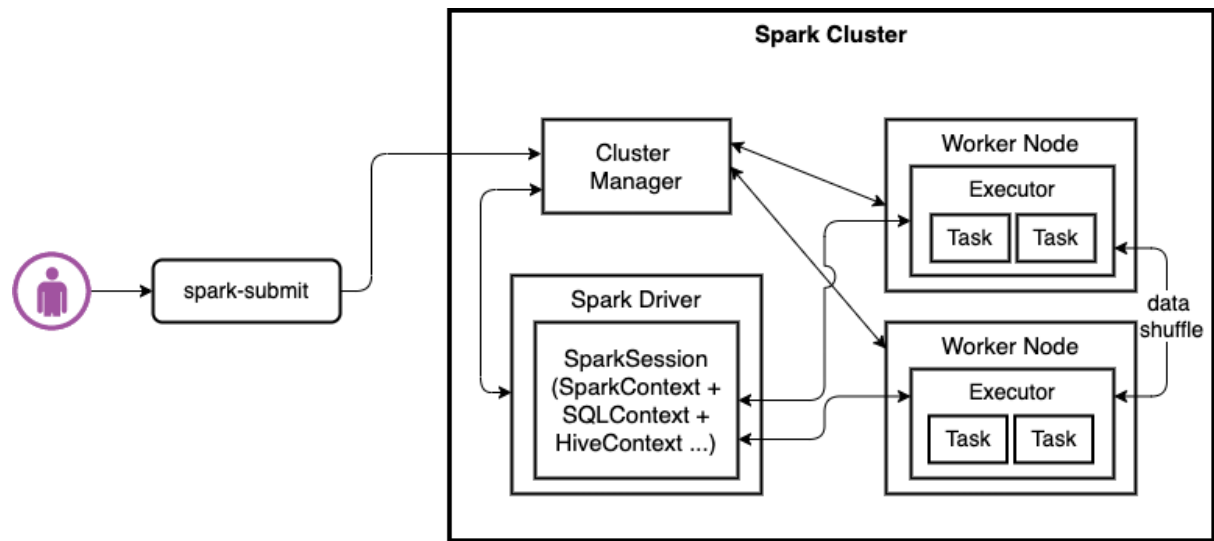


## Chapter 1: Data Management – Introduction and Concepts



## Chapter 2: Introduction to Important AWS Glue Features

```
# col_name      data_type      comment
id              bigint
country        struct<name:string,capital:string,countrycode:string,phonecode:string>

# Partition Information
# col_name      data_type      comment
currency       string
```

