# Stream Processing and Big Data Platforms

*Hong-Linh Truong*
*Department of Computer Science*
*linh.truong@aalto.fi, https://rdsea.github.io*

CS-E4640 Big Data Platforms, Spring 2025, Hong-Linh Truong
05/03/2025

**A"**

Aalto University
School of Science

# Learning objectives

- **Understand fundamental concepts and techniques in stream processing in big data**
- **Able to design streaming analytics in big data platforms and applications**
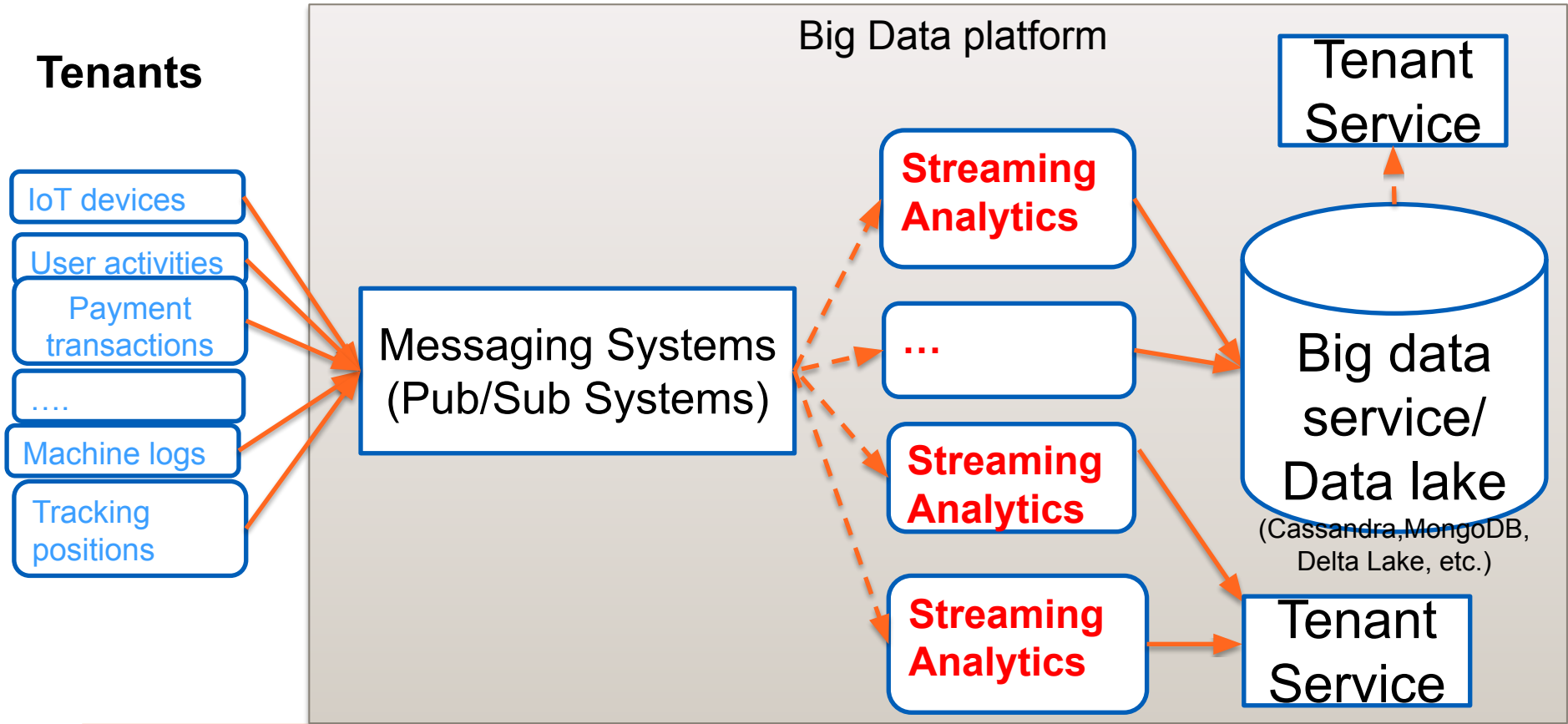- **Able to select and use common stream processing frameworks**

# Big data at large-scale: the big picture in this course



Today

**Data sources** (sensors, files, database, queues, log services)

**Messaging/Ingest systems** (e.g., Kafka, Pulsar, AMQP, MQTT, NATS, Kinesis, Nifi, Google PubSub, Azure IoT Hub)

**Stream processing/ML systems** (e.g. Flink, Kafka KSQL, Spark, Redpanda, Google Dataflow, Azure Stream)

**Analytics/ML Systems** (e.g., Azure Synapse Analytics, BigQuery, Redshift, ClickHouse)

**Operation/Management/ Business Services**

**Storage/Database/Data Lake** (S3, Minio, HDFS, DuckDb, CockroachDB, Cassandra, MongoDB, Elastic Search, Chroma, Weaviate, InfluxDB, Druid, Hudi, Iceberg, Delta Lake, etc.)

