



eyouth

هيئة الاتصالات
وتقنيولوجيا المعلومات

Harnessing the Data Flood: A Real-Time Pipeline for Smart Building Management

A Data Engineering Project for the Digital Egypt Pioneers Initiative

Date : 12/11/2025

Presenters:

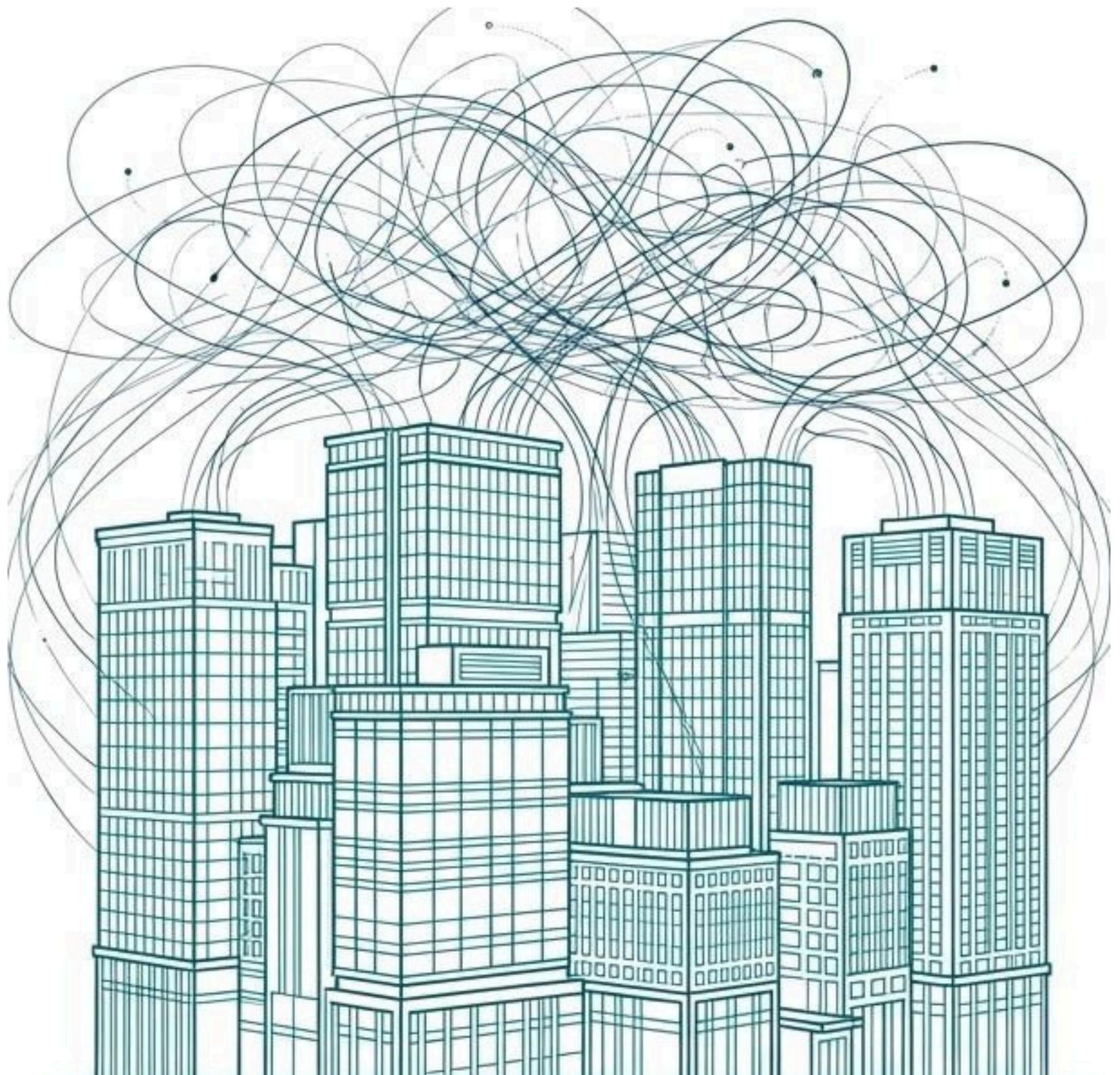
Mohamed Ahmed Abdelaziz
Rodina Amr
Hanin Salem
Youssef Tamer
Ganna Ehab
Mohamed Ibrahim

The Modern Building's Dilemma: Drowning in Data, Starving for Insight

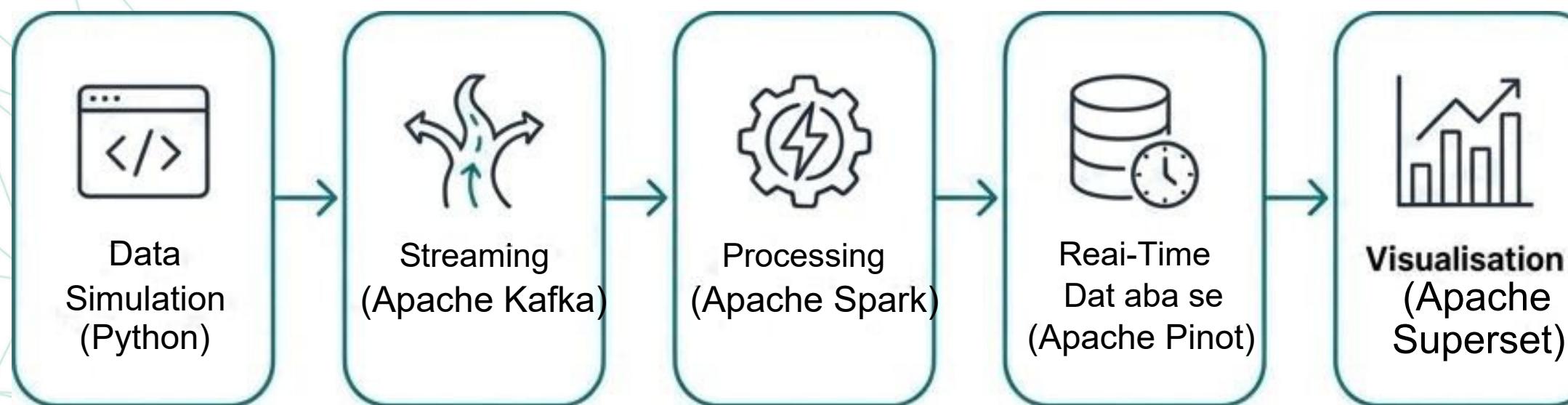
Modern buildings generate a massive, continuous stream of sensor data, including temperature, humidity, CO₂, motion, and energy consumption.

Without a proper data pipeline, this valuable information becomes messy, unorganised, and difficult to analyse in real time.

This directly leads to significant operational inefficiency, unnecessary energy waste, and poor performance decisions by facility managers.



Our Solution: An End-to-End, Real-Time Analytics Pipeline



End-to-End Real-Time Flow:
Continuous ingestion, processing, and visualisation with no batch delays.



