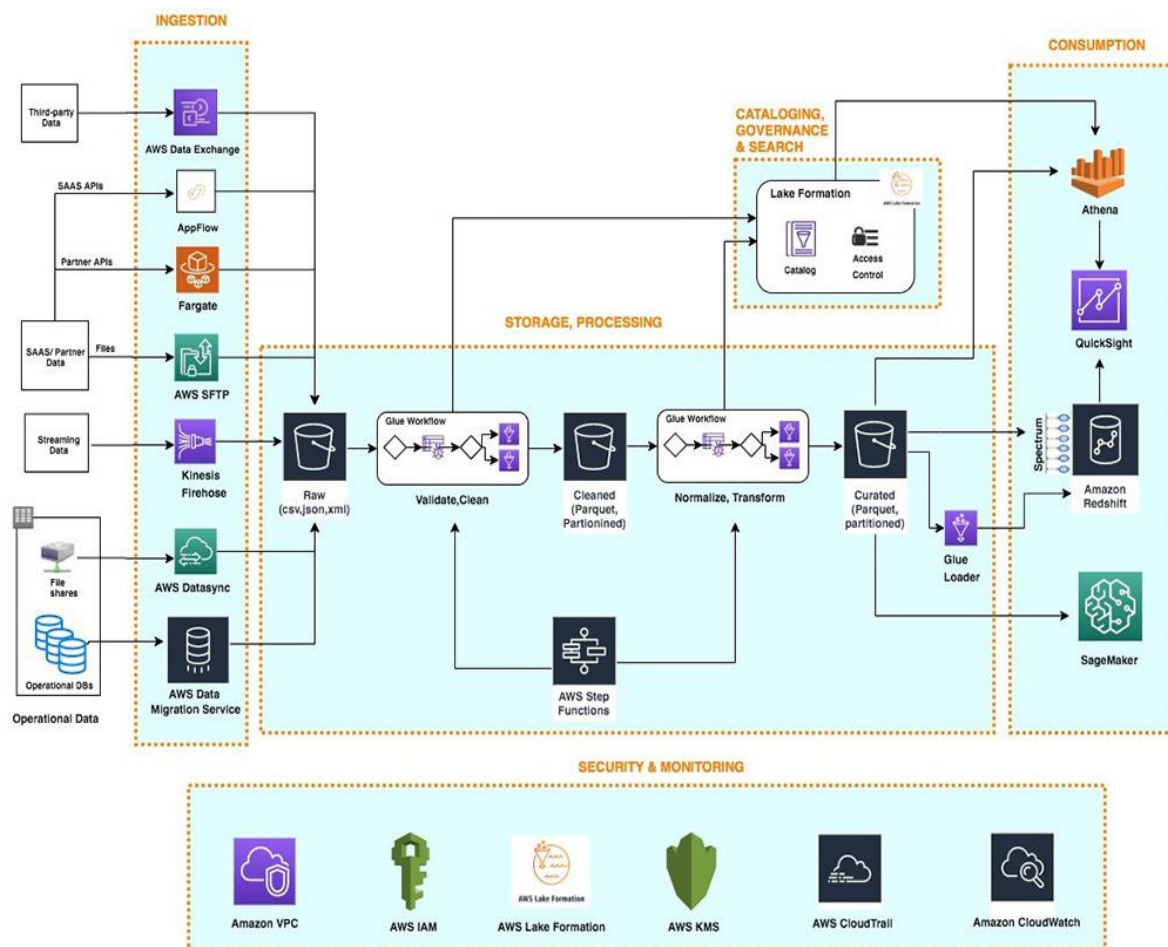# CREATE A SERVERLESS IOT DATA PROCESSING



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 lake centric analytics architecture:



