# RADICALBIT

## radicalbit.io

## Blending Event Stream Processing with Machine Learning using the Kafka Ecosystem

Data Council, Barcelona, Oct 2nd, 2019

# DISCLAIMER

A bit of me.

## Andrea Spina
*Head of R&D* @Radicalbit

[andrea.spina@radicalbit.io](mailto:andrea.spina@radicalbit.io)

@Spina89

andrea-spina

RADICALBIT
radicalbit.io

# RADICALBIT PRODUCTS

*Radicalbit is a highly specialized software firm, founded in Milan, in 2015, focused on the design and development of products dedicated to Event Stream Processing solutions, daily working to combine streaming technologies, Machine Learning and AI with a self-service approach.*

## Radicalbit Natural Analytics

RNA is a platform offering the most advanced self-service capabilities for Data Integration, Data Governance, Data Preparation and Data Visualization over streaming based architectures. It offers a complete set of features aimed to manage every step of the Data Lifecycle: from ingestion to visualization.

## Natural Series Database

NSDb is a storage solution conceived having streaming real-time analytics in mind. It fits perfectly the read side of Kappa Architectures (or for systems based on Command Query Responsibility Segregation pattern). The idea is to store metrics and to bind directly the incoming indexed data to the final users, thanks to pushing technologies like WebSocket.

# DISCLAIMER (AGAIN)

During this talk, you're going to listen about some **buzzwords**

- Event Stream Processing
  - Machine Learning

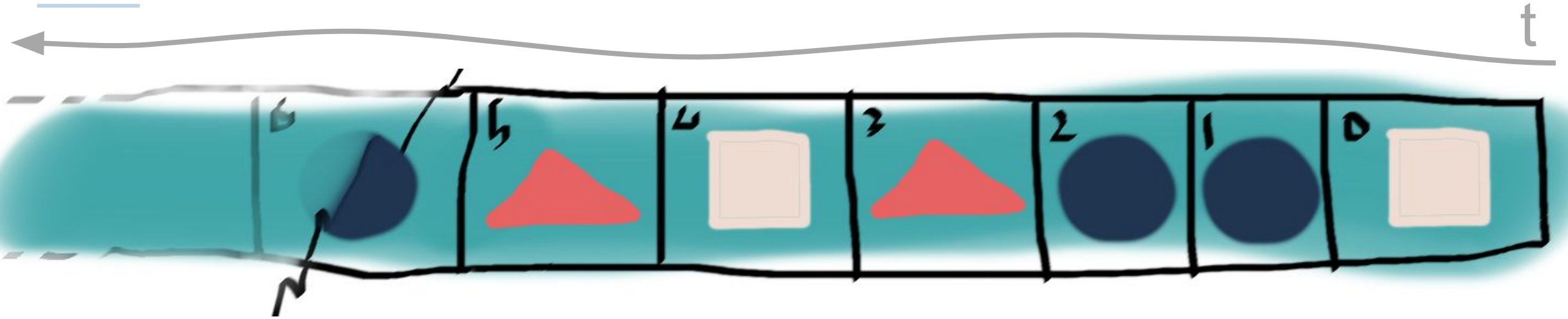You might also hear about topics you already know, and a few you might not ;)

- Lambda v.s. Kappa architectures
  - Machine Learning Logistics
  - Online Machine Learning

RADICALBIT
radicalbit.io

# AGENDA

**RADICALBIT**
radicalbit.io

# Events Stream Processing

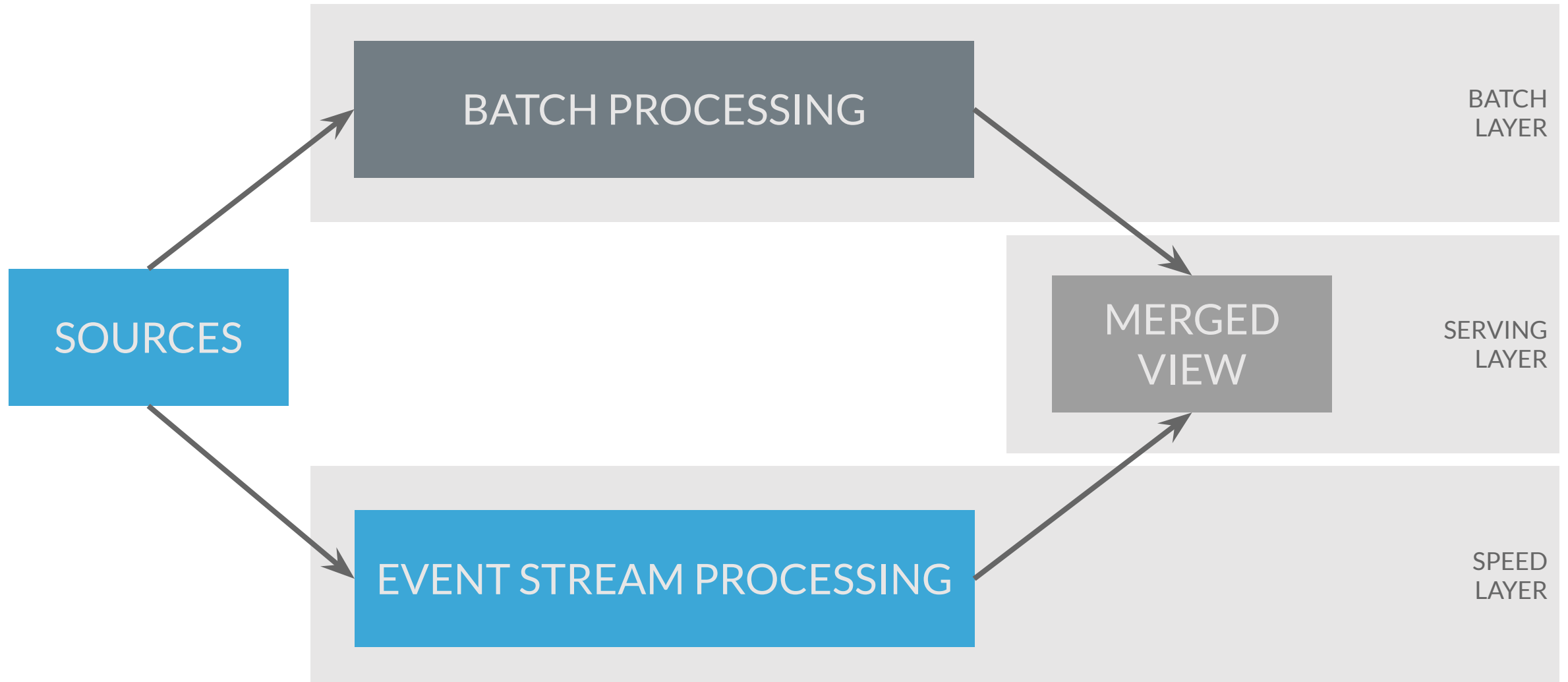# Data Streams



Unbounded

Immutable

Unknown

Batch is stream-*able*

Storable

Transformable

# LAMBDA ARCHITECTURE



SOURCES

BATCH PROCESSING — BATCH LAYER

EVENT STREAM PROCESSING — SPEED LAYER

MERGED VIEW — SERVING LAYER

# LAMBDA ARCHITECTURE ISSUE



BATCH PROCESSING

BATCH LAYER

SOURCES

$$$!

MERGED VIEW

SERVING LAYER

EVENT STREAM PROCESSING

SPEED LAYER

$$$!