Low-Latency Analytics: A powerful combination of Apache Pinot and Superset enables immediate querying and insights.



Scalable & Flexible: The architecture is designed to handle both streaming and batch workloads effectively.



Actionable Insights: Unlocks critical smart-building metrics for energy, air quality, and occupancy management.



The Blueprint: Our Technology Stack

Data Streaming



Apache Kafka

For high-throughput, fault-tolerant real-time data ingestion.

Data Processing



Apache Spark

For scalable batch and structured streaming data processing, cleaning, aggregation,

Real-Time Analytics Database



Apache Pinot

A real-time, distributed OLAP datastore, purpose-built for low-latency queries on fresh data.



Trino

A high-performance query engine that acts as the vital connector between Pinot and Superset.

Data Visualisation



Apache Superset

For creating interactive, real-time dashboards and visualising key performance indicators.

Languages & Supporting Tools



Python

Used for data simulation and processing logic.



SQL

The language for querying data in Pinot via Trino.



GitHub

For version control and collaborative development.

Anatomy of the Pipeline: Data Flow & Architecture

Detailed Data Flow

1.  A Python script simulates and generates real-time sensor data.
2.  Data is published to the Kafka topic `bms_data`.
3.  An Apache Spark streaming job consumes from `bms_data`, performs cleaning and aggregation, and publishes the processed data to a new Kafka topic, `"spark_data"`.
4.  The processed data is ingested from `'spark_data'` into Apache Pinot columnar tables.
5.  Trino queries Pinot's tables, exposing the data for visualisation.
6.  Superset dashboards query Trino to display real-time KPIs.

Database Architecture Rationale

Apache Pinot Rationale

Chosen for its columnar storage and optimisation for real-time OLAP queries, handling high-throughput streaming ingestion with powerful aggregation and filtering capabilities.

Key Entities

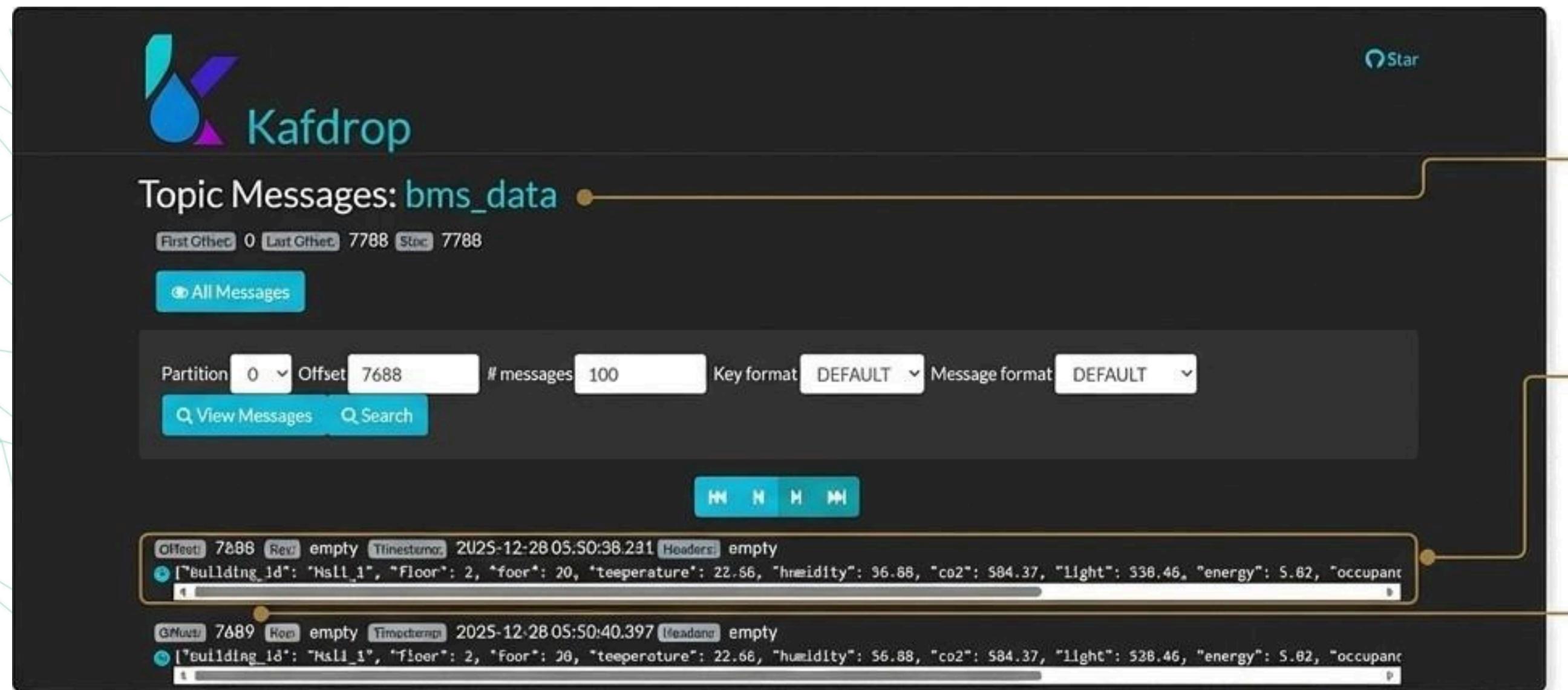
The data model centres on Sensors (temperature, humidity, CO₂, etc.) which produce timestamped Readings .





The Data's Journey, Part Real-Time Ingestion

The pipeline begins with a Python simulator generating continuous sensor readings. These raw data points are immediately published as messages to the Apache Kafka topic 'bns_data'. The screenshot below shows a live view of these messages arriving in the topic, confirming successful data capture at the entry point of our system.



The screenshot shows the Kafdrop interface with the topic 'Topic Messages: bns_data'. The interface includes filters for Partition (0), Offset (7688), # messages (100), Key format (DEFAULT), and Message format (DEFAULT). Below these controls, two messages are displayed:

- Message 1: Offset 7688, Time 2025-12-28 05:50:38.291, Headers empty. Payload: {"Building_Id": "Hall_1", "Floor": 2, "Floor": 20, "temperature": 22.66, "humidity": 56.88, "co2": 584.37, "light": 528.46, "energy": 5.82, "occupant": 1}
- Message 2: Offset 7689, Time 2025-12-28 05:50:40.397, Headers empty. Payload: {"Building_Id": "Hall_1", "Floor": 2, "Floor": 20, "temperature": 22.66, "humidity": 56.88, "co2": 584.37, "light": 528.46, "energy": 5.82, "occupant": 1}

The topic name:
"bns data".

The message payload,
showing fields like
"building_id",
"floor", "temperature",
"co2", and "energy".

