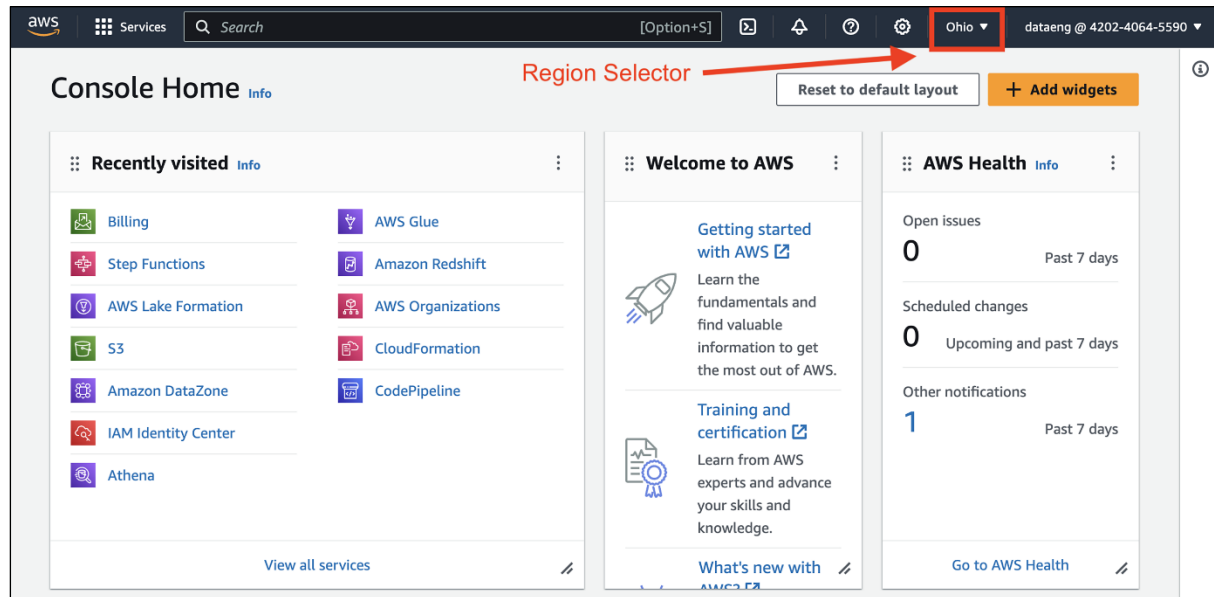


Chapter 1: An Introduction to Data Engineering



User name

dataeng

The user name can have up to 64 characters. Valid characters: A-Z, a-z, 0-9, and + = , . @ _ - (hyphen)

☒ Provide user access to the AWS Management Console - *optional*

If you're providing console access to a person, it's a [best practice](#) to manage their access in IAM Identity Center.

Are you providing console access to a person?

User type

☐ Specify a user in Identity Center - Recommended

We recommend that you use Identity Center to provide console access to a person. With Identity Center, you can centrally manage user access to their AWS accounts and cloud applications.

☒ I want to create an IAM user

We recommend that you create IAM users only if you need to enable programmatic access through access keys, service-specific credentials for AWS CodeCommit or Amazon Keyspaces, or a backup credential for emergency account access.

Console password

☐ Autogenerated password

You can view the password after you create the user.

☒ Custom password

Enter a custom password for the user.

- Must be at least 8 characters long
- Must include at least three of the following mix of character types: uppercase letters (A-Z), lowercase letters (a-z), numbers (0-9), and symbols ! @ # \$ % ^ & * () _ + - (hyphen) = [] { } | ' "

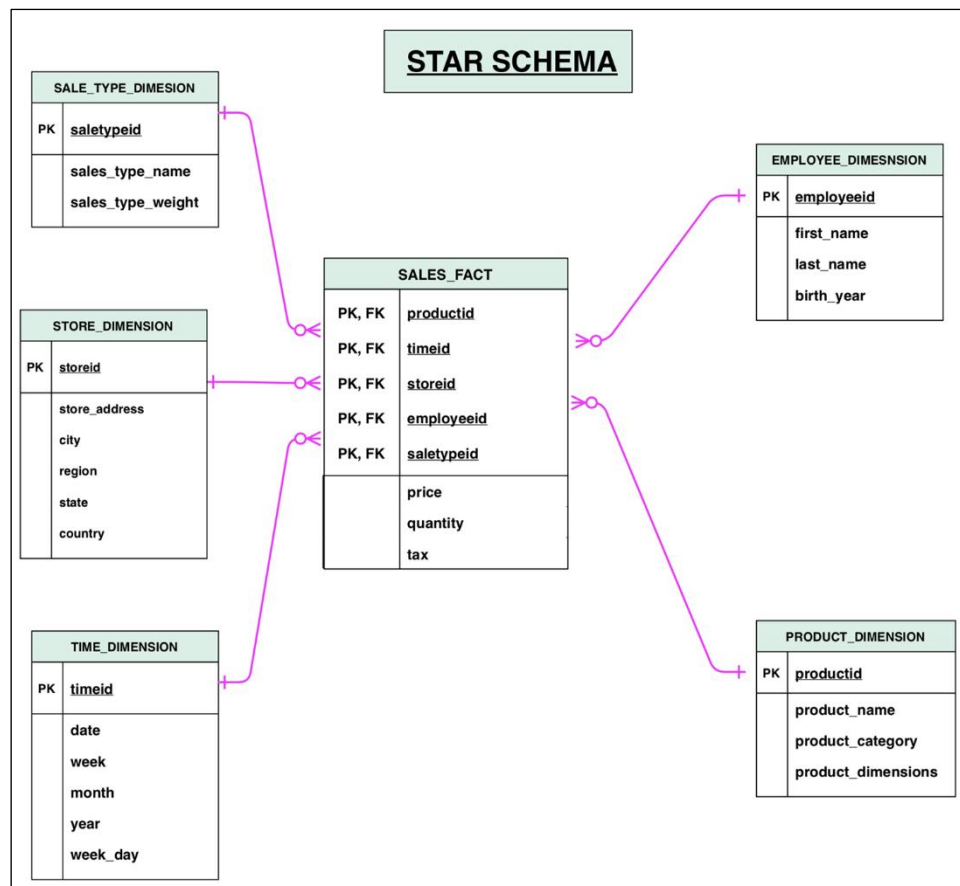
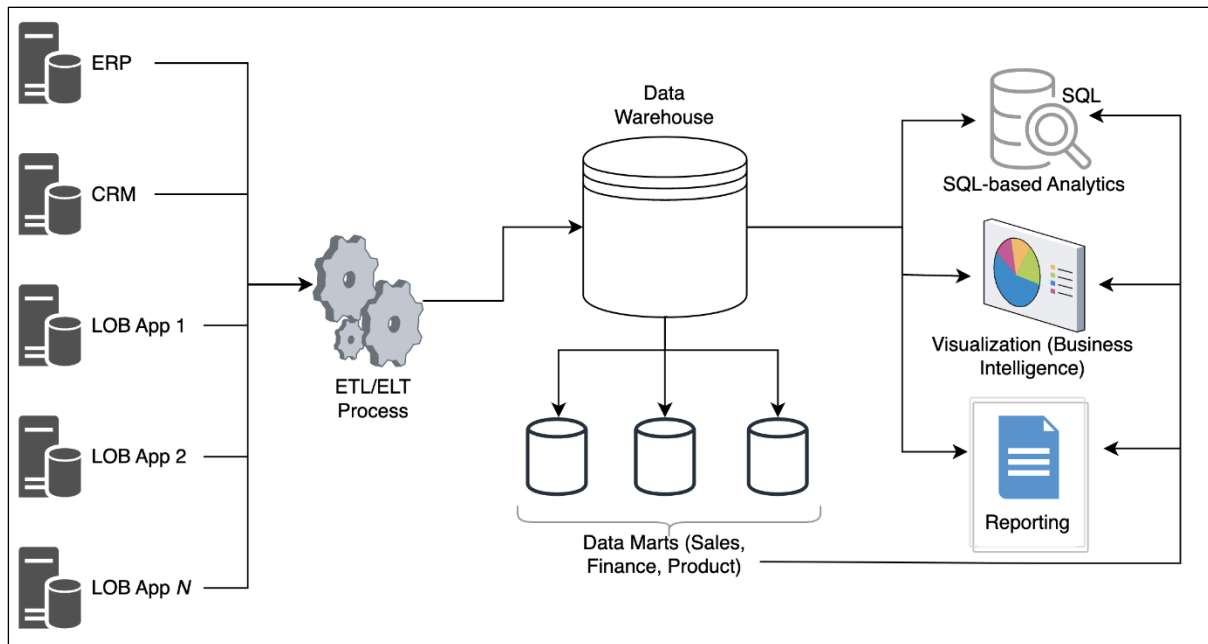
☐ Show password

☐ Users must create a new password at next sign-in - Recommended

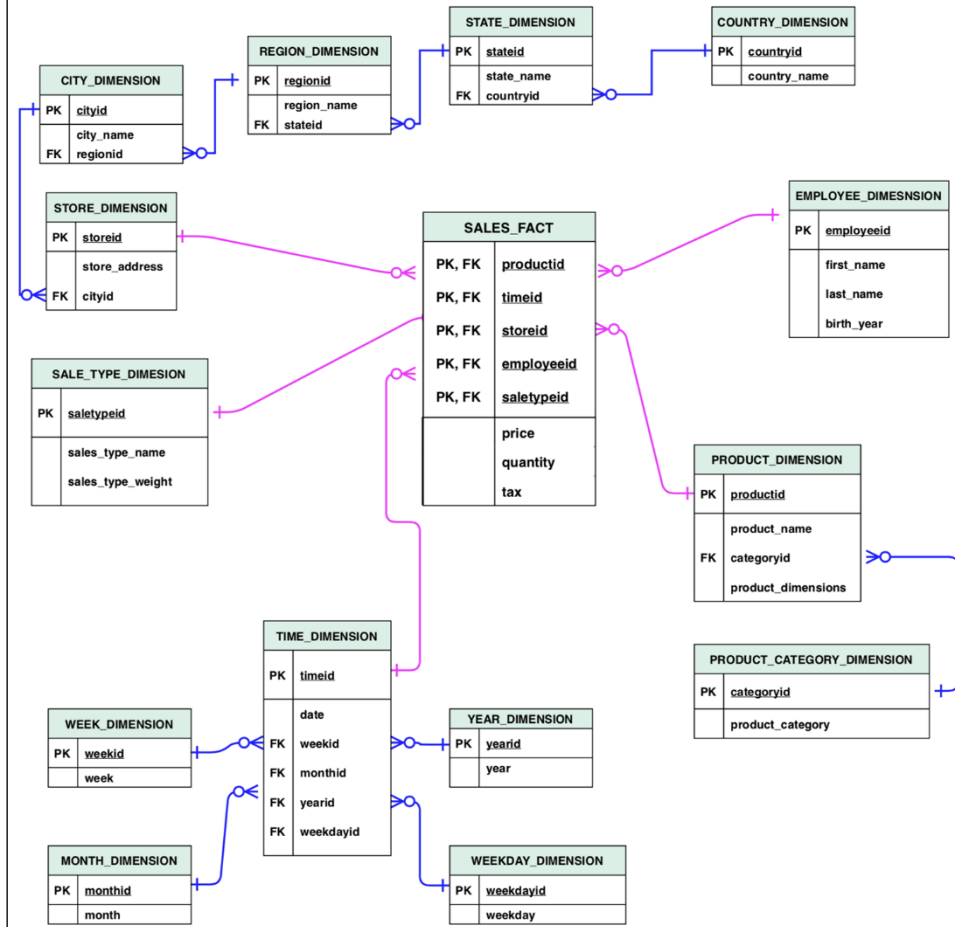
Users automatically get the [IAMUserChangePassword](#) policy to allow them to change their own password.

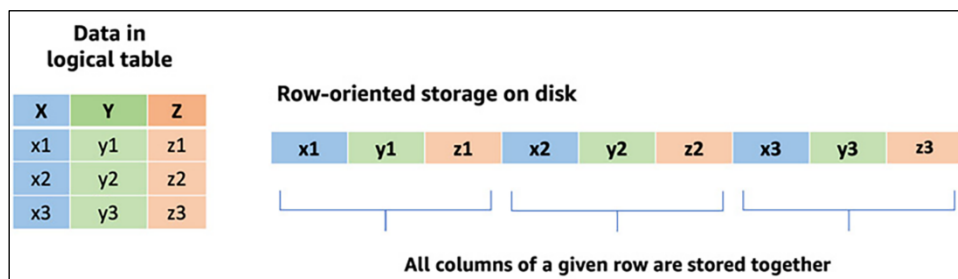
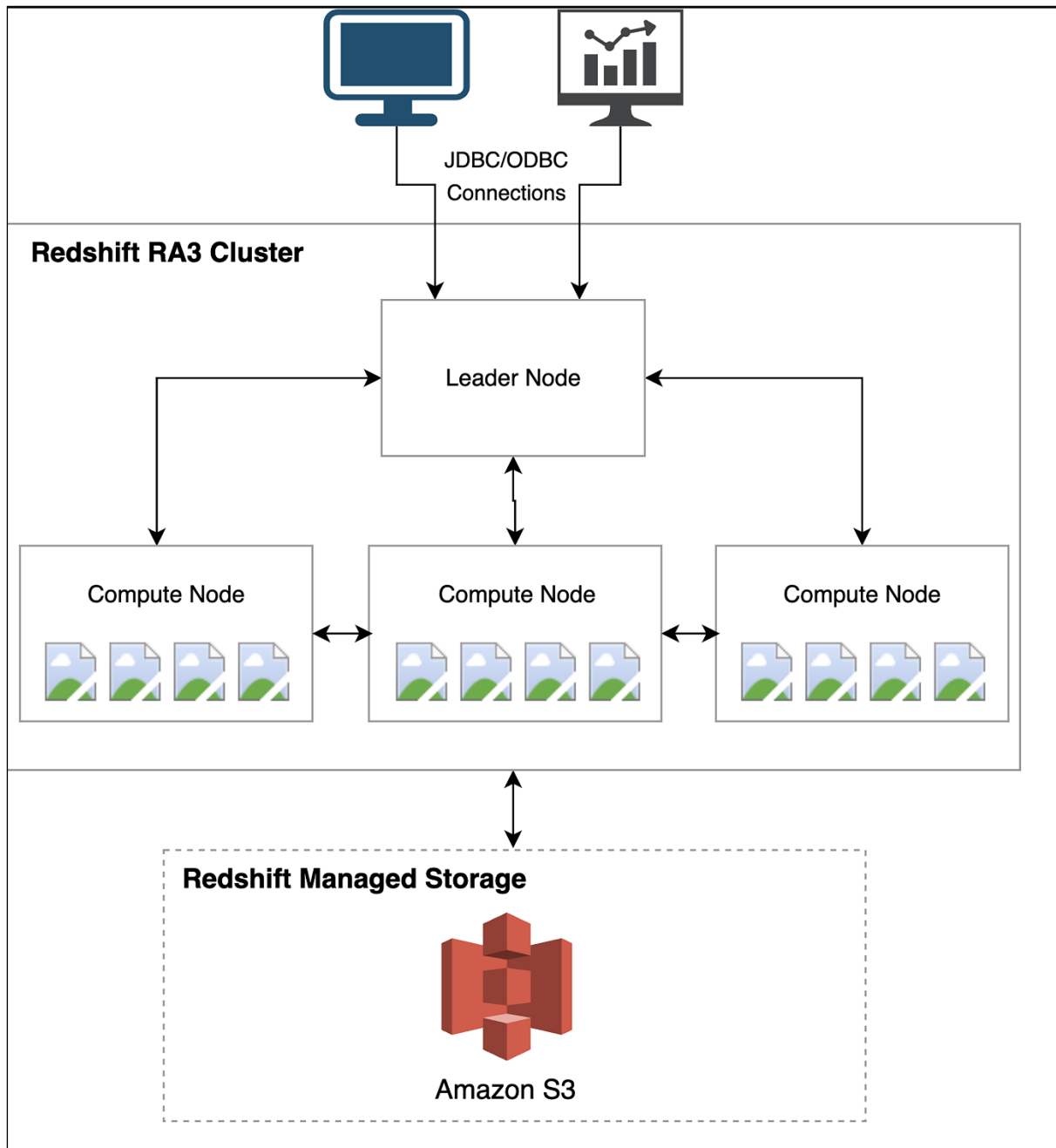
If you are creating programmatic access through access keys or service-specific credentials for AWS CodeCommit or Amazon Keyspaces, you can generate them after you create this IAM user. [Learn more](#)

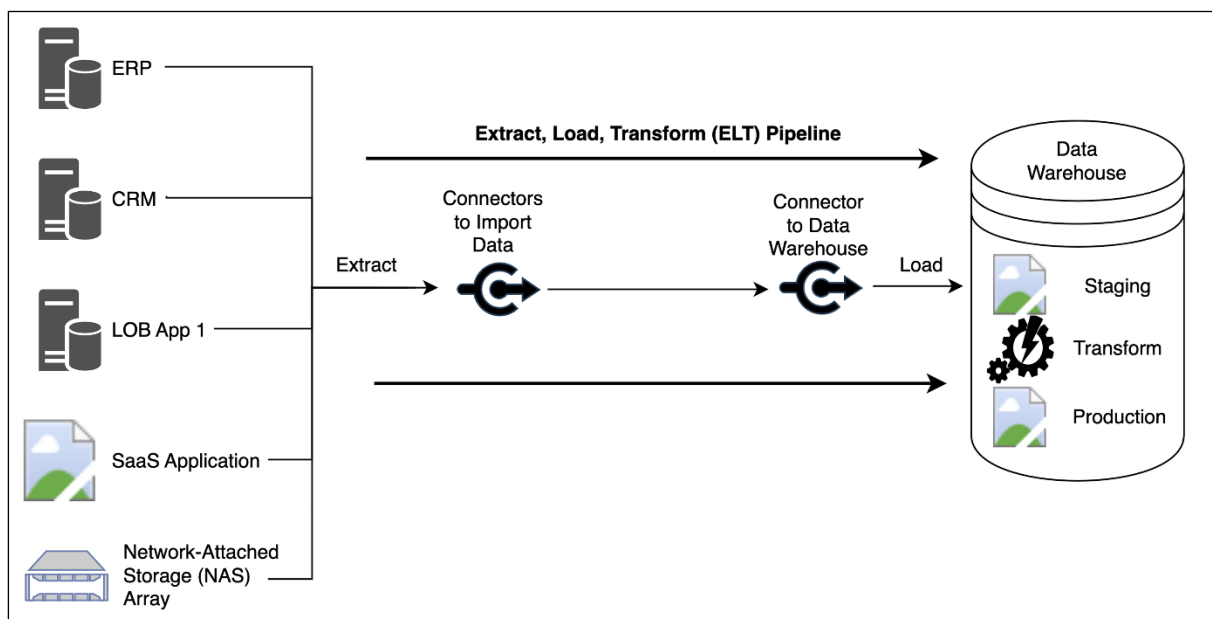
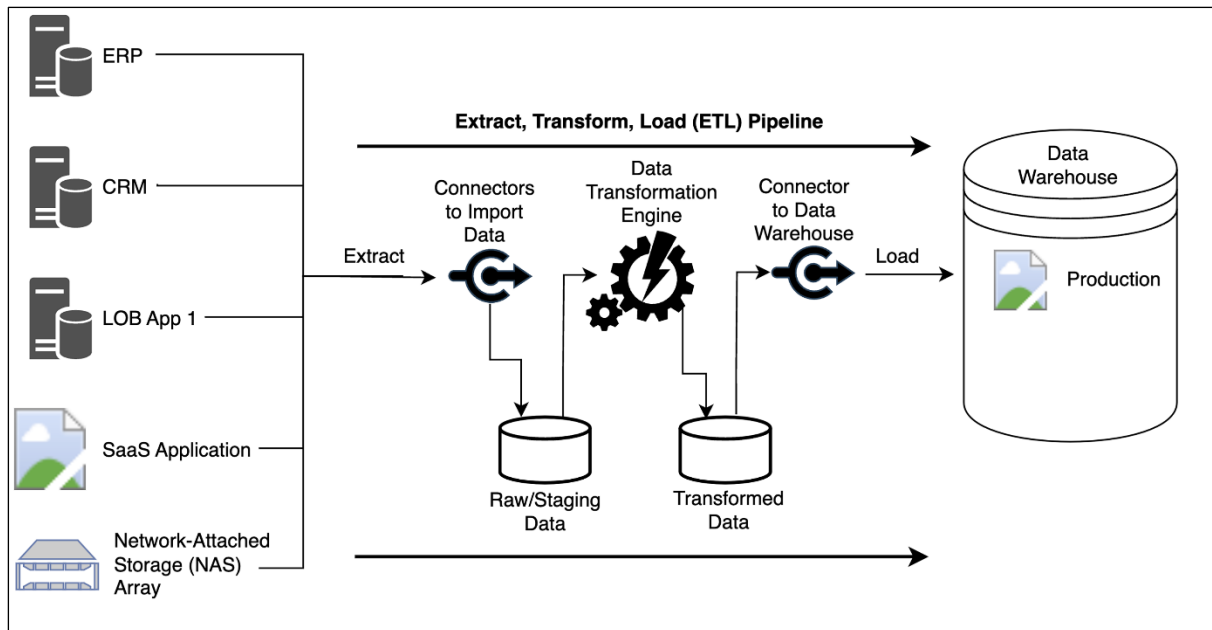
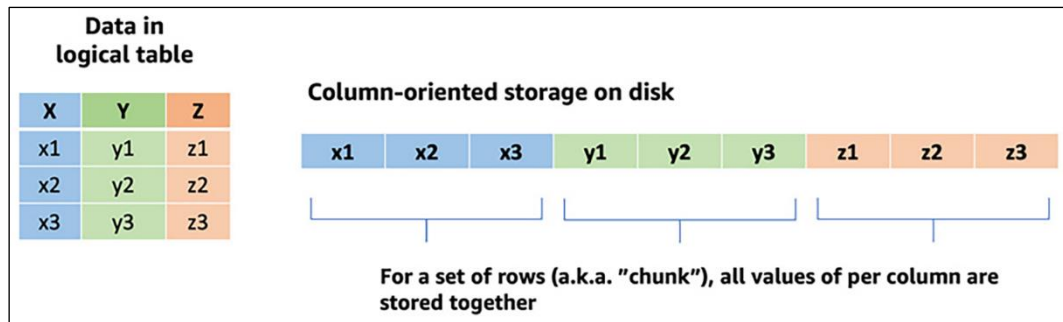
Chapter 2: Data Management Architectures for Analytics

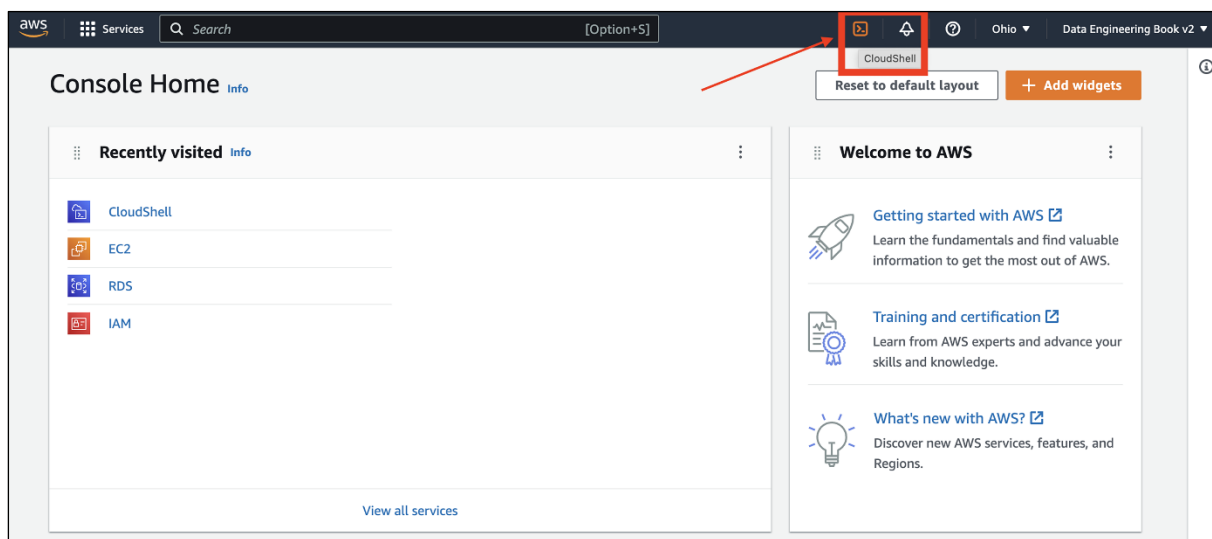
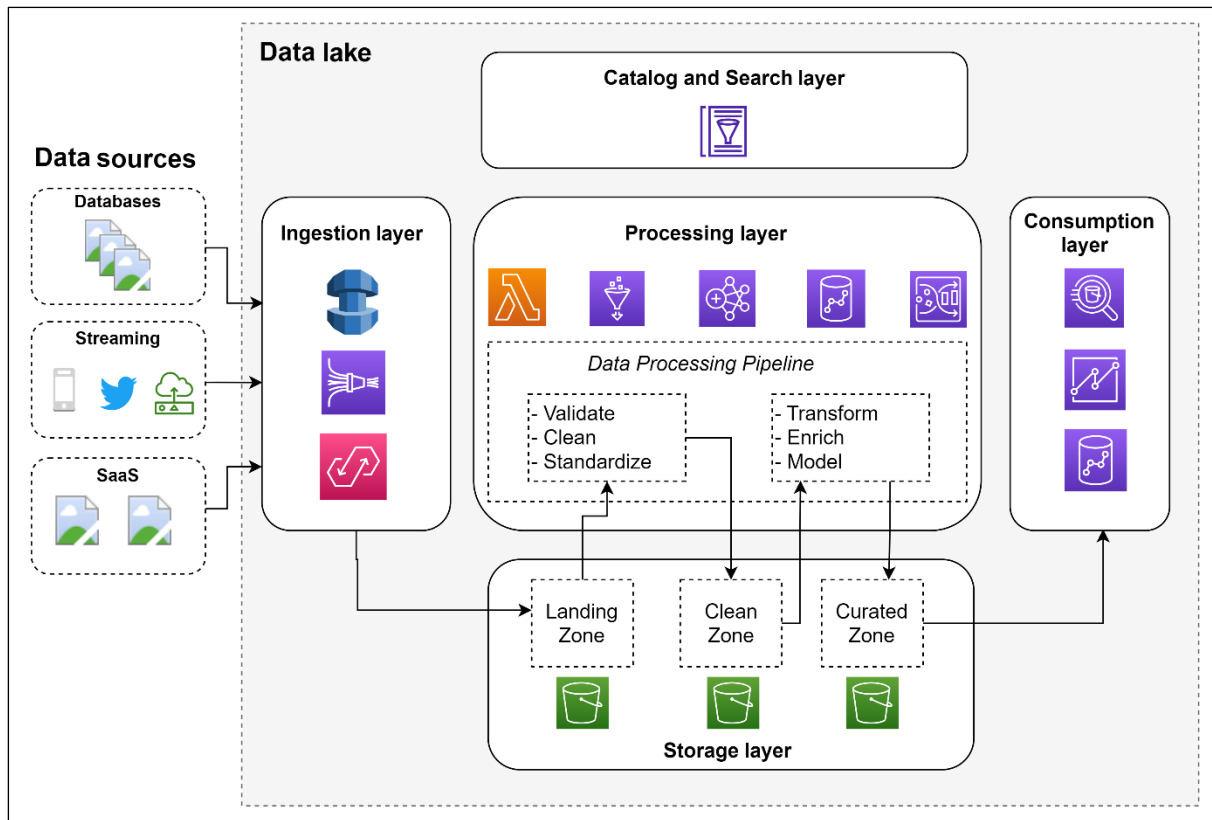


SNOWFLAKE SCHEMA











AWS CloudShell

us-east-2

S3()

S3()

NAME

s3 -

DESCRIPTION

This section explains prominent concepts and notations in the set of high-level S3 commands provided.