# KAPPA ARCHITECTURE



**Log Store / Replay Layer**

Topic_n

Topic_n+1

**Streaming processing system**

Job_Version_n

Job_Version_n+1

**Serving DB**
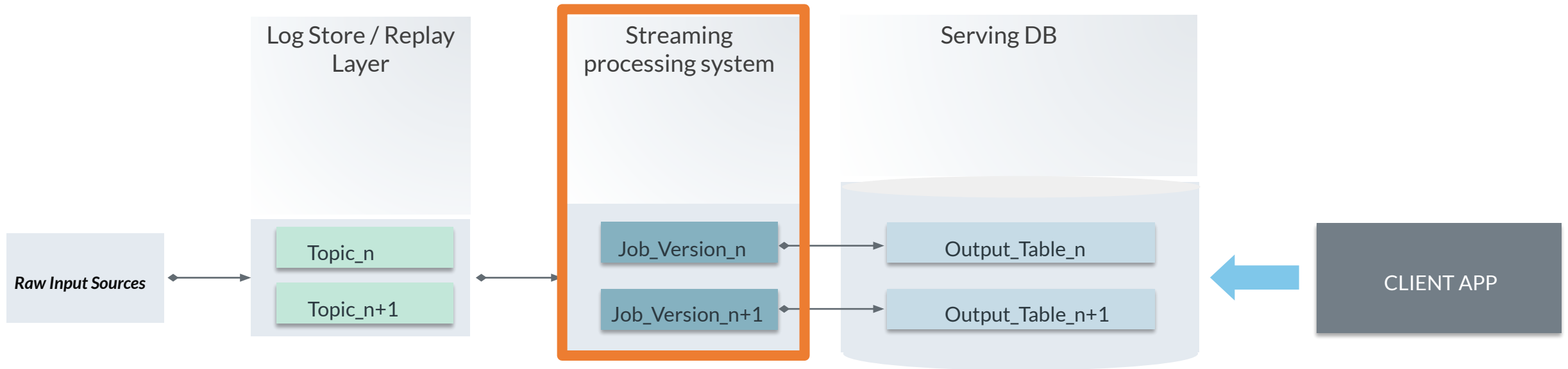
Output_Table_n

Output_Table_n+1

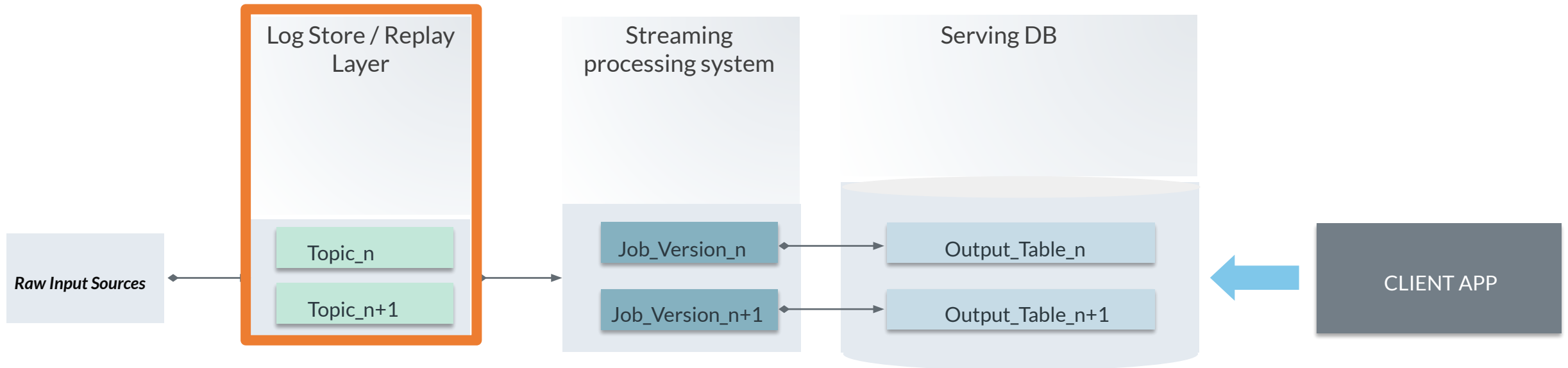*Raw Input Sources*

CLIENT APP

# KAPPA ARCHITECTURE: REQUIREMENTS

1. Low Latency / High throughput

# KAPPA ARCHITECTURE: REQUIREMENTS

1. Low Latency / High throughput
2. Agile *data-reprocessing* method

RADICALBIT
radicalbit.io

# KAPPA ARCHITECTURE: REQUIREMENTS

1. Low Latency / High throughput
2. Agile *data-reprocessing* method
3. Long-time retention message system

| Log Store / Replay Layer | Streaming processing system | Serving DB |
|---|---|---|
| Topic_n | Job_Version_n | Output_Table_n |
| Topic_n+1 | Job_Version_n+1 | Output_Table_n+1 |

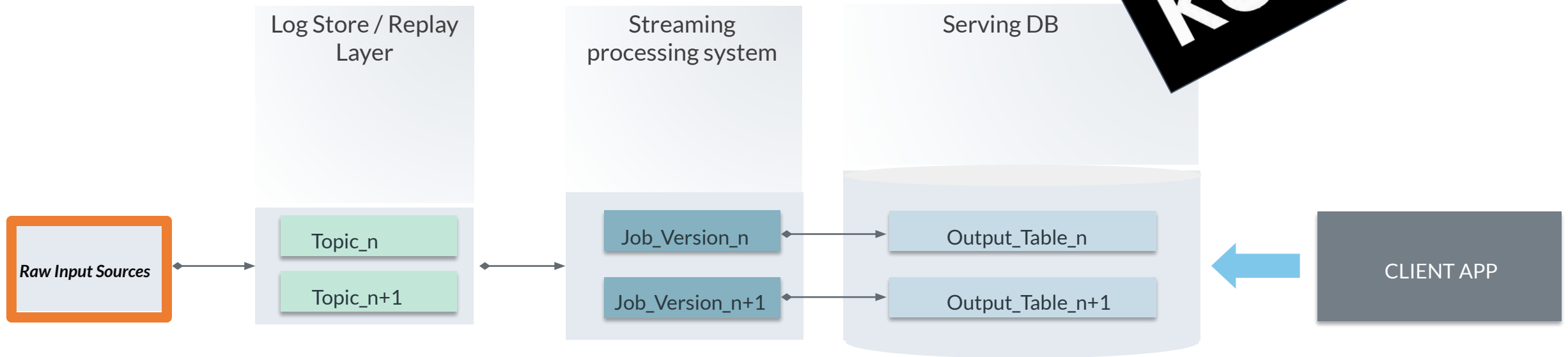**Raw Input Sources**

CLIENT APP

RADICALBIT
radicalbit.io

# KAPPA ARCHITECTURE: REQUIREMENTS

1. Low Latency / High throughput
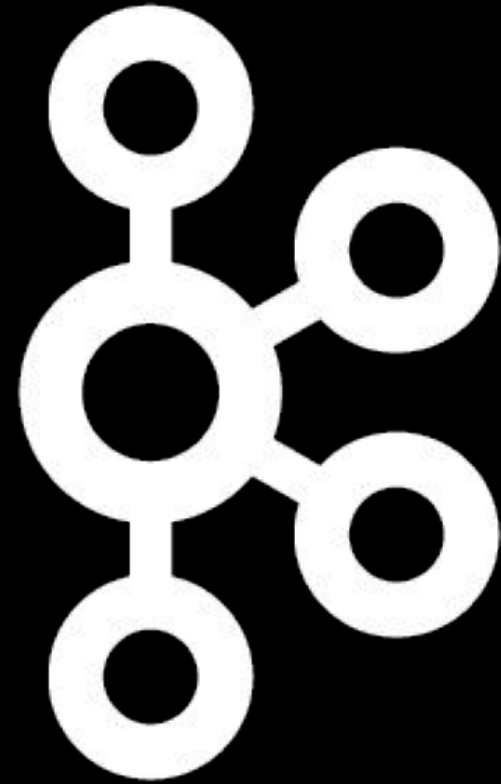2. Agile *data-reprocessing* method
3. Long-time retention message system



| Log Store / Replay Layer | Streaming processing system | Serving DB |
|---|---|---|