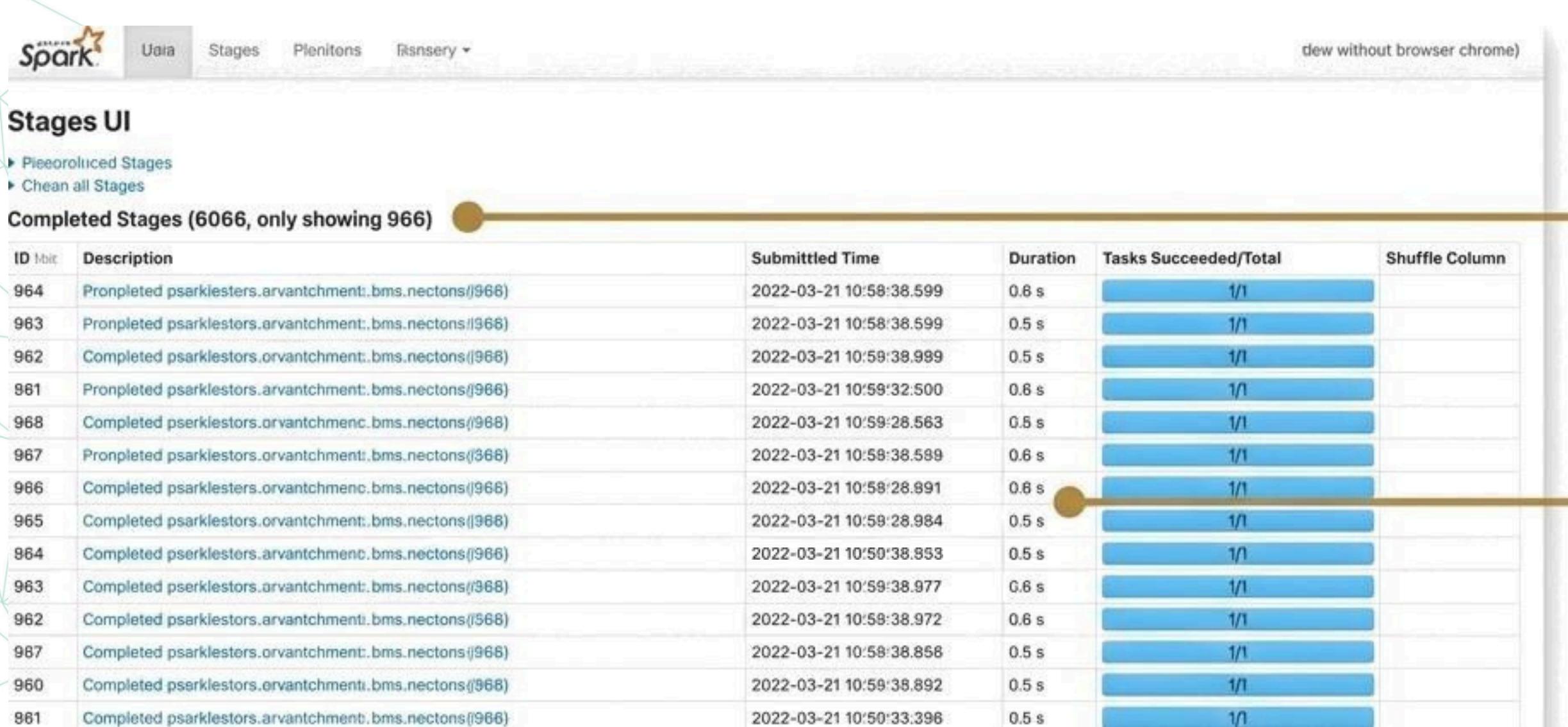
The constantly
increasing offset
number, indicating a
live stream.

The Data's Journey, Part 2: Stream Processing with Spark

An Apache Spark Streaming application continuously consumes raw data from the "bms_data" Kafka topic.

Spark jobs perform essential cleaning, transformation, and aggregation logic.

The Spark UI confirms that these processing stages are completing successfully and efficiently in near real-time, preparing the data for storage.

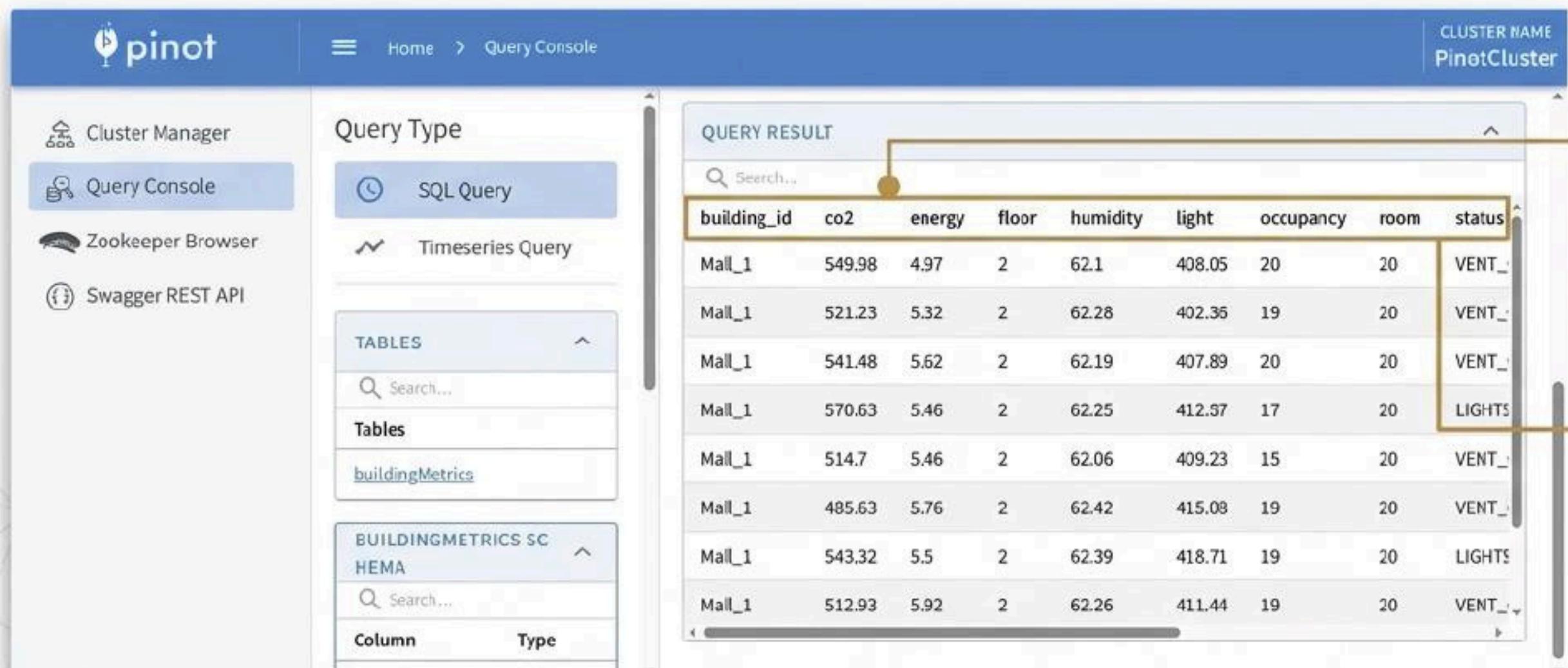


A list of thousands of successfully executed micro-batches, proving continuous processing.