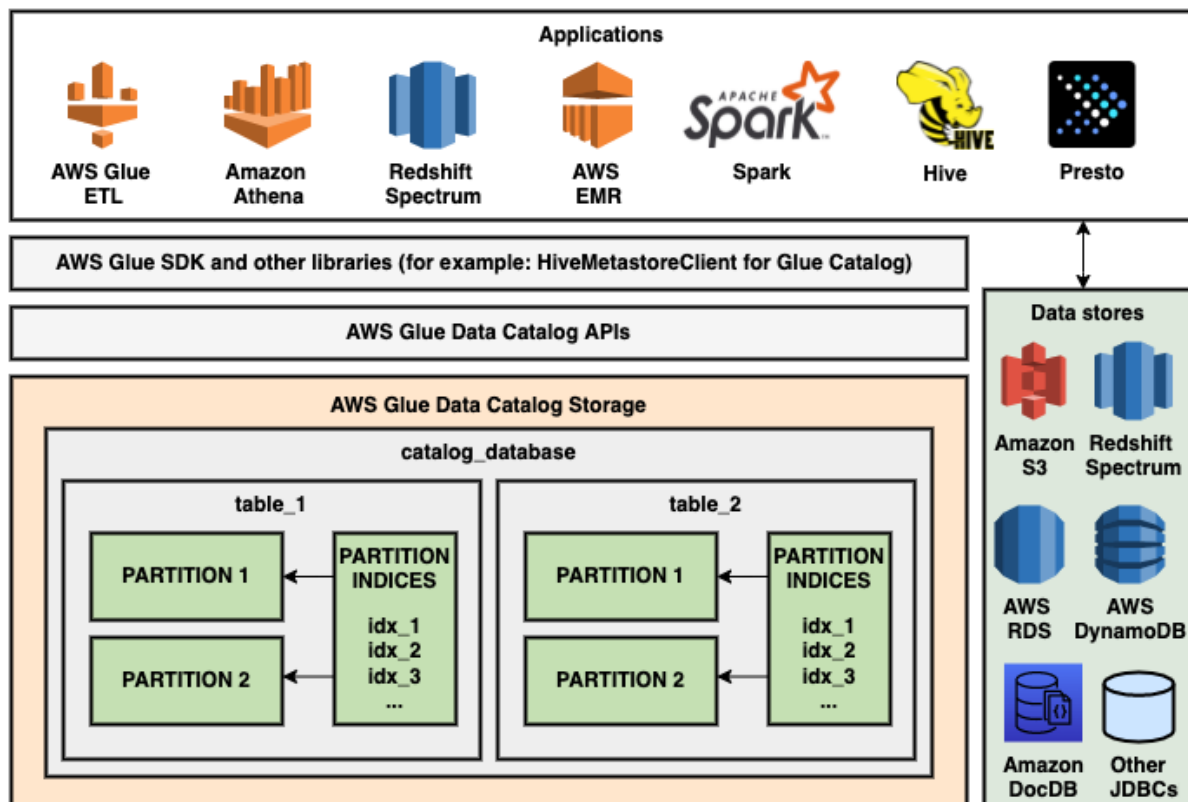
Table Schema

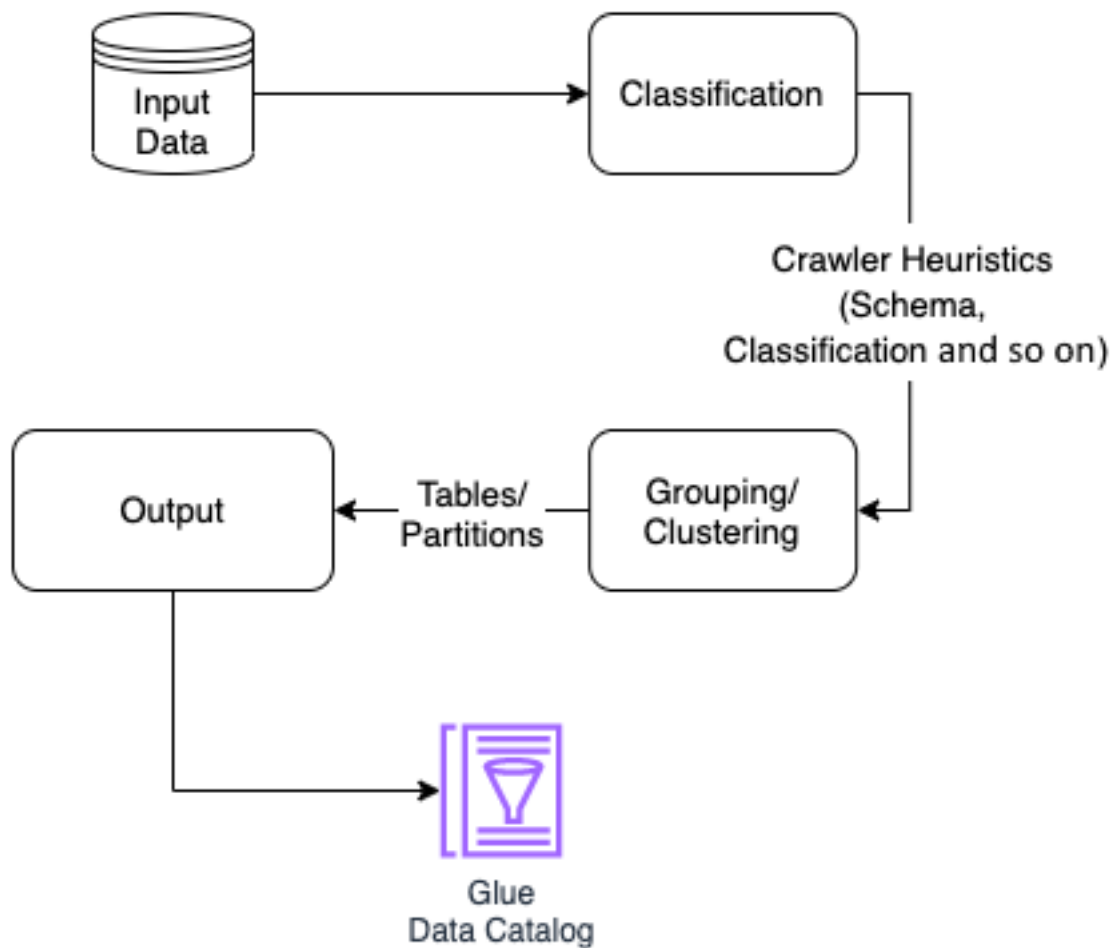
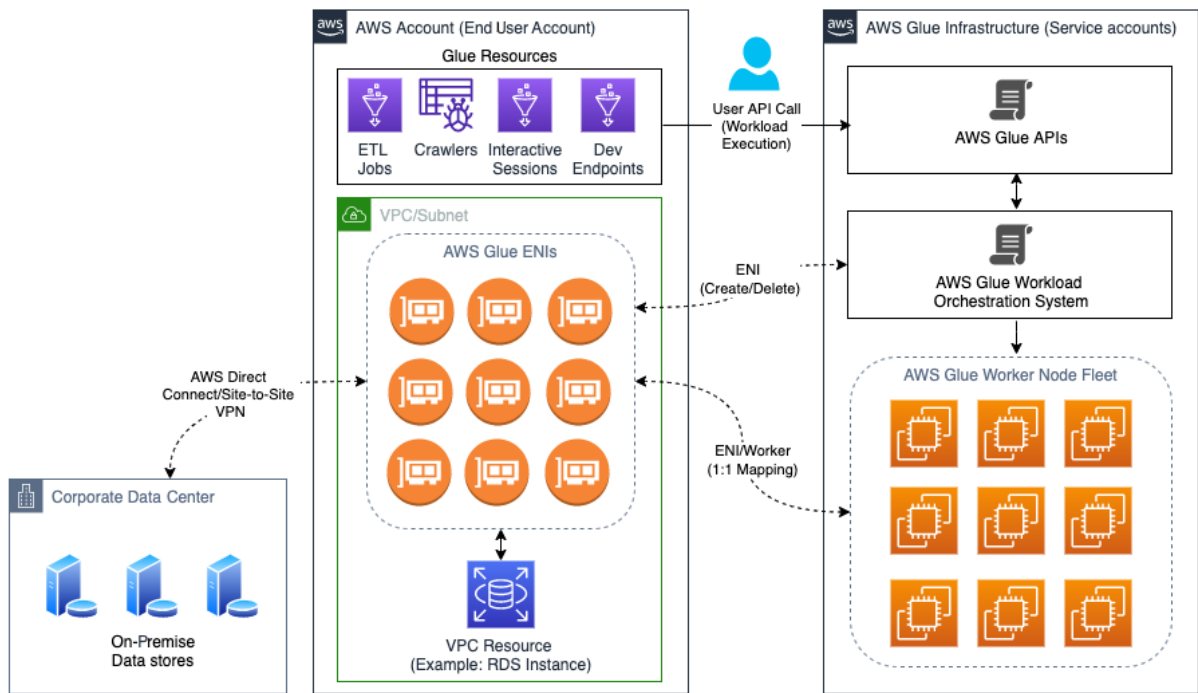
```
# Detailed Table Information
Database:      default
Owner:         hadoop
CreateTime:    Tue Jul 12 22:01:56 UTC 2022
LastAccessTime: UNKNOWN
Protect Mode:  None
Retention:     0
Location:      s3://BUCKET_NAME/S3_PREFIX
Table Type:    EXTERNAL_TABLE
Table Parameters:
  EXTERNAL          TRUE
  classification    parquet
  compressionType   none
  transient_lastDdlTime 1657663316
  typeOfData        file
```

