# Building Adaptive Data Mining Models on Streaming Data in Real-Time, an Outlook on Challenges, Approaches and Ongoing Research

## Marine Perception Research Department of the German Research Center for Artificial Intelligence (DFKI)

Frederic Stahl

www.dfki.de/map

Marie-Curie-Str. 1
26129 Oldenburg
Germany

map-info@dfki.de

# Data never sleeps!

- Forbes: 2.5 quintillion bytes of data created every day.

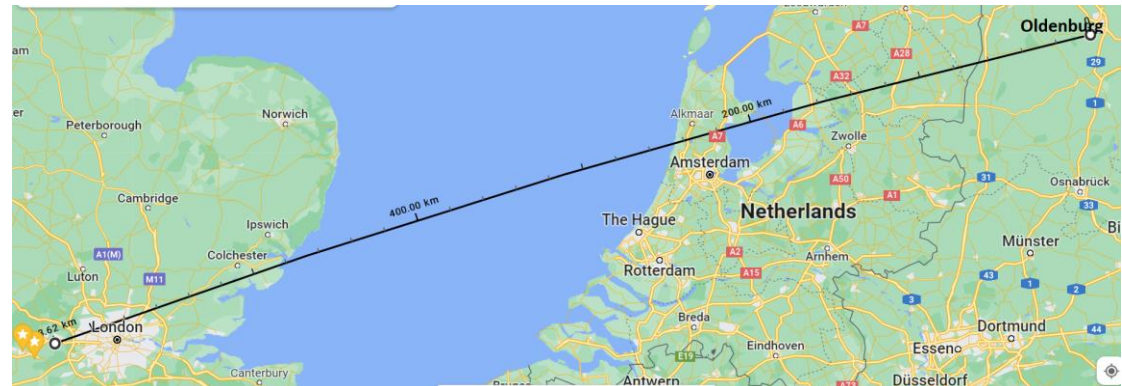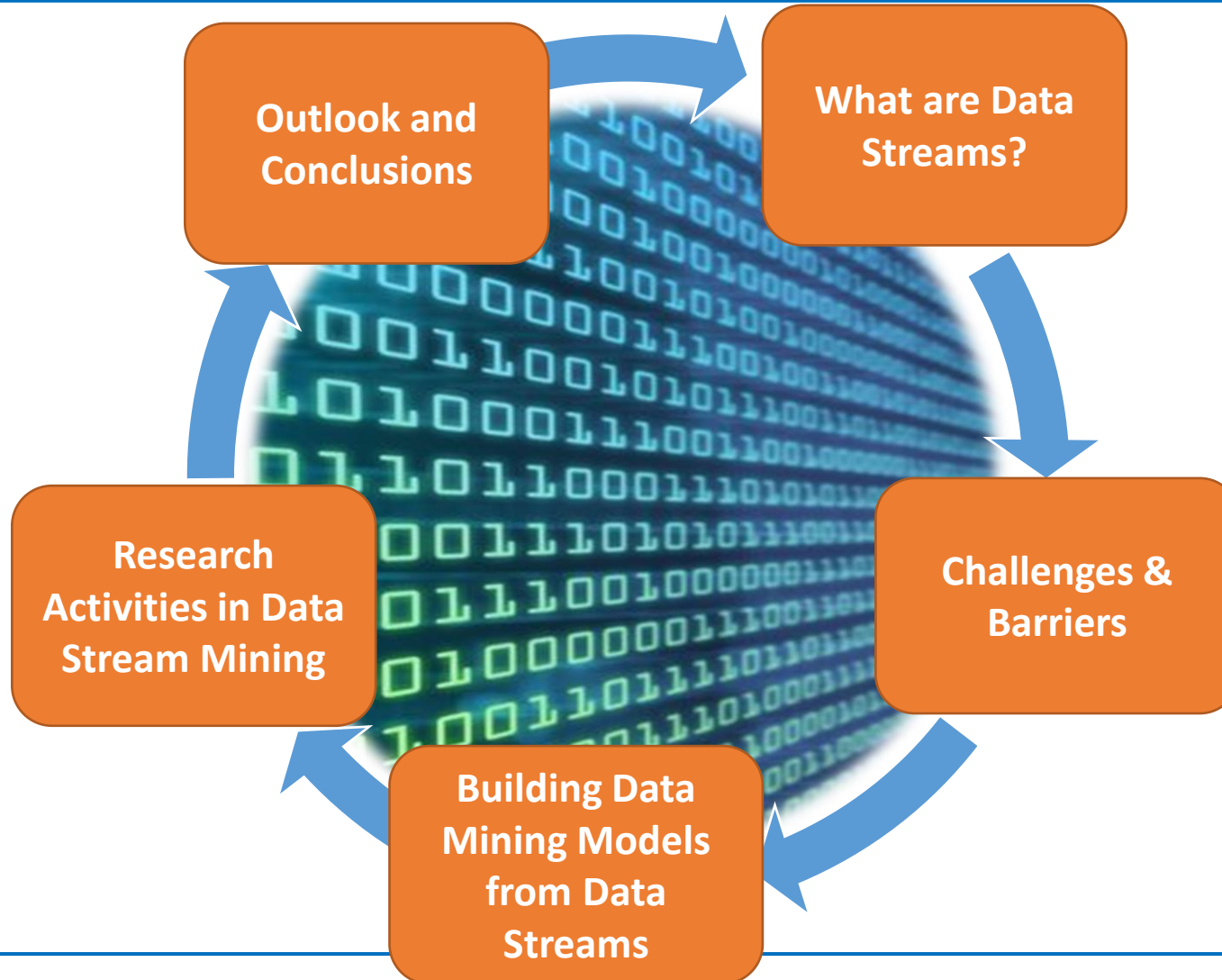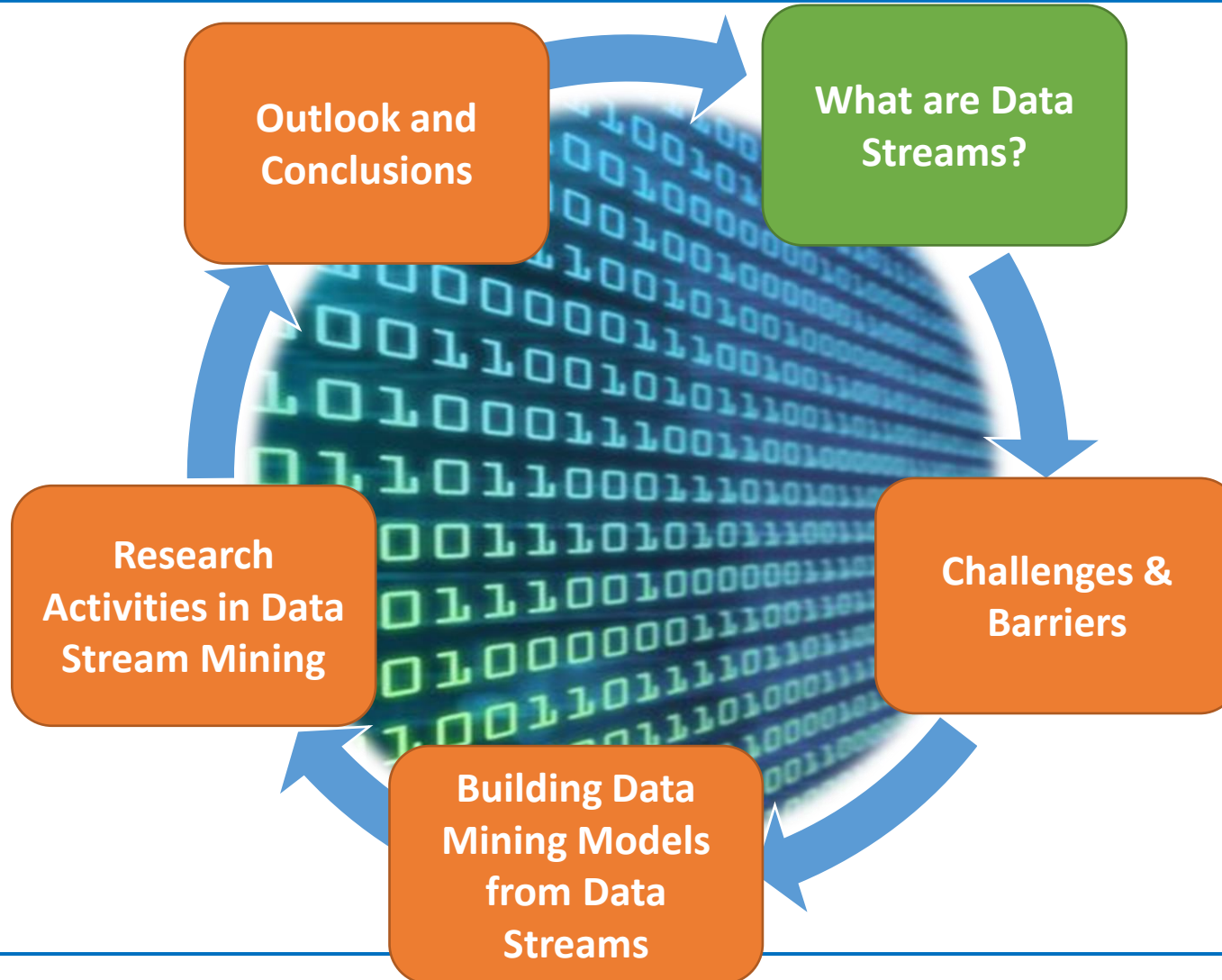- That's about 100 million Blue-ray discs or about 530 million DVD discs.

# How much Data Is created Every day?

- That's about 100 million Blue-ray each 25 GB discs.
- Each disc is 1.2mm thick
    - $\Rightarrow$ This stacks to **120 km!**
    - $\Rightarrow$ Distance Oldenburg to Hamburg!

- Or in DVDs (4.7 GB each disc)
- Each disc is 1.2mm thick
    - $\Rightarrow$ This stacks to **630 km!**
    - $\Rightarrow$ Distance Oldenburg to London/Reading!

To make sense of this real-time data,
analytics methods that never sleep
are required!

# Outline



What are Data Streams?

Challenges & Barriers

Building Data Mining Models from Data Streams

Research Activities in Data Stream Mining

Outlook and Conclusions

Outlook and Conclusions

What are Data Streams?

