The slide or note you've shared highlights **three major complexities in stream processing** — data, systems, and workloads. Let's break each one down with context and examples:

---

## 1. COMPLEX DATA

✅ **Challenges:**

- **Diverse data formats**: Data might arrive in different formats like:

    - **JSON** – readable structured data

    - **Avro** – compact binary format for serialization

    - **Binary** – unreadable raw data that needs decoding

- **Dirty data**: Incoming streams may have:

    - Missing fields, corrupt records, or wrong types (e.g., timestamp is a string instead of a number)

- **Late or out-of-order data**:

    - Events might not arrive in the order they were generated.

    - For example, a sensor might send data with timestamps like `10:00`, `10:02`, `10:01`.

🛠️ **What this means for stream processing:**

You need tools that can **parse multiple formats**, **clean or validate data in real-time**, and **handle event-time vs. processing-time** (using watermarks or windowing logic).

---

## 2. COMPLEX SYSTEMS

✅ **Challenges:**

- **Diverse storage systems**:

○ A stream processor might pull data from **Kafka**, but enrich it using **RDBMS**, store results in **S3**, and send alerts via **Kinesis**.

- **System failures**:

    ○ Streams need to be **fault-tolerant** — what happens if Kafka goes down? Or your job crashes?

    ○ You must support **exactly-once or at-least-once** semantics depending on your use case.

🛠️ **What this means:**

Your processing engine (like **Flink**, **Spark Streaming**, or **Kafka Streams**) needs to be able to **connect to various sources/sinks**, **recover from crashes**, and **maintain consistency**.

---

## 3. COMPLEX WORKLOADS

✅ **Challenges:**

- **Streaming + interactive queries**:

    ○ Users might want dashboards that run **real-time + historical** queries together.

    ○ For example, show today's order rate **vs.** average order rate over the last week.

- **Machine learning**:

    ○ Stream processing might involve:

        ■ Real-time predictions (e.g., fraud detection)

        ■ Online learning (models updated continuously)

        ■ Model inference in the stream

🛠️ **What this means:**

You need systems that allow **low-latency queries**, support **caching**, and can **integrate ML pipelines** either via APIs or embedding models into the stream logic.

---

📌 **Summary Table:**

| Complexity | Challenges | Needs |
|---|---|---|
| Complex Data | Dirty, out-of-order, multi-format data | Schema handling, cleansing, windowing |
| Complex Systems | Many storage layers, system failures | Integration, checkpointing, exactly-once semantics |
| Complex Workloads | Real-time + batch, ML inference | Low-latency queries, scalable ML integration |

Would you like a visual diagram or example pipeline showing how these come together in a real system (e.g., fraud detection or real-time analytics)?

This slide titled **"Anatomy of a Streaming Query"** explains the **source** part of a streaming pipeline in **Apache Spark Structured Streaming**.

---

◆ **Code Breakdown:**

```python
CopyEdit
spark.readStream
    .format("kafka")
    .option("subscribe", "input")
    .load()
```

✅ **What it does:**

- **`spark.readStream`**: Tells Spark to read data as a **streaming source**, rather than a static DataFrame.

- **`.format("kafka")`**: Specifies that the source is Kafka.

- **`.option("subscribe", "input")`**: Subscribes to a Kafka topic named `"input"`.

- **`.load()`**: Actually loads the stream and returns a streaming DataFrame.

---

◆ **Right-hand Side Explanation:**

**Source – Key Concepts:**

- **"Specify one or more locations to read data from"**
  You can stream from a single source or multiple (e.g., Kafka, files, sockets).

- **"Built-in support for Files/Kafka/Socket, pluggable."**
  Spark comes with built-in connectors for:

    ○ **Kafka** (common for real-time logs, events)

    ○ **Files** (like CSV or JSON in a folder being populated over time)

    ○ **Socket** (like streaming text from a TCP server)

    ○ You can plug in other custom sources (e.g., MQTT, Kinesis, Pub/Sub).

**"Can include multiple sources of different types using `union()`"**
If needed, you can combine streams from multiple sources (like Kafka + file-based ingestion) using:

```python
CopyEdit
df1.union(df2)
```

- 

---

🧠 **Real-world Use Case:**

Imagine you're building a real-time analytics dashboard:

- You want to **read events from Kafka** (user clicks).

- You also want to **read config files** from S3 that update periodically.

- Spark lets you read both as streaming sources and **merge them** with `union()` to power downstream analytics.

## Code Breakdown:

```python
CopyEdit
spark.readStream
    .format("kafka")
    .option("subscribe", "input")
    .load()
    .groupBy('value.cast("string") as key')
    .agg(count("*") as 'value')
```

✅ **What it does:**

- **Reads from Kafka** like in the previous slide.

