# Building a Real-Time IoT Data Pipeline
# AWS Kinesis, Lambda, S3, Quick Sight
# [Bootcamp - AWS Data Engineering Project - 4]

**Project Overview:** Build a real-time data pipeline to monitor factory machinery using IoT sensor data. The pipeline should ingest real-time sensor streams, detect anomalies, store data in a layered data lake (Bronze, Silver, Gold), and provide operational dashboards and alerts.

The pipeline provides:

- Real-time anomaly detection – Immediate alerts for abnormal machine conditions

- Historical analysis – Trend identification for predictive maintenance

- Operational efficiency – Automated data processing and storage

- Decision support – Visual dashboards for equipment health monitoring

The solution leverages a scalable, serverless architecture to process streaming IoT data while optimizing costs.

## IAM Roles:
1. Managed Flink:
   a. AmazonKinesisFullAccess
   b. CloudWatchLogsFullAccess
   c. AmazonKinesisAnalyticsFullAccess
   d. AWSGlueConsoleFullAccess

2. **SNS:**
   a. AmazonKinesisFullAccess
   b. CloudWatchLogsFullAccess
   c. AmazonKinesisAnalyticsFullAccess

      d. AmazonSNSFullAccess

      e. AWSLambda_FullAccess

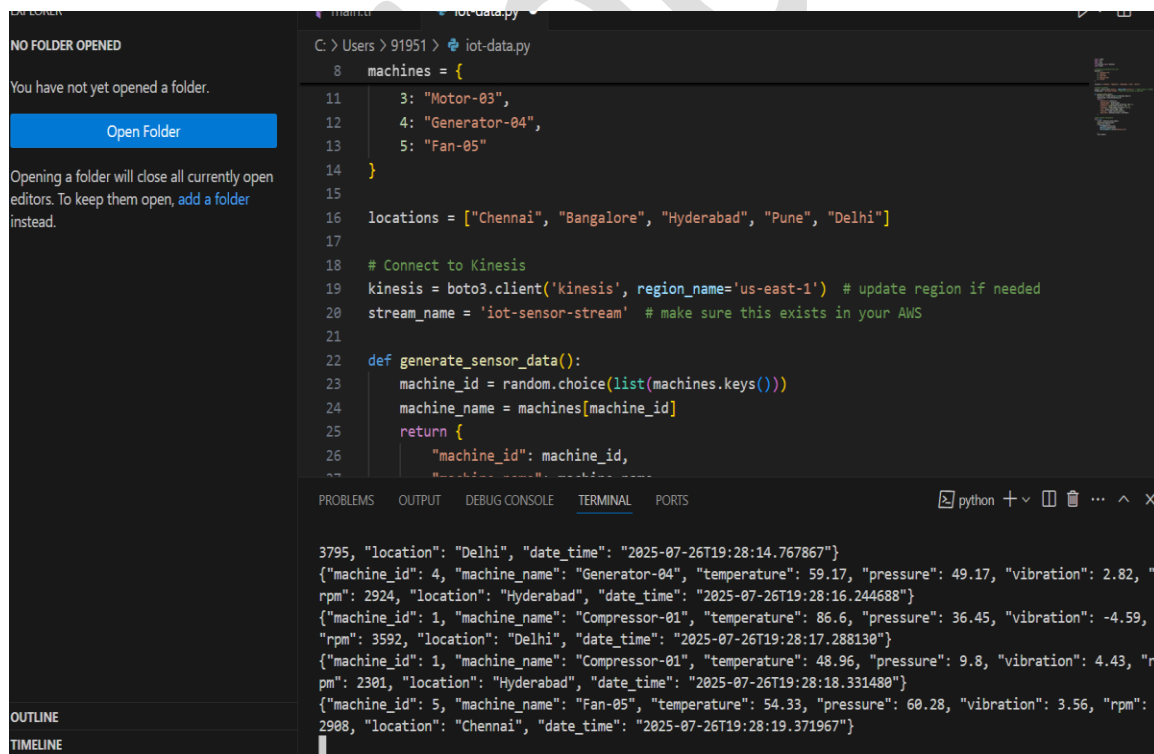      f. CloudWatchEventsFullAccess

3. **Glue**:
   a. AmazonAthenaFullAccess
   b. AmazonS3FullAccess
   c. AWSGlueServiceRole
   d. AWSQuicksightAthenaAccess
   e. CloudWatchEventsFullAccess

## Data Ingestion - Kinesis

Simulate IoT sensor data and send it to an AWS Kinesis Data Stream (`sensor-stream`) using a Python generator script. Each record includes temperature, pressure, vibration, rpm, location, and timestamp.