**Batch data processing/Distributed Query/ML systems** (e.g., Hadoop, Airflow, Spark, Presto)

**Elastic Edge-Cloud Infrastructures** (VMs, dockers, Kubernetes, OpenStack elastic resource management tools, storage)

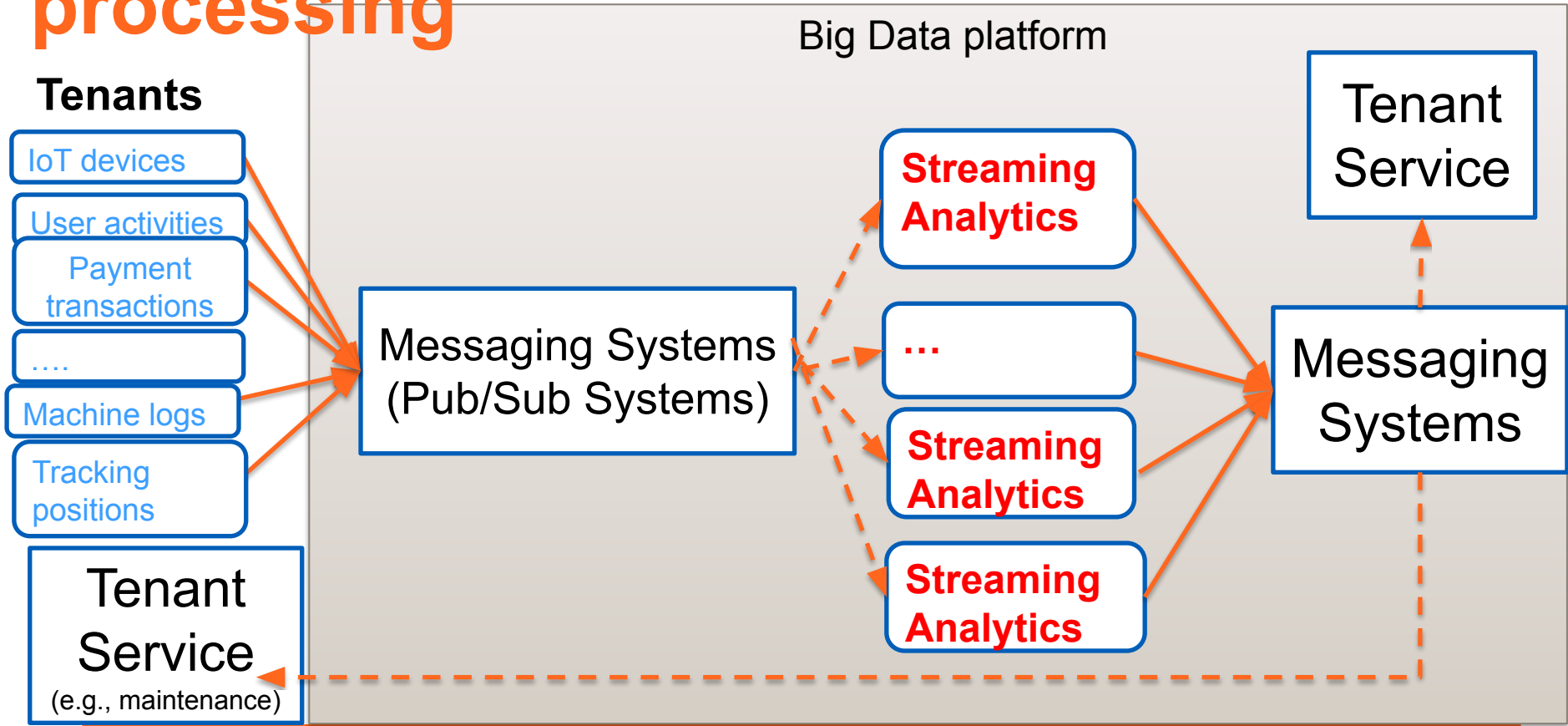# Stream analytics for data in motion

# Stream processing in big data

- **Big data coming from streams at near real-time**
  - the data element/unit may be "small" but voluminous and delivered in a near real-time manner
  - high and volatile throughput, but low processing time expected
  - not just *"take a record and store it into a database"*
- **Require large-scale computing infrastructures and many other platform services**
  - *task parallelism*: multiple tasks for processing data
  - *data parallelism*: data is partitioned into concurrent/parallel data streams ⇒ distributed, parallel processing tasks
  - *stateful analytics:* processing needs state information across multiple data records and time

# Near real-time streaming  data processing

# Near real-time streaming data processing



**Big Data platform**

Tenants

- IoT devices
- User activities
- Payment transactions
- ....
- Machine logs
- Tracking positions

Tenant Service
(e.g., maintenance)

Messaging Systems
(Pub/Sub Systems)

**Streaming Analytics**

...

**Streaming Analytics**

**Streaming Analytics**

Messaging Systems

Tenant Service

# Example in the cloud – Azure stream analytics for stream processing and big data platforms



**Figure source:** https://docs.microsoft.com/en-us/azure/stream-analytics/stream-analytics-introduction

Known public cloud services: Amazon Kinesis, Google Dataflow, Alibaba Cloud DataHub

# From complex event processing (CEP) in the age of enterprise computing



Our practices focus on modern technologies like: Apache Flink, Apache Spark and Arroyo, which are used intensively in business systems and big cloud platforms

# Stream processing and big data platforms

- **Stream processing services as a component of data platforms**
  - a big data technology for pre-processing, ingestion and high-level analytics, including near-real time machine learning
- **Stream processing services as data platforms**
  - a big data platform offers mainly stream processing services for streaming analytics
  - analytics on the fly as the first class feature
  - e.g., IoT analytics, e-commerce user activities, fraud detection, real time AI/ML