If you are looking for the low level S3 commands for the CLI, please see the [s3api command reference page](#).

Path Argument Type

Whenever using a command, at least one path argument must be specified. There are two types of path arguments: LocalPath and S3Uri.

LocalPath: represents the path of a local file or directory. It can be written as an absolute path or relative path.

S3Uri: represents the location of a S3 object, prefix, or bucket. This must be written in the form `s3://mybucket/mykey` where `mybucket` is the specified S3 bucket, `mykey` is the specified S3 key. The path argument must begin with `s3://` in order to denote that the path argument refers to a S3 object. Note that prefixes are separated by forward slashes. For example, if the S3 object `myobject` had the prefix `myprefix`, the S3 key would be `myprefix/myobject`, and if the object was in the bucket `mybucket`, the S3Uri would be `s3://mybucket/myprefix/myobject`.

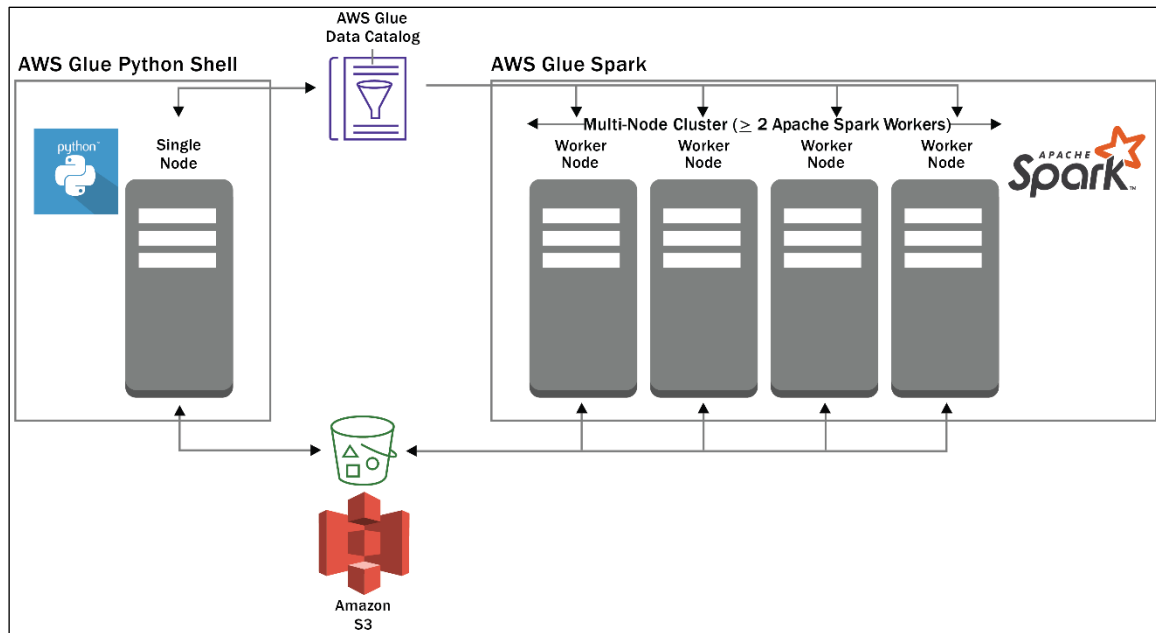
S3Uri also supports S3 access points. To specify an access point, this value must be of the form `s3://<access-point-arn>/<key>`. For example if the access point `myaccesspoint` to be used has the ARN: `arn:aws:s3:us-west-2:123456789012:accesspoint/myaccesspoint` and the object being accessed has the key `mykey`, then the S3URI used must be:

:

Feedback

Language

Chapter 3: The AWS Data Engineer's Toolkit



Amazon S3 > dataeng- > hr/ > employee/

employee/ Copy S3 URI

Objects Properties

Objects (20)

Objects are the fundamental entities stored in Amazon S3. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Refresh Delete Actions Create folder Upload

< 1 > Settings

<input type="checkbox"/>	Name	Type	Storage class
<input type="checkbox"/>	LOAD00000001.csv	csv	Standard
<input type="checkbox"/>	LOAD00000002.csv	csv	Standard
<input type="checkbox"/>	LOAD00000003.csv	csv	Standard
<input type="checkbox"/>	LOAD00000004.csv	csv	Standard
<input type="checkbox"/>	LOAD00000005.csv	csv	Standard
<input type="checkbox"/>	LOAD00000006.csv	csv	Standard
<input type="checkbox"/>	LOAD00000007.csv	csv	Standard
<input type="checkbox"/>	LOAD00000008.csv	csv	Standard
<input type="checkbox"/>	LOAD00000009.csv	csv	Standard
<input type="checkbox"/>	LOAD00000010.csv	csv	Standard
<input type="checkbox"/>	LOAD00000011.csv	csv	Standard

Name

employee

Description

Database

hr

Classification

csv

Location

s3://dataeng-temp/hr/employee/

Connection

Deprecated

No

Last updated

Wed Dec 09 21:46:50 GMT-500 2020

Input format

org.apache.hadoop.mapred.TextInputFormat

Output format

org.apache.hadoop.hive ql.io.HiveIgnoreKeyTextOutputFormat

Serde serialization lib

org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe

Serde parameters

field.delim ,

skip.header.line.count1sizeKey6300objectCount20UPDATED_BY_CRAWLERhr-employee-crawler

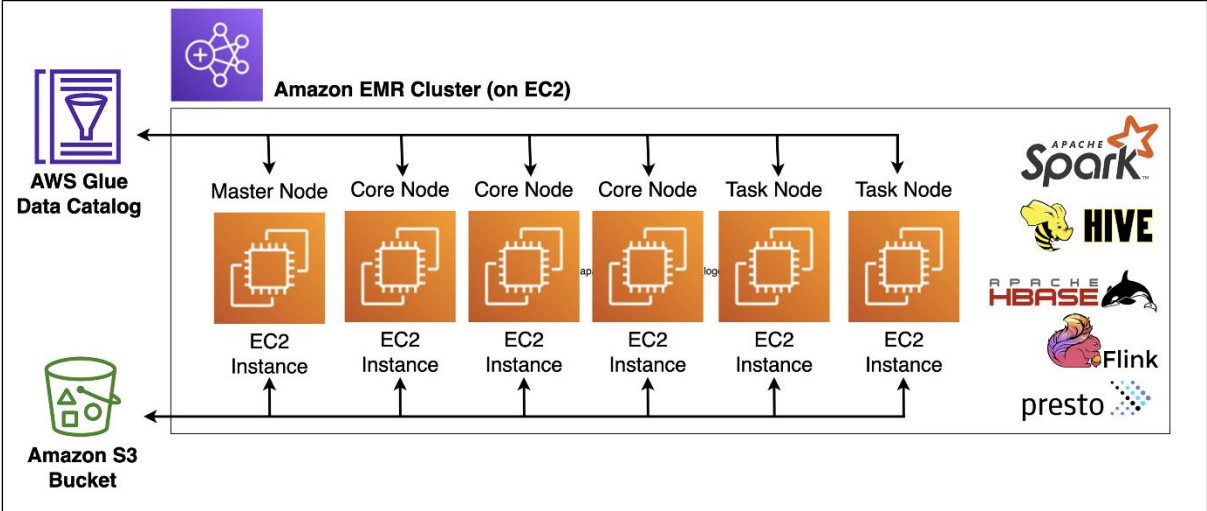
Table properties

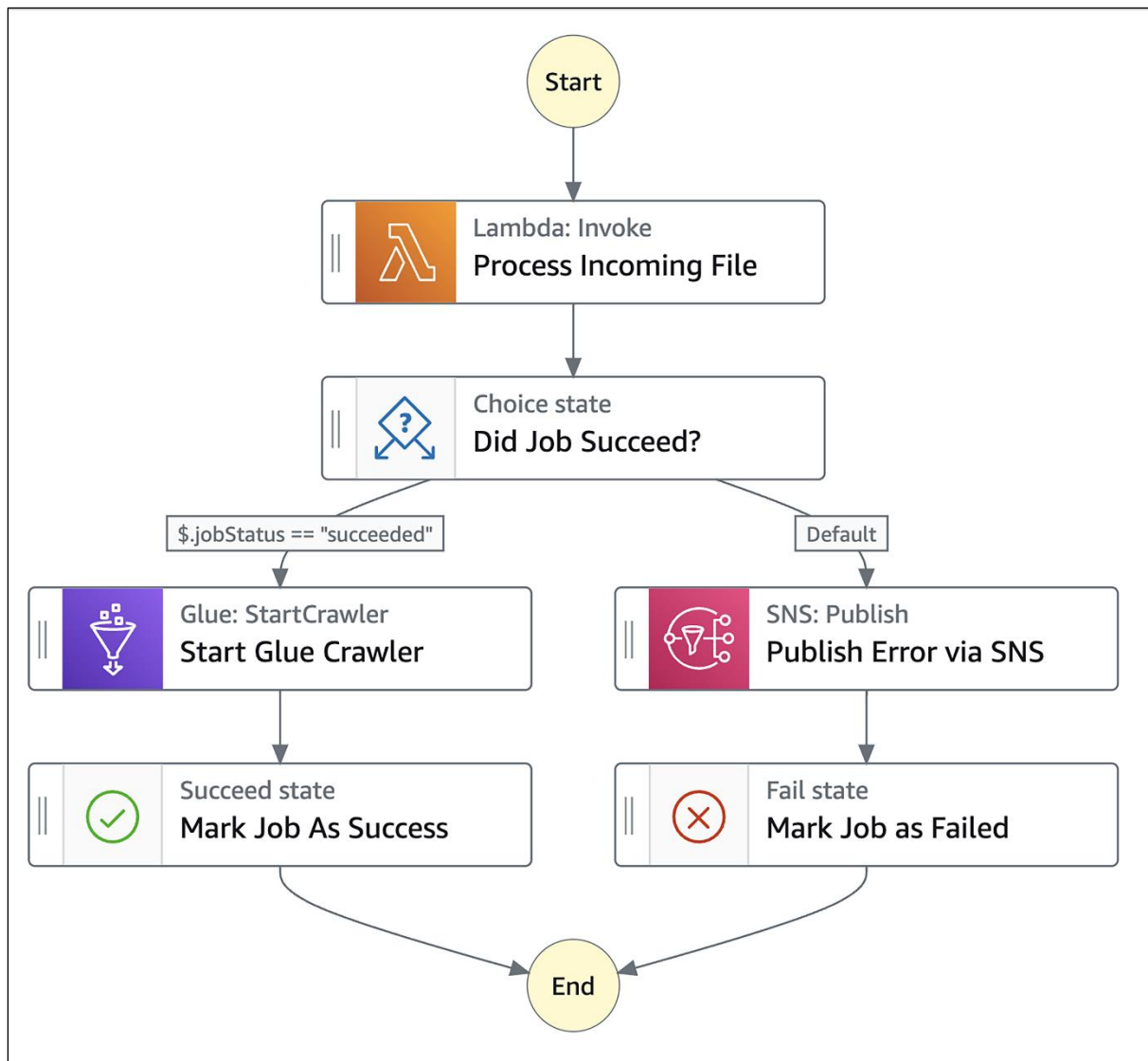
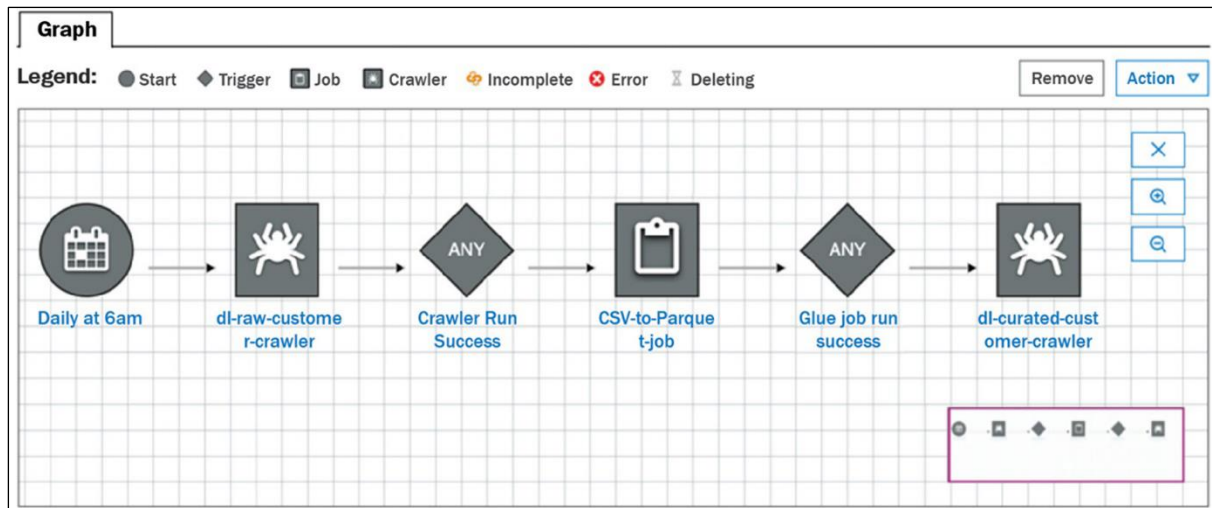
CrawlerSchemaSerializerVersion1.0recordCount40averageRecordSize156CrawlerSchemaDeserializerVersion1.0compressionTypenonecolumnsOrderedtrueareColumnsQuotedfalsedelimiter ,typeOfDatafile

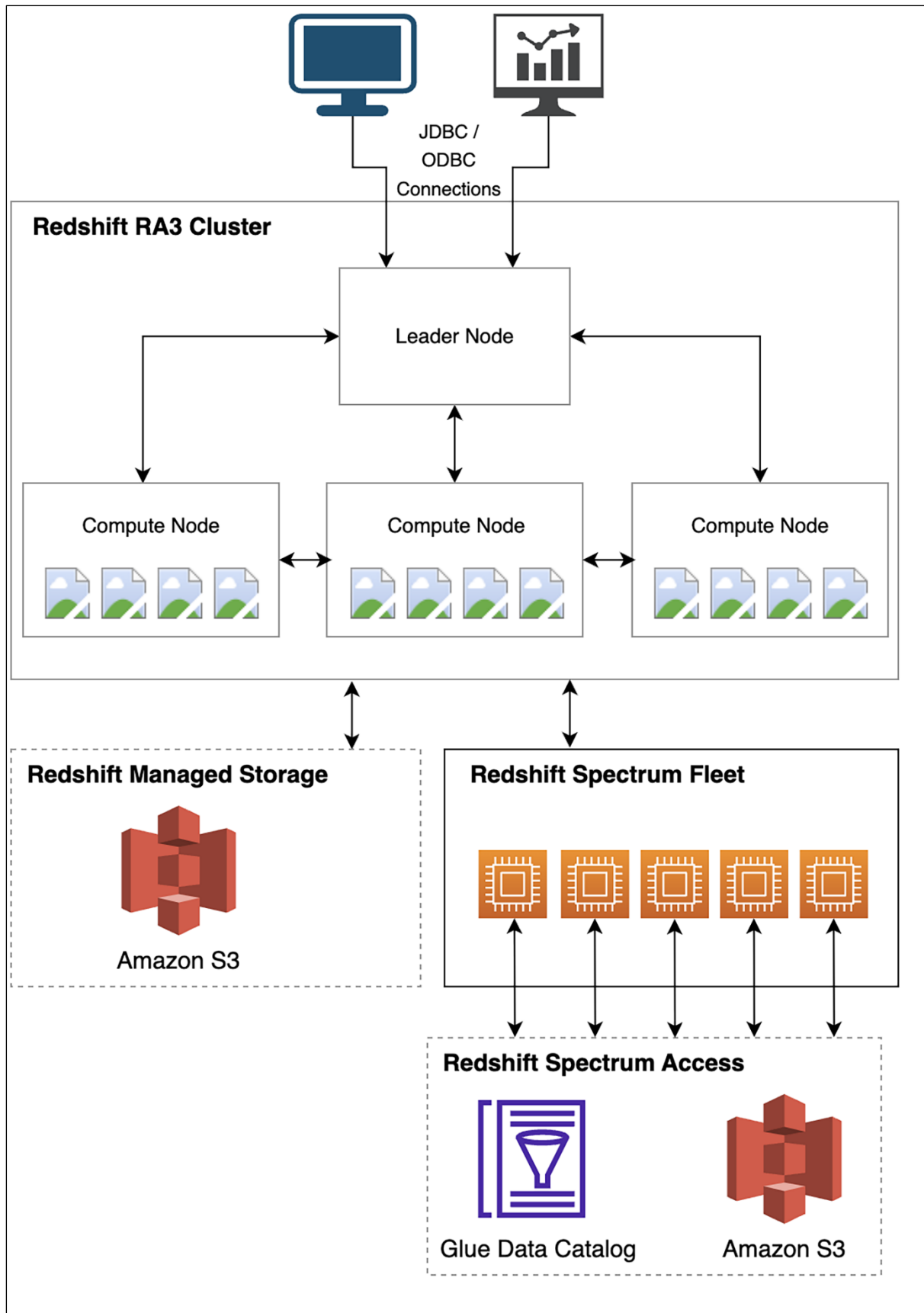
Schema

Showing: 1 -

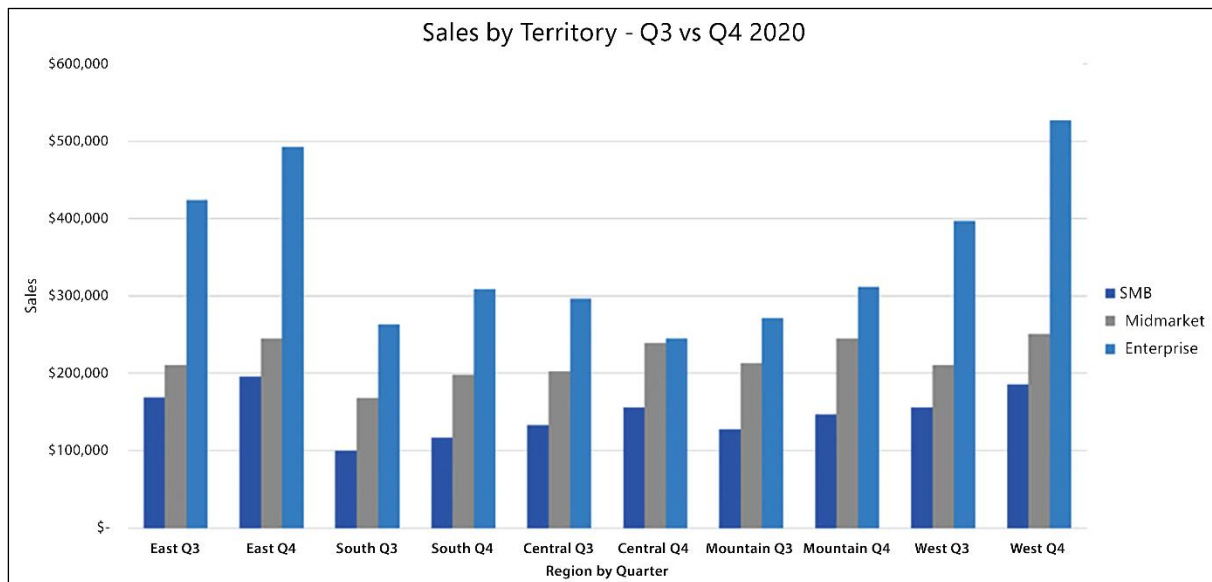
	Column name	Data type	Partition key	Comment
1	emp_id	bigint		
2	last_name	string		
3	first_name	string		
4	hire_date	bigint		
5	street_address	string		







Sales Data by Territory and Segment			
Territory	SMB	Midmarket	Enterprise
East Q3	\$ 168,778	\$ 210,696	\$ 423,875
East Q4	\$ 196,254	\$ 244,995	\$ 492,878
South Q3	\$ 99,361	\$ 168,572	\$ 263,119
South Q4	\$ 116,895	\$ 198,320	\$ 309,552
Central Q3	\$ 132,882	\$ 203,082	\$ 296,332
Central Q4	\$ 156,332	\$ 238,920	\$ 245,000
Mountain Q3	\$ 127,699	\$ 213,247	\$ 271,440
Mountain Q4	\$ 146,780	\$ 245,112	\$ 312,000
West Q3	\$ 156,147	\$ 210,558	\$ 396,885
West Q4	\$ 185,889	\$ 250,664	\$ 526,995



Create layer

Layer configuration

Name

awsSDKpandas219_python39

Description - *optional*

AWS SDK for Pandas, v2.19.0 for Python 3.9

- ☒ Upload a .zip file
☐ Upload a file from Amazon S3

 Upload

aws wrangler-layer-2.19.0-py3.9.zip
49.75 MB



For files larger than 10 MB, consider uploading using Amazon S3.

Compatible architectures - *optional* [Info](#)

Choose the compatible instruction set architectures for your layer.

- ☐ x86_64
☐ arm64

Compatible runtimes - *optional* [Info](#)

Choose up to 15 runtimes.

Runtimes



Python 3.9



License - *optional* [Info](#)

Function name

Enter a name that describes the purpose of your function.

CSVtoParquetLambda

Use only letters, numbers, hyphens, or underscores with no spaces.

