



POLITECNICO DI MILANO
SCHOOL OF INDUSTRIAL AND INFORMATION ENGINEERING
ACADEMIC YEAR 2024-2025

Streaming Data Analytics

Professor: Emanuelle Della Valle

Last updated: November 4, 2024

**This document is intended for educational purposes only.
These are unreviewed notes and may contain errors.
Made by Roberto Benatuil Valera**

Contents

1	Introduction	5
1.1	Enabling Continuity	5
1.1.1	Sensors and Actuators	5
1.1.2	Connectivity	6
1.1.3	Streaming Data Engineering	6
1.1.4	Streaming Data Science	6
1.2	Importance of Streaming Data Analytics	7
2	A Tale of Four Streams	9
2.1	Types of data streams	9
2.2	Time models	9
2.2.1	Stream-only time model	9
2.2.2	Absolute time model	10
2.2.3	Interval-based time model	10
2.2.4	Trade-off: latency vs expressiveness	10
2.3	Architectures for stream processing	11
2.3.1	Event Based Systems (EBS)	11
2.3.2	Data Stream Management Systems (DSMS)	12
2.3.3	Complex Event Processing (CEP)	14
2.3.4	Event-driven Architecture	16
3	EPL and Esper	17
3.1	What is Esper?	17
3.1.1	Esper features	17
3.2	EPL in a nutshell	18
3.3	EPL syntax: fire alarms example	18
3.3.1	Creating a schema	18
3.3.2	Creating a simple query	19
3.3.3	Defining windows	20
3.3.4	Output stream	22
3.3.5	Complex event patterns	23
3.3.6	Composing queries	27
3.4	Joins in EPL	28
3.4.1	Stream to stream joins	28

3.4.2	Table to table joins	29
3.4.3	Stream to table joins	30
3.5	Contexts in EPL	31
3.5.1	Context by key	32
3.5.2	Context by start/end conditions	32
3.5.3	Context by time	33
3.5.4	Comparison: Windows vs. Contexts	33
4	Kafka and Spark	35
4.1	Introduction: scaling stream ingestion and processing	35
4.1.1	Latency vs throughput	35
4.1.2	Data/message	36
4.2	Apache Kafka	37
4.2.1	Kafka: high-level vs system view	37
4.2.2	Kafka: how it works	39
4.2.3	Distributed consumption	40
4.3	Apache Spark	40
4.3.1	Key concepts in Spark	41
4.3.2	Resilient Distributed Datasets (RDDs)	41
4.3.3	Spark at work	42
4.3.4	DataFrames in Spark	44
4.4	Spark Structured Streaming	44
4.4.1	Programming model	45
4.4.2	Creating streaming DataFrames	45
4.4.3	Writing streaming DataFrames	47
4.4.4	Window operations on Event Time	48
4.4.5	Joins	48
4.4.6	Late arrivals	49
4.4.7	Unsupported operations	51

Chapter 1

Introduction

In streaming data engineering, continuity of data analysis is a key concept. To process large amounts of data, we have two main approaches:

- Batch processing: this is the traditional approach where data is collected and stored, then processed in a batch to obtain insights.
- Stream processing: in this approach, data is processed as it arrives, allowing for real-time insights.

In some cases, batch processing is enough, like when we need to process data that is not time-sensitive, such as historical data. However, in many cases, we need to process data in real-time, like in fraud detection or monitoring systems like fire alarms. In these cases, stream processing is the way to go. But, how do we enable this continuity of data flow?

1.1 Enabling Continuity

To enable continuity of data processing, we need four main components:

- Sensors and actuators
- Connectivity
- Streaming data engineering
- Streaming data science

1.1.1 Sensors and Actuators

Sensors are devices that collect data from the environment, like temperature sensors, cameras, etc. Actuators are devices that can act on the environment, like turning on a light, opening a door, etc. These devices are the first step in the data processing pipeline. They are a way of giving the computers the means to gather data from the environment and act on it.

[...] sensor technology enable computers to observe, identify and understand the world - without the limitations of human-entered data. - Kevin Ashton, the brander of "Internet of Things"

1.1.2 Connectivity

This is the enabling and constraining factor of the Internet of Things. It is the way that sensors and actuators communicate with the rest of the system. This can be done through a variety of means, like WiFi, Bluetooth, etc.

In this case, the distance and the use case are the main factors that determine the type of connectivity that we need. For example, if we need to connect a sensor to a computer that is far away, we might need to use a cellular network, while if we need to connect a sensor to a computer that is close, we might use Bluetooth.

Also, bandwidth is a factor that we need to consider. If we need to send a lot of data, we might need a high bandwidth connection, while if we need to send a small amount of data, we might use a low bandwidth connection. Another important factor is the latency of the connection. If we need to send data in real-time, we need a low latency connection, while if we don't need to send data in real-time, we might use a high latency connection.

Note that most of the times, it comes down to a trade-off between these factors, so we need to choose the right connectivity for our use case.

1.1.3 Streaming Data Engineering

This is the main topic of this chapter. This area focuses on the systems and infrastructure that we need to manage various types of data streams efficiently. We have four main concepts:

- Event-based systems: to tame myriads of tiny flows of data
- Data Stream Management Systems (DSMS): to handle continuous massive flows of unstoppable data
- Complex Event Processing (CEP): to manage continuous numerous flows of data that can turn into a torrent
- Event-driven architecture: to tame the forming of an immense delta made of myriads of flows of any size and speed.

We will cover these concepts in more detail in the next sections.

1.1.4 Streaming Data Science

This is the area that focuses on the algorithms and techniques that we need to analyze data streams efficiently. Note that changes in the data stream can happen at any time, and they cause ML models to lose accuracy. This is why we need to use techniques to adapt to these changes. The three main components of this section are:

- Time series analytics: to explain the past and forecast the future of a continuous flow of data without assuming data independence

- Streaming machine learning: to learn one data at a time from a continuous flow of data without assuming identically distributed data
- Continual learning: to do a long-life learning from a sequence of experiences without forgetting past knowledge

This is the second topic of this course, and we will cover it in more detail in the next chapters.

1.2 Importance of Streaming Data Analytics

One idea: value is about time. Traditionally, data processing creates value by providing insights that are generated with months or years of data gathering. However, in many cases, we need to process data in real-time to create value.

Continuous streaming analysis creates value by providing insights that are generated within seconds or minutes, allowing for real-time decision making, in other words, reactive applications that can respond to events as they happen.

Streaming data analytics is important, as it can:

- Provide real-time insights
- Reduce costs
- Improve efficiency
- Create innovative products
- Generate new revenue streams

Chapter 2

A Tale of Four Streams

2.1 Types of data streams

Data streams come in many forms and sizes. They can be classified into four main categories for the sake of this course:

- Myriads of tiny flows that you can collect, for example, coming from physical or software sensors
- Continuous massive flows that you cannot stop, for example, from telecommunications and utilities monitoring
- Continuous numerous flows that can turn into a torrent, like physical or cyber alarms
- Myriads of flows of any size and speed forming an immense delta, like physical or software actuators

2.2 Time models

There are 3 main type of time models for streaming data analytics:

- Stream-only time model
- Absolute time model
- Interval-based time model

2.2.1 Stream-only time model

In this model, time is defined by the order of the events in the stream. This model is used when we only care about the order of the events, and not about the time that they happened.

Therefore, in this model there is no notion of time, so this limits the type of queries that we can do. This reduces the expresiveness of the model, but it also reduces the latency of the system.

2.2.2 Absolute time model

In this model, time is defined by the time that the events happened. This model is useful when we need to know the exact time that the events happened, not only the order in which they happened.

Therefore, in this model we can do more complex queries, like time-based queries or window-based queries. This increases the expressiveness of the model, but it also increases the latency of the system.

2.2.3 Interval-based time model

In this model, time is defined by intervals. This model is useful when we need to know the time that the events happened, but we also need to group the events in intervals.

Therefore, in this model we can do more complex queries, not only time-based queries, but also group the events in specific intervals. This increases the expressiveness of the model, but it also increases the latency of the system.

2.2.4 Trade-off: latency vs expressiveness

There is a trade-off between latency and expressiveness. Let us see the following graph:

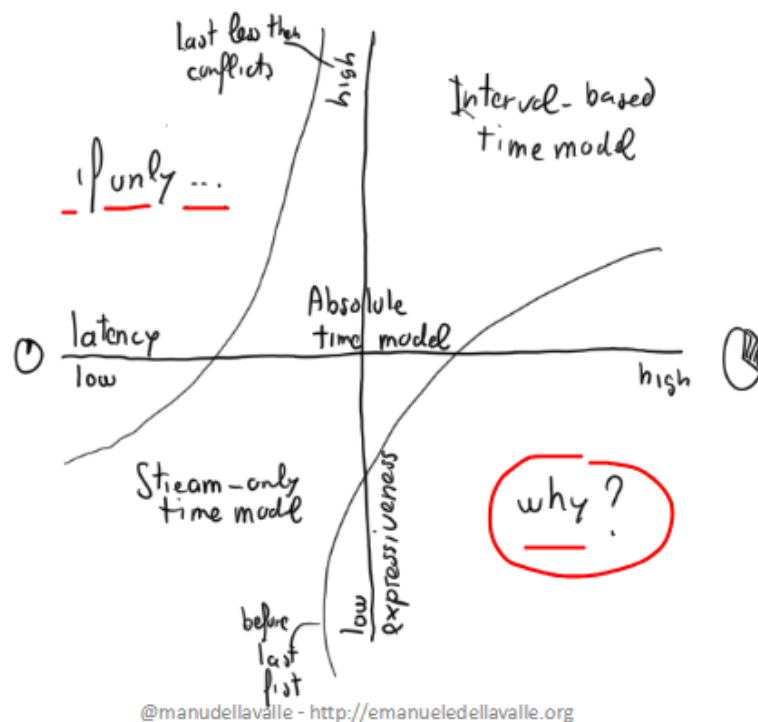


Figure 2.1: Trade-off between latency and expressiveness

2.3 Architectures for stream processing

There are four main architectures for stream processing:

- Event-based systems
- Data Stream Management Systems
- Complex Event Processing
- Event-driven architecture

2.3.1 Event Based Systems (EBS)

An Event-Based System is a software architecture paradigm in which data flows and operations are driven by events rather than following a pre-defined control flow. In other words, in an EBS, actions are executed in response to specific events that occur, rather than as a part of a predetermined logical sequence.

Two main concepts:

- Event: an immutable and append-only stream of "business facts" that are generated by sensors, applications, or users.
- Decoupling: means decentralizing the freedom to act, adapt and change

There are two main types of event-based systems:

- Content-based systems
- Topic-based systems

Content-based systems

Content-based systems filter and route messages based on the actual content within the messages. In this model, subscribers specify their interest through content attributes, and the system delivers only messages that match these attributes.