The low duration of each stage underscores the high performance of the processing layer.

The Data's Journey, Part 3: Storing for Low-Latency Analytics

Cleaned and aggregated data from Spark is ingested into Apache Pinot. Pinot stores the data in a columnar format, highly optimised for the fast analytical queries required by our dashboards. The Pinot Query Console below shows the structured `buildingMetrics` table, confirming that processed data is available and ready for immediate querying.



The screenshot shows the Apache Pinot Query Console interface. On the left, there's a sidebar with links to Cluster Manager, Query Console (which is selected), Zookeeper Browser, and Swagger REST API. Under 'Query Console', there are tabs for 'SQL Query' (selected) and 'Timeseries Query'. Below these are sections for 'TABLES' (listing 'Tables' and 'buildingMetrics') and 'BUILDINGMETRICS SCHEMA' (listing 'Column' and 'Type'). The main area is titled 'QUERY RESULT' and shows a table with the following data:

building_id	co2	energy	floor	humidity	light	occupancy	room	status
Mall_1	549.98	4.97	2	62.1	408.05	20	20	VENT_
Mall_1	521.23	5.32	2	62.28	402.35	19	20	VENT_
Mall_1	541.48	5.62	2	62.19	407.89	20	20	VENT_
Mall_1	570.63	5.46	2	62.25	412.37	17	20	LIGHTS
Mall_1	514.7	5.46	2	62.06	409.23	15	20	VENT_
Mall_1	485.63	5.76	2	62.42	415.08	19	20	VENT_
Mall_1	543.32	5.5	2	62.39	418.71	19	20	LIGHTS
Mall_1	512.93	5.92	2	62.26	411.44	19	20	VENT_

Annotations on the right side of the screenshot explain the highlighted elements:

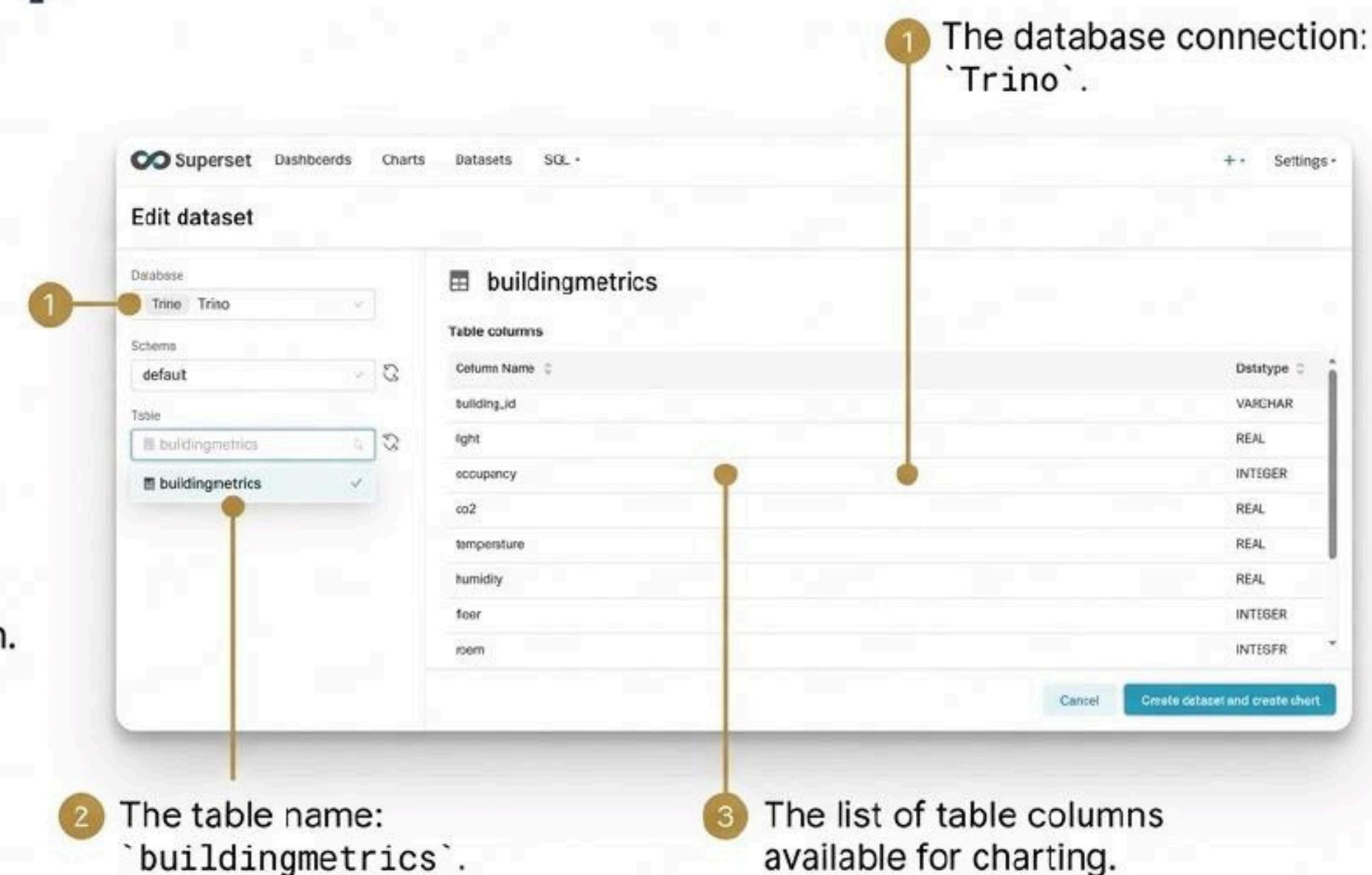
- A callout points to the 'buildingMetrics' table in the 'TABLES' section with the text: 'The selected table: "buildingMetrics".'
- A callout points to the table header in the 'QUERY RESULT' section with the text: 'The query result, showing structured columns like "building_id", "co2", "energy", "occupancy", etc.'



The Data's Journey, Part 4: Connecting to Superset via Trino

To make the data in Pinot available for visualisation, we use Trino as a federated query engine. Apache Superset connects to our Trino server, which in turn queries the Pinot database. This architecture allows Superset to leverage Pinot's speed while using standard SQL.

The screenshot shows the `buildingmetrics` dataset successfully configured within Superset, ready for dashboard creation.