Raw Input Sources

Topic_n
Topic_n+1

Job_Version_n
Job_Version_n+1

Output_Table_n
Output_Table_n+1

CLIENT APP

# KAFKA

*" Apache Kafka is a distributed streaming platform.*

- *Publish and subscribe to streams of records, similar to a message queue or enterprise messaging system*

- *Store streams of records in a fault-tolerant durable way*
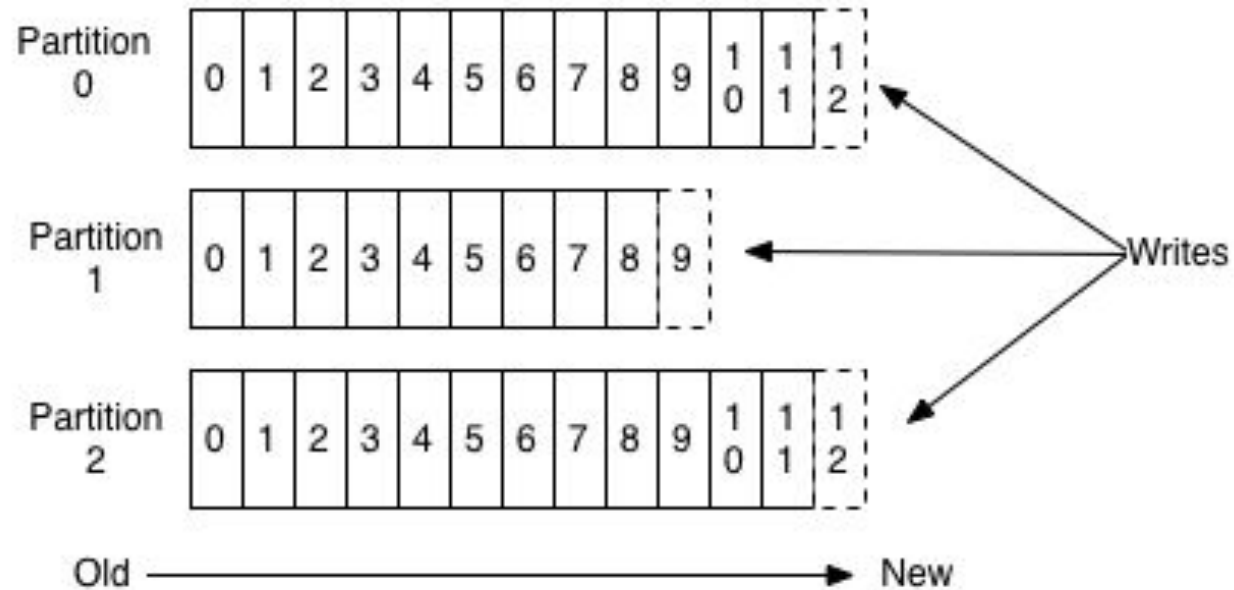
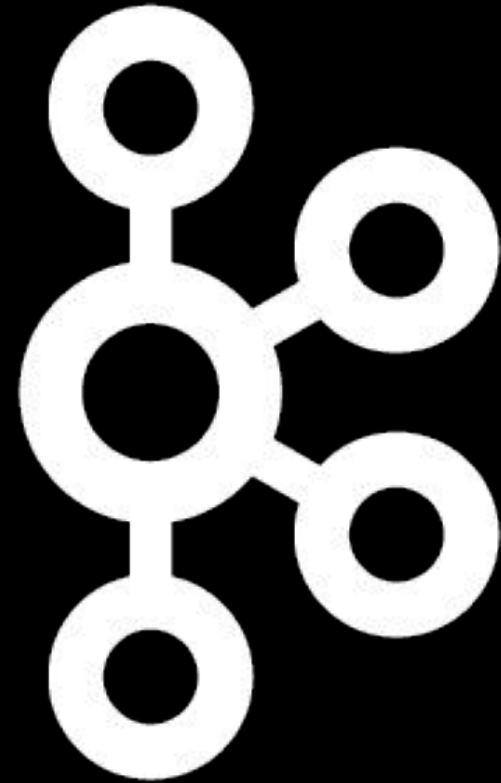- *Process streams of records as they occur "*

*Data reprocessing means "**resetting offsets**"*

## Anatomy of a Topic

# KAFKA MODULES

1. Consumer API
2. Producer API
3. Connect API
4. **Streams API**

*"Kafka at the core of tens of thousands production use-cases"*
**Jay Kreps, Kafka Summit - New York, 2019**

RADICALBIT
radicalbit.io

# RNA AND THE KAPPA ARCHITECTURE

Radicalbit platform has been optimized to take full advantage of Kafka core features such as Kafka Connect, the Schema registry, and Kafka Streams but can be used to manage data pipelines also over Apache Flink or Spark Streaming with code portability

# Blending Machine Learning with Streaming

RADICALBIT
radicalbit.io

# TARGETED ML TASKS

- Models Serving
- Online Machine Learning

RADICALBIT
radicalbit.io

# Streaming Models Serving

*and the magic of machine learning logistics*

RADICALBIT
radicalbit.io

# STREAMING MODELS SERVING

- Serve models in a event stream processing architecture

- It's a **Machine learning logistics** issue (1)

  ○ Organisations need Data Scientists and Data Engineers

  ○ New Tools make it harder (2)

- Fragmented solution space

  ○ Framework based: Tensorflow Serving, Spark, Openscoring

  ○ Cloud based: Google, IBM, MS Azure, Amazon

(1)    Ted Dunning & Ellen Friedman - Machine Learning Logistics - OREILLY

(2)    Boris Lublinsky - Serving Machine Learning Models - OREILLY

# SOLUTIONS TO ML ENTROPY

1. **STANDARD BASED**

   Define a **youNameIt**-independent format to represent a wide range of ML models

   - PMML (PFA) - traditional learning

   - ONNX - deep learning

   - MLEAP - not a STD

```
<LocalTransformations>
    <DerivedField name="GroupedInputColor" dataType="integer"
        optype="categorical">
        <MapValues outputColumn="group" defaultValue="3" mapMissingTo="4">
            <FieldColumnPair column="color" field="VarInputColor" />
            <TableLocator>
                <Extension extender="ADAPA" name="TABLE_NAME"
                    value="GroupedInputColor" />
            </TableLocator>
        </MapValues>
    </DerivedField>
    <DerivedField name="VarColorGroup_1" dataType="integer" optype="categorical">
        <NormDiscrete value="1" field="GroupedInputColor" method="indicator" />
    </DerivedField>
    <DerivedField name="VarColorGroup_2" dataType="integer" optype="categorical">
        <NormDiscrete value="2" field="GroupedInputColor" method="indicator" />
    </DerivedField>
    <DerivedField name="VarColorGroup_3" dataType="integer" optype="categorical">
        <NormDiscrete value="3" field="GroupedInputColor" method="indicator" />
    </DerivedField>
    <DerivedField name="VarColorGroup_4" dataType="integer" optype="categorical">
        <NormDiscrete value="4" field="GroupedInputColor" method="indicator" />
    </DerivedField>
</LocalTransformations>
```

RADICALBIT
radicalbit.io

# SOLUTIONS TO ML ENTROPY

2.  **CONTAINER BASED**

    Creating containers wrapping environments natively aimed at models deployment

    Exposing a communication protocol for serving (usually a REST endpoint)

    - Seldon core[1]

    - Clipper[2]

    - MLFlow[3]

(1)   https://www.seldon.io/open-source/
(2)   http://clipper.ai/
(3)   https://mlflow.org/