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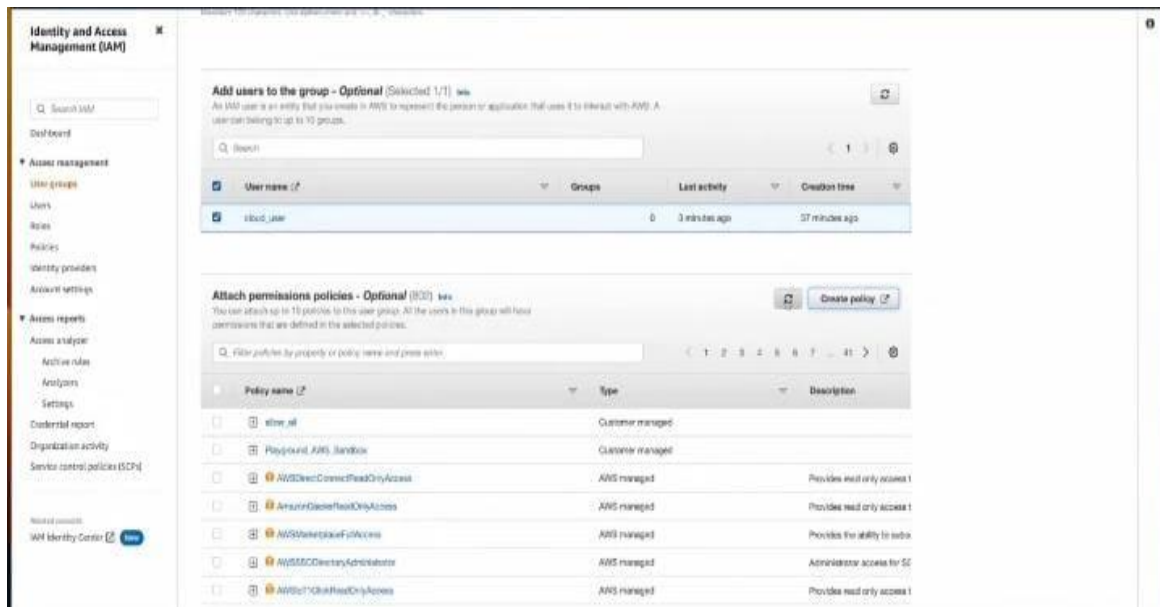
It provides the following key benefit:

- Easy configuration-driven use.

- Freedom from infrastructure management.
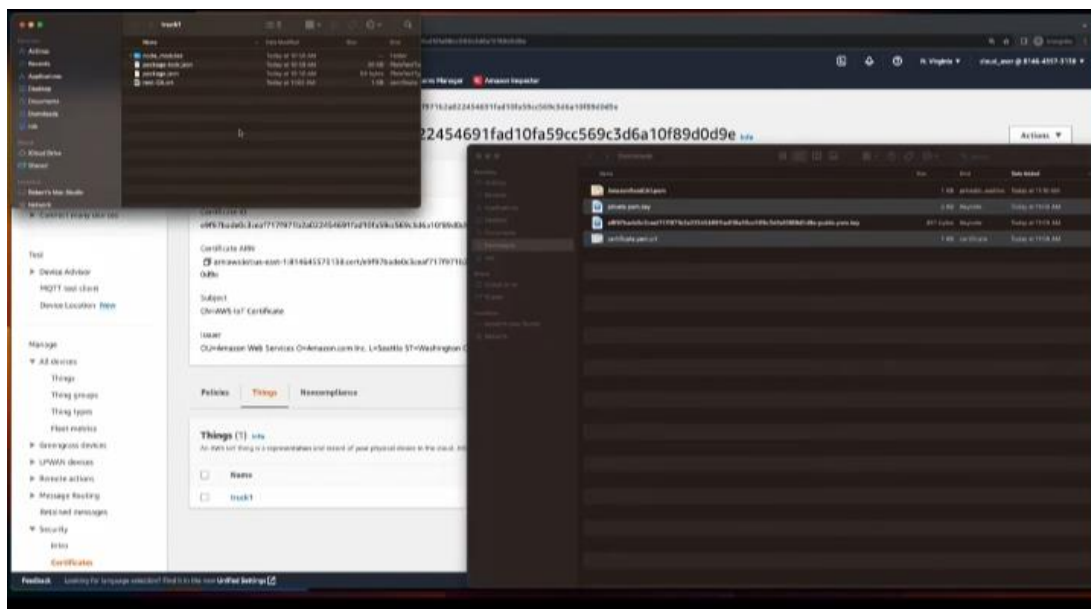- Pay-per-use pricing model.