1 The database connection: 'Trino'.

2 The table name: 'buildingmetrics'.

3 The list of table columns available for charting.

Column Name	Datatype
building_id	VARCHAR
light	REAL
occupancy	INTEGER
co2	REAL
temperature	REAL
humidity	REAL
floor	INTEGER
room	INTEGER



Energy Team Dashboard

Draft

hanin salim

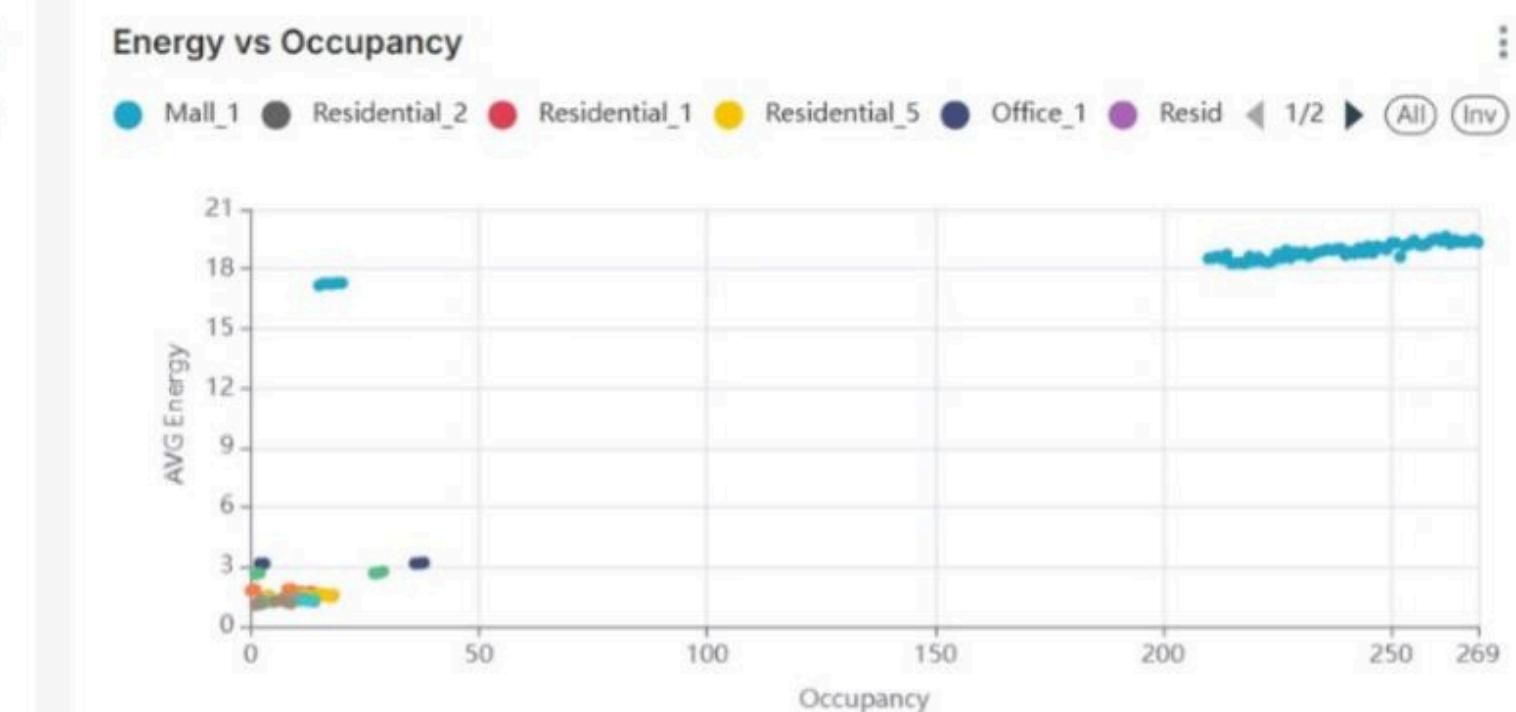
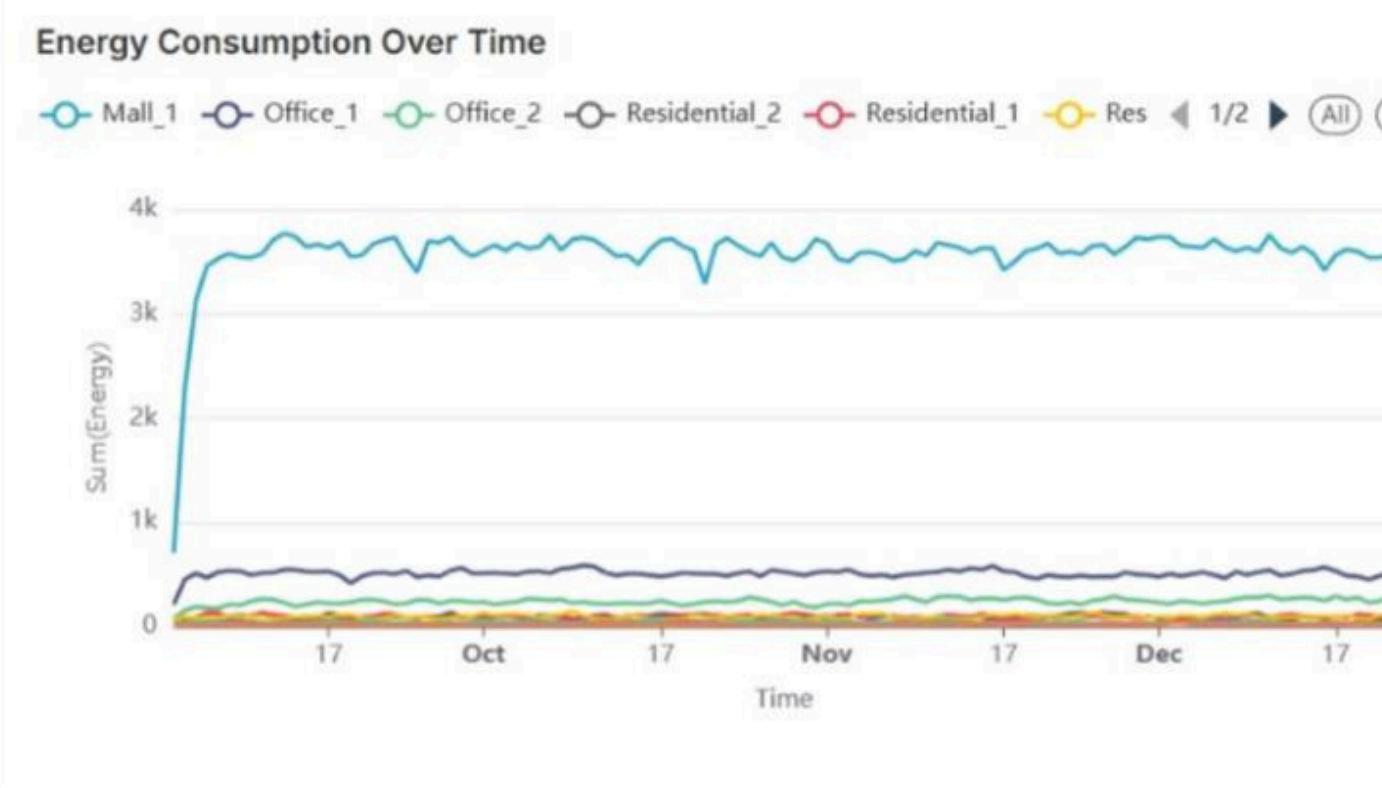
a minute ago

