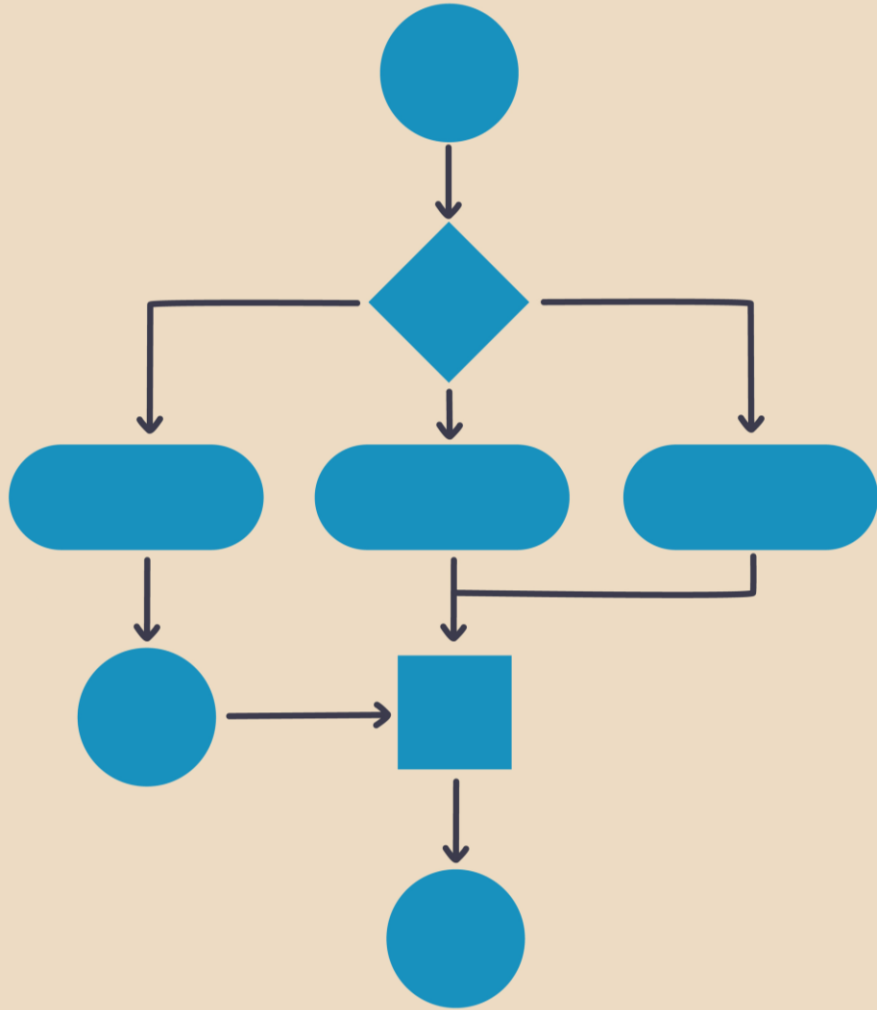# Orchestration of Data Pipeline

# Orchestration of Data Pipeline



- Automation for data movement from multiple storage locations into a centralised repository.



Extraction | Staging Area | Transform | Load | Data Warehouse | Analytics

- May include various tools, services, triggers, paths, sub-paths, schedules.
- May involve sequential, parallel processing steps.
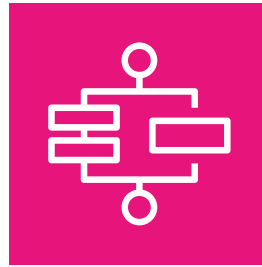- Should be a robust system that can handle failures, retries.
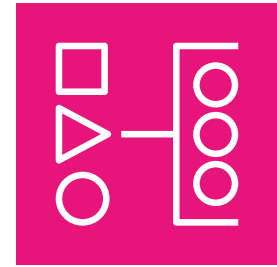
Orchestration Service in AWS