This type of systems are mainly used as an academic effort, as they are hard to implement in practice.

Topic-based systems

Topic-based systems organize messages by predefined topics. In this model, publishers send messages to a specific topic, and subscribers receive all messages sent to the topic they have subscribed to.

This type of systems are the most common in practice, and they are widely available for industrial applications. Some examples are Apache Kafka, MQTT, among others.

Kafka: "the" event-based system

Apache Kafka is a distributed event streaming platform that is capable of handling trillions of events a day. It is used by thousands of companies to build real-time streaming data pipelines and applications.

In this system, topics are partitioned into kafka brokers. Producers share messages over the partitions on a certain topic via hash partitioning or round-robin partitioning. Different consumers can read messages from the same topic, and partitions guarantee parallel reads.

2.3.2 Data Stream Management Systems (DSMS)

Data Stream Management Systems (DSMS) are systems that are designed to handle continuous massive flows of unstopable data. They are used to process and analyze data streams in real-time.

Data streams are sequences (sometimes unbounded) of time-varying data elements. They represent an almost continuous flow of information, with the most recent information being more relevant, as it describes the current state of a dynamic system.

The nature of streams require a paradigmatic change, from persistent data (one-time semantics) to transient data (continuous semantics).

DSMS semantics

This syem uses continuous queries registered over streams that are observed through a window, like the following picture:

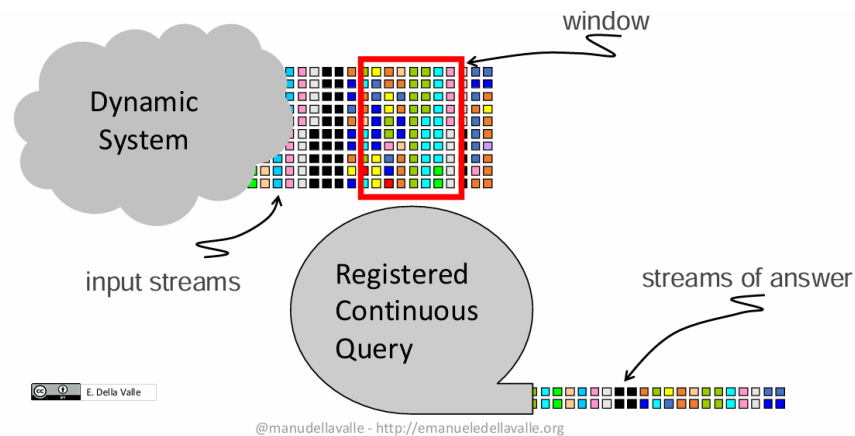


Figure 2.2: DSMS semantics

We can do different operations over the data stream. For example, we can do incremental operations, like:

→ Filtering (more generally, mapping)

- Grouping
- Count, sum, min, max, means
- Variance and standard deviation

Specially, for variance calculation, we can use the Welford's algorithm, which is an online algorithm for calculating the variance of a sample. It is more numerically stable than the standard formula for variance, and it is also more efficient. The algorithm is as follows:

Algorithm 1 Welford's algorithm

```

1:  $n \leftarrow 0$ 
2:  $\mu \leftarrow 0$ 
3:  $M2 \leftarrow 0$ 
4:
5: procedure UPDATE( $x$ )
6:    $n \leftarrow n + 1$ 
7:    $\delta \leftarrow x - \mu$ 
8:    $\mu \leftarrow \mu + \frac{\delta}{n}$ 
9:    $\delta_2 \leftarrow x - \mu$ 
10:   $M2 \leftarrow M2 + \delta \cdot \delta_2$ 
11: end procedure

```

However, there are some operations that we are not able to exactly calculate with continuous semantics, like the median, mode, top-k, distinct count, etc. For some of these operations, we can use approximations.

Window types

There are three main types of windows in DSMS:

- Sliding window: a window that slides over the stream
- Tumbling window: a window that is fixed in time
- Session window: a window that groups events that are close in time

Note that we are not specifying the time model. This is because DSMS can work with any time model, as long as it is specified in the query. For example, if we are working with the stream-only time model, we define physical windows. E.g.: for a tumbling window, every 10 events; for a sliding window, consider the last 10 events every 5 events; Note that session windows cannot be specified.

If we consider an absolute time model, we define logical windows. E.g.: for a tumbling window, every 10 seconds; for a sliding window, consider the last 10 seconds every 5 seconds; for a session window, every 10 seconds of inactivity.

An example of a windowed aggregation query in DSMS is the following:

```

1 SELECT STREAM
2     HOP_END(rowtime, INTERVAL '1' HOUR, INTERVAL '3' HOUR),
3     COUNT(*),
4     SUM(units)
5 FROM Orders
6 GROUP BY
7     HOP(rowtime, INTERVAL '1' HOUR, INTERVAL '3' HOUR);

```

Streaming joins

A streaming join is an operation that combines two or more streams based on a common key or specific attribute.

In traditional databases, combining rows from different tables is a basic operation, but in stream processing is hard because the data is not static but in constant motion.

Streams are data that flows over time. Therefore, joins consider the temporal synchronization of the data via time windows to limit the join to events that occur within a certain time interval.

An example:

```

1 SELECT STREAM o.rowtime o.productId, o.orderId,
2                 s.rowtime AS shipTime
3 FROM Orders AS o JOIN Shipments AS s
4     ON o.orderId = s.orderId AND
5         s.rowtime BETWEEN
6             o.rowtime AND
7             o.rowtime + INTERVAL '1' HOUR;

```

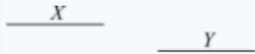

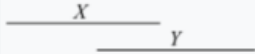
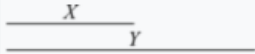
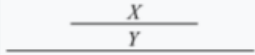
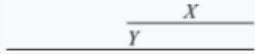
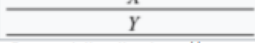
2.3.3 Complex Event Processing (CEP)

Complex Event Processing (CEP) is a method of tracking and analyzing streams of information from different sources to identify patterns and relationships. It is used to process and analyze data streams in real-time.

CEPs add the ability to deploy rules that describe how composite events can be generated from primitive ones. Typical CEP rules search for sequences of events that match a pattern, for example: Raise C if $A \rightarrow B$. Time is a key factor in CEP, as it is used to determine the order of events and the time window in which the events must occur.

CEP semantics

Complex Event Processing (CEP) uses subsets of Allen's interval algebra to define the temporal relationships between events. The main operators are:

Relation	Illustration	Interpretation
$X < Y$ $Y > X$		X takes place <u>before</u> Y
$X m Y$ $Y mi X$		X meets Y (<i>i</i> stands for <i>Inverse</i>)
$X o Y$ $Y oi X$		X overlaps with Y
$X s Y$ $Y si X$		X starts Y
$X d Y$ $Y di X$		X during Y
$X f Y$ $Y fi X$		X finishes Y
$X = Y$		X is equal to Y

@manudellavalle - <http://emanueledellavalle.org>

Figure 2.3: Allen's interval algebra

Event Processing Language (EPL)

Event Processing Language (EPL) is a language that is used to define rules in CEP. It has the capability to describe complex event patterns and relationships, to create queries that monitor real-time data streams.

It is supported by many CEP engines. The main engine we are going to use in this course is Esper, although there are others like Oracle CEP.

An example of a query in EPL is the following:

```

1 insert into alertIBM
2 select *
3 from pattern [
4     every (
5         StockTick(name="IBM", price > 100)
6         ->
7         (StockTick(name="IBM", price < 100)
8           where timer:within(60 seconds))
9     )
10 ];

```

This query alerts on each IBM stock tick with a price greater than 100, followed by a tick with a price less than 100 within 60 seconds.

2.3.4 Event-driven Architecture

Event-driven architecture is a software architecture paradigm promoting the production, detection, consumption of, and reaction to events. An event can be defined as a significant change in state.

In this paradigm, the system only worries about writing the events to the event log, and the rest of the system can react to these events. This allows for a more decoupled system, where each consumer can read react to the events asynchronously.

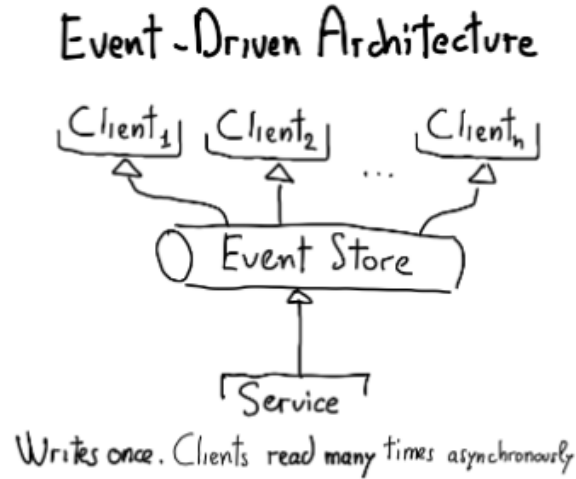


Figure 2.4: Event-driven architecture

EDAs have two main functionalities:

- Notifications, which is just a stream of observations
- Data replication, which is a stream of changes

Chapter 3

EPL and Esper

3.1 What is Esper?

Esper is an open-source software product for Complex Event Processing (CEP) and Stream Analytics, developed by EsperTech. It turns a large amount of incoming event series or streams into actionable insights in real-time. It is a lightweight, high-performance, and scalable engine that can process millions of events per second.

Esper provides a SQL-like language called EPL (Event Processing Language) to define queries and patterns over event streams. It also provides a Java API to integrate with Java applications.

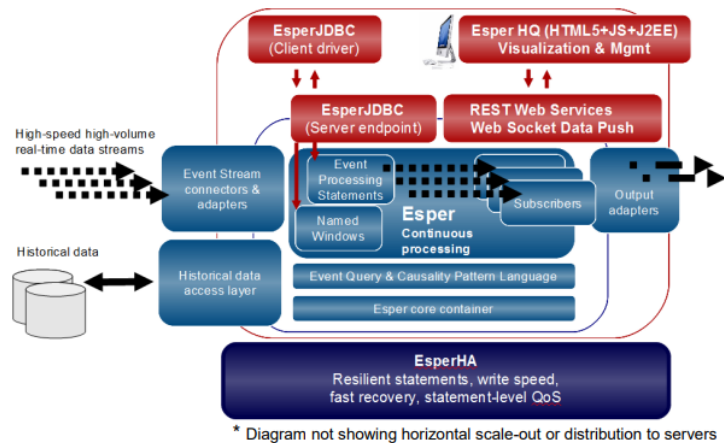


Figure 3.1: Esper Architecture

3.1.1 Esper features

- It is implemented as a Java library, so it can be easily integrated with Java applications.
- It provides a SQL-like language called EPL to define queries and patterns over event streams.
- it is designed for performance and scalability, and it can process millions of events per second (high throughput and low latency).