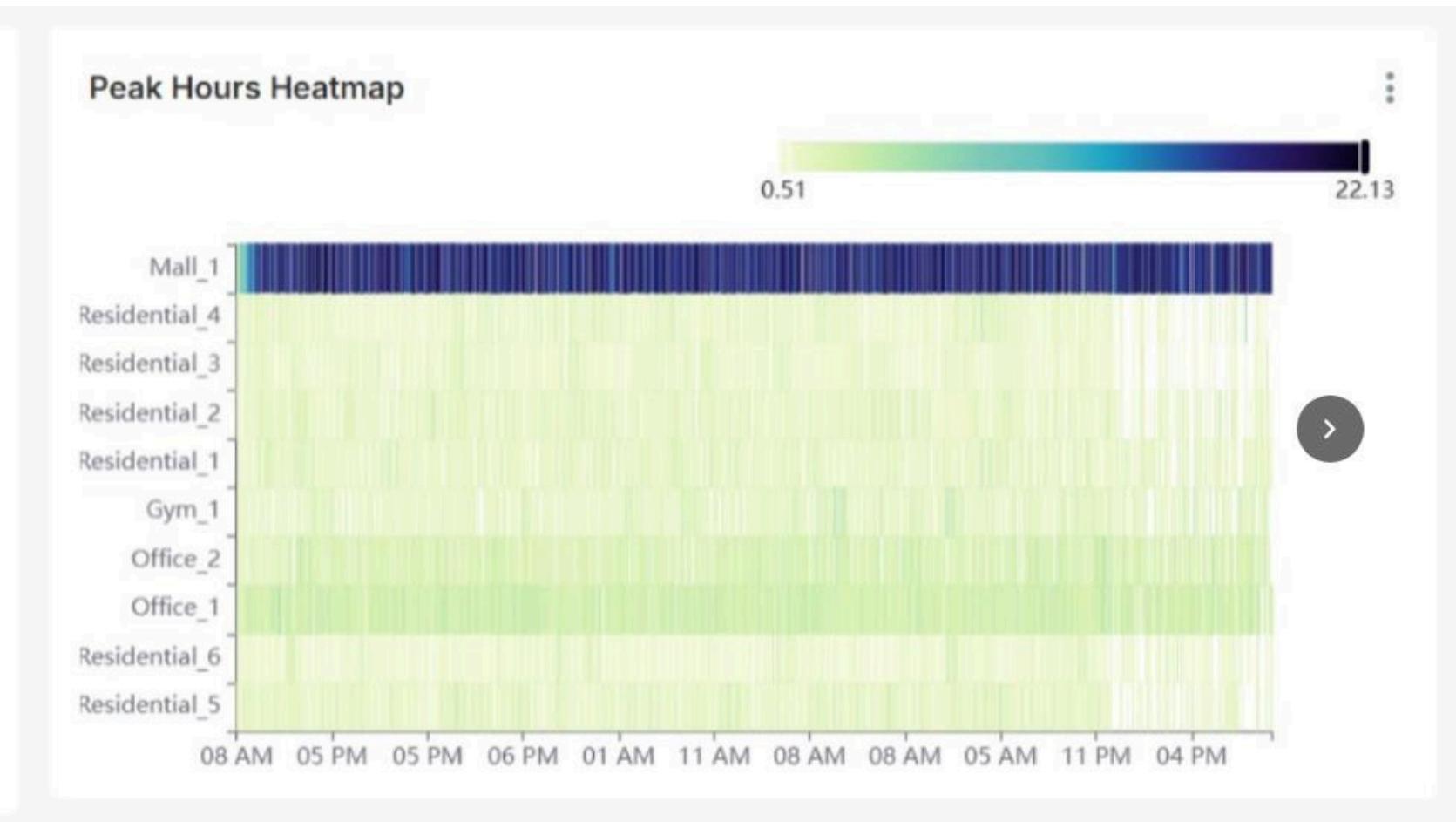
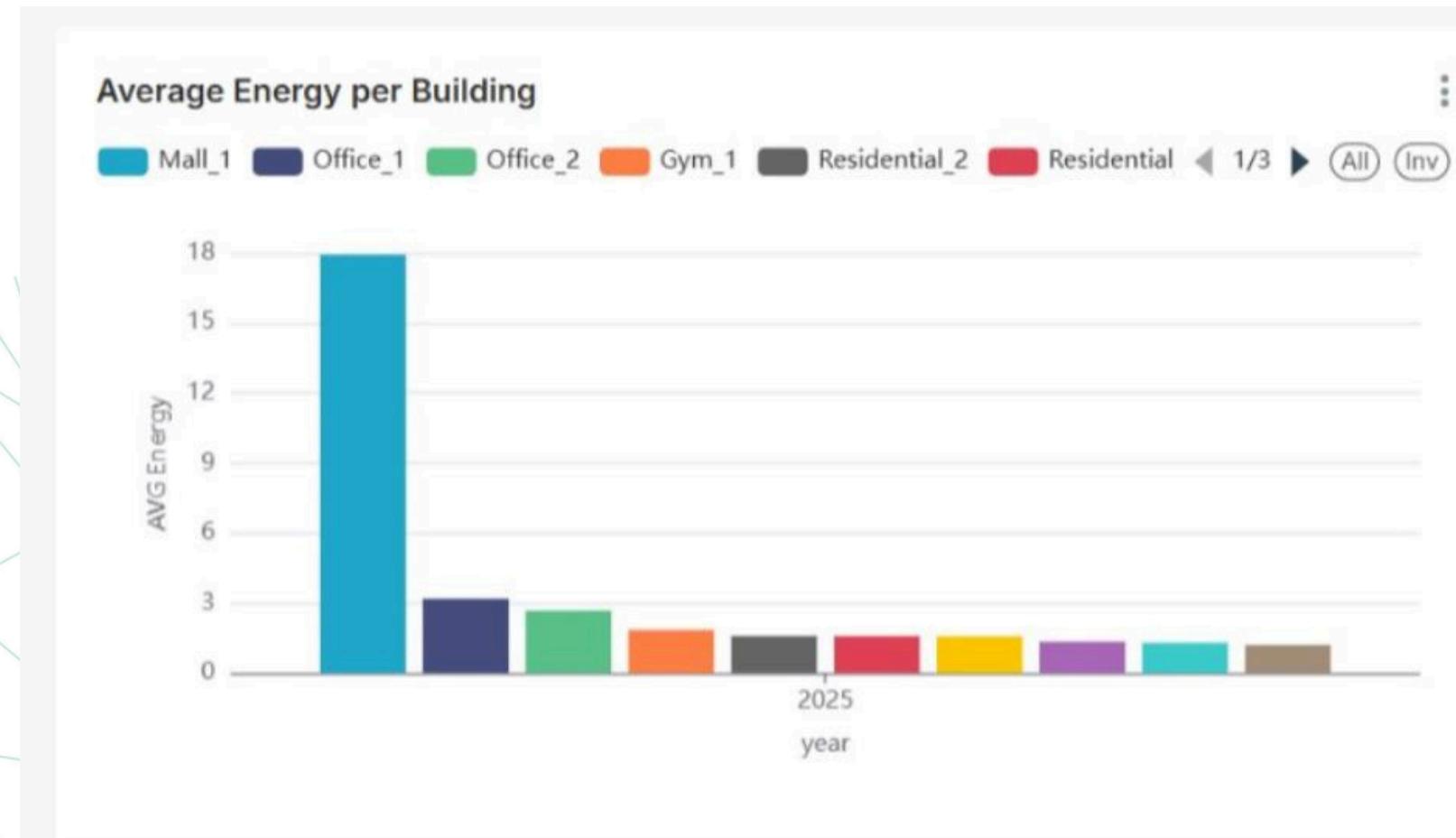
Edit dashboard