Runtime

Info

Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

Python 3.9

Architecture

Info

Choose the instruction set architecture you want for your function code.

☒ x86_64

☐ arm64

Permissions

Info

By default, Lambda will create an execution role with permissions to upload logs to Amazon CloudWatch Logs. You can customize this default role later when adding triggers.

▼ Change default execution role

Execution role

Choose a role that defines the permissions of your function. To create a custom role, go to the IAM console.

☐ Create a new role with basic Lambda permissions

☒ Use an existing role

☐ Create a new role from AWS policy templates

Existing role

Choose an existing role that you've created to be used with this Lambda function. The role must have permission to upload logs to Amazon CloudWatch Logs.

DataEngLambdaS3CWGlueRole

View the DataEngLambdaS3CWGlueRole role on the IAM console.

Lambda > Layers > Add layer

Add layer

Function runtime settings

Runtime

Python 3.9

Architecture

x86_64

Choose a layer

Layer source

Info

Choose from layers with a compatible runtime and instruction set architecture or specify the Amazon Resource Name (ARN) of a layer version. You can also create a new layer.

☐ AWS layers

Choose a layer from a list of layers provided by AWS.

☒ Custom layers

Choose a layer from a list of layers created by your AWS account or organization.

☐ Specify an ARN

Specify a layer by providing the ARN.

Custom layers

Layers created by your AWS account or organization that are compatible with your function's runtime.

awsSDKpandas219_python39

Version

1

Cancel

Add

Add trigger

Trigger configuration

 **S3**
aws storage

Bucket

Please select the S3 bucket that serves as the event source. The bucket must be in the same region as the function.

dataeng-landing-zone-



Event type

Select the events that you want to have trigger the Lambda function. You can optionally set up a prefix or suffix for an event. However, for each bucket, individual events cannot have multiple configurations with overlapping prefixes or suffixes that could match the same object key.

All object create events

Prefix - *optional*

Enter a single optional prefix to limit the notifications to objects with keys that start with matching characters.

e.g. images/

Suffix - *optional*

Enter a single optional suffix to limit the notifications to objects with keys that end with matching characters.

.csv

Lambda will add the necessary permissions for Amazon S3 to invoke your Lambda function from this trigger. [Learn more](#) about the Lambda permissions model.



Recursive invocation

If your function writes objects to an S3 bucket, ensure that you are using different S3 buckets for input and output. Writing to the same bucket increases the risk of creating a recursive invocation, which can result in increased Lambda usage and increased costs. [Learn more](#)

- ☒ I acknowledge that using the same S3 bucket for both input and output is not recommended and that this configuration can cause recursive invocations, increased Lambda usage, and increased costs.

Cancel

Add

Chapter 4: Data Governance, Security, and Cataloging

AWS Glue > Tables > csvtoparquet

csvtoparquet

Last updated (UTC)
March 26, 2023 at 14:07:08

Version 1 (Current version)

Actions

Table overview

Data quality New

Table details

Advanced properties

Name
csvtoparquet

Description
-

Database
cleanzonedb

Classification
parquet

Location
s3://dataeng-clean-zone-gse23/cleanzonedb/csvtoparquet

Connection
-

Deprecated
-

Last updated
March 26, 2023 at 14:07:08

Input format
org.apache.hadoop.hive.qLio.parquet.Ma
predParquetInputFormat

Output format
org.apache.hadoop.hive.qLio.parquet.Ma
predParquetOutputFormat

Serde serialization lib
org.apache.hadoop.hive.qLio.parquet.ser
de.ParquetHiveSerDe

Schema

Partitions

Indexes

Schema (2)

View and manage the table schema.

Filter schemas

< 1 > ⚙

#	Column name	Data type	Partition key	Comment
1	name	string	-	-
2	favorite_num	int	-	-

AWS Lake Formation > Tables > csvtoparquet

csvtoparquet

Version 1 (Current version)

Actions

Compare versions

Drop table

View properties

Table details

Edit table

Database
cleanzonedb

Description
-

Governance
Disabled

Location
s3://dataeng-clean-zone-gse23/cleanzonedb/csvtoparquet

Data format
parquet

Compaction Status
-

Connection
-

Last updated
Sunday, March 26, 2023 at 2:07 PM UTC

▶ Advanced table properties

Schema

Edit schema

Find Columns

< 1 > ⚙

#	Column Name	Data type	Partition key	Comment	LF-Tags
1	name	string	-	-	-
2	favorite_num	int	-	-	-

Create policy

1

2

3

A policy defines the AWS permissions that you can assign to a user, group, or role. You can create and edit a policy in the visual editor and using JSON. [Learn more](#)

Visual editor

JSON

[Import managed policy](#)

```
28     "glue:CreatePartition",
29     "glue:DeletePartition",
30     "glue:BatchDeletePartition",
31     "glue:UpdatePartition",
32     "glue:GetPartition",
33     "glue:GetPartitions",
34     "glue:BatchGetPartition"
35 ],
36 "Resource": [
37     "arn:aws:glue::*:catalog",
38     "arn:aws:glue::*:database/cleanzonedb",
39     "arn:aws:glue::*:database/cleanzonedb*",
40     "arn:aws:glue::*:table/cleanzonedb/*"
41 ]
42 },
43 }
```

Security: 0 Errors: 0 Warnings: 0 Suggestions: 1

Visual editor

JSON

[Import managed policy](#)

```
39     "arn:aws:glue::*:database/cleanzonedb",
40     "arn:aws:glue::*:table/cleanzonedb/*"
41 ]
42 },
43 {
44     "Effect": "Allow",
45     "Action": [
46         "s3:GetBucketLocation",
47         "s3:GetObject",
48         "s3:ListBucket",
49         "s3:ListBucketMultipartUploads",
50         "s3:ListMultipartUploadParts",
51         "s3:AbortMultipartUpload",
52         "s3:PutObject"
53     ],
54     "Resource": [
55         "arn:aws:s3:::dataeng-clean-zone-gse23/*"
56     ]
57 },
58 {
59     "Effect": "Allow",
60     "Action": [
61         "s3:GetBucketLocation",
```

Security: 0 Errors: 0 Warnings: 0 Suggestions: 1

Welcome to Lake Formation

The first step in creating your data lake in Lake Formation is defining one or more administrators. Administrators have full access to the Lake Formation console, and control the initial data configuration and access permissions.

Choose the initial administrative users and roles

You may add yourself and/or other principals.

☒ Add myself

AWS account: 42[REDACTED]30

☐ Add other AWS users or roles

Select additional IAM users and roles to be data lake administrators.

Cancel

Get started

AWS Lake Formation

Dashboard

Data catalog

Databases

Tables

Data filters

Data sharing New

Settings 1

Register and ingest

Data lake locations

Blueprints

Crawlers 2

Jobs 2

Permissions

Administrative roles and tasks

LF-Tags

LF-tag permissions

Data lake permissions

AWS Lake Formation > Permissions

Data permissions for database cleanzonedb (2)

Revoke

Grant

Filter permissions by property or value

< 1 ... >

	Principal	Principal type	Resource type	Database	Table	Resource
<input type="radio"/>	DataEngLambdaS3CWGlueRole	IAM role	Database	cleanzonedb	-	cleanzonedb
<input type="radio"/>	IAMAllowedPrincipals	Group	Database	cleanzonedb	-	cleanzonedb

Amazon Athena > Query editor

Editor

Recent queries

Saved queries

Settings

Workgroupprimary

Data

↺

↻

Data source

AwsDataCatalog

Database

cleanzonedb

Tables and views

Create

⚙

Filter tables and views

▼ Tables (1)

↺ 1 ↻

📦 csvtoparquet

⋮

► Views (0)

↺ 1 ↻

Query 1

⋮

+

▼

1

select * from cleanzonedb.csvtoparquet

SQL

Ln 1, Col 1

⌵

⌴

⚙

Run again

Explain

Cancel

Clear

Create

Reuse query results

up to 60 minutes ago

Query results

Query stats

Completed

Time in queue: 102 ms

Run time: 357 ms

Data scanned: 0.19 KB

Results (8)

Copy

Download results

Search rows

↺ 1 ↻ ⚙

#

▼

name

▼

1

Gareth

2

Tracy

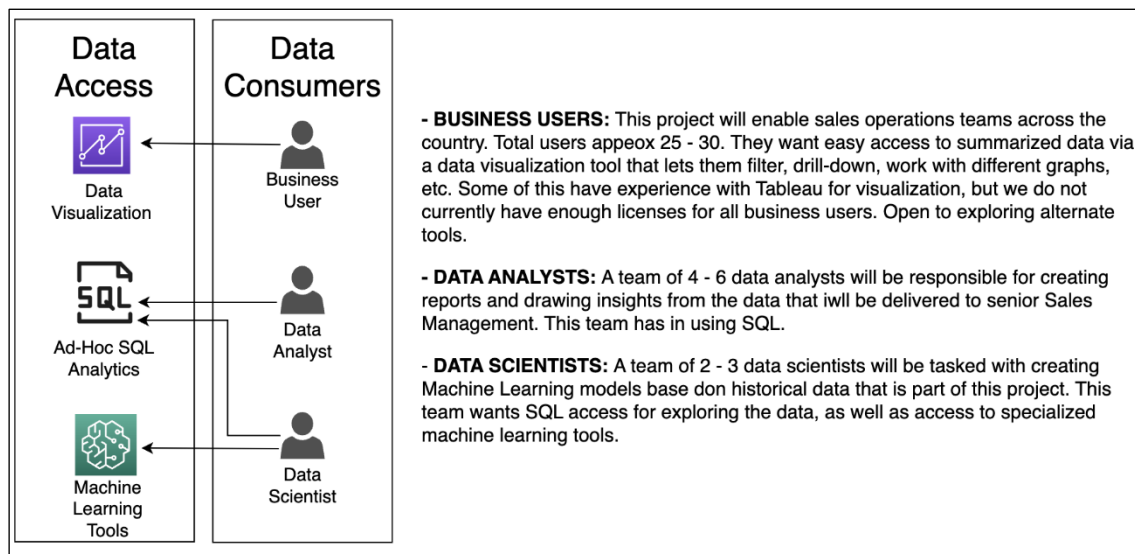
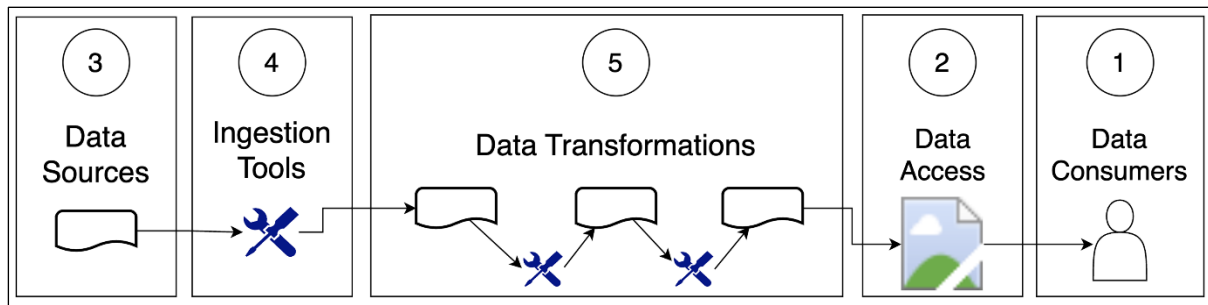
3

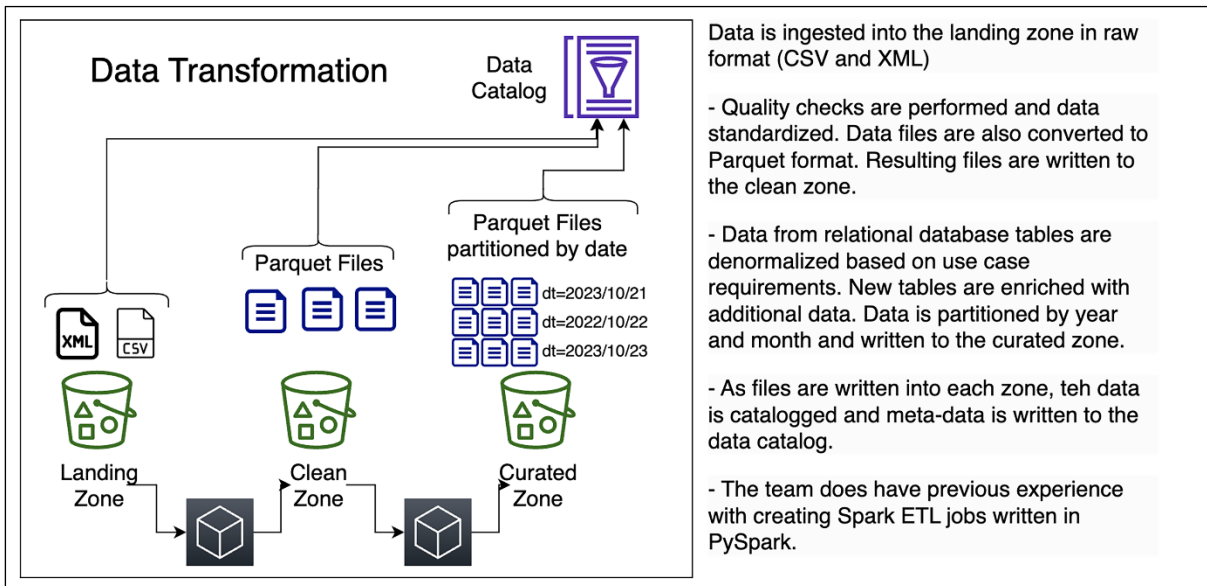
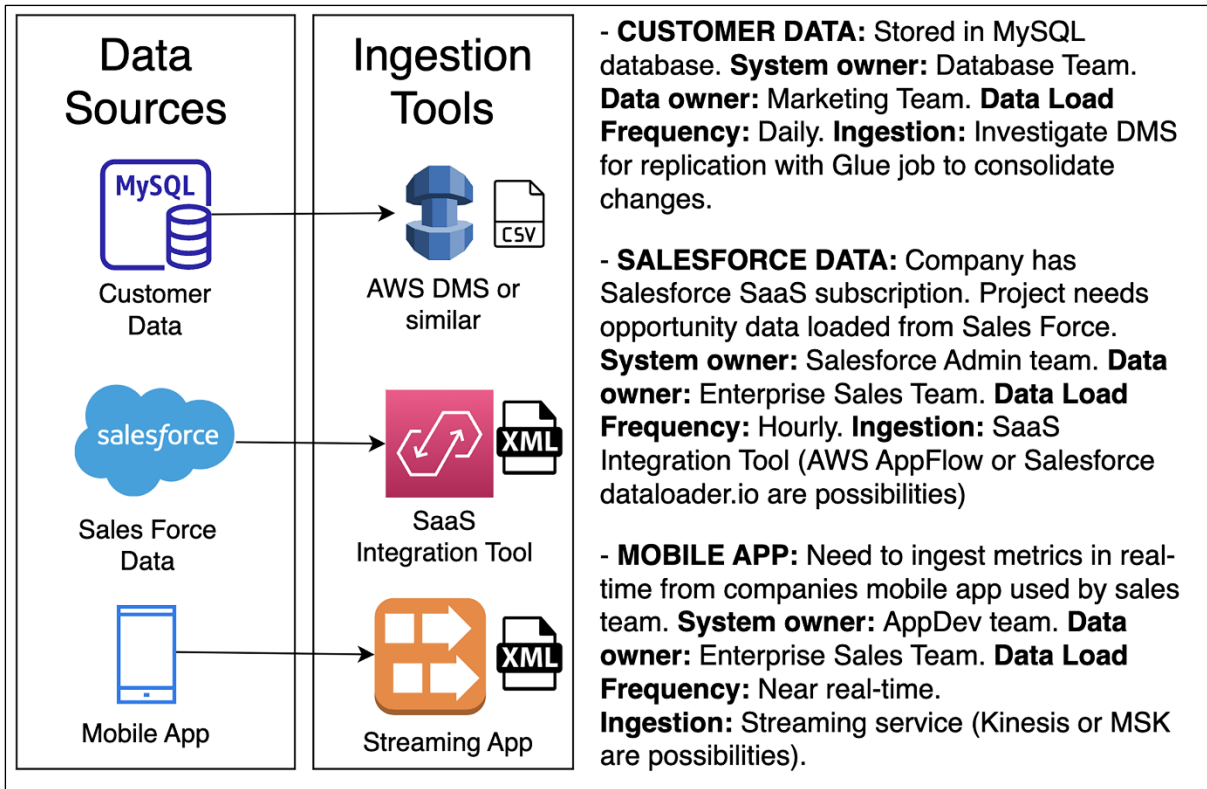
Chris

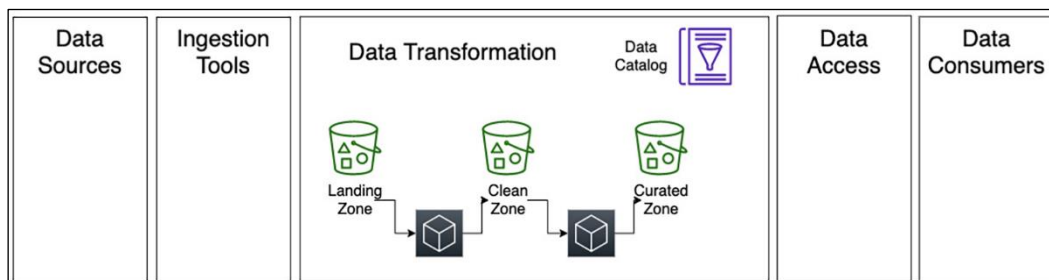
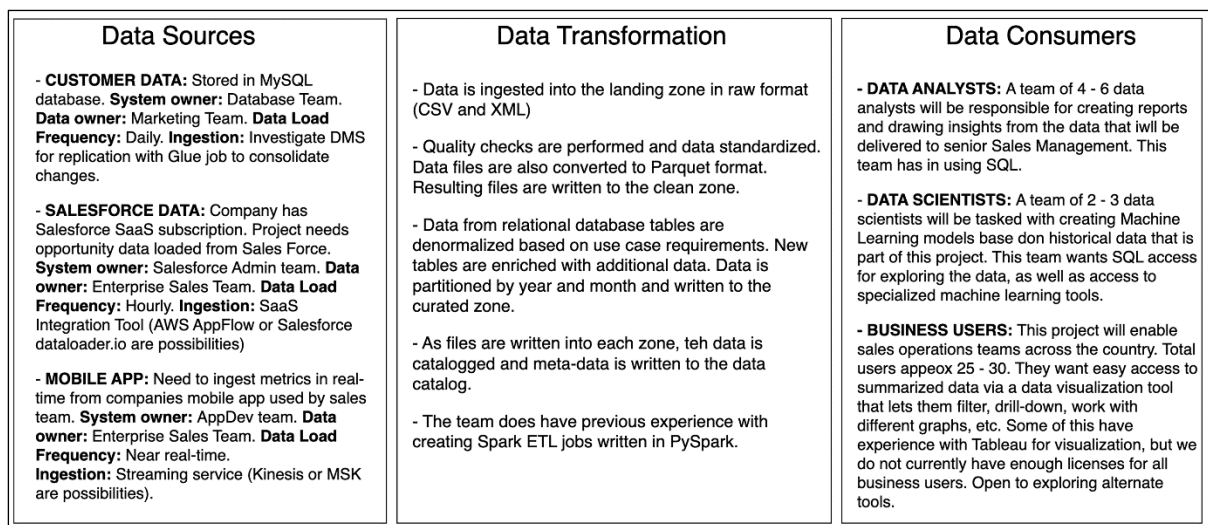
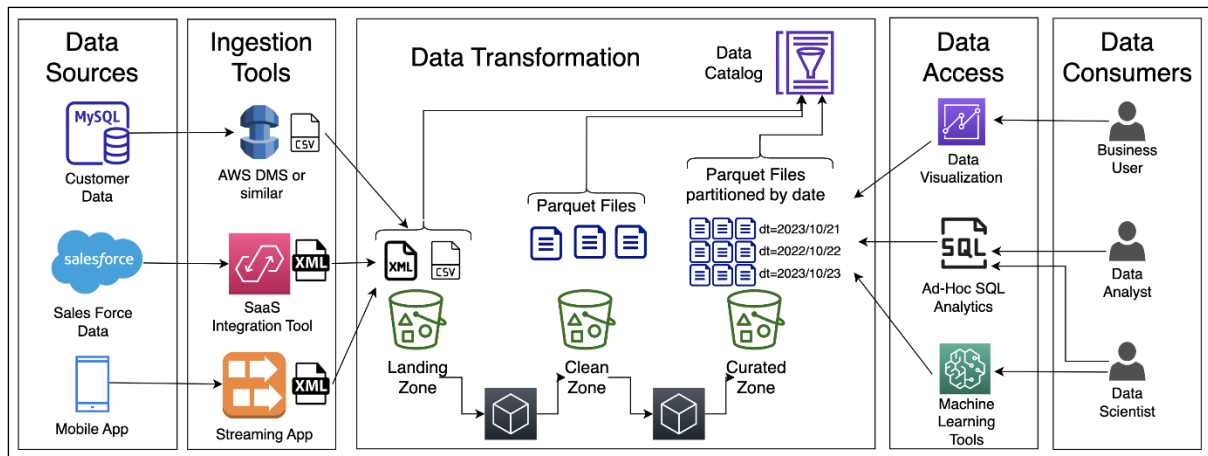
4

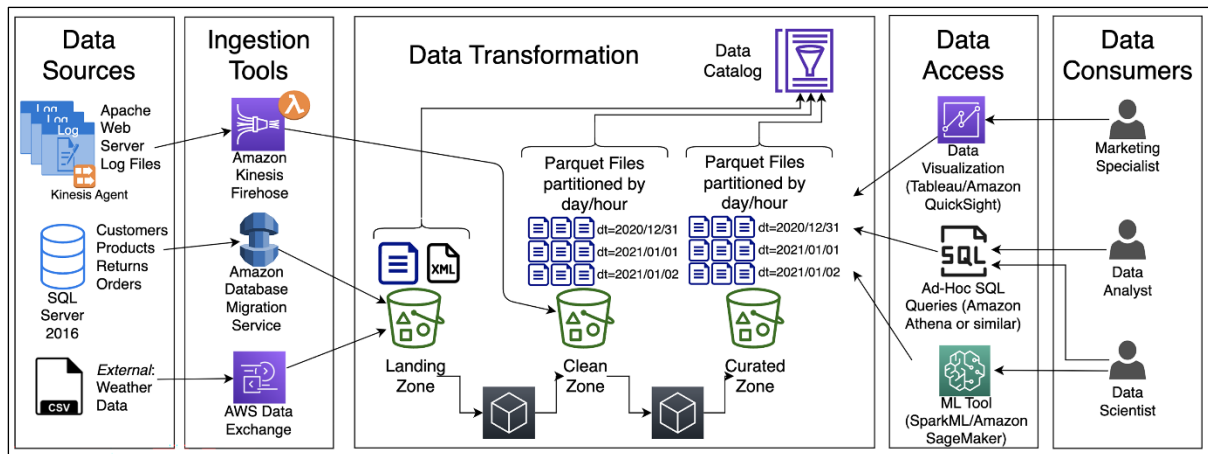
Emma

Chapter 5: Architecting Data Engineering Pipelines









Data Sources	Data Transformation	Data Consumers
<p>- Apache Web Server Log Files: From 4 Apache web servers. System Owner: Natalie Rabinovich. Data Owner: Marketing. Ingestion: Could use Kinesis Agent to transform to JSON and send to Kinesis Firehose. Firehose does validation (using Lambda function) and transforms to Parquet format. Could write direct to clean zone, partitioned by day (yyyy/mm/dd).</p> <p>- Databases: Customers, Products, Returns, Orders on SQL Server 2016 Enterprise Edition. System Owner: Owen McClave. Data Owner: Sales Team. Potentially use Amazon DMS to replicate to Amazon S3 raw zone in Parquet format.</p> <p>- Weather Data: External data source available via subscription. Data Owner: Marketing. Ingestion: Available from AWS Data Exchange marketplace. Lambda function can load data into Amazon S3 raw zone when available.</p>	<p>- Raw Zone: Database and weather data replicated into raw zone. When files ingested triggers Lambda function to perform data quality checks and then loads into Clean Zone partitioned by yyyy/mm/dd.</p> <p>- Clean Zone: Web server log files loaded directly into clean zone after Kinesis Firehose uses a Lambda function to perform data quality checks. Firehose configured to write to clean zone partitioned by yyyy/mm/dd. Database and weather files loaded from raw zone after data quality checks, and partition by yyyy/mm/dd.</p> <p>- Curated Zone: Database files denormalized, enriched (with weather data potentially), other business logic added. Partitioned by either day (databases, weather) or hour (web server log files)</p>	<p>- Marketing Specialists: Want to use business intelligence (visualization) tool to view up-to-date website analytics (ad-campaign referrals, coupon redemption, heatmap showing activity by geographic location). Refresh on at least hourly basis. Analytics team generally uses Tableau, but marketing team does not have licenses. Open to other BI tools.</p> <p>- Data Analysts: Responsible for creating reports and insights using SQL queries. Database and weather data could be refreshed daily, but they would need web server clickstream log files refreshed at least hourly.</p> <p>- Data Scientists: Need ad-hoc SQL access to databases, weather and web server log files. They currently use SparkML on-premises, but open to new cloud based tools that may make speed up delivery and collaboration for their machine learning products.</p>

Chapter 6: Ingesting Batch and Streaming Data

Food_Code	Display_Name	Portion_Display_Name	Total Calo
71411000	Potato skin with cheese & bacon	order (10 halves)	1667.4
24301010	Roasted duck	duck half	1283.52
21103120	Breaded fried steak (eat lean & fat)	large steak	1069.2
28141010	Fried chicken frozen meal	large meal (16 oz)	1024.92
27347100	Chicken or turkey pot pie	16-ounce pie (Hungry Man)	976.1
58200100	Wrap sandwich (meat, vegetables, rice)	wrap	818.37
21103120	Breaded fried steak (eat lean & fat)	medium steak	801.9
58106730	Meat & veggie pizza, thick crust	small pizza (8" across)	798.64
24401010	Roasted Cornish game hen	hen	792.54
58106530	Meat pizza, thick crust	small pizza (8" across)	785.4

53 lines (45 sloc) | 1.67 KB

RawBlame

```
1 ----
2 AWSTemplateFormatVersion: 2010-09-09
3 Description: Chapter 6 - Data Engineering with AWS
4 Parameters:
5   DBPassword:
6     Type: String
7     NoEcho: true
8     Description: The database admin account password
9     MinLength: 8
10    AllowedPattern: ^[a-zA-Z0-9]*$
11    ConstraintDescription: Password must contain only alphanumeric characters.
12  LatestAmiId:
13    Type: 'AWS::SSM::Parameter::Value<AWS::EC2::Image::Id>'
14    Default: '/aws/service/ami-amazon-linux-latest/al2023-ami-kernel-6.1-x86_64'
15 Resources:
16   MySQLInstance:
```

CloudFormation > Stacks > dataeng-aws-chapter6-mysql-ec2

Stacks (1)

Filter by stack name

Active View nested

Stacks

dataeng-aws-chapter6-mysql-ec2

2023-04-09 21:55:11 UTC-0400

CREATE_COMPLETE

dataeng-aws-chapter6-mysql-ec2

DeleteUpdateStack actionsCreate stack

Stack infoEventsResourcesOutputsParametersTemplateChange sets

Events (7)

Search events

Timestamp	Logical ID	Status	Status reason
2023-04-09 22:03:01 UTC-0400	EC2Instance	CREATE_COMPLETE	-
2023-04-09 22:02:54 UTC-0400	EC2Instance	CREATE_IN_PROGRESS	Resource creation Initiated
2023-04-09 22:02:52 UTC-0400	EC2Instance	CREATE_IN_PROGRESS	-
2023-04-09 22:02:48 UTC-0400	MySQLInstance	CREATE_COMPLETE	-
2023-04-09 21:55:19	MySQLInstance	CREATE_IN_PROGRESS	Resource creation

Review and create [Info](#)

Review the permissions, specify details, and tags.

