

BIG DATA STREAM MINING WITH APACHE SAMOA

Albert Bifet @abifet
DBWeb Team
INFRES Department

10 November 2015



MOTIVATION

REAL TIME ANALYTICS

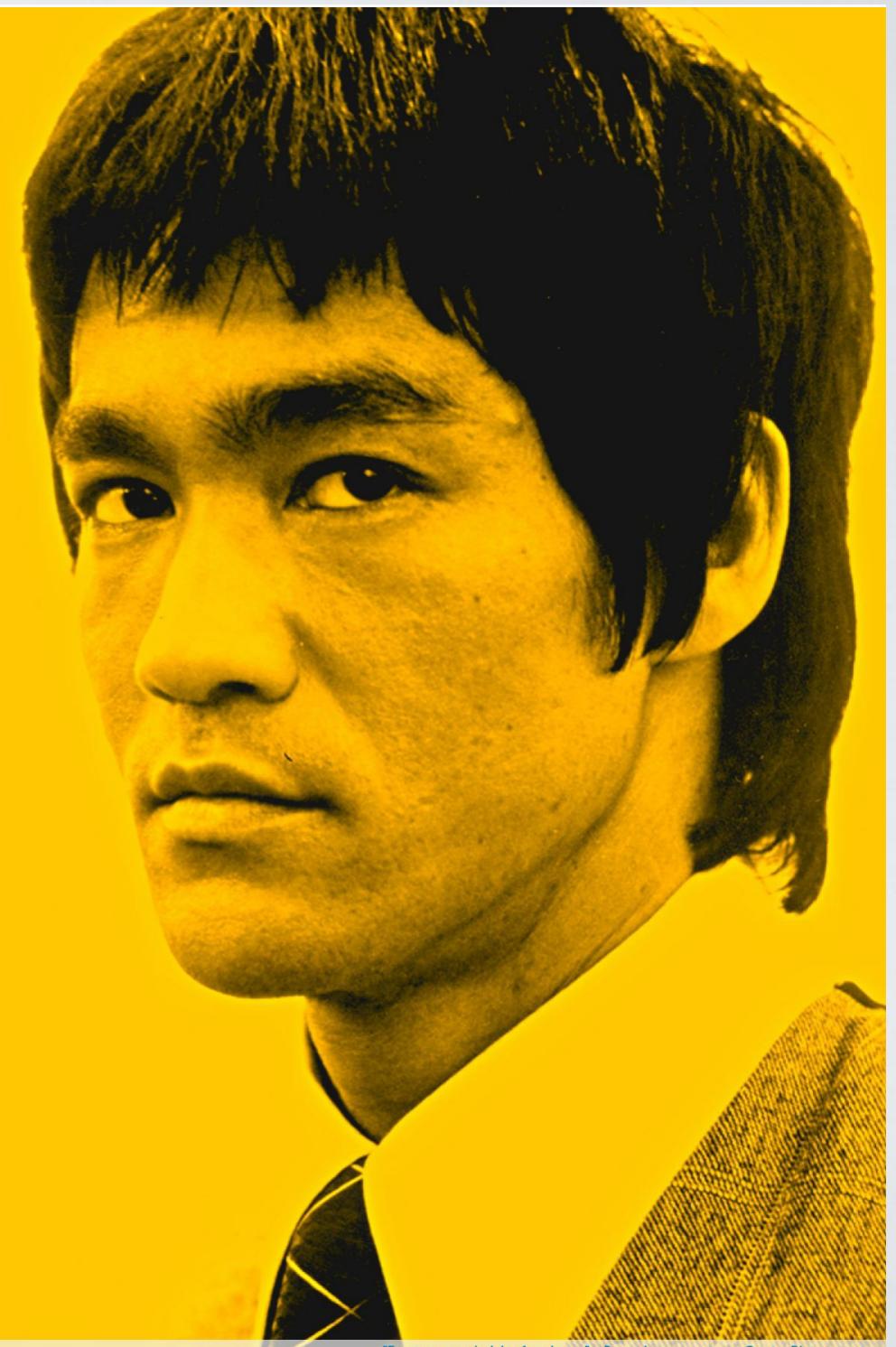


REAL TIME ANALYTICS

**Empty your mind, be formless.
Shapeless, like water.
If you put water into a cup,
it becomes the cup.
You put water into a bottle
and it becomes the bottle.
You put it in a teapot,
it becomes the teapot.
Now, water can flow or
it can crash.**

李小龍

Bruce Lee

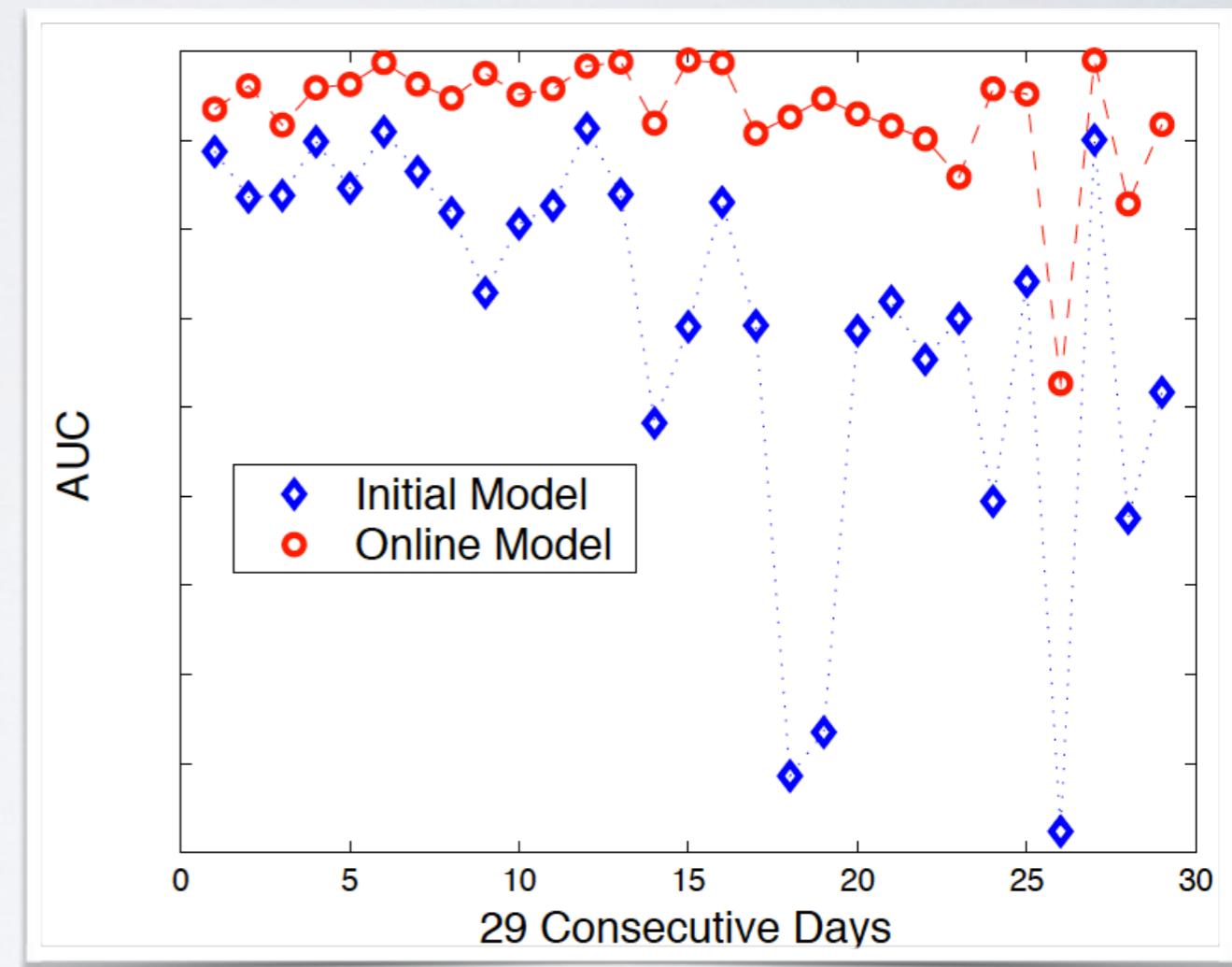


APACHE SA(MOA) VISION

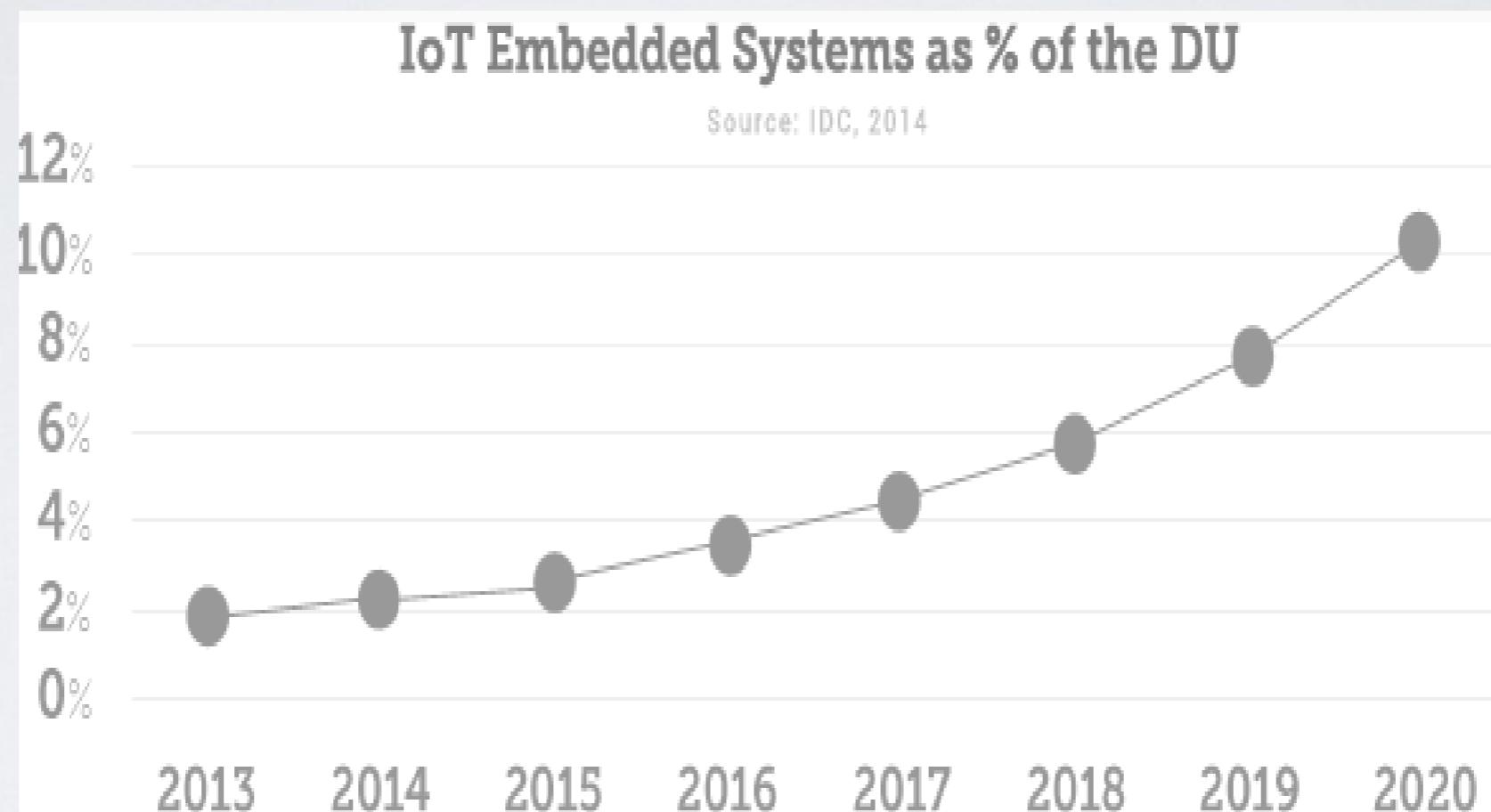
- Data Stream mining platform
 - Library of state-of-the-art algorithms for practitioners
 - Development and collaboration framework for researchers
- Algorithms & Systems