- **Casts the Kafka `value` to a string**, and **groups** the data by that value (acting as a key).

- **Aggregates** (counts) how many times each key appears in the stream.

This transformation helps perform **real-time analytics**, like counting how often a message appears in the stream.

---

🔹 **Right-hand Side Explanation:**

**Transformation – Key Concepts:**

- **"Using DataFrames, Datasets and/or SQL"**
  Spark allows you to write transformations using:

  - **DataFrame API**: like `.groupBy().agg()`

  - **Datasets** (typed APIs in Scala/Java)

  - **SQL**: by registering temporary views and running SQL queries.

- **"Catalyst figures out how to execute the transformation incrementally"**
  Spark's **Catalyst optimizer**:

  - Automatically plans how to run the transformation.

- ○ Ensures it only processes **new data (micro-batches)**, not redoing everything.

- ● **"Internal processing always exactly-once"**
  Spark guarantees **exactly-once semantics**:

  - ○ Each event is processed **only once**, even during recovery from failure.

  - ○ Checkpointing and write-ahead logs make this possible.

---

## 🧠 Real-world Use Case:

Let's say you're building a real-time **sentiment dashboard**:

- ● Users are sending reviews via Kafka.

- ● You want to **count how many reviews are positive, negative, or neutral** in real-time.

- ● You can group by the `sentiment` label and count them — exactly like this slide shows.

### 🔷 Code Block (Left Side):

```python
CopyEdit
input = spark.readStream
    .format("kafka")
    .option("subscribe", "topic")
    .load()

result = input \
    .select("device", "signal") \
    .where("signal > 15")

result.writeStream \
    .format("parquet") \
    .start("dest-path")
```

### ✅ What this does:

- ● **Reads** streaming data from a Kafka topic.

- **Selects** the columns `"device"` and `"signal"`.

- **Filters** only the rows where signal > 15.

- **Writes** the result to a Parquet file in `dest-path`.

---

### ◆ Center: Spark Planning Breakdown

**1. Logical Plan**

- This is what Spark generates from your **DataFrame operations**.

- It includes abstract steps:

  - **Read from Kafka**

  - **Project** (select `device`, `signal`)

  - **Filter** (`signal > 15`)

  - **Write to Kafka (or sink)**

**2. Optimized Physical Plan**

- Spark's **Catalyst optimizer** converts the logical plan into an efficient execution plan.

- Optimization includes:

  - **Operator fusion** (combining steps)

  - **Code generation** (for speed)

  - **Off-heap memory use**

- The plan includes:

  - Kafka as **source**

  - **Optimized operators**

    ○   Kafka (or other) as **sink**

---

◆ **Right: Incremental Execution**

🔁 **At each time step `t=1, t=2, t=3...`**

- Spark **processes only the new data** arriving in the stream (called a **micro-batch**).

- It applies the same optimized plan on just the new input:

    ○   Filters rows

    ○   Projects columns

    ○   Writes results

- Spark creates a **new execution plan for each micro-batch**, enabling continuous and efficient streaming

◆ **Code Breakdown (Left Side):**

```python
CopyEdit
spark.readStream
    .format("kafka")
    .option("subscribe", "input")
    .load()
    .groupBy('value.cast("string") as key')
    .agg(count("*") as 'value')
    .writeStream
    .format("kafka")
    .option("topic", "output")
```

✅ **What it does:**

- Reads from a Kafka topic called `"input"`.

- Groups and counts messages by their value.

- Writes the output to a **Kafka sink**, publishing results to a topic named `"output"`.

---

◆ **Right Side Explanation – Sink Concepts:**

✅ **What is a Sink?**

A **sink** is the destination where **each processed micro-batch** is sent. Spark supports a variety of sinks like:

- **Kafka** – to publish back into a stream

- **File systems** – to store in formats like CSV, Parquet, etc.

- **Console** – for debugging

- **Memory** – for in-app querying (not for production)

---

🧠 **Key Points:**

- **"Accepts the output of each batch."**
  Spark processes streaming data in **micro-batches**, and each batch's result is sent to the sink.

- **"Transactional and exactly once (when supported)."**
  Some sinks like **files (e.g., Parquet)** and **Kafka (with idempotent writes)** support **exactly-once delivery semantics** to avoid duplicates or data loss.

- **"Use foreach to execute arbitrary code."**
  If you need to:

    - Send results to a **custom API**

    - Write to **external databases**

    - Trigger **notifications**
      You can use `.foreach()` to run **custom code** for each row or batch.