# 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.
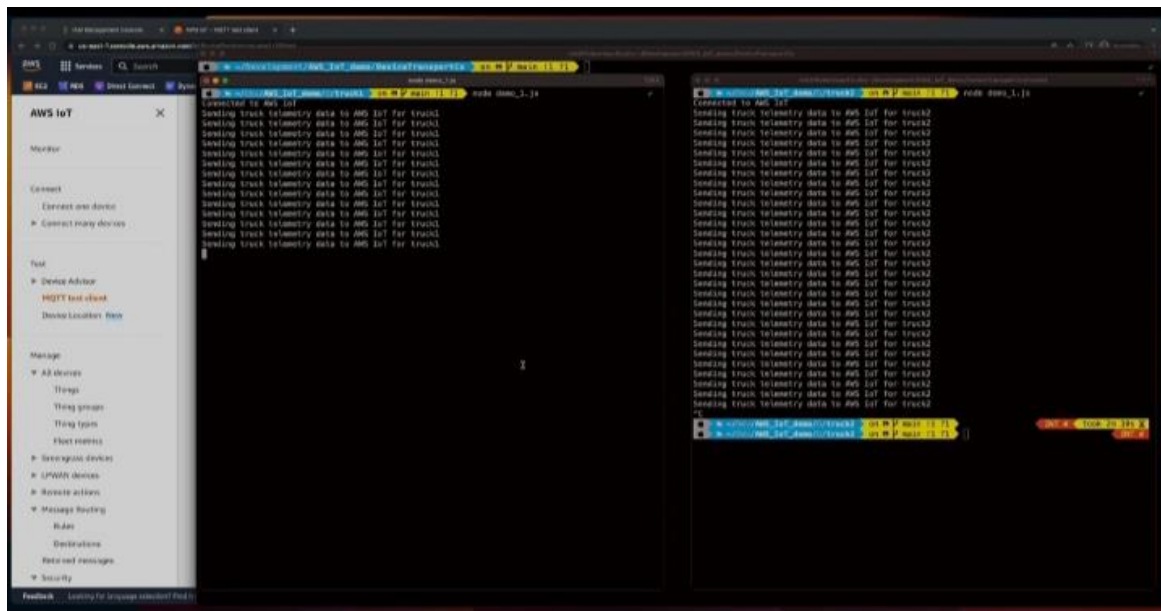


# STORAGE LAYER:

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.
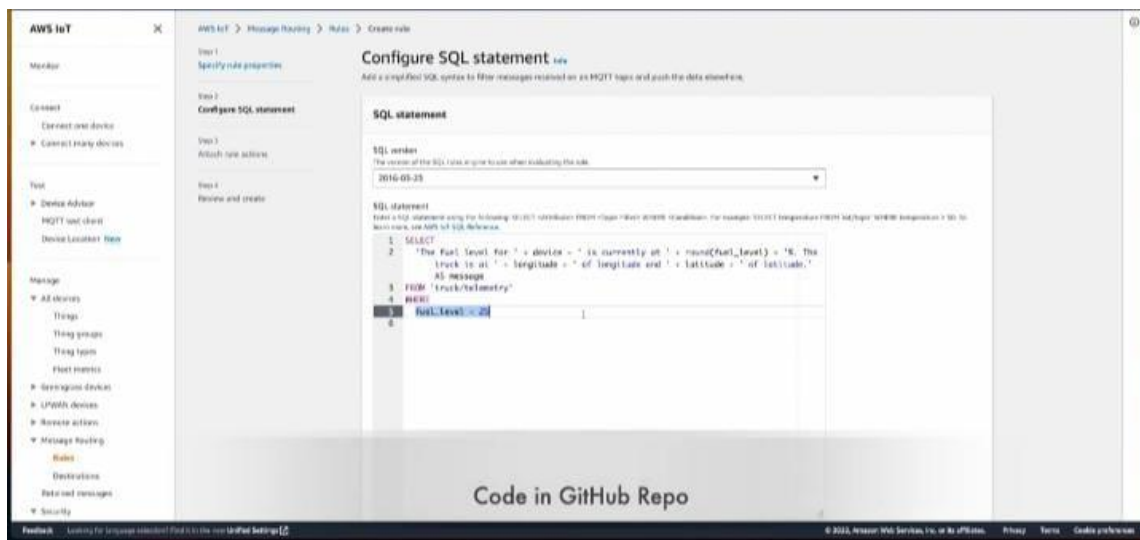
# Cataloging and search layer:

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.
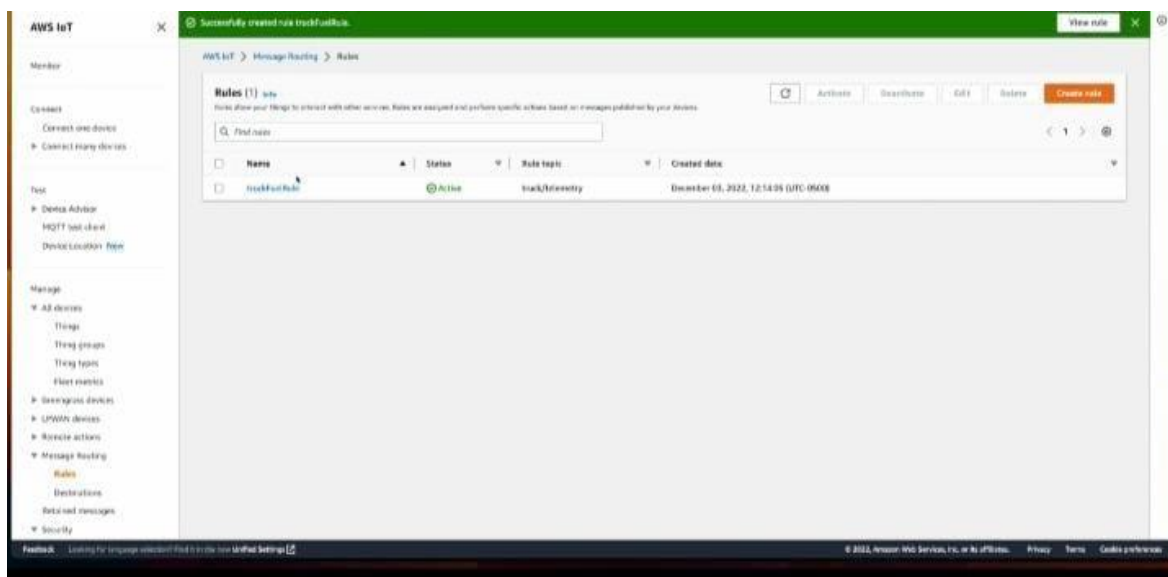
# Processing layer:

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)
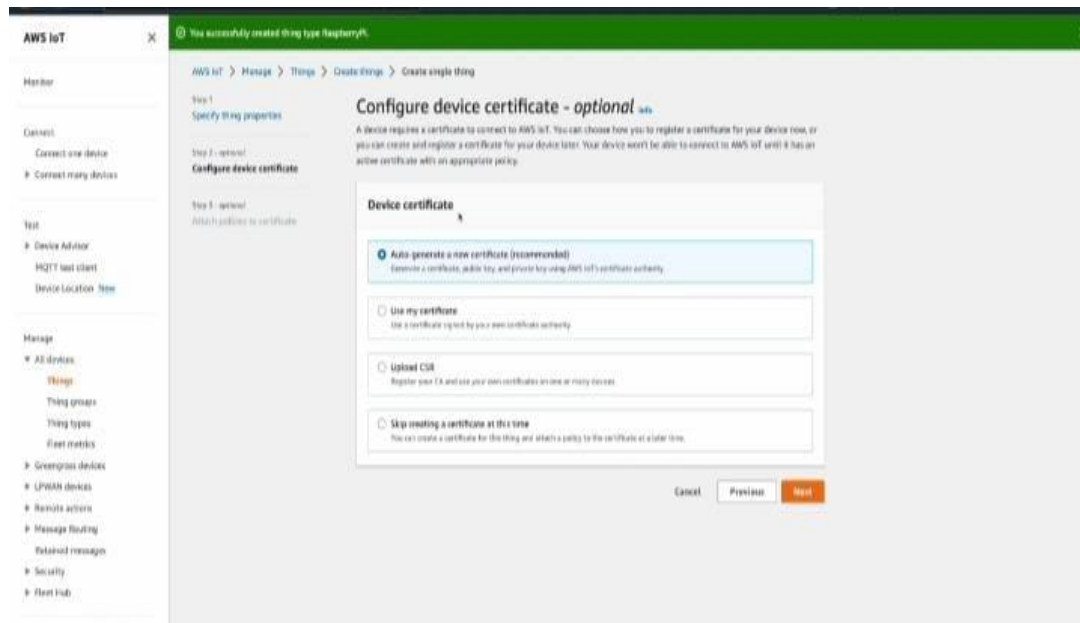
Code in GitHub Repo