# Stream Processing – key concepts

# Common building blocks

- **The way to connect data to streams and obtain data records (messages) from the streams**
  - focusing very much on *connector concepts* and well-defined message structures (JSON, Avro, customized binary format, etc.)
  - connectors implement complex data handling mechanisms (low level session management, message retainment, delivery quality of service)
- **The way to specify/program the "analytics" logic**
  - *analytics functions, statements* and how they are glued together to process flows of messages
  - high-level, easy to use
- **The distributed engine to process analytics tasks**
  - handle complex task processing atop multiple compute nodes
- **The way to push the result to external components (sink databases, new streams, files)**

# Data stream programming
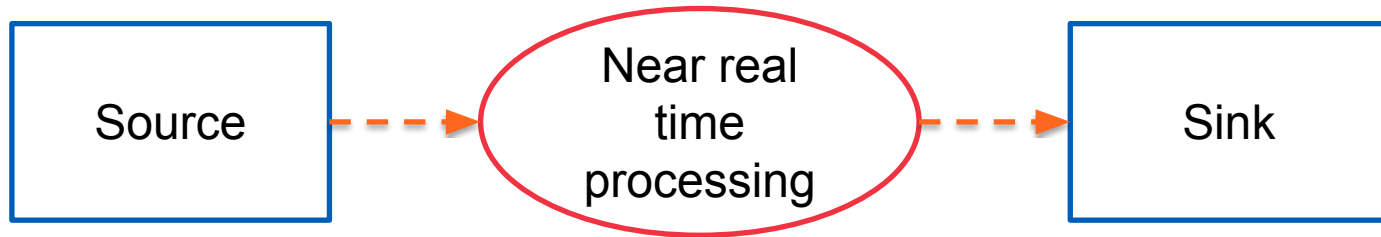
Data stream: a sequence/flow of data units

Data units are defined by applications: a data unit can be data described by a primitive data type or by a complex data type, a serializable object, etc.

Streaming data: produced by (near)realtime data sources as well as (big) static data sources ⇒ *unbounded* and *bounded*

- Examples of data streams
  - Continuous media (e.g., video for video analytics)
  - Discrete media (e.g., stock market events, twitter events, system monitoring events, comments, notifications, log records)

# Messages of events/data records

- **Messages encapsulating real-world events, data records and other types of data**
- **Data to be sent/processed can be in a simple or complex structure**

Source → Near real time processing → Sink

**We focus on unbounded discrete messages of data**

Aalto University
School of Science

# Message representations and streams

- **Data Sources**
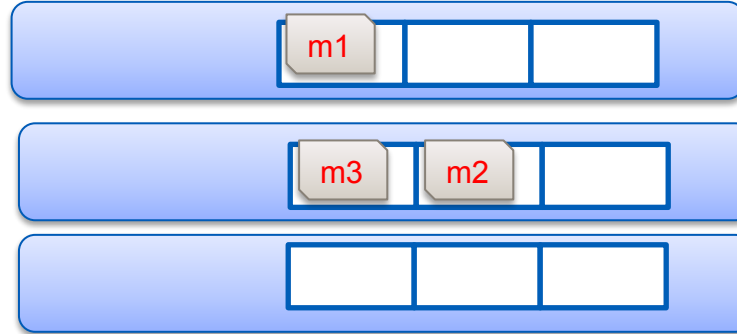  - via message brokers, databases, websocket, different IO adapters/connectors, etc.
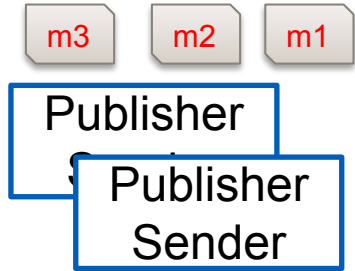- **Data Sinks**
  - messaging systems, databases, file storage/systems (S3, HDFS), etc.
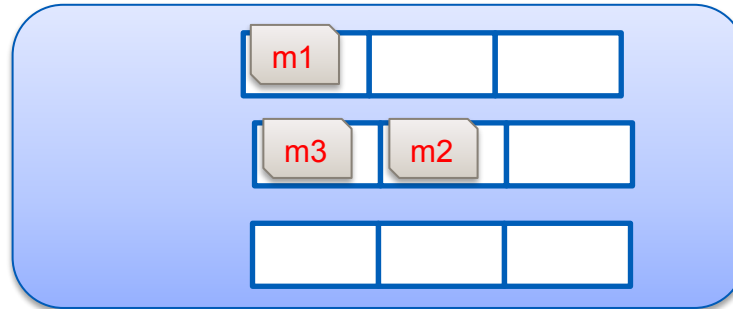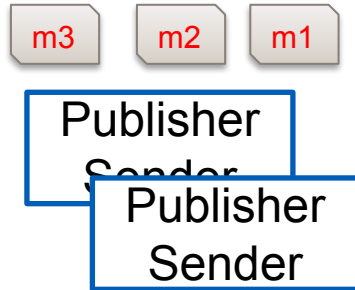- **Data representations**
  - text/CSV, JSON, Arvo format, etc.
  - serialization and deserialization (short name: SerDe) are required
  - data format validation
  - data schema registry for registered schemas

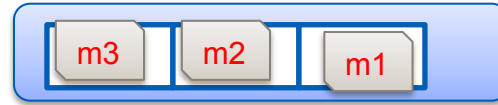# Publisher view: how messages are published