---

## 💡 Example Use Cases for Sinks:

- Write aggregated metrics to a **dashboard via Kafka**

- Archive processed logs to **Amazon S3** using the file sink

- Store real-time fraud detection results into a **NoSQL DB** using `foreachBatch`



```
spark.readStream
  .format("kafka")
  .option("subscribe", "input")
  .load()
  .groupBy('value.cast("string") as key')
  .agg(count("*") as 'value')
  .writeStream
  .format("kafka")
  .option("topic", "output")
  .trigger("1 minute")
  .outputMode("update")
  .option("checkpointLocation", "…")
  .start()
```

## 🧠 Real-World Example:

Let's say you're monitoring sensor data from IoT devices:

- You count how often a sensor reports a specific status (like "ALERT").

- You group and aggregate per device and send that to Kafka every minute.

Without checkpointing:

- If the job fails at 10:03, all progress is lost.

- When restarted, it reprocesses from the beginning (duplicates results, re-counts data).

With checkpointing:

- If the job fails at 10:03, the last checkpoint (say at 10:02) allows it to resume from there.

- It ensures **exactly-once** output.

---

◆ **Best Practices:**

- Always set a **valid, persistent location** for `"checkpointLocation"` (e.g., `/mnt/checkpoints/iot-query`)

- Don't share a checkpoint directory across different queries

- Make sure the location is accessible to all nodes in the cluster

This slide dives deeper into **fault tolerance with checkpointing** in Spark Structured Streaming, explaining how it ensures **exactly-once guarantees** across failures or restarts.

---

◆ **Core Concept: Checkpointing**

**Definition (from slide):**

> *"Checkpointing – tracks progress (offsets) of consuming data from the source and intermediate state."*

✅ **What gets tracked?**

- **Offsets**: The last successfully read position from a data source (e.g., Kafka offset).

- **State**: If you're doing stateful operations (like `groupBy` + `agg`), intermediate aggregation values are stored.

- **Metadata**: Info about execution plans, progress, schema, etc.

These are **persisted as JSON** in the **checkpoint directory** you define using:

.option("checkpointLocation", "/some/path/")

---

### ◆ Visual Explanation (Right side)

The timeline at **t = 1, 2, 3** shows:

- New data being processed in each micro-batch.

- After processing, **metadata and results are recorded into a write-ahead log**.

- In case of failure (e.g., crash after t=2), Spark:

  - Reloads the **checkpoint data**

  - Avoids reprocessing t=1 and t=2

  - Resumes cleanly from **t=3**

This is how Spark achieves **end-to-end exactly-once guarantees**.

---

### 🧠 Real-World Example:

Imagine you're processing **credit card transactions** in real time to detect fraud.

**Use case:**

- Stream transaction data from Kafka.

- Use stateful logic: group by `card_id` and keep a running sum of spend in the last 10 minutes.

- Send alerts when suspicious thresholds are exceeded.

**Without checkpointing:**

- If Spark crashes, **state is lost**.

- It may miss alerts or double-send them (if it reprocesses old Kafka records).

**With checkpointing:**

- Spark saves:

  - Which transactions were already read (Kafka offsets)

  - Intermediate state (running sums)

- After a crash, Spark:

  - Picks up where it left off

  - Keeps your **fraud detection logic consistent**

---

◆ **Key Takeaways:**

| Feature | Why It Matters |
| --- | --- |
| **Offsets** tracked | Prevents reprocessing of old data |
| **State** tracked | Ensures aggregations/window logic can continue |
| **Metadata as JSON** | Human-readable and easy to debug |
| **Resilience after crash** | Recovers query without duplicating output |
| **Exactly-once guarantee** | Critical for finance, alerts, analytics |

---

Would you like to go over how `outputMode` or `trigger` settings affect checkpointing next?

```scala
val rawData = spark.readStream
  .format("kafka")
  .option("kafka.boostrap.servers",...)
  .option("subscribe", "topic")
  .load()

val parsedData = rawData
  .selectExpr("cast (value as string) as json"))
  .select(from_json("json", schema).as("data"))
  .select("data.*")

val query = parsedData.writeStream
  .option("checkpointLocation", "/checkpoint")
  .partitionBy("date")
  .format("parquet")
  .start("/parquetTable")
```

databricks

Great! You're referring to **Kafka source configuration options** used in **Spark Structured Streaming** via `.option(...)` when reading from Kafka.

