 - identify the compatibility of two concepts
 - maintain the pool of learners

Methods

Design components:

- **DD**: Drift Detection
- **DP**: Drift Prediction
- **LM**: Labels Missing
- **LN/A**: Labels Not Available
- **MetaL**: Meta Learning
- **MetaF**: Meta Features
- **Clust**: Clustering
- **CEqSim**: Conceptual Equivalence/Concept Similarity
- **Ens**: Ensemble
- **CPool**: Concept Pool.

Please refer to section 3 of the survey for more information

Methods

Explicit Handling of Recurrences

(model for each data batch)

Meta Learning

(act as a wrapper algorithm to determine the best model/s for the current concept)

Clustering

Drift Prediction

Meta Features

Sec	Method	Year	DD	DP	LM	LN/A	MetaL	MetaF	Clust	CEqSim	Ens	CPool
3.1	LEARN++*	2011									X	X
	PMRCD	2012									X	X
	Dynse	2018									X	X
	ASE	2017								X	X	
	GraphPool	2018								X	X	X
3.2	RCD	2013	X				X			X	X	X
	CPF	2016	X				X			X		X
	ECPF	2019	X				X			X		X
	PEARL	2022	X				X				X	X
3.3	REDLLA	2012	X		X				X	X		X
	SUN	2012	X		X				X	X		X
	ContexTrac	2012	X						X	X		X
	ESCR	2021	X		X				X	X	X	
	CDMSE	2021	X		X				X		X	
	CCP	2010						X	X		X	
	UClust	2020	X		X	X		X	X	X		X
	CDCMS	2020	X						X		X	
3.4	MM-PRec	2016	X	X			X			X		X
	PCCF	2016		X								
	BLPA	2017	X	X								
	CPRD	2019	X	X								
	ProSeed	2016	X	X								
	ProChange	2018	X	X	X							
	MDP	2018	X	X				X	X			
	Nacre	2021	X	X							X	
3.5	SELeCT	2022						X				X
	FiCSUM	2023	X					X				

Design components: **DD**: Drift Detection, **DP**: Drift Prediction, **LM**: Labels Missing, **LN/A**: Labels Not Available, **MetaL**: Meta Learning, **MetaF**: Meta Features, **Clust**: Clustering, **CEqSim**: Conceptual Equivalence/Concept Similarity, **Ens**: Ensemble, **CPool**: Concept Pool.

Please refer to section 3 of the survey for more information

Evaluation

- Under recurrent concept drift
 - Model performance
 - Drift Detection performance

Evaluation

Relative Performance

- **compares** the performance of classifier A **against a baseline classifier B .**
- **at instance i :** $\log(B_{error_i}/A_{error_i})$ [1]
- **Cumulative Accuracy Gain:** $\sum (accuracy(A) - accuracy(B))$ [2]

[1] Joao Gama, Raquel Sebastiao, and Pedro Pereira Rodrigues. On evaluating stream learning algorithms. Machine learning, pages 317–346, 2013.

[2] Ocean Wu, Yun Sing Koh, Gillian Dobbie, and T Lacombe. Probabilistic exact adaptive random forest for recurrent concepts in data streams. Int. J. Data Sci. Anal., pages 1–16, 2022.

Evaluation

Model Selection for Each Concept

- measures the strength of the relationship between each $\langle model, context \rangle$ pair [1].
 - **context**: an underlying condition that results in a concept
 - **measures the context linkage** for model reuse (from a model pool)

Evaluation

Drift Detection on **Synthetic Data** (drift points known in advance)

- **True drift** points are **compared** to the **drift detected points** to detect **Type I (FP)** and **Type II (FN)** errors. [1-4]

[1] Robert Anderson, Yun Sing Koh, and Gillian Dobbie. Predicting concept drift in data streams using metadata clustering. In IJCNN, pages 1–8. IEEE, 2018.

[2] Alexandr Maslov, Mykola Pechenizkiy, Indre Z'liobaite, and Tommi Ka'rkka'inen. Modelling recurrent events for improving online change detection. In SDM, pages 549–557. SIAM, 2016.

[3] David Tse Jung Huang, Yun Sing Koh, Gillian Dobbie, and Russel Pears. Detecting volatility shift in data streams. In ICDM, pages 863–868. IEEE, 2014.

[4] Yun Sing Koh, David Tse Jung Huang, Chris Pearce, and Gillian Dobbie. Volatility drift prediction for transactional data streams. In ICDM, pages 1091– 1096. IEEE, 2018.