Policy details

Policy name

Enter a meaningful name to identify this policy.

DataEngDMSLandingS3BucketPolicy

Maximum 128 characters. Use alphanumeric and '+=, @-_' characters.

Description - optional

Add a short explanation for this policy.

Maximum 1,000 characters. Use alphanumeric and '+=, @-_' characters.

i This policy defines some actions, resources, or conditions that do not provide permissions. To grant access, policies must have an action that has an applicable resource or condition. For details, choose **Show remaining**. [Learn more](#)

Permissions defined in this policy [Info](#)

Permissions defined in this policy document specify which actions are allowed or denied. To define permissions for an IAM identity (user, user group, or role), attach a policy to it

Edit

Search

Allow (1 of 384 services)

Show remaining 383 services

Service	Access level	Resource	Request condition
S3	Limited: Read, List, Permissions management, Tagging, Write	Multiple	None

Use case

Allow an AWS service like EC2, Lambda, or others to perform actions in this account.

Common use cases

- ☐ EC2
Allows EC2 instances to call AWS services on your behalf.
- ☐ Lambda
Allows Lambda functions to call AWS services on your behalf.

Use cases for other AWS services:

DMS

- ☒ DMS
Allows Database Migration Service to call AWS services on your behalf.

Cancel

Next

IAM > Roles > DataEngDMSLandingS3BucketRole

DataEngDMSLandingS3BucketRole

Delete

Allows Database Migration Service to call AWS services on your behalf.

Summary

Edit

Creation date

April 09, 2023, 22:25 (UTC-04:00)

ARN

arn:aws:iam::420240645590:role/DataEngDMSLandingS3BucketRole

Last activity

None

Maximum session duration

1 hour

Permissions

Trust relationships

Tags

Access Advisor

Revoke sessions

Endpoint configuration

Endpoint identifier [Info](#)

A label for the endpoint to help you identify it.

s3-landing-zone-sakila-csv

Descriptive Amazon Resource Name (ARN) - optional

A friendly name to override the default DMS ARN. You cannot modify it after creation.

Friendly-ARN-name

Target engine

The type of database engine this endpoint is connected to.

Amazon S3

Service access role ARN

Role that can access target

arn:aws:iam::26:role/DataEngDMSLandingS3BucketRole

Bucket name

The name of an Amazon S3 bucket where DMS will read the files from

dataeng-landing-zone

Bucket folder

The Amazon S3 bucket path where the CSV files can be found

sakila-db

▼ Endpoint settings

Define additional specific settings for your endpoints using wizard or editor. [Learn more](#)

☒ Wizard

Enter endpoint settings using the guided user interface.

☐ Editor

Enter endpoint settings in JSON format.

Endpoint settings

Setting

Value - A value is required

Q AddColumnName

X

Q True

X

Remove

Add new setting

☐ Use endpoint connection attributes

Chapter 7: Transforming Data to Optimize for Analytics

Customer_ID	Last_Name	First_Name	Address_Street	Address_City	Address_State	Phone_Number	Sales_Person_ID
1	Smith	Jonathan	123 Main Street	Springville	MA	555-943-1987	2
2	Mendez	Bruno	5449 South West Street	Jersey	PA	555-615-1609	3
3	Sachdeva	Viyoma	94 Midland Avenue	Oxford	NJ	555-664-0464	1

Sales_Person_ID	Last_Name	First_Name	Territory_Code
1	Taylor	Chris	95
2	Williams	Carmen	42
3	Kelly	Michael	23

Customer_ID	Last_Name	First_Name	Address_Street	Address_City	Address_State	Phone_Number	Sales_Person_Last	Sales_Person_First
1	Smith	Jonathan	123 Main Street	Springville	MA	555-943-1987	Williams	Carmen
2	Mendez	Bruno	5449 South West Street	Jersey	PA	555-615-1609	Kelly	Michael
3	Sachdeva	Viyoma	94 Midland Avenue	Oxford	NJ	555-664-0464	Taylor	Chris

Untitled job

Job has not been saved Try new UI

Actions Save Run

Visual Script Job details Runs Data quality New Schedules Version Control

+
Data source - S3 bucket
S3 – Film-Category

Data source properties - S3 Output schema Data preview

Name

S3 source type Info

- ☐ S3 location
Choose a file or folder in an S3 bucket.
- ☒ Data Catalog table

Database
Choose a database.

► Use runtime parameters

Table

► Use runtime parameters

Untitled job

Job has not been savedTry new UI

ActionsSaveRun

VisualScriptJob detailsRunsData quality NewSchedulesVersion Control

+

Data source - S3 bucket

S3 - Film-Category

Data source - S3 bucket

S3 - Film

Data source properties - S3

Output schema

Data preview

Name

S3 - Film

S3 source type

Info

☐ S3 location

Choose a file or folder in an S3 bucket.

☒ Data Catalog table

Database

Choose a database.

sakila

► Use runtime parameters

Table

film

► Use runtime parameters

Untitled job

Job has not been savedTry new UI

ActionsSaveRun

VisualScriptJob detailsRunsData quality NewSchedulesVersion Control

+

Data source - S3 bucket

S3 - Film-Category

Data source - S3 bucket

S3 - Film

Transform - ApplyMapping

Renamed keys for Join

Transform - Join

Join

Transform

Output schema

Data preview

Name

Join

Node parents

Choose which nodes will provide inputs for this one.

Choose one or more parent node

S3 - Film

Renamed keys for Join

Join type

Select the type of join to perform.

Left join

Select all rows from the left dataset and the rows that meet the join condition from the rig...

Join conditions

Select a field from each parent node for the join condition.

S3 - Film

Renamed keys for Join

film_id

=

right_film_id

Data target properties - S3

Output schema

Data preview

Parquet

After you save your job, it will use Glue Studio's optimized Parquet writer.

Compression Type

Snappy

S3 Target Location

Choose an S3 location in the format s3://bucket/prefix/object/ with a trailing slash (/).

s3://dataeng-curved-zone-gse23/streaming/streaming

View

Browse S3

Data Catalog update options

[Info](#)

Choose how you want to update the Data Catalog table's schema and partitions. These options will only apply if the Data Catalog table is an S3 backed source.

☐ Do not update the Data Catalog

☒ Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions

☐ Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions

Database

Choose the database from the AWS Glue Data Catalog.

curatedzonedb

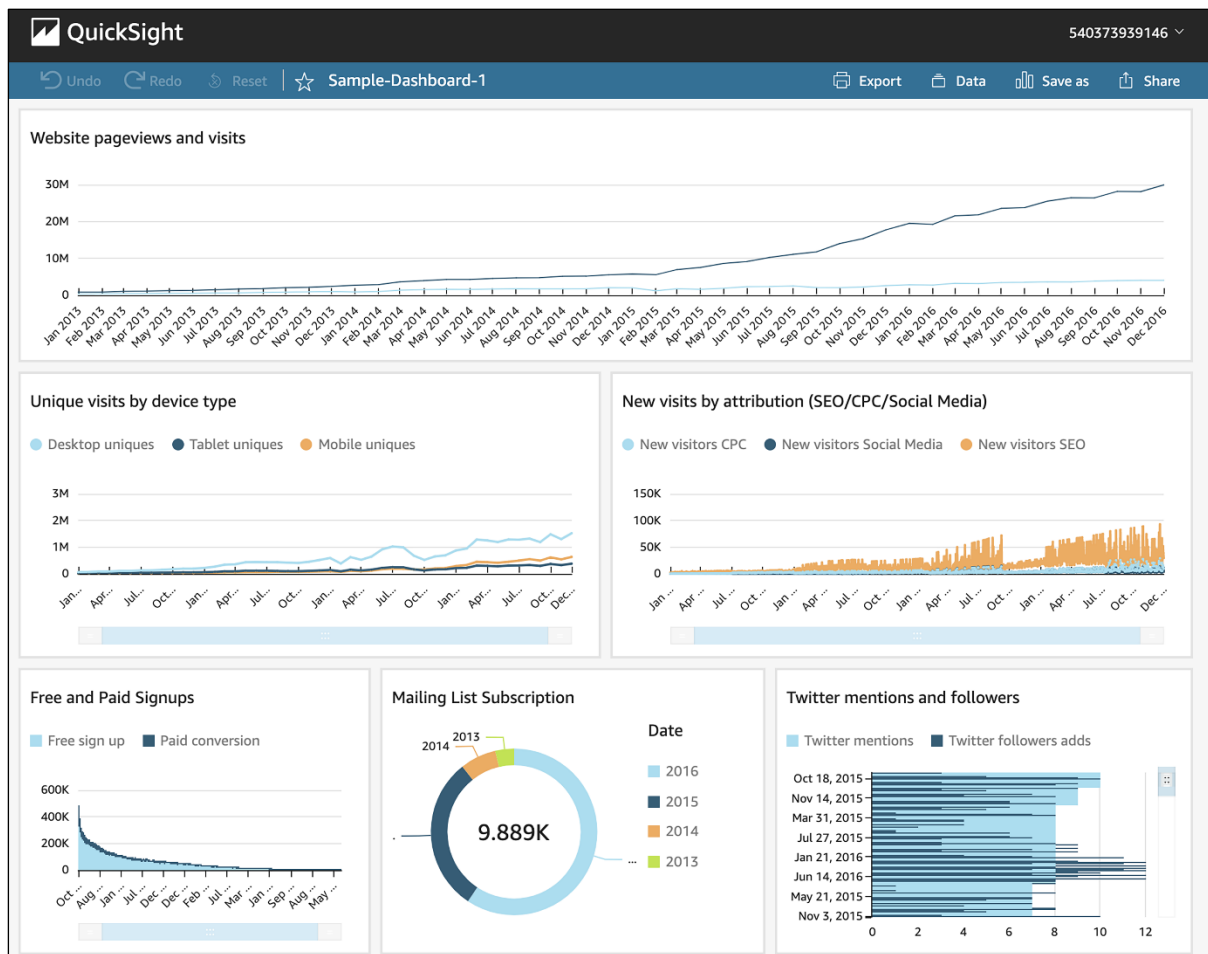
Use runtime parameters

Table name

Enter a table name for the AWS Glue Data Catalog.

streaming_films

Chapter 8: Identifying and Enabling Data Consumers



demo-project

Dataset: streaming-films | Sample: First n sample (500 rows)

Viewing 19 columns | 500 rows

timestamp	eventtype	# film_id_streaming	distributor
2021-02-15T00:03:29-...	trailer	6	vudo
2021-02-15T00:06:07-...	rent	6	fandango now
2021-02-15T00:05:14-...	buy	5	microsoft
All other values		483	
2021-02-15T00:04:08-05:00	rent	9	google play
2021-02-15T00:04:56-05:00	buy	9	google play
2021-02-15T00:05:54-05:00	trailer	9	apple itunes
2021-02-15T00:06:27-05:00	trailer	9	amazon prime
2021-02-15T00:03:14-05:00	trailer	11	youtube
2021-02-15T00:05:58-05:00	trailer	11	apple itunes
2021-02-15T00:03:02-05:00	trailer	32	microsoft
2021-02-15T00:04:16-05:00	rent	32	fandango now
2021-02-15T00:05:40-05:00	trailer	32	youtube
2021-02-15T00:07:34-05:00	...	40	...

Zoom: 130%

DATASETS

PROJECTS

RECIPES

DQ RULES

JOBS

WHAT'S NEW

Dataset name

customer-dataset

The dataset name must contain 1-255 characters. Valid characters are alphanumeric (A-Z, a-z, 0-9), hyphen (-), period (.), and space.

Connect to new dataset

File upload

Data lake/data store

Amazon S3

Database connections

Amazon Redshift

JDBC

AWS Glue Data Catalog

Data Catalog S3 tables

Data Catalog Redshift tables

Data Catalog RDS tables

All AWS Glue tables

Others

Amazon AppFlow

AWS Data Exchange

External data connections

Snowflake

AWS Account

Current AWS account

420240645590

Another AWS account

Your source from Data Catalog

Permission from AWS Lake Formation will apply to datasets with this icon.

AWS Glue databases

Search databases by name

Database name	Description	Created on
cleanzonedb		2 months ago February 28, 2023, 10:16:19 pm
curatedzonedb		12 days ago April 17, 2023, 9:11:49 pm
sakila		20 days ago April 9, 2023, 10:57:56 pm
streaming_db		18 days ago April 11, 2023, 9:22:23 pm

DATASETS

PROJECTS

RECIPES

DQ RULES

JOBS

WHAT'S NEW

DataBrew > Projects > Create project

Create project

Project details

Project name

customer-mailing-list

The project name must contain 1-255 characters. Valid characters are alphanumeric (A-Z, a-z, 0-9), hyphen (-), period (.), and space.

Recipe details

Attached recipe

Create new recipe

Recipe name

customer-mailing-list-recipe

The recipe name must contain 1-255 characters. Valid characters are alphanumeric (A-Z, a-z, 0-9), hyphen (-), period (.), and space.

Import steps from recipe

Import recipe steps from an existing recipe into your project. The existing recipe that you chose will not be edited.

Select a dataset

Select the dataset that you want to work on

My datasets

Your imported datasets

Sample files

Explore example files for your dataset

New dataset

Import new dataset

Find datasets

Dataset name	Data type	Source	Create date
address-dataset	Data Catalog table	Data Catalog	5 days ago April 24, 2023, 10:48:44 pm
customer-dataset	Data Catalog table	Data Catalog	5 days ago April 24, 2023, 10:47:53 pm

Recipe (4)

customer-mailing-list-recipe

Working version

Publish

More

Applied steps (4) | [Clear all](#)

1. Left join address-dataset

2. Change format of first_name to Capital case

3. Change format of last_name to Capital case

4. Change format of email to Lowercase

DATASETS

PROJECTS

RECIPES

JOBS

WHAT'S NEW

Created recipe job "mailing-list-job".

[DataBrew](#) > [Jobs](#) > mailing-list-job

mailing-list-job

Dataset: customer-dataset

Project: customer-mailing-list

Recipe: customer-mailing-list-recipe

▶ Run job

Actions ▼

OPEN PROJECT

Job run history

Job details

Data lineage

4 Recipe

Last job run 4 minutes, no job runs scheduled

Job run history

Stop job run

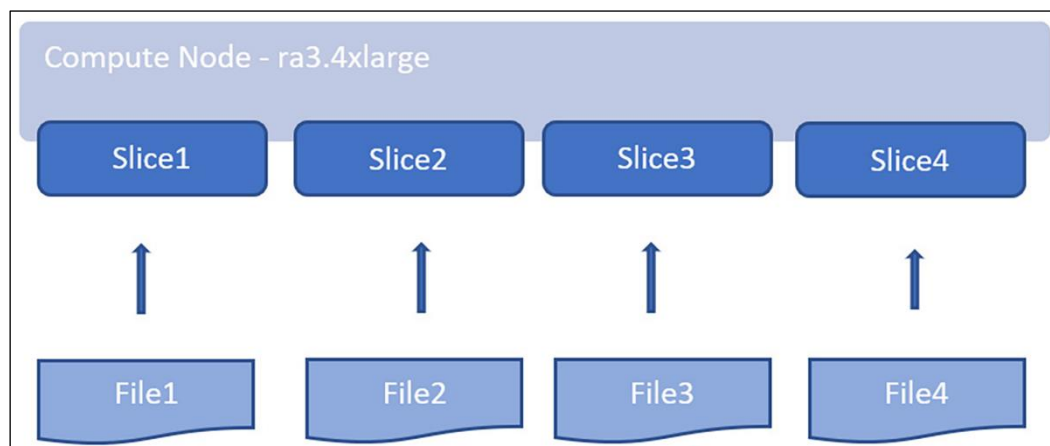
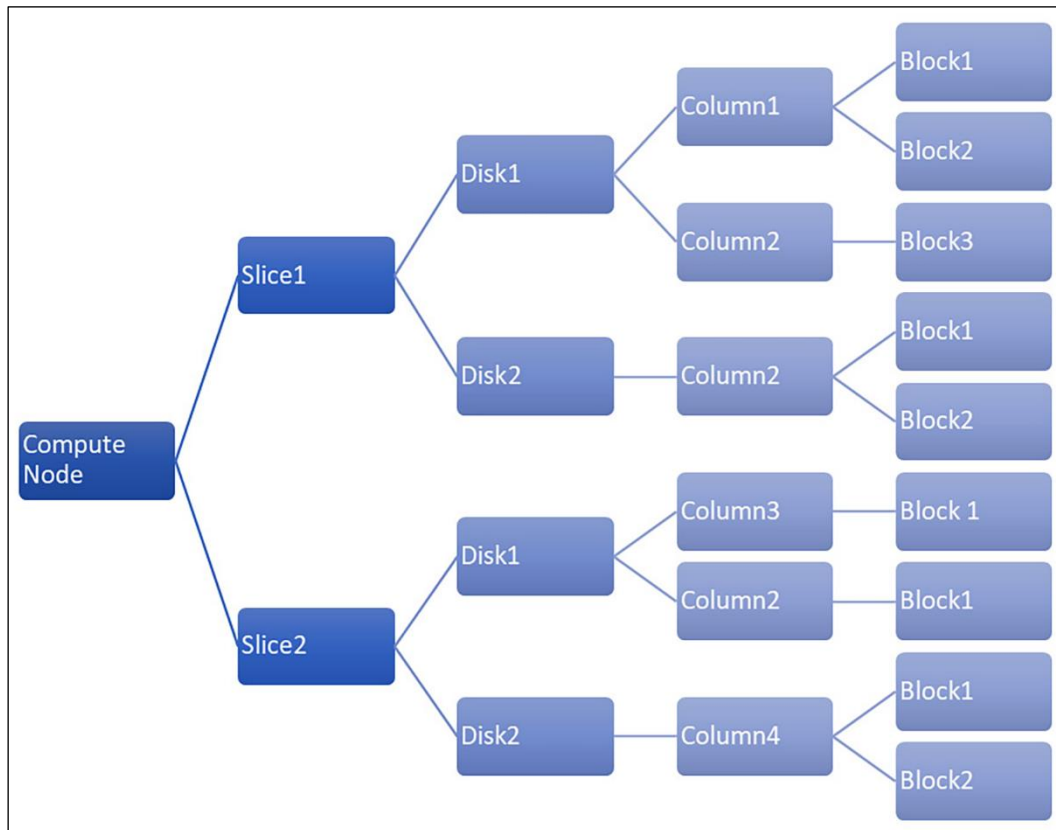
Actions ▼

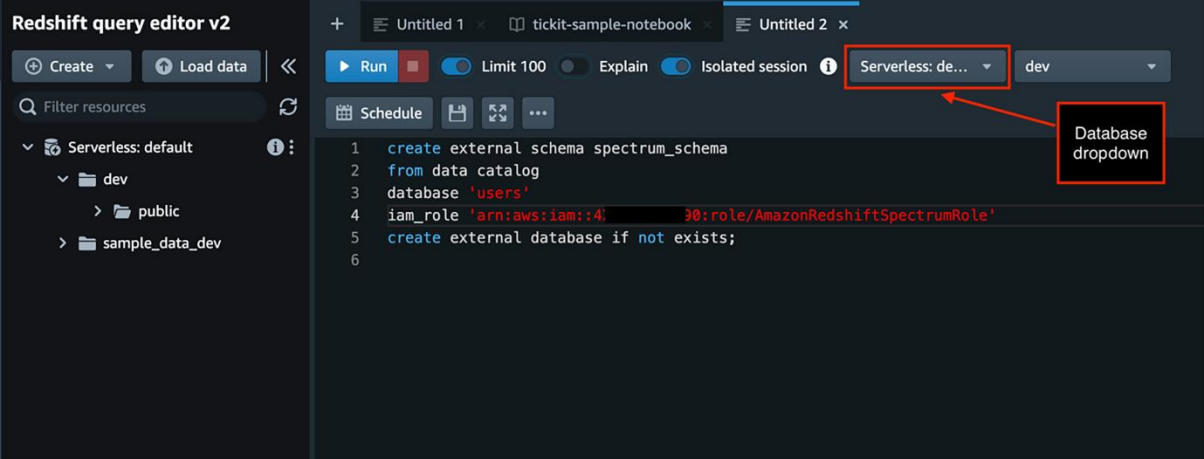
Show all ▼

< 1 >

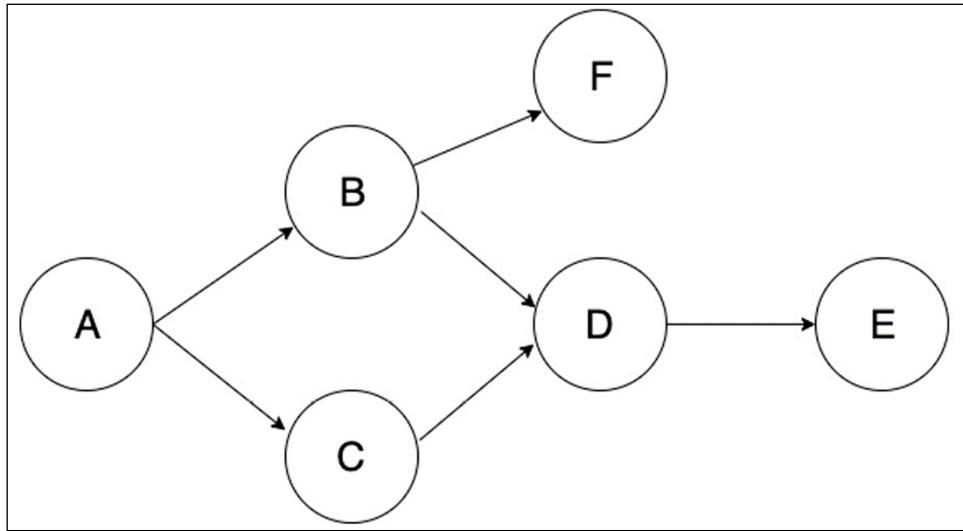
Job run ID	Last job run status	Run time	Output	Summary
mailing-list-job_2023-04-29-20:18:43	Succeeded	1 minute, 24 seconds	1 output	

