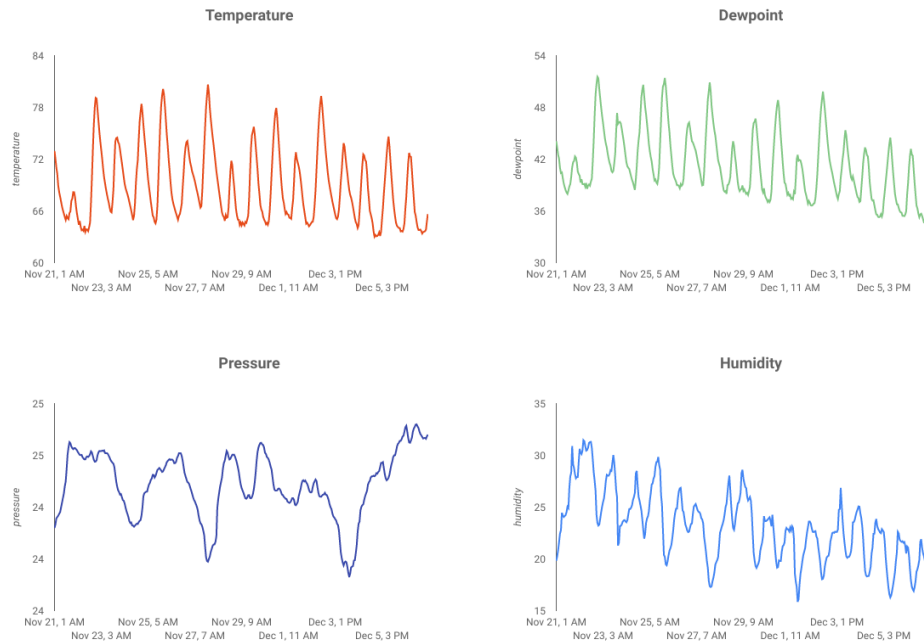# Serverless IOT Data Processing

**PROBLEM DEFINITION:**

While multi-tier applications consisting of web, application server and database are foundational to web development and are the starting point for many websites, success will often bring challenges around scalability, integration and agility. For example, how can data be handled in real-time and how can it be distributed to multiple key business systems? These issues coupled with the demands of internet-scale applications drove the need for a distributed messaging systemand gave rise to an architectural pattern of using data pipelines to achieve resilient, real-time systems. As a result, understanding how to publish real-time data to build a data pipeline are crucial skills for developer and architect alike.

In this code lab, we are going to build a weather data pipeline that starts with an Internet of Things (IoT) device, utilizes a message queue to receive and deliver data. Leverages a serverless function to move the data to a data warehouse and then creates a dashboard that displays the information. A Raspberry Pi with a weather sensor will be used for the IoT device and several components of the Google Cloud Platform will be from the data pipeline. Building out the Raspberry Pi, while beneficial, is an optional portion of this code lab.

Temperature

Dewpoint

Pressure

Humidity



Google Cloud Platform

Cloud Pub/Sub → Cloud Functions → BigQuery → Data Studio

**DESIGN THINKING**

→Create Pub/sub

   Deploy a Google Cloud Function.

→Leverage Google Big Query.

→Dashboard using Google Data Studio.

→If we build out the IoT sensor, we will also learn how to utilize the Google Cloud SDK and how to secure remote access calls to the Google Cloud Platform.

If we want to build the IoT sensor portion of this code lab instead of leveraging sample data and a script, we will also need the following

→Raspberry Pi Zero W with power supply, SD memory card and case.

→USB card reader.

→USB hub (to allow for connecting a keyboard and mouse into the sole USB port on the Raspberry Pi)

→Female-to-female breadboard wires.

→GPIO Hammer Headers.

→BME280 sensor.

**CREATE A BIG QUERY TABLE:**

Big Query is a serverless, highly scalable, low-cost enterprise data warehouse and will be an ideal option to store data being streamed from IoT devices while also allowing an analytics dashboard to query the information. Let's create a table that will hold all the IoT weather data. Select Big Query from the Cloud console. **This will open Big Query in a new window** (don't close the original window as you'll need to access it again).