IMPORTANCE

- Example: spam detection in comments on Yahoo News
- Trends change in time
- Need to retrain model with new data



INTERNET OF THINGS



- EMC Digital Universe, 2014

BIG DATA STREAM

- Volume + Velocity (+ Variety)
- Too large for single commodity server main memory
- Too fast for single commodity server CPU
- A solution should be:
 - Distributed
 - Scalable



BIG DATA PROCESSING ENGINES

- Low latency



S4 distributed stream
computing platform

APACHE
STORM™
Distributed • Resilient • Real-time

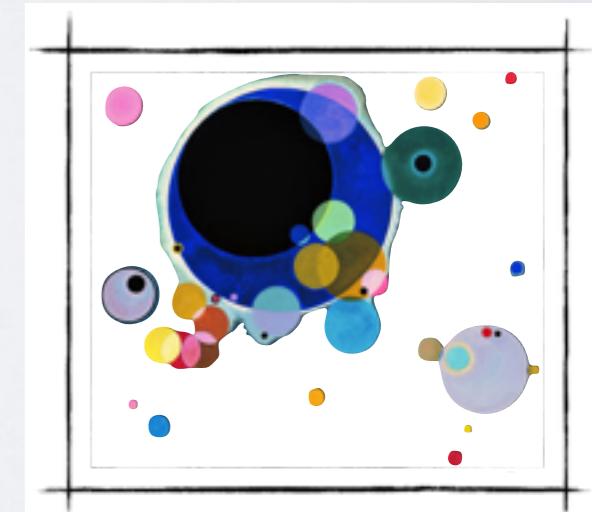
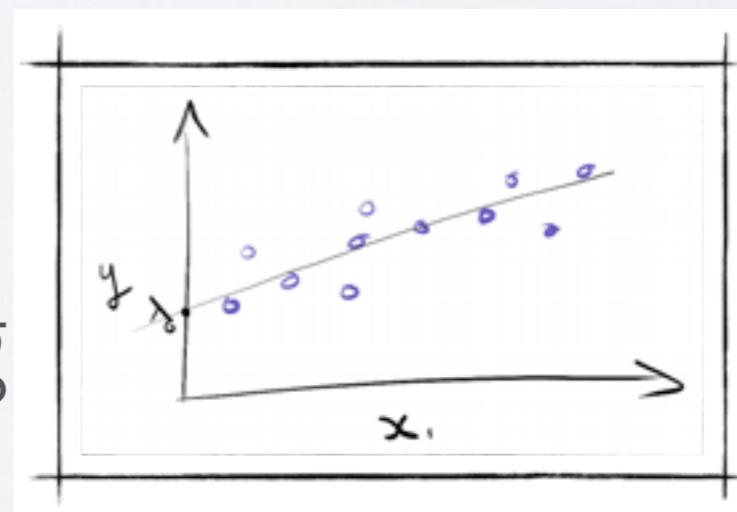
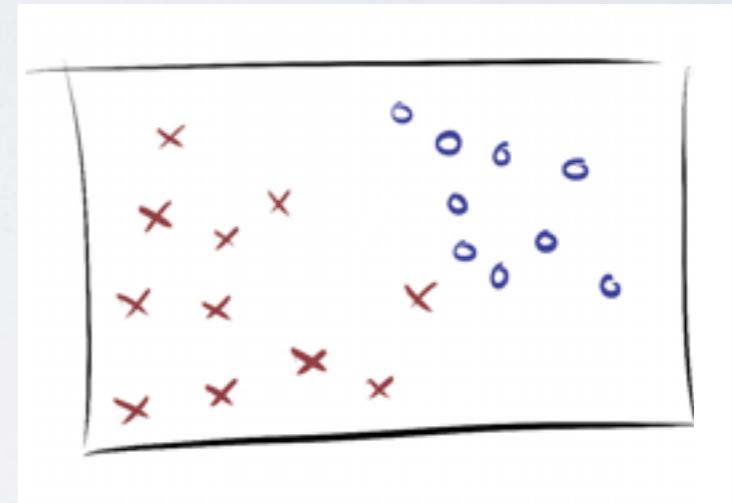
samza

- High Latency (Not real time)

Spark Streaming

MACHINE LEARNING

- Classification
- Regression
- Clustering
- Frequent Pattern Mining



WHAT IS MOA?

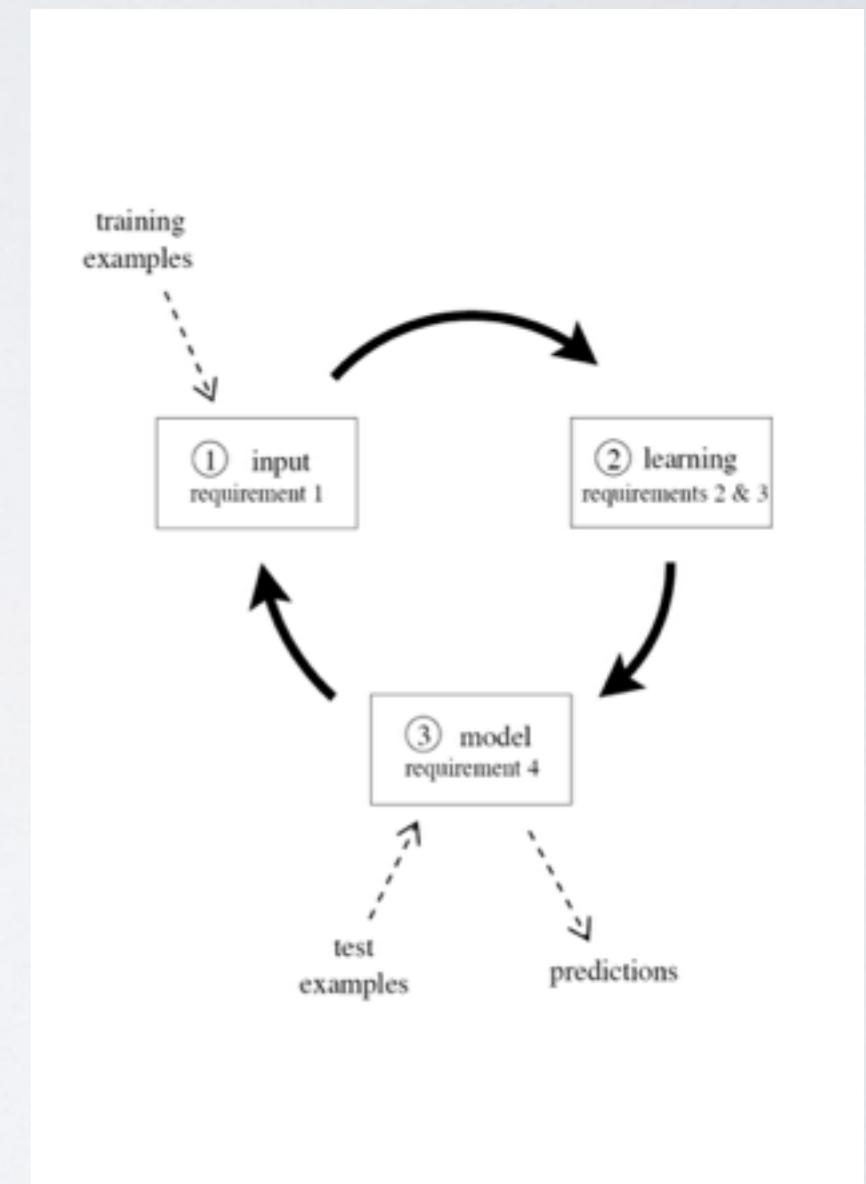
MOA

- {M}assive {O}nline {A}nalysis is a framework for online learning from data streams.
- It is closely related to WEKA
- It includes a collection of offline and online as well as tools for evaluation:
 - classification, regression
 - clustering, frequent pattern mining
- Easy to extend, design and run experiments