RADICALBIT
radicalbit.io

# PROS and CONS[1]

## STANDARD BASED

PROS Performance, flexibility, many people are happy

CONS Adoption, algorithms

## CONTAINER BASED

PROS Repeatability, adoption is not a problem, everybody is happy
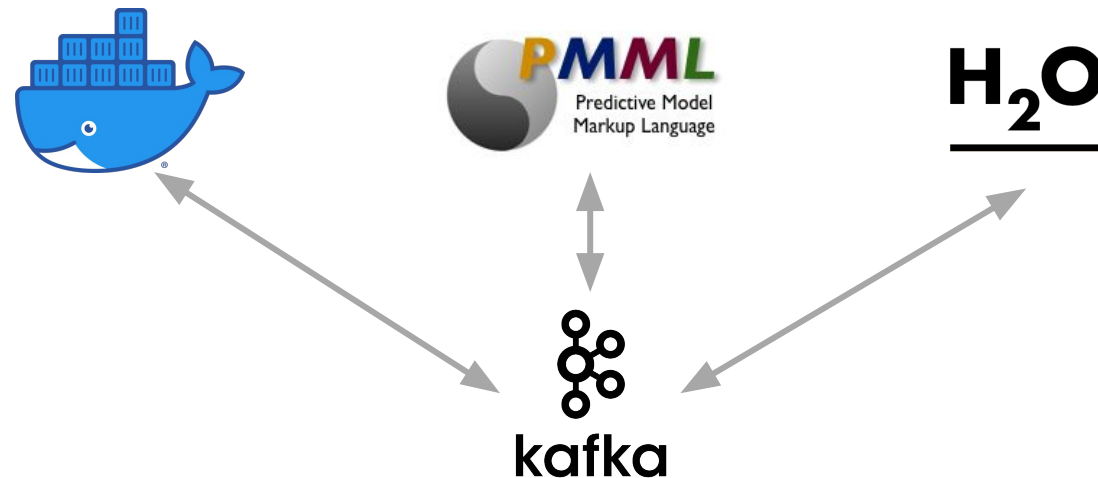
CONS Performance depends on systems, devops competence

(1)   https://qconsp.com/sp2018/system/files/presentation-slides/qconsp18-deployingml-may18-npentreath.pdf

# SERVING AS A SERVICE WITH KAFKA

## The goal

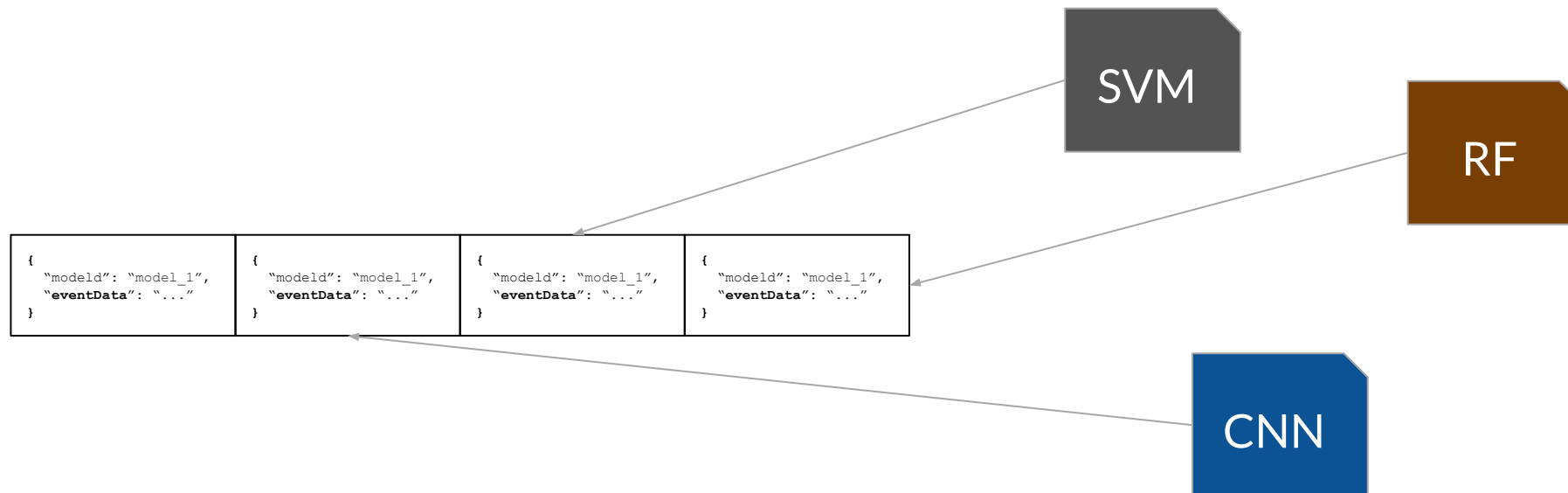Attempting to serve seamless **Standards**, **Containers**, and **Tools** using Kafka

- No constraints about models deployment (it has not to be even a ML model!)

- It potentially has not to be even a ML model!

Our predictive **k**-pipelines shall:

- **dynamically serve the evolution of trained models**
  - models often change in behavior during their long-lasting lifetime
  - updates

- apply simultaneously multiple models against the same stream, the same model to many streams

SVM

RF

```
{
  "modelId": "model_1",
  "eventData": "..."
}
```
```
{
  "modelId": "model_1",
  "eventData": "..."
}
```
```
{
  "modelId": "model_1",
  "eventData": "..."
}
```
```
{
  "modelId": "model_1",
  "eventData": "..."
}
```

CNN

RADICALBIT
r a d i c a l b i t . i o

# KAFKA STREAMS APIs

- Kafka Streams is not a DSPE, is a library[1]

- By Kafka Streams APIs, users define a processor topology

- Two API levels

    - Kafka Streams DSL

    - Processor API

(1)   https://kafka.apache.org/23/documentation/streams/

# SERVING AS A SERVICE: KS-H$_2$O Example

- Gartner 2019 magic quadrant for Machine Learning

- Most of the code is *open source*

- High support for algorithms

- H$_2$O flow

Main features

- well-built Rest API layer

- **POJO** and **MOJO** formats + client library

Figure 1. Magic Quadrant for Data Science and Machine Learning Platforms

Source: Gartner (January 2019)

```
{
    …
    "id": "unsupervised_cusomers_1",
    "algorithm": "kmeans",
    "format": "mojo",
    "exp_date": null,
    "more_info": " ... "
    …
}
```