Total energy
529k
Total energy Consumption

Average Energy Per Room
6.58
Average Energy Per Room

Max Energy Peak
29.64
Max Energy Peak







Safety & Indoor Air Quality Dashboard

Draft

hanin salem

19 minutes ago

Edit dashboard

...

Average Co2 level

920.5

Average Co2 Level

Average Temperature

22.69

Average Temperature

Average Humidity

61.41

Average Humidity

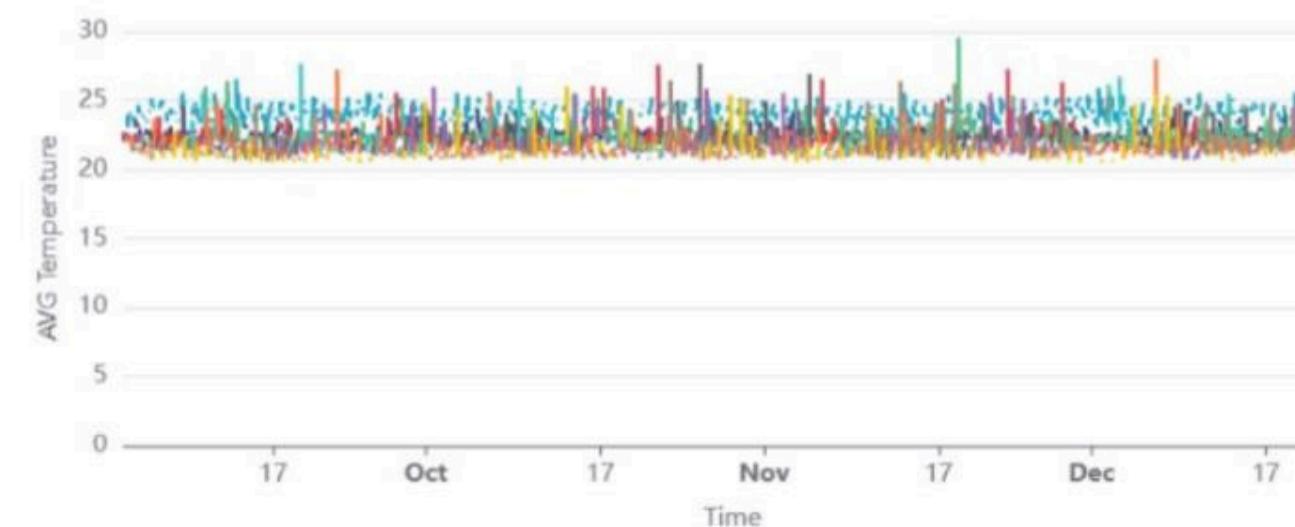
Max CO2 Today

3.14k

Max CO2 Today

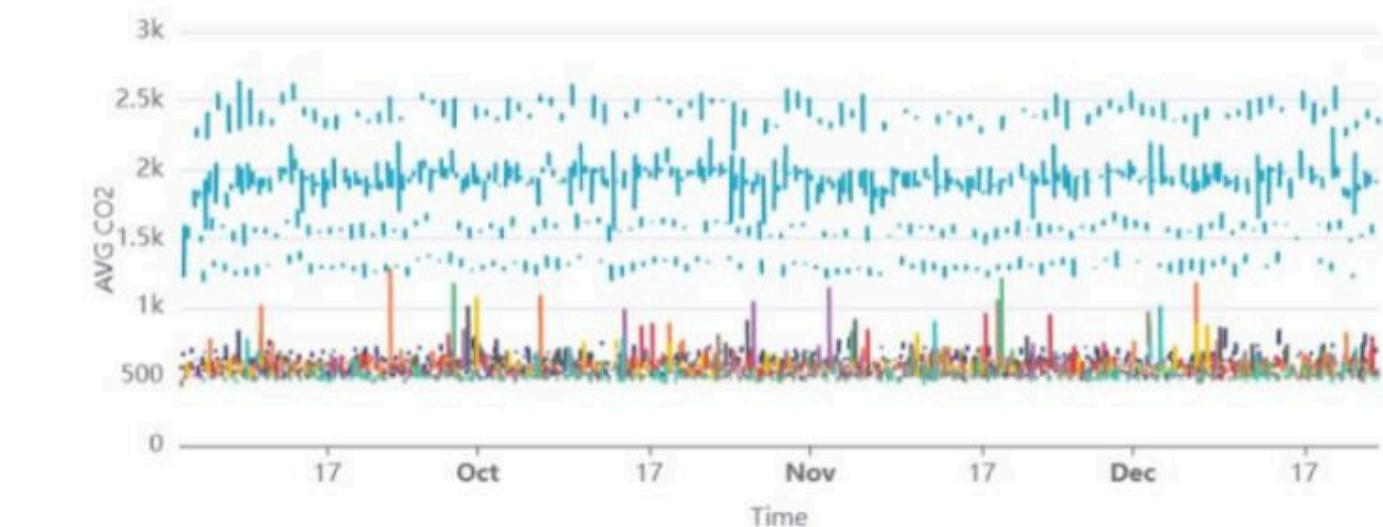
Temperature Trend

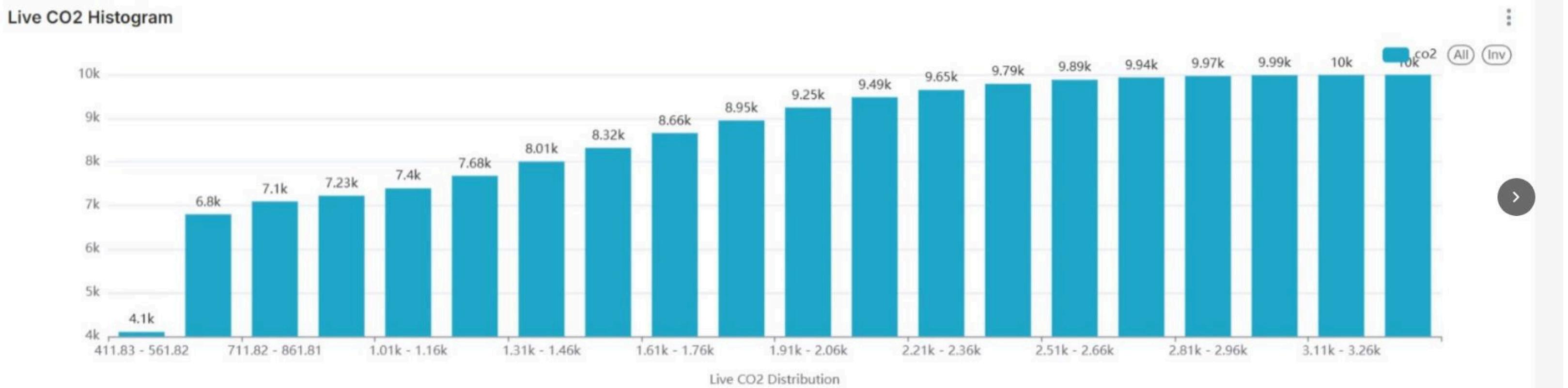
Mall_1 Office_1 Office_2 Residential_2 Residential_5 Res ▶ 1/2 (All) Inv

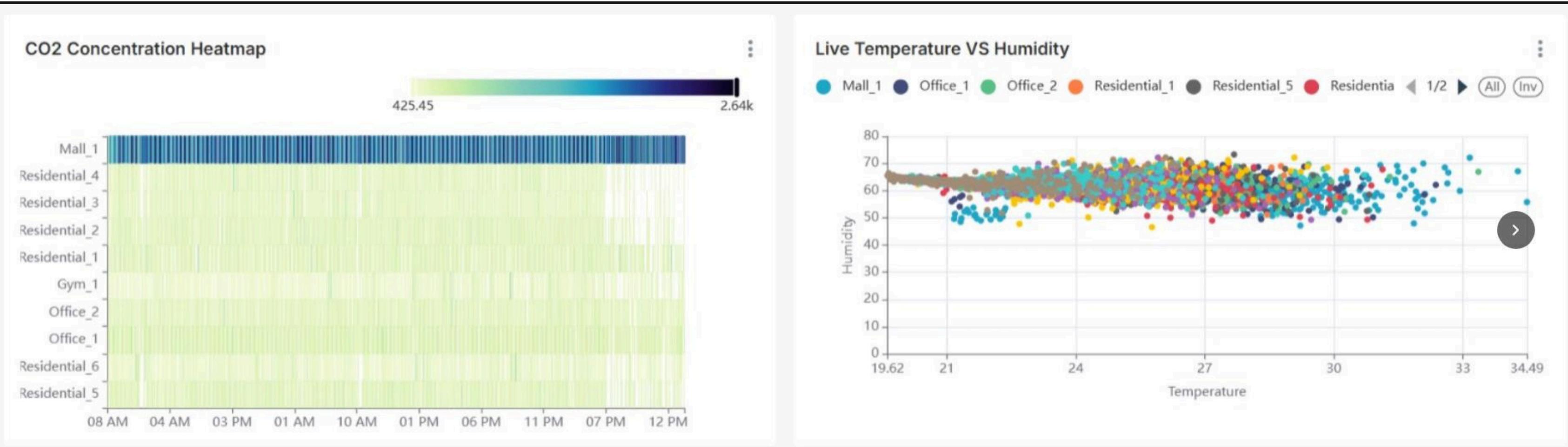


Live CO2 Trend

Mall_1 Office_1 Office_2 Residential_1 Residential_5 Res ▶ 1/2 (All) Inv







Empowering the Facility Manager: From Data to Decisions

Feature



Real-time monitoring of energy, COC, and occupancy.



Automated alerts for abnormal readings (e.g., high CO₂).



Historical trend analysis for all metrics.

Benefit



Faster, data-driven operational decisions.



Proactive air quality management and improved occupant comfort.



Root cause analysis of inefficiencies.

Outcome



Optimised energy consumption and reduced operational costs.



Healthier and more productive building environments.



Prevention of systemic waste and equipment malfunction.

Project Execution: Status, Testing, and Quality Assurance

Current Status

Pipeline is **fully functional** in beta.

Data flows correctly from Python simulation through to the Superset dashboards.



Unit Testing: Individual Spark processing functions were thoroughly tested for correctness.



Integration Testing: The end-to-end connection between Kafka, Spark, Pinot, and Superset was validated.