Messaging system

m3  m2  m1

Publisher Sender

Publisher Sender

| m1 | | |
| m3 | m2 | |
| | | |

**topic=queue; no partition**

m3  m2  m1

Publisher Sender

Publisher Sender

| m1 | | |
| m3 | m2 | |
| | | |

Topic and topic partitions
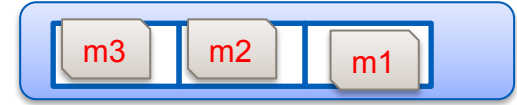
**topic = n partitions = n queues**

Aalto University
School of Science

# Handling messages for consumption (processing)

Messages in systems

simple straight forward



different mechanisms for "routing"

fant-out/broadcast

complex mapping/routing

**messages for consumption**
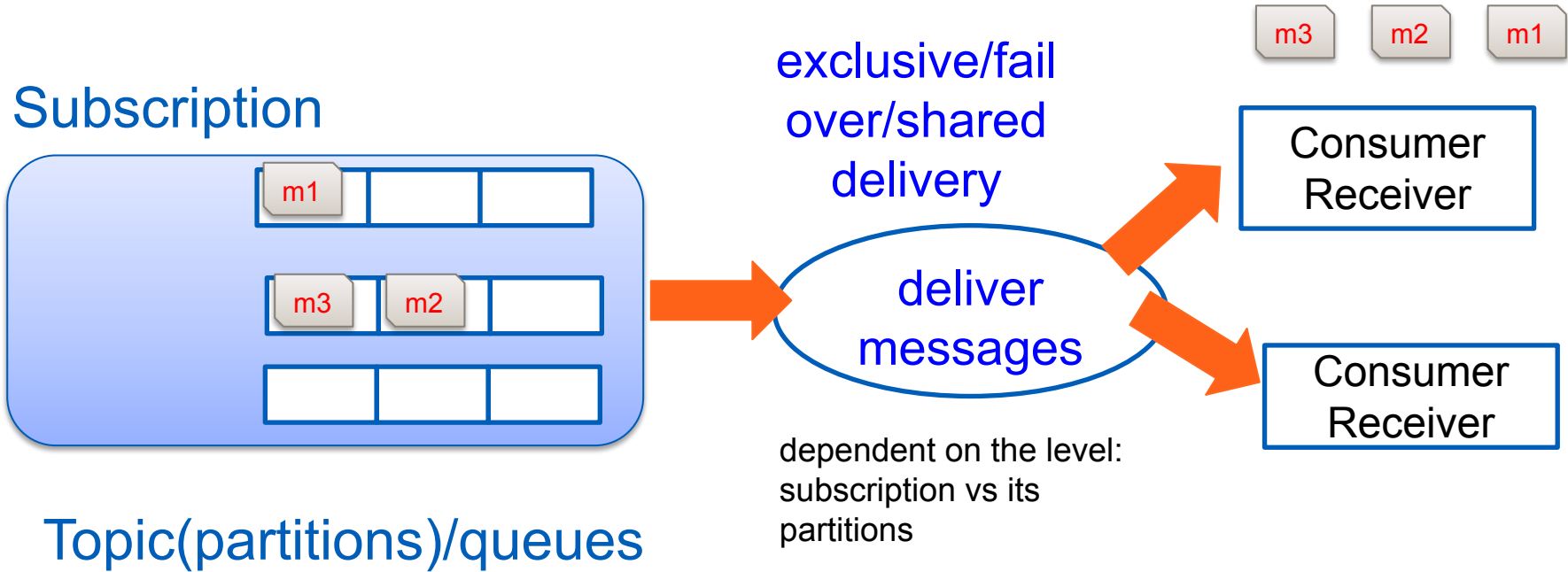
allowing parallel processing of the same messages

parallel processing of different messages

# Consumer view in accessing messages: subscription and delivery



Subscription

exclusive/fail over/shared delivery

deliver messages

dependent on the level: subscription vs its partitions

m3   m2   m1

Consumer Receiver

Consumer Receiver

m1

m3   m2

Topic(partitions)/queues

# Some key issues

- **Data order & delivery**
  - late data, out of order data
- **Times associated with messages and processing**
- **Data parallelism**
  - key-based data processing
- **Task parallelism**
  - stateful vs stateless processing