Data format,  
Location and  
other information

```
# Storage Information
SerDe Library:  org.apache.hadoop.hive ql.io.parquet.serde.ParquetHiveSerDe
InputFormat:    org.apache.hadoop.hive ql.io.parquet.MapredParquetInputFormat
OutputFormat:   org.apache.hadoop.hive ql.io.parquet.MapredParquetOutputFormat
Compressed:     No
Num Buckets:    -1
Bucket Columns: ☐
Sort Columns:   ☐
Storage Desc Params:
  serialization.format 1
```

Input/Output format,  
Serialization and  
Deserialization Library (SerDe)  
information





## Chapter 3: Data Ingestion

```
mysql> SELECT argument FROM mysql.general_log WHERE argument LIKE '%city%' AND  
command_type = 'Prepare' ORDER BY argument;
```

argument
SELECT * FROM (select * from city WHERE ID % 10 = 0) as city
SELECT * FROM (select * from city WHERE ID % 10 = 0) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 1) as city
SELECT * FROM (select * from city WHERE ID % 10 = 1) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 2) as city
SELECT * FROM (select * from city WHERE ID % 10 = 2) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 3) as city
SELECT * FROM (select * from city WHERE ID % 10 = 3) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 4) as city
SELECT * FROM (select * from city WHERE ID % 10 = 4) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 5) as city
SELECT * FROM (select * from city WHERE ID % 10 = 5) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 6) as city
SELECT * FROM (select * from city WHERE ID % 10 = 6) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 7) as city
SELECT * FROM (select * from city WHERE ID % 10 = 7) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 8) as city
SELECT * FROM (select * from city WHERE ID % 10 = 8) as city WHERE 1=0
SELECT * FROM (select * from city WHERE ID % 10 = 9) as city
SELECT * FROM (select * from city WHERE ID % 10 = 9) as city WHERE 1=0

20 rows in set (0.20 sec)



# Chapter 4: Data Preparation

The screenshot illustrates the process of creating a recipe for data preparation in Data Studio. It is divided into two main sections: the top section shows the initial dataset view, and the bottom section shows the recipe configuration.