- It can interact with static/historical data, and it can be used to correlate real-time data with historical data.
- It has a configurable push and pull communication.

3.2 EPL in a nutshell

EPL is a SQL-like language that allows you to define queries and patterns over event streams. It is used to define the logic that Esper uses to process events.

As in the DSMS approach, it offers stream-to-relation operators, e.g., windowing, relation-to-relation operators, like joins, projections, selections, and aggregation; and relation-to-stream operators to produce output. As in the CEP approach, it offers a comprehensive set of operators to define and detect patterns in event streams.

EPL is similar to SQL, as seen in some of its syntax (e.g., `SELECT`, `FROM`, `WHERE`, `GROUP BY`, etc.). However, it is based on event streams and views, rather than tables. Here, views define the data that is available for queries, they can represent windows over streams, and they can also sort events, derive statistics from event attributes, group events, etc.

3.3 EPL syntax: fire alarms example

To learn more about its syntax, let us consider the following example. Suppose we have a set of fire alarms that send events when they detect smoke. We want to detect when a fire is happening by correlating the events from the fire alarms. Specifically, to detect a fire, we need to receive a smoke event followed by a temperature > 50 event, within 2 minutes, at the same sensor.

In this section, we will learn how to:

- Define event types as schemas
- Create simple queries
- Define different types of windows
- Personalize the output stream
- Recognize complex patterns of events

3.3.1 Creating a schema

In EPL, we can define schemas to define the structure of the events that we are going to process. They are useful to register and recognize specific event types when creating queries, especially the most complex ones.

To define a schema, we mainly use the `CREATE SCHEMA` statement, although we can also use the runtime configuration API `addEventType` method. We use the following syntax:

```

1 create schema --keyword to create a schema
2 schema_name [as...] --name of schema (you can add an alias)
3 (
4     field_name data_type [...], --field name and data type
5     [inherits inherited_schema_name] --it can inherit from another schema
6 );

```

In our example, we need to describe 3 types of events: smoke events, temperature events, and fire events. For that, we can define the following schemas:

```

1 create schema SmokeSensorEvent (
2     sensor string,
3     smoke boolean,
4 );
5
6 create schema TemperatureSensorEvent (
7     sensor string,
8     temperature double,
9 );
10
11 create schema FireEvent (
12     sensor string,
13     smoke boolean,
14     temperature double,
15 );

```

3.3.2 Creating a simple query

In EPL, we can define queries in the same way we do in SQL. We can use the **SELECT** statement to select fields, the **FROM** statement to select from a stream, the **WHERE** statement to filter the data, and so on. Let's see the following example:

```

1 @Name('Query_0')
2 select *
3 from TemperatureSensorEvent
4 where temperature > 50;

```

In this query, we are selecting all fields from the **TemperatureSensorEvent** stream where the temperature is greater than 50. However, in this case, the query is getting executed every time a new event arrives, so it is not efficient. We can consider using an event-based system style, like this:

```

1 @Name('Query_0_BIS')
2 select *
3 from TemperatureSensorEvent(temperature > 50);

```

In this way, the events are filtered before being processed by the query, which is a little more efficient.

We could also want to obtain the results of each sensor separately. For that, we can use the GROUP BY statement, as follows:

```
1 @Name('Query_1')
2 select sensor, avg(temperature)
3 from TemperatureSensorEvent
4 group by sensor;
```

In this query, we are selecting the sensor and the average temperature from the corresponding stream, grouping by sensor. This way, we can obtain the average temperature of each sensor, observed in the whole stream.

3.3.3 Defining windows

In EPL, we can define windows to group events in time or in the number of events. There are 2 main types of windows, regarding the chosen time model: logical windows, that are based on absolute time, and physical windows, that are based on the occurrence of events.

From each type, we can define 3 subtypes: tumbling windows, sliding windows, and hopping windows. We will see how to define each one of them.

Tumbling windows

Tumbling windows are the simplest type of windows. They divide the stream into fixed-size, non-overlapping windows. This size may be measured in absolute time (logical tumbling windows) or in the number of events (physical tumbling windows).

To use a logical tumbling window, we use the following syntax:

```
1 @Name('Query_2')
2 select sensor, avg(temperature)
3 from TemperatureSensorEvent.win:time_batch(4 seconds)
4 group by sensor;
```

In this query, we are selecting the sensor and the average temperature from the corresponding stream, grouping by sensor, and using a logical tumbling window of 4 seconds. This way, we can obtain the average temperature of each sensor, observed in the last 4 seconds, every 4 seconds.

We can also use physical tumbling windows, as follows:

```
1 @Name('Query_3')
2 select sensor, avg(temperature)
3 from TemperatureSensorEvent.win:length_batch(4)
4 group by sensor;
```

In this way, we can obtain the average temperature of each sensor, observed in the last 4 events, every 4 events.

Sliding windows

Sliding windows are a more complex type of windows. They are defined by every event that arrives in the stream, and they can overlap. The slide size is determined either by the time (logical) or by a number of events (physical).

To use a sliding window, we use the following syntax:

```
1 @Name('Query_4')
2 select sensor, avg(temperature)
3 from TemperatureSensorEvent.win:time(4 seconds)
4 group by sensor;
```

In this query, we obtain the average temperature of each sensor, observed in the last 4 seconds, every time an event arrives.

If we want to use a physical sliding window, we can use the following syntax:

```
1 @Name('Query_5')
2 select sensor, avg(temperature)
3 from TemperatureSensorEvent.win:length(4)
4 group by sensor;
```

In this query, we obtain the average temperature of each sensor, observed in the last 4 events, every time an event arrives.

Hopping windows

Hopping windows are a combination of tumbling and sliding windows. They are defined by a fixed-size window that slides by a fixed-size increment. This type of windows tend to overlap. The window size is determined either by the time (logical) or by a number of events (physical).

To use a hopping window, we use the following syntax:

```
1 @Name('Query_6')
2 select sensor, avg(temperature)
3 from TemperatureSensorEvent.win:time(4 seconds)
4 group by sensor;
5 output snapshot every 2 seconds;
```

Notice that to use a hopping window, we need to specify the hopping length by modifying the output stream. In this query, we obtain the average temperature of each sensor, observed in the last 4 seconds, every 2 seconds.

Now, let's deepen our knowledge about how to personalize the output stream.

3.3.4 Output stream

In EPL, we can declare output policies to define how the output stream is going to be generated. We use the `OUTPUT` statement, followed by the desired policy. Also, with the `EVERY` keyword, we can specify the frequency of the output.

There are 4 main types of output policies: snapshot, all, first, and last. Let's see how to use each one of them:

- Snapshot: it generates a snapshot of the current state of the window.
- All: it generates all the events in the window that matches the query.
- First: it generates the first event in the window that matches the query.
- Last: it generates the last event in the window that matches the query.

Let's see an example of this:

```
1 @Name('Query_7_SNAP')
2 select sensor, avg(temperature)
3 from TemperatureSensorEvent.win:time(4 seconds)
4 group by sensor;
5 output snapshot every 2 seconds;
6
7 @Name('Query_7_ALL')
8 select sensor, avg(temperature)
9 from TemperatureSensorEvent.win:time(4 seconds)
10 group by sensor;
11 output all every 2 seconds;
12
13 @Name('Query_7_FIRST')
14 select sensor, avg(temperature)
15 from TemperatureSensorEvent.win:time(4 seconds)
16 group by sensor;
17 output first every 2 seconds;
18
19 @Name('Query_7_LAST')
20 select sensor, avg(temperature)
21 from TemperatureSensorEvent.win:time(4 seconds)
22 group by sensor;
23 output last every 2 seconds;
```

In this example, we are obtaining the average temperature of each sensor, observed in the last 4 seconds, every 2 seconds, using the snapshot, all, first, and last output policies.

The main difference between `SNAPSHOT` and `ALL` is that the former generates a snapshot of the current state of the window, while the latter generates all the events in the window that matches the query, no matter if they are a null event.

3.3.5 Complex event patterns

In EPL, we can define complex event patterns to detect specific sequences of events in the stream. We use the **PATTERN** statement, followed by the desired pattern. Let's see an example of this:

```

1 @Name('Query_8')
2 select *
3 from pattern [
4     a=SmokeSensorEvent(smoke=true) ->
5     b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
6 ];

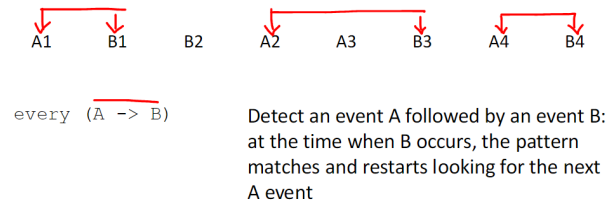
```

This query is searching for a smoke event followed by a temperature > 50 event. Note that this query only matches the first occurrence of the pattern. If we want to match all the occurrences, we should use the **EVERY** keyword.

EVERY clause

The **EVERY** clause is used to re-start the pattern matching process after a successful or failed match. It has different approaches:

- **EVERY (A -> B)**: it matches the pattern every time the event A is followed by the event B. The moment the event B is matched, the pattern matching process is re-started, and the event A is expected again.



B1	{A1, B1}
B3	{A2, B3}
B4	{A4, B4}

Figure 3.2: Pattern matching with **EVERY (A -> B)**

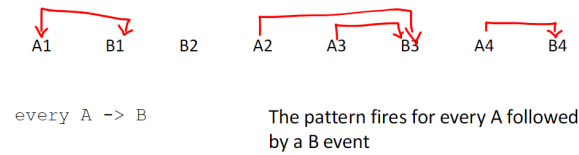
Let's see an example of this:

```

1 @Name('Query_9')
2 select *
3 from pattern [
4     every (
5         a=SmokeSensorEvent(smoke=true)
6         ->
7         b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
8     )
9 ];

```

- **EVERY A -> B**: it matches the pattern every time the event A is followed by the event B. In this case, the pattern matching process fires every time an event A is matched, and successfully returns when the event B is matched, so it will return every A followed by B.



B1	{A1, B1}
B3	{A2, B3}, {A3, B3}
B4	{A4, B4}

Figure 3.3: Pattern matching with **EVERY A -> B**

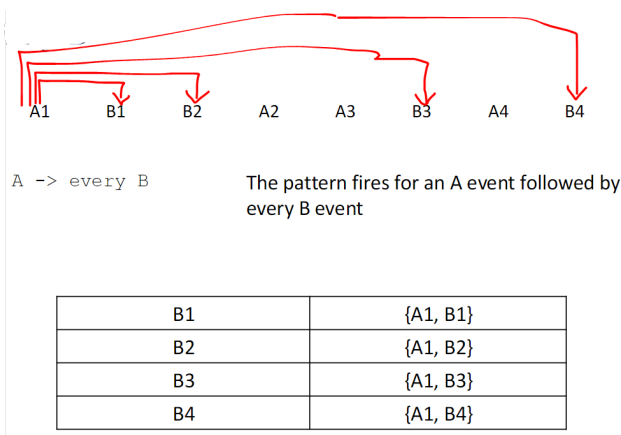
Let's see an example of this:

```

1 @Name('Query_10')
2 select *
3 from pattern [
4     every a=SmokeSensorEvent(smoke=true)
5     ->
6     b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
7 ];

```

- **A -> EVERY B**: it matches the pattern every time the event A is followed by the event B. In this case, once the event A is matched, this pattern will be matched every time the event B is encountered, so it will return every B following that A.