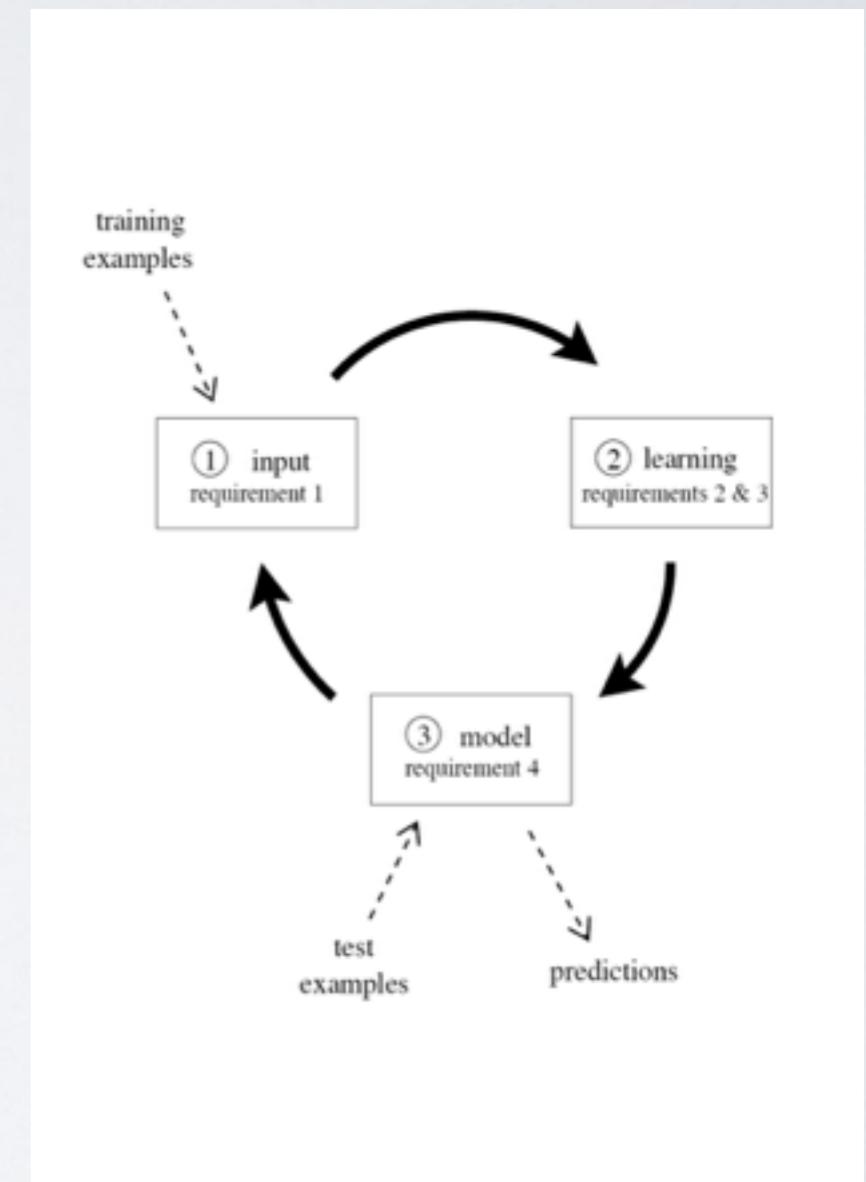
STREAM SETTING

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



STREAM EVALUATION

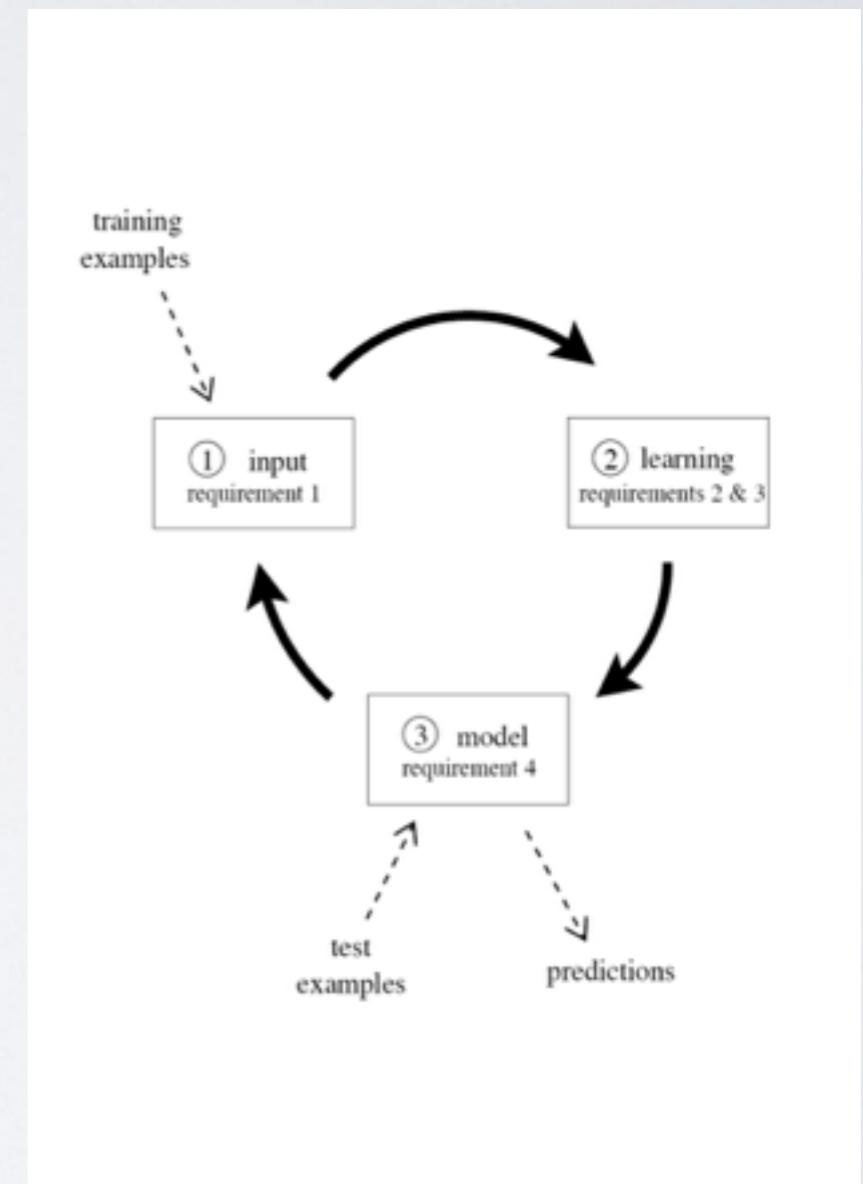
- Holdout Evaluation
- Interleaved Test-Then-Train or Prequential



STREAM EVALUATION

Holdout an independent test set

- Apply the current decision model to the test set, at regular time intervals
- The loss estimated in the holdout is an unbiased estimator

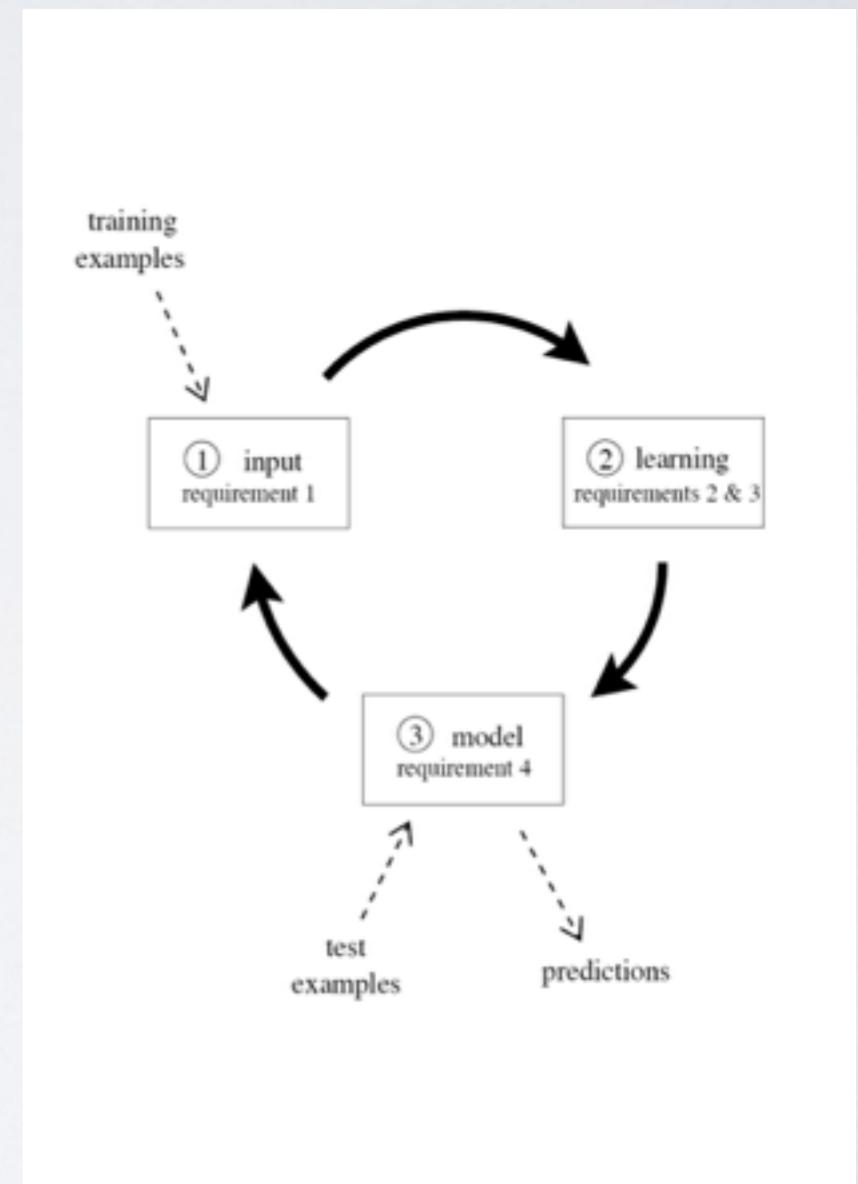


STREAM EVALUATION

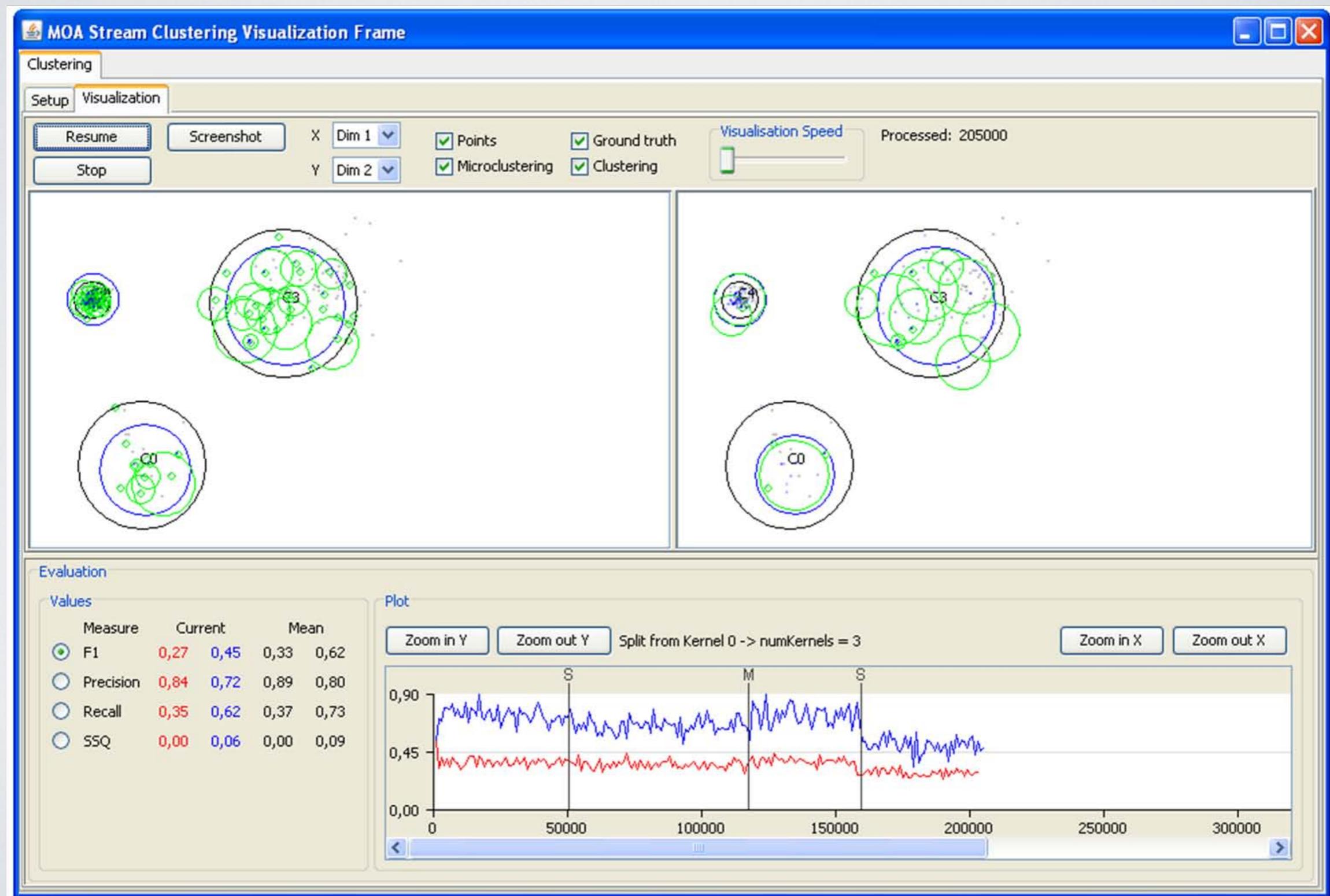
Prequential Evaluation

- The error of a model is computed from the sequence of examples.
- For each example in the stream, the actual model makes a prediction based only on the example attribute-values.