Open Source Software & Benchmark Datasets

- Most methods have **custom** open source implementations
- **Traditional streaming datasets**
 - Real world: may not know the reoccurrence
 - Synthetic : reproducibility
- **CapyMOA recurrent concept API**

Please refer to: Section 5 & 6 of the survey for more information

Delayed & Unsupervised Drift Detection

Delayed & Unsupervised Drift Detection

- Most concept drift detection algorithms are applied to the **univariate stream of correct/incorrect** classifier predictions
- Such strategies require that labeled data is available as soon as possible to respond to concept drifts in a timely fashion
- Despite their intrinsic differences, most drift detectors **trigger when the observed model's predictive performance starts to degrade**

Delayed & Unsupervised Drift Detection

Terminology

Delayed drift detection: The label will arrive at some point in the future, and it will be used for feeding the learner with a **delayed** univariate stream of correct/incorrect predictions

Unsupervised drift detection: The label **will not arrive**, thus the detection should be based on the input data or the output of the learner itself

Delayed Drift Detection (Example)

- Experiment with data generated using the AGRRAWAL generator with **3 abrupt concept drifts** (at instances 25, 000, 50, 000, and 75, 000).
- Used an ensemble algorithm capable of detecting and adapting to changes by resetting base models whenever **changes** are **detected** on their **univariate stream of correct/incorrect** predictions.
- Figure depicts the amount of concept drifts detected (y-axis) over the processing of 100,000 instances with and without delayed labelling.
- The detections for the “No delay” experiment shows a high rate of detection immediately after the concept drifts, except for a few arbitrary drift signals in-between the concept drifts.

Delayed Drift Detection (Example)

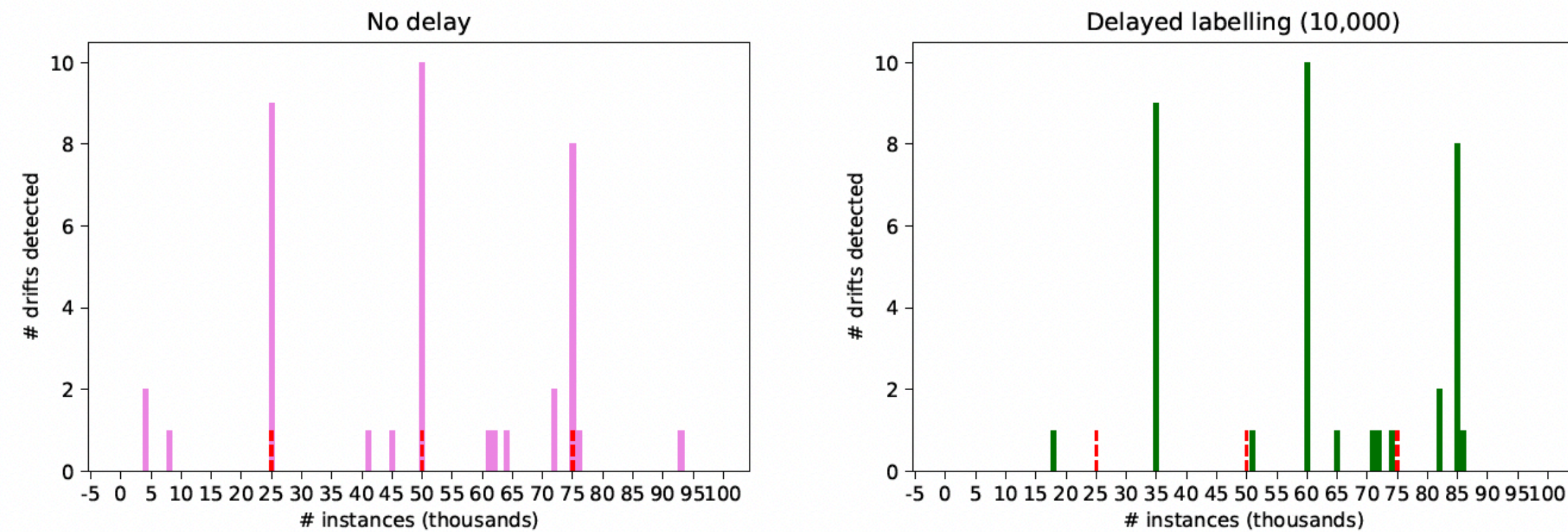


Fig. 7. Drifts detected by a 10 learner SRP model using ADWIN on AGRAWAL with and without labelling delay. Red dotted vertical lines indicate the location of concept drifts.