1. Create a S3 bucket and folders bronze, gold and silver.
2. "myp4-hema-b", objects: silver, gold and bronze.
3. Generate a stream in your local machine and stream into kinesis.

**Stream Processing - Apache Flink (Managed)**

Use Apache Flink Studio Notebook to define real-time transformation on the incoming stream.

 Key transformations:

- Calculate average metrics per minute

- Flag high temperatures

- Output streaming results to new Kinesis streams.



**Data Lake Design (Bronze - Silver - Gold)**

• Bronze: Raw data is ingested into S3 via Kinesis Firehose (prefix: `bronze/`)

• Silver: Processed stream with aggregated averages (prefix: `silver/`)

• Gold: Final Delta Lake output written using AWS Glue jobs with PySpark.

Zeppelin    Notebook ▾                                    Search        ● Configuration ▾

iot-sensor  ▷ ⋈ ▤ ✐ ⎘ ⬇ ⬆ ⟳   🔍 🗑                      ⌨ ⚙ 🔒  default ▾  Actions for project4 ▾

Took 2 hrs 3 min 33 sec. Last updated by anonymous at July 26 2025, 5:55:51 PM. (outdated)

```
%flink.ssql

CREATE TABLE rolling_avg (
machine_id INT,
machine_name STRING,
avg_temp DOUBLE,
avg_presure DOUBLE,
avg_vibration DOUBLE,
avg_rpm DOUBLE,
temperature_alert STRING,
window_end TIMESTAMP(3)
) WITH (
'connector' = 'kinesis',
'stream' = 'rolling-avg-stream',
'aws.region' = 'us-east-1',
'format' = 'json'
);

Table has been created.
```
FINISHED ▷ ⋈ ▦ ⚙

Took 1 sec. Last updated by anonymous at July 26 2025, 4:32:44 PM.

---



Amazon S3 > Buckets > myp4-hema-b > bronze/ > 2025/ > 07/ > 26/

**Amazon S3**                    **26/**                                    ⧉ Copy S3 URI
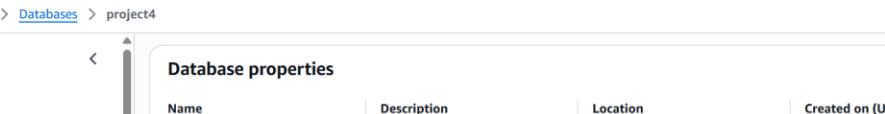
General purpose buckets          Objects    Properties
Directory buckets
Table buckets                    **Objects (1)**
Vector buckets  Preview          ⟳  ⧉ Copy S3 URI  ⧉ Copy URL  ⬇ Download  Open ↗  Delete  Actions ▾
Access Grants
Access Points (General Purpose   Create folder  ⬆ Upload
Buckets, FSx file systems)       Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For
Access Points (Directory Buckets) others to access your objects, you'll need to explicitly grant them permissions. Learn more ↗
Object Lambda Access Points
Multi-Region Access Points       🔍 Find objects by prefix                                  < 1 > ⚙
Batch Operations
IAM Access Analyzer for S3       | ☐ | Name | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
                                 | ☐ | 📁 20/ | Folder | - | - | - |

---



AWS Glue > Databases > project4

**AWS Glue**                     **Database properties**

Getting started                  **Name**        **Description**    **Location**    **Created on (UTC)**
ETL jobs                         project4        -                  -               July 26, 2025 at 19:34:56
  Visual ETL
  Notebooks
  Job run monitoring             **Tables (3)**              Last updated (UTC)  ⟳  Delete  Add tables using crawler  Add table
Data Catalog tables              View and manage all available tables.  July 26, 2025 at 20:26:43
Data connections
Workflows (orchestration)        🔍 Filter tables                                    < 1 > ⚙
Zero-ETL integrations  New
▼ Data Catalog                   | ☐ | Name ▲ | Database ▽ | Location ▽ | Classific... ▽ | Depreca... ▽ | View data | Data quality | Colur |
  Databases                      | ☐ | iot_stream | project4 | - | - | - | - | | View data qualit | View |
  Tables                         | ☐ | iot_stream_data | project4 | - | - | - | - | | View data qualit | View |
Stream schema registries         | ☐ | iot_stream1 | project4 | - | - | - | - | | View data qualit | View |
  Schemas
Connections

```
%flink.ssql
select * from iot_stream1
```
FLINK JOB  RUNNING 0%

| machine_id | machine_name | temperature | pressure | vibration | rpm | location | date_time |
|---|---|---|---|---|---|---|---|
| 1 | Compressor-01 | 78.29 | 71.61 | -0.8 | 3901 | Delhi | 2025-07-26 19:52:46.580499 |
| 1 | Compressor-01 | 71.11 | 36.19 | 4.47 | 3411 | Pune | 2025-07-26 19:52:49.694421 |
| 1 | Compressor-01 | 75.48 | 39.2 | 1.07 | 2809 | Pune | 2025-07-26 19:52:51.775273 |
| 1 | Compressor-01 | 12.48 | 20.09 | -3.03 | 3370 | Delhi | 2025-07-26 19:52:52.813347 |
| 1 | Compressor-01 | 101.58 | 71.5 | -1.86 | 3610 | Bangalore | 2025-07-26 19:53:00.207648 |

Started a minute ago.