### Create Table

**Source Data**   ○ Create from source   ● Create empty table

**Destination Table**

| Table name | weath | . | weatherDataTable | ? |

| Table type | Native table ⇕ | ? |

**Schema**

| Name | Type | Mode | |
|------|------|------|---|
| sensorID | STRING ⇕ | NULLABLE ⇕ | ✕ |
| timecollected | TIMESTAMP ⇕ | NULLABLE ⇕ | ✕ |
| zipcode | INTEGER ⇕ | NULLABLE ⇕ | ✕ |
| latitude | FLOAT ⇕ | NULLABLE ⇕ | ✕ |
| longitude | FLOAT ⇕ | NULLABLE ⇕ | ✕ |
| temperature | FLOAT ⇕ | NULLABLE ⇕ | ✕ |
| humidity | FLOAT ⇕ | NULLABLE ⇕ | ✕ |
| dewpoint | FLOAT ⇕ | NULLABLE ⇕ | ✕ |
| pressure | FLOAT ⇕ | NULLABLE ⇕ | ✕ |

Add Field                                    Edit as Text

**Options**

| Partitioning | None ⇕ |
| Encryption Type | Default ⇕ |

Create Table

**CREATE A PUB/SUB TOPIC:**

Cloud pub/sub is a simple, reliable, scalable foundation for stream analytics and event-driven computing systems. As a result, it is perfect for handling incoming IOT messages and then allowing downstream systems to process them.

If you are still in the window for Big Query, switch back to the Cloud Console. If you closed the Cloud Console, go to https://console.cloud.google.com

From the Cloud Console, select Pub/Sub and then Topics.

If you see an Enable API prompt, click the Enable API button.

After the key upload is complete, it should appear in the Cloud Storage browser.



**CREATE A CLOUD FUNCTION:**

Cloud computing has made possible fully serverless models of computing where logic can be spun up on-demand in response to events originating from anywhere. For this lab, a Cloud Function will start each time a message is published to the weather topic, will read the message and then store it in Big Query.



**START THE DATA PIPELINE:**

Might need to enable compute API.

**CHECK THAT DATA IS FLOWING:**

**CLOUD FUNCTION:**

Ensure that the Cloud Function is being triggered by Pub/Sub

gcloud beta functions logs read function-weatherPubSubToBQ

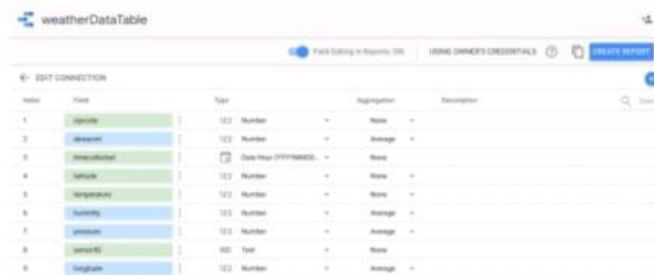**CREATE A DATA STUDIO DASHBOARD:**

Google data studio turns our data into informative dashboards and reports that are easy to read, easy to share, and fully customizable.



Onboarding new data or building new analytics pipelines in traditional analytics architectures typically requires extensive coordination across business, data engineering, and data science and analytics teams to first negotiate requirements, schema, infrastructure capacity needs, and workload management.

For many use cases today however, business users, data scientists, and analysts are demanding easy, frictionless, self-service options to build end-to-end data pipelines because it's hard and inefficient to predefine constantly changing schemas and spend time negotiating capacity slots on shared infrastructure. The exploratory nature of machine learning (ML) and many analytics tasks means you need to rapidly ingest new

datasets and clean, normalize, and feature engineer them without worrying about operational overhead when you must think about the infrastructure that runs data pipelines.