RADICALBIT
radicalbit.io

1 - to - 1 Bind

Model Repository Server → Control Stream

RADICALBIT
radicalbit.io

# KSH$_2$O - FEEDING A METADATA TABLE

**GlobalKTable**

```
{
    ...
    "id": "unsupervised_cusomers_1",
    "algorithm": "kmeans",
    "format": "mojo"
    "exp_date": null
    "more_info": " ... "
    ...
}
```
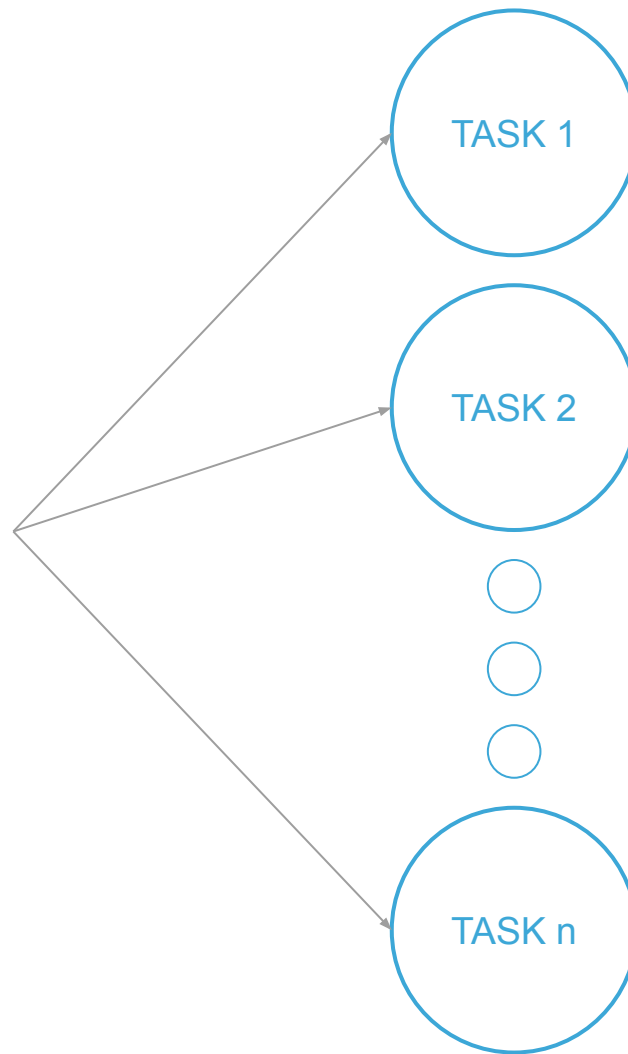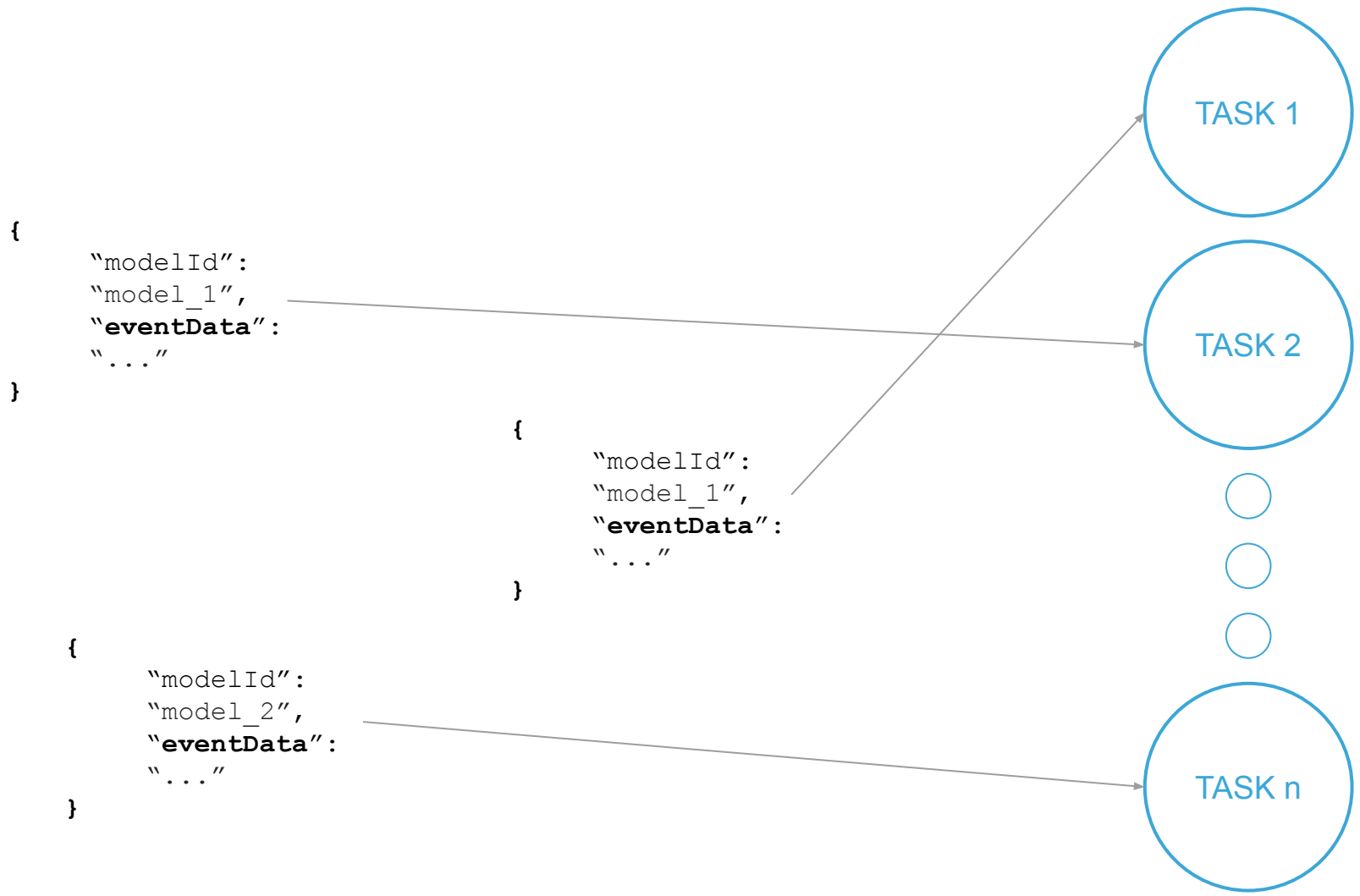
TASK 1

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

TASK 2

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

TASK n

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

RADICALBIT
radicalbit.io

# KSH$_2$O - THE DATA STREAM

```
{
    "modelId": "ID1",
    "eventData":
    "..."
}
```

POJO model

I'M A TASK!

| ID1 | /../.. | m1 | more |
|-----|--------|----|------|
| ID2 | /../.. | m2 | more |

H$_2$O

**T1**

**RESTORED**

| ID1 | /../.. | m1 | more |
|-----|--------|----|------|
| ID2 | /../.. | m2 | more |

**T2**

**RESTORED**

| ID1 | /../.. | m1 | more |
|-----|--------|----|------|
| ID2 | /../.. | m2 | more |

**T3**

**RESTORED**

| ID1 | /../.. | m1 | more |
|-----|--------|----|------|
| ID2 | /../.. | m2 | more |

**T4**

**RESTORED**

| ID1 | /../.. | m1 | more |
|-----|--------|----|------|
| ID2 | /../.. | m2 | more |

```
{
    "modelId": "ID1",
    "eventData": "..."
}
```

**H$_2$O**

**POJO**

**T2**

| ID1 | /../.. | m1 | more |
|-----|--------|----|------|
| ID2 | /../.. | m2 | more |

On *restore*, lazy uploading applies models' recovering

RADICALBIT
radicalbit.io

# THE KS-H$_2$O ARCHITECTURE

Configuration
- **Empty predictions** strategy
- **NaN** management
- **UDF**
- ...

```
{
    …
    "id":
    "unsupervised_cusomers_1",
    "algorithm": "kmeans",
    "format": "mojo"
    "exp_date": null
    "more_info": " ... "
    …
}
```

**TASK 1**

**POJO**

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

**TASK 2**

**MOJO**

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

```
{
    "modelId": "model_1",
    "eventData": "..."
}
```
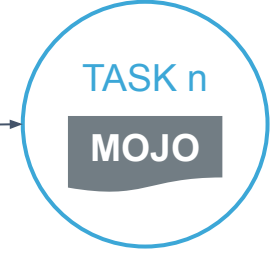
**TASK n**

**MOJO**

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

```
{
    "modelId": "model_2",
    "eventData": "..."
}
```

```
{
    "UserDefinedOutput":
    { … },
    "prediction": {
    /*PREDICTION OBJECT*/
    }
}
```

**ks-h2o -** https://github.com/radicalbit

RADICALBIT
radicalbit.io

```
{
     …
     "id":
     "unsupervised_cusomers_1",
     "algorithm": "kmeans",
     "format": "mojo"
     "exp_date": null
     "more_info": " ... "
     …
}
```

TASK 1

**POJO**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

# First generalisation

Given a control message, how to build the shared state

RADICALBIT
radicalbit.io

# Second generalisation

Given the record to score, how to build the model

| | |
|---|---|
| { | |
| "modelId": "model_2", | |
| "eventData": "…" | |
| } | |

TASK n

**MOJO**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

RADICAL BIT
radicalbit.io

# Third generalisation

## Given the record to score, implement the scoring method

{

    "modelId": "model_2",
    "eventData": "..."