# Orchestration Service in AWS

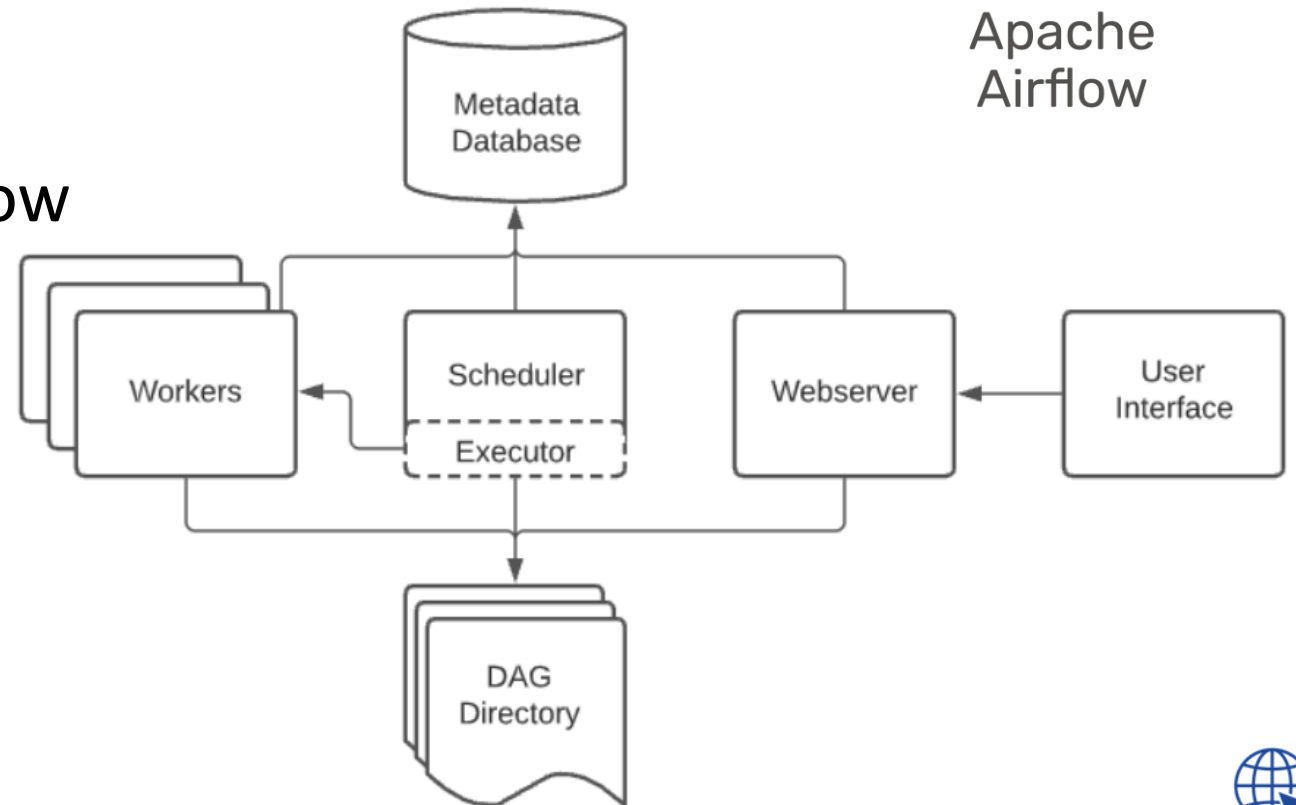AWS Glue
Workflows

AWS Step
Functions

Amazon Managed
Workflows for
Apache Airflow
(Amazon MWAA)

# Apache Airflow

- Apache Airflow is an open-source platform for developing, scheduling, and monitoring batch-oriented workflows.

- Key Use Cases for Apache Airflow
  - ETL
  - AI/ML
  - DevOps

- Main Components of Apache Airflow
  - Meta Database
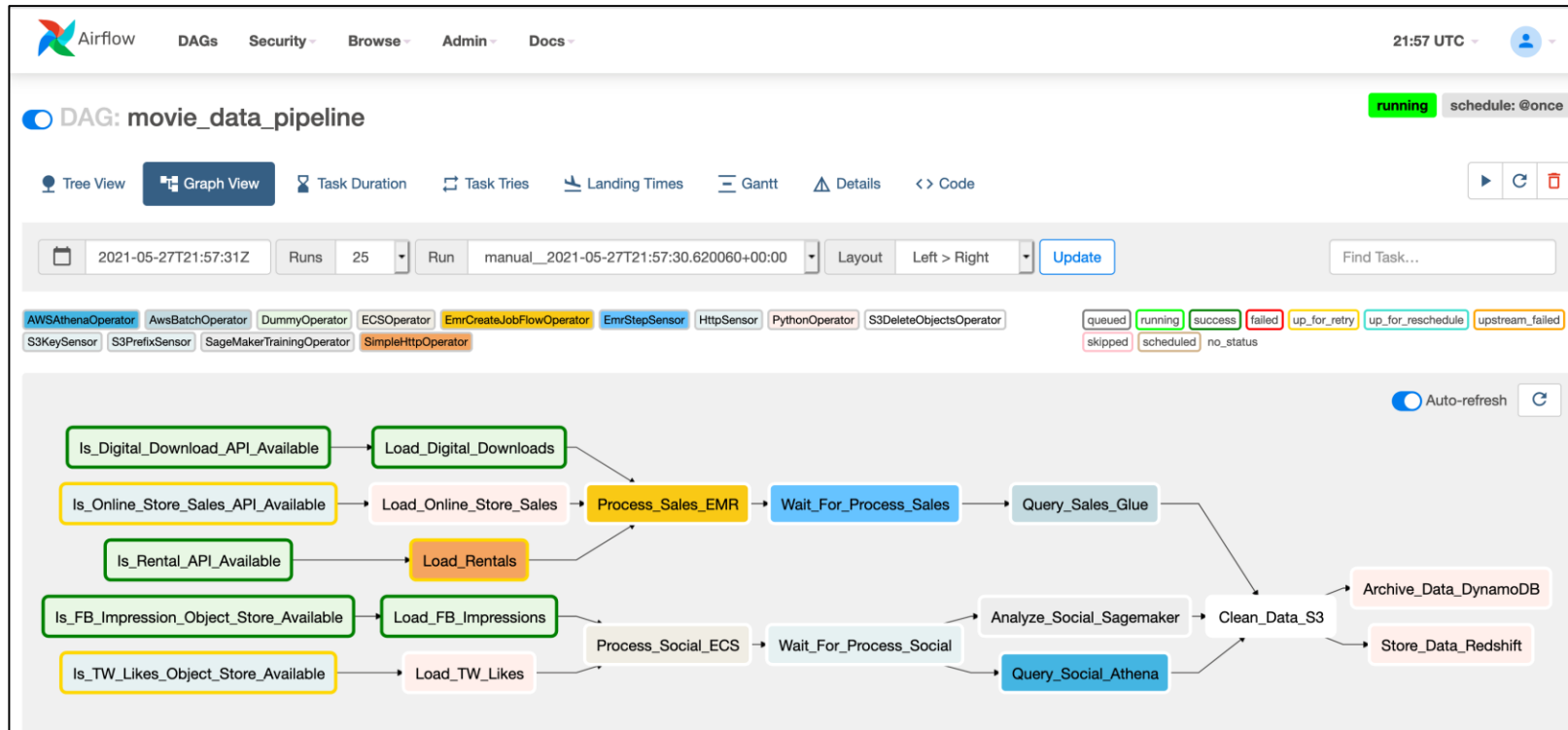  - Scheduler
  - Worker
  - Web Server

# Apache Airflow

- In Airflow, a DAG – or a Directed Acyclic Graph – is a workflow: a collection of all the tasks you want to run, organized to reflect their relationships and dependencies.