**Zeppelin**   Notebook ▾        Search        ● Configuration ▾

## iot-sensor  ▷ ✕ 圓 ✎ ⊕ ⬇ ⬆ ⟳  Q  🗑        ⌨ ⚙ 🔒  default ▾  Actions for project4 ▾

```
CREATE TABLE iot_stream1 (
    machine_id INT,
    machine_name STRING,
    temperature DOUBLE,
    pressure DOUBLE,
    vibration DOUBLE,
    rpm INT,
    location STRING,
    date_time TIMESTAMP(3),
    temp_flag AS CASE
        WHEN temperature > 90 THEN 'High'
        ELSE 'Normal'
    END,
    WATERMARK FOR date_time AS date_time - INTERVAL '5' SECOND
)
WITH (
    'connector' = 'kinesis',
    'stream' = 'iot-sensor-stream',
    'aws.region' = 'us-east-1',
    'format' = 'json',
    'scan.stream.initpos' = 'LATEST',
    'json.timestamp-format.standard' = 'ISO-8601'
);


Table has been created.
```

Took 2 sec. Last updated by anonymous at July 26 2025, 3:51:40 PM.

**Zeppelin**   Notebook ▾        Search        ● Configuration ▾
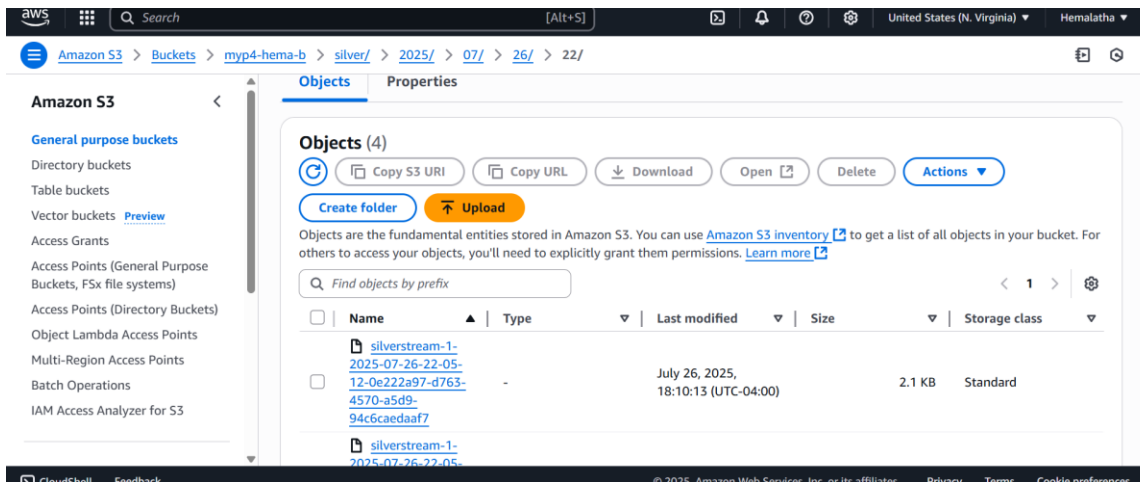
## iot-sensor  ▷ ✕ 圓 ✎ ⊕ ⬇ ⬆ ⟳  Q  🗑        ⌨ ⚙ 🔒  default ▾  Actions for project4 ▾

```
TUMBLE(date_time, INTERVAL '1' MINUTE),
machine_id,machine_name;
```
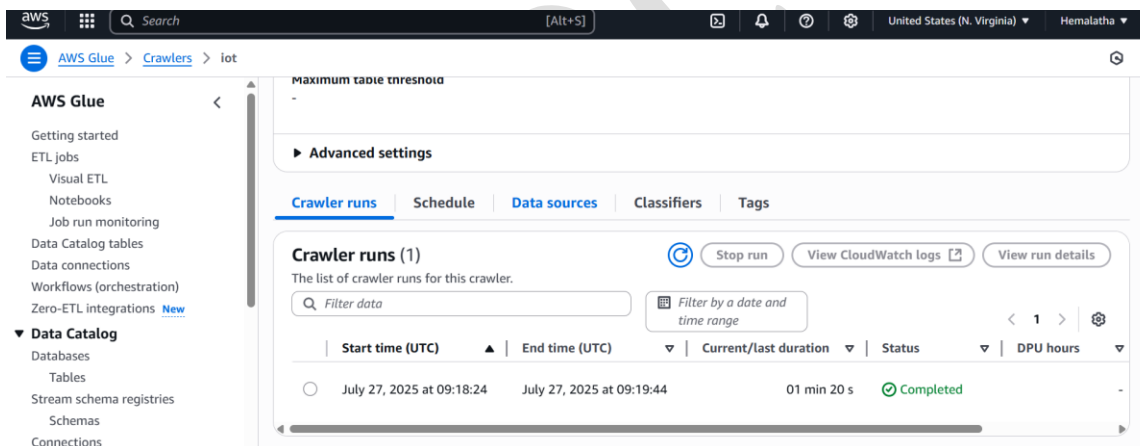Started 14 minutes ago.

```
%flink.ssql
select * from rolling_avg
```
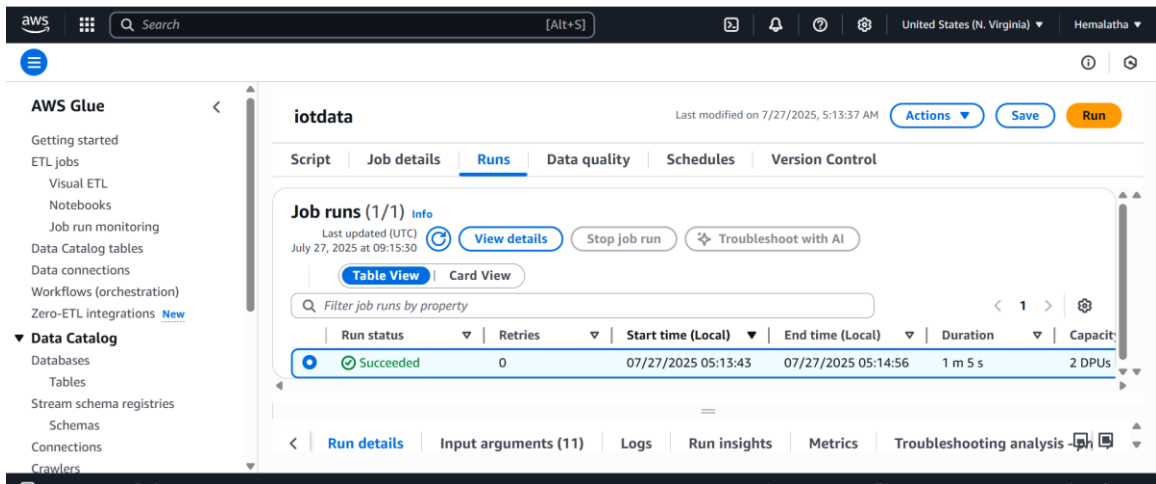FLINK JOB  RUNNING 0%