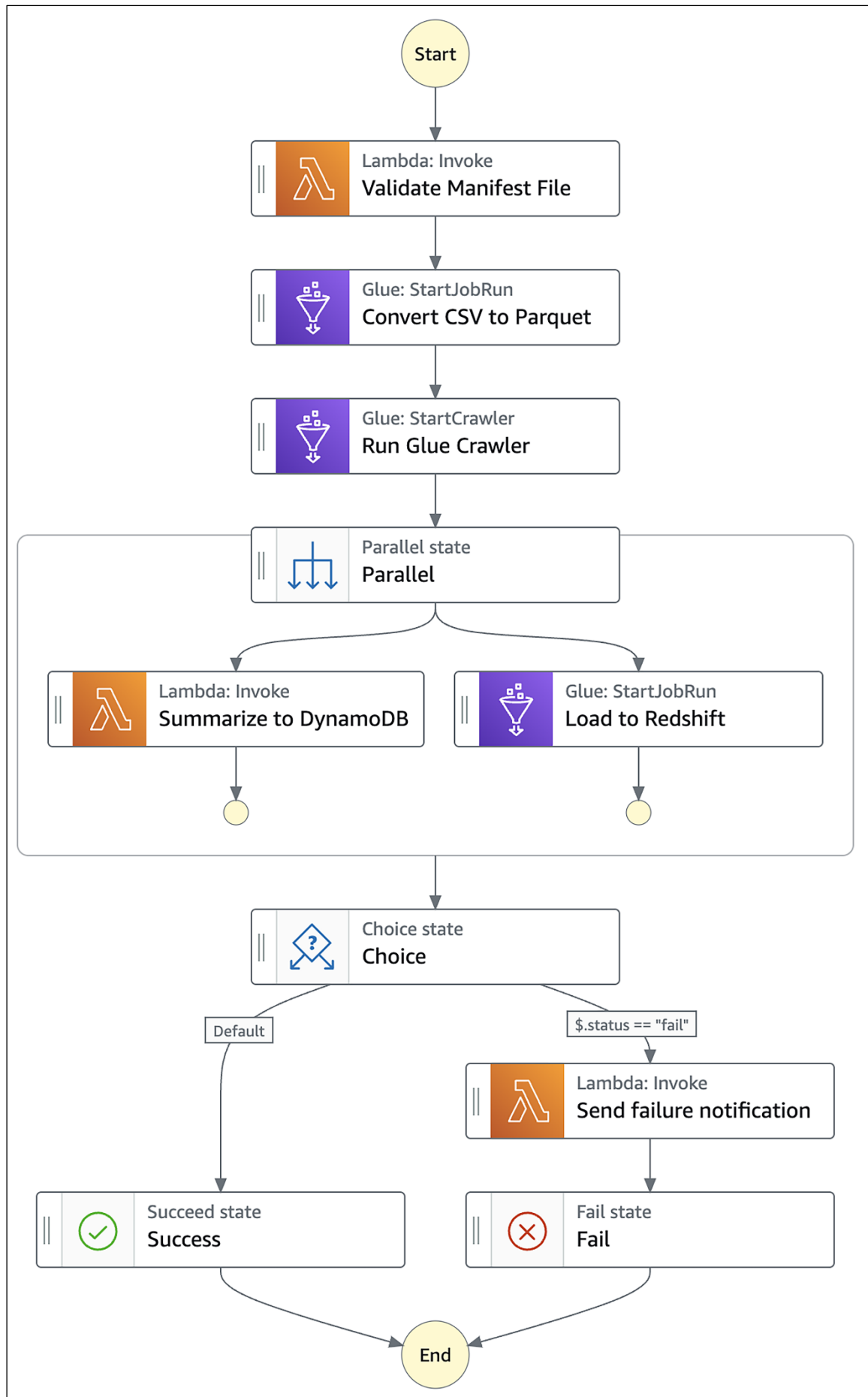
Chapter 9: A Deeper Dive into Data Marts and Amazon Redshift





Chapter 10: Orchestrating the Data Pipeline





Criteria	AWS Step Functions	Amazon Managed Workflows for Apache Airflow (MWAA)
Short description	Serverless AWS native orchestration service	Managed AWS service for open source Apache Airflow
Graphical pipeline development	Yes	No
Graphical run visualization	Yes	Yes
Error and retry single step	Yes	Yes
Re-run from failed step	Custom workaround	Yes
Open source community support	No	Yes
Cost	Usage-based cost that depends on the complexity of the workflow	Constant base infrastructure cost, plus worker costs that can scale up and down
Scalability	Highly scalable, fully automatic	Highly scalable, managed by user or autoscaling groups, and can be configured
Infrastructure management	No infrastructure management or provisioning as everything handled by AWS	Requires making choices about infrastructure, but AWS manages the infrastructure and software
Language for pipeline development	JSON (or use of visual designer)	Python
Serverless/managed	Serverless	Managed
Integration	Seamlessly integrates with AWS services and manual integration with non-AWS services	Strong integration support for many AWS services, as well as extensive third-party services

MyStateMachine-k4l5ohf02

DesignCodeConfig

Workflow not createdCancelActionsCreate

UndoRedoZoom inZoom outCenterDuplicateDeleteFeedback

Search

ActionsFlowPatternsInfo

MOST POPULAR

AWS Lambda Invoke

Amazon SNS Publish

Amazon ECS RunTask

AWS Step Functions StartExecution

AWS Glue StartJobRun

COMPUTE

Amazon Data Lifecycle Manager

Amazon EBS

Amazon EC2

AWS EC2 Instance Connect

Elastic Inference

Start

Lambda: Invoke

Check file extension

End

Check file extension

Definition

ConfigurationInputOutputError handling

During execution, the Task state calls an API and the response goes into the task result. The result can be manipulated with filters before it is passed as output to the next state. Info

Lambda:Invoke task result example

A read-only example of the kind of task result to expect from this API:

```
{  "ExecutedVersion": "$LATEST",  "Payload": {    "foo": "bar",    "colors": [      "red",      "blue",      "green"    ],    "car": {      "year": 2000    }  }}
```

☐ Transform result with ResultSelector - optional

Use the ResultSelector filter to construct a new JSON object using parts of the task result. Info

☐ Add original input to output using ResultPath - optional

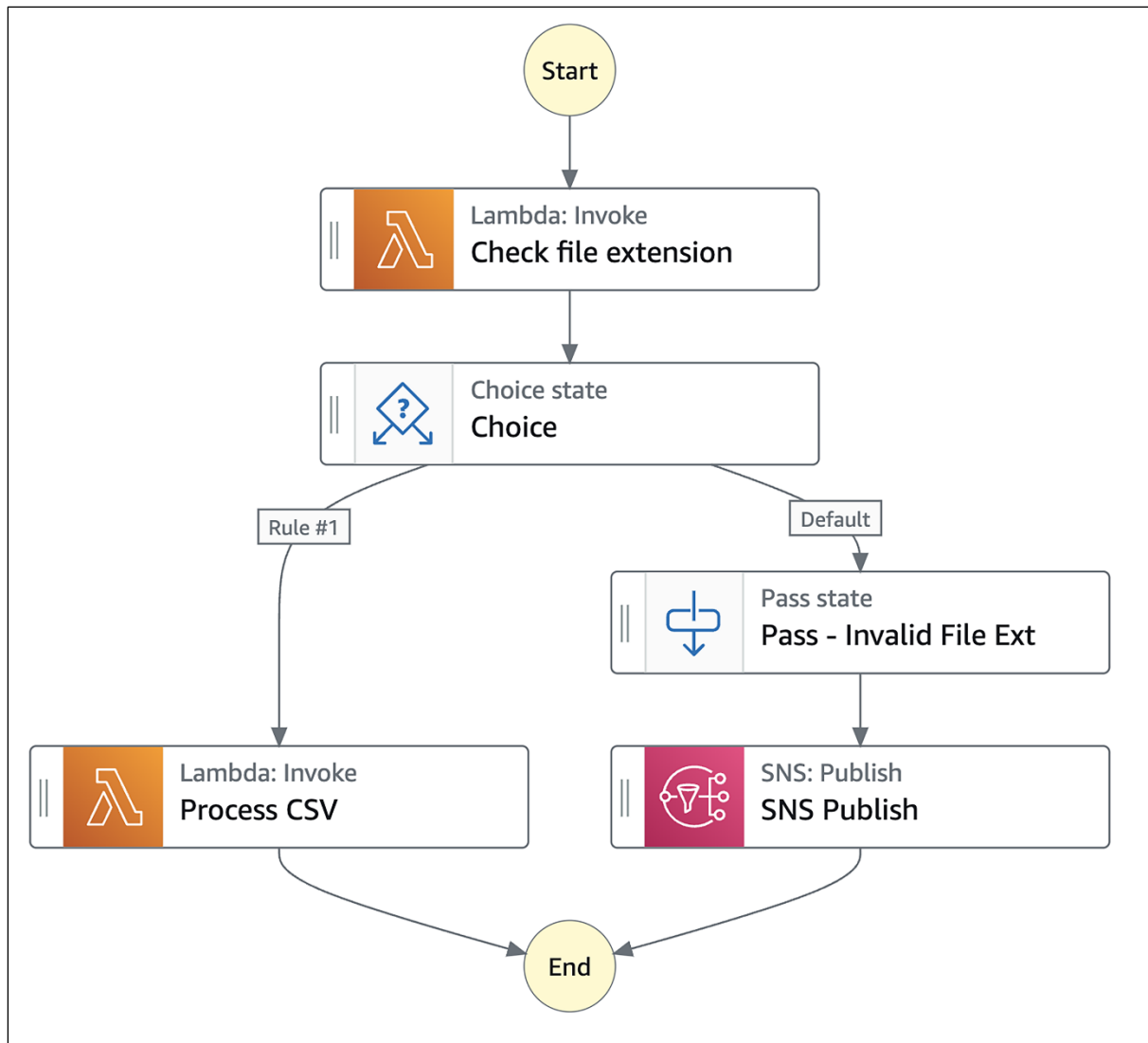
By default, a state sends its task result as output. Use the ResultPath filter to include the original input in the state's output. Info

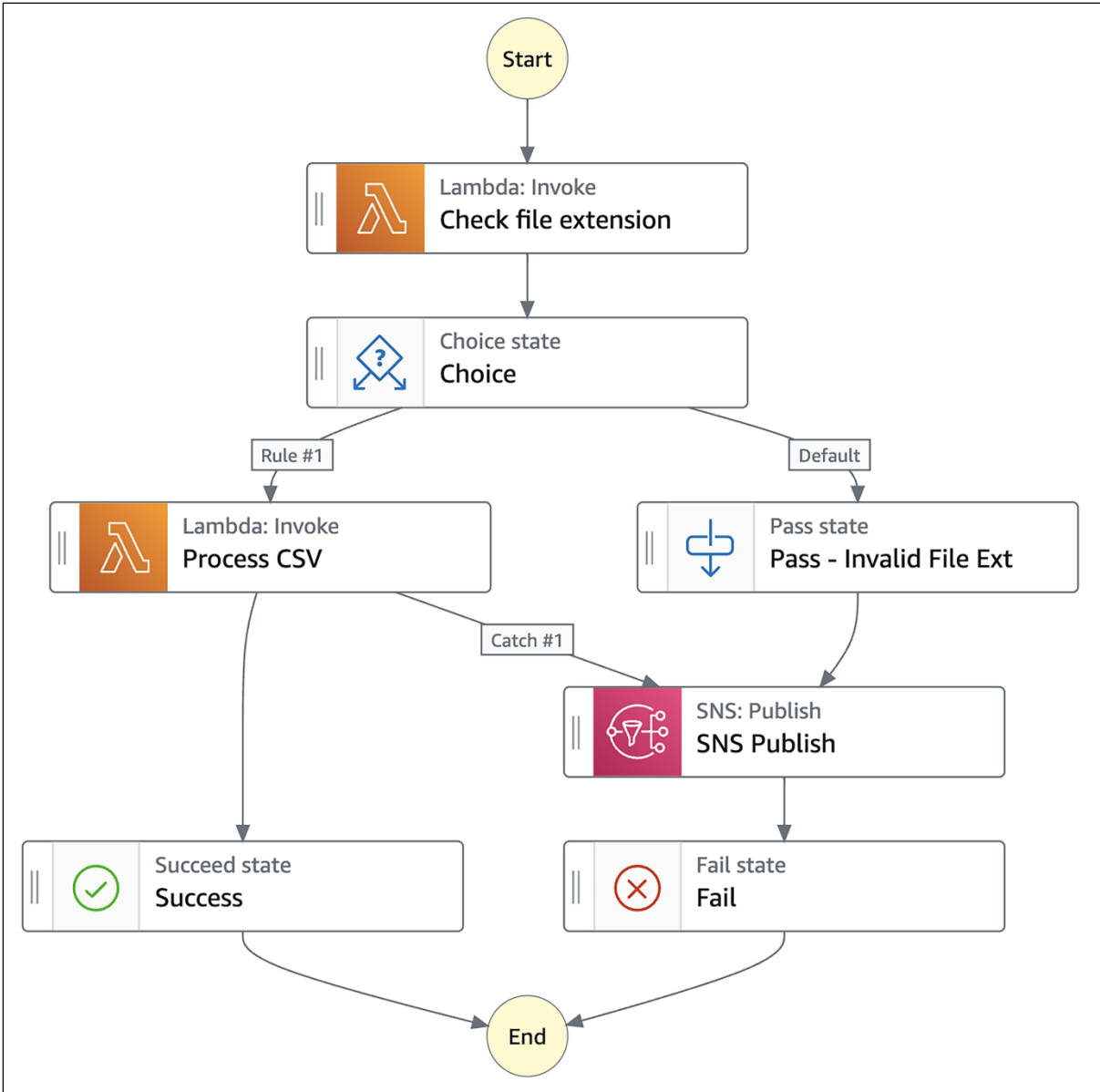
☒ Filter output with OutputPath - optional

Use the OutputPath filter to select a portion of the effective output to pass to the next state. Info

\$Payload

Must use valid JSONPath syntax, and point to an existing key-value pair in the state





Amazon EventBridge > Rules > dataeng-s3-trigger-rule

dataeng-s3-trigger-rule

EditDisableDeleteCloudFormation Template

Rule details

Rule name

dataeng-s3-trigger-rule

Status

Enabled

Event bus name

default

Type

Standard

Description

Rule ARN

arn:aws:events:us-east-2:45590:rule/dataeng-s3-trigger-rule

Event bus ARN

arn:aws:events:us-east-2:45590:event-bus/default

Event pattern

Targets

Monitoring

Tags

Event pattern

```
1 {
2   "source": ["aws.s3"],
3   "detail-type": ["Object Created"],
4   "detail": {
5     "bucket": {
6       "name": ["dataeng-clean-zone-gse23"]
7     },
8     "object": {
9       "key": [{
10        "prefix": "chapter10"
11      }]
12    }
13  }
14 }
```

Edit

Copy

Graph view

Actions

+

Q

🎯

✂

Start

Check file extension

Choice

Process CSV

Pass - Invalid File Ext

SNS Publish

Success

Fail

End

In progress

Failed

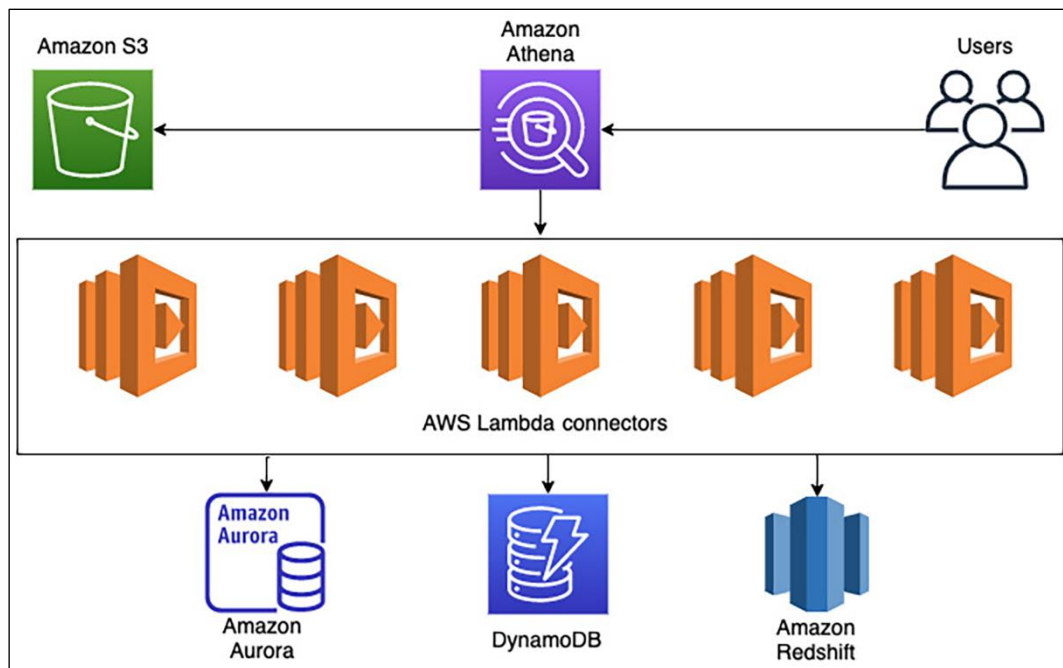
Caught error

Canceled

Succeeded

```
graph TD; Start((Start)) --> Check[Check file extension]; Check --> Choice[Choice]; Choice --> Process[Process CSV]; Choice --> Pass[Pass - Invalid File Ext]; Process --> Success[Success]; Process --> SNS[SNS Publish]; SNS --> Fail[Fail]; Success --> End((End)); Fail --> End;
```


Chapter 11: Ad Hoc Queries with Amazon Athena



Amazon Athena ×

Query editor
Notebook editor [New](#)
Notebook explorer [New](#)

▼ **Jobs**
Workflows
Powered by Step Functions

▼ **Administration**
Workgroups
Capacity reservations [New](#)
Data sources

○ Turn on compact mode

Amazon Athena > Workgroups

Workgroups (1) [Info](#)
Use workgroups to separate users, teams, applications, workloads, and to set limits on amount of data for each query or the entire workgroup process. You can also view query-related metrics in AWS CloudWatch.

Actions ▼ **Create workgroup**

Filter workgroups

	Name	Description	Analytics engine	Engine updates	Created on	Status
<input type="radio"/>	primary	-	Athena engine vers...	Automatic	2023-02-28T22:19...	Turned On

Amazon Athena > Query editor

Editor Recent queries Saved queries Settings Workgroup datalake-user-sand...

Data

Data source: AwsDataCatalog Database: curatedzonedb

Tables and views: **Tables (2)** film_category streaming_films **Views (0)**

Query 1 :

```
1 SELECT category_name,
2 count(category_name) streams
3 FROM streaming_films
4 GROUP BY category_name
5 ORDER BY streams DESC
```

SQL Ln 5, Col 22

Run again Explain Cancel Clear Create

Reuse query results up to 60 minutes ago

Query results **Query stats**

Completed Time in queue: 112 ms Run time: 456 ms Data scanned: 3.63 KB

Results (16) Copy Download results

Search rows

#	category_name	streams
1	Sports	662
2	Foreign	648
3	Children	603
4	Documentary	558

Amazon Athena > Query editor

Editor Recent queries **Saved queries** Settings Workgroup datalake-user-sand...

Data

Data source: AwsDataCatalog Database: curatedzonedb

Tables and views: **Tables (2)** film_category streaming_films **Views (0)**

Query 1 :

```
1 SELECT category_name,
2 count(category_name) streams
3 FROM streaming_films
4 GROUP BY category_name
5 ORDER BY streams DESC
```

SQL Ln 1, Col 1

Run again Explain Cancel Clear Create

Reuse query results up to 60 minutes ago

Query results **Query stats**

Completed Time in queue: 115 ms Run time: 444 ms Data scanned: 3.63 KB

Results (16) Copy Download results

Search rows

#	category_name	streams
1	Sports	662

Amazon Athena > Query editor

Editor

Recent queries

Saved queries

Settings

Workgroup datalake-user-sand...

Recent queries (1/12)

Q Search recent queries

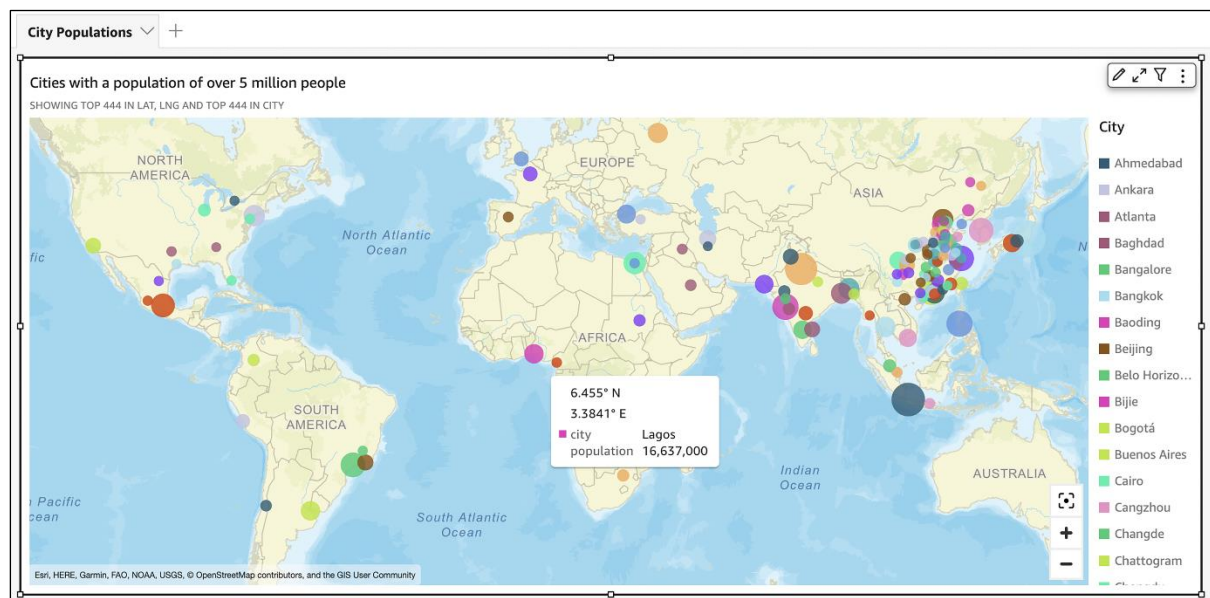
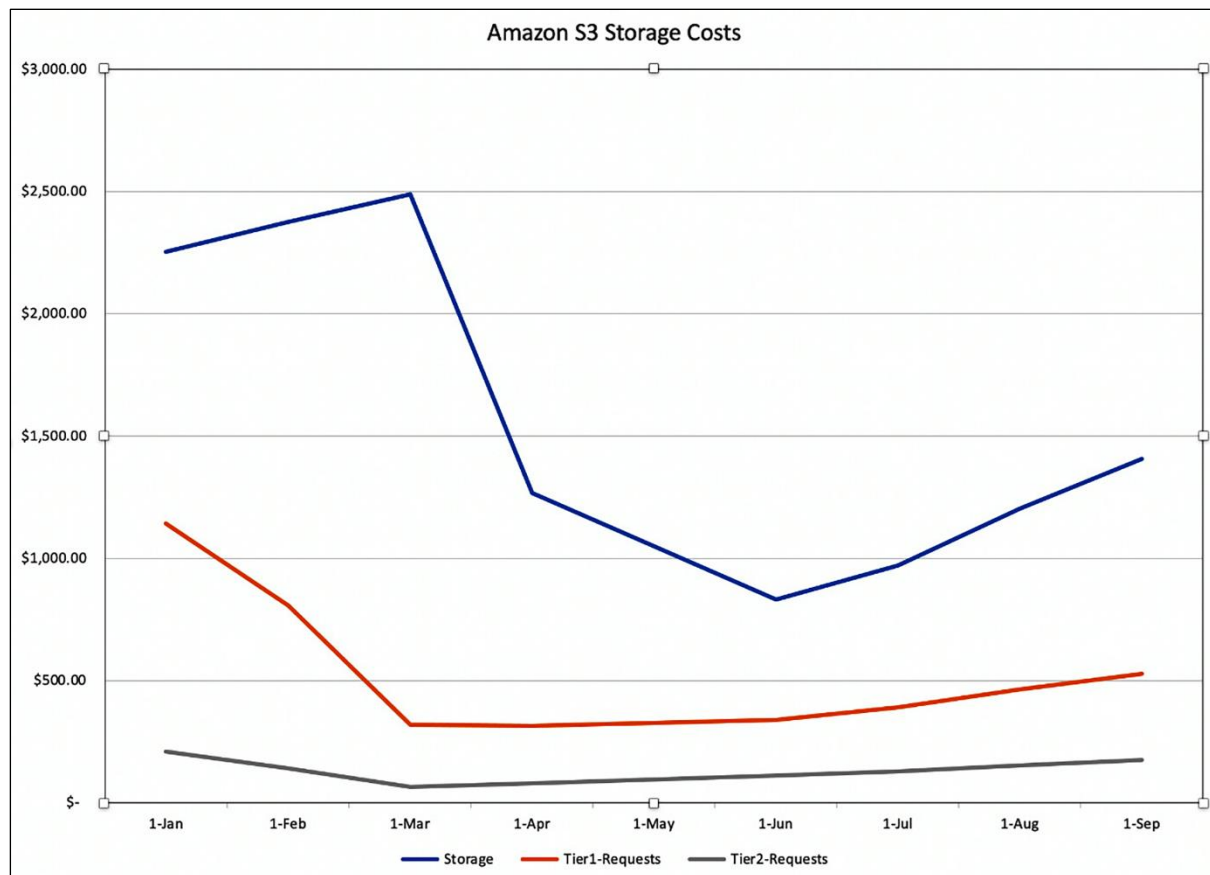
< 1 >