**Serverless data analytics architecture:**

To compose the layers described in our logical architecture, we introduced a reference architecture that uses serverless and managed services.
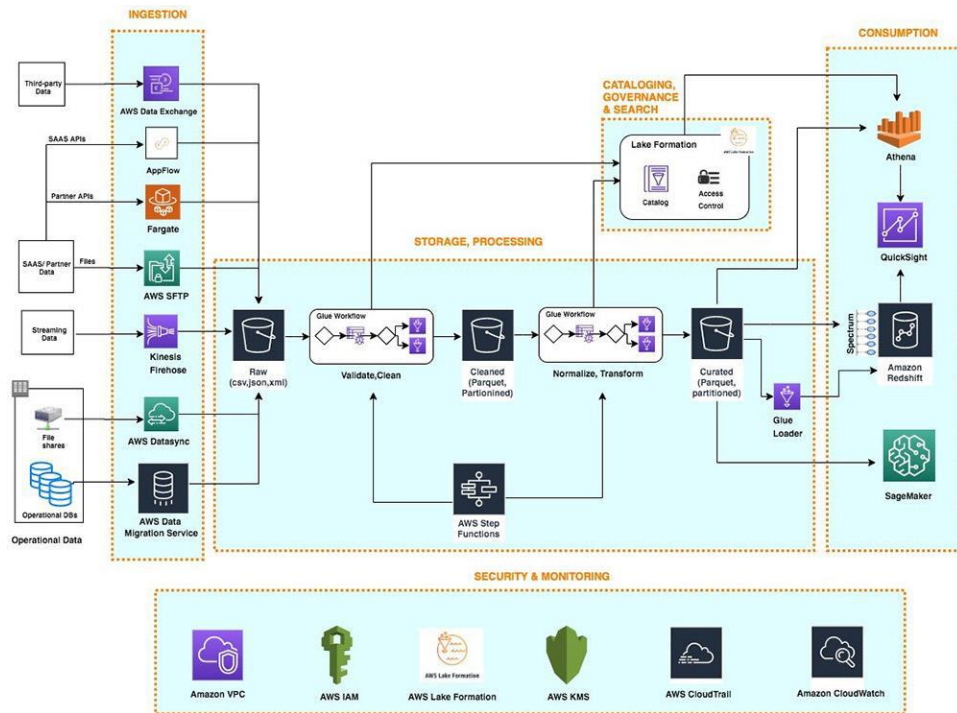
**It provides the following key benefit:**

•

**Easy configuration-driven use.**

•

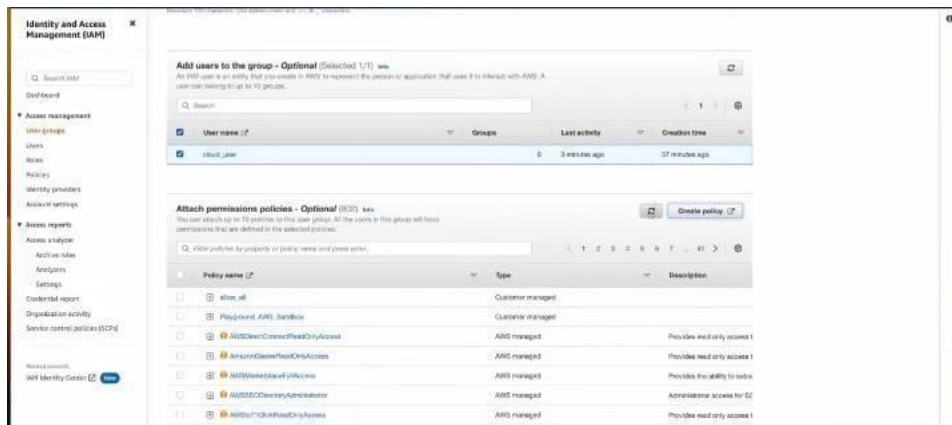**Freedom from infrastructure management.**
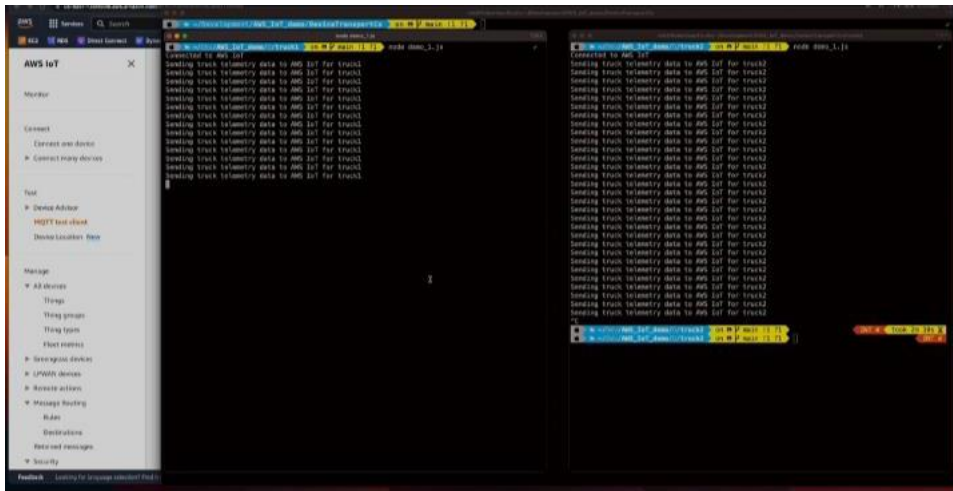
•