$$S = \sum_{i=1}^n L(y_i, \hat{y}_i).$$



CLUSTERING



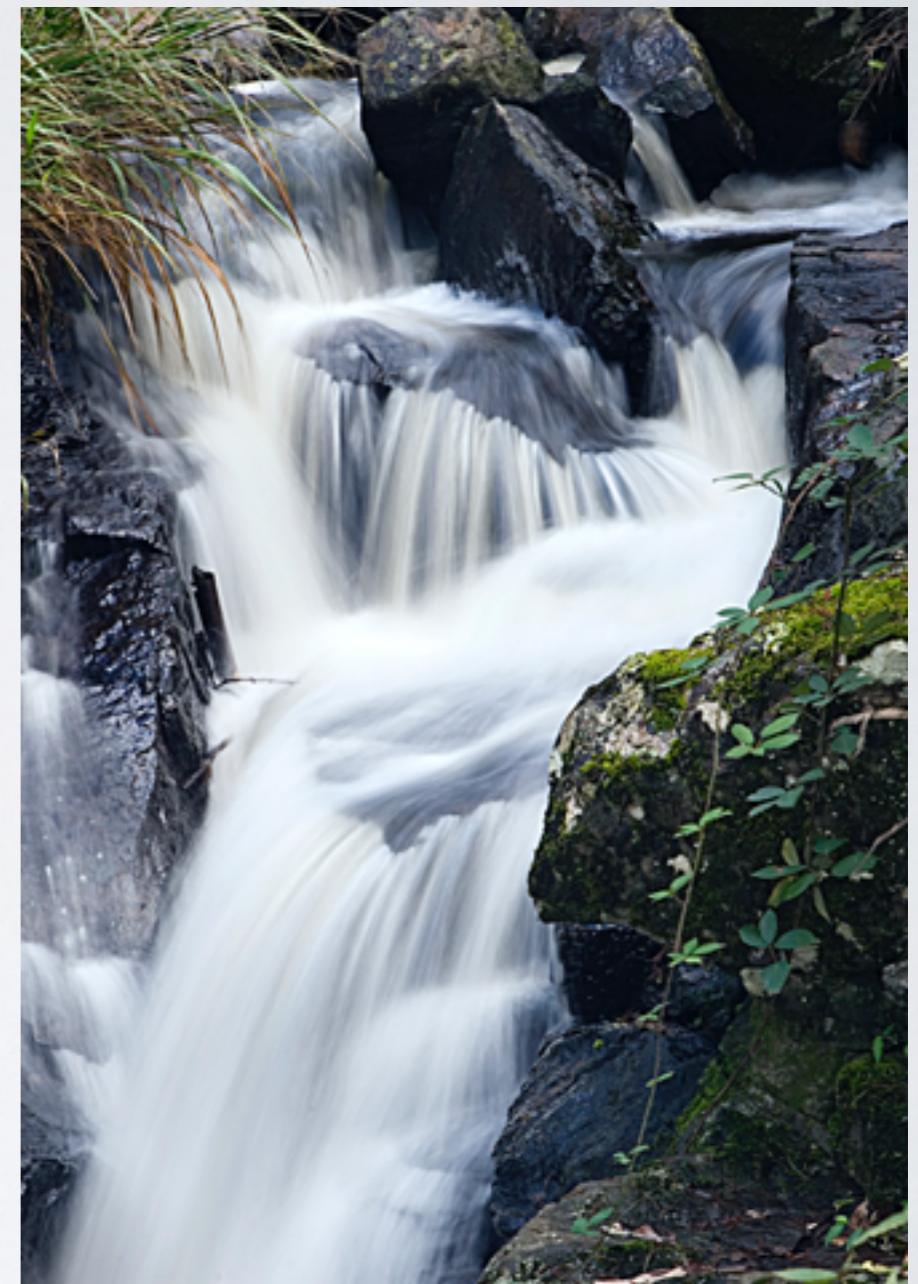
COMMAND LINE

- ```
java -cp .:moa.jar:weka.jar -javaagent:sizeofag.jar moa.DoTask "EvaluatePeriodicHoldoutTest -l DecisionStump -s generators.WaveformGenerator -n 100000 -i 100000000 -f 1000000" > dsresult.csv
```
- This command creates a comma separated values file:
  - training the DecisionStump classifier on the WaveformGenerator data,
  - using the first 100 thousand examples for testing,
  - training on a total of 100 million examples,
  - and testing every one million examples

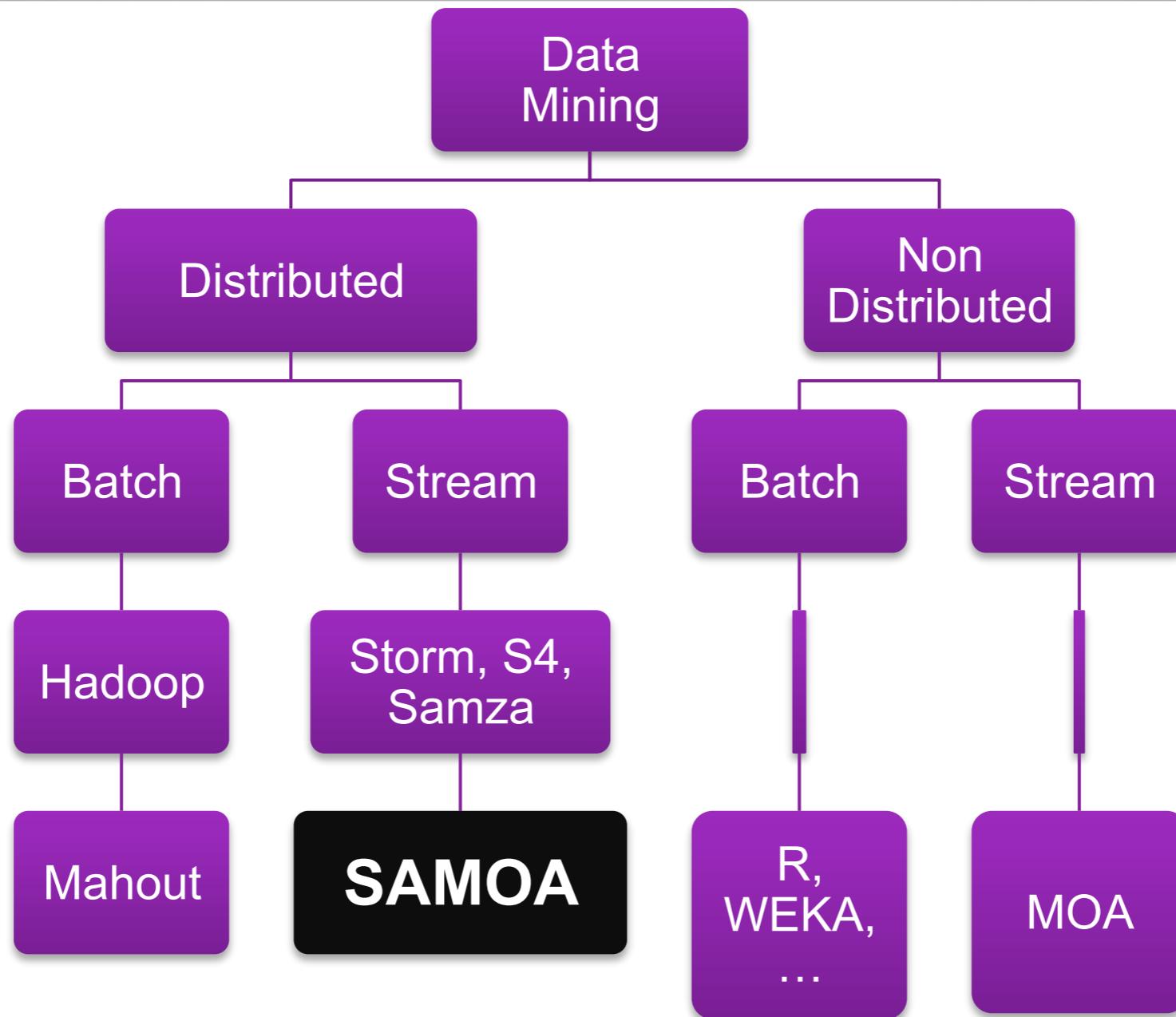
# WHAT IS APACHE SAMOA?