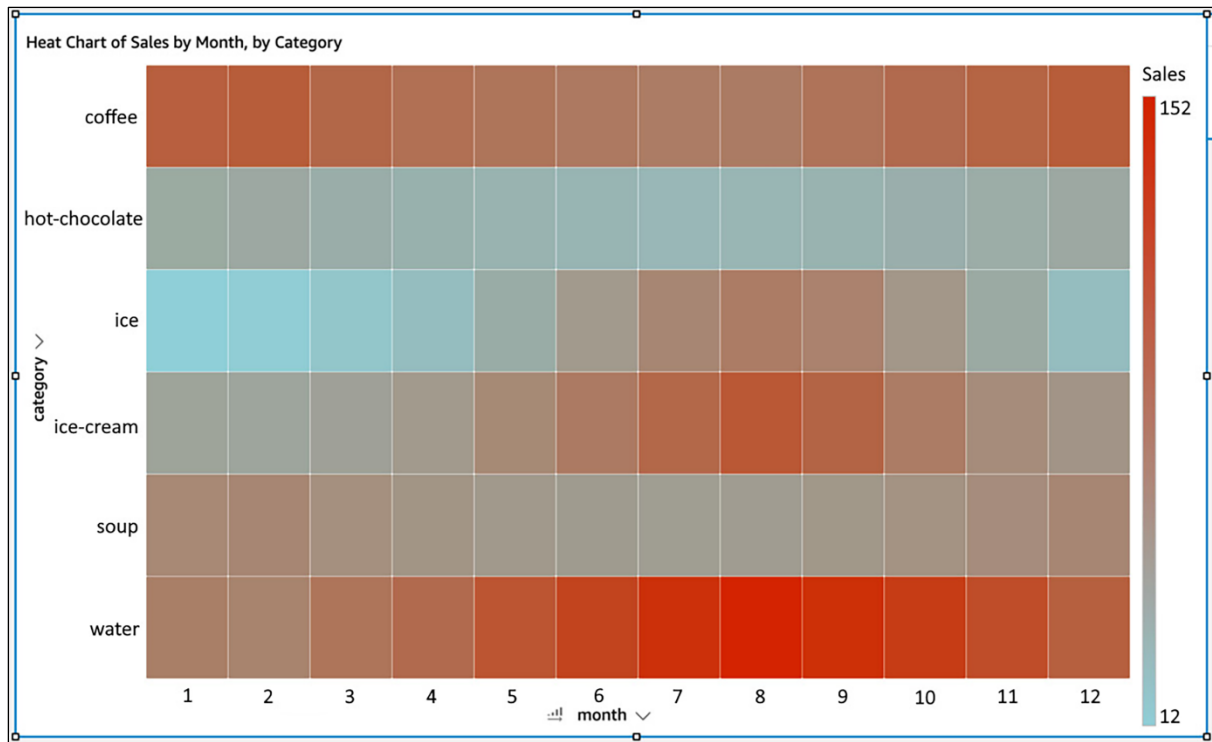
Download results

Download CSV

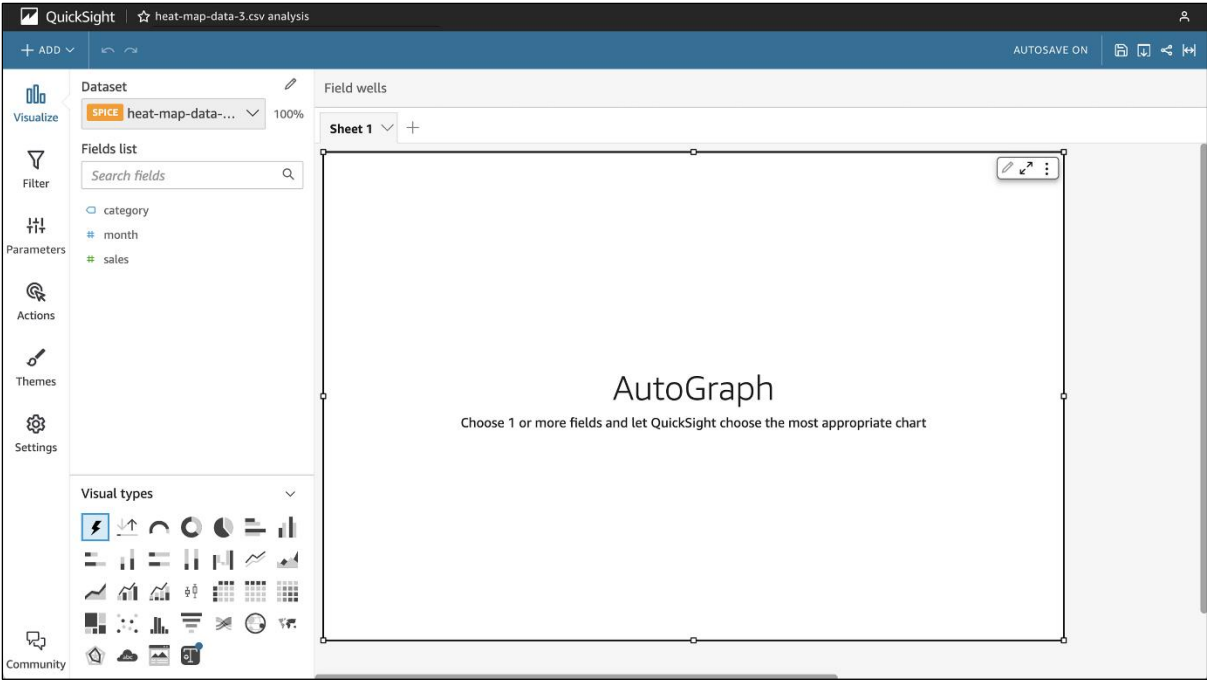
Execution ID	Query	Start time	Status	Run time	Cache	Data sc...	Query engine versi...	Encryption
490a4d28-5dfa...	select name from film_cat...	2023-06-11T21...	Failed	150 ms	-	0 MB	Athena engine version 3	SSE_S3
e135fae1-2980...	SELECT state, count(state)...	2023-06-11T21...	Completed	455 ms	-	8.65 KB	Athena engine version 3	SSE_S3
74acb79b-6b6...	SELECT category_name, c...	2023-06-11T21...	Completed	444 ms	-	3.63 KB	Athena engine version 3	SSE_S3
ef508c93-d889...	SELECT category_name, c...	2023-06-11T21...	Completed	194 ms	Result reuse	0 MB	Athena engine version 3	SSE_S3
390b4229-a90...	SELECT category_name, c...	2023-06-11T21...	Completed	170 ms	Result reuse	0 MB	Athena engine version 3	SSE_S3
f7f0b4ed-2216...	SELECT category_name, c...	2023-06-11T21...	Completed	192 ms	Result reuse	0 MB	Athena engine version 3	SSE_S3
968b8529-3ea...	SELECT category_name, c...	2023-06-11T21...	Completed	182 ms	Result reuse	0 MB	Athena engine version 3	SSE_S3
98b25897-d28...	SELECT category_name, c...	2023-06-11T21...	Completed	172 ms	Result reuse	0 MB	Athena engine version 3	SSE_S3
69267b6a-9ab...	SELECT category_name, c...	2023-06-11T21...	Completed	424 ms	-	3.63 KB	Athena engine version 3	SSE_S3

Chapter 12: Visualizing Data with Amazon QuickSight





Upload a file (.csv, .tsv, .clf, .elf, .xlsx, .json)	Salesforce Connect to Salesforce	S3 Analytics	S3
Athena	RDS	Redshift Auto-discovered	Redshift Manual connect
MySQL	PostgreSQL	ORACLE	SQL Server
Aurora	MariaDB	Presto	Spark
Teradata Provided by Teradata	Snowflake	AWS IoT Analytics	Amazon OpenSearch Ser... Successor to Amazon Elasticsearch Ser...
Timestream	Exasol	Databricks	GitHub
Twitter	Jira	ServiceNow	Adobe Analytics



KPI's ▾ +

Sales Revenue - Current vs Target

Current

78,520

Target goal

100,000

78.52%

78.52%



New Customers - Current vs Target

Current

1,350

Target goal

1,500

90%

90%



Customer Cancellations - Current vs Max Target

Current

268

Target goal

300

89.33%

89.33%



Edition	<input checked="" type="radio"/> Enterprise	<input type="radio"/> Enterprise + Q Learn more
Team trial for 30 days (4 authors)*	FREE	FREE
Author per month (yearly)**	\$18	\$28
Author per month (monthly)**	\$24	\$34
Readers (pay-per-Session)	\$0.30 / session (max \$5)****	\$0.30 / session (max \$10)****
Additional SPICE per month	\$0.38 per GB	\$0.38 per GB
QuickSight Q regional fee	N/A	\$250 / mo / region
Natural language query with QuickSight Q	N/A	INCLUDED
Single Sign On with SAML or OpenID Connect	✓	✓
Connect to spreadsheets, databases & business apps	✓	✓
Access data in Private VPCs	✓	✓
Row-level security for dashboards	✓	✓
Secure data encryption at rest	✓	✓
Connect to your Active Directory	✓	✓
Use Active Directory groups***	✓	✓
Send email reports	✓	✓
Embed QuickSight	✓	✓
Capacity-based pricing	✓	✓
Supported regions	Learn more	Learn more

* Trial authors are auto-converted to month-to-month subscription upon trial expiry

** Each additional author includes 10GB of SPICE capacity

*** Active Directory groups are available in accounts connected to Active Directory

**** Sessions of 30-minute duration. Total charges for each reader are capped at \$5 per month. [Conditions apply.](#)

[Sign up for Standard Edition here.](#)

[QuickSight Standard Sign-up](#)

Create your QuickSight account

Standard

[Back](#)

Authentication method

- ☐ Use IAM federated identities & QuickSight-managed users
Authenticate with single sign-on (SAML or OpenID Connect), AWS IAM credentials, or QuickSight credentials
- ☒ Use IAM federated identities only
Authenticate with single sign-on (SAML or OpenID Connect) or AWS IAM credentials

QuickSight region

Select a region

US East (Ohio)

Account info

QuickSight account name

You will need this for you and others to sign in

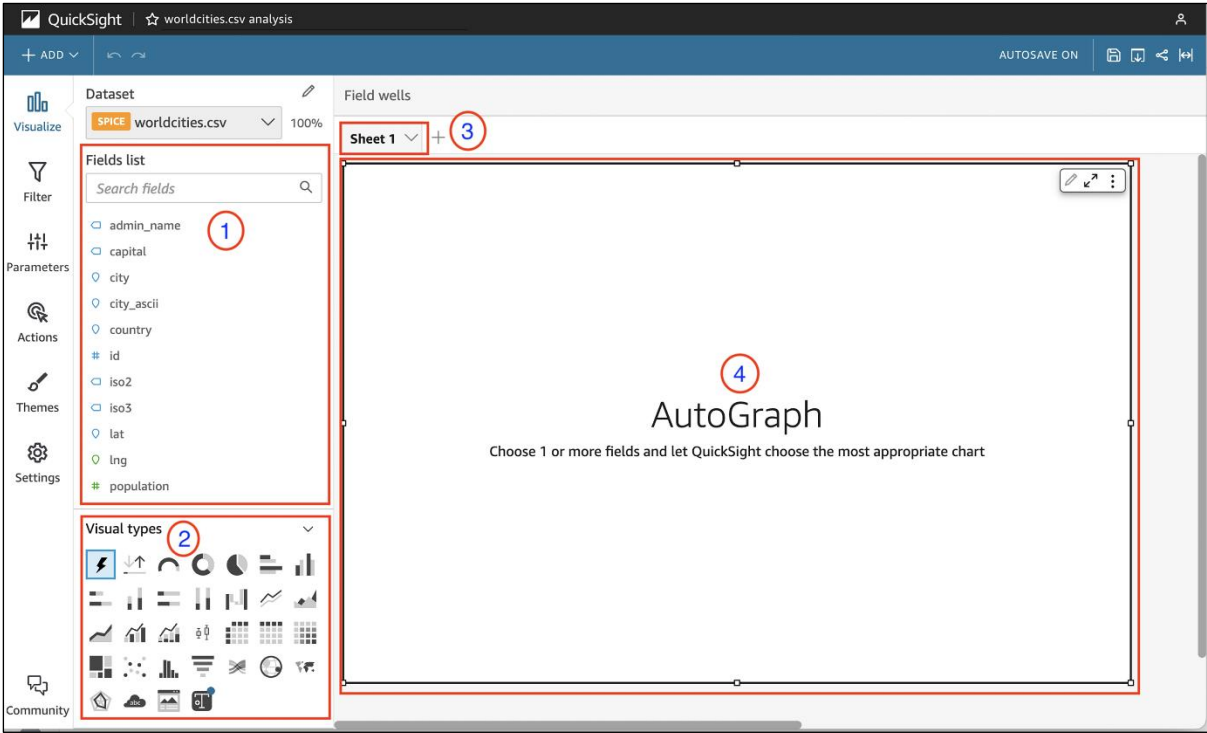
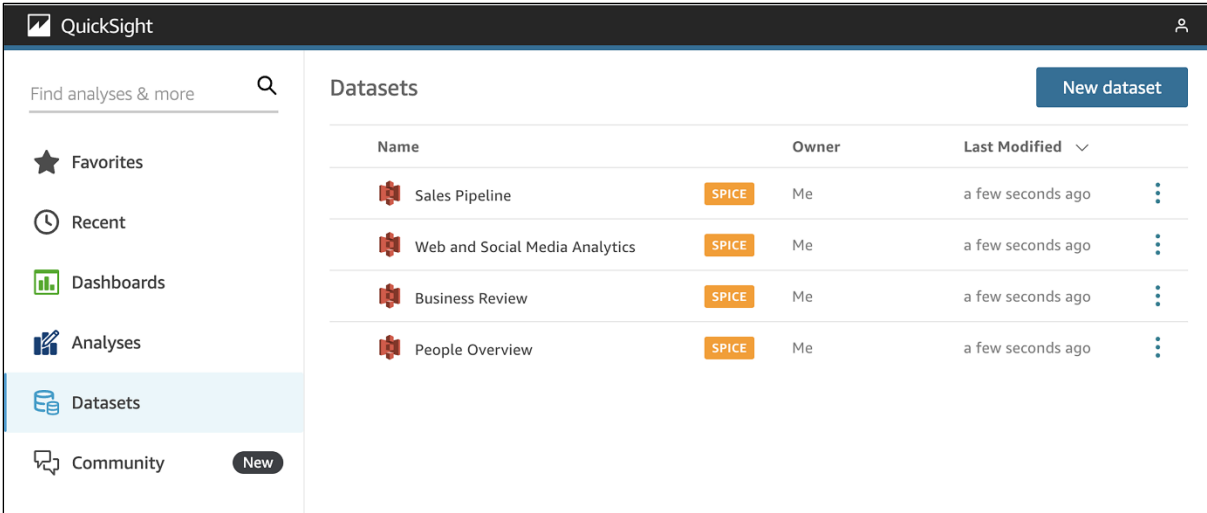
data-

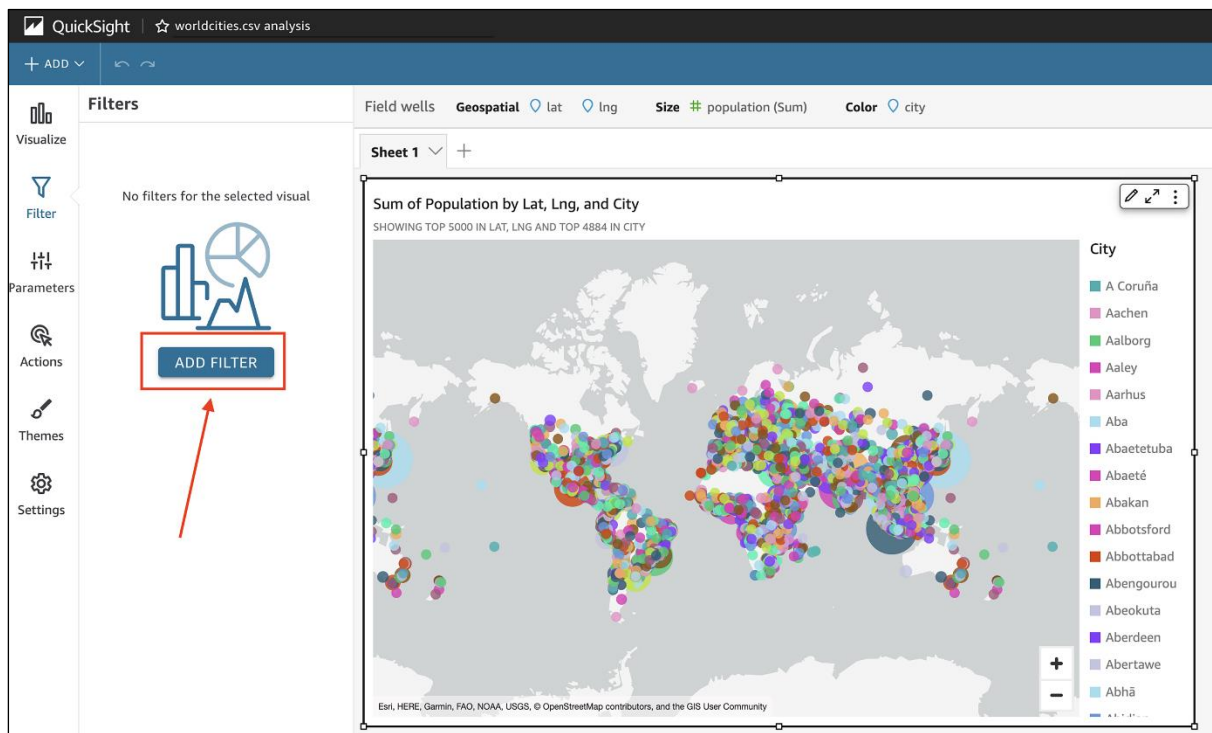
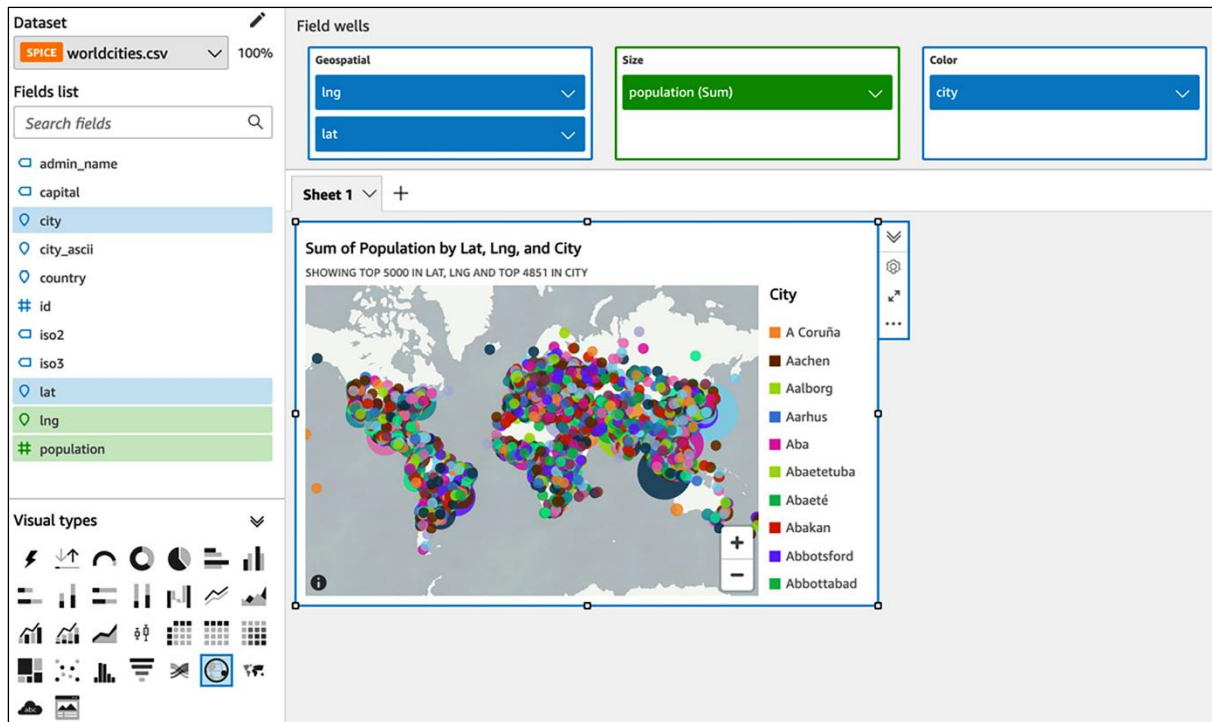
Notification email address

For QuickSight to send important notifications

gare

l.com





< Edit filter



Applied to



Only this visual



population



Equals - none

Aggregation

Sum



Filter condition

Greater than or equal to



Use parameters

Minimum value

3000000

Null options

Include nulls



OR

ADD FILTER CONDITION

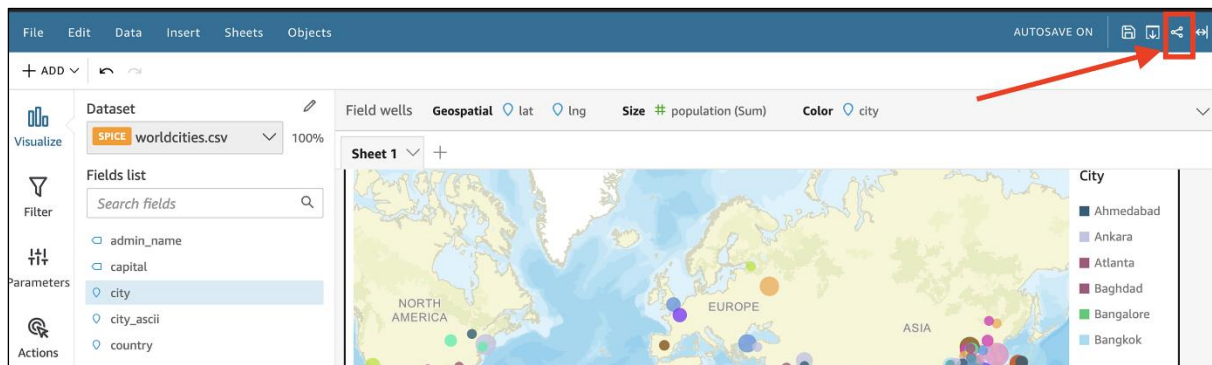
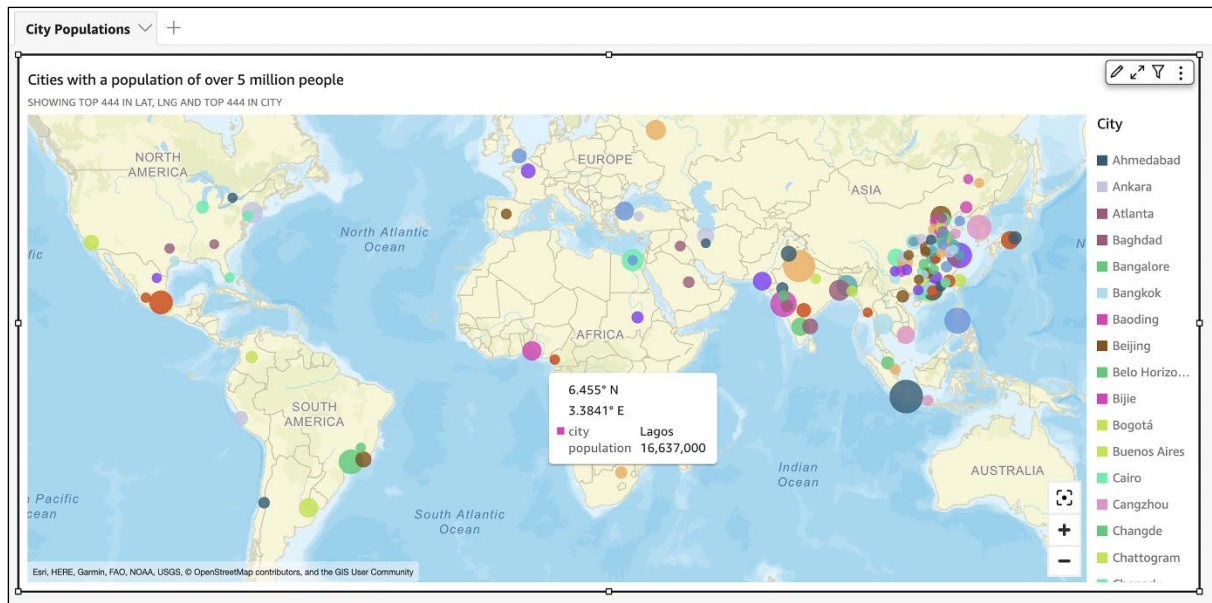


Note: There are limitations on how you can group filters.