# 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.
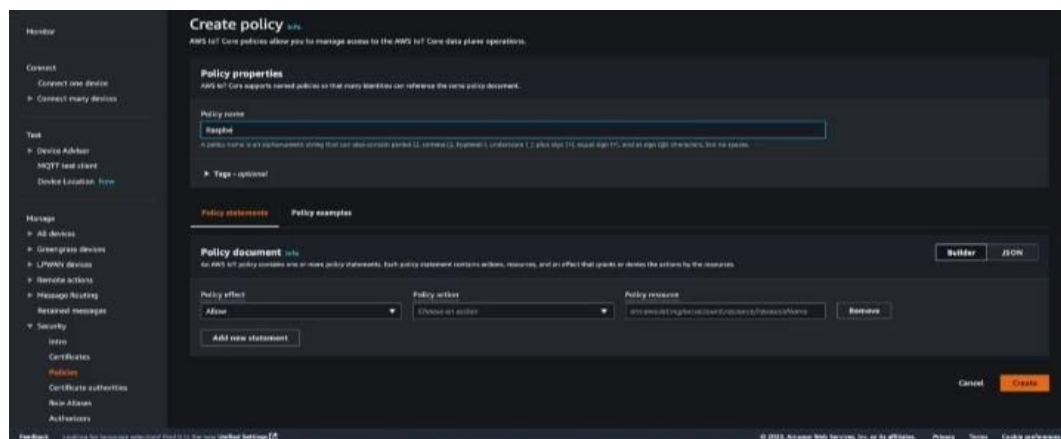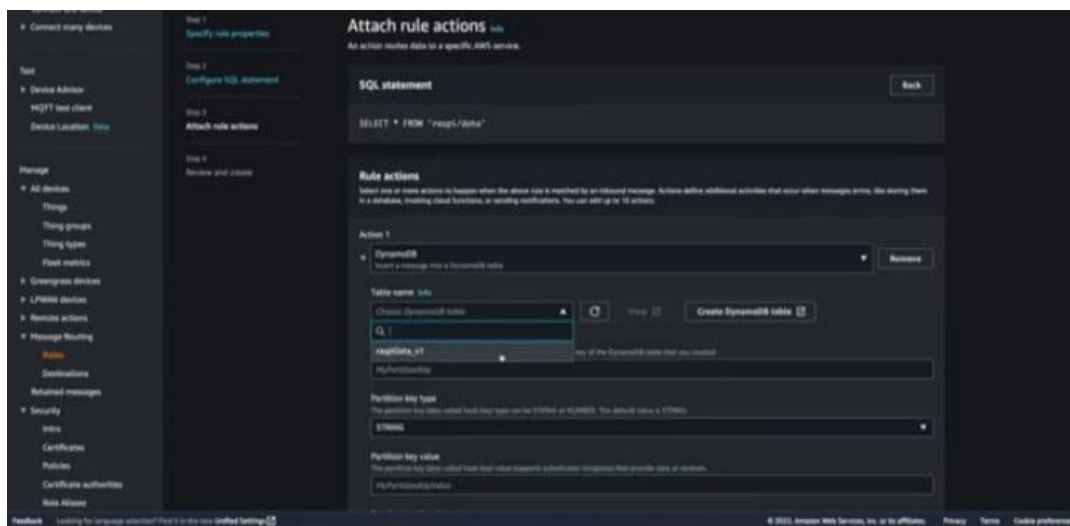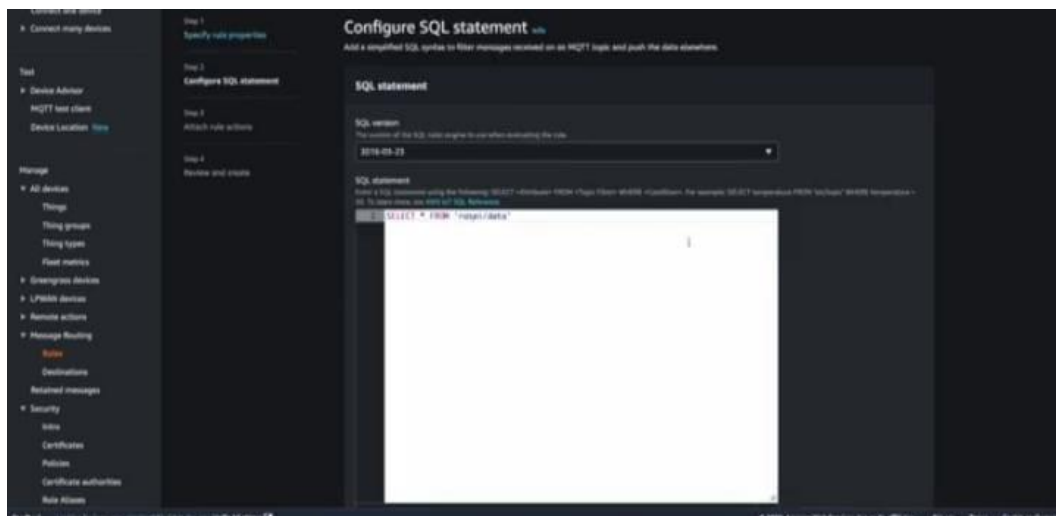


# Security and governance layer:

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.