- A DAG is defined in a Python script, which represents the tasks and their dependencies as code.
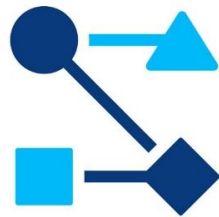
# Challenges with Self-Managed Apache Airflow
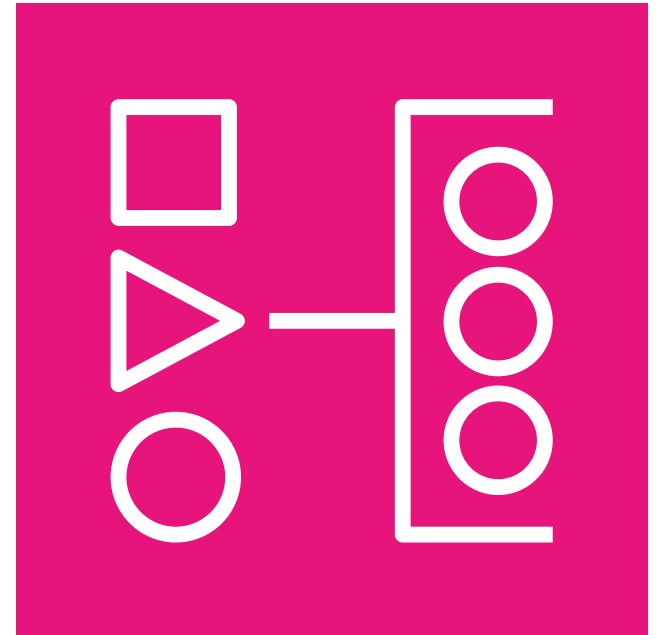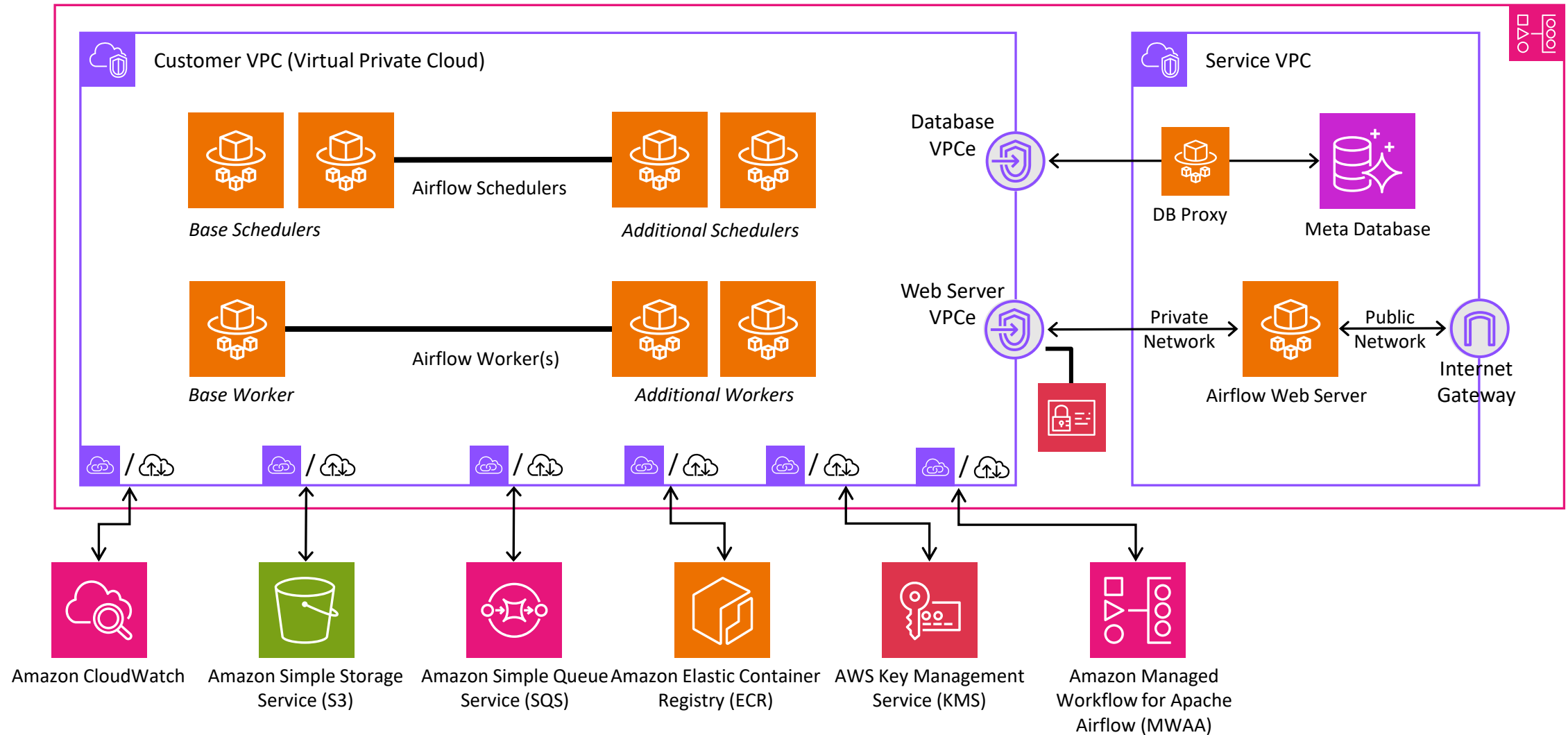
Setup

Scaling

Security

Upgrade

Maintenance

# Amazon Managed Workflows for Apache Airflow (MWAA)

- A managed service for Apache Airflow, making it easy for data engineers and data scientists to invoke data processing workflows on AWS.

- MWAA is 100% Open Source. It is not a Fork or "branch" of Airflow.