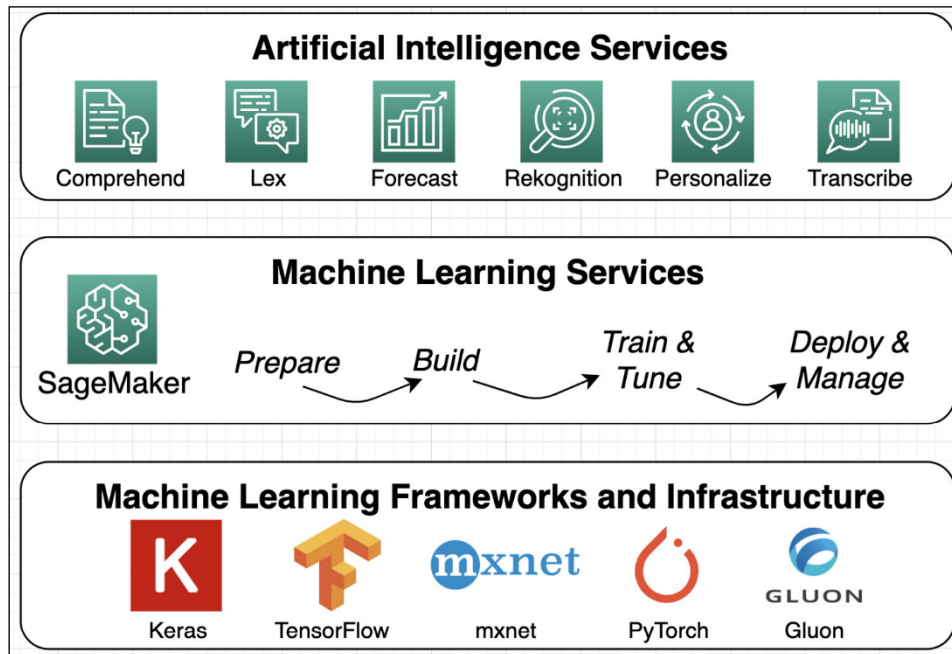
[Learn more](#)

APPLY

DELETE FILTER



Chapter 13: Enabling Artificial Intelligence and Machine Learning



DATE	REF NO	DESCRIPTION	CHARGES
4/15/2019	2559498	GUEST ROOM	\$179.00
4/15/2019	2559498	STATE TAX	\$10.74
4/15/2019	2559498	CITY TAX	\$16.11
4/16/2019	2559777	C3 FOOD DRINK	\$7.00
4/16/2019	2559811	VS	(\$212.85)
BALANCE			\$0.00

Hilton Honors(R) stays are posted within 72 hours of checkout. To check your earnings or book your next stay at more than 4,000 hotels and resorts in 100 countries, please visit [Honors.com](https://honors.com)

Thank you for choosing Doubletree! Come back soon to enjoy our warm chocolate chip cookies and relaxed hospitality. For your next trip visit us at doubletree.com for our best available rates!

DATE	REF NO	DESCRIPTION	CHARGES	
4/15/2019	2559498	GUEST ROOM	\$179.00	
4/15/2019	2559498	STATE TAX	\$10.74	
4/15/2019	2559498	CITY TAX	\$16.11	
4/16/2019	2559777	C3 FOOD DRINK	\$7.00	
4/16/2019	2559811	VS	(\$212.85)	



Amazon SQS > Queues > Create queue

Create queue

Details

Type

Choose the queue type for your application or cloud infrastructure.

☒ Standard [Info](#)

At-least-once delivery, message ordering isn't preserved

- At-least once delivery
- Best-effort ordering

☐ FIFO [Info](#)

First-in-first-out delivery, message ordering is preserved

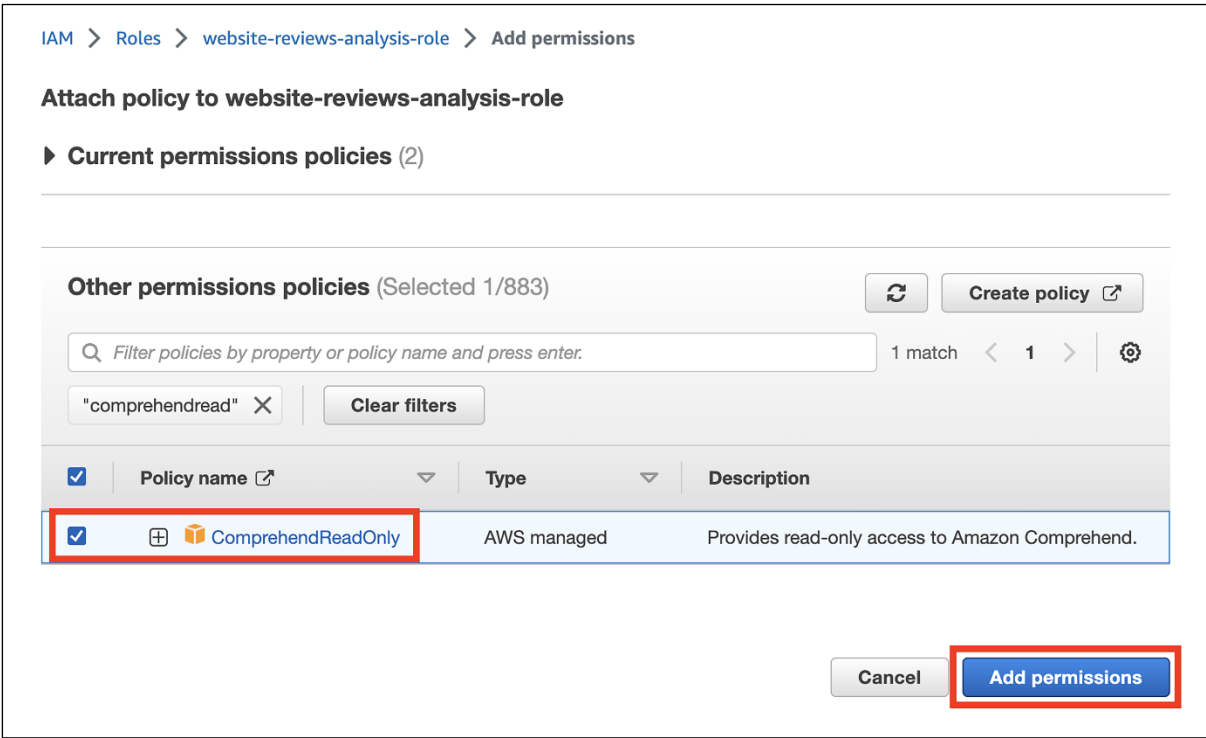
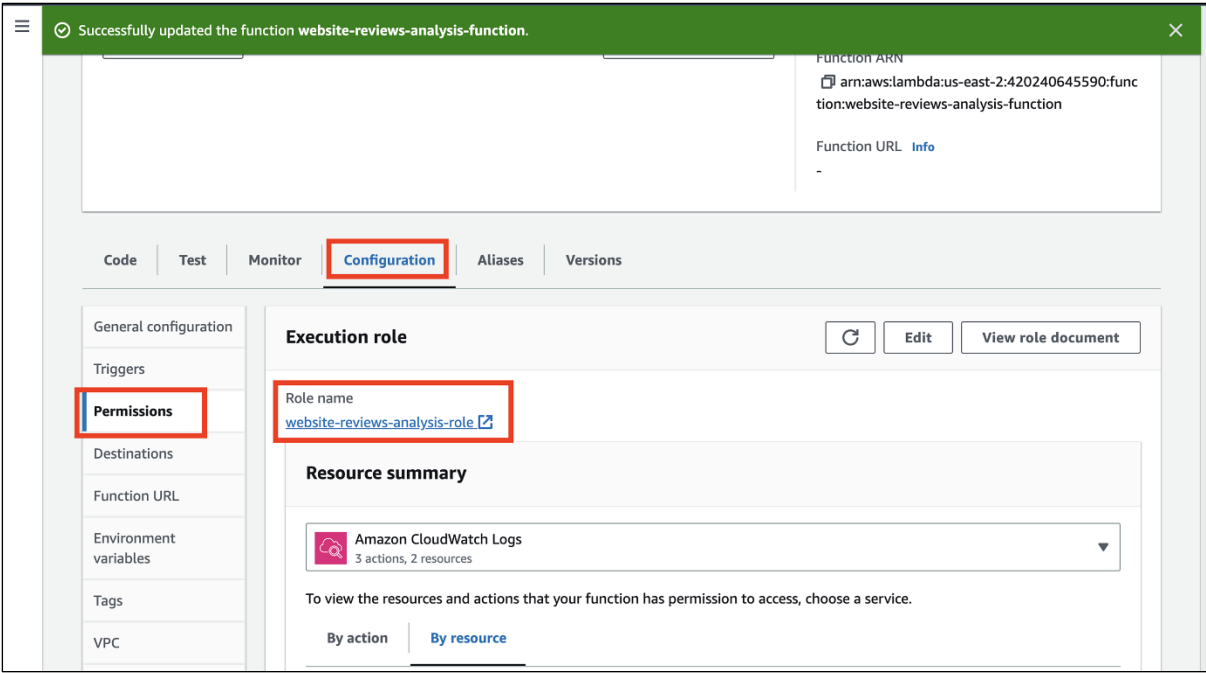
- First-in-first-out delivery
- Exactly-once processing

You can't change the queue type after you create a queue.

Name

website-reviews-queue

A queue name is case-sensitive and can have up to 80 characters. You can use alphanumeric characters, hyphens (-), and underscores (_).



Amazon SQS > Queues > website-reviews-queue

website-reviews-queue

Edit Delete Purge **Send and receive messages** Start DLQ redrive

Details [Info](#)

Name website-reviews-queue	Type Standard	ARN arn:aws:sqs:us-east-2:123456789012:website-reviews-queue
Encryption Amazon SQS key (SSE-SQS)	URL https://sqs.us-east-2.amazonaws.com/123456789012/website-reviews-queue	Dead-letter queue -

► More

SNS subscriptions **Lambda triggers** Dead-letter queue Monitoring Tagging Access policy Encryption Dead-letter queue redrive tasks

Lambda triggers (1) [Info](#)

[Refresh](#) [View in Lambda](#) [Delete](#) [Configure Lambda function trigger](#)

Search triggers

UUID	ARN	Status	Last modified
ba81ed70-3c5f-4b05-bb6d-9b58796f2fac	arn:aws:lambda:us-east-2:123456789012:function:website-reviews-analysis-function	Enabled	7/6/2023, 11:38:56 PM

CloudWatch > Log groups > /aws/lambda/website-reviews-analysis-function > 2023/07/09/[\$LATEST]718f3579ce364c9fb733ffcc957913c7

Log events

You can use the filter bar below to search for and match terms, phrases, or values in your log events. [Learn more about filter patterns](#)

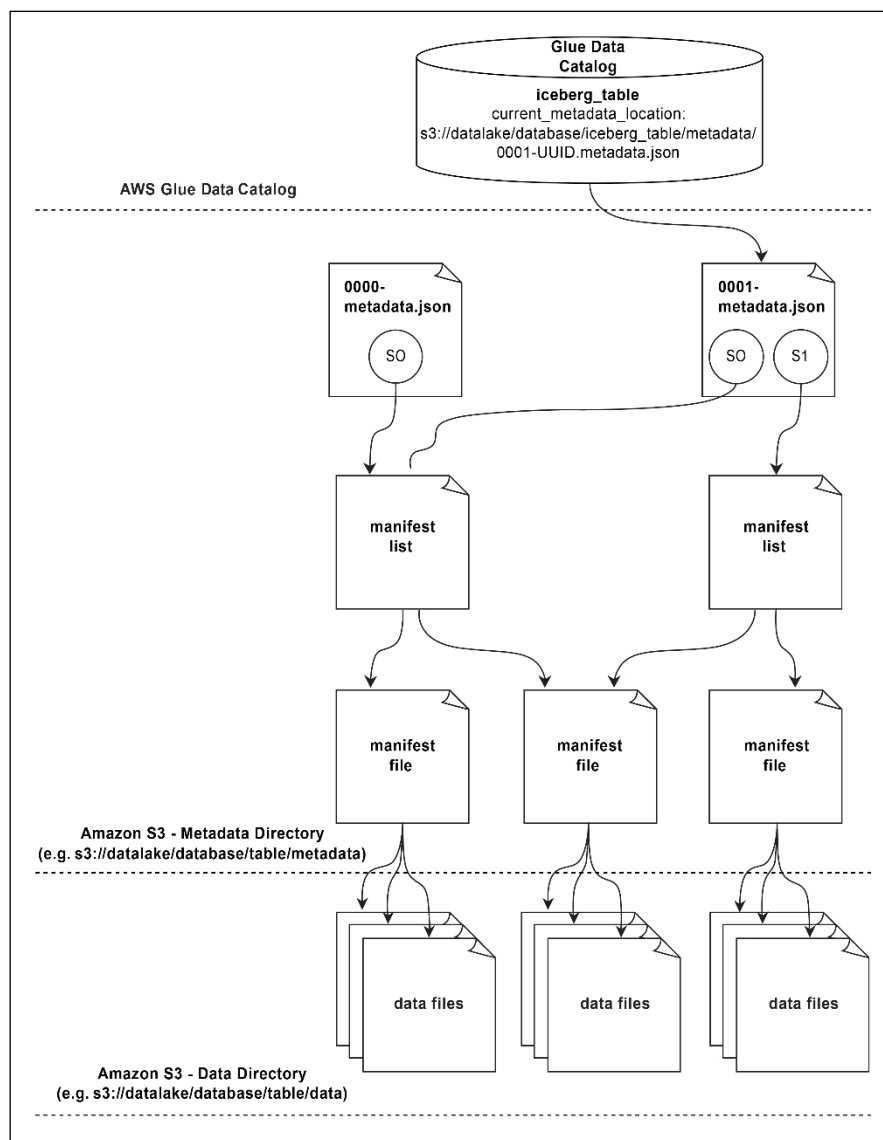
[Refresh](#) Actions ▼ Start tailing Create metric filter

Filter events Clear 1m 30m 1h 12h Custom [Display](#) ▼

Timestamp	Message
	No older events at this moment. Retry
2023-07-09T14:51:24...	INIT_START Runtime Version: python:3.10.v5 Runtime Version ARN: arn:aws:lambda:us-east-2::runtime:51b59a...
2023-07-09T14:51:25...	START RequestId: 344249c8-6797-5c5d-8318-14ed6475c331 Version: \$LATEST
2023-07-09T14:51:25...	I recently stayed at the Kensington Hotel in downtown Cape Town and was very impressed. The hotel is bea...
2023-07-09T14:51:25...	Calling DetectSentiment
2023-07-09T14:51:25...	SENTIMENT: POSITIVE
2023-07-09T14:51:25...	SENTIMENT SCORE: {'Positive': 0.999713122844696, 'Negative': 2.825235787895508e-05, 'Neutral': 0.0002193...
2023-07-09T14:51:25...	Calling DetectEntities
2023-07-09T14:51:25...	ENTITY: Kensington Hotel, ENTITY TYPE: ORGANIZATION
2023-07-09T14:51:25...	ENTITY: Cape Town, ENTITY TYPE: LOCATION
2023-07-09T14:51:25...	ENTITY: Elizabeth's Kitchen, ENTITY TYPE: ORGANIZATION
2023-07-09T14:51:25...	END RequestId: 344249c8-6797-5c5d-8318-14ed6475c331

Chapter 14: Building Transactional Data Lakes

	Copy-On-Write (COW)	Merge-On-Read (MOR)
Action on record update/delete	New version of affected file is created containing newly updated records, or skipping deleted records	Deleted and updated rows are written to a deletion tracking file, and updated records are written to a new file
Action when table is read	The metadata for the table tracks the files that have the latest data, and only those files are read	The metadata for the table tracks the original file, the file tracking deleted records, and files containing new/updated data. The query engine needs to merge each of these files to query the data.
Optimized for reads or writes?	Optimized for table reads	Optimized for table writes



AWS Glue > Tables > streaming_films_ib

streaming_films_ib

Last updated (UTC)
July 23, 2023 at 18:33:55

Version 1 (Current version)

Actions

Table overview

Data quality New

Table details

Advanced properties

Serde parameters (0)

Key	Value
No parameters	
No parameters to display.	

Table properties (3)

Key	Value
metadata_location	s3://dataeng-curved-zone-gse23/iceberg/streaming_films/metadata/00001-ec5cdb77-7b7b-4277-abae-
previous_metadata_location	s3://dataeng-curved-zone-gse23/iceberg/streaming_films/metadata/00000-b118a8c2-eeee-485f-8f79-5
table_type	ICEBERG

Objects (9)

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 others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create Folder

Upload

Find objects by prefix

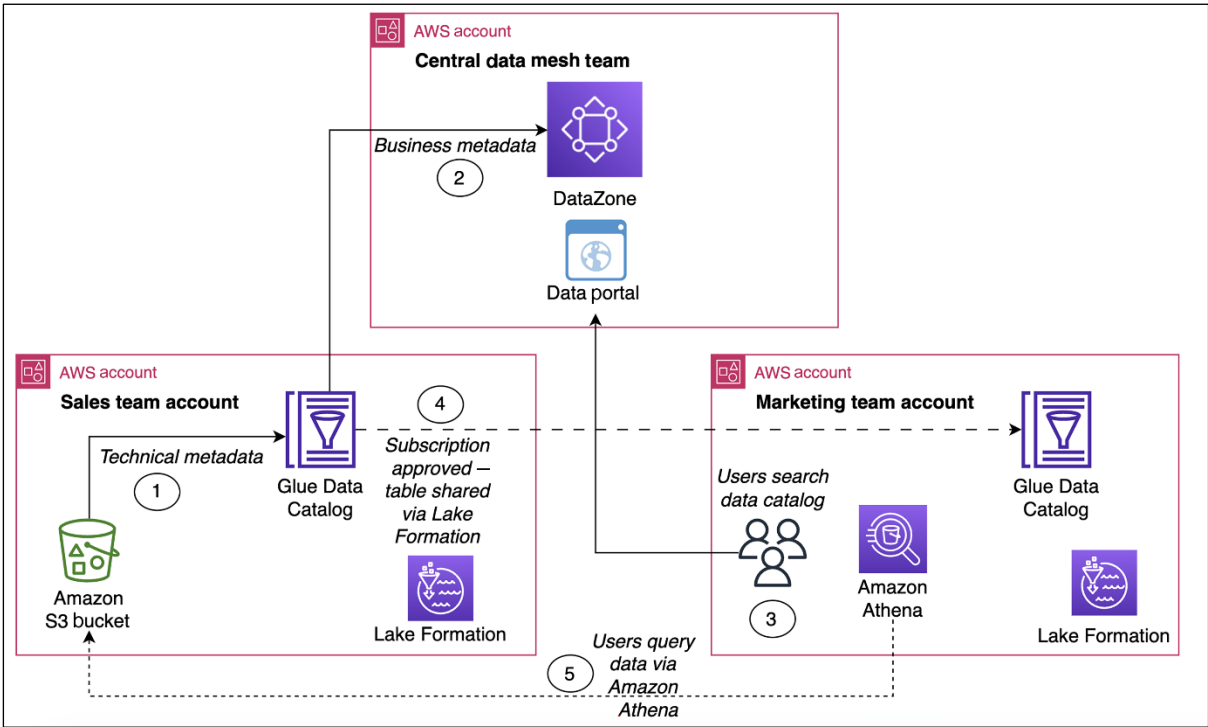
	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	00000-76233ac9-4abc-49e6-ad89-60713d3ade15.metadata.json	json	July 30, 2023, 11:58:16 (UTC-04:00)	2.9 KB	Standard
<input type="checkbox"/>	00001-8e506fa1-0d90-4060-b423-7f178349b5d6.metadata.json	json	July 30, 2023, 11:58:44 (UTC-04:00)	4.1 KB	Standard
<input type="checkbox"/>	00002-5b02dd23-5973-45b2-a1be-5c665f881698.metadata.json	json	July 30, 2023, 11:58:45 (UTC-04:00)	8.4 KB	Standard
<input type="checkbox"/>	00003-7c53f5a1-2aea-4e5f-bcf8-da11f0df93c5.metadata.json	json	July 30, 2023, 12:10:13 (UTC-04:00)	9.4 KB	Standard
<input type="checkbox"/>	20230730_155839_00131_46438-r71566d0-2724-4678-939f-9ebbaebe1859.stats	stats	July 30, 2023, 11:58:45 (UTC-04:00)	44.1 KB	Standard
<input type="checkbox"/>	2f8d6261-6ad4-48a6-a206-5ecfaf3fa843-m0.avro	avro	July 30, 2023, 11:58:44 (UTC-04:00)	30.2 KB	Standard
<input type="checkbox"/>	5362d55d-8199-4e68-bee1-6152b6683aef-m0.avro	avro	July 30, 2023, 12:10:13 (UTC-04:00)	30.3 KB	Standard
<input type="checkbox"/>	snap-2240417827427115165-1-2f8d6261-6ad4-48a6-a206-5ecfaf3fa843.avro	avro	July 30, 2023, 11:58:44 (UTC-04:00)	4.2 KB	Standard
<input type="checkbox"/>	snap-3424169349391369617-1-5362d55d-8199-4e68-bee1-6152b6683aef.avro	avro	July 30, 2023, 12:10:13 (UTC-04:00)	4.2 KB	Standard