**Top Section: Dataset View**

The dataset is named "glue-UN-general-assembly-resolutions" and is sampled from the "resolution" dataset. It contains 15 columns and 500 rows. The columns are: # assembly\_session, # vote\_id, ABC resolution, # amendment, # colonization, # human\_rights, and # isra. The interface includes a sidebar with navigation options like DATASETS, PROJECTS, RECIPES, DQ RULES, JOBS, and WHAT'S NEW. The main area displays a table of data with various statistics (Distinct, Unique, Total, Min, Median, Mean, Mode, Max) and a histogram for each column.

**Bottom Section: Recipe Configuration**

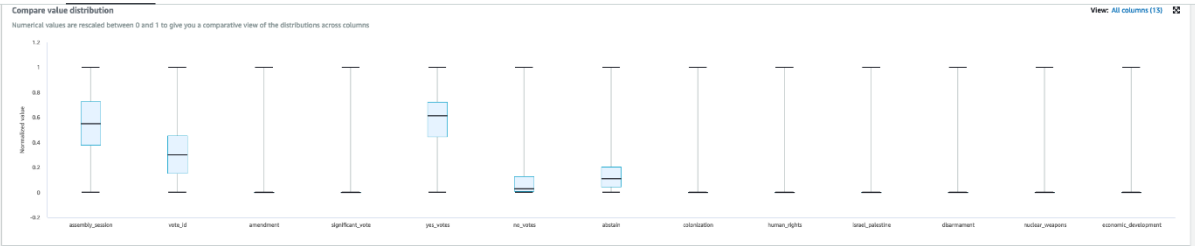
The recipe is named "glue-UN-general-assembly-resolutions-recipe" and is in the "Working version" state. It shows the applied steps (3) and the source dataset. The steps are:

1. Delete empty rows with missing values in amendment
2. Fill with last valid value
3. Fill with most frequent value

The recipe is configured to use the "GRID" view. The "Column details" panel on the right shows the statistics for the columns and provides recommendations for handling missing values, such as "Delete rows with missing values" or "Fill with last valid value".

**Annotations:**

- 1: Points to the "RECIPE" button in the top right corner.
- 2: Points to the "Publish" button in the recipe configuration panel.
- 3: Points to the "Download as JSON" button in the recipe configuration panel.
- 4: Points to the "Create Job" button in the top right corner.



**Columns summary (15)**

Column name	assembly_session	vote_id	resolution	amendment	vote_date	significant_vote	yes_votes	no_votes	abstain
Column type	# Integer	# Integer	# String	# Integer	# String	# Integer	# Integer	# Integer	# Integer
Data quality	100% Valid	100% Valid	97% Valid 0% Invalid 3% Missing	92% Valid 0% Invalid 48% Missing	100% Valid 0% Invalid	100% Valid 0% Invalid	100% Valid 0% Invalid	100% Valid 0% Invalid	100% Valid 0%
Distribution									
Box plot									
Total valid	5429 (100%)	5429 (100%)	5278 (97%)	2544 (52%)	5429 (100%)	5429 (100%)	5429 (100%)	5429 (100%)	5429 (100%)
Total missing	0 (0%)	0 (0%)	151 (3%)	2385 (48%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Distinct values	70 (1%)	5429 (100%)	5270 (97%)	2 (0%)	780 (14%)	2 (0%)	189 (3%)	77 (1%)	95 (2%)

**Data quality rules (2)**

Expand all | Collapse all | Find

ALL (2) | SUCCEEDED (1) | FAILED (1) | ERROR (0) | DISABLED (0)

☒ **united-nations-resolutions-ruleset** 2 rules

- Check Dataset For Duplicate Rows**  
Check if **dataset** has duplicate rows count <= 0
- Check All Columns For Missing Values**  
Check if **all columns** has missing values == 0%  

SUCCEEDED 7 columns | FAILED 8 columns | ERROR 0 columns

**Check All Columns For Missing Values**

Check if **all columns** has missing values == 0%

**Failed**

SUCCEEDED 7 columns | FAILED 8 columns | ERROR 0 columns

**Columns (15)**

Find

ALL (15) | SUCCEEDED (7) | FAILED (8) | ERROR (0)

- # no\_votes
- # vote\_id
- # colonization
- # abstain
- # vote\_date
- # significant\_vote
- # israel\_palestine
- # economic\_development

Visual Script Job details Runs Schedules

Source

Transform

Target

Undo

Redo

Remove

Node properties

Transform

Output schema

Data preview

Data source - Data Catalog  
hudi\_nyc\_yellow\_trip...Transform - Filter  
FilterTransform - Custom code  
Custom Transform

Schema Info Edit

Key

vendorid

tpep\_pickup\_datetime

tpep\_dropoff\_datetime

passenger\_