B1	{A1, B1}
B2	{A1, B2}
B3	{A1, B3}
B4	{A1, B4}

Figure 3.4: Pattern matching with **A -> EVERY B**

Let's see an example of this:

```

1 @Name('Query_11')
2 select *
3 from pattern [
4     a=SmokeSensorEvent(smoke=true)
5     ->
6     every b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
7 ];

```

- **EVERY A -> EVERY B**: it matches the pattern every time the event A is followed by the event B. In this case, the pattern matching process fires every time an event A is matched, and for each A, it will return every B that follows that A.

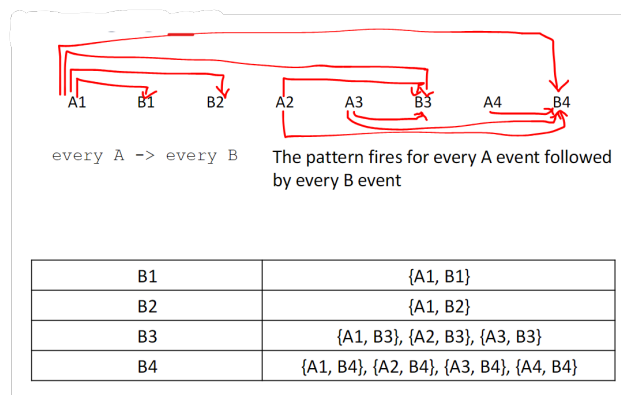


Figure 3.5: Pattern matching with **EVERY A -> EVERY B**

Let's see an example of this:

```

1 @Name('Query_12')
2 select *
3 from pattern [
4     every a=SmokeSensorEvent(smoke=true)
5     ->
6     every b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
7 ];

```

Note that the **EVERY** clause restarts the pattern matching process multiple times. This can be very resource-intensive, so it should be used with caution. Also, we don't always want to match every pattern in the stream, but only in an specific window, or in an specific context. For that, we have pattern guards.

Pattern guards

Pattern guards are used to restrict the pattern matching process to a specific context. There are 2 main types of pattern guards: **TIMER:WITHIN** and **AND NOT**. Let's see how to use each one of them:

- **TIMER:WITHIN:** it restricts the pattern matching process to a specific time window. It is used to specify a time limit for the pattern matching process. Let's see an example of this:

```

1 @Name('Query_13')
2 select *
3 from pattern [
4     a=SmokeSensorEvent(smoke=true)
5     ->
6     b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
7     where timer:within(2 minutes)
8 ];

```

In this query, we are searching for a smoke event followed by a temperature > 50 event, within 2 minutes. This way, we are restricting the pattern matching process to a specific time window.

- **AND NOT:** it restricts the pattern matching process to exclude specific events. It is used to specify that an event should not be matched. Let's see an example of this:

```

1 @Name('Query_14')
2 select *
3 from pattern [
4     a=SmokeSensorEvent(smoke=true)
5     ->
6     b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
7     and not c=FireEvent(sensor = a.sensor)
8 ];

```

In this query, we are searching for a smoke event followed by a temperature > 50 event, and excluding the case when a fire event is detected at the same sensor. This way, we are restricting the pattern matching process to exclude specific events.

Operators precedence

For all these operations, we need to consider the operators precedence. The operators are evaluated in the following order:

Order	Operator	Description	Example
1	guard postfix	where timer:within and while (expression) (incl. withinmax and plug-in pattern guard)	MyEvent where timer:within(1 sec) a=MyEvent while (a.price between 1 and 10)
2	unary	every, not, every-distinct	every MyEvent timer:interval(5 min) and not MyEvent
3	repeat	[num], until	[5] MyEvent [1..3] MyEvent until MyOtherEvent
4	and	and	every (MyEvent and MyOtherEvent)
5	or	or	every (MyEvent or MyOtherEvent)
6	followed-by	->	every (MyEvent -> MyOtherEvent)

Figure 3.6: Operators precedence

3.3.6 Composing queries

Typically, queries in EPL are placed in a network. In this way, we can make queries interact with each other, and we can define the order in which they are executed. This allows us to define complex queries by composing simpler ones.

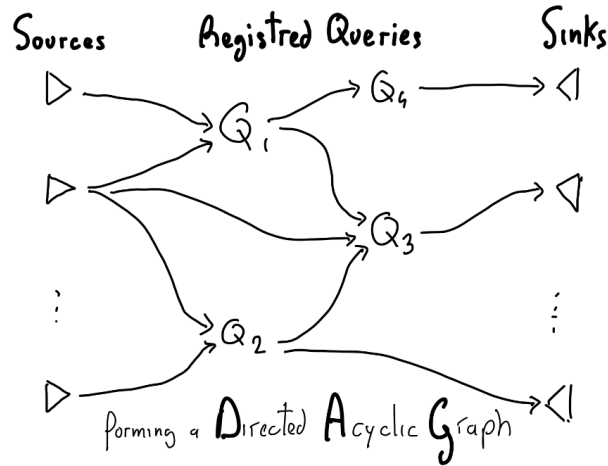


Figure 3.7: Query network

We use the `INSERT INTO` statement to insert the results of a query into a stream. In this way, we can use the results of a query as input for another query. Let's see an example of this:

```

1 @Name('Query_15')
2 insert into FireEvent
3 select a.sensor, a.smoke, b.temperature
4 from pattern [
5     every (
6         a=SmokeSensorEvent(smoke=true)
7         ->
8         b=TemperatureSensorEvent(temperature > 50, sensor = a.sensor)
9         where timer:within(2 minutes)
10    )
11 ];

```

In this query, we are searching for a smoke event followed by a temperature > 50 event, within 2 minutes. If the pattern is matched, the results are inserted into the `FireEvent` stream. Then, we can use this query:

```

1 @Name('Query_16')
2 select count(*)
3 from FireEvent.win:time(10 minutes);

```

In this query, we are counting the number of fire events detected in the last 10 minutes. This way, we can use the results of the first query as input for the second query.

3.4 Joins in EPL

In SQL, joins are performed on tables, which are static data structures. The semantics of SQL joins are based on relational data model and queries that work with complete data sets. These joins do not consider the time dimension of the data.

In EPL, joins are performed on event streams, which are dynamic data structures. The semantics of EPL joins are based on event streams and queries that work with continuous data. These joins intrinsically consider the time dimension of the data.

Generally, joins are classified in 3 main types:

- Inner join: it returns only the rows that have matching values in both tables.
- Left/right join: it returns all the rows from the left/right table, and the matched rows from the other table.
- Full join: it returns all the rows when there is a match in one of the tables.

Note that this is true for SQL and EPL. Also based on the structure on which the join is performed, in EPL we can distinguish between 3 types of joins:

- Stream to stream joins
- Table to table joins
- Stream to table joins (and viceversa)

Note that complex event pattern matching (CEP), seen in the previous section, is a form of join, where the join condition is based on the temporal order of the events.

3.4.1 Stream to stream joins

Stream to stream joins are used to join 2 event streams. In each stream, we specify a window to define the time frame in which the join is performed. For example, if we use a window of 5 seconds, we will only consider the events that arrive in the last 5 seconds. Its syntax is as follows:

```
1 select *
2 from View#time(9 sec) as v
3     inner join
4     Click#time(9 sec) as c
5     on v.id = c.id;
```

In this query, we are joining the **View** and **Click** streams on the **id** field. We are using a window of 9 seconds for each stream. This way, we are only considering the events that arrive in the last 9 seconds.

Consider that the queries are only executed when an event arrives in the stream, and only returns the new results that match the join condition. The following figure describes the process of a stream to stream (inner) join:

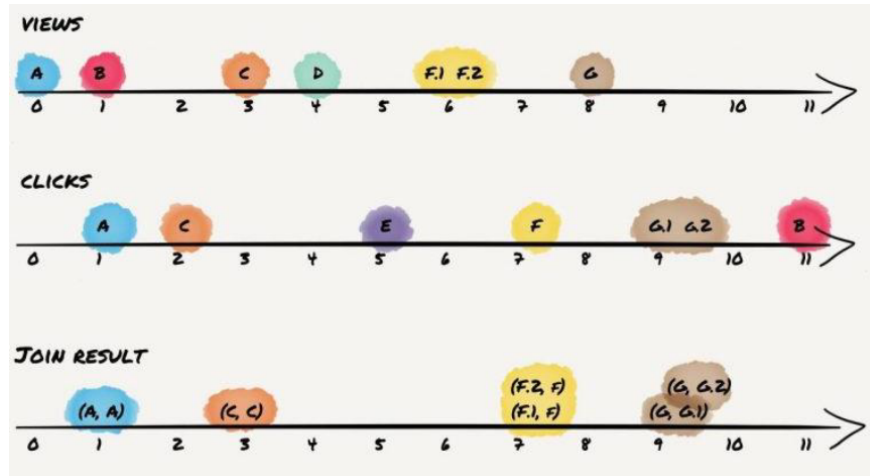


Figure 3.8: Stream to stream join

Note that we can also perform left/right outer joins and full outer joins. The syntax is similar to the inner join, but we use the `left outer join`, `right outer join`, and `full outer join` statements, respectively.

Joins among streams have profoundly different semantics, because they are designed to run on real-time data streams, essentially unbounded record sequences. Events are joined only if they arrive within a specific time window. This temporal aspect is central and intrinsic to the join semantics of streams. The concept of time is crucial, as Joins depend on the time synchronization of the events in the streams, which introduces a significant difference with respect to traditional SQL joins.

3.4.2 Table to table joins

EPL offers 3 ways to build a table from an event stream:

- **KEEP ALL:** this window is one that retains all arriving events. However, take should be taken to remove events from the window, in a timely manner.
- **UNIQUE:** this window is one that retains only the most recent event for each unique key.
- **CREATE TABLE with INSERT INTO:** this window is one that retains all arriving events, but it is not a window in the sense of a time window.

We will mainly use the **UNIQUE** window to build tables from event streams. This window is one that retains only the most recent event for each unique key. Its syntax is as follows:

```

1 select *
2 from View#unique(id) as v
3     inner join
4     Click#unique(id) as c
5     on v.id = c.id;
```

In this query, we are joining the `View` and `Click` streams on the `id` field. We are using the `unique` window to build tables from the event streams. This way, we are only considering the most recent event for each unique key, in this case, the `id` field.