# STREAMING MODEL

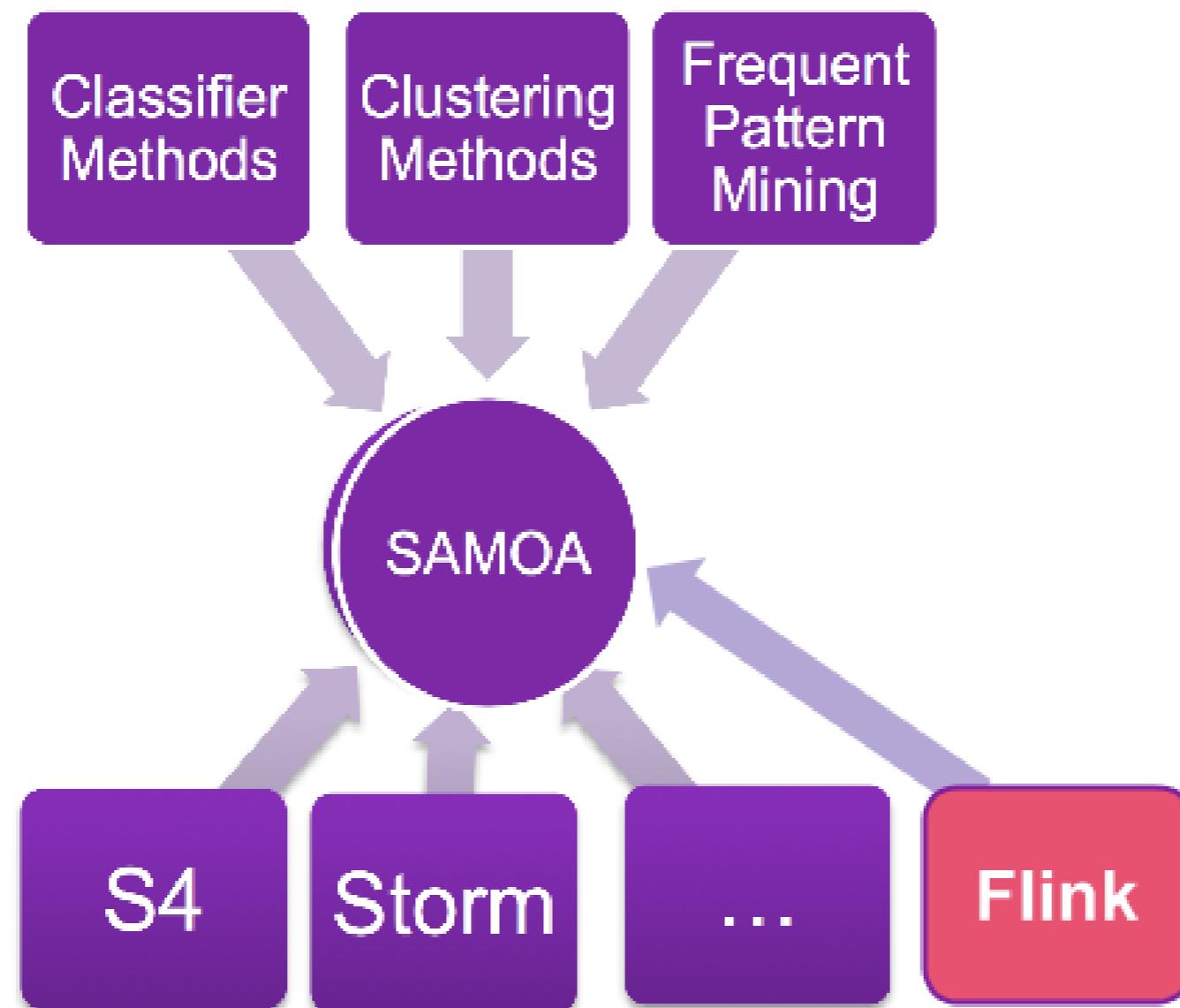
- Sequence is potentially infinite
- High amount of data, high speed of arrival
- Change over time (concept drift)
- Approximation algorithms  
(small error with high probability)
  - Single pass, one data item at a time
  - Sub-linear space and time per data item



# TAXONOMY

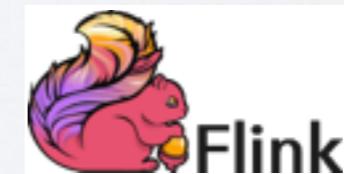
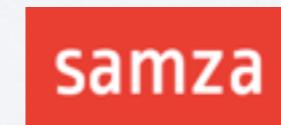
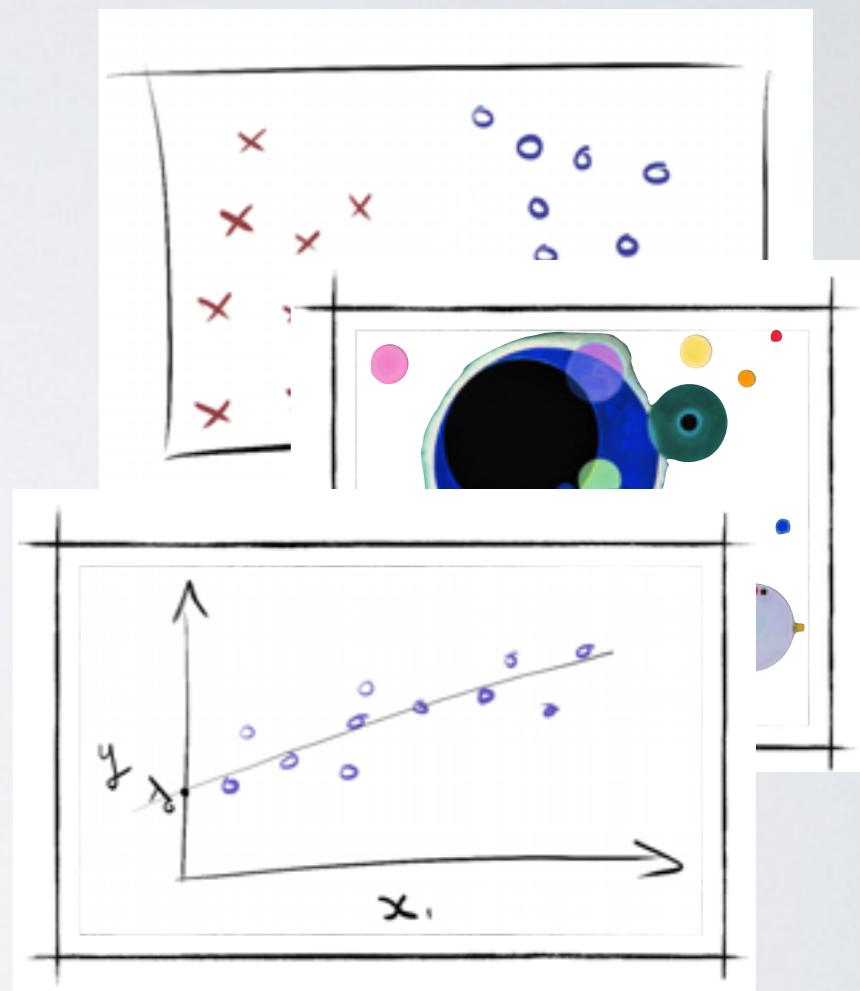


# ARCHITECTURE



# STATUS

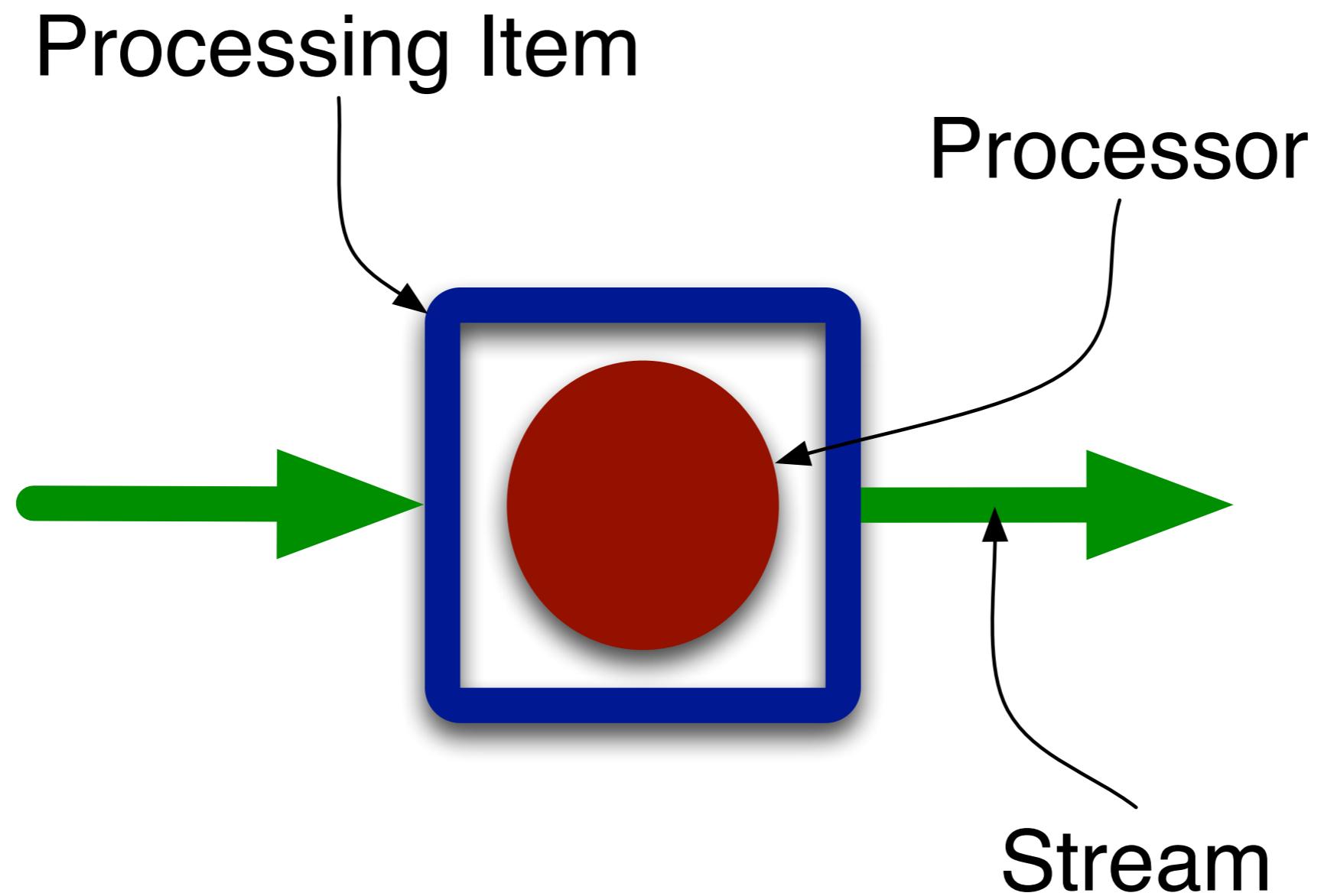
- Parallel algorithms
  - Classification (Vertical Hoeffding Tree)
  - Clustering (CluStream)
  - Regression (Adaptive Model Rules)
- Execution engines



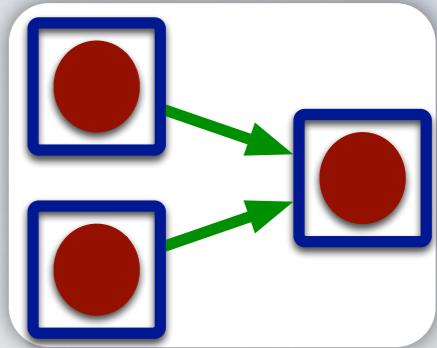
# IS SAMOA USEFUL FOR YOU?