Let's walk through this block **line by line**, explaining:

- **How to use each option**

- **What it means**

- **What the syntax does**

---

## ✅ How? (Basic connection to Kafka)

```scala
```

```scala
.option("kafka.bootstrap.servers", "broker1:9092,broker2:9092")
```

- **Required**: This is the **Kafka cluster** your stream will connect to.
- Provide a **comma-separated list** of Kafka brokers.
- This does **not** mean your app connects to all brokers — Kafka will handle that internally.

---

## ✅ What? (How to specify topics or partitions)

### Option 1: **Fixed list of topics**
```scala
.option("subscribe", "topic1,topic2,topic3")
```

- Use when you **know exactly** which topics you want to read from.

- Spark reads **all partitions** from these topics.

---

## Option 2: Dynamic topic pattern
```
.option("subscribePattern", "topic.*")
```

- Use a **regex pattern** to match topics dynamically.

- Spark will pick up any topic that matches `topic.*` — useful for logs, metrics, etc.

- **Topics created later** that match the pattern will be picked up too.

---

## Option 3: Specific partitions (manual assignment)
```
.option("assign", """{"topicA":[0,1]}""")
```

- JSON string specifying **which partitions to read** from which topic.

- Use this to control **exact partition-level consumption**.

- This **disables consumer group balancing** — you're fully in control.

---

# ✅ Where? (Offsets — where to start reading from)

## Option: Set starting offset

.option("startingOffsets", "latest")

- Controls **where Spark starts reading from**:

  - `"latest"`: Only read **new messages** (messages produced **after** job starts).

  - `"earliest"`: Read **everything** from the beginning.

## Advanced: Specific offsets per partition

.option("startingOffsets", """{"topicA":{"0":23,"1":345}}""")

- Start reading from **offset 23 on partition 0** and **345 on partition 1** of topicA.

- Useful for **resuming** from a specific point.
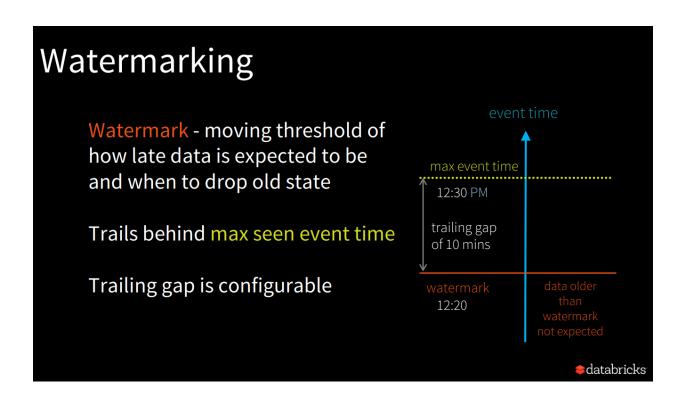
---

# 🧠 Example Code Putting It All Together

val df = spark.readStream
  .format("kafka")
  .option("kafka.bootstrap.servers", "localhost:9092")
  .option("subscribe", "logs,metrics")
  .option("startingOffsets", "earliest")
  .load()

✅ This connects to Kafka at `localhost:9092`, subscribes to `logs` and `metrics`, and starts reading from the **beginning**.

---

## 🔁 Summary Table

| Option | Purpose | Value Type |
| --- | --- | --- |
| `kafka.bootstrap.servers` | Connect to Kafka cluster | Comma-separated |
| `subscribe` | Read from fixed topics | Comma-separated |
| `subscribePattern` | Read from topics matching a regex | Regex string |
| `assign` | Read specific topic partitions | JSON string |
| `startingOffsets` | Where to begin reading messages | "latest", "earliest", or JSON |

Let me know if you want to go over **endingOffsets**, checkpoint integration, or auto offset commits next!



# Event time Aggregations

### Windowing is just another type of grouping in Struct. Streaming

number of records every hour

```
parsedData
    .groupBy(window("timestamp","1 hour"))
    .count()
```

avg signal strength of each device every 10 mins

```
parsedData
    .groupBy(
        "device",
        window("timestamp","10 mins"))
    .avg("signal")
```

## Support UDAFs!

This slide introduces **Arbitrary Stateful Operations** using Spark's `mapGroupsWithState` (introduced in **Spark 2.2**), which allows **fine-grained control over per-key state** in streaming applications.