# Key issues in streaming data ingestion

Consumer/Ingestion Pipeline

Data producer



End-to-end delay

Message arrival orders

Aalto University
School of Science

**Without a timestamp associated to a message, do we know the delay or out of order?**
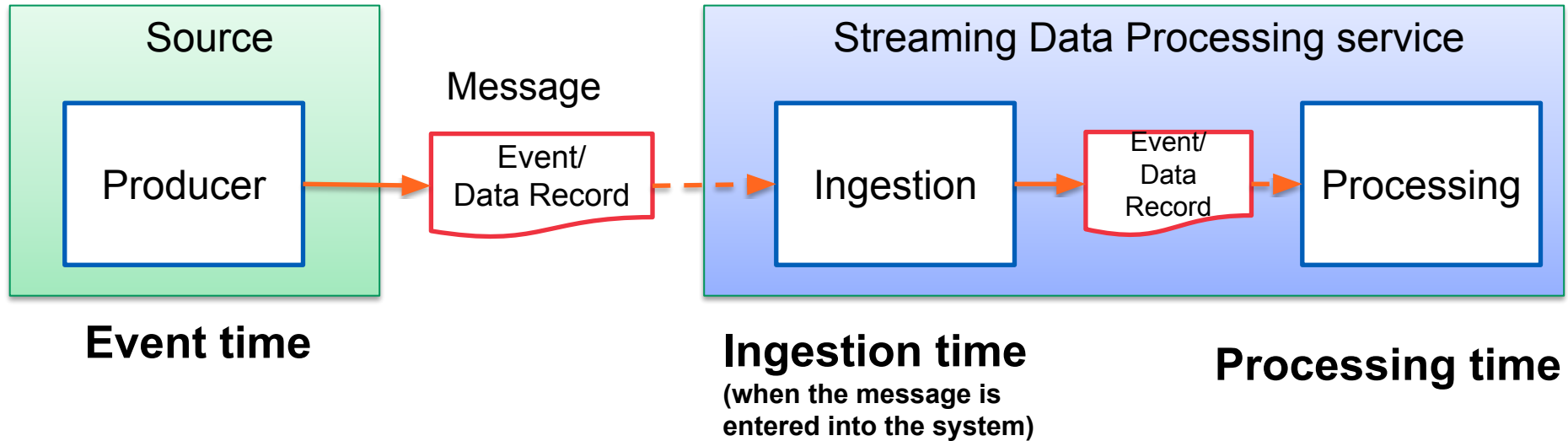
**What is the consequence of delay/out of order for processing?**

Aalto University
School of Science

# Key issues in streaming data: the notion of times

**Times associated with data and processing**



**Which time is important for analytics (from business viewpoint)?**

# Data parallelism: partition stream data based on some keys for analytics

Stream IN

IF...

m3 m2 m1

m1

m3 m2

Stream OUT

Can be processed in parallel using multiple nodes

With **keyed data**: enable parallel processing based on the keys

```
12  */
13  public class BTSAlarmEvent {
14      public String station_id;          keys?
15      public String datapoint_id;
16      public String alarm_id;
17      public Date event_time;
18      public float value;
19      public float valueThreshold;
20      BTSAlarmEvent() {
21
```

# Windows of data

Window is used to group data for processing:

Which constraints are used to determine a window?

m1000 .. .. m2 m1 ... m4

**a stream of events**

sliding/tumble window size: period of time or number of events/records

Arrival order

# Windowing

- **Windows size:**
  - time or number of records

- **Tumbling window:**
  - identified by size, no gap between windows

- **Sliding window:**
  - identified by size and a sliding internal

- **Session Window:**
  - identified by "gap" between windows (e.g., the gap of events is used to mark "sessions")

# Functions applied to Windows of data

If we

- specify a set of conditions ⇒ windows will be created according to the conditions to store message in corresponding windows

then we can

- Apply functions to messages in  the window that match these conditions

Task parallelism: we can have a lot of such functions executed in parallel in multiple compute nodes

# Functions

- **Can be simple or complex!**
  - built-in and user-defined functions
- **Core for analytics and ML**
- **Examples**
  - individual threshold/alarm based alerting, atypical events monitoring
  - data rollup
  - anomaly detection based on statistical functions, like quantile/T-digest, ...
  - real time AI/machine learning

# Example

**Monitoring working hours of (taxi/truck) drivers (assume events about pickup/drop captured at near real-time):**

- Windows: 12 hours

- Partitioning data/Keyed streams: licenseID

- Function: determine working and break times and check with the law/regulation

**Source:**
https://www.infoworld.com/article/3293426/how-to-build-stateful-streaming-applications-with-apache-flink.html

What if events/records come late into the windows?

Do we need to deal with late, out of order events/records?

*correctness and completeness* issues

Aalto University
School of Science

# Support lateness

- **Identify timestamp of events/data records**
- **Identify watermark in streams**
  - a watermark is a timestamp
  - a watermark indicates that no events which are older that the watermark should be processed
  - enable the delay of processing functions to wait for late events
- **Using watermark to ignore late data ⇒ computing under "incompleteness assumption"**

| m.. | m.. | w | m3 | m2 | m1 |

**watermark**

**Aalto University
School of Science**

# Delivery guarantees

**Exactly once? at least once? or at-most-once**

**End-to-end?**

**message delivery**

**result delivery**

m[e1,e2,..en]

r([e1,e2,..en])

Stream IN → Stream processing → Stream OUT

**What if the stream processing fails and restarts**

# Message and processing guarantees

- **Message guarantees are the job of the broker/messaging system**

- **Processing guarantees are the job of the stream processing frameworks**

- **They are highly connected if messaging systems and processing frameworks are tightly coupled (e.g., Kafka case)**

**Aalto University**
**School of Science**

# End-to-end exactly once

- **Exactly once for processing is not enough**
- **Messaging systems must support**
  - redeliver messages/data, recoverable data
- **Sink and output must support exactly one**
  - idempotent results, roll back
- **Coordination among various components**

**Further reading:**
**https://flink.apache.org/features/2018/03/01/end-to-end-exactly-once-apache-flink.html**
**https://www.confluent.io/blog/simplified-robust-exactly-one-semantics-in-kafka-2-5/**
**https://docs.microsoft.com/en-us/azure/hdinsight/spark/apache-spark-streaming-exactly-once**

# Performance metrics



**Response time**

Job submission

Batch processing

(e.g., MapReduce/Hive, Spark SQL)

Job completion

m3 → Stream processing → r(m3)