Research Activities in Data Stream Mining

Challenges & Barriers

Building Data Mining Models from Data Streams
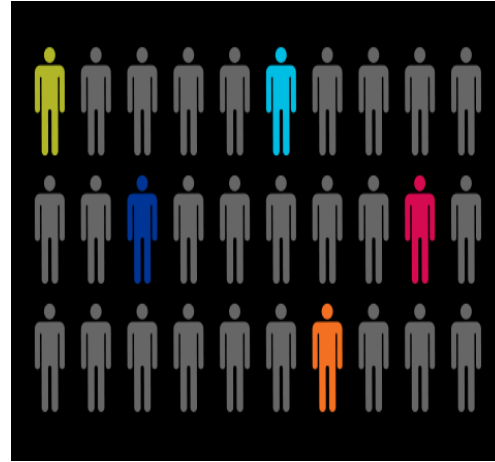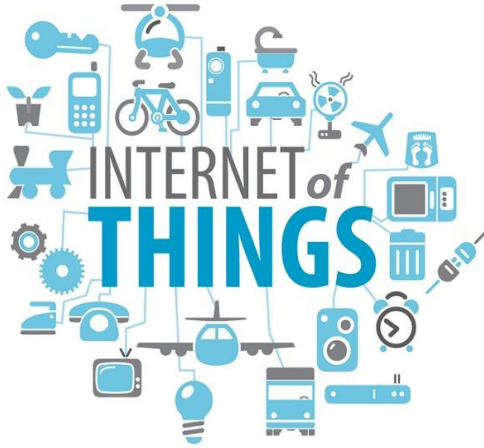
# Sources of Data Streams



**Internet of Things**

- By year-end 2039, IoT devices worldwide are forecasted to almost triple from 9.7 billion in 2020 to 29 billion in 2030 [1]

[1] statistica. (2020). Number of Internet of Things (IoT) connected devices worldwide from 2019 to 2021, with forecasts from 2022 to 2030

# Sources of Data Streams



**Internet of Things**

- By year-end 2039,  IoT devices worldwide are forecasted to almost triple from 9.7 billion in 2020 to 29 billion in 2030 [1]
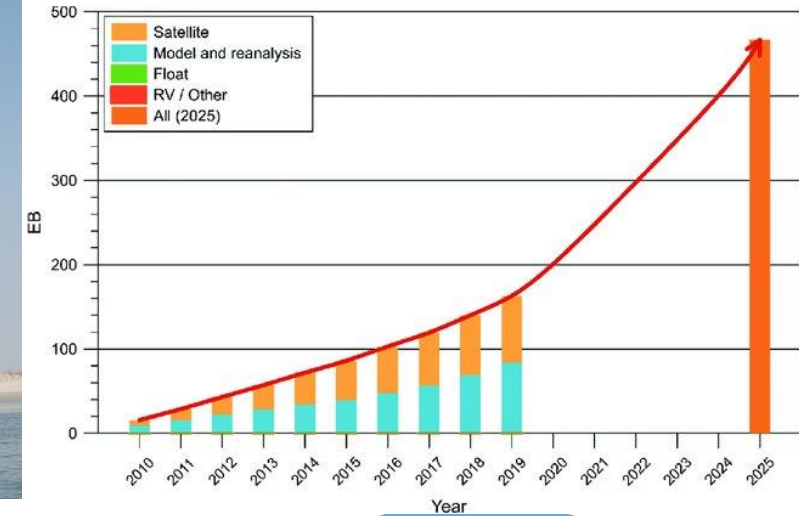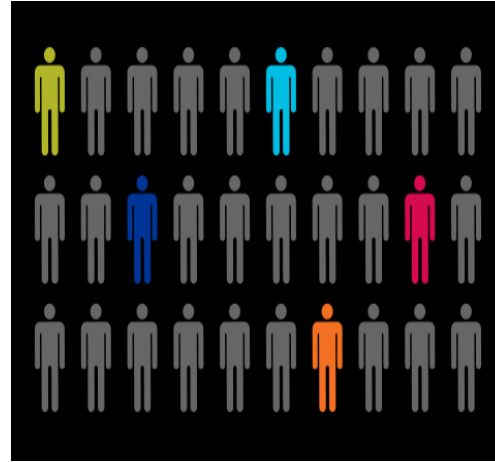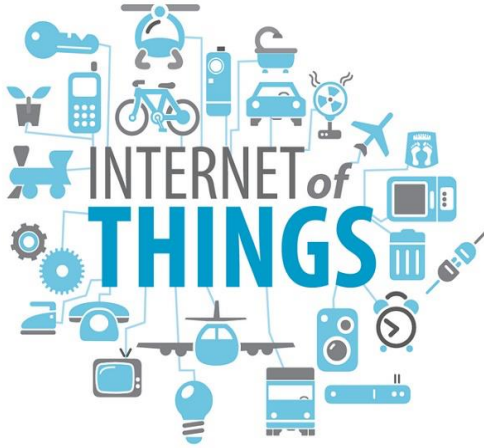
**Personalisation**

- Facebook:
  - 1.91 billion active users every day [2]
  - 4.75 billion pieces of content shared

[1] statistica. (2020). Number of Internet of Things (IoT) connected devices worldwide from 2019 to 2021, with forecasts from 2022 to 2030

[2] Noyes, A. and Noyes, D. (2014). The Top 20 Valuable Facebook Statistics - Updated October 2014 –     Zephoria Inc.. [online] Zephoria Inc. Available at: https://zephoria.com/social-media/top-15-valuable-facebook-statistics/ [Accessed 2022].

# Sources of Data Streams









## Internet of Things

- By year-end 2039, IoT devices worldwide are forecasted to almost triple from 9.7 billion in 2020 to 29 billion in 2030 [1]

## Personalisation

- Facebook:
  - 1.91 billion active users every day [2]
  - 4.75 billion pieces of content shared

## Marine Sciences

- Distribution of ocean science data acquired in the past decade, based on publicly available data from the internet (CC BY 4.0) [3]
- Expected to reach almost 500 Exabytes by the year 2025

[1] statistica. (2020). Number of Internet of Things (IoT) connected devices worldwide from 2019 to 2021, with forecasts from 2022 to 2030

[2] Noyes, A. and Noyes, D. (2014). The Top 20 Valuable Facebook Statistics - Updated October 2014 – Zephoria Inc.. [online] Zephoria Inc. Available at: https://zephoria.com/social-media/top-15-valuable-facebook-statistics/ [Accessed 2022].

[3] Qian, C., Huang, B., Yang, X. and Chen, G., 2022. Data science for oceanography: From small data to big data. Big Earth Data, 6(2), pp.236-250.

# Static versus Streaming Data

A data stream is *a continuous, rapid flow of data that challenges our state-of-the-art processing and communication infrastructure.*

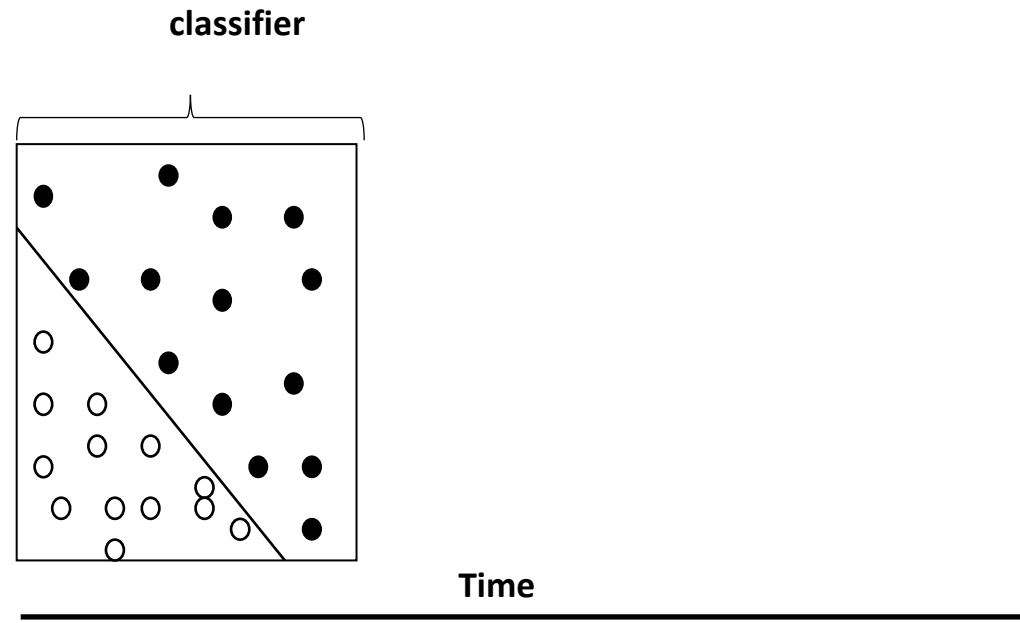| Static Data | Streaming Data |
|---|---|
| • Historical data | • Often live, real-time data feed |
| • Randomly accessible | • Sequentially accessed |
| • Secondary storage | • Limited memory requirements |
| • No/low processing latency criticality | • High processing latency criticality |
| • Assumption of pre-processed dataset | • Assumption of inaccurate raw data |

**Volume** and **Velocity**

# Concept Drift

- Underlying concept defining the knowledge being learned, begins to shift over time.

- Concept change is unforeseen and unpredictable.

- Concepts from the past may re-occur in the future.

- Concept drift exists in real-life problems:
  - Seasonal weather
  - Stock market rallies because of breaking news
  - etc.