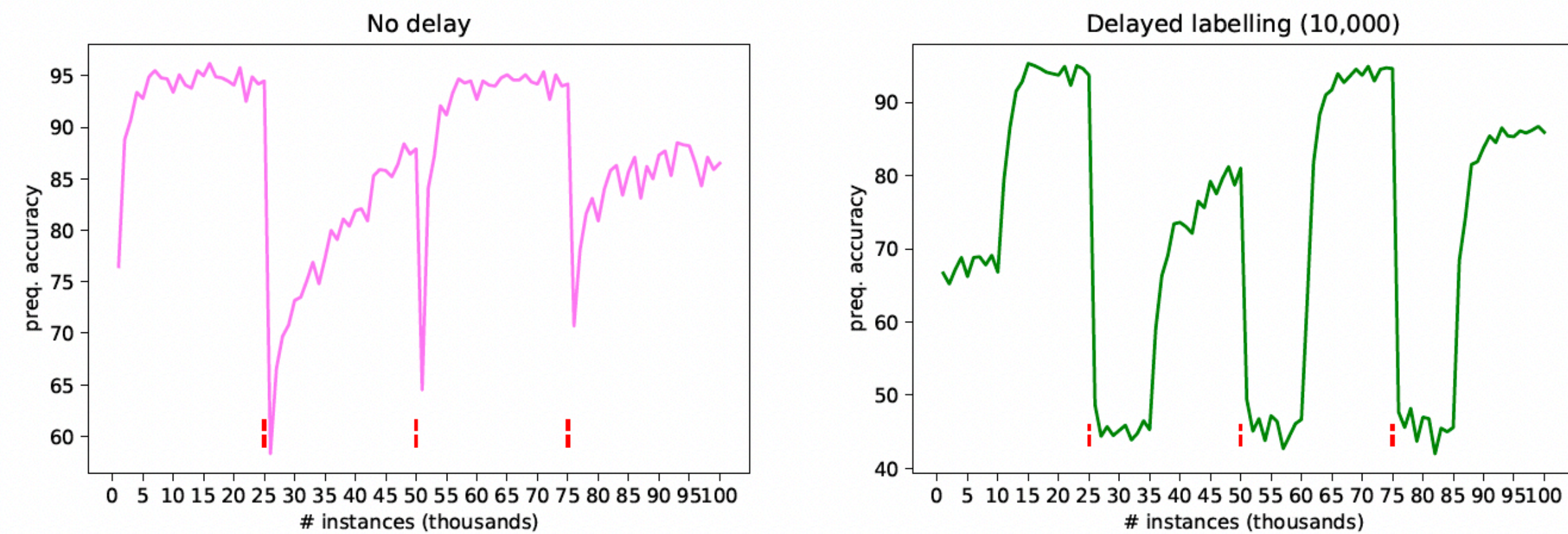


Fig. 8. Accuracy by a 10 learner SRP model using ADWIN on AGRAWAL with and without labelling delay. Red dotted vertical lines indicate the location of concept drifts.

STUDD: Unsupervised Concept Drift Detection using a Student–Teacher Approach

- Detecting concepts drifts in the absence of labeled data
- **Procedure.** A predictive model (teacher) is built using an initial batch of labelled training data. The teacher's predictions are used as class labels to train a surrogate model (student), which will learn to mimic the teacher. A drift detection algorithm is used to identify variations in the mimicking error of the student.
- **Hypothesis.** If the mimicking error increases, then it means that a concept drift has occurred.

Conclusions

Conclusions

- Practical aspects w.r.t. CD: simulate, evaluate, utilise
- Opportunities in identifying drifts on partially and delayed labeled settings
- Opportunities w.r.t. recurrent drifts in the intersection with Online Continual Learning
 - **Datasets** (OCL -> Recurrent SL)
 - **Model pool management** techniques (OCL <- Recurrent SL)

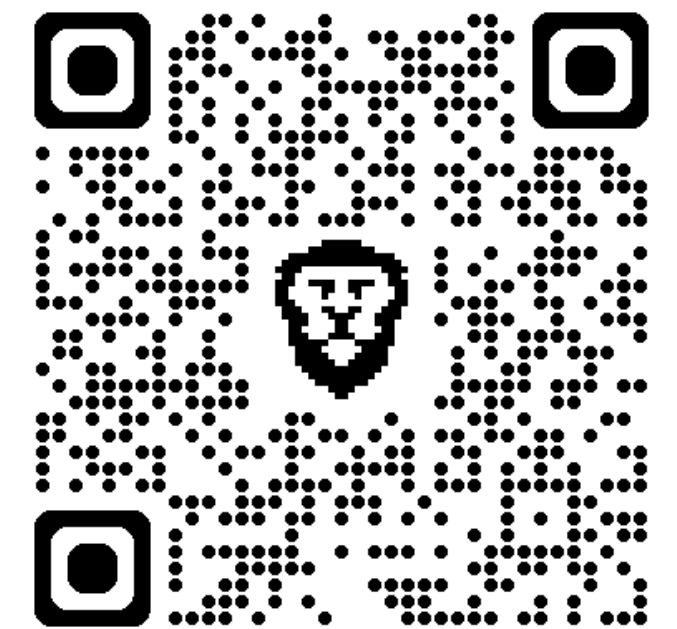
Thank you!

- Consider trying **CapyMOA** for your drift detection needs!

Contact: heitor.gomes@vuw.ac.nz



<https://discord.gg/RekJArWKNZ>



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