Functional Testing: Dashboards were verified to ensure they display accurate, real-time KPIs that match the source data.

Rigorous Testing Phases

QA Feedback Highlights

- Dashboards are responsive with low-latency updates.
- Displayed KPIs are accurate and reliable.

The Engineering Team and Collaborative Process

Team Members & Roles



Mohamed Ahmed: Team Leader / DataEngineer (Pipeline Architecture & DataGeneration)



Hanin Mohamed: Streaming Specialist (Kafka Operations)



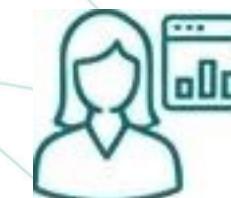
Rodina Amr: Data Processing Specialist (Spark Batch 8 Streaming)



Mohamed Ibrahim: Documentation & QA (Testing & Reporting)



Youssef Tamer: Database Specialist (Pinot Setup & Performance)



Ganna Ehab: Visualisation Specialist (Superset Dashboards)

Collaboration Methods



Process: Agile methodology with regular team syncs.



Tools: GitHub for source control, VS Code for development, and Slack/Teams for communication.



Note: The team used Tailscale VPN to create a secure network, allowing all components to connect seamlessly without public internet exposure.

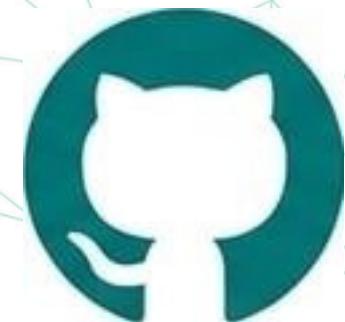


Conclusion: Transforming Raw Data into a Powerful Decision-Making Tool

This project successfully demonstrates the design and implementation of an end-to-end, real-time data pipeline for smart building management. We have built a robust, scalable system that converts a chaotic stream of sensor data into clear, actionable insights, empowering facility managers to optimise energy, improve occupant comfort, and reduce waste.

For Further Exploration:

The complete source code and technical documentation are available for review.



GitHub Repository: github.com/MAAbdelaziz22 DEPI Project



Questions?