# Concept Drift (cont.)
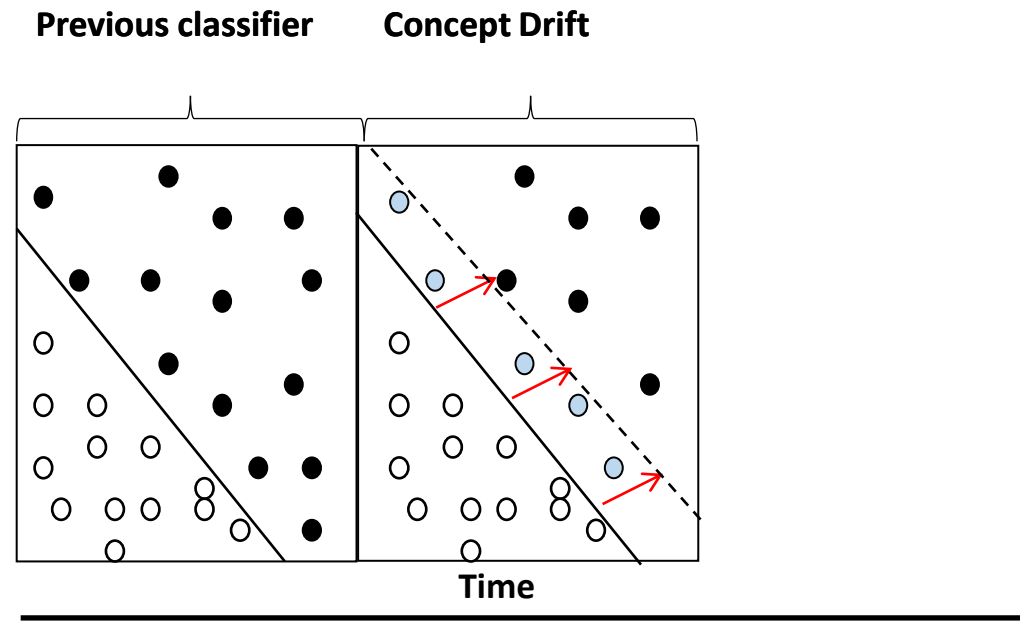
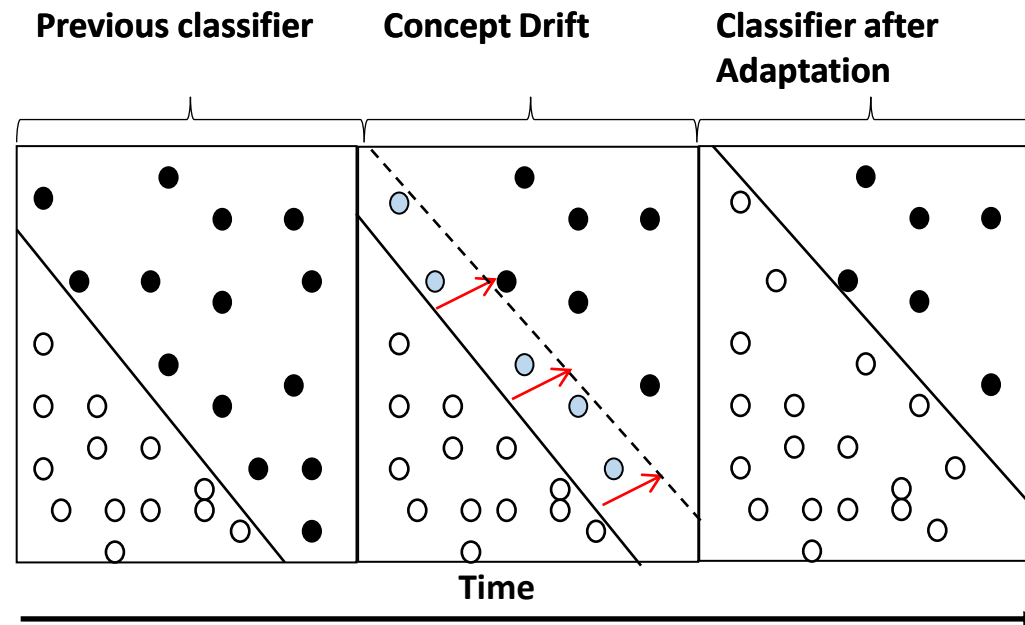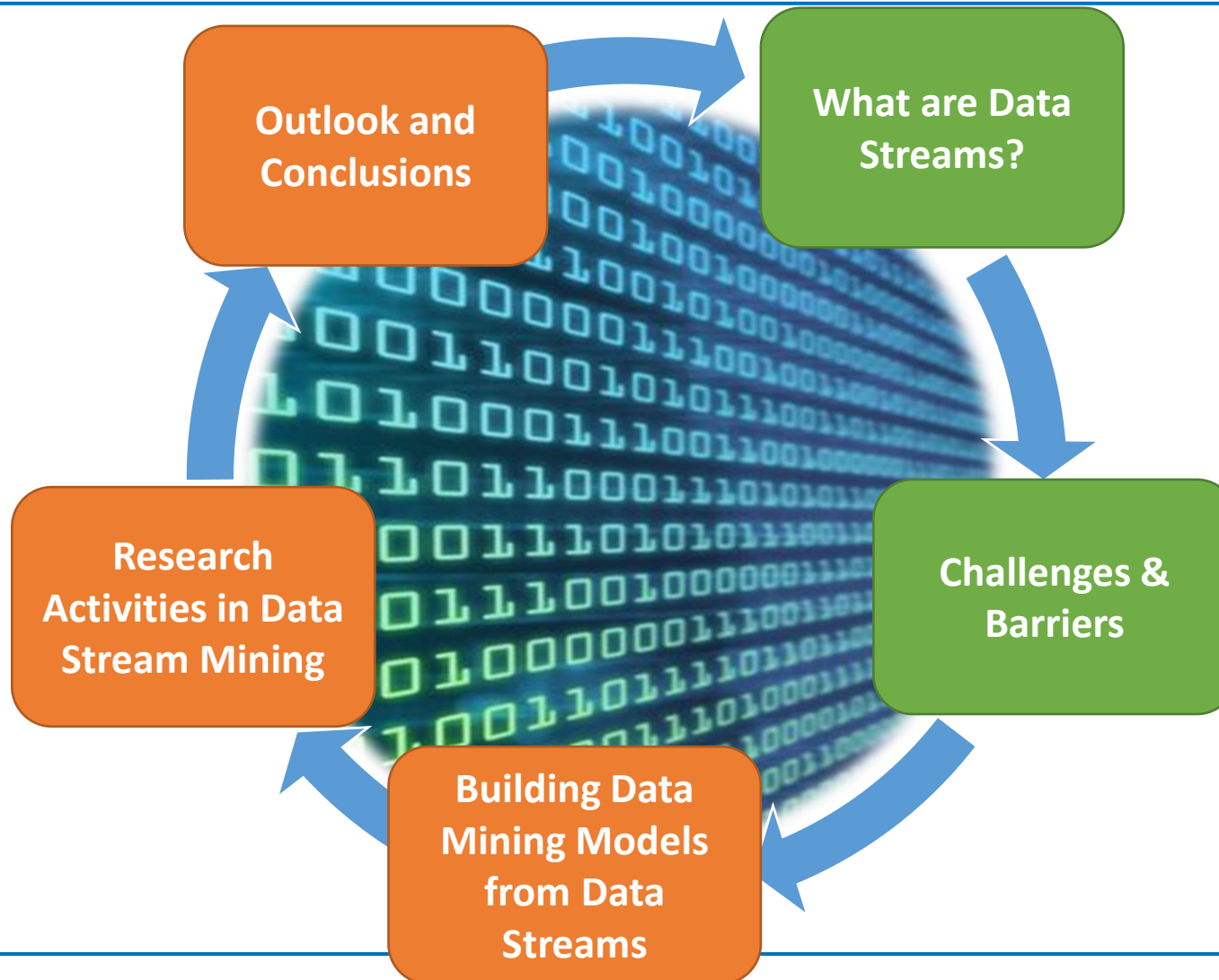Concept shift/drift: changes mining set statistics

- A model should always reflect the time-changing concept.
- Render previously learned models inaccurate or invalid.
- Robustness and adaptability: quickly recover/adjust after concept changes.

**classifier**

**Time**

# Concept Drift (cont.)

Concept shift/drift: changes mining set statistics

- A model should always reflect the time-changing concept.
- Render previously learned models inaccurate or invalid.
- Robustness and adaptability: quickly recover/adjust after concept changes.
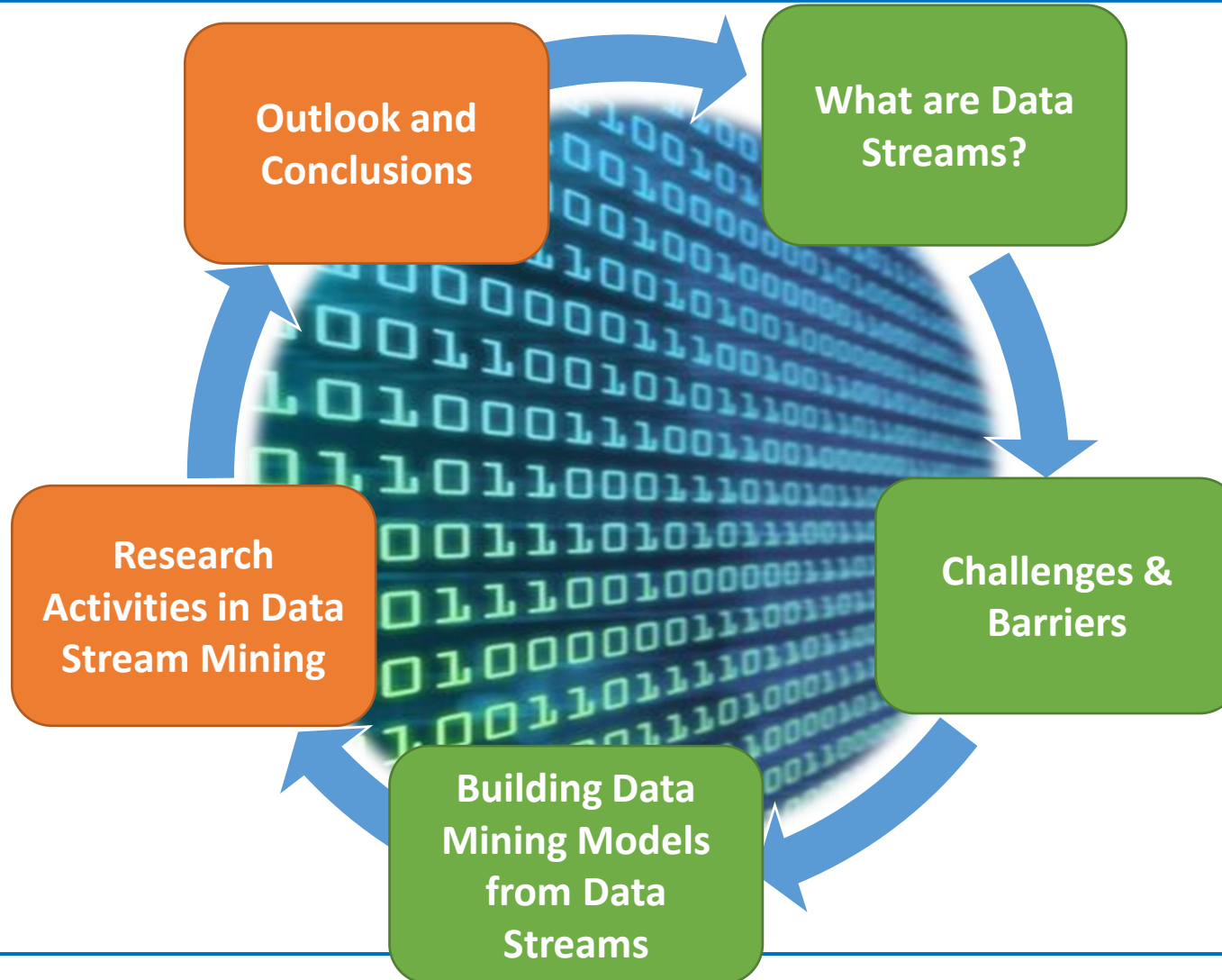
**Previous classifier**      **Concept Drift**

**Time**

# Concept Drift (cont.)

Concept shift/drift: changes mining set statistics

- A model should always reflect the time-changing concept.
- Render previously learned models inaccurate or invalid.
- Robustness and adaptability: quickly recover/adjust after concept changes.

# The Data Tsunami

| Challenges | Barriers |
|---|---|
| 1) Data generated at a fast rate (Velocity), at potentially large and unknown quantities (Volume) | 1) Limited scalable (parallel) real-time high throughput data stream mining algorithms |
| 2) Concept Drift (changes of pattern encoded in in the data over time) | 2) Different and changing types of concept drift |
| 3) Modelling real-time analytics workflows from streaming data | 3) Lack of customisable pre-processing techniques |
| 4) Multi-modality of data sources (text, video/images, unstructured) | 4) Different time stamps but co-occurring data items |
| 5) Class label sparsity: adapting predictive models | 5) Supervised algorithms not applicable in many cases |
| 6) Explaining Concept Drift | 6) Lack of drift detectors explaining concept drift |

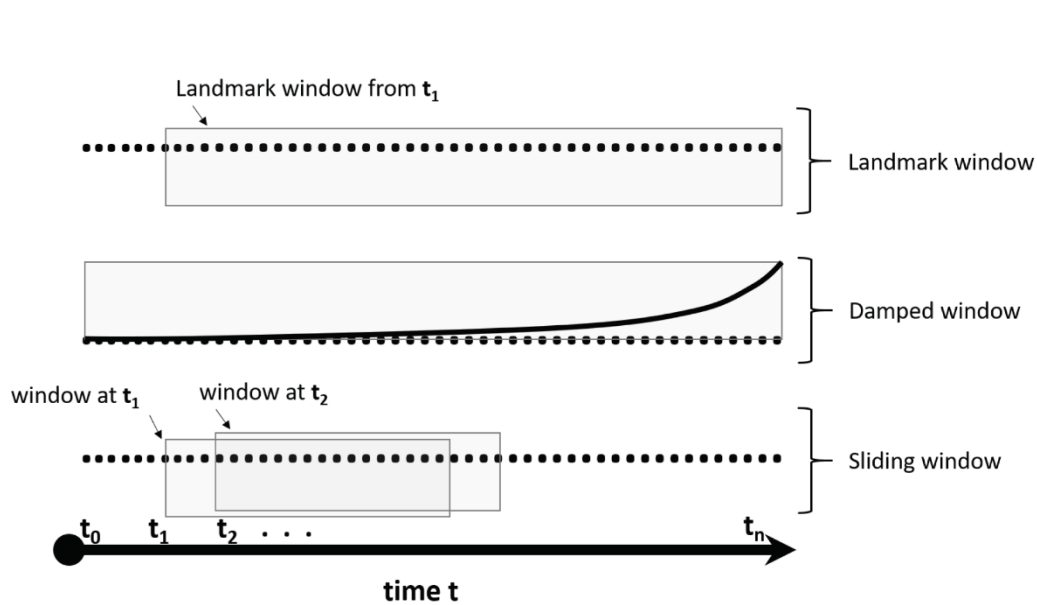# Methods: Windowing approaches to induce data mining models
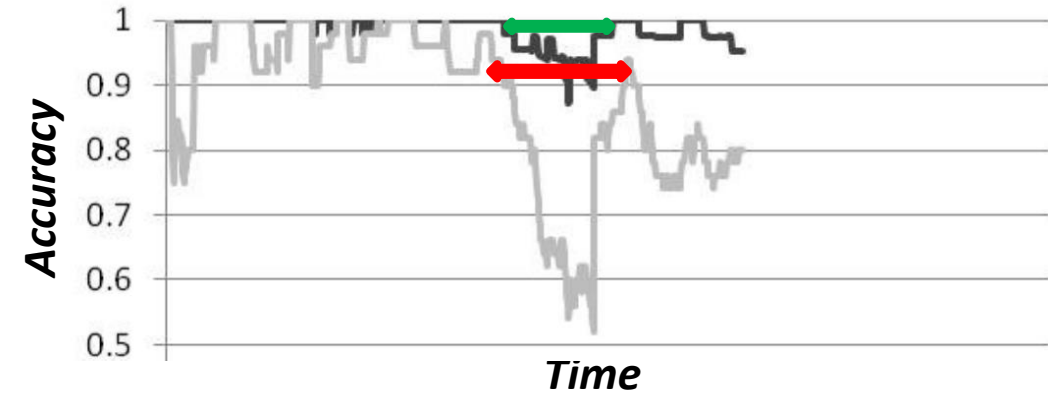
## 1) Create time windows