Note that we are not actually joining tables, but event streams. The tables are built from the event streams, and the join is performed on these tables. So we call them "tables" as they seem to inherit the unique key constraint from the SQL tables. The following figure describes the process of a table to table (inner) join:

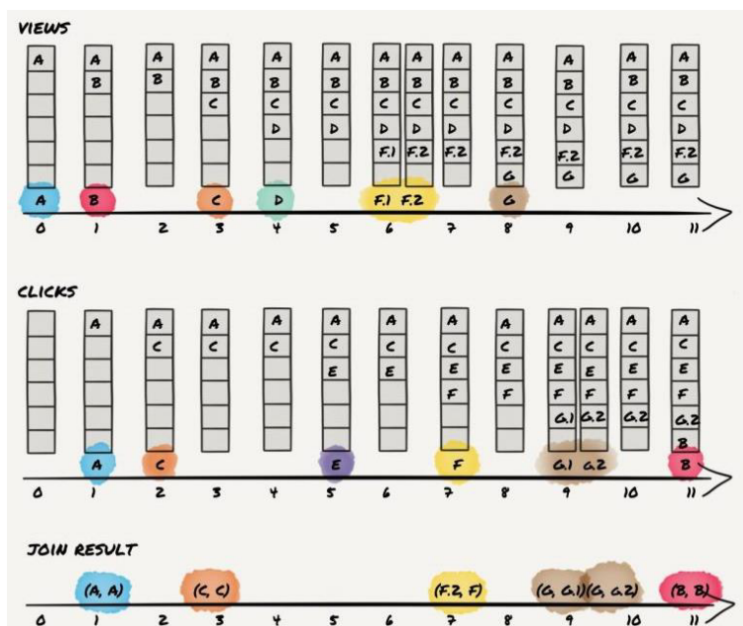


Figure 3.9: Table to table join

Note that the events persist over time, only being replaced by newer events with the same key. This is a key difference with respect to stream to stream joins, where the events are only considered within a specific time window.

We can also perform left/right outer joins and full outer joins. The syntax is similar to the inner join, but we use the `left outer join`, `right outer join`, and `full outer join` statements, respectively.

3.4.3 Stream to table joins

These type of joins introduce a new keyword: `UNIDIRECTIONAL`. This keyword is used in the `JOIN` clause to identify streams that provide the events to execute the join. If the keyword is present for a stream, all other streams in the `FROM` clause become passive streams. That means that when a new event arrives or leaves a data window of a passive stream, the join does not generate new results.

Therefore, the `UNIDIRECTIONAL` keyword makes the stream-to-table join asymmetric: only the input from the stream with the `UNIDIRECTIONAL` keyword will trigger the join. And because the

join is not-windowed, this stream is stateless; thus join lookups from the table to the stream are not possible.

The syntax for a stream to table join is as follows:

```
1 select *
2 from View as v
3     unidirectional inner join
4     Click#unique(id) as c
5     on v.id = c.id;
```

In this query, we are joining the **View** stream with the **Click** table on the **id** field. We are using the **unique** window to build the table from the **Click** stream. This way, we are only considering the most recent event for each unique key, in this case, the **id** field. In this case, a new click event won't trigger the join, but a new view event will, and because the stream is stateless, join lookups from the **Click** table to the **View** stream are not possible.

The following figure describes the process of a stream to table (inner) join:

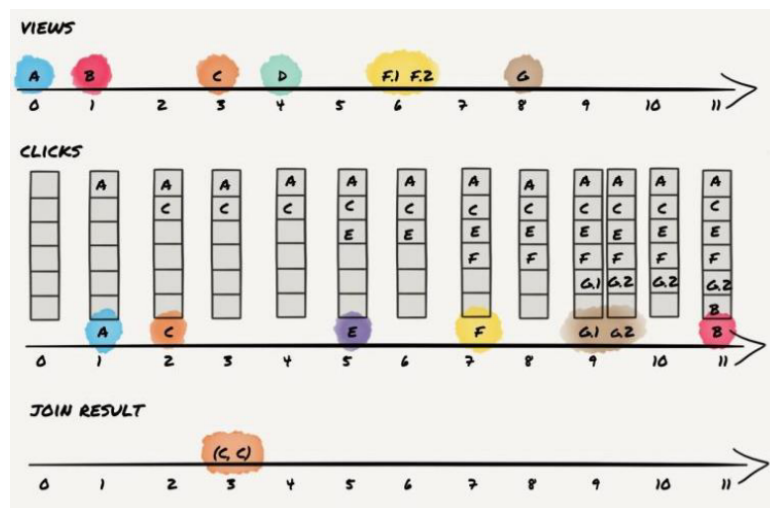


Figure 3.10: Stream to table join

We can also perform left/right outer joins and full outer joins. The syntax is similar to the inner join, but we use the `left outer join`, `right outer join`, and `full outer join` statements, respectively.

3.5 Contexts in EPL

In EPL, we can also define contexts. A context enables more flexible and stateful processing in scenarios requiring grouping or conditions that go beyond simple event windows. In fact, contexts define a new type of window (different than tumbling and hopping windows), which we will name a session window.

Why use contexts? Here is a list of advantages:

- Better control: context gives fine-grained control over how events are processed.
- Stateful processing: handles cases where the state is important, e.g., sessions or transactions.
- Efficient partitioning: it helps reduce complexity in event, partitioning by user, session, or custom logic.

We have three main types of contexts:

- Context by key: partitions events based on a key, e.g., `userId`.
- Context by start/end conditions: defined by specific events or triggers.
- Context by time: defines temporal windows (generalizes logical windows, e.g., it allows to declare session windows).

3.5.1 Context by key

Context by key is used to partition events based on a key. It is useful when we want to group events by a specific attribute, e.g., `userId`. Its syntax is as follows:

```
1 create context UserSessionContext
2 partition by userId from UserActionEvent;
```

In this context, we are partitioning events by the `userId` attribute from the `UserActionEvent` stream. This way, we are grouping events by the `userId` attribute. Note that this process is highly parallelizable, as each partition can be processed independently.

3.5.2 Context by start/end conditions

Context by start/end conditions is defined by specific events or triggers. It is useful when we want to group events based on specific conditions, e.g., when a session starts or ends. Its syntax is as follows:

```
1 create context OrderContext
2 initiated by OrderEvent(orderStatus = 'NEW')
3 terminated by OrderEvent(orderStatus = 'COMPLETED');
```

In this context, we are defining a session window that starts when an `OrderEvent` with `orderStatus = 'NEW'` arrives and ends when an `OrderEvent` with `orderStatus = 'COMPLETED'` arrives. This way, we are grouping events based on the start and end conditions.

Note that we could use this context along with the context by key to partition events based on a key and start/end conditions. Look at the following example:

```
1 create context OrderContext;
2 partition by userId from UserActionEvent;
3 initiated by UserActionEvent(action = 'start_order');
4 terminated by UserActionEvent(action = 'complete_order');
```


3.5.3 Context by time

Context by time allows us to track actions within fixed periods of time. They are triggered by some condition and last for a specific time. Its syntax is as follows:

```
1 create context TimeBatchContext
2 partition by userId from UserActionEvent;
3 initiated by UserActionEvent(action = 'start_order');
4 terminated after 10 minutes;
```

In this context, we are defining a session window that starts when an `UserActionEvent` with `action = 'start_order'` arrives and ends after 10 minutes. This way, we are grouping events based on the time window. Note that we are also partitioning events by the `userId` attribute.

3.5.4 Comparison: Windows vs. Contexts

Windows and contexts are both used to group events, but they have different semantics. Let us look at the following figure that shows the main differences between windows and contexts:

Windows

- operate over a **fixed** time frame or number of events
- are stateless
- reset after each use

Contexts

- **provide more flexible segmentation,**
- have states that persist across event streams
- persist until terminated, providing a more robust grouping for complex scenarios

Figure 3.11: Windows vs. Contexts

Chapter 4

Kafka and Spark

4.1 Introduction: scaling stream ingestion and processing

To process a high volume of data, we need systems that are horizontally scalable. This means that we can add more machines to the system to increase its capacity. In this chapter, we will discuss two such systems: Apache Kafka and Apache Spark, and how they can be used together to process large volumes of data.

In general, the roles of Kafka and Spark are as follows:

- Kafka is used to collect and store data in real-time. It manages the ingestion of data from multiple sources and stores it in a distributed manner, making it available for processing by other systems.
- Spark is used to process the data stored in Kafka, and generate insights. It is a distributed computing system that can process large volumes of data in parallel.

Before getting into the details of Kafka and Spark, let's discuss some important concepts that are used in these systems:

- Latency vs throughput
- Data/message

4.1.1 Latency vs throughput

We refer to latency as the time taken for a message to travel from the producer to the consumer. More generally, latency is the amount of time needed to complete a task. Throughput, on the other hand, is the number of messages that can be processed in a given time period.

We can see this with a visual example: consider a highway with cars moving at different speeds. The latency is the time taken for a car to travel from one end of the highway to the other, while the throughput is the number of cars that can pass through the highway in a certain time period, say an hour.

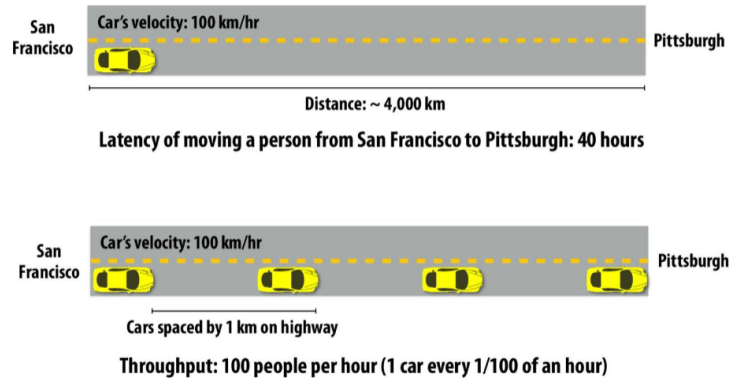


Figure 4.1: Latency vs throughput

We could improve the throughput by adding more lanes to the highway, but this would not necessarily reduce the latency. To improve the latency, we would need to reduce the distance between the two ends of the highway, or increase the speed of the cars. Notice that the latter would also improve the throughput, as more cars would be able to pass through the highway in a given time period.

