# IoT Sensor Data Processor

## Short Description

This project demonstrates how to build a complete ETL (Extract, Transform, Load) pipeline using AWS Glue and Amazon S3 to process IoT sensor data. The pipeline reads raw sensor logs (temperature, humidity, timestamps), cleans invalid records, performs transformations, and calculates hourly average temperatures. The cleaned and aggregated data is stored in Parquet format for efficient querying and future analytics.

```
TOOLS & TECHNOLOGIES USED:
----------------------------------------------------------------------
| Tool/Technology   | Description                                              |
|-------------------|----------------------------------------------------------|
| AWS S3            | Cloud storage to store raw, processed, and aggregated sensor data. |
| AWS Glue Studio   | Used to create custom ETL jobs with PySpark script editor.    |
| AWS Glue Catalog  | Metadata store that keeps schema info about the raw data.     |
| PySpark           | Distributed data processing framework used to transform the data. |
| Apache Parquet    | Columnar data format for efficient storage and querying.      |
| GitHub            | Version control and collaboration platform for storing project scripts. |

----------------------------------------------------------------------
S3 BUCKET STRUCTURE:
----------------------------------------------------------------------
s3://iot-sensor-data-satyam/

| Folder Name | Description                                              |
|-------------|----------------------------------------------------------|
| raw/        | Contains raw sensor data in CSV or JSON format.          |
| processed/  | Stores cleaned and validated data in Parquet format.     |
| aggregated/ | Contains hourly average temperature data in Parquet format. |

----------------------------------------------------------------------
DATA SAMPLE (RAW CSV):
----------------------------------------------------------------------
| sensor_id | timestamp           | temperature | humidity |
|-----------|---------------------|-------------|----------|
| sensor-1  | 2025-06-01 10:00:00 | 25.3        | 60       |
| sensor-2  | 2025-06-01 10:05:00 | -100        | 65       |
| sensor-3  | 2025-06-01 10:10:00 | 30.1        | 70       |
| sensor-4  | 2025-06-01 10:15:00 | 45.0        | 75       |
| sensor-5  | 2025-06-01 10:20:00 | 151         | 66       |
```

## ETL Workflow

### Step 1: Raw Data Ingestion

- Upload sensor logs to: s3://iot-sensor-data-satyam/raw/

### Step 2: Cleaning & Transformation (ETL Job 1)

- Read data from AWS Glue Data Catalog
- Filter out invalid temperature readings (below -50 or above 150)
- Convert timestamp to datetime format
- Cast humidity from long to double
- Save cleaned data to s3://iot-sensor-data-satyam/processed/ as Parquet

**Step 3: Aggregation (ETL Job 2)**

- Read cleaned data from processed folder
- Extract hour from each timestamp
- Compute average temperature for each hour
- Save output to s3://iot-sensor-data-satyam/aggregated/ in Parquet

```
OUTPUT SAMPLE (CLEANED):
------------------------------------------------------------
| sensor_id | timestamp           | temperature | humidity |
|-----------|---------------------|-------------|----------|
| sensor-1  | 2025-06-01 10:00:00 | 25.3        | 60       |
| sensor-3  | 2025-06-01 10:10:00 | 30.1        | 70       |
| sensor-4  | 2025-06-01 10:15:00 | 45.0        | 75       |


------------------------------------------------------------

REPOSITORY STRUCTURE:
------------------------------------------------------------
IOT-SENSOR/
├── iot_sensor_etl_job.py        # Glue script for cleaning raw data
├── iot_sensor_aggregate_job.py  # Glue script for hourly aggregation
├── read_parquet.py              # Script to read Parquet files locally
└── README.md                    # Documentation
```

## Project Benefits

- **Performance**: Fast reads and queries via Parquet
- **Cost-Efficient**: Only valid, transformed data is saved
- **Flexible**: Fully custom ETL logic using PySpark
- **Scalable**: Easily expandable to more sensors or metrics
- **Analytics Ready**: Compatible with Athena, QuickSight, or Python tools

## Use Cases

- Real-time environment monitoring
- Industrial IoT sensor data analysis
- Smart home automation
- Weather pattern tracking
- Anomaly detection in sensor readings

## Project Repository

Explore the full project code and scripts here:

GitHub Repository: https://github.com/Satyam25613/IOT-SENSOR

aws    Q Search    [Alt+S]    United States (N. Virginia) ▾    SATYAM ▾

S3

**AWS Glue**

Getting started
ETL jobs
  Visual ETL
  Notebooks
  Job run monitoring
Data Catalog tables
Data connections
Workflows (orchestration)
Zero-ETL integrations  New
▼ Data Catalog
Databases
  Tables
Stream schema registries
  Schemas
Connections
Crawlers
  Classifiers
Catalog settings
▶ Data Integration and ETL
▶ Legacy pages

What's New ↗
Documentation ↗

**iot_sensor_etl_job**    Last modified on 6/9/2025, 7:25:22 PM    Actions ▾    Save    Run

Script    **Job details**    Runs    Data quality    Schedules    Version Control

**IAM Role**
Role assumed by the job with permission to access your data stores. Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job.