Source: Stahl, F., Le, T., Badii, A., Gaber, M.M. (2021) A frequent pattern conjunction Heuristic for rule generation in data streams. Information 12(1) (2021), ISSN 2078-2489, doi: 10.3390/info12010024

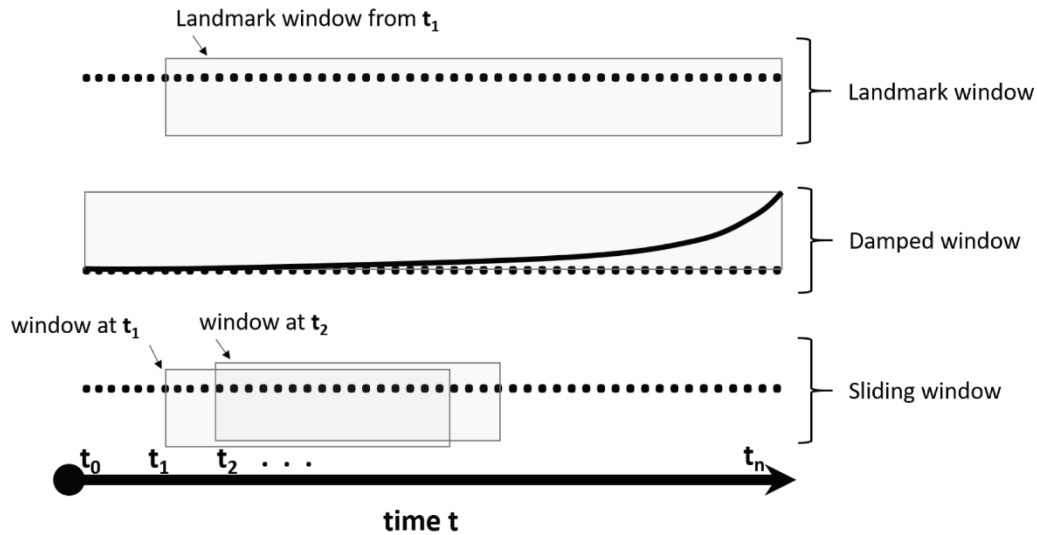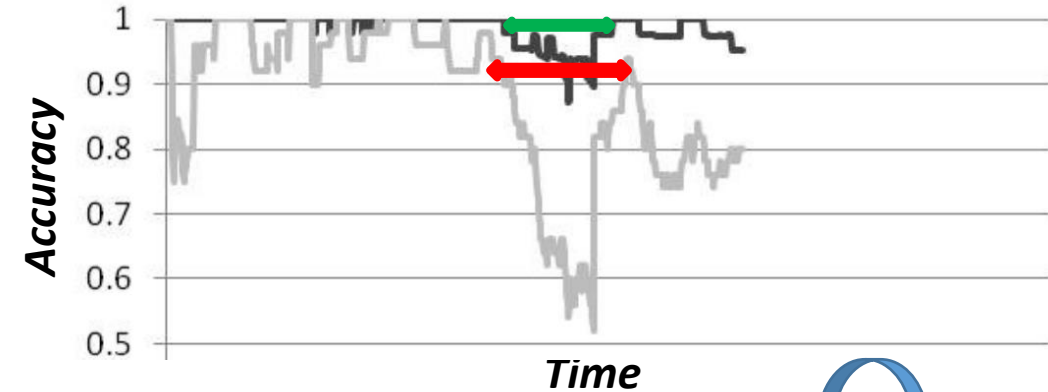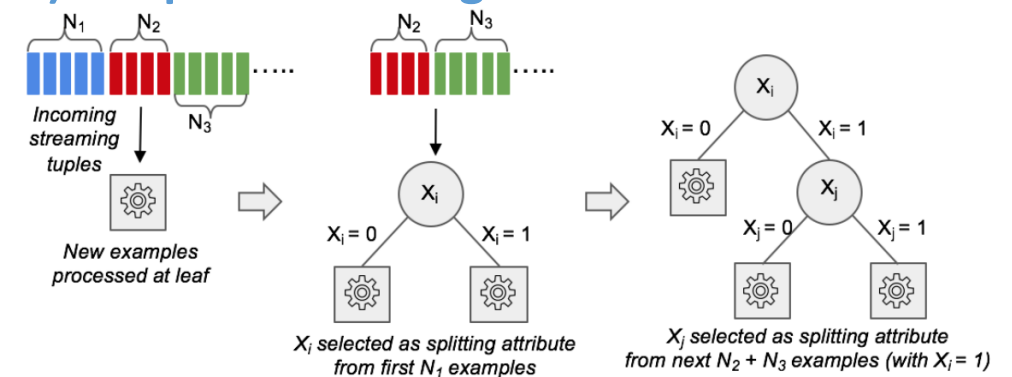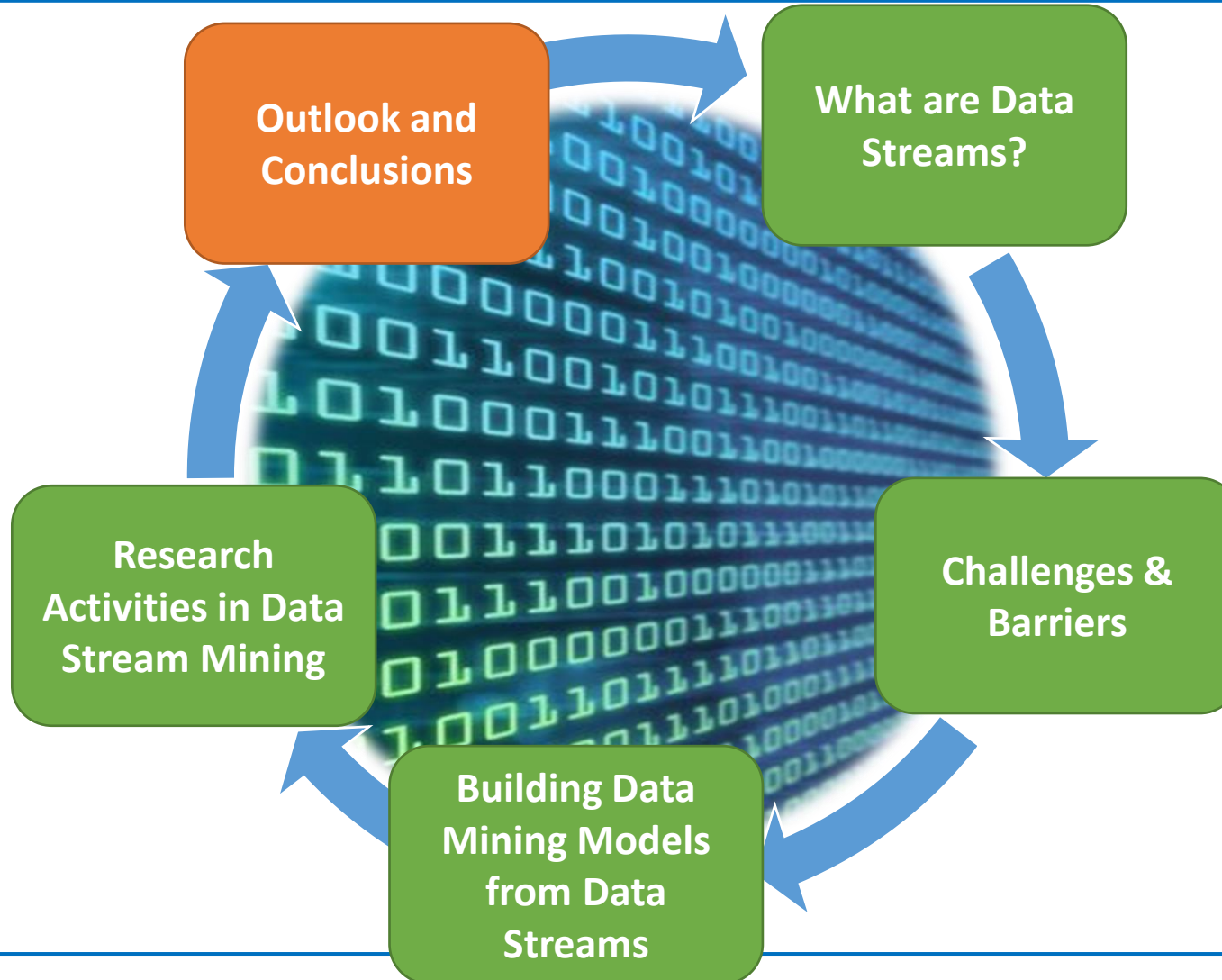# Methods: Windowing approaches to induce data mining models

## 2) Detect concept drift

## 1) Create time windows



Landmark window from **t₁**

Landmark window

Damped window

window at **t₁**    window at **t₂**

Sliding window

t₀  t₁  t₂  . . .                                          tₙ

**time t**



*Accuracy*

1
0.9
0.8
0.7
0.6
0.5

*Time*

Source: Stahl, F., Le, T., Badii, A., Gaber, M.M. (2021) A frequent pattern conjunction Heuristic for rule generation in data streams. Information 12(1) (2021), ISSN 2078-2489, doi: 10.3390/info12010024

# Methods: Windowing approaches to induce data mining models

## 1) Create time windows



Landmark window from $t_1$

Landmark window

Damped window

window at $t_1$    window at $t_2$

Sliding window

$t_0$  $t_1$  $t_2$  . . .    $t_n$

time t

Source: Stahl, F., Le, T., Badii, A., Gaber, M.M. (2021) A frequent pattern conjunction Heuristic for rule generation in data streams. Information 12(1) (2021), ISSN 2078-2489, doi: 10.3390/info12010024

## 2) Detect concept drift



Accuracy

Time

## 3) Adapt data mining model



$N_1$  $N_2$

Incoming streaming tuples

$N_3$

New examples processed at leaf

$X_i$ selected as splitting attribute from first $N_1$ examples

$N_2$  $N_3$

$X_i$

$X_i = 0$    $X_i = 1$

$X_j$ selected as splitting attribute from next $N_2 + N_3$ examples (with $X_i = 1$)