}

TASK n

**MOJO**

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

RADICALBIT
radicalbit.io

Configuration
- Empty predictions strategy
- NaN management
- UDF
- ...

```
{
    …
    "id":
    "unsupervised_cusomers_1",
    "algorithm": "kmeans",
    "format": "mojo"
    "exp_date": null
    "more_info": " ... "
    …
}
```

TASK 1
**POJO**

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

TASK 2
**MOJO**

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

```
{
    "modeId": "model_1",
    "eventData": "..."
}
```

TASK n
**MOJO**

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

```
{
    "modeId": "model_2",
    "eventData": "..."
}
```

```
{
    "UserDefinedOutput":
    { … },
    "prediction": {
    /*PREDICTION OBJECT*/
    }
}
```

https://github.com/radicalbit

**ks-h2o -** https://github.com/FlinkML/flink-jpmml
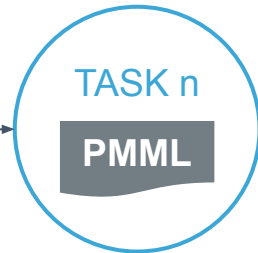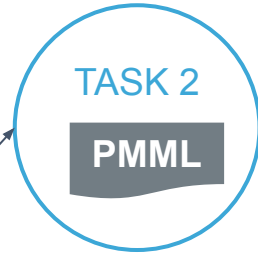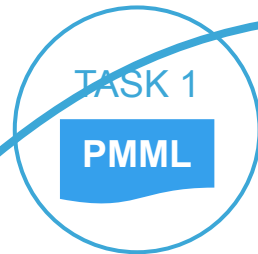
# SERVING STANDARD MODELS DEPLOYMENTS

```
{
    …
    "id":
    "unsupervised_cusomers_1",
    "algorithm": "svm",
    "format": "PMML"
    "model_path": /opt/mdls
    "more_info": " ... "
    …
}
```

Configuration
- **Empty predictions** strategy
- **NaN** management
- **UDF**
- **...**

**TASK 1**

**PMML**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

**TASK 2**

**PMML**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

```
{
    "UserDefinedOutput":
    { … },
    "prediction": {
    /*PREDICTION OBJECT*/
    }
}
```

```
{
    "modelId": "model_1",
    "eventData": "..."
}
```

**TASK n**

**PMML**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

```
{
    "modelId": "model_2",
    "eventData": "..."
}
```

**flink-jpmml -** https://github.com/FlinkML/flink-jpmml

43

RADICALBIT
radicalbit.io

# SERVING CONTAINERIZED MODELS DEPLOYMENTS

```
{
    …
    "id": "rnn-1",
    "algorithm": "net",
    "format": "Tensorflow"
    "serve_url":
    https://serve.com/opt/mdls
    …
}
```

## Configuration
- **Empty predictions** strategy
- **NaN** management
- **UDF**
- …

**TASK 1**
**REST client**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

**TASK 2**
**REST**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

```
{

    "modelId": "model_1",
    "eventData": "..."

}
```

**TASK n**
**REST**

| unsupervised_customers_1 | kmeans | mojo | … |
|---|---|---|---|
| deep_net | deeplearning | pojo | … |

```
{

    "modelId": "model_2",
    "eventData": "..."

}
```
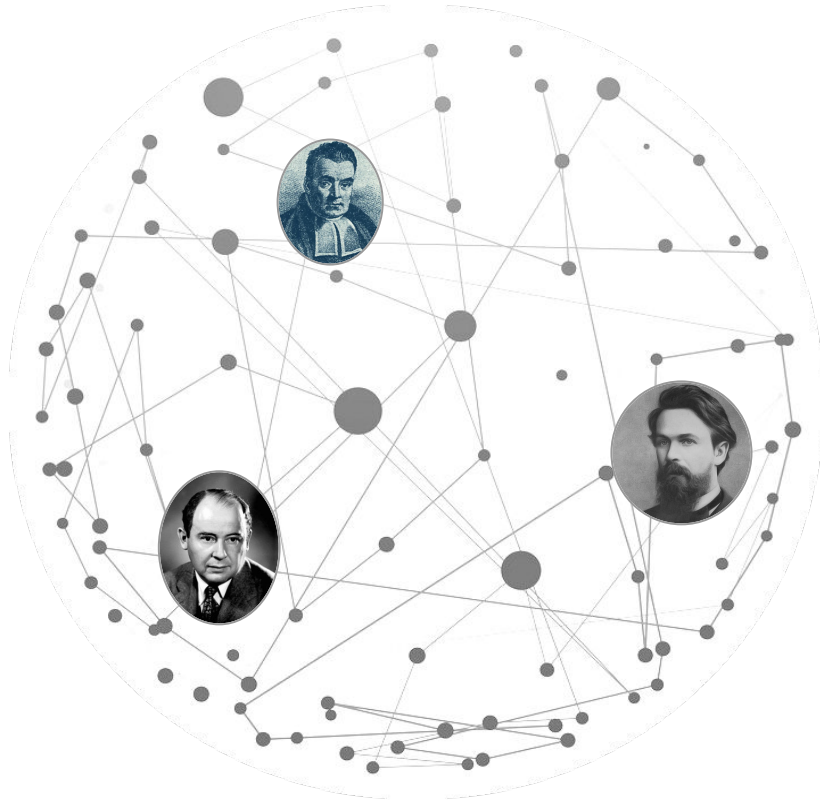
```
{

    "UserDefinedOutput":
    { … },
    "prediction": {
    /*PREDICTION OBJECT*/
    }

}
```

44

# KS-OML

*Online Machine Learning using Kafka Streams*

RADICALBIT
radicalbit.io

# MACHINE LEARNING, TRAINING

| DIM | WEIGHT | COLOR |
|-----|--------|-------|
| 3 | 0.8 | R |
| 1 | 0.2 | b |
| 1 | 1.2 | y |
| 4 | 2.1 | g |
| 3 | 0.9 | r |
| 2 | 1.0 | r |
| 12 | 0.2 | b |
| 1 | 0.3 | g |
| 1 | 0.4 | y |
| 3 | 0.1 | b |
| 3 | 0.2 | g |
| 4 | 2.0 | r |
| 4 | 3.1 | c |
| 3 | 0.8 | R |
| 1 | 0.2 | b |
| 1 | 1.2 | y |

RADICALBIT
radicalbit.io

1.  Computational Model



TASK

RADICALBIT
radicalbit.io

1. Computational Model



TASK

# KS-OML - ONLINE LEARNING CHALLENGES

2. Evolving Data - or *the long-standing issue of the Concept Drift*

# ONLINE MACHINE LEARNING STATE

Growing **academic** interest

Online Machine Learning in Big Data Streams, 2018, András Benczúr, Levente Kocsis, Róbert Pálovics

Still *rare* productionized OML implementations

The value

1. When the ability **to fast adapt** is more important than the best performance
   *Newspaper domestic affairs drift example*