- Only if you need to deal with:
  - **Large** fast data
  - Evolving process (model updates)
- What is happening now?
  - Use feedback in real-time
  - Adapt to changes faster

# ML DEVELOPER API



# ML DEVELOPER API



```
TopologyBuilder builder;
Processor sourceOne = new SourceProcessor();
builder.addProcessor(sourceOne);
Stream streamOne = builder.createStream(sourceOne);
```

```
Processor sourceTwo = new SourceProcessor();
builder.addProcessor(sourceTwo);
Stream streamTwo = builder.createStream(sourceTwo);
```

```
Processor join = new JoinProcessor();
builder.addProcessor(join)
.connectInputShuffle(streamOne)
.connectInputKey(streamTwo);
```

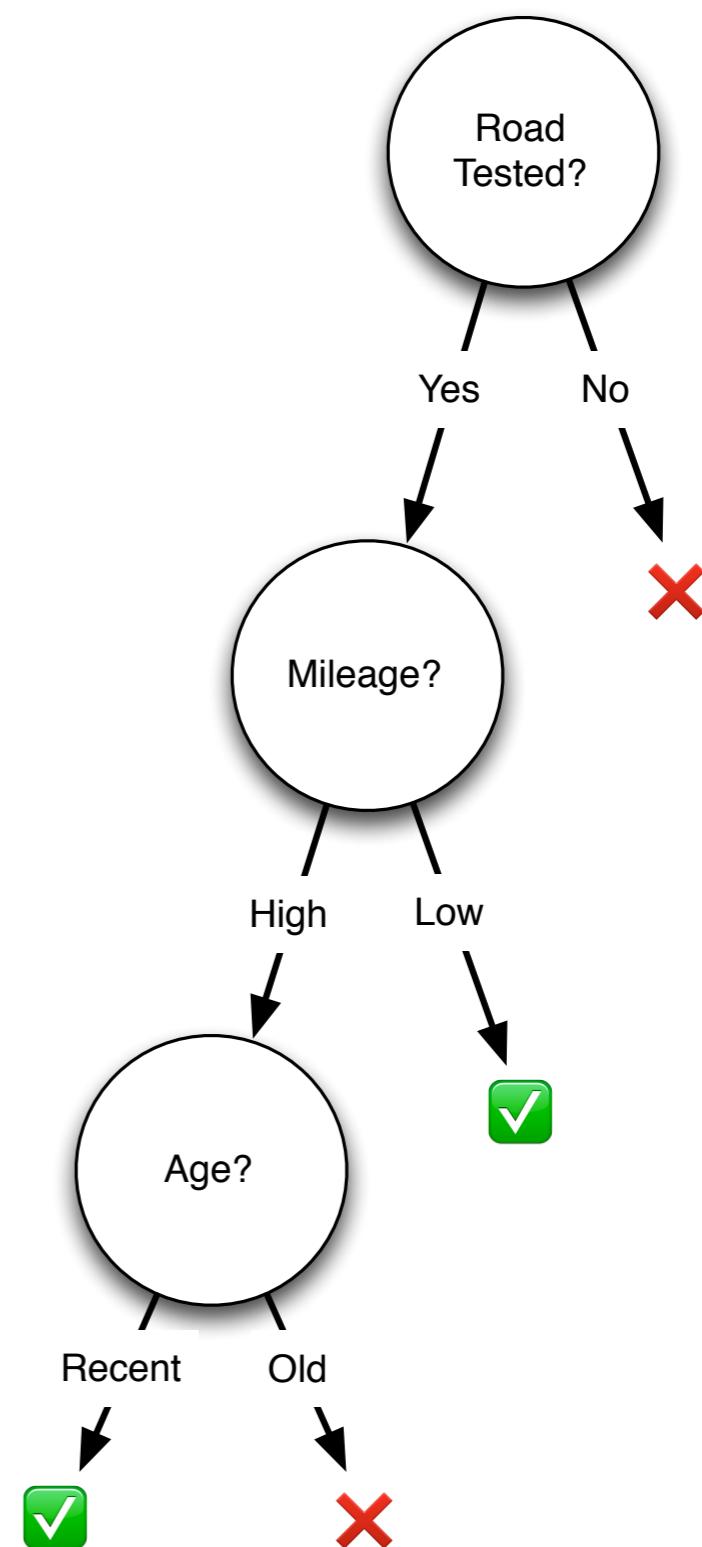
# VERTICAL HOEFFDING TREE (VHT)

# DECISION TREE

- Nodes are tests on attributes
- Branches are possible outcomes
- Leafs are class assignments



## Car deal?



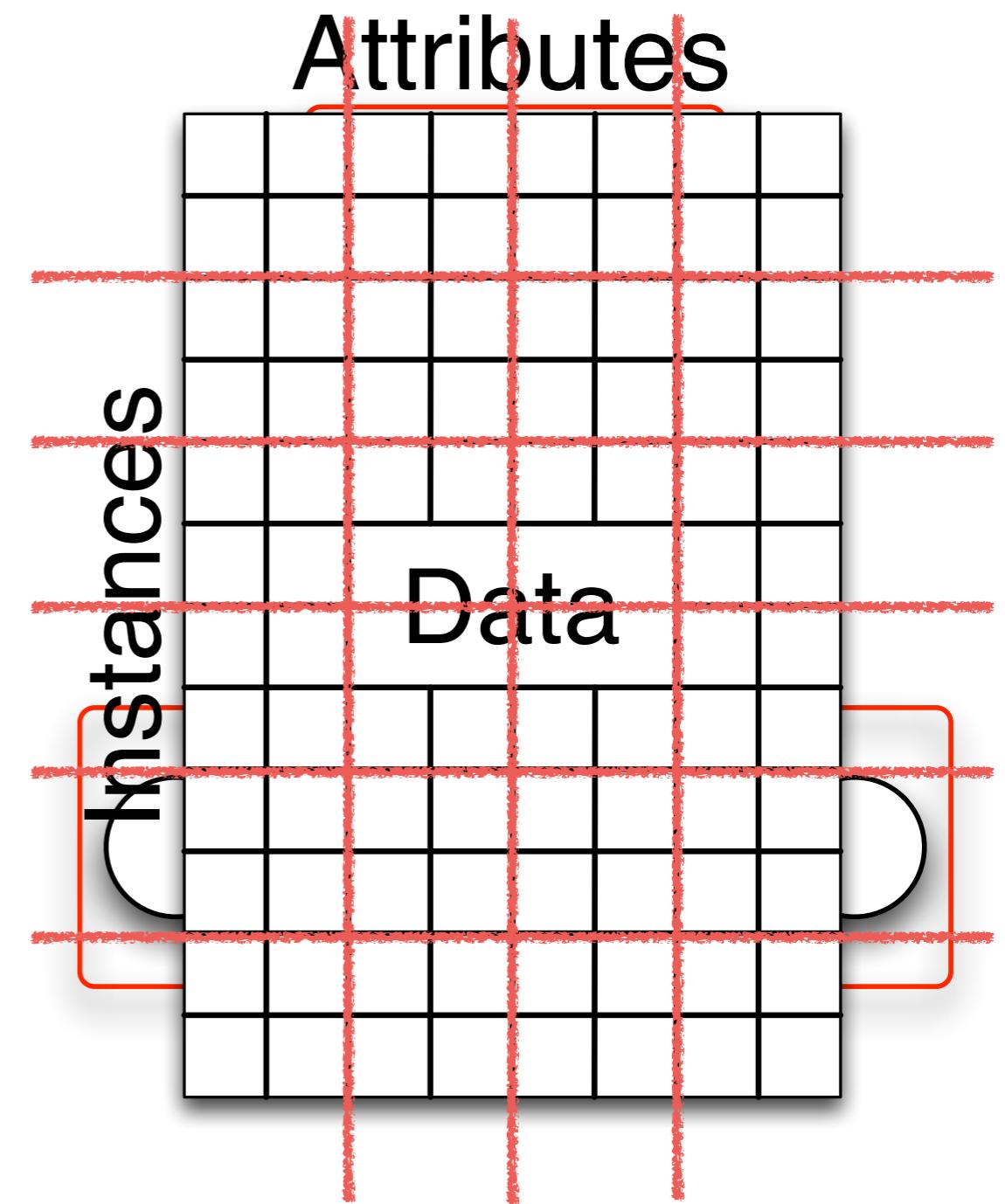
# HOEFFDING TREE

- Sample of stream enough for near optimal decision
- Estimate merit of alternatives from prefix of stream
- Choose sample size based on statistical principles
- When to expand a leaf?
  - Let  $x_1$  be the most informative attribute,  
 $x_2$  the second most informative one
  - Hoeffding bound: split if

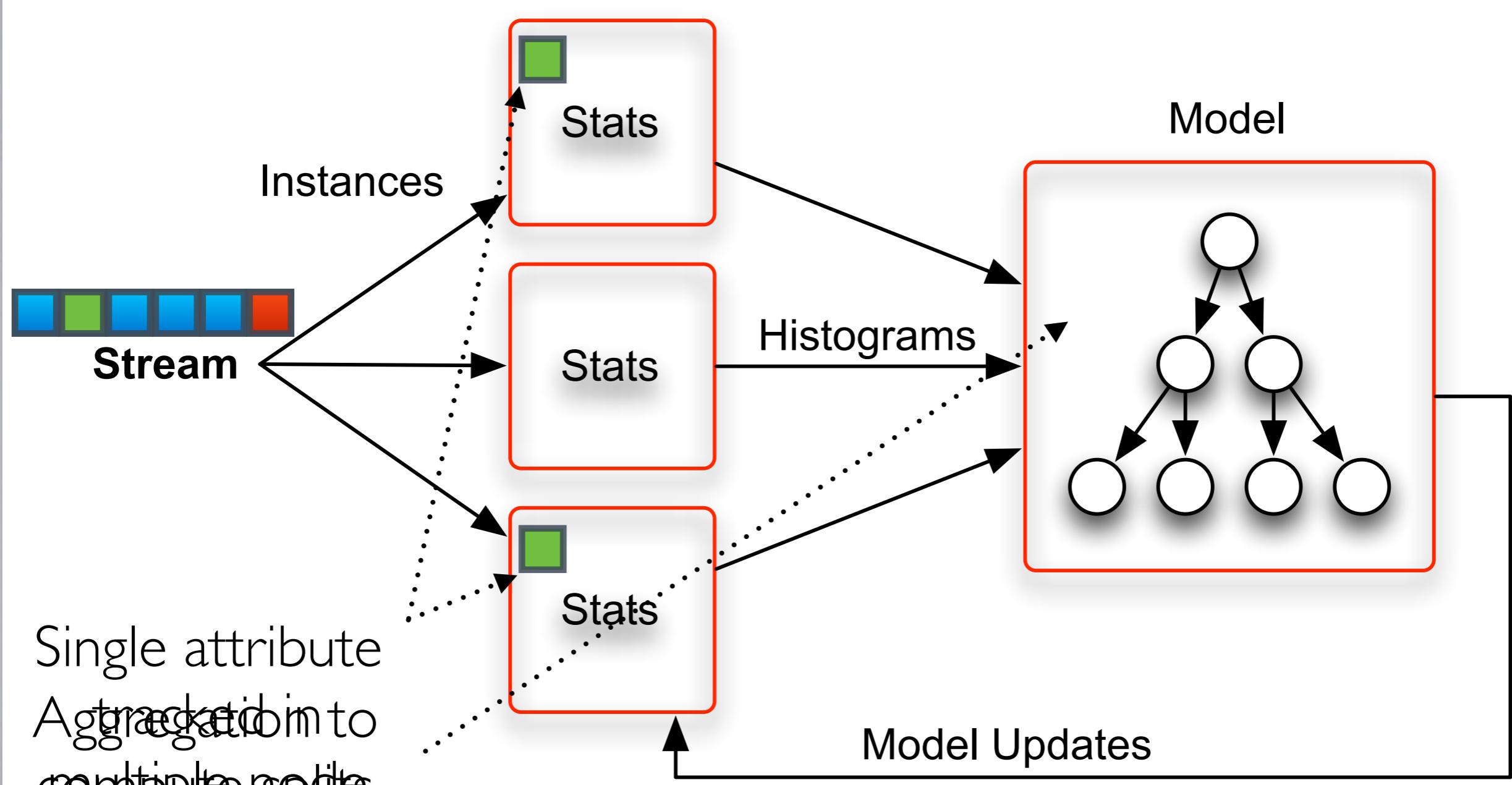
$$\Delta G(x_1, x_2) > \epsilon = \sqrt{\frac{R^2 \ln(1/\delta)}{2n}}$$