- How Amazon MWAA Helps?
  - Deployments and Operations
    - Easy to Set Up and Maintain
  - Availability and Sizing
    - Multi-AZ for HA with Airflow on ECS Fargate
  - Scaling
    - Auto Scaling and Celery Executor
  - Security
    - IAM and VPC

# Amazon MWAA Architecture

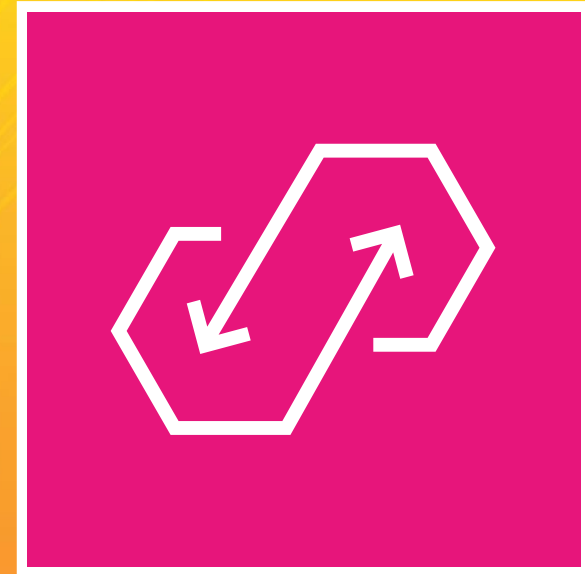| Comparison | AWS Glue Workflow | AWS Step Function | Amazon Managed Workflows for Apache Airflow (Amazon MWAA) |
|---|---|---|---|
| Terminology | Workflows | State Machine | Directed Acyclic Graph (DAG) |
| Deployment | Serverless | Serverless | Runs on managed servers |
| Authoring | Visual Designer SDK, CDK | Visual Designer SDK, CDK, State Language | Python |
| Type | Purpose built | General purpose | General purpose |
| High availability | Built-in | Built-in | Configurable across AZs |
| Integration | Primarily integrate with Glue services (crawlers, jobs) | Out-of-box integration with various AWS Services | Can be integrated with AWS Services and third-party tools as Apache Hadoop, Presto, Hive and Spark |
| Suitability | Glue specific orchestration | Integration with AWS Services | Migration from On-prem Airflow Prioritization of open source |
| Pricing | No charges for workflows. You pay of underlying Glue services. | Charged based on state transition (standard workflow) or duration (express workflow) | Infrastructure charges (compute and storage) |

Amazon AppFlow

# Amazon AppFlow

- A fully managed integration service that enables you to securely transfer data between SaaS applications like Salesforce, SAP, Zendesk, Slack, and ServiceNow, and AWS services like Amazon S3 and Amazon Redshift.

| NO / LOW CODE | COST SAVINGS | SPEED & AGILITY | SECURE & SCALABLE |
|---|---|---|---|

**Source**
Ingest data from supported SaaS application

| Salesforce | Marketo |
|---|---|
| SAP | CircleCI |
| Google Analytics | And many more… |

| Mask fields | Catalog data | Map and merge |
|---|---|---|

Amazon AppFlow

| Filter & validate | Partition, aggregate | Add formulas |
|---|---|---|

**Destination**
Transfer data to supported destinations for processing or storage

| Amazon Redshift |
|---|
| Amazon S3 |
| And many more… |

# How it works?

## Creating a flow

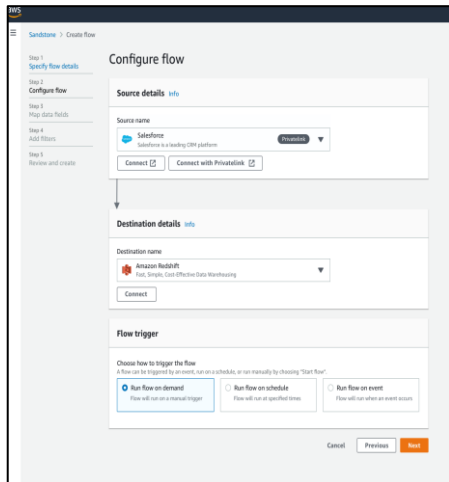**Select source and destination; login** → **Specify flow trigger** → **Map fields from source to destination** → **Add filters, validation, transformation** → **Activate your flow or run it at the click of a button**
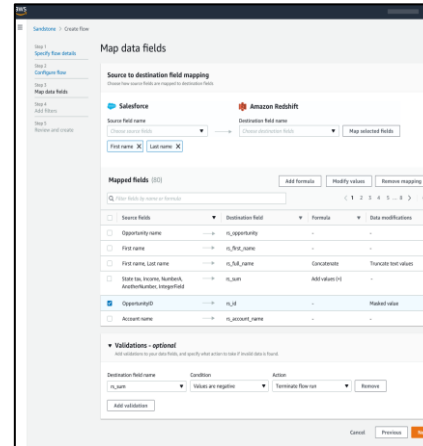
**Event Based**

**Scheduled (1 minute minimum)**

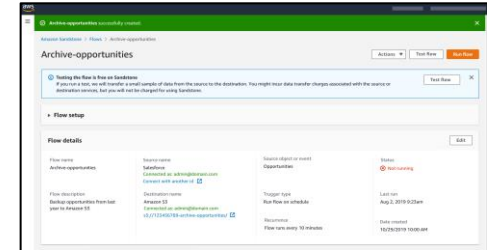**On-demand (run flow immediately)**

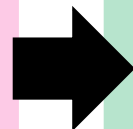**Filter records**

**Mask sensitive fields**

**Combine fields to create new ones**
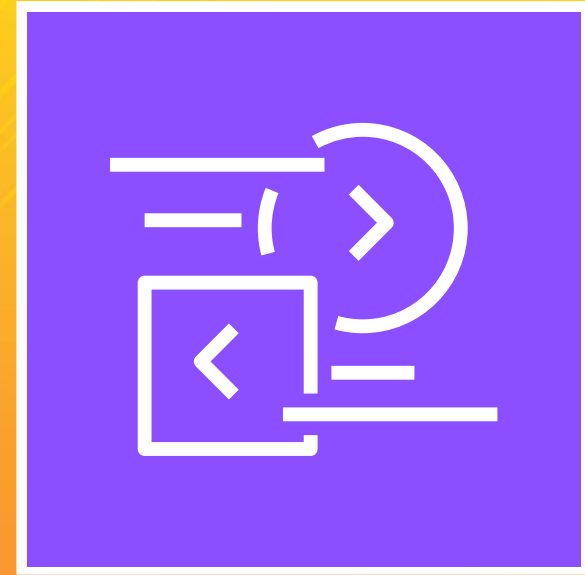
**Validate records**

**Truncate fields**

User defines the data flow requirements using UI → Amazon AppFlow provisions compute, storage and networking resource and executes the flow