2. When keeping data offline is not possible
   *healthcare data, not reachable data*

3. Resource savings

RADICALBIT
radicalbit.io

# KS-OML

*Online Machine Learning on Kafka architecture*

## Main objective

all-contained operator with configurable **algorithms suite**

**First Implementation:** *passive-aggressive algorithm*, Daniele Tria

**Second Implementation:** *soft confidence-weighted algo*, Seyedmasih Hosseinimotlagh

(1)  **Online Passive-Aggressive Algorithms -** Crammer, Dekel, Keshet, Shalev-Shwartz, Singer

http://jmlr.csail.mit.edu/papers/volume7/crammer06a/crammer06a.pdf

(2)  **Soft confidence-weighted Algorithms -** Wang, Zaho, Hoi

https://arxiv.org/ftp/arxiv/papers/1206/1206.4612.pdf

RADICALBIT
r a d i c a l b i t . i o
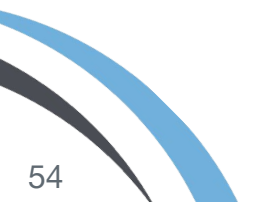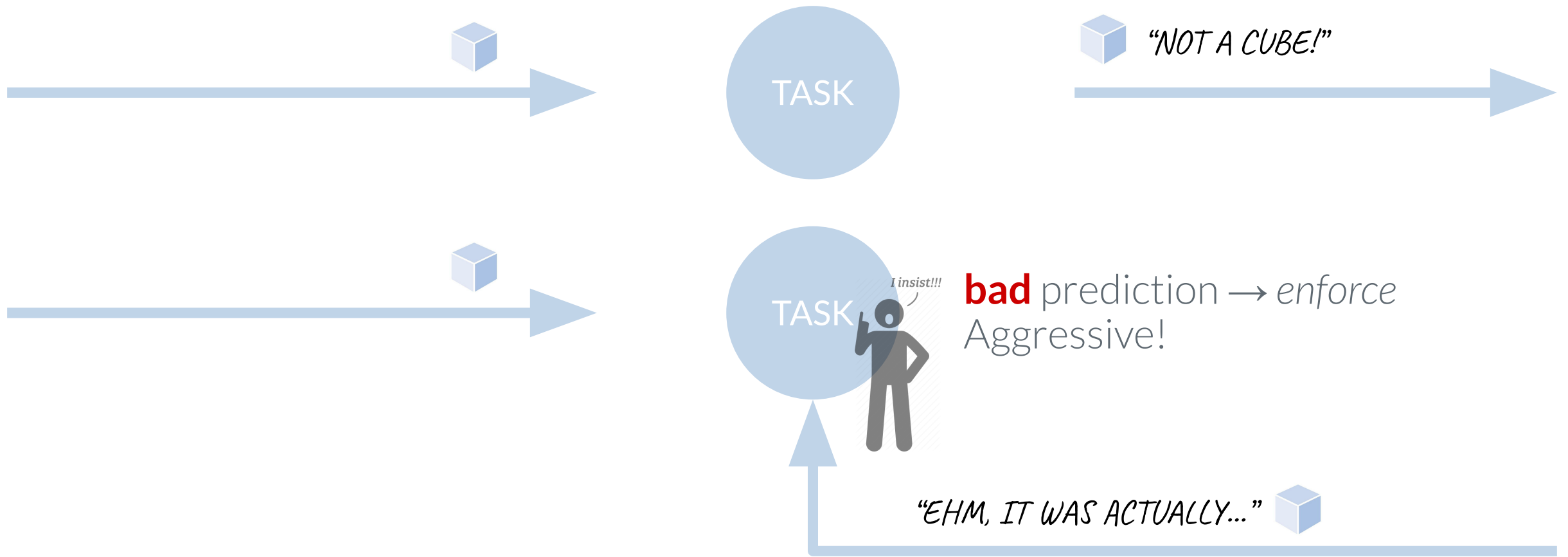
# KS-OML - PASSIVE AGGRESSIVE ALGORITHM

- *margin* based algorithm
- able to solve **binary** class, **multiclass** and **regression** problems
- *feedback* concept

- for binary class, given
  - $y_t$ the true label
  - $\mathbf{x}_t$ the feature vector
  - $\mathbf{w}_t$ the weighted vector of the model

- $y_t(\mathbf{x}_t \cdot \mathbf{w}_t) = margin$     *margin > 0 → correct prediction*     $sign(\mathbf{x}_t \cdot \mathbf{w}_t) = y_t$

- why *passive-aggressive* ?
  - *if margin ≥ +1 → do nothing*
  - *else → enforce the margin*

RADICAL BIT
r a d i c a l b i t . i o

# KS-OML - HIGH-LEVEL WORKFLOW



"CUBE!"

TASK

TASK    *I insist!!!*    good prediction
**margin < 1** → *enforce*
Aggressive!

"YOU WERE RIGHT!"

# KS-OML - HIGH-LEVEL WORKFLOW

TASK

"CUBE!"

TASK

good prediction
**margin > 1** $\rightarrow$ *do nothing*
Passive

"YOU WERE RIGHT!"

# KS-OML - INPUT STREAMS

1. Main Event Stream :     `UnlabeledData`
2. Feedback Event Stream :   `LabeledData`

- *Connected* by the same Scala case class
- written in the same topic "data"

topic "data"

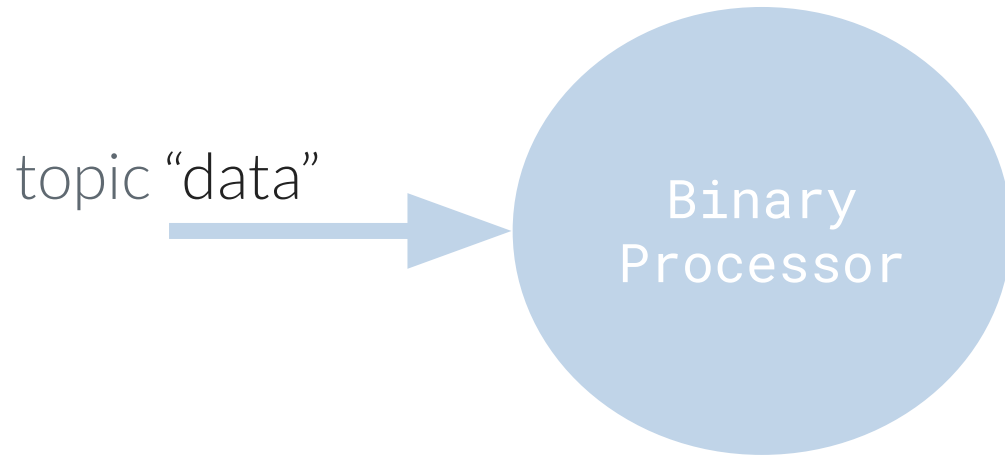`UnlabeledData`

`LabeledData`

topic "data"

Binary
Processor

event Store

model Store

When `UnlabeledData`

- get model from store
- compute prediction
- store event by hashing with prediction
- emit prediction

When `LabeledData`

- get event and model from stores
- check loss function *margin*
- eventually update the model
- delete the event from store
- emit again if required

RADICALBIT
radicalbit.io

| | BANKNOTES DATASET (binary) | IRIS DATASET (multiclass) |
|---|---|---|
| First RUN | - | - |
| Second RUN | accuracy: 0.949671 precision: 0.97872 | accuracy: 1.0 precision: 1.0 |
| Third RUN | accuracy: 0.950747 precision: 0.98391 | accuracy: 1.0 precision: 1.0 |

RADICALBIT
radicalbit.io

1. Feedback algos are good if you get the feedback a.s.a.p

2. Passive Aggressive is good when you can suffer of *cold start*

3. Passive Aggressive is *adaptive*

4. It works. Cool!

RADICALBIT
radicalbit.io

# Pouring the blend

## Model Serving

Streams are the perfect fit

Kafka is a natural solution for distribution and performance - but you need to tackle Kappa challenges!