**Service time/latency**

Aalto University
School of Science

# Latency and throughput

- **Service latency**
  - subseconds! e.g., milliseconds
  - max, min or percentile? $\Rightarrow$ up to application requirements
- **Throughput**
  - how many messages can be processed per second?
- **Goal: low latency and high throughput!**

# Structure of streaming data processing programs (1)

- **We have multiple streams of data, different functions for processing data, multiple computing nodes**

- **Data exchange between tasks**
  - links in task graphs reflect data flows

- **Stream processing**
  - centralized or distributed (in terms of computing resources)
  - simple functions vs complex ones

# Structure of streaming data processing programs (2) - examine a simple example

```
124
12
126         WAIT AND PROCESS DATA
127         '''
128     while True:
129         '''
130         Receive the data from source
131         '''
132         msg = consumer.receive()
133         '''
134         when should we do this?
135         consumer.acknowledge(msg)
136         '''
137         try:
138             '''
139             MAIN TRANSFORMATION, HERE IS WITH A FUNCTION
140             '''
141             ## assume that the selected data schema is json
142             result =dt_process_json_style(msg,op_processor)
143             ##store the result to the right data sink
144             dt_store_to_sink(result)
145
146         except Exception as ex:
147             logging.warn(f'{ex}')
148             logging.info("Continue to wait")
149
```

**How to handle possible errors**

Note: Example with an external Pulsar consumer for data transformation

# Structure of streaming data processing programs (3)



```
119  //     parse the data, determine alert and return the alert in a json string
120  //     Apply function on data stream
121        DataStream<String> alerts = btsdatastream
122  //          .flatMap(new BTSParser()
123              .flatMap(new BTS_Trend_Parser()
124              /*
125                  Another example is to have:
126                  new FlatMapFunction<String, BTSAlarmEvent>() {
127                      @Override
128                      public void flatMap(String valueString, Collector<BTSAlarmEvent> out) {
129                          String[] record = valueString.split(",");
130                          ....
131                          out.collect(...);
132                      }
133              })
134              */
135              )
136              //.setParallelism(5)    // uncomment this line to scale the Parser stream and set the valu
137              .keyBy(new AlarmKeySelector()
138              /* another way is to have:
139                  new KeySelector<BTSAlarmEvent, String>() {
140                      public String getKey(BTSAlarmEvent btsalarm) { return btsalarm.station_id; }
141                  }
142              */
143              )
144              .window(SlidingProcessingTimeWindows.of(Time.seconds(60), Time.seconds(5)))
145  //          .window(SlidingEventTimeWindows.of(Time.minutes(5), Time.seconds(5))) // set the window si
146  //          .process(new MyProcessWindowFunction()).setParallelism(1);
147              .process(new TrendDetection()).setParallelism(1);
148  //.setParallelism(5);  // uncomment this line to scale the stream processing and set the value for
```
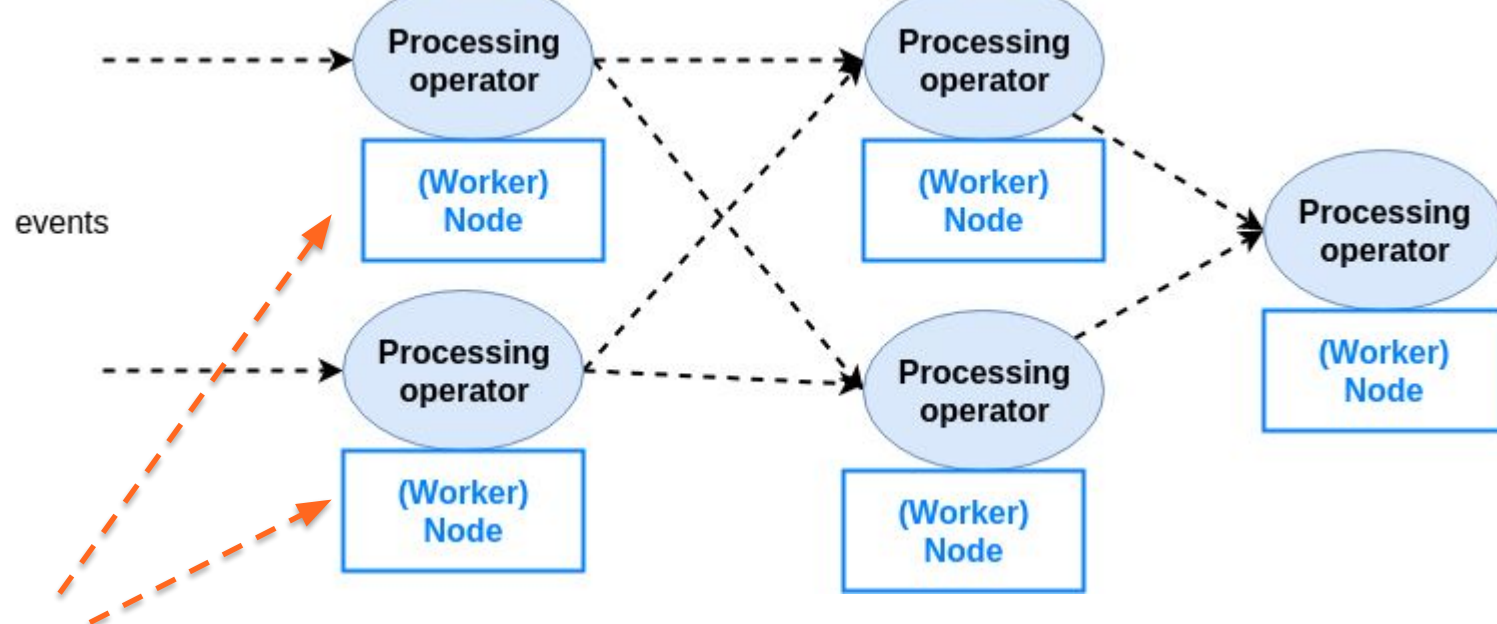
- **Dataflows:**
  - *Data source operators: represent  sources of streams*
  - *Processing operators:  represent processing functions*