AWS Data Exchange

# How we use Mobile applications?



App Publishers
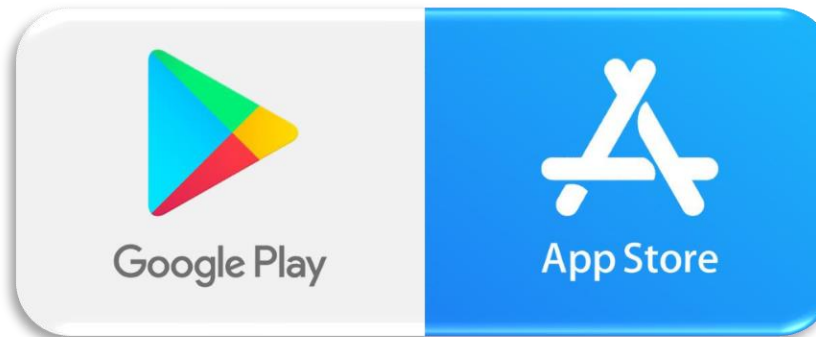
App Consumers

| Mobile client | Mobile client | Mobile client | Mobile client | Mobile client | Mobile client |

# AWS Data Exchange

| Financial Services | Healthcare & Life Sciences | Marketing & Advertising | Media & Entertainment | Technology | Travel & Hospitality | CPG & Retail |
|---|---|---|---|---|---|---|