$X_i$

$X_i = 0$    $X_i = 1$

$X_j$

$X_j = 0$    $X_j = 1$

Source: Domingos and Hulten, 2000] Pedro M. Domingos and Geoff Hulten. Mining high-speed data streams. In SIGKDD, pages 71–80, 2000

# MC-NN: Micro Cluster based Neareast Neighbour, basic Approach

**Objective:** Develop a scalable predictive Data Stream classification

# MC-NN: Micro Cluster based Neareast Neighbour, basic Approach

**Objective:** Develop a scalable predictive Data Stream classification

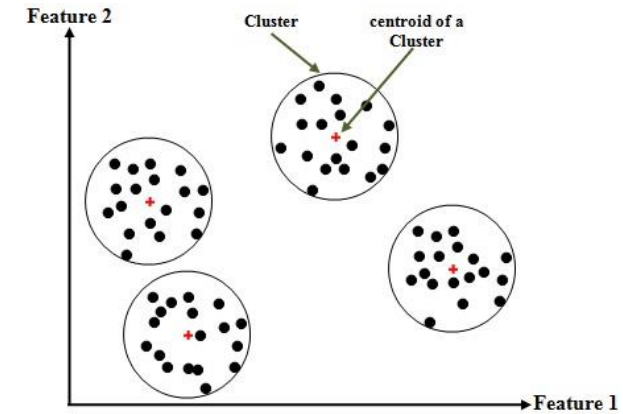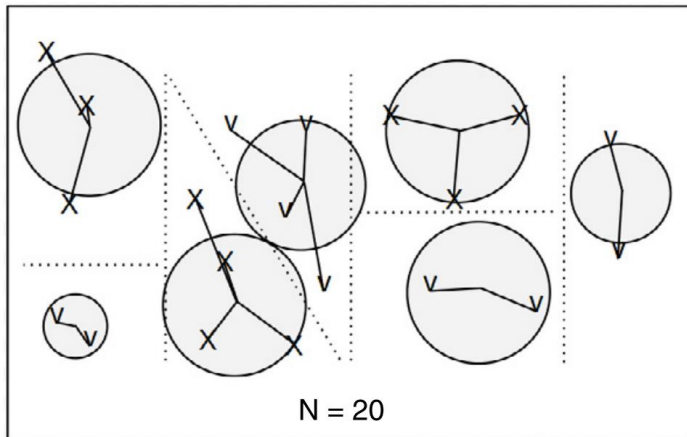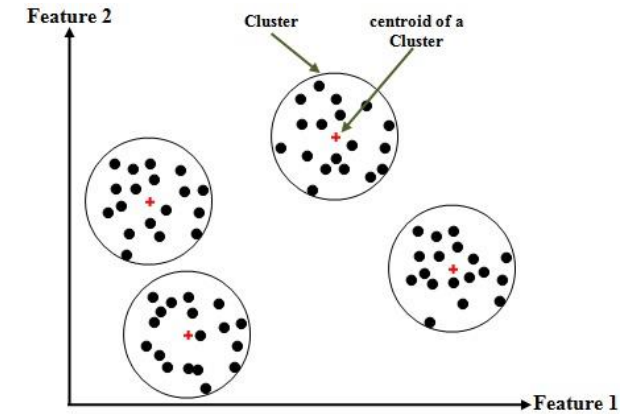**1) Initialising Micro-Clusters and maintenance statistics**
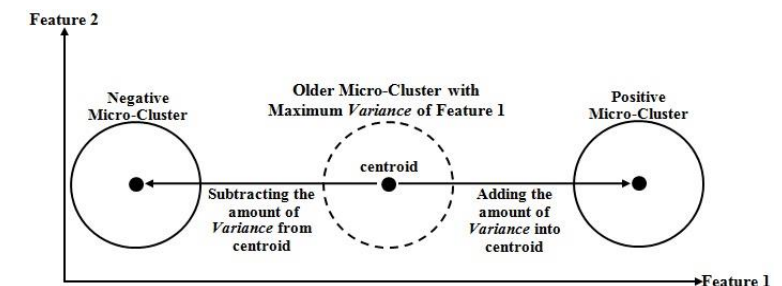


N = 20

$$< CF2^x, CF1^x, CF1^t, n, CL, \epsilon, \Theta, \alpha, \Omega >$$

$$centroid(x) = \frac{CF1^x}{n}$$

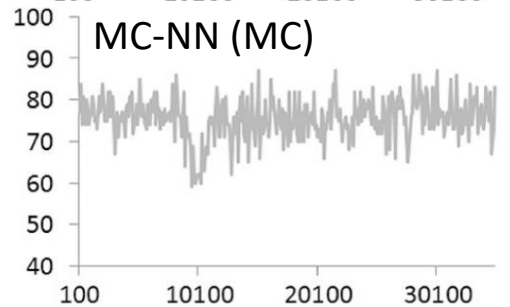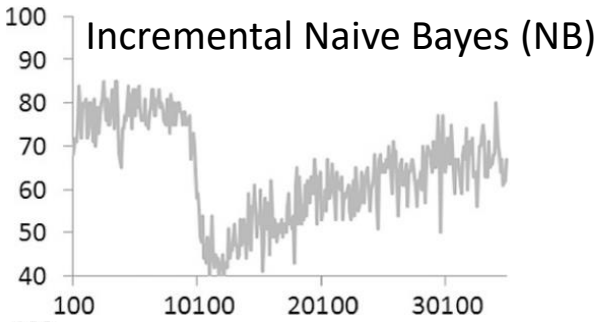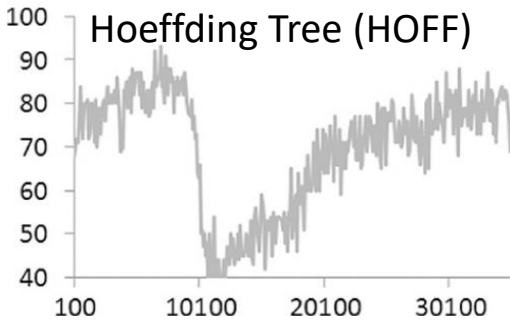$$Variance[x] = \sqrt{\left(\frac{CF2^x}{n}\right) - \left(\frac{CF1^x}{n}\right)^2}$$

- Initially a fixed number of Micro-Clusters is randomly initialised.

- Only components outlined in the table are stored.

- These can be used to calculate the clusters centroid and boundary (variance).

# MC-NN: Micro Cluster based Nearest Neighbour, basic Approach

**Objective:** Develop a scalable predictive Data Stream classification

**1) Initialising Micro-Clusters and maintenance statistics**



N = 20

$$< CF2^x, CF1^x, CF1^t, n, CL, \epsilon, \Theta, \alpha, \Omega >$$

$$centroid(x) = \frac{CF1^x}{n}$$

$$Variance[x] = \sqrt{\left(\frac{CF2^x}{n}\right) - \left(\frac{CF1^x}{n}\right)^2}$$

- Initially a fixed number of Micro-Clusters is randomly initialised.
- Only components outlined in the table are stored.
- These can be used to calculate the clusters centroid and boundary (variance).

**2) Absorbing new data stream instances**

# MC-NN: Micro Cluster based Neareast Neighbour, basic Approach

**Deutsches Forschungszentrum für Künstliche Intelligenz GmbH**

**Objective:** Develop a scalable predictive Data Stream classification

## 1) Initialising Micro-Clusters and maintenance statistics



N = 20

$$< CF2^x, CF1^x, CF1^t, n, CL, \epsilon, \Theta, \alpha, \Omega >$$

$$centroid(x) = \frac{CF1^x}{n}$$

$$Variance[x] = \sqrt{\left(\frac{CF2^x}{n}\right) - \left(\frac{CF1^x}{n}\right)^2}$$

- Initially a fixed number of Micro-Clusters is randomly initialised.
- Only components outlined in the table are stored.
- These can be used to calculate the clusters centroid and boundary (variance).

## 2) Absorbing new data stream instances

**EPSRC** Engineering and Physical Sciences Research Council



## 3 ) Splitting and removing of Micro-Clusters

# MC-NN: Results

**Adaptating to Concept Drift**


Hoeffding Tree (HOFF)


Incremental Naive Bayes (NB)


MC-NN (MC)

# MC-NN: Results
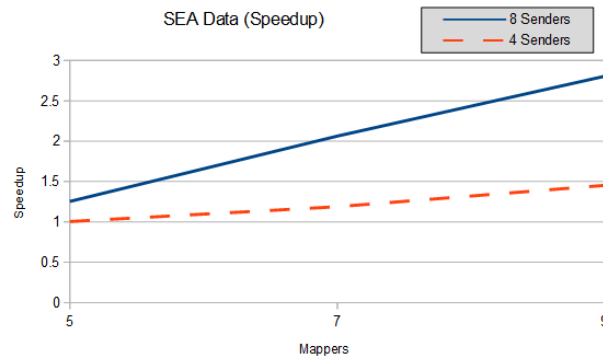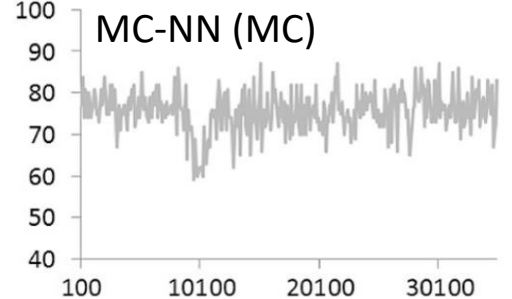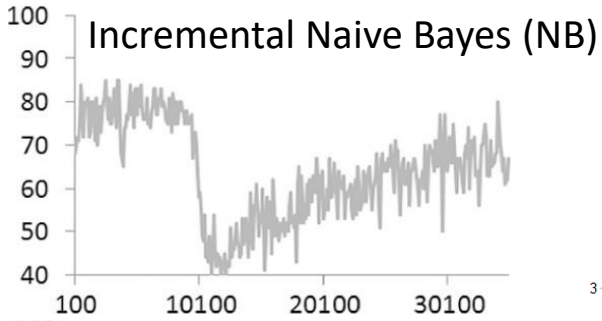
**Adaptating to Concept Drift**
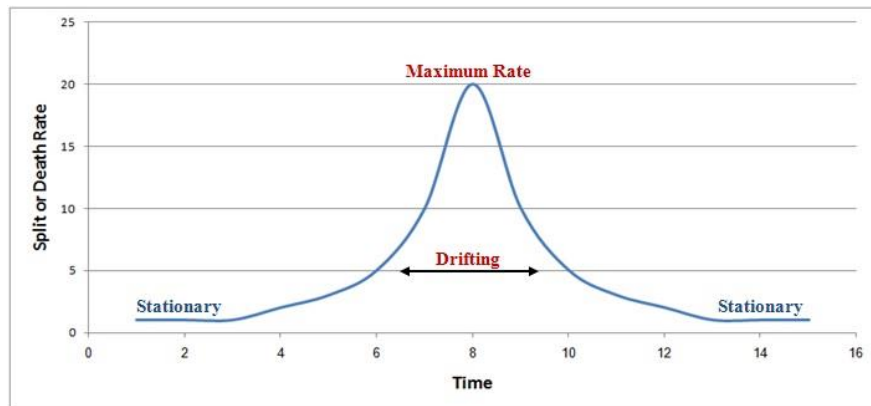
**Scalability through parallelisation**



Hoeffding Tree (HOFF)

Incremental Naive Bayes (NB)

MC-NN (MC)

**Deutsches Forschungszentrum für Künstliche Intelligenz GmbH**

## Adaptating to Concept Drift

Hoeffding Tree (HOFF)

Incremental Naive Bayes (NB)

MC-NN (MC)