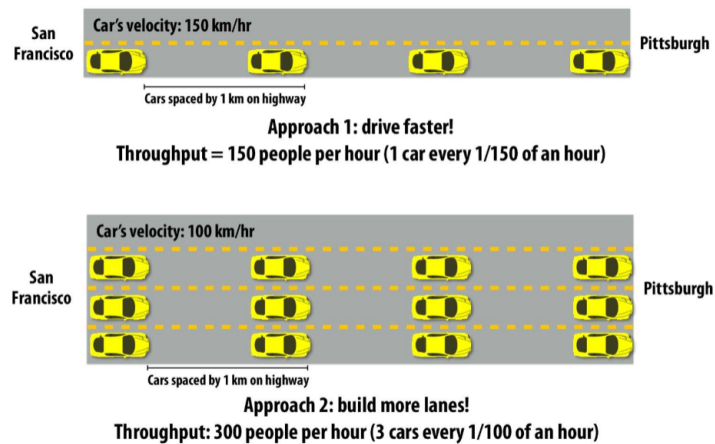


Figure 4.2: Improving latency and throughput

In the context of data processing, we usually want to minimize the latency and maximize the throughput. This is because we want to process data as quickly as possible, and we want to be able to process a large volume of data in a given time period.

4.1.2 Data/message

In the context of data processing, the data that comes in is usually in the form of messages. A message is a unit of data that is sent from a producer to a consumer. For example, in a messaging system like Kafka, a message could be a log entry, a sensor reading, or a user event. Messages can be of different types and sizes, and can contain different types of data. For example, a message

could be a JSON object, a binary blob, or a string.

Note that the way messages are represented can have a big impact on the performance of the system. Following the previous example, imagine that the messages are cars on a highway, and the people in the cars are the data. If the cars are small and can carry only one person, then we would need a lot of cars to transport a large number of people. On the other hand, if the cars are large and can carry many people, then we would need fewer cars to transport the same number of people. The same is true for messages: if the messages are small, then we would need more of them to transport the same amount of data, which would increase the overhead of the system. On the other hand, if the messages are large, then we would need fewer of them to transport the same amount of data, which would reduce the overhead of the system.

In general, the larger the data points per message, the higher the throughput, but the higher the latency. This is because larger messages take longer to process, but fewer of them are needed to transport the same amount of data. So, it is needed to find a balance between the size of the messages and the performance of the system.

Compression and binary encoding are the typical ways to increase the data per message, given a fixed message size.

4.2 Apache Kafka

Apache Kafka is a distributed streaming platform that is used for building real-time data pipelines and streaming applications. It is designed to be scalable, fault-tolerant, and highly available. Kafka is used for collecting, storing, and processing large volumes of data in real-time. It is used by many companies, including LinkedIn, Netflix, Uber, and Airbnb, to process data at scale.

4.2.1 Kafka: high-level vs system view

Let us take a conceptual view of Kafka. In general, Kafka is composed of the following components:

- **Producer:** A producer is a system that sends data to Kafka. It sends the messages based on a certain topic.
- **Consumer:** A consumer is a system that reads data from Kafka. It consumes the messages from certain topics that subscribes to.
- **Topics:** A topic is a category of messages in Kafka. It is a way to organize the messages in Kafka. Producers send messages to topics, and consumers read messages from topics.
- **Messages:** A message is a unit of data in Kafka. It is a key-value pair that is sent from a producer to a consumer. Messages are stored in topics and can be read by consumers.
- **Cluster:** A Kafka cluster manages a set of topics.

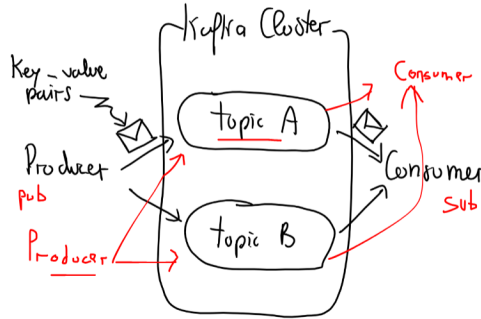


Figure 4.3: Kafka high-level view

This is a high-level view of Kafka. In practice, Kafka uses a distributed architecture to achieve scalability, fault-tolerance, and high availability. So, for that, we need to introduce the following concepts:

- **Brokers:** A broker is a Kafka server that stores the messages. A Kafka cluster is composed of multiple brokers. Each broker stores a subset of the messages in the cluster.
- **Partitions:** A partition is a unit of parallelism in Kafka. A topic is divided into multiple partitions, and each partition is stored on a different broker. This allows Kafka to scale out by distributing the load across multiple brokers.

In general, a broker receives messages from producers, stores them in partitions, and serves them to consumers. A production Kafka cluster typically consists of multiple brokers, each storing multiple partitions of multiple topics. Each broker can handle a certain amount of traffic, and the cluster as a whole can handle a large volume of data.

To ensure fault-tolerance, Kafka uses replication. Each partition is replicated across multiple brokers, so that if one broker fails, the data is still available on other brokers. This allows Kafka to provide high availability and durability.

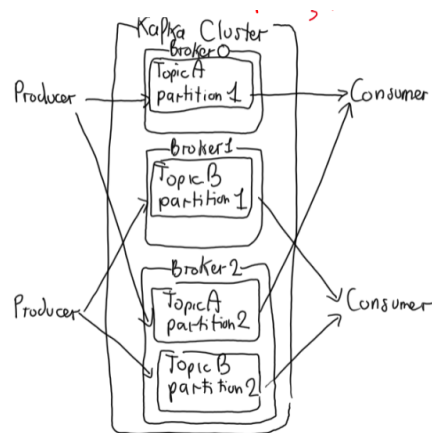


Figure 4.4: Kafka architecture

4.2.2 Kafka: how it works

Let us describe in details how Kafka works. In general, Kafka works as follows:

1. A producer sends a message to a topic. The producer specifies the topic and the message to be sent.
2. The message is sent to a broker. The broker receives the message and stores it in a partition. The partition is chosen based on the key of the message. If no key is specified, then round-robin is used to choose the partition.
3. The message is replicated to other brokers, to ensure fault-tolerance. The number of replicas is configurable, and the replication factor determines how many replicas are created.
4. Each partition is stored on the disk of the broker. The messages are stored in a log-structured format, which allows for fast reads and writes. Each message is assigned an offset, which is a unique identifier for the message. This log is called the **commit log**.

Commit log

Offset	0	1	2	3	4	5	6	7	8
key	K ₁	K ₂	K ₁	K ₃	K ₄	K ₅	K ₅	K ₂	K ₆
Value	V ₂	V ₃	V ₄	V ₅	V ₆	V ₇	V ₈	V ₉	V ₁₀

Figure 4.5: Commit log

5. A consumer reads messages from a topic. The consumer specifies the topic and the partition to read from. The consumer reads messages in order of their offset, starting from the last offset read.

Note that each broker can configure its log retention policy, which determines how long messages are kept in the commit log. This allows Kafka to manage the storage of messages and prevent the log from growing indefinitely. In general, the default retention policy is to keep messages for a certain amount of time, after which they are deleted, but we can also configure the retention policy to be based on the size of the log, or the number of messages in the log.

Also, when deleting a log, there is a concept of **compaction**, which is a process that removes duplicate messages from the log. This is useful when we want to keep only the latest version of a message, and discard older versions. Compaction is useful for maintaining the state of a system, and for ensuring that the log does not grow too large.

When making replicas, Kafka uses a leader-follower model. Each partition has one leader and multiple followers (replicas). The leader is responsible for handling reads and writes, while the followers are responsible for replicating the data. If the leader fails, one of the followers is elected as the new leader, and the system continues to operate.

4.2.3 Distributed consumption

Kafka allows for distributed consumption, which means that multiple consumers can read messages from a topic in parallel. This allows Kafka to scale out by distributing the load across multiple consumers.

Consumers can be part of a consumer group, which is a group of consumers that work together to read messages from a topic. Each consumer in the group reads messages from a different partition, so that the load is distributed evenly across the consumers. If a consumer fails, then the partitions that it was reading from are reassigned to other consumers in the group, so that the system continues to operate.

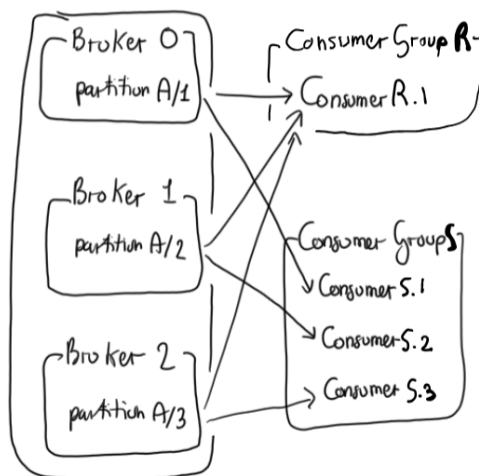


Figure 4.6: Consumer group

In general, Kafka provides at-least-once delivery guarantees, which means that messages are delivered to consumers at least once. This is because Kafka stores messages in the commit log, and consumers can re-read messages if needed. However, Kafka does not provide exactly-once delivery guarantees, which means that messages can be delivered more than once.

4.3 Apache Spark

Apache Spark is an unified analytics engine for big data processing, with built-in modules for streaming, SQL, machine learning, and graph processing. It is designed to be fast, easy to use, and general-purpose. Spark is used by many companies, including Netflix, Uber, and Airbnb, to process data at scale.

Spark is built around the concept of Resilient Distributed Datasets (RDDs), which are fault-tolerant collections of objects that can be operated on in parallel. RDDs are the fundamental data structure in Spark, and are used to represent data that is distributed across a cluster of machines.

4.3.1 Key concepts in Spark

Let us discuss some key concepts in Spark:

- **Use RAM instead of disk:** Spark is designed to use RAM instead of disk for storing intermediate results. This allows Spark to be much faster than traditional MapReduce systems, which use disk for storing intermediate results.

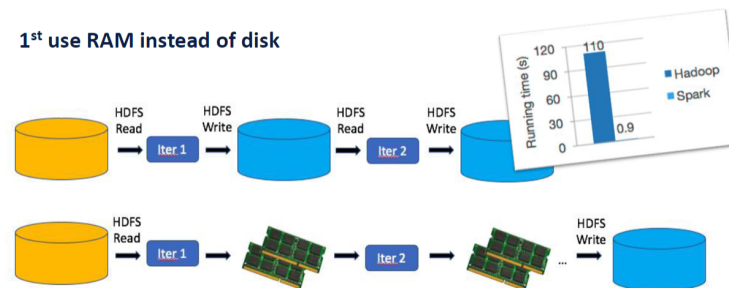


Figure 4.7: Spark in-memory processing

- **Ease of use:** Spark provides a high-level API that makes it easy to write distributed programs. The API is available in multiple languages, including Scala, Java, Python, and R.
- **Generality:** Spark is a general-purpose computing engine that can be used for a wide range of applications, including batch processing, real-time processing, machine learning, and graph processing. This makes Spark a versatile tool that can be used for many different use cases.