**Pay-per-use pricing model.**

## INGESTION LAYER:

The ingestion layers in our serverless architecture are composed of a set of purpose-built services to enable data ingestion from a variety of sources.
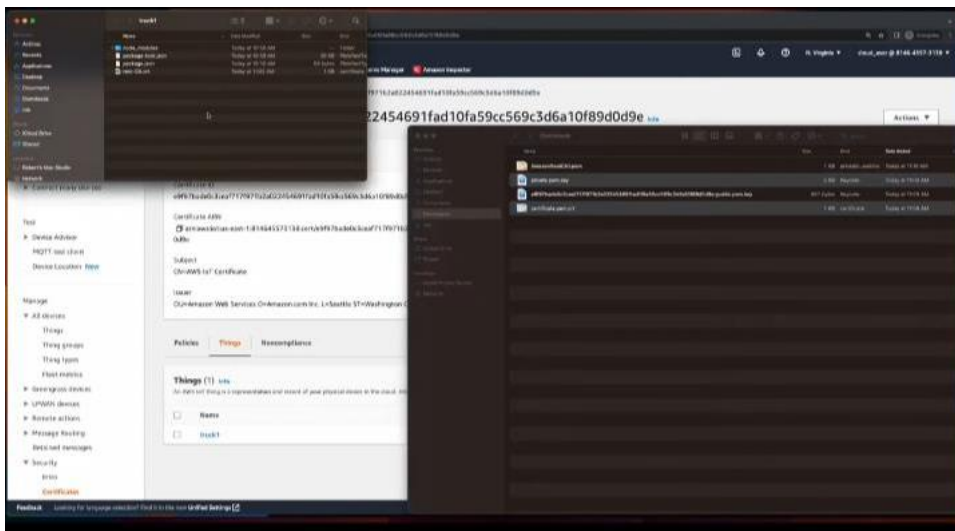


## STORAGE LAYER:

Data of any structure and any format can be stored as S3 objects without needing to predefine any schema. This enables services in the ingestion layer to quickly land a variety of source data into the data lake in its original source format.

**Cataloging and search layer:**

A data lake typically hosts a large number of datasets, and many of these datasets have evolving schema and new data partitions. A central Data Catalog that manages meta data for all the datasets in the data lake is crucial to enabling self-service discovery of data in the data lake. Additionally, separating metadata from data into a central schema enables schema-on-read for the processing and consumption layer components.
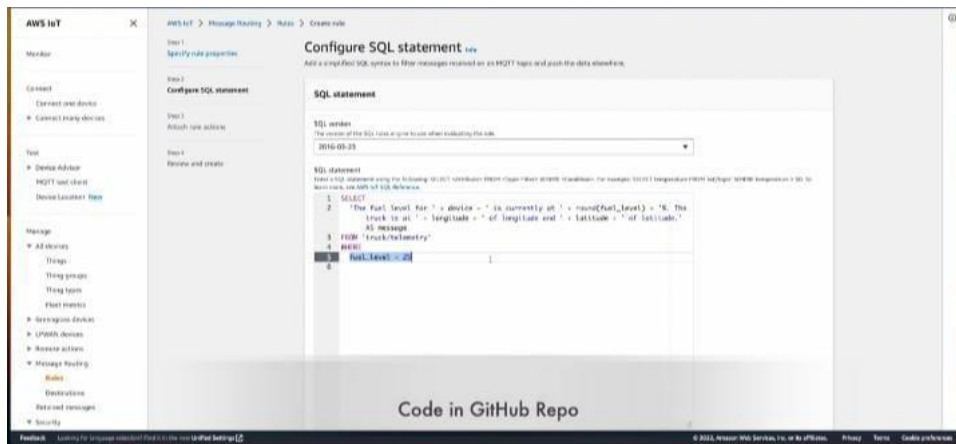


