

Project 2

Stream Processing System — Report

Abril Cano Castro - Emanuele Cimino - Jose Luiz

0.1 Purpose

The aim of this project is to write a stream processing system, starting from a message flow of `<key, value>` pairs.

0.2 Sensor-data producer

To begin, we implemented a simple Kafka producer that creates the initial topic (if needed) and configures the required properties.

Table 1: Kafka Producer Properties: meaning and purpose.

Property	Meaning	Effect / Why Used
Connection & Serialization		
<code>bootstrap.servers</code>	List of Kafka brokers to connect to.	The producer contacts these brokers to discover the cluster and send data.
<code>key.serializer</code>	Class to serialize the message key.	Converts the key into bytes before sending (here: <code>StringSerializer</code>).
<code>value.serializer</code>	Class to serialize the message value.	Converts the value into bytes before sending (here: <code>StringSerializer</code>).
Exactly-Once Semantics (EOS) / Reliability		
<code>enable.idempotence</code>	Enables idempotent producer (no duplicates).	Guarantees no duplicate messages even on retries.
<code>acks</code>	Number of replicas that must acknowledge a message.	<code>all</code> waits for all in-sync replicas (strongest delivery guarantee).
<code>retries</code>	Maximum number of retries on failure.	<code>Integer.MAX_VALUE</code> retries indefinitely until success.
Throughput / Performance		
<code>linger.ms</code>	Time to wait before sending a batch.	Allows messages to accumulate (20 ms) for larger batches and better throughput.
<code>batch.size</code>	Maximum batch size in bytes for one partition.	32 KB per batch → fewer network requests.

Property	Meaning	Effect / Why Used
<code>delivery.timeout.ms</code>	Total time allowed for message delivery.	120 s before giving up, including retries.
<code>request.timeout.ms</code>	Timeout for a single request to a broker.	30 s; if the broker does not respond, retry.
<code>transaction.timeout.ms</code>	Max allowed time for an open transaction.	180 s; the broker aborts if the transaction is not completed.

0.3 Sensor data consumer and output producer

Here we list the properties for the consumer of produced data and the producer of aggregated values, after the transformations performed with Akka.

Table 2: Kafka Consumer and Producer Properties used in the system.

Property	Meaning	Effect / Why Used
Consumer Properties		
<code>bootstrap.servers</code>	List of Kafka brokers to connect to.	Tells the consumer which brokers to contact to fetch data.
<code>group.id</code>	Identifier of the consumer group.	Allows multiple consumers to share the same topic load, coordinating partitions among them.
<code>key.deserializer / value.deserializer</code>	Classes used to deserialize key/value from bytes to objects.	Here: <code>StringDeserializer</code> converts <code>byte[]</code> to <code>String</code> .
<code>auto.offset.reset</code>	What to do if there is no committed offset.	<code>earliest</code> starts from the beginning of the topic; <code>latest</code> starts from new messages.
<code>enable.auto.commit</code>	Enables automatic offset commits by the consumer.	Disabled (<code>false</code>) to allow transactional commits together with producer output.
<code>isolation.level</code>	Controls which records are visible to consumers in case of transactions.	<code>read_committed</code> reads only messages from committed transactions (EOS support).
<code>max.poll.records</code>	Maximum number of records returned in a single poll.	Limits batch size per poll (e.g., 500).
<code>max.poll.interval.ms</code>	Maximum delay between two <code>poll()</code> calls before the consumer is considered dead.	300 s by default; prevents rebalancing if processing takes too long.
Producer Properties		
<code>bootstrap.servers</code>	List of Kafka brokers to connect to.	Producer uses these brokers to send data.
<code>key.serializer / value.serializer</code>	Classes to serialize key/value to <code>byte[]</code> .	Here: <code>StringSerializer</code> converts <code>String</code> to <code>byte[]</code> .

Property	Meaning	Effect / Why Used
<code>enable.idempotence</code>	Enables idempotent producer (no duplicates).	Required for exactly-once semantics and safe retries.
<code>acks</code>	Number of replicas that must acknowledge a message.	<code>all</code> = wait for all in-sync replicas (strongest durability).
<code>retries</code>	Maximum number of retries on failure.	<code>Integer.MAX_VALUE</code> = retry indefinitely.
<code>linger.ms</code>	Time to wait before sending a batch.	Accumulates more messages (20 ms) → better batching and throughput.
<code>batch.size</code>	Max batch size per partition.	32 KB → fewer network requests.
<code>compression.type</code>	Compression algorithm for batches.	<code>lz4</code> gives fast compression and lower network usage.
<code>delivery.timeout.ms</code>	Total allowed time to deliver a message.	120 s including retries, before failure.
<code>request.timeout.ms</code>	Timeout for a single request to the broker.	30 s; if exceeded, a retry is triggered.
<code>transaction.timeout.ms</code>	Max allowed time for an open transaction.	180 s; broker aborts if transaction not finished.

In fact, this class also initializes cluster sharding, which is responsible for creating and routing messages to the `SensorActorSupervisor` actors that perform the aggregate calculations. It also creates a main `AggregateSink` actor, which collects all the results and publishes the aggregated data to the appropriate Kafka topic.

0.4 Akka actors

Table 3: Overview of Akka actors used for data aggregation.

Actor	Responsibility	Notes
<code>AggregateSink</code>	Collects the results produced by all sensor actors and acts as the final aggregation point.	Publishes the aggregated data to the correct Kafka topic.
<code>SensorActorSupervisor</code>	Supervisor actor for each sensor type. Manages window state, triggers aggregation, and forwards results to the sink.	Instantiated and sharded by Akka Cluster Sharding.
<code>TempActor</code>	Processes temperature readings.	Maintains a sliding/tumbling window for temperature data and performs aggregation (e.g., average).
<code>HumActor</code>	Processes humidity readings.	Similar to <code>TempActor</code> but for humidity data.

Actor	Responsibility	Notes
WindActor	Processes wind readings.	Can compute average/max wind speed over time windows.
AirActor	Processes air quality readings.	Aggregates AQI or pollutant levels over time.
SensorMessageExtractor	Extracts <code>entityId</code> and <code>shardId</code> from incoming messages.	Required by Akka Cluster Sharding to route messages to the correct actor instance.

Finally, we implemented a simple consumer for the aggregated data.