Growing desire of an unique abstraction (both low and high level)

## Online Learning

Increasing interest, industry still immature

Global shared streams, or states, even stores are fundamental to Machine Learning

RADICALBIT
radicalbit.io

# THANKS!

## Office hours 12.45 - 13.30
### Your questions are welcome!

info@radicalbit.io

RADICALBIT
r a d i c a l b i t . i o

# REFERENCES

https://medium.com/value-stream-design/online-machine-learning-515556ff72c5
https://medium.com/analytics-vidhya/data-streams-and-online-machine-learning-in-python-a382e9e8d06a

RADICALBIT
radicalbit.io

# Bonus Slides

RADICALBIT
radicalbit.io

# MACHINE LEARNING

| DIM | WEIGHT | COLOR |
|-----|--------|-------|
| 3 | 0.8 | "R" |

IT'S A "CUBE"!

| DIM | WEIGHT | COLOR |
|-----|--------|-------|
| 3 | 0.8 | R |
| 1 | 0.2 | b |
| 1 | 1.2 | y |
| 4 | 2.1 | g |
| 3 | 0.9 | r |
| 2 | 1.0 | r |
| 12 | 0.2 | b |
| 1 | 0.3 | g |
| 1 | 0.4 | y |
| 3 | 0.1 | b |
| 3 | 0.2 | g |
| 4 | 2.0 | r |
| 4 | 3.1 | c |
| 3 | 0.8 | R |
| 1 | 0.2 | b |
| 1 | 1.2 | y |

# MACHINE LEARNING, TRAINING

| DIM | WEIGHT | COLOR |
|-----|--------|-------|
| 3 | 0.8 | R |
| 1 | 0.2 | b |
| 1 | 1.2 | y |
| 4 | 2.1 | g |
| 3 | 0.9 | r |
| 2 | 1.0 | r |
| 12 | 0.2 | b |
| 1 | 0.3 | g |
| 1 | 0.4 | y |
| 3 | 0.1 | b |
| 3 | 0.2 | g |
| 4 | 2.0 | r |
| 4 | 3.1 | c |
| 3 | 0.8 | R |
| 1 | 0.2 | b |
| 1 | 1.2 | y |

# BLENDING ISSUES - MAIN GOALS

Main goal is introducing the above mentioned features in a native event stream platform, whereby:

- data is not finite and is *unknown*

- domain semantic changes over time

- processing logic might change over time

- applications evolve dynamically

- ...

RADICALBIT
radicalbit.io

GlobalKTable

LeftJoin

TASK 1

TASK 2

TASK n

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

| unsupervised_customers_1 | kmeans | mojo | ... |
|---|---|---|---|
| deep_net | deeplearning | pojo | ... |

```
{
    "modeld":
    "model_1",
    "eventData":
    "..."
}
```

```
{
    "modeld":
    "model_1",
    "eventData":
    "..."
}
```

```
{
    "modeld":
    "model_2",
    "eventData":
    "..."
}
```

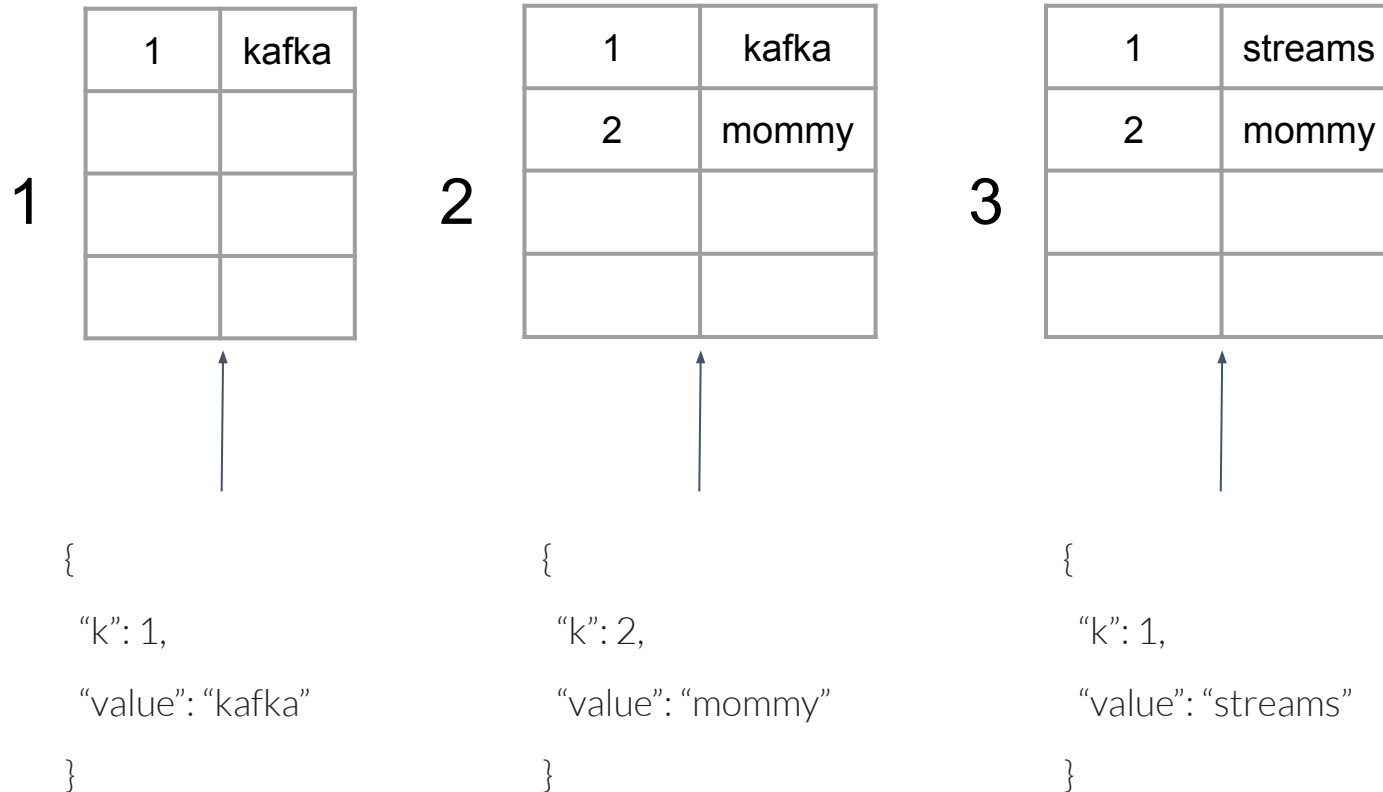# KSH$_2$O - KTABLE

- KTable is referring to a stream as of a **changelog**

| | 1 | kafka |
|---|---|---|
| **1** | | |
| | | |
| | | |

| | 1 | kafka |
|---|---|---|
| | 2 | mommy |
| **2** | | |
| | | |

| | 1 | streams |
|---|---|---|
| | 2 | mommy |
| **3** | | |
| | | |

```
{
  "k": 1,
  "value": "kafka"
}
```

```
{
  "k": 2,
  "value": "mommy"
}
```

```
{
  "k": 1,
  "value": "streams"
}
```

# KSH$_2$O - KTABLE

- GlobalKTable is a KTable that is global in terms of **topic supervision**

**KTable**

( 1, kafka )

| 1 | kafka |
|---|-------|
|   |       |
|   |       |
|   |       |

Partition 1

→

APP 1

| 1 | streams |
|---|---------|
| 2 | mommy   |
|   |         |
|   |         |

Partition 2

( 2, mommy )

( 1, streams )

APP 2

- GlobalKTable is a KTable that is global in terms of **topic supervision**

**KTable**

( 1, kafka )

| 1 | kafka |
|---|---|
|   |   |
|   |   |
|   |   |

APP 1

Partition 1

Partition 2

( 2, mommy )

( 1, streams )

| 1 | streams |
|---|---|
| 2 | mommy |
|   |   |
|   |   |

APP 2

**GlobalKTable**

( 1, kafka )

Partition 1

| 1 | streams |
|---|---|
| 2 | mommy |
|   |   |
|   |   |

APP 1

Partition 2

( 2, mommy )

( 1, streams )

| 1 | streams |
|---|---|
| 2 | mommy |
|   |   |
|   |   |

APP 2

Why don't partitioning accordingly model stream and data stream?

# KSH$_2$O - DEMO

# KSH$_2$O

## OML tools

- ## Apache SAMOA
  Large-Scale Learning from Data Streams with Apache SAMOA, 2018
  Nicolas Kourtellis, Gianmarco De Francisci Morales, and Albert Bifet

- ## side ML libraries on Apache Flink, Apache Spark, Apache Storm

| Multi-class problem GENERATED DATASET | KS - OML Accuracy | Python "batch" implementation Accuracy |
|---|---|---|
| First RUN | - | 0.905 |
| Second RUN | 0.99749374 | 0.99 |
| Third RUN | 0.9987469 | 0.995 |

RADICALBIT
radicalbit.io