# PARALLEL DECISION TREES

- Which kind of parallelism?
  - Task
  - Data
    - Horizontal
    - Vertical

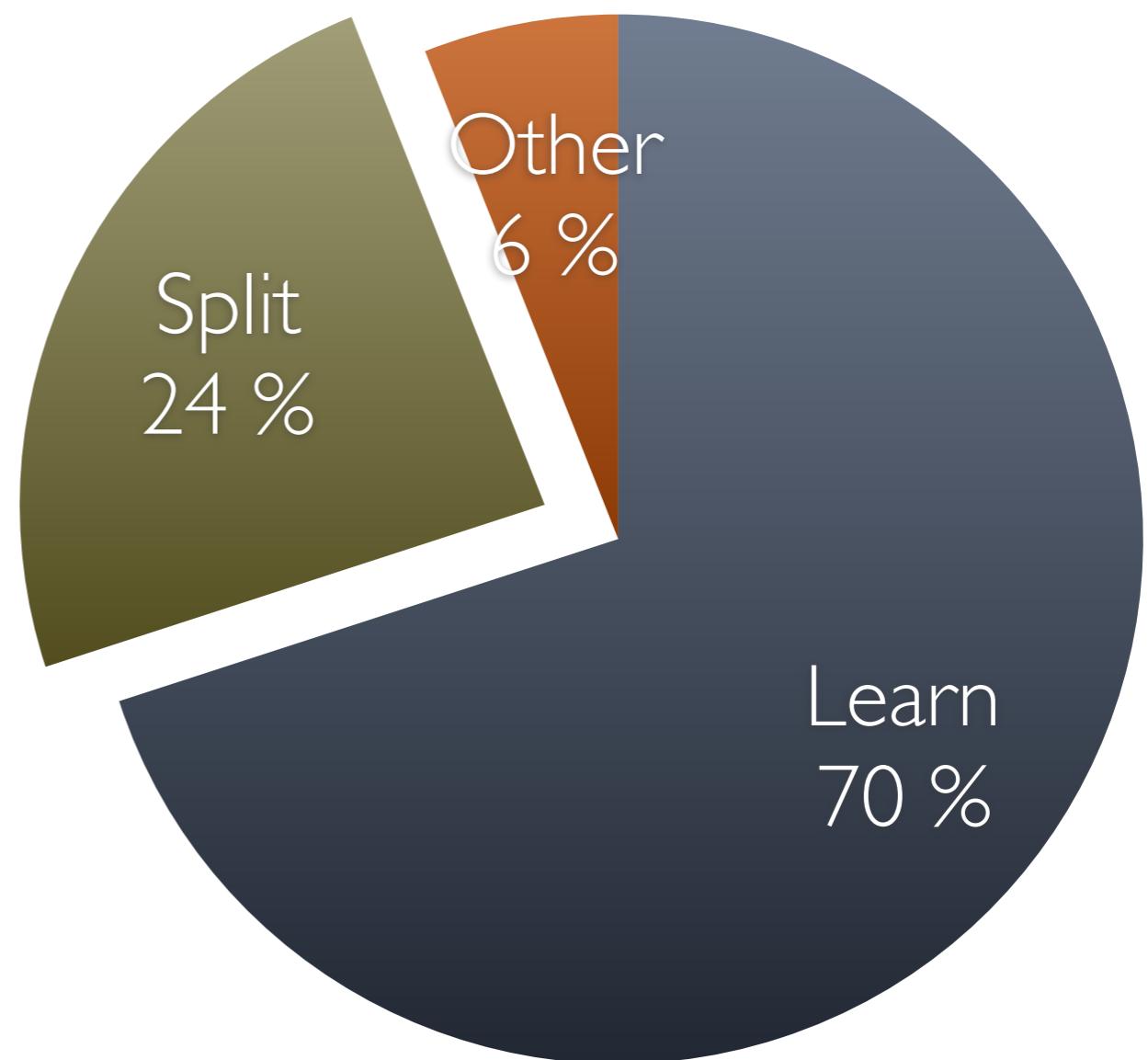


# HORIZONTAL PARALLELISM

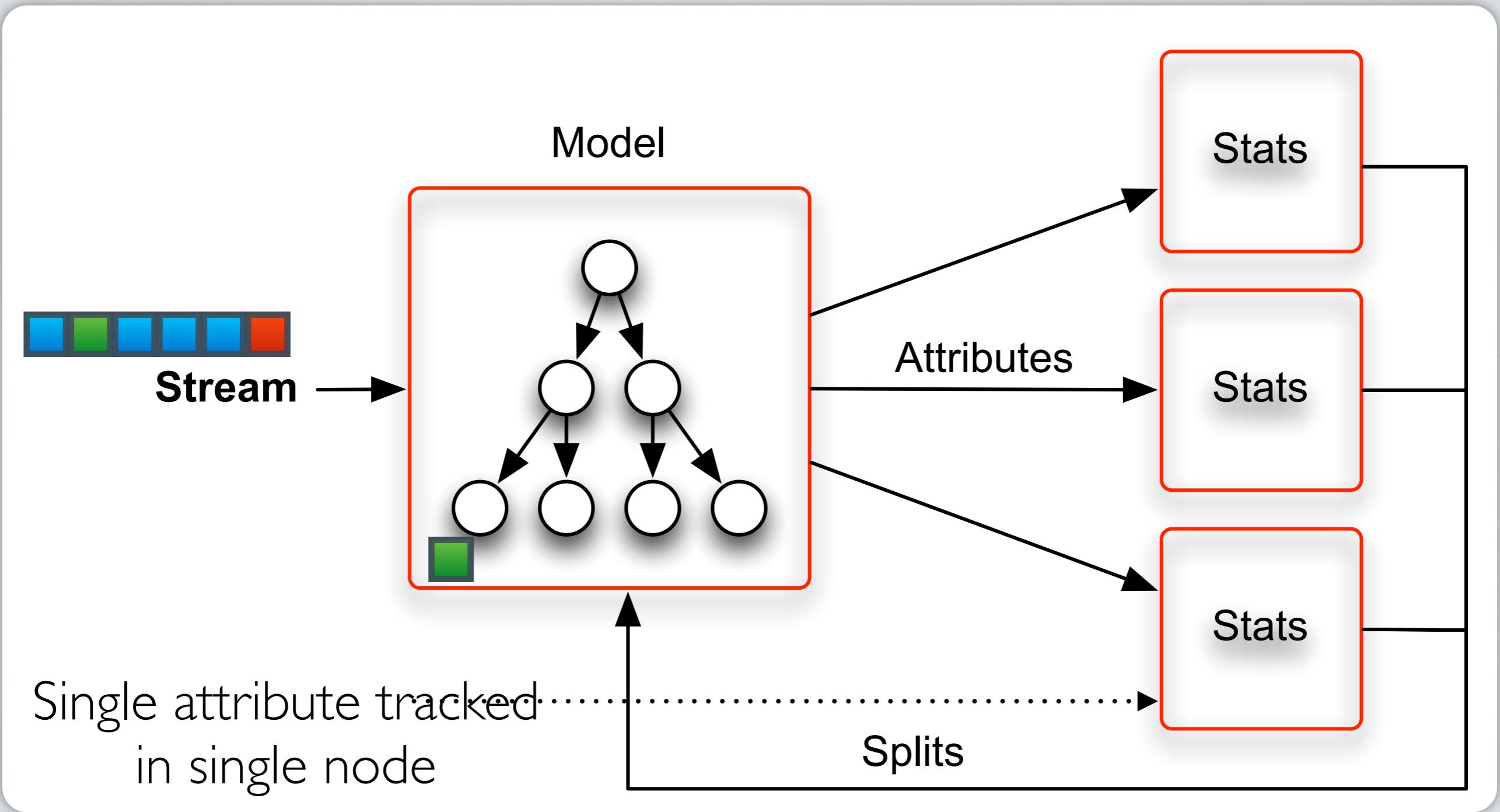


# HOEFFDING TREE PROFILING

CPU time for training  
100 nominal and 100  
numeric attributes



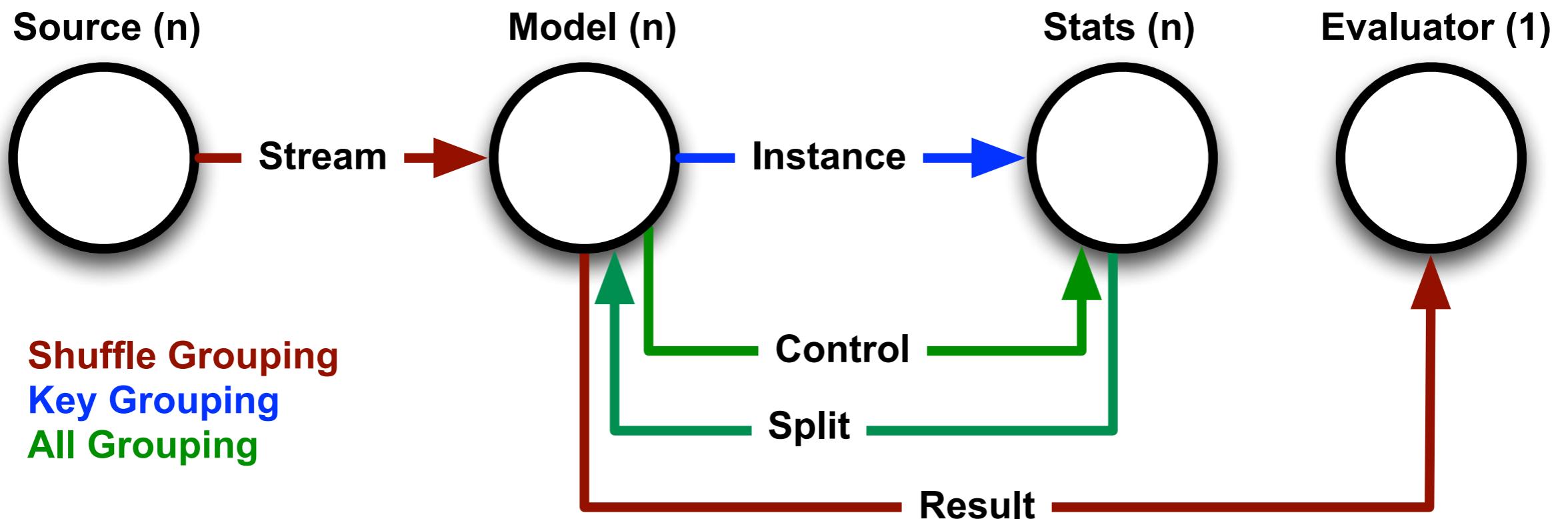
# VERTICAL PARALLELISM



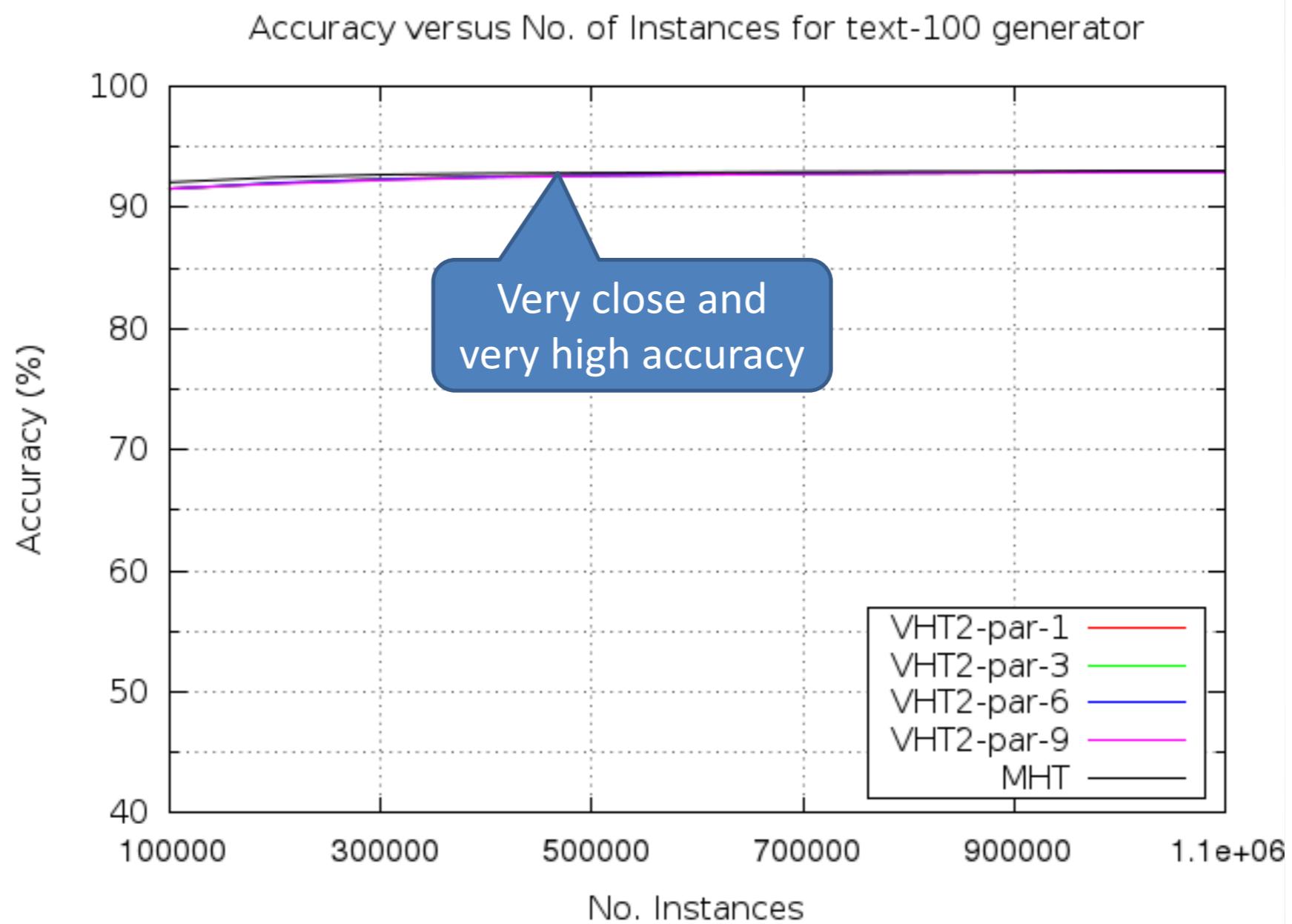
# ADVANTAGES OF VERTICAL

- High number of attributes => high level of parallelism  
(e.g., documents)
- Vs task parallelism
  - Parallelism observed immediately
- Vs horizontal parallelism
  - Reduced memory usage (no model replication)
  - Parallelized split computation

# VERTICAL HOEFFDING TREE

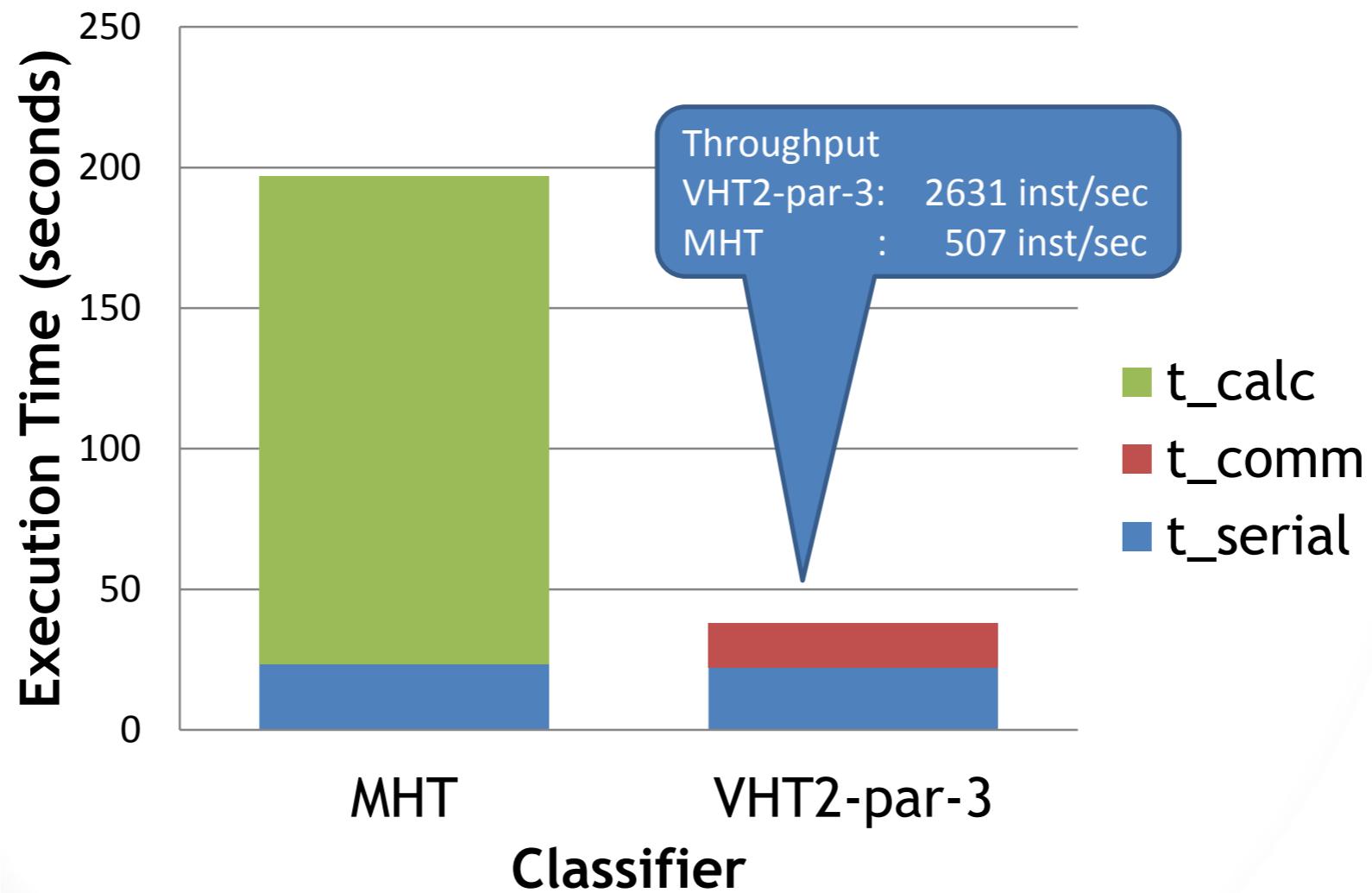


# ACCURACY



# PERFORMANCE

## Profiling Results for text-10000 with 100000 instances



# SUMMARY

- Streaming is an important V of Big Data