**Aalto University
School of Science**

# Distributed processing topology in a cluster

**A graph of tasks (running operators); all tasks are running**



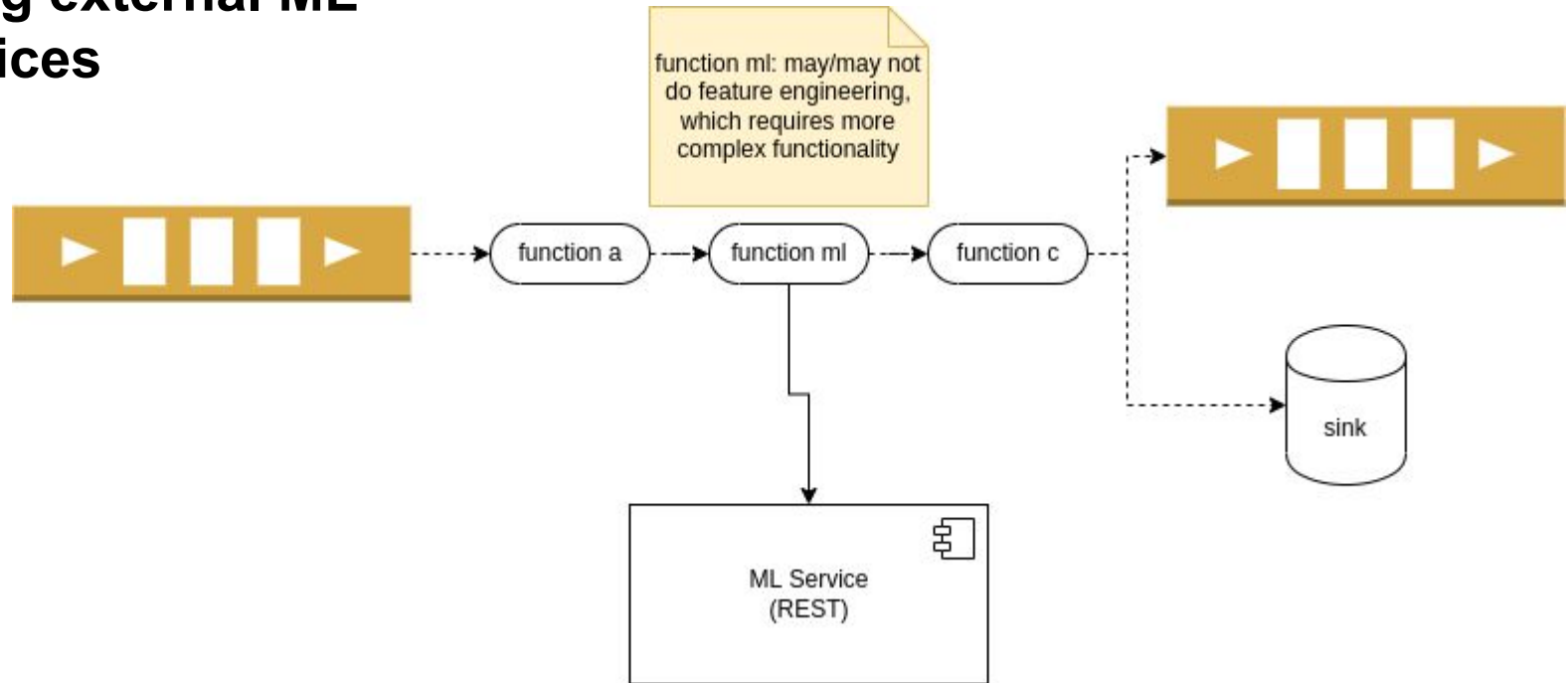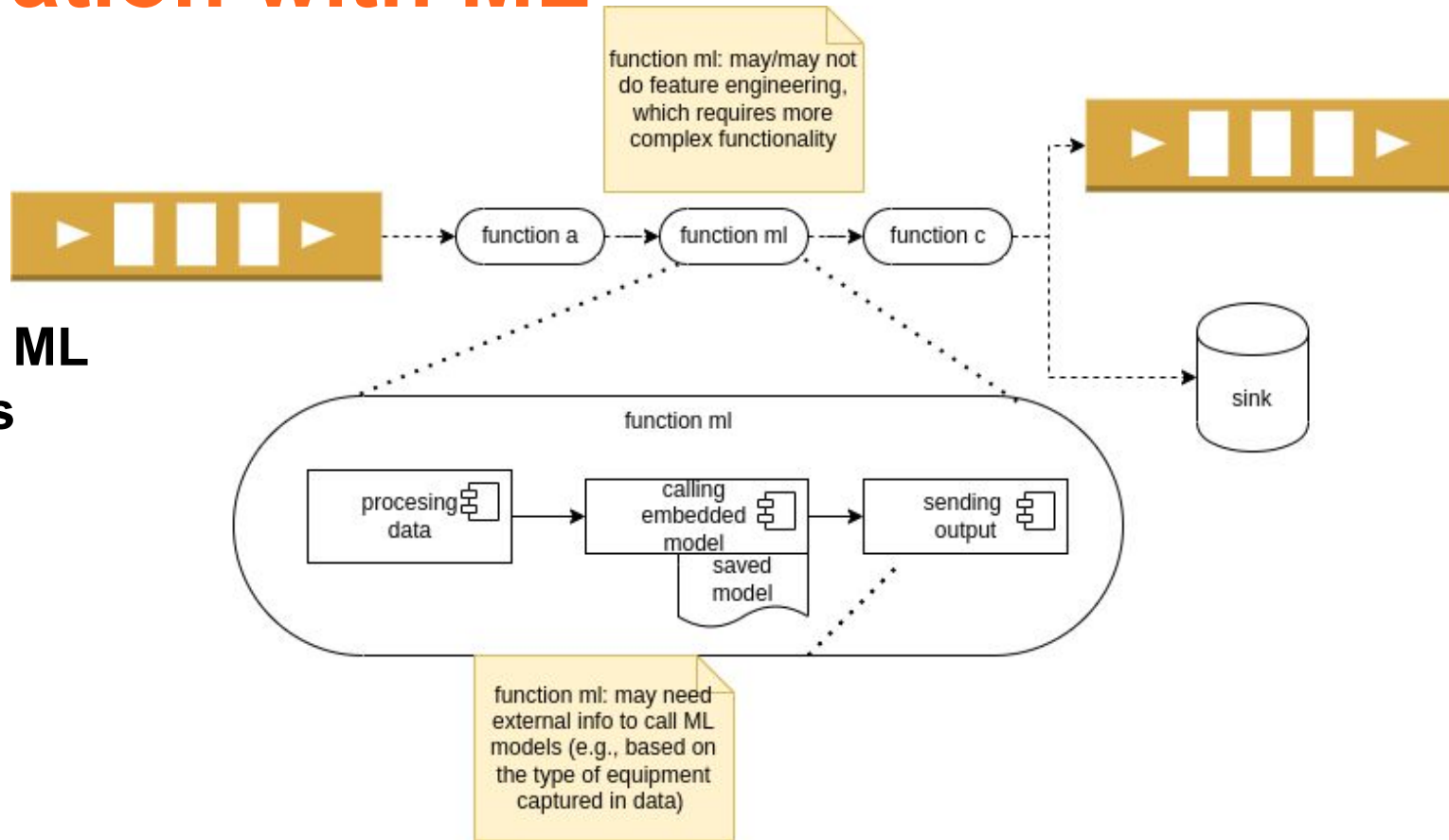**Nodes of a cluster (VMs, containers, Kubernetes)**