3,500+ data sets from 300+ data providers      **Data Publishers**

**AWS Data Exchange**

Data Files, Data Tables, Data APIs

**Data Consumers**

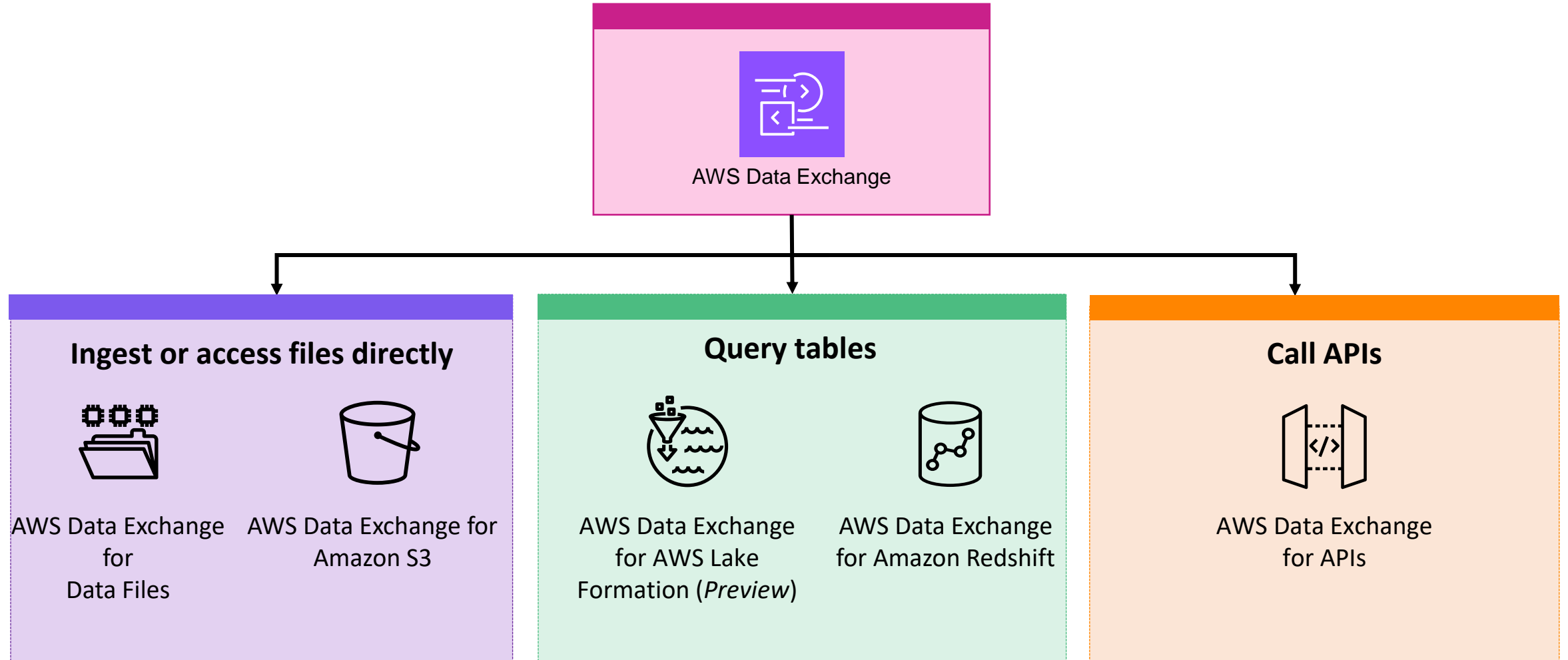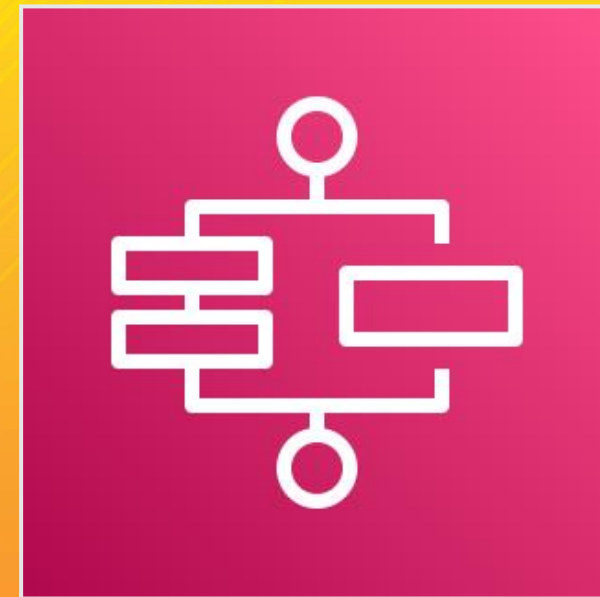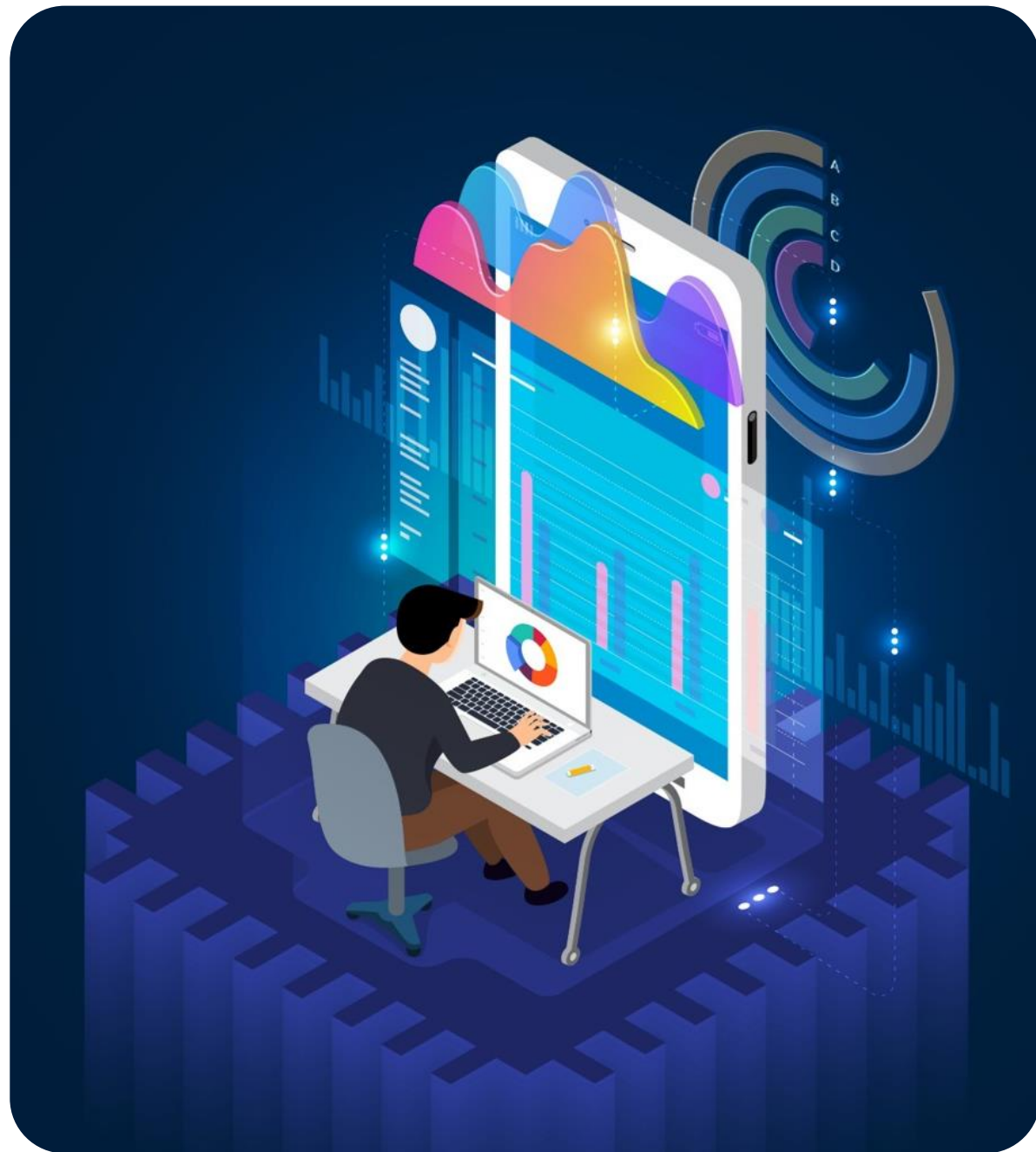| Automatically export new or updated data sets to Amazon S3 | Query data directly from vendor's databases | Use AWS-native authentication and governance, AWS SDKs |
|---|---|---|