## Scalability through parallelisation

Hadoop Node

SAMZA Container

KAFKA Input Stream → MCNN Classifier → KAFKA Output Stream

SEA Data (Speedup)
- 8 Senders
- 4 Senders

Speedup HYP
- 8 Senders
- 4 Senders

## Long term adaptation

**EPSRC** Engineering and Physical Sciences cil

Accuracy vs Time 10,000,000 Instances 10% Noise

| | |
|---|---|
| NB | (10.12) |
| HOFF | (29.44) |
| KNN 250 | (311.0) |
| MC2 | (11.7) |
| MC10 | (9.72) |

Accuracy vs Time 10,000,000 Instances 10% Noise

| | |
|---|---|
| NB | (10.17) |
| HOFF | (27.1) |
| KNN 250 | (332.0) |
| MC2 | (11.61) |
| MC10 | (9.91) |

# Using MC-NN for Explaining Concept Drift through Feature Tracking
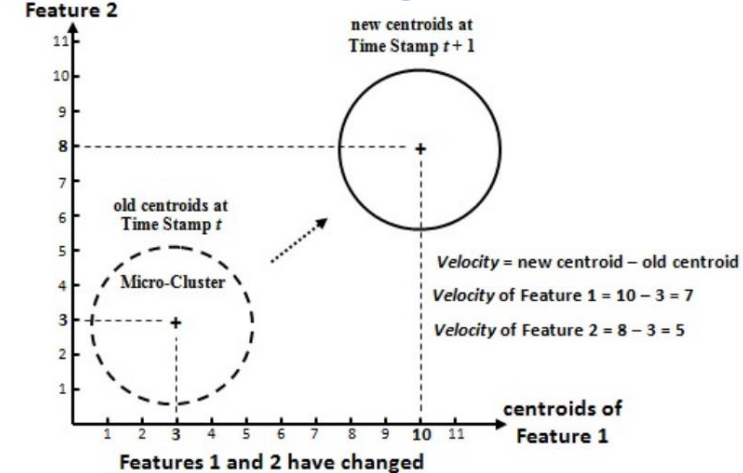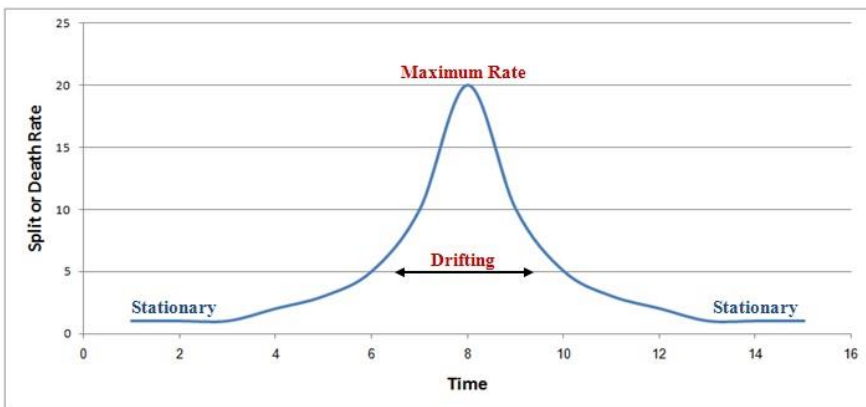
**Measuring centroid velocity**



**Measuring split & death rate**

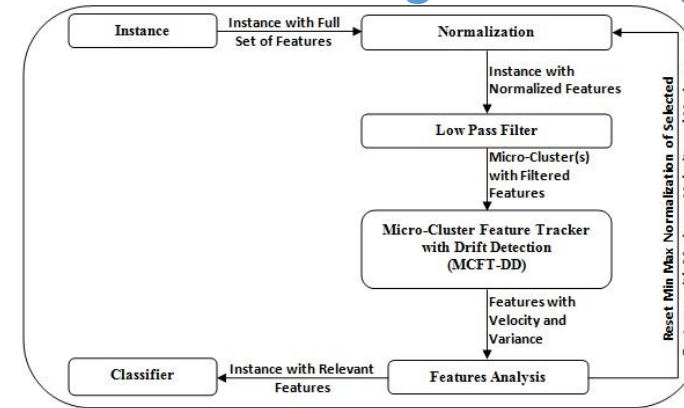# Using MC-NN for Explaining Concept Drift through Feature Tracking



**Measuring centroid velocity**

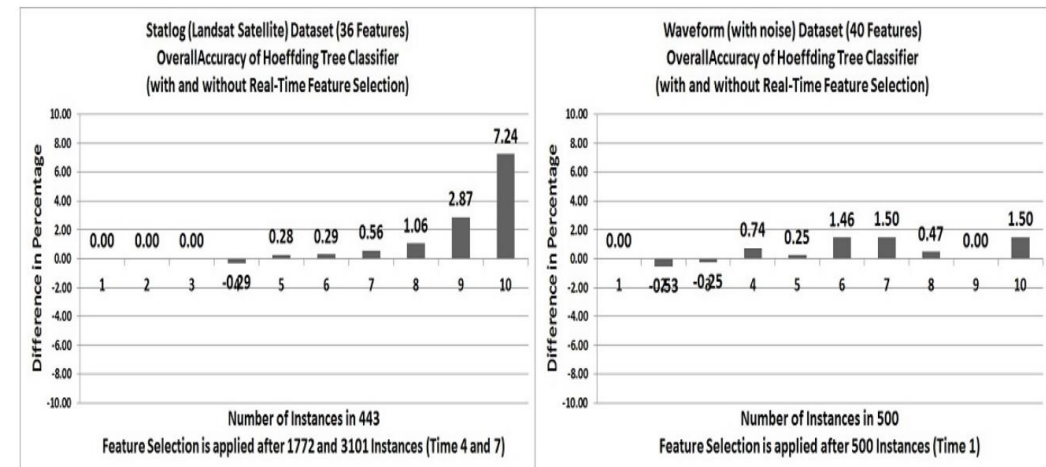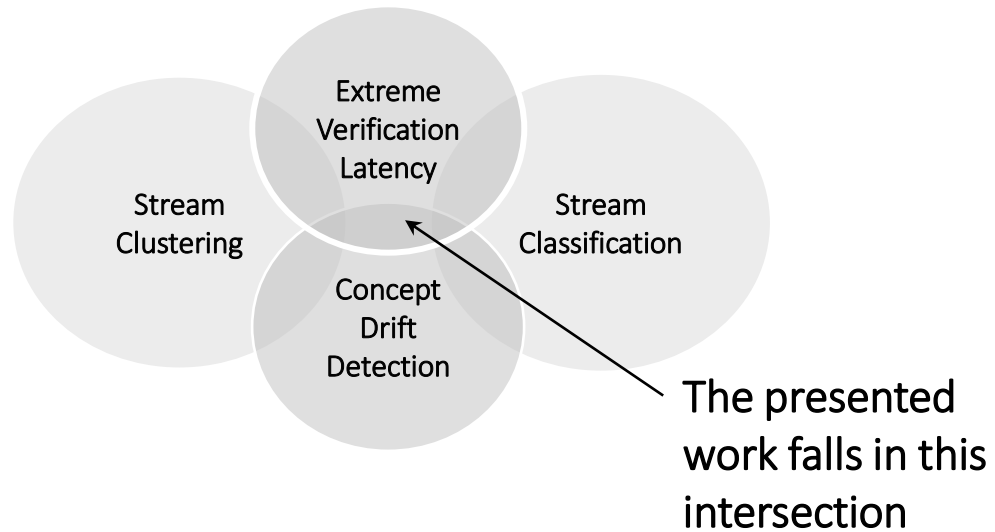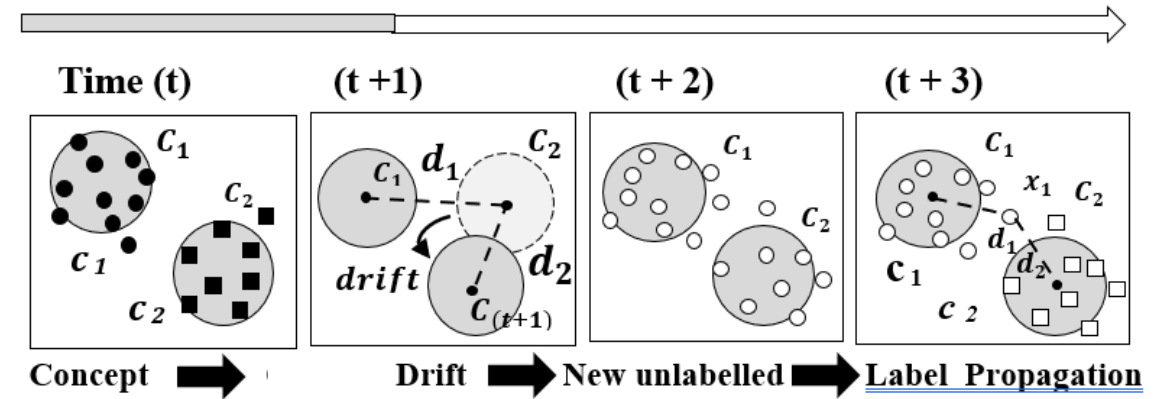**Measuring split & death rate**

**Feature tracking and ranking**
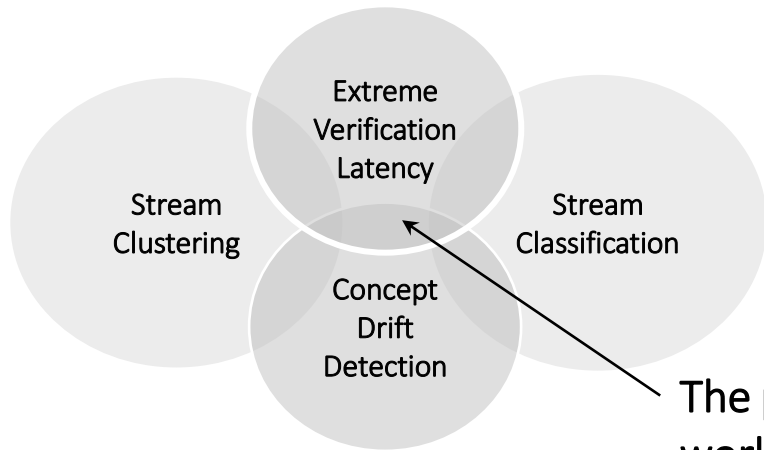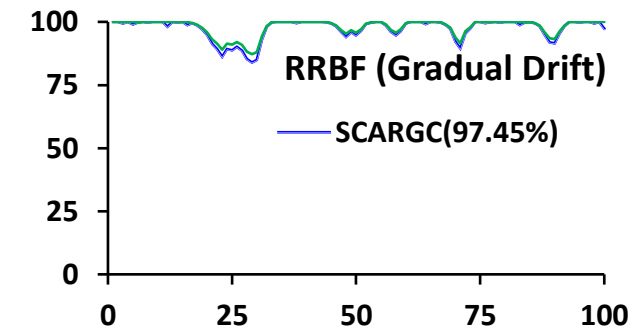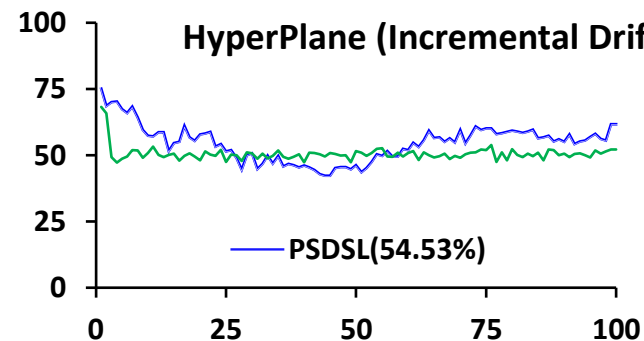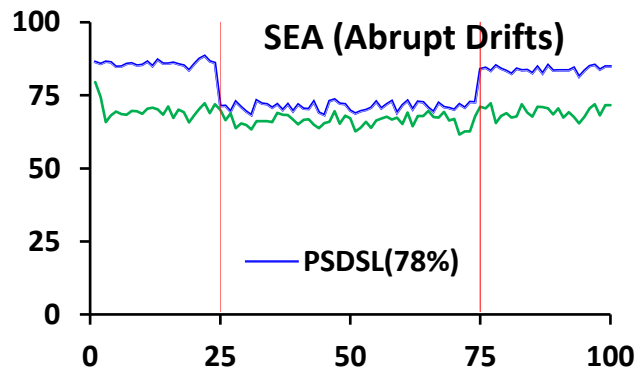
**Results**