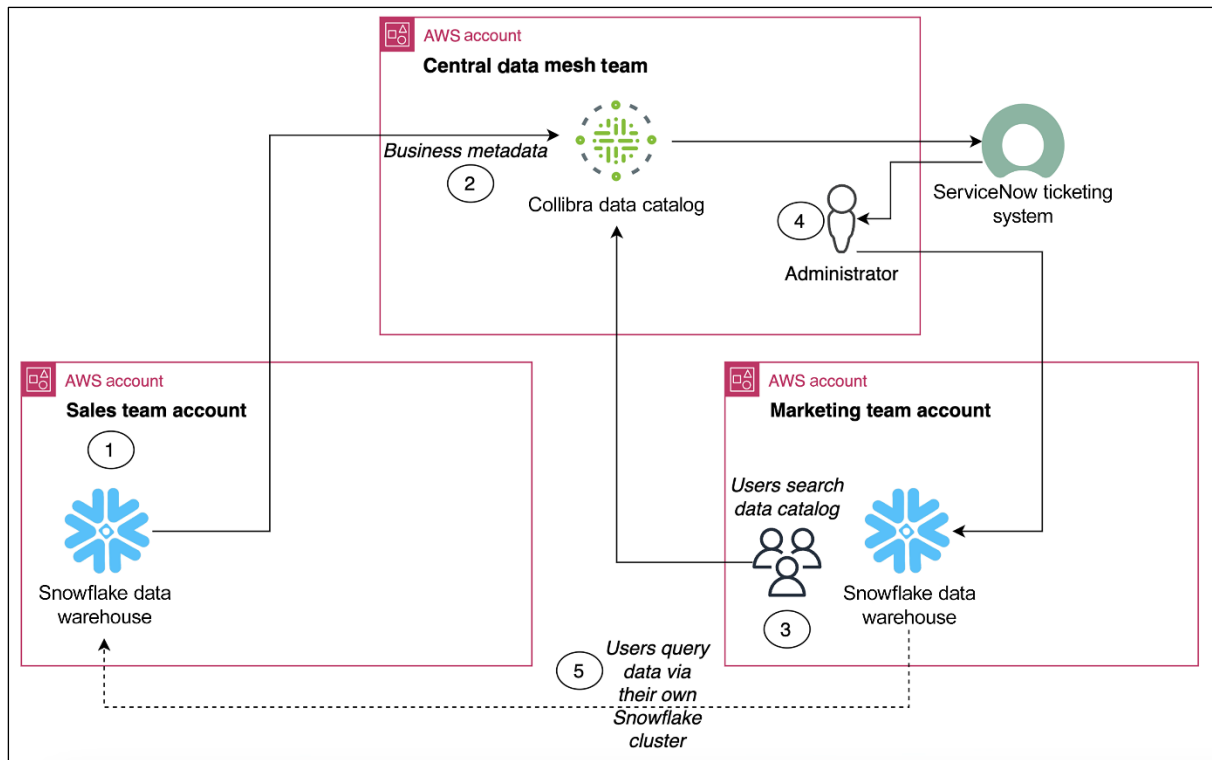
```
d6ec39.metadata.json  {} 00002-1d77459b-273a-4e6d-adf6-f3c62379a6f9.metadata.json × 申  〇  〇  ...
{} 00002-1d77459b-273a-4e6d-adf6-f3c62379a6f9.metadata.json > {} refs > {} main > abc type
130   "current-snapshot-id" : 1954287477388955304,
131   "refs" : {
132     "main" : {
133       "snapshot-id" : 1954287477388955304,
134       "type" : "branch"
135     }
136   },
137   "snapshots" : [ {
138     "sequence-number" : 1,
139     "snapshot-id" : 8054068847429778299,
140     "timestamp-ms" : 1690137234995,
141     "summary" : {
142       "operation" : "append",
143       "trino_query_id" : "20230723_183350_00145_y5gpr",
144       "added-data-files" : "127",
145       "added-records" : "8550",
146       "added-files-size" : "665754",
147       "changed-partition-count" : "16",
148       "total-records" : "8550",
149       "total-files-size" : "665754",
150       "total-data-files" : "127",
151       "total-delete-files" : "0",
152       "total-position-deletes" : "0",
153       "total-equality-deletes" : "0"
154     },
155     "manifest-list" : "s3://dataeng-curated-zone-gse23/iceberg/streaming_films/m
156     "schema-id" : 0
157   }, {
158     "sequence-number" : 2,
159     "snapshot-id" : 1954287477388955304,
160     "parent-snapshot-id" : 8054068847429778299,
161     "timestamp-ms" : 1690143975094,
162     "summary" : {
163       "operation" : "delete",
164       "trino_query_id" : "20230723_202612_00037_gpms6",
165       "deleted-data-files" : "11",
166       "deleted-records" : "558",
167       "removed-files-size" : "54404",
168       "changed-partition-count" : "1",
169       "total-records" : "7992",
170       "total-files-size" : "611350",
171       "total-data-files" : "116",
```

Ln 134, Col 24 Spaces: 2 UTF-8 LF {} JSON

Chapter 15: Implementing a Data Mesh Strategy

Field Name	Field Type	Description
Sales Region	<i>Country</i> Business Glossary	The sales region for this dataset
Sales Quarter	<i>Quarter</i> Business Glossary	The calendar quarter that this dataset covers
Dataset Support	String	Provide the email address to use for queries about this dataset
Last Audit Date	Date	Provide the date that this dataset was last audited





Summary		
Description -	Status ✔ Available	Created on October 08, 2023, 18:09 (UTC-04:00)
Data portal URL https://d[redacted].us-east-2.on.aws	IAM Identity Center Enabled	ARN arn:aws:datazone:us-east-2:4[redacted]:domain/d[redacted] x5

Film Analysis Project

OVERVIEWDATAENVIRONMENTSMEMBERS

films-CuratedZoneDB

ACTIONS

RUN

Source type: AWS Glue • Created at: 10/15/2023, 3:08:56 PM

No description

DATA SOURCE RUNS

DATA SOURCE DEFINITION

DETAILS

SCHEDULE

Date and time

REFRESH

STATUS: ALL

10/15/2023, 3:09:54 PM

Completed

Data source run activities list

Run type	Duration	Added	Updated	Unchanged	Failed
On demand	00:00:02	3	0	0	0

ASSET STATUS: ALL

Asset name	Database name	Status	More info
streaming_films	curatedzonedb	Successfully created	-
film_category	curatedzonedb	Successfully created	-
category_streams	curatedzonedb	Successfully created	-

FILM CATALOG P...

Search Assets

CatalogDomain

Ffilm-catalog-t...

Film Catalog Project

OVERVIEWDATAENVIRONMENTSMEMBERS

Automated Metadata Generation

REJECT ALL

ACCEPT ALL

Green icons indicate automatically generated metadata suggestions for the data asset. Click on these icons to edit, accept, or reject each suggestion. You also have the option to select 'Accept All' or 'Reject All' for all auto-generated suggestions related to the asset. [Learn More](#)

Film Category

ACTIONS

PUBLISH ASSET

Technical name: film_category • Asset Type: amazon.datazone.GlueTableAssetType

No description

BUSINESS METADATA

SCHEMA

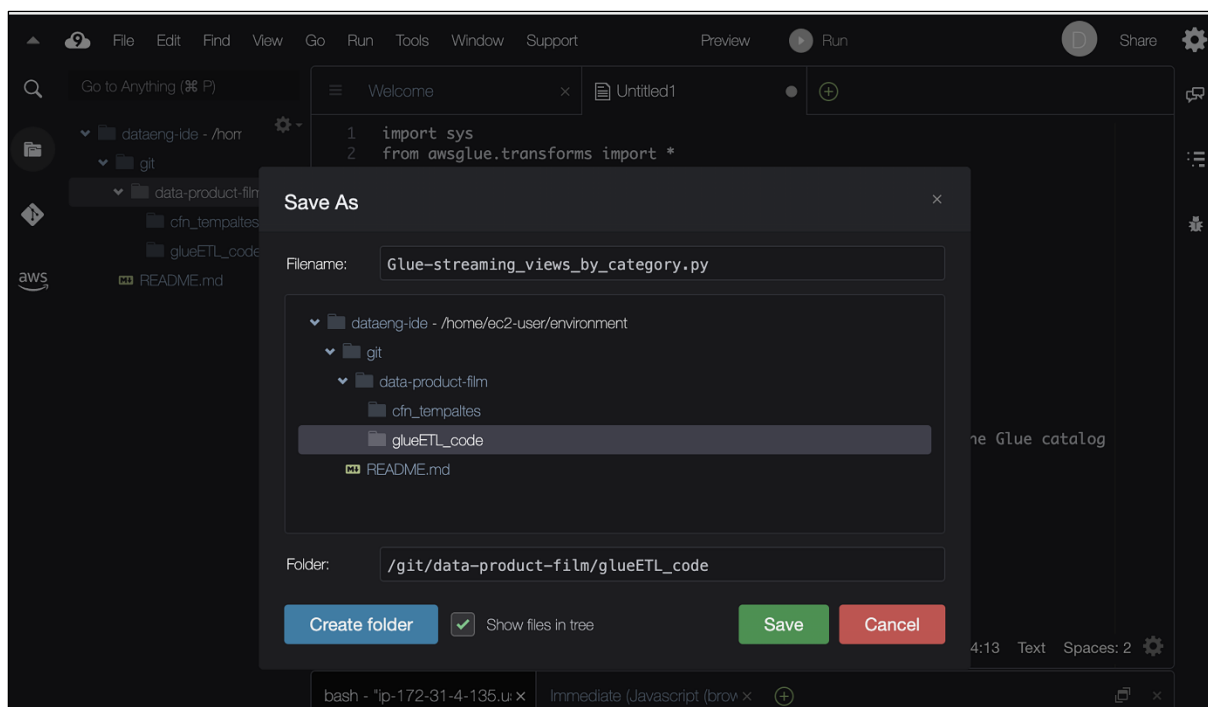
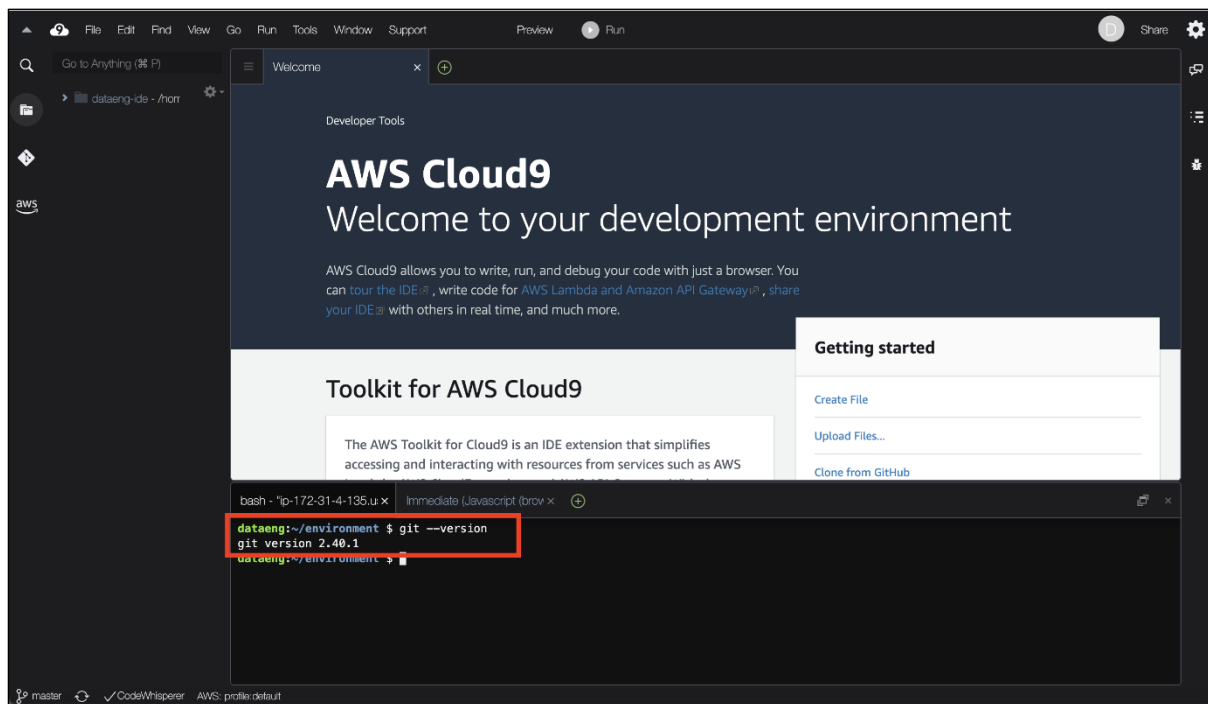
HISTORY

Use this section to view or modify the columnar business metadata of this asset.

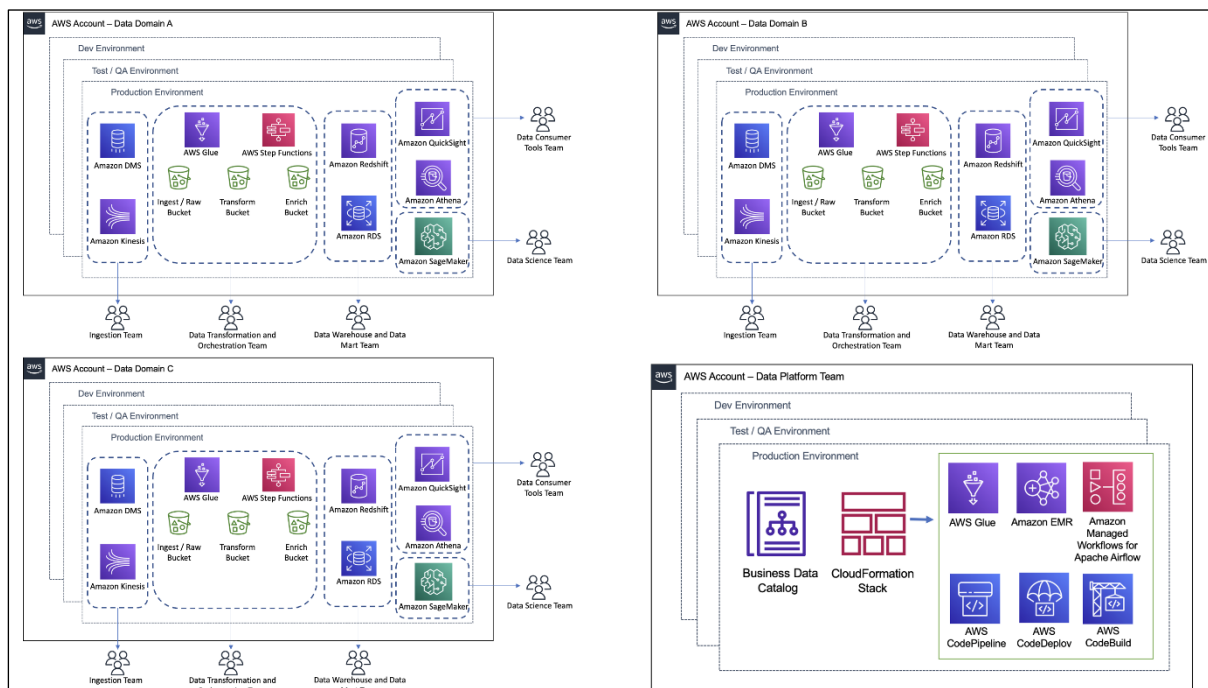
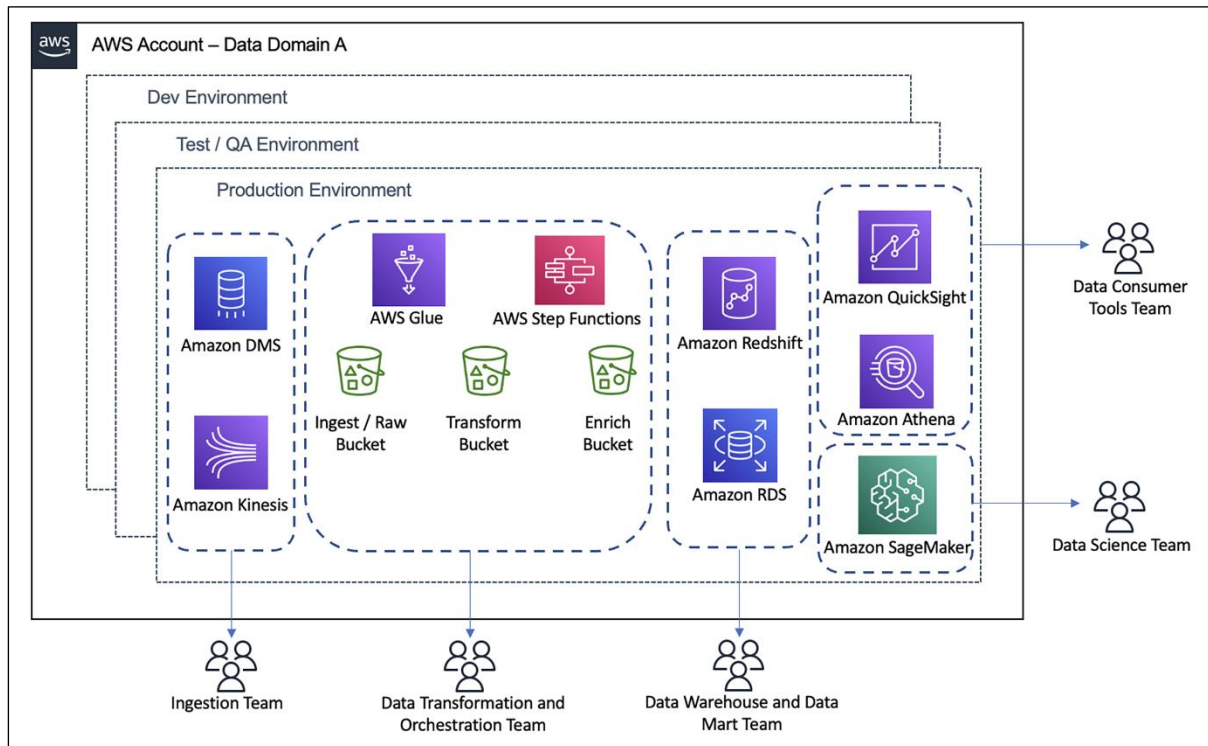
Filter

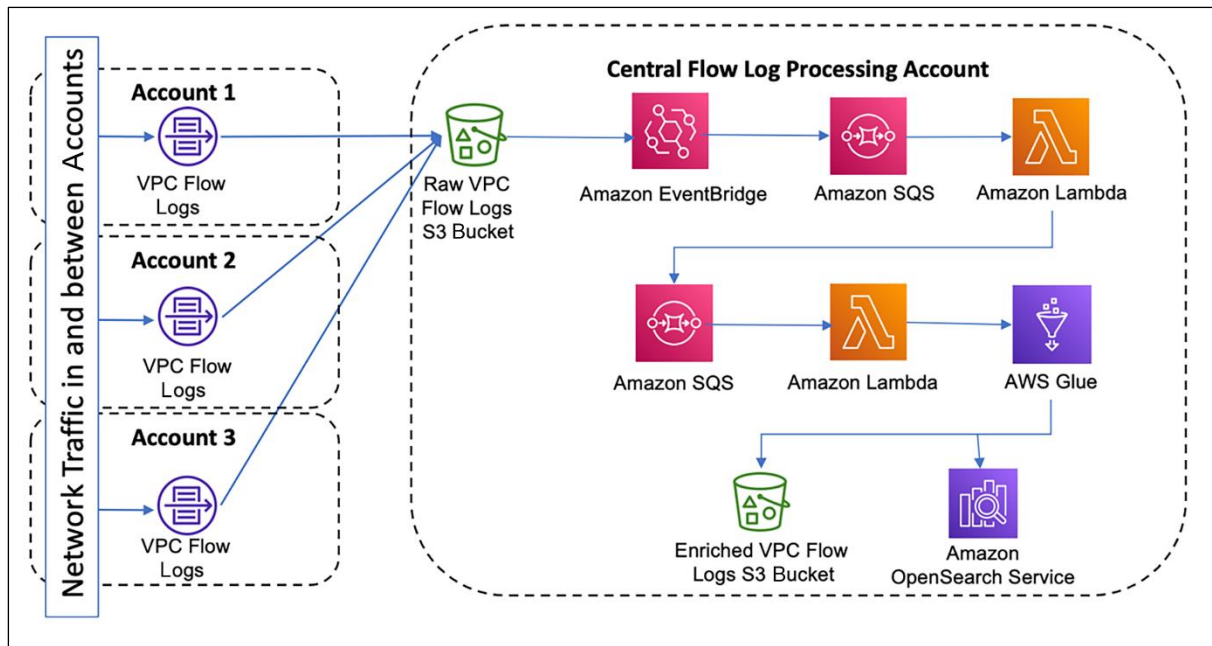
Name	Data type	Description	Terms	Primary key	Sort key
Category ID	bigint				
Category Name	string				
Film ID	bigint				

Chapter 16: Building a Modern Data Platform on AWS



Chapter 17: Wrapping Up the First Part of Your Learning Journey





AWS Billing Dashboard

AWS Billing Dashboard

Info

Page refresh time: Sunday, September 24, 2023 at 12:15:33 PM EDT

Unhide

AWS summary

Info

Viewing an overview of your AWS costs.

Current month's total forecast

Info

USD 16.00

Total number of active services

15

Current MTD balance

USD 12.69

Total number of active AWS accounts

1

Prior month for the same period with trend

USD 9.54

↑ 33.0%

Total number of active AWS Regions

5

Highest cost

Info

Viewing highest service spend.

Highest cost type

Highest service spend

Service name

QuickSight

Trend compared to prior month

↑ 3.2%

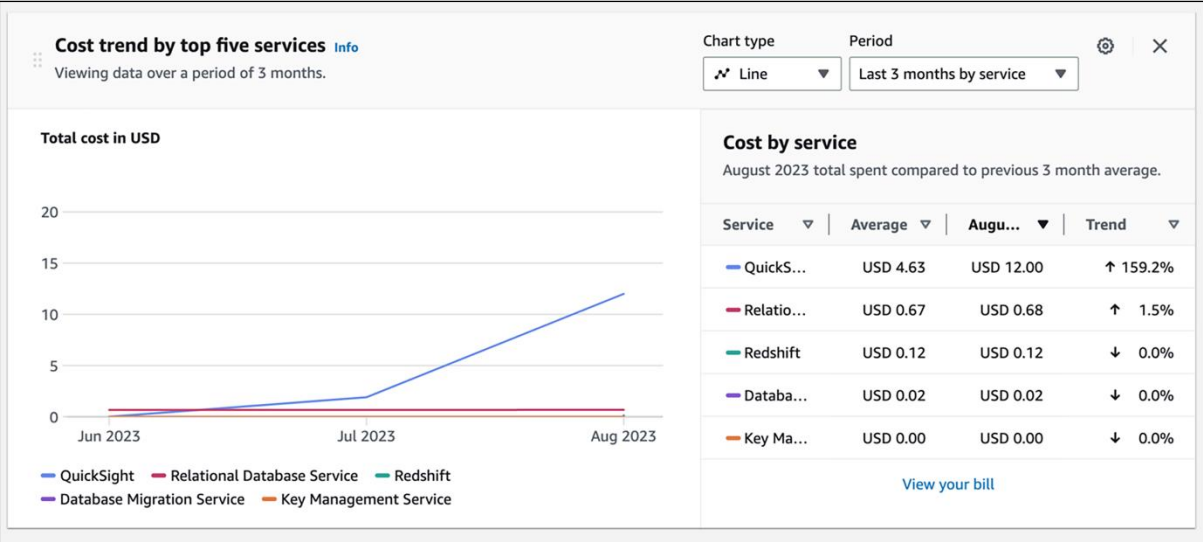
Current MTD balance

USD 9.32

Prior month for the same period

USD 9.03

View your bill



Amazon Web Services, Inc. charges by service [Info](#) [Expand all](#)

Total active services: **15** Total pre-tax service charges in USD: **USD 12.69**

Filter by service name or region name

Description	Usage Quantity	Amount in USD
QuickSight		USD 9.32
Glue		USD 2.10
Redshift		USD 0.70
US East (N. Virginia)		USD 0.35
Amazon Redshift USE1-RMS:Serverless		USD 0.35
Storage charges with Redshift managed storage	14.439 GB-Mo	USD 0.35
US East (Ohio)		USD 0.35
Amazon Redshift USE2-RMS:Serverless		USD 0.35
Storage charges with Redshift managed storage	14.518 GB-Mo	USD 0.35

[Option+S]

Introducing the new AWS account page experience
We've redesigned the AWS account page. [Let us know what you think.](#)

[AWS Billing](#) > Account

Account [Info](#)

This page now uses granular permission.
Legacy format permissions have been deprecated and will no longer be supported. This page has been updated based on your granular permissions. If you're missing permissions, ask your administrator to update your granular permissions using the policies tool. [Learn more.](#)

Account settings

Account ID: 42

- Account
- Organization
- Service Quotas
- Billing Dashboard
- Security credentials

[Sign out](#)

[Edit](#)

Close account

Please review this [important account closure guidance](#). Specifically:

- **Agreement termination:** Closing your account will serve as your notice of termination of the [AWS Customer Agreement](#) (or any other AWS agreement governing this account) for this account.
- **Billing:** You remain responsible for all [outstanding fees and charges](#), including [this month's usage](#) and [active subscriptions](#) (such as [Reserved Instances](#)).
- **Reactivation:** You may reopen your AWS account for 90 days after closure. If you reopen your account, you may be charged for any [active](#) resources. After 90 days your account will be permanently closed, any remaining content will be deleted, and unused credits will be lost.
- **GovCloud:** Closing this account will also close any linked [GovCloud accounts](#).

If you are experiencing unexpected charges, review how to [troubleshoot common changes](#), including Free Tier changes, unwanted resources and unauthorized activity, without closing your account.

Close account