**Processing layer:**

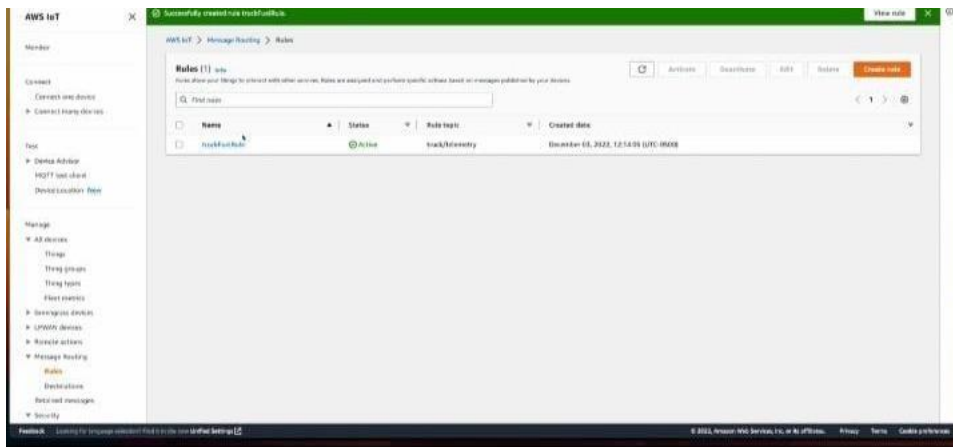The processing layer in our architecture is composed of two types of components:

•Components used to create multi-step data processing pipelines

•Components to orchestrate data processing pipelines on schedule or in response to event triggers (such as ingestion of new data into the landing zone)
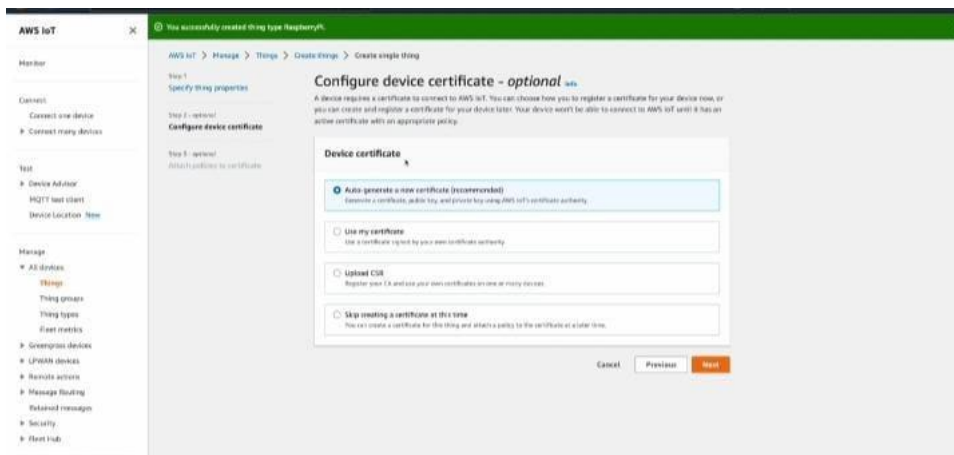
Code in GitHub Repo