# MC-NN: Unsupervised Classification



The presented work falls in this intersection

# MC-NN: Unsupervised Classification

Stream Clustering

Extreme Verification Latency

Stream Classification

Concept Drift Detection

The presented work falls in this intersection

**Data Stream** (Few labelled)

| Time (t) | (t +1) | (t + 2) | (t + 3) |

$c_1$ · $c_2$ · $c_1$ · $c_2$

$c_1$ · $d_1$ · $c_2$ · $drift$ · $d_2$ · $C_{(t+1)}$

$c_1$ · $c_2$

$c_1$ · $x_1$ · $c_2$ · $d_1$ · $d_2$ · $c_1$ · $c_2$
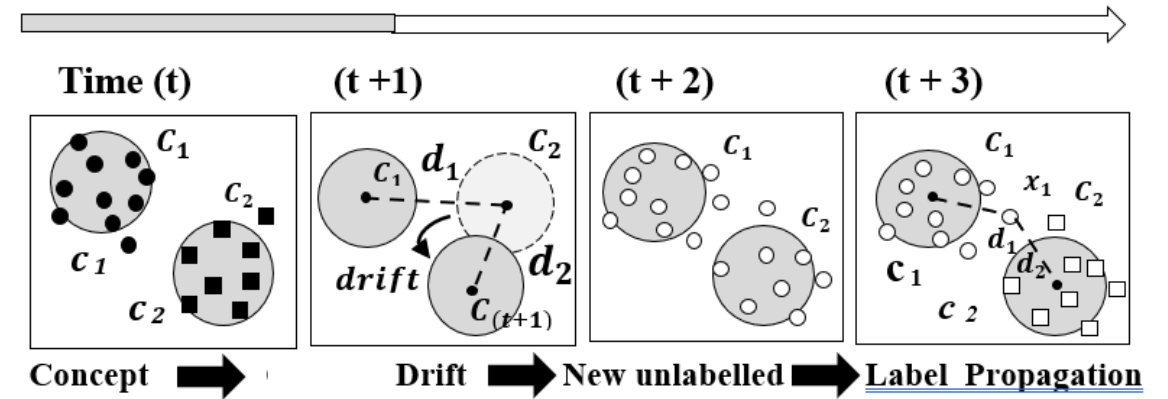
**Concept** ➤ **Drift** ➤ **New unlabelled** ➤ **Label Propagation**

# MC-NN: Unsupervised Classification



The presented work falls in this intersection

**Data Stream** (Few labelled)




SEA (Abrupt Drifts)
PSDSL(78%)


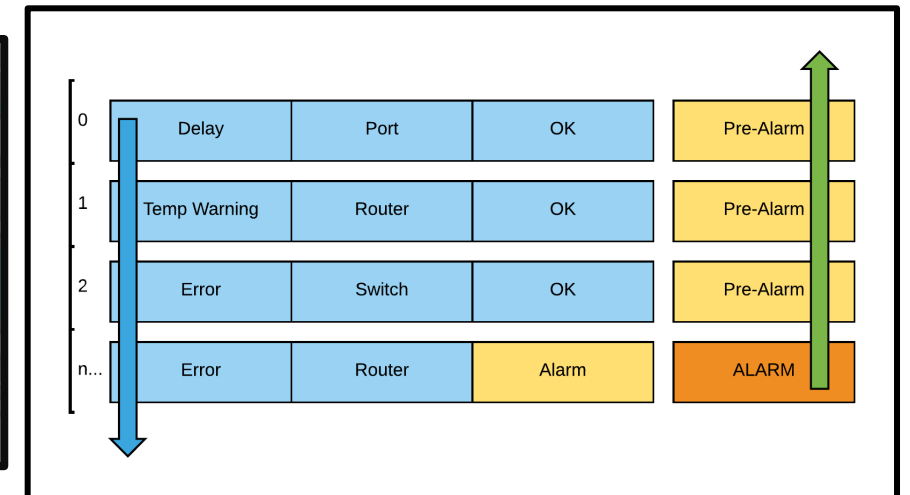HyperPlane (Incremental Drift)
PSDSL(54.53%)


RRBF (Gradual Drift)
SCARGC(97.45%)
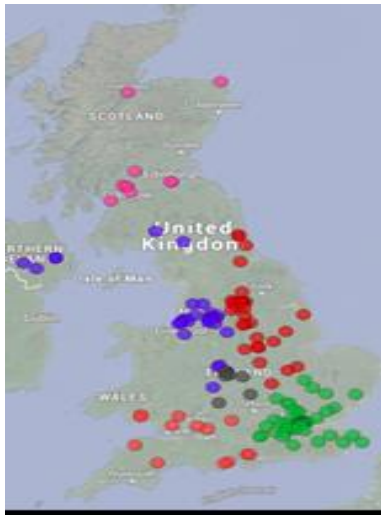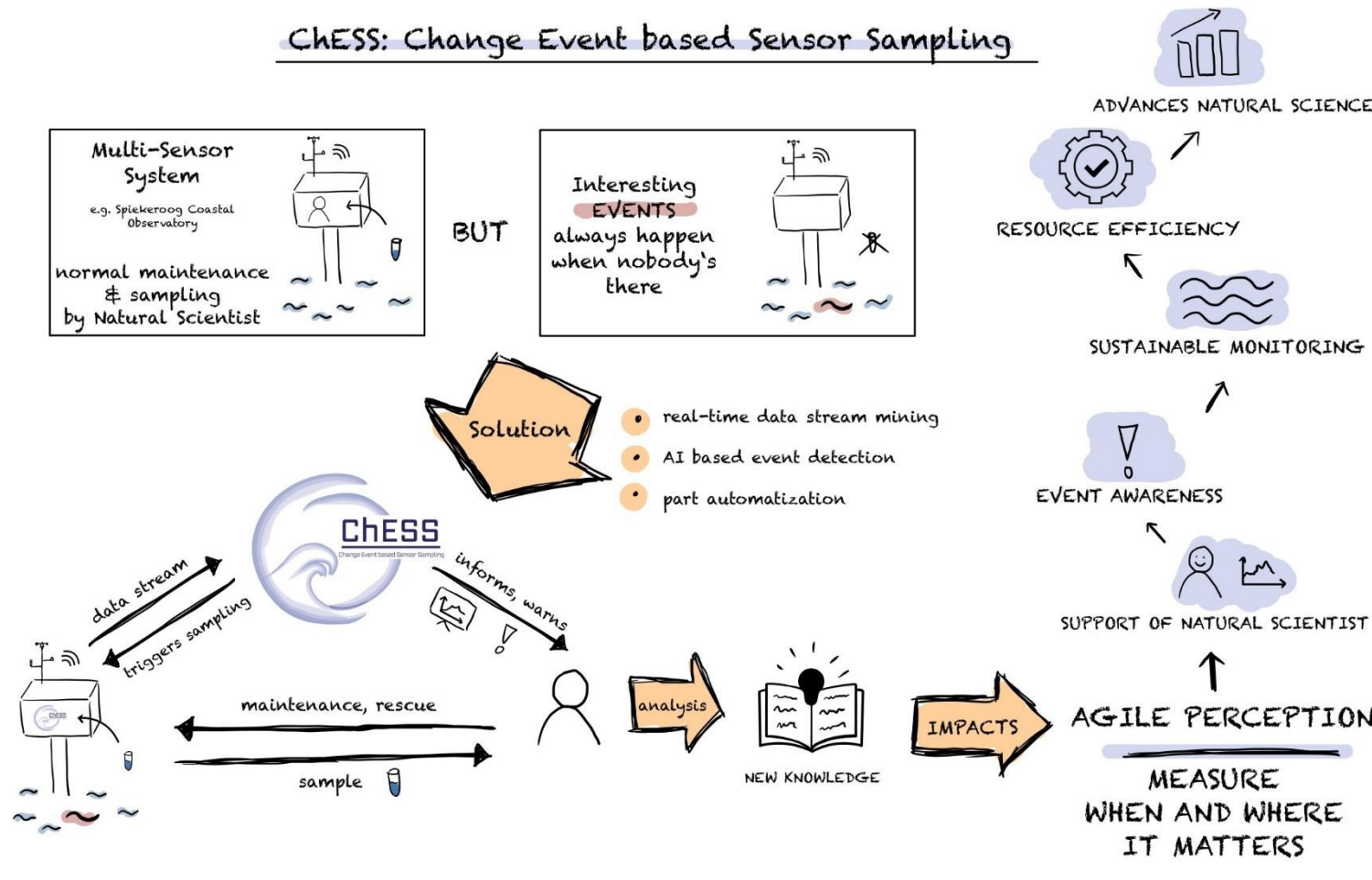
# Applications: BT (completed)

**Example Problem: Real-time Network Alarm Forecasting**

- Increasing reliance on Telecommunication services for business and personal use

- Telecommunication Networks have a great deal of redundancy (99.999% availability), however, the "last mile" is often a single point of failure

- Network devices emit different events data at different frequencies under different conditions. Yet they may be linked.
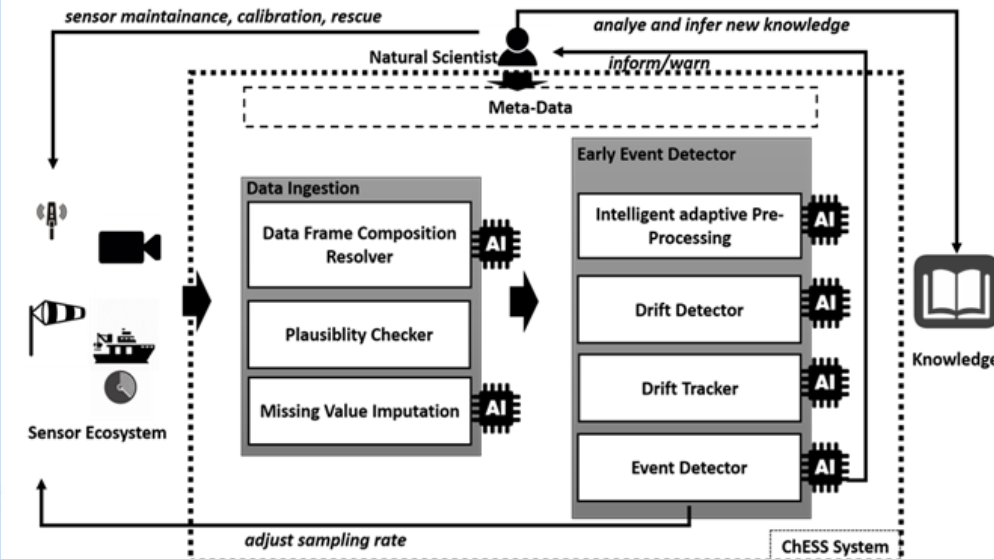
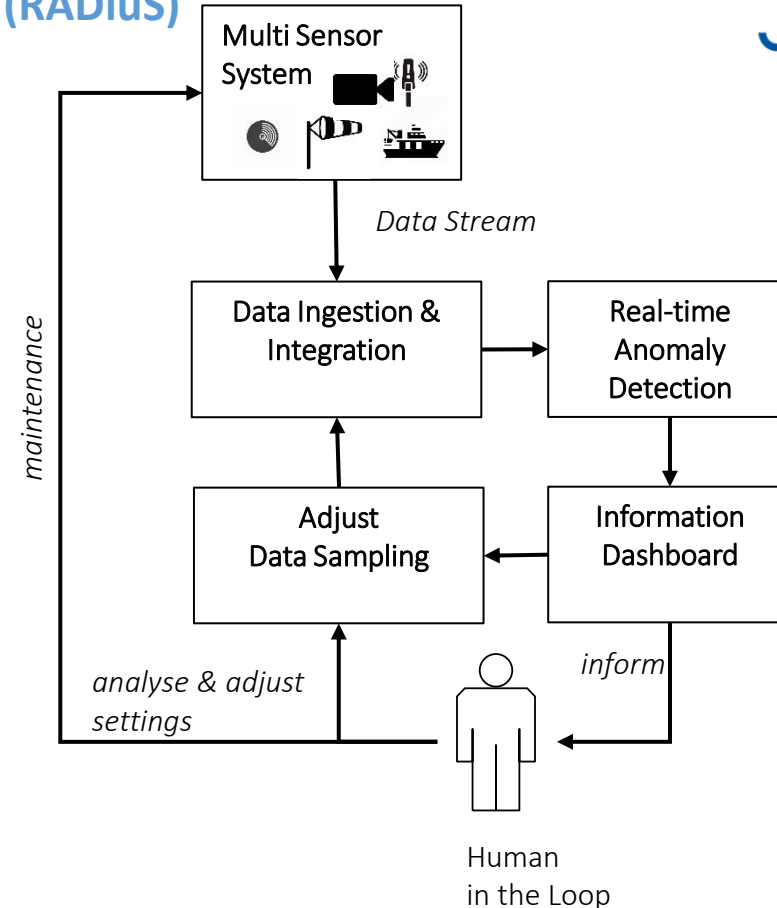# Systems Development: ChESS (ongoing)