- Mining big data streams is an open field
- MOA: Massive Online Analytics
  - Available and open-source <http://moa.cms.waikato.ac.nz/>
- SAMOA: A Platform for Mining Big Data Streams
  - Available and open-source (incubating @ASF)  
<http://samoa.incubator.apache.org>

# OPEN CHALLENGES

- Distributed stream mining algorithms
  - Active & semi-supervised learning + crowdsourcing
  - Millions of classes (e.g., Wikipedia pages)
  - Multi-target learning
- System issues (load balancing, communication)
- Programming paradigms and abstractions

# SAMOA TEAM



Albert  
Bifet



Gianmarco  
De Francisci Morales



Nicolas  
Kourtellis



Matthieu  
Morel

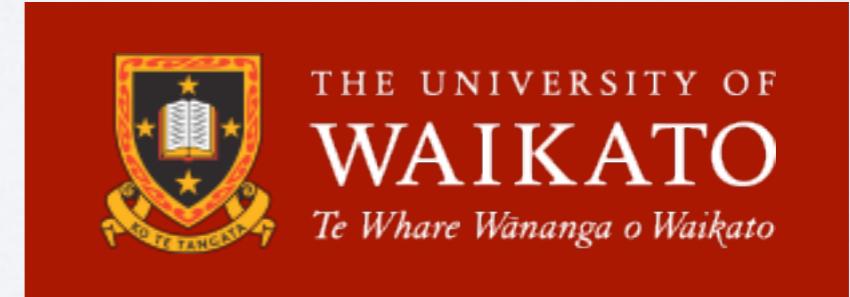


Arinto  
Murdopo



Olivier  
Van Laere

# SUPPORTING ORGANISATIONS



# THANKS!



**@ApacheSAMOA**

<https://samoa.incubator.apache.org>