| machine_id | machine_name | avg_temp | avg_presure | avg_vibration | avg_rpm | temperature_alert | window_end |
|---|---|---|---|---|---|---|---|
| 1 | Compressor-01 | 59.00357142857143 | 47.12928571428572 | 0.5428571428571428 | 2952.0 | Normal | 2025-07-26 22:02:00.000 |
| 1 | Compressor-01 | 67.20727272727272 | 43.838181818181816 | 0.4345454545454546 | 3012.0 | Normal | 2025-07-26 22:03:00.000 |
| 2 | Pump-02 | 65.14500000000001 | 62.896 | 0.4860000000000001 | 2754.0 | Normal | 2025-07-26 22:02:00.000 |
| 2 | Pump-02 | 57.138 | 49.111 | -1.5510000000000002 | 3004.0 | Normal | 2025-07-26 22:03:00.000 |
| 3 | Motor-03 | 56.99857142857143 | 32.79857142857143 | -1.612857142857143 | 2884.0 | Normal | 2025-07-26 |

## ETL with AWS Glue

Configure Glue jobs to read from Silver S3 (Parquet format), perform transformations, and write Delta format output to gold layer in S3.

## Alerts via Lambda + SNS

Deploy a Lambda function triggered by Kinesis Data Stream to monitor `temperature > 90°C` and send real-time alerts to email via SNS.

**Visualization with Athena + QuickSight**

Use Glue Crawlers to catalog Delta Lake outputs. Connect Athena tables to QuickSight for building visual dashboards (e.g., average temp, RPM, anomaly trends).

**Summary**

The end-to-end pipeline enables real-time monitoring, alerting, storage, and insight generation for smart factory operations. It demonstrates scalable design using cloud-native tools and Delta Lake architecture.