**Consumption layer:**

The consumption layer in our architecture is composed using fully managed, purpose-built, analytics services that enable interactive SQL, BI dashboarding, batch processing, and ML.
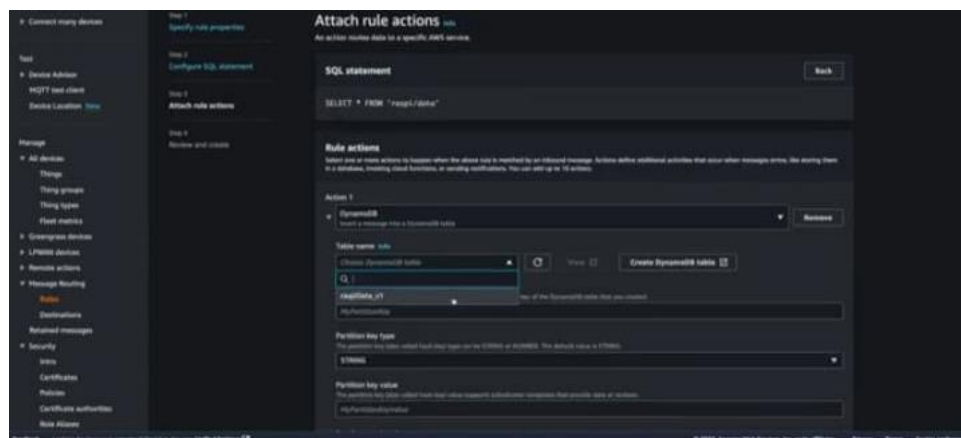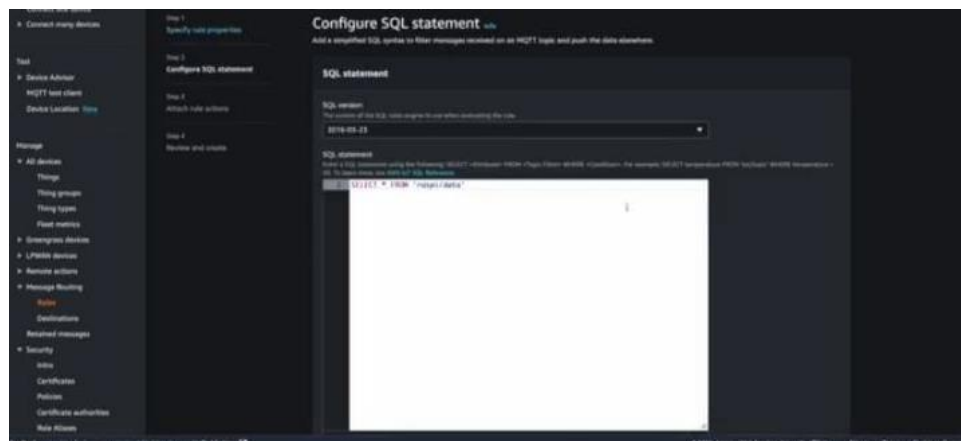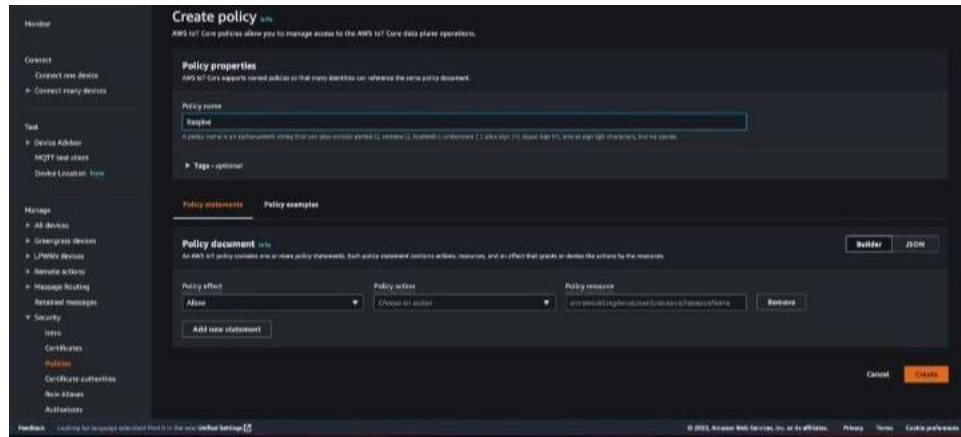


**Security and governance layer:**

Components across all layers of our architecture protect data, identities, and processing resources by natively using the following capabilities provided by the security and governance layer.

**Additional considerations:**

In this post, we talked about ingesting data from diverse sources and storing it as S3 objects in the data lake and then using it to process ingested datasets until they're in a consumable state.







**Conclusion:**

Serverless and managed services, you can build a modern, low-cost data lake centric analytics architecture in days. A decoupled, component-driven architecture allows you to start small and quickly

add new purpose-built components to one of six architecture layers to address new requirements and data sources.