Aalto University
School of Science

# Integration with ML

**From outside**

Aalto University
School of Science

# Integration with ML

**Using external ML services**

# Integration with ML

**Embedded ML capabilities**



function ml: may/may not do feature engineering, which requires more complex functionality

function a → function ml → function c

sink

function ml

procesing data → calling embedded model → sending output

saved model

function ml: may need external info to call ML models (e.g., based on the type of equipment captured in data)

**Aalto University**
**School of Science**

# Distributed, composable processing topologies in cross distributed sites

# Common concepts in existing frameworks - programming level

- **How to write streaming program?**
- **With programming languages**
  - low level APIs
  - DSL
  - Java, Scala, Python (Spark, Flink, Kafka)
- **High-level data models**
  - KSQL
- **Flow/pipeline description**
  - Node-RED/GUI-based flow editors

# Common concepts in existing frameworks - key common concepts

- **Abstraction of streams**
- **Connector library for data sources/sinks**
  - very important for application domains
- **Runtime elasticity**
  - add/remove (new) operators
  - add/remove underlying computing nodes
- **Fault tolerance**

Aalto University
School of Science

# Where do you find most of concepts that we have discussed

- **Apache Storm**
    - *https://storm.apache.org/*

- **Apache Spark  (Structured Streaming)**
    - *https://spark.apache.org/*

- **Apache Kafka Streams and KSQL**
    - strongly bounded to Kafka messaging

- **Apache Flink (Stream Analytics)**
    - native, clustered, better data sources/sinks

- **Apache Beam (https://beam.apache.org/)**
    - unifying programming models for batch and stream processing

# Practical learning paths

- **Path 1: if you don't have a preference and need challenges**
  - Apache Flink Stream API (e.g., with RabbitMQ/Kafka connectors)
- **Path 2: many of you have worked with Kafka**
  - Kafka Streams DSL (everything can be done with Kafka)
- **Path 3: for those of you who are working with Spark (and Python is the main programming language)**
  - Apache Spark Structured Streaming
- **Path 4: for those who deal with MQTT brokers**
  - Apache Storm (but also Kafka, …): Spout and Bolt API or Stream API

# Summary

- **Focus:**
  - Practical programming with one of the stacks:
    - *Apache Flink Stream API (with different connectors)*
    - *Kafka Streams*
  - Check the common concepts in other tools/systems
- **Action:**
  - Work on use cases where you can use stream analytics (as a user/developer) $\Rightarrow$ there are many interesting analytics
  - Provision services for stream processing (as a platform)

# Thanks!

**Hong-Linh Truong**
**Department of Computer Science**

**rdsea.github.io**