Let's break it down with simple language and an example 👇

---

## ◆ What is `mapGroupsWithState`?

### → A powerful Spark API for stateful stream processing

It allows you to:

- Define your own **custom state** per key (like a counter, a timer, or any object)

- Update, access, or remove this state as new data arrives

- Set **timeouts per key** to clean up unused state

---

# 🔧 Syntax Overview from the Slide

ds.groupByKey(_.id)
  .mapGroupsWithState(timeoutConf)(mappingWithStateFunc)

| Component | Role |
|---|---|
| `ds` | A dataset (structured streaming input) |
| `groupByKey(_.id)` | Group data by some key (e.g., user ID, sensor ID, etc.) |
| `mapGroupsWithState` | Apply a **user-defined function** that manages state per group |
| `timeoutConf` | How long to retain state per key (event-time or processing-time) |
| `mappingWithStateFunc` | Your custom function that defines how to update and return output |

---

## 🔹 The Function Signature

def mappingWithStateFunc(
  key: K,
  values: Iterator[V],
  state: GroupState[S]
): U = {
  // Your logic here
}

## Parameters:

- `key`: the current group key (e.g., a user ID)

- `values`: the new events for this key (batch of events)

- `state`: the **current state** for this key (you can read, update, or remove it)

## Return:

- A result (U) for each group, like a transformed row or alert

---

## 🧠 **Real-World Example: User Session Timeout**

Say you're tracking user activity:

- Each message has `userId`, `action`, `timestamp`

- You want to maintain a **session window** per user

- If no activity from a user for 30 minutes, mark the session as closed

You can do this with `mapGroupsWithState`:

- Store last action timestamp as state

- If new event arrives after timeout → emit session close

- Otherwise → update session state

---

## ✅ **Key Features Highlighted in the Slide**

| Feature | What it means |
|---|---|
| `mapGroupsWithState` | Custom function per key with control over state |
| User-defined state | You choose what to store (count, last timestamp, etc.) |
| Per-key timeout | Clean up unused state after a configurable period |
| Event-time or processing-time | Timeouts can be based on **event timestamp** or system clock |
| Supports Scala and Java | Works with typed Datasets in both languages |

---

## 📌 **Summary**

- Use `mapGroupsWithState` when you need **custom logic** and **flexible state handling**.

- It's ideal for:

    - **Sessionization**

    - **Fraud detection**

    - **Last-known state tracking**

    - **Custom time-based alerts**

Let me know if you'd like to see a full Scala example with `GroupState` logic!

Great! This slide introduces **three powerful operations** in Spark Structured Streaming:

---

## 🔶 1. Streaming Deduplication

### ✅ Purpose:

To **remove duplicate records** from a stream — for example, when:

- Data has retry logic or network issues

- You don't want to process the same event twice

### ✅ Code:

parsedData.dropDuplicates("eventId")

- This keeps **only the first record** per unique `eventId`

- Internally, Spark stores a state of seen `eventId`s to avoid re-processing

### ⚠️ Important:

You should combine this with **watermarking**:

```
parsedData
  .withWatermark("timestamp", "10 minutes")
  .dropDuplicates("eventId")
```

💡 This allows Spark to **forget old IDs after 10 minutes** and not hold memory forever.

---

## 🔶 2. Stream–Batch Joins

### ✅ Purpose:

Join **real-time data** with **static or slowly-changing data**, like:

- Enriching clickstream logs with user profile info

- Joining incoming orders with product catalog

### ✅ Code:

```
val batchData = spark.read
  .format("parquet")
  .load("/additional-data")

parsedData.join(batchData, "device")
```

- `batchData` is **loaded once** (not continuously)

- Spark joins every incoming streaming record with the batch table on `device`

### ⚠️ Tip:

The batch table should be small enough to **broadcast** to executors.

---

## 🔶 3. Stream–Stream Joins

### ✅ Purpose:

Join **two real-time data streams** together, like:

- Joining sensor readings with device status updates

- Matching orders with payments in real time

## ✅ Current Options:

- Use `mapGroupsWithState` for custom logic

- (Back when this slide was made) direct built-in support was *coming soon*
  → **Now it exists!** Spark supports **stream-stream joins** with watermark-based logic.

## ✅ Example (modern syntax):

```
stream1
  .withWatermark("ts1", "10 minutes")
  .join(
    stream2.withWatermark("ts2", "10 minutes"),
    expr("stream1.id = stream2.id AND ts1 BETWEEN ts2 - interval 5 minutes AND ts2 + interval
5 minutes")
  )
```

## 📌 Summary Table

| Operation | Description | Key Spark Features Used |
|---|---|---|
| Streaming Deduplication | Removes duplicates using a unique key | `dropDuplicates()`, Watermark |
| Stream–Batch Join | Enriches stream with static reference data | `join()` |
| Stream–Stream Join | Combines two streams using keys + time conditions | `join()` + Watermarks + Event-time windows |

Let me know if you'd like to see examples of **handling late data in joins** or how **stream–batch joins** are optimized!