# Applications: Intelligent Maintenance of Costal Environments (just starting)
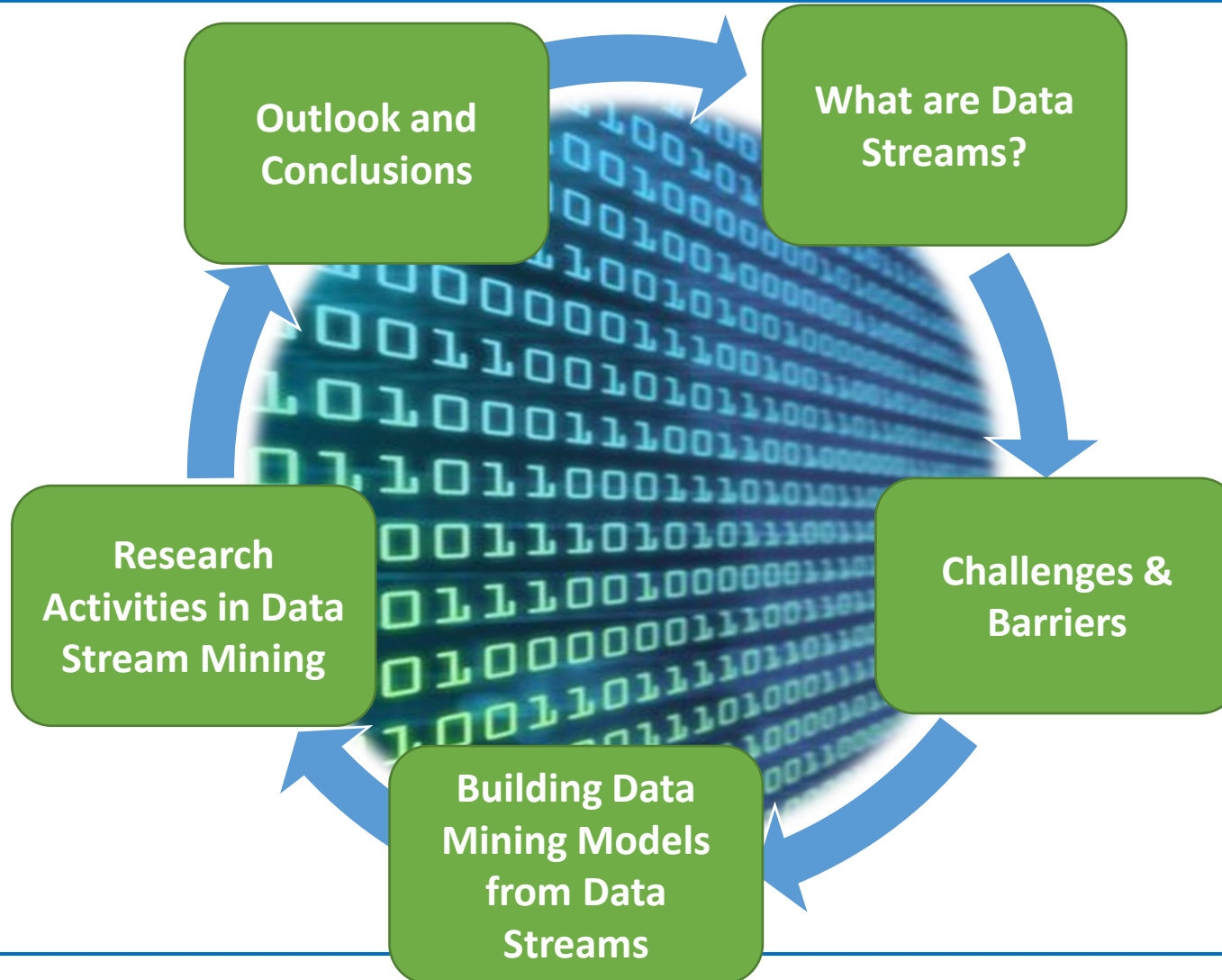
**Ad-hoc data acquisition mesh for enhanced versatile explorations of waters**



@ Prof. Jan Schulz

**Real-time anomaly detection in aquatic system (RADiuS)**

# The Data Tsunami



## Challenges

1) Data generated at a fast rate (Velocity), at potentially large and unknown quantities (Volume)
2) Concept Drift (changes of pattern encoded in in the data over time)
3) Modelling real-time analytics workflows from streaming data
4) Multi-modality of data sources (text, video/images, unstructured)
5) Class label sparsity: adapting predictive models
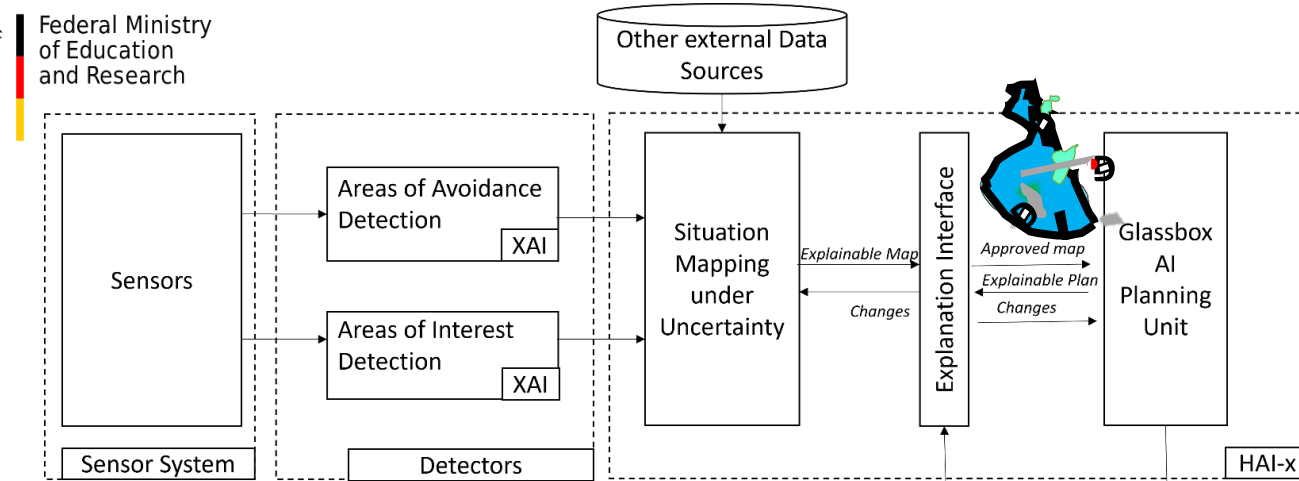6) Explaining Concept Drift

## Barriers

1) Limited scalable (parallel) real-time high throughput data stream mining algorithms
2) Different and changing types of concept drift
3) Lack of customisable pre-processing techniques
4) Different time stamps but co-occurring data items
5) Supervised algorithms not applicable in many cases
6) Lack of drift detectors explaining concept drift

Deutsches Forschungszentrum für Künstliche Intelligenz GmbH

## Hybride AI Explainer (HAI-x): weed harvester scenario
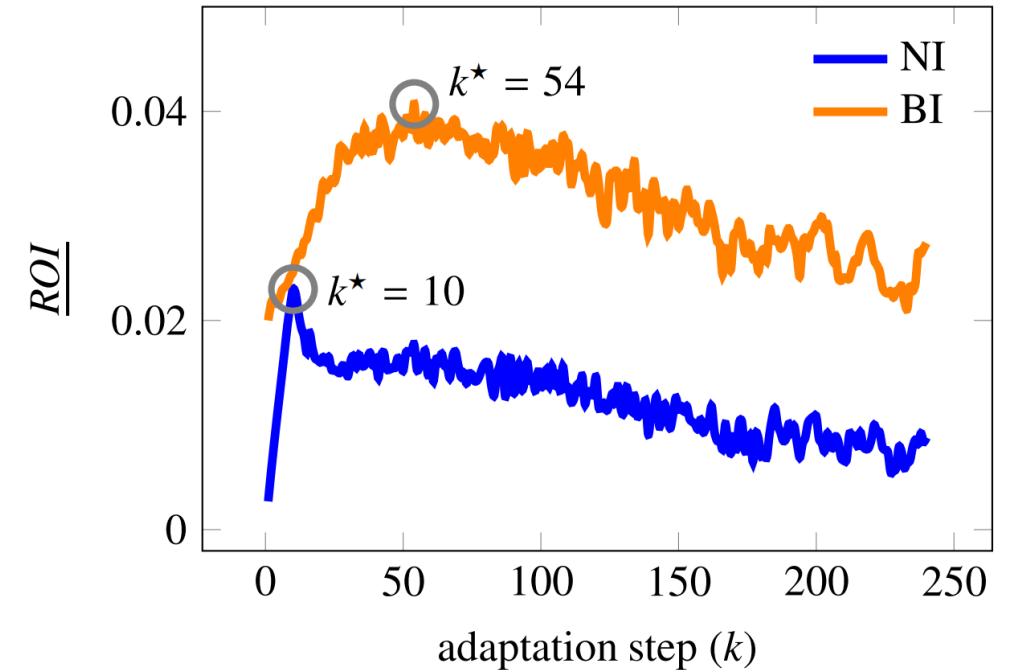


## Utility of Adaptation: when is it worth updating your model?

<u>ROI</u> is return on employing an adaptive predictor as compared to keeping a fixed nonadaptive model

# Thank you!



- Marine Perception Team at DFKI
- Prof. Oliver Zielinski
- Prof. Lars Nolle
- The team of the Department of Computer Science at the University of Reading
- Prof. Giuseppe Di Fatta
- UK and ISG Teams
- Prof. Mohamed Gaber
- My wife Laura
- And others.

# Questions?