# AWS Data Exchange

- AWS Data Exchange makes it easy to ingest third-party data into your data warehouses and data lakes

- It allows you to quickly and easily start using third-party data in your applications, analytics, and machine learning models.

- Data providers benefits
  - Package and publish data sets inside free or paid products for the price and terms that you control.
  - Data owners can curate data in S3, local files, Redshift datashares, APIs, or Lake Formation tables/databases.
  - Reach millions of AWS customers

- Data subscribers benefits
  - Quickly find diverse data in one place and analyze data as it's published
  - Spend less time and effort to use data in production.

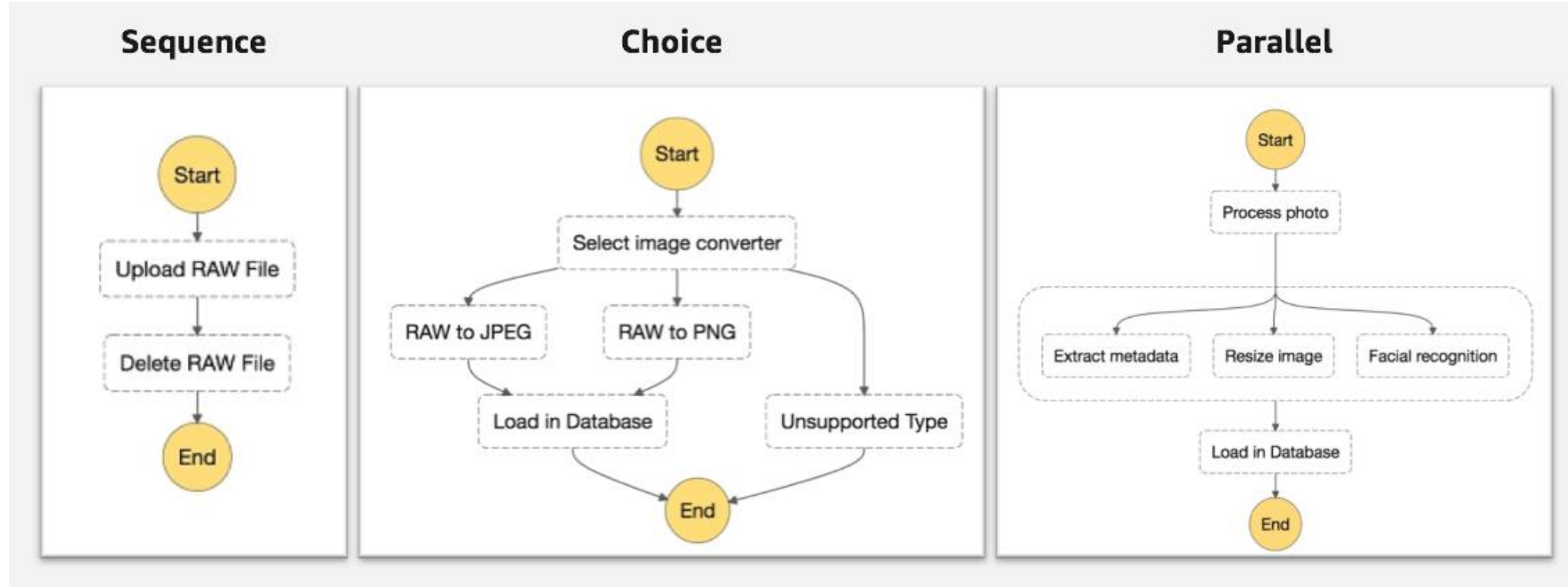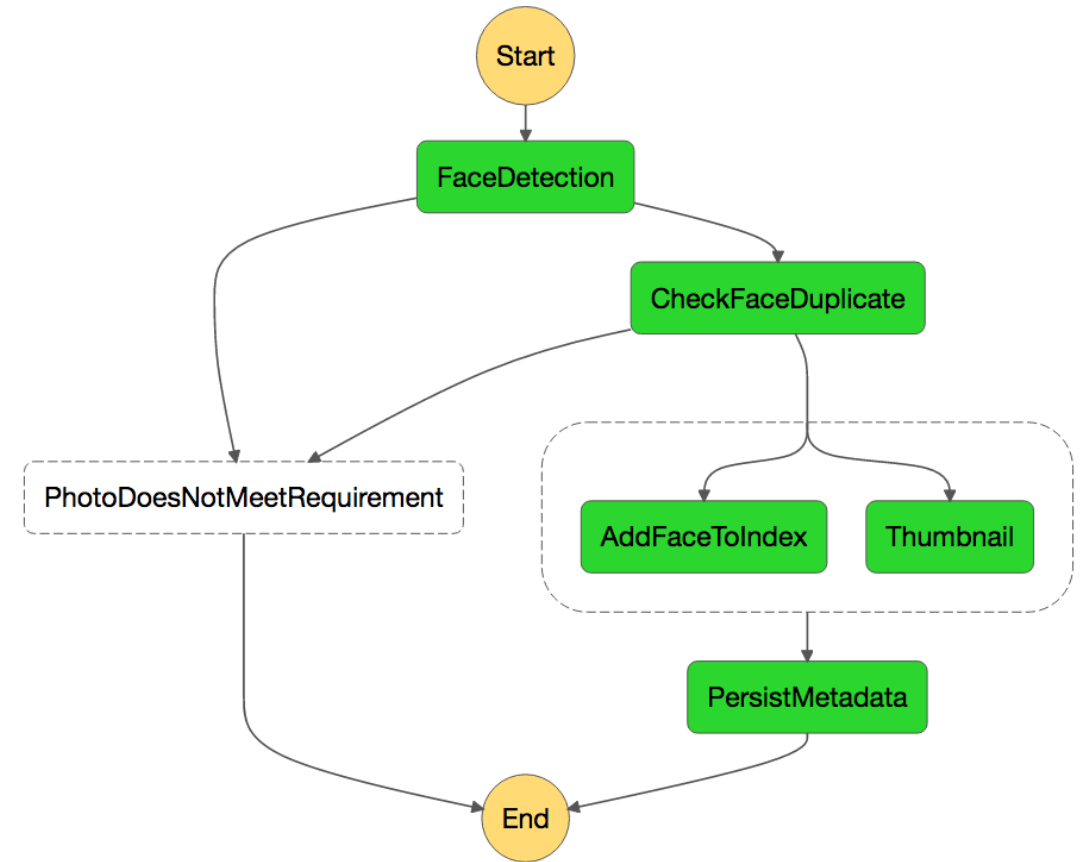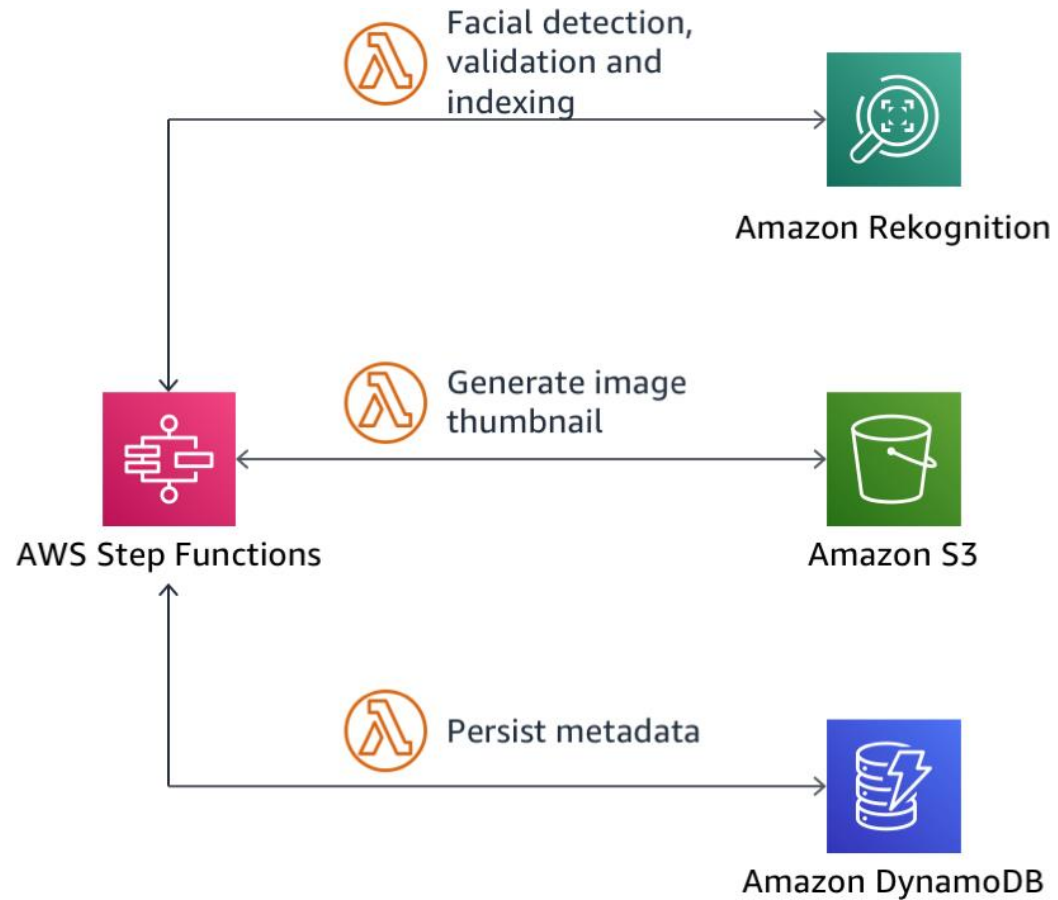# Ways to license third-party data faster and easier