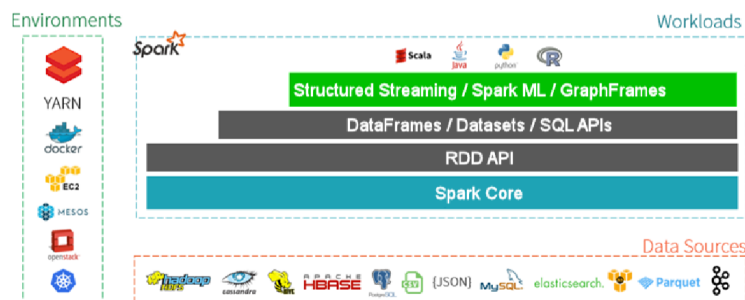


Figure 4.8: Spark generality

4.3.2 Resilient Distributed Datasets (RDDs)

RDDs are the fundamental data structure in Spark. They are fault-tolerant collections of objects that can be operated on in parallel. RDDs are immutable, which means that they cannot be changed once they are created. Instead, transformations are applied to RDDs to create new RDDs.

RDDs can be created from a variety of data sources, including HDFS, S3, Cassandra, and Kafka. Once created, RDDs can be transformed using a variety of operations, including map, filter, reduce, and join. These operations are applied in parallel across the cluster, so that the processing is distributed across multiple machines.

Transformations and actions

In Spark, there are two types of operations that can be applied to RDDs: transformations and actions. Transformations are operations that create a new RDD from an existing RDD, while actions are operations that return a value to the driver program.

Some common transformations in Spark include:

- **map**: Applies a function to each element in the RDD.
- **filter**: Filters the elements in the RDD based on a predicate.
- **reduce**: Aggregates the elements in the RDD using a commutative and associative function.
- **join**: Joins two RDDs based on a key.

Some common actions in Spark include:

- **collect**: Returns all the elements in the RDD to the driver program.
- **count**: Returns the number of elements in the RDD.
- **take**: Returns the first n elements in the RDD.
- **saveAsTextFile**: Saves the RDD to a text file.

4.3.3 Spark at work

Let us see how Spark can be used to process data. A Spark cluster consists of a driver program and multiple worker nodes. The driver program is responsible for coordinating the execution of the job, while the worker nodes are responsible for executing the tasks.

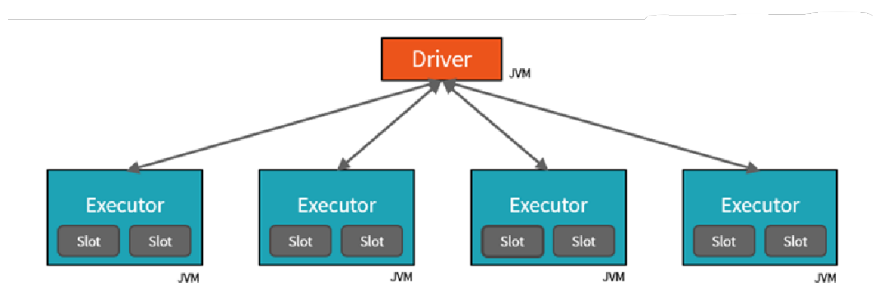


Figure 4.9: Spark components

When a job is submitted to the Spark cluster, the driver program breaks it down into stages, which are then executed in parallel across the worker nodes. Each stage consists of tasks, which are executed in parallel across the worker nodes. The tasks read data from RDDs, apply transformations to the data, and write the results to new RDDs.

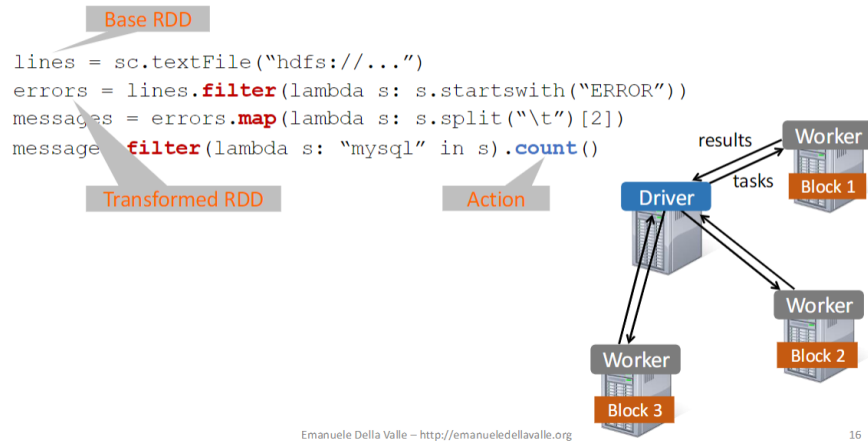


Figure 4.10: Spark execution model

In general, Spark uses a lazy evaluation model, which means that transformations are not executed immediately. Instead, they are queued up and executed when an action is called. This allows Spark to optimize the execution of the job, by combining multiple transformations into a single stage, and reducing the number of shuffles that are needed.

Spark also provides fault-tolerance through lineage, which is a record of the transformations that were applied to an RDD. If a partition of an RDD is lost, Spark can recompute it by replaying the transformations that were applied to the RDD. This allows Spark to recover from failures and continue processing the data.

Spark also provides caching, which allows RDDs to be stored in memory for faster access. This is useful when an RDD is used multiple times in a job, as it avoids the need to recompute the RDD each time it is used. Caching can be done at different levels, including memory only, disk only, and memory and disk.

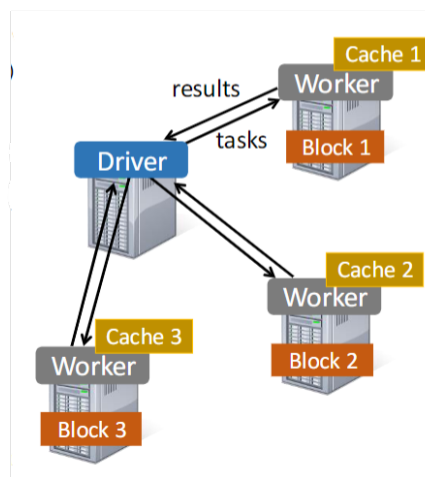


Figure 4.11: Spark caching

4.3.4 DataFrames in Spark

DataFrames are a high-level API for working with structured data in Spark. They provide a more user-friendly interface than RDDs, and are optimized for performance. DataFrames are built on top of RDDs, and provide a more structured way to work with data.

DataFrames can be created from a variety of data sources, including JSON, CSV, and Parquet files. Once created, DataFrames can be queried using SQL, and transformed using a variety of operations, including select, filter, groupBy, and join. These operations are optimized for performance, and are executed in parallel across the cluster.

Let us see an example of how to create a DataFrame in Spark:

```
1 from pyspark.sql import SparkSession
2
3 # Create a Spark session
4 spark = SparkSession.builder.appName("example").getOrCreate()
5
6 # Create a DataFrame from a JSON file
7 df = spark.read.json("data.json")
8
9 # Show the first 5 rows of the DataFrame
10 df.show(5)
11
12 # Filter the DataFrame
13 filtered_df = df.filter(col("age") > 30)
14
15 # Create a temporary view to query the DataFrame using SQL
16 filtered_df.createOrReplaceTempView("filtered_df")
17
18 # Query the DataFrame using SQL
19 query = "SELECT * FROM filtered_df"
20 result = spark.sql(query)
```

Catalyst and Tungsten

DataFrames in Spark are built on top of two key components: Catalyst and Tungsten. Catalyst is a query optimizer that optimizes the execution of DataFrame operations. It uses a rule-based optimizer to transform the query plan and generate an optimized plan. Tungsten is a memory-centric execution engine that optimizes the execution of DataFrame operations. It uses code generation and memory management techniques to improve the performance of the operations.

Together, Catalyst and Tungsten provide a high-performance engine for working with structured data in Spark. They optimize the execution of DataFrame operations, and provide a more user-friendly interface for working with data.

4.4 Spark Structured Streaming

Spark Structured Streaming provides a fast, scalable, fault-tolerant, end-to-end exactly-once stream processing engine, without the user having to reason about streaming. It is built on the Spark SQL

engine, and provides a high-level API for working with streaming data. Structured Streaming allows users to write streaming queries in the same way as batch queries, and provides the same fault-tolerance and exactly-once guarantees as batch processing.

4.4.1 Programming model

Spark Structured Streaming treats a live data stream as an unbounded table, which is continuously updated as new data arrives. As a result of this idea, users can logically express streaming computation as a Directed Acyclic Graph (DAG) of standard batch-like queries on static tables. Spark will physically run this DAG incrementally and continuously on the streaming data.

It uses the model of a DSMS (Data Stream Management System), as we saw in previous chapters. In every micro-batch, Spark will execute the DAG incrementally to update the final result. For each query, it will maintain a state, which is updated in each micro-batch. We can see this in the following figure:

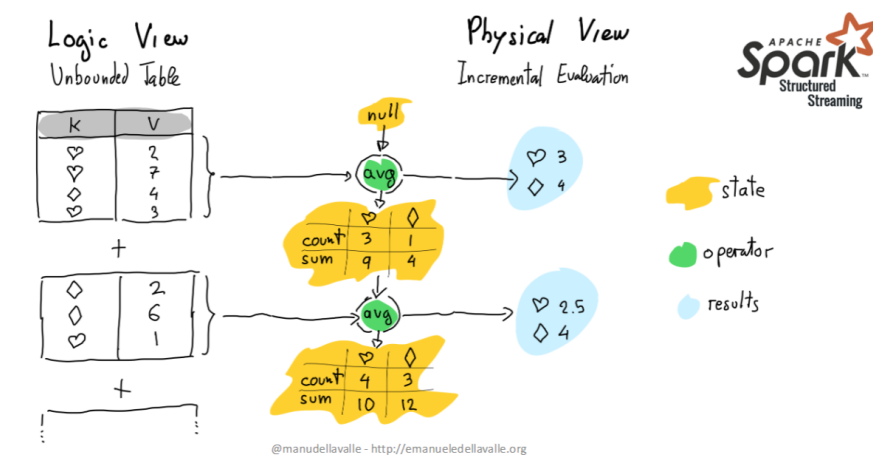


Figure 4.12: Spark incremental evaluation

It is important to note that Spark Structured Streaming does not materialize the entire table, but only the incremental changes to the table. It will read the latest available data from the stream, process it incrementally to update the result, and then discard the data. It will only keep around the minimal intermediate state data as required to update the result, for example, the intermediate counts and average for the above example.

This model is very powerful, as it allows users to express complex streaming computations in a simple and declarative way, and provides the same fault-tolerance and exactly-once guarantees as batch processing. It also allows users to easily integrate streaming and batch processing, as the same code can be used for both.

4.4.2 Creating streaming DataFrames

To create a streaming DataFrame, we can use the `readStream` method of the `SparkSession` object. It can receive as input the following sources:

- File sources, i.e., a folder in which new files get dropped.
- Kafka sources, i.e., a Kafka topic.
- Socket sources, i.e., a TCP socket. This is only for testing purposes, as it does not provide fault-tolerance.