AWSGlueServiceRole-SensorData    ▾    ↻

**Type**
The type of ETL job. This is set automatically based on the types of data sources you have selected.

Spark    ▾

**Glue version**    Info

Glue 5.0 - Supports spark 3.5, Scala 2, Python 3    ▾

**Language**

Python 3    ▾

**Worker type**
Set the type of predefined worker that is allowed when a job runs.

G 1X
(4vCPU and 16GB RAM)    ▾

**Automatically scale the number of workers**
☐ AWS Glue will optimize costs and resource usage by dynamically scaling the number of workers up and down throughout the job run. Requires Glue 3.0 or later.

© 2025, Amazon Web Services, Inc. or its affiliates.    Privacy    Terms    Cookie preferences

CloudShell    Feedback

---

aws    Q Search    [Alt+S]    United States (N. Virginia) ▾    SATYAM ▾

S3

**AWS Glue**

Getting started
ETL jobs
  Visual ETL
  Notebooks
  Job run monitoring
Data Catalog tables
Data connections
Workflows (orchestration)
Zero-ETL integrations  New
▼ Data Catalog
Databases
  Tables
Stream schema registries
  Schemas
Connections
Crawlers
  Classifiers
Catalog settings
▶ Data Integration and ETL
▶ Legacy pages

What's New ↗
Documentation ↗

**iot_sensor_etl_job**    Last modified on 6/9/2025, 7:25:22 PM    Actions ▾    Save    Run

**Script**    Job details    Runs    Data quality    Schedules    Version Control

**Script**  Info

```python
1   import sys
2   from awsglue.transforms import *
3   from awsglue.utils import getResolvedOptions
4   from pyspark.context import SparkContext
5   from pyspark.sql.functions import col, to_timestamp, hour, dayofmonth, avg
6   from awsglue.context import GlueContext
7   from awsglue.job import Job
8   from awsglue.dynamicframe import DynamicFrame
9
10  args = getResolvedOptions(sys.argv, ['JOB_NAME'])
11
12  sc = SparkContext()
13  glueContext = GlueContext(sc)
14  spark = glueContext.spark_session
15  job = Job(glueContext)
16  job.init(args['JOB_NAME'], args)
17
18  # 1. Read cleaned data from S3 (processed folder)
19  input_path = "s3://iot-sensor-data-satyam/processed/"
```

Python    Ln 1, Col 1    ⊗ Errors: 0    ⚠ Warnings: 0

© 2025, Amazon Web Services, Inc. or its affiliates.    Privacy    Terms    Cookie preferences

CloudShell    Feedback

---

## Screenshot 2 — AWS Glue Studio Runs Editor

Runs - Editor - AWS Glue Studio

us-east-1.console.aws.amazon.com/gluestudio/home?region=us-east-1#/editor/job/%09iot_sensor_etl_job/runs

aws | Search [Alt+S] | United States (N. Virginia) ▼ | SATYAM ▼

S3

### AWS Glue
- Getting started
- **ETL jobs**
  - Visual ETL
  - Notebooks
  - Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations New
- **Data Catalog**
  - Databases
  - Tables
  - Stream schema registries
    - Schemas
  - Connections
  - Crawlers
    - Classifiers
  - Catalog settings
- Data Integration and ETL
- Legacy pages

What's New ↗
Documentation ↗

# iot_sensor_etl_job

Last modified on 6/9/2025, 7:25:22 PM | Actions ▼ | Save | Run

Script | Job details | **Runs** | Data quality | Schedules | Version Control

## Job runs (1/3) Info
Last updated (UTC) June 10, 2025 at 18:06:31

View details | Stop job run | ✨ Troubleshoot with AI | **Table View** | Card View

Filter job runs by property

< 1 >

| | Run status ▼ | Retries ▼ | Start time (Local) ▼ | End time (Local) ▼ | Duration ▼ | Capacity (D... ▼ | Worker type ▼ | Glue version |
|---|---|---|---|---|---|---|---|---|
| ● | ✓ Succeeded | 0 | 06/09/2025 19:27:23 | 06/09/2025 19:29:15 | 1 m 37 s | 10 DPUs | G.1X | 5.0 |
| ○ | ✓ Succeeded | 0 | 06/09/2025 19:04:20 | 06/09/2025 19:06:06 | 1 m 34 s | 10 DPUs | G.1X | 5.0 |
| ○ | ✗ Failed | 0 | 06/09/2025 18:56:36 | 06/09/2025 18:58:24 | 1 m 38 s | 10 DPUs | G.1X | 5.0 |

Run details | Input arguments (9) | Logs | Run insights | Metrics | Troubleshooting analysis - preview | Spark UI

| Job name | Start time (Local) | Glue version | Last modified on (Local) |
|---|---|---|---|
| iot_sensor_etl_job | 06/09/2025 19:27:23 | 5.0 | 06/09/2025 19:29:15 |
| Id | End time (Local) | Worker type | Log group name |
| jr_9751ab070056f2020b22c0aa76987da820 | 06/09/2025 19:29:15 | G.1X | /aws-glue/jobs |

CloudShell  Feedback  © 2025, Amazon Web Services, Inc. or its affiliates.  Privacy  Terms  Cookie preferences

26°C Haze | Search | ENG IN | 23:36 10-06-2025