AWS Data Exchange

## Ingest or access files directly

AWS Data Exchange for Data Files

AWS Data Exchange for Amazon S3

## Query tables

AWS Data Exchange for AWS Lake Formation (*Preview*)

AWS Data Exchange for Amazon Redshift

## Call APIs

AWS Data Exchange for APIs
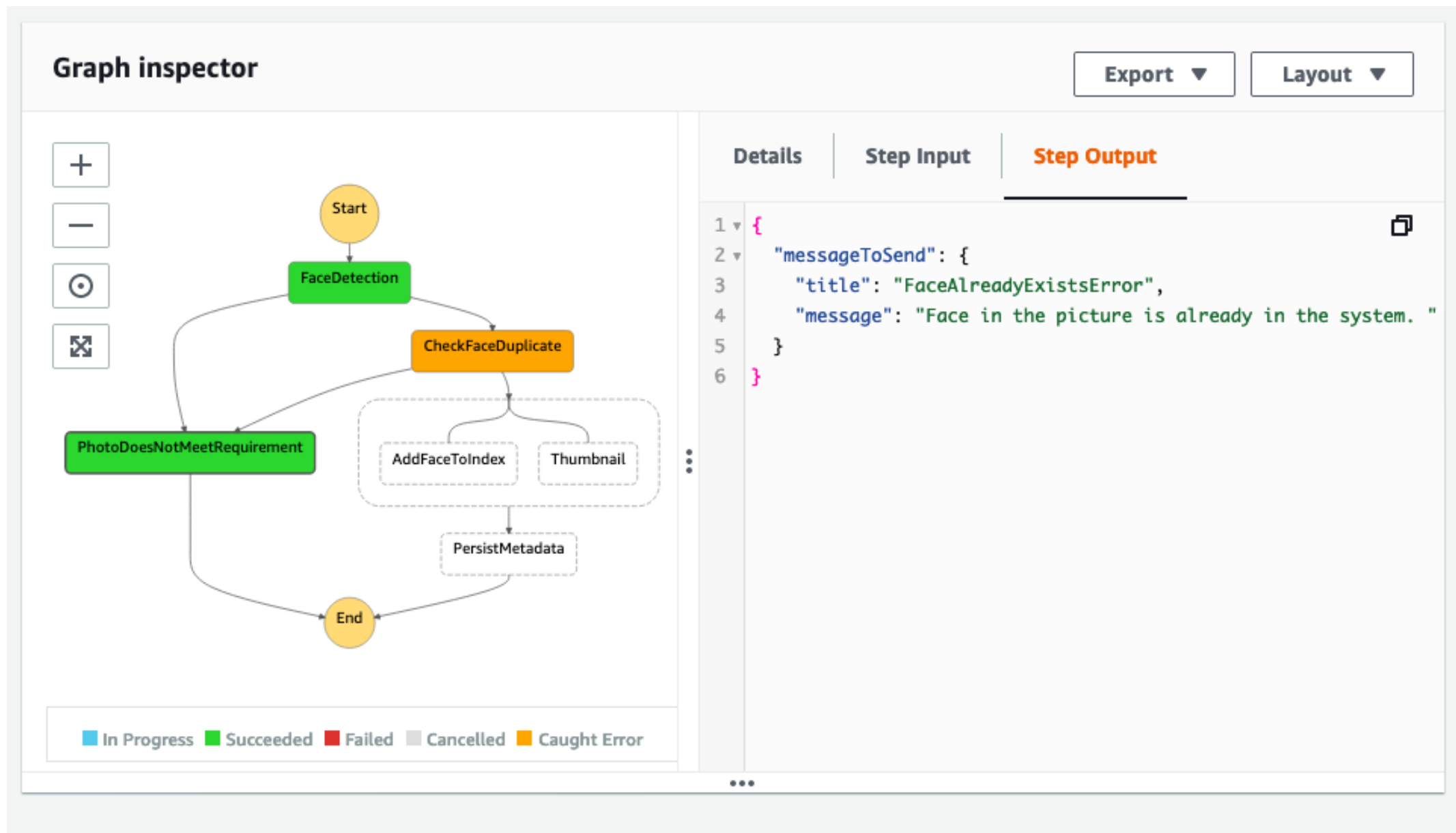
AWS Step Function

# AWS Step Functions

- A workflow for Lambda Function
- Visual Interface to create and run workflows
- The output of one step acts as an input to the next.

# Serverless Image Processing Workshop

# Serverless Image Processing Workshop

# AWS Lambda Power Tuning

# Things Go Better With Step Functions