Let us see an example of how to create a streaming DataFrame in Spark, reading from a Kafka topic:

```

1 from pyspark.sql import SparkSession
2
3 # Create a Spark session
4 spark = SparkSession.builder.appName("example").getOrCreate()
5
6 # Create a streaming DataFrame from a Kafka topic
7 raw_df = spark.readStream.format("kafka") \
8     .option("kafka.bootstrap.servers", "localhost:9092") \
9     .option("startingOffsets", "earliest") \
10    .option("subscribe", "topic") \
11    .load()
12
13 # Show the first 5 rows of the streaming DataFrame
14 raw_df.show(5)

```

In this example, we create a streaming DataFrame from a Kafka topic. We specify the Kafka bootstrap servers, the starting offsets, and the topic to subscribe to. We then show the first 5 rows of the streaming DataFrame.

Note that when extracting data from Kafka, our `raw_df` DataFrame will have the following schema:

- **key**: The key of the message, as a binary blob.
- **value**: The value of the message, as a binary blob.
- **topic**: The topic of the message, as a string.
- **partition**: The partition of the message, as an integer.
- **offset**: The offset of the message, as an integer.
- **timestamp**: The timestamp of the message, as a timestamp.
- **timestampType**: The timestamp type of the message, as an integer.

Here, the column that contains the actual data is the **value** column. We can extract the data from this column using the **select** method, as we will see in the following example:

```

1 # Extract the data from the value column
2 casting = ("CAST(key AS STRING)", "CAST(value AS STRING)")
3 data_df = raw_df.selectExpr(*casting)

```

In this example, we extract the data from the `key` and `value` columns, and cast them to strings. This allows us to work with the data as strings, which is more convenient for processing. But the data of value will be a stringified JSON object, so we need to parse it to a DataFrame. We can do this as follows:

```
1 from pyspark.sql.functions import from_json
2
3 # Parse the data from the value column
4 schema = "id STRING, name STRING, age INT"
5 parsed_df = data_df.select(from_json("value", schema).alias("data"))
6
7 # Obtain the final DataFrame
8 final_df = parsed_df.select("data.*")
```

In this example, we parse the data from the `value` column using the `from_json` function, and specify the schema of the data as a string. We then select the fields of the data, and obtain the final DataFrame. This allows us to work with the data as a structured DataFrame, which is more convenient for processing.

4.4.3 Writing streaming DataFrames

We have to remember that Spark is a lazy evaluation engine, so we need to call an action to start the execution of the streaming DataFrame. The most common action to start the execution of a streaming DataFrame is the `writeStream` method. This method allows us to write the streaming DataFrame to a sink, which can be a file, a Kafka topic, or a console.

To do so, we need to specify the following options:

- **Output sink:** The sink to write the data to.
- **Output mode:** The mode to use when writing the data, i.e., append, complete, or update.
- **Query name:** (Optional) Name of the query, which can be used to identify the query in the Spark UI.
- **Trigger interval:** (Optional) The interval at which to trigger the query, i.e., processing time or event time.
- **Checkpoint location:** (Optional) The location to store the checkpoint data, which is used to recover from failures.

We have 3 different output modes:

- **Append:** Only the new rows added to the result table since the last trigger will be written to the sink. This is only supported for queries where rows added to the result table are never updated. For example, queries with only `select`, `where`, `map`, `filter`, `flatMap`, `join`, etc.
- **Complete:** The entire result table will be written to the sink after every trigger. This is dangerous, as it can cause the sink to grow indefinitely, and can lead to out-of-memory errors. It is only supported for aggregation queries.
- **Update:** Only the rows in the result table that were updated since the last trigger will be written to the sink. Supported for aggregation queries and few others, not for joins.

4.4.4 Window operations on Event Time

Spark Structured Streaming supports only logical windows on event time. It allows for declaring tumbling, hopping and session windows, and it treats windows as a particular grouping criteria. For example, suppose we use the example of the fire alarms. We have:

```
1 temperature_sdf.groupBy(window("TS", "1 minutes", "30 seconds"), "SENSOR")
```

This will group the data in 1-minute windows, with a 30-second slide, and by sensor. This will create a new column with the window start and end times, and group the data by this column.

It is important to note that Spark Structured Streaming does not support physical windows on event time, as it is a continuous processing engine. This means that it does not have the concept of windows in the traditional sense, and does not have the notion of window boundaries. Instead, it treats windows as a logical grouping criteria, and processes the data continuously as it arrives.

4.4.5 Joins

Spark Structured Streaming supports joins between two streaming DataFrames. The following table describes the types of joins that are supported:

Inner	Supported , optionally specify watermark on both sides + time constraints for state cleanup
Left Outer	Conditionally supported, must specify watermark on right + time constraints for correct results, optionally specify watermark on left for all state cleanup
Right Outer	Conditionally supported, must specify watermark on left + time constraints for correct results, optionally specify watermark on right for all state cleanup
Full Outer	Conditionally supported, must specify watermark on one side + time constraints for correct results, optionally specify watermark on the other side for all state cleanup
Left Semi	Conditionally supported, must specify watermark on right + time constraints for correct results, optionally specify watermark on left for all state cleanup

Figure 4.13: Spark joins

Note that, comparing Spark Structured Streaming with EPL, Spark does not support the "followed by" operator (\rightarrow). This is because Spark is a continuous processing engine, and does not have the concept of events that are followed by other events. To achieve the same result, we need to use a stream-to-stream join with temporal constraints on the events timestamps. For example:

```
1 join_sdf = smoke_events.join(  
2     high_temp_events, expr("""  
3         (sensorTemp == sensorSmoke) AND  
4         (tsTemp > tsSmoke) AND  ## this is the EPL's -> operator  
5         (tsTemp < tsSmoke + interval 2 minute)  ## this is the timer:within clause  
6     """)  
7 )
```


Note that this will not tame the torrent effect, as it will report every event that matches the condition, but this is expected. These two joins are equivalent:

```
1 # Spark
2 join_sdf = A.join(B, expr("""
3 (idA == idB) AND
4 (tsB > tsA ) AND
5 (tsB < tsA + interval 2 minute )
6 """))
```

```
1 # EPL
2 every x=A ->
3     every B(id=x.id)
4     where timer:within(2 minutes)
```

Spark Structured Streaming also supports stream-to-table joins, which allow us to join a streaming DataFrame with a static DataFrame. This is useful when we have a static lookup table that we want to join with a streaming DataFrame. For example, suppose we have a static DataFrame with the sensor information, and we want to join it with the streaming DataFrame with the temperature data. We can do this as follows:

```
1 # Create a static DataFrame with the sensor information
2 sensor_df = spark.read.csv("sensor.csv", header=True)
3
4 # Join the streaming DataFrame with the static DataFrame
5 join_sdf = temperature_sdf.join(sensor_df, "SENSOR")
```

4.4.6 Late arrivals

Late arrivals refer to events that arrive after the expected time window in which they were supposed to be processed. This happens because in real-time streaming, data may experience delays due to temporarily disconnection, network high latency, processing issues, etc.

To handle late arrivals, Spark Structured Streaming provides a watermarking mechanism. Watermarking allows us to specify a threshold on event time, after which we consider the data to be late. This allows us to handle late arrivals by specifying a window on event time, and discarding data that arrives after the window has closed.

To use watermarking, we need to specify the following options:

- **Event time column:** The column that contains the event time.
- **Watermark delay:** The threshold on event time, after which we consider the data to be late.

The system will track the processing time P as the maximum seen event time (last event time), and if the difference between P and the event time of the event is smaller than the watermark delay, the event is considered on time. Otherwise, it is considered late, and it is discarded.

To use watermarking, we need to specify the event time column and the watermark delay in the query. For example, suppose we have a streaming DataFrame with the temperature data, and we want to handle late arrivals. We can do this as follows:

```

1 # Create a watermark column
2 watermarked_df = temperature_sdf.withWatermark("TS", "2 minutes")
3
4 # Group the data by window and sensor
5 windowed_df = watermarked_df.groupBy(window("TS", "1 minutes", "30 seconds"), "
  SENSOR")

```

In this example, we create a watermark column on the TS column with a delay of 2 minutes. This allows us to handle late arrivals by specifying a window on event time, and discarding data that arrives after the window has closed. We then group the data by window and sensor, and process the data accordingly.

The following graph shows how watermarking works, in a windowed groped aggregation with update mode:

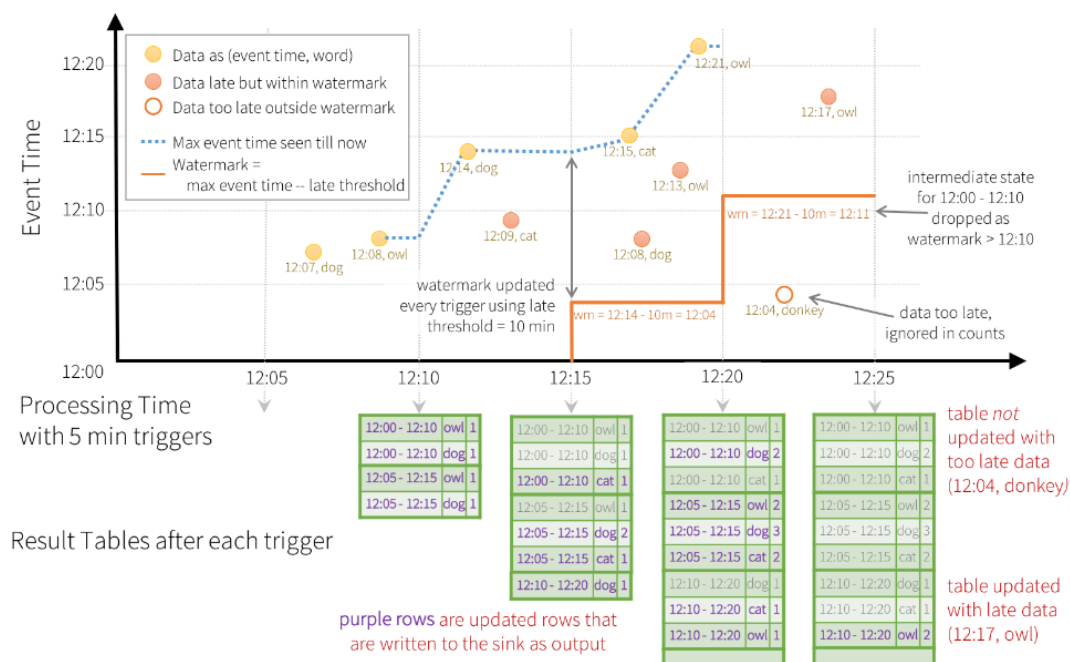


Figure 4.14: Spark watermarking in update mode

And the following graph shows how watermarking works, in a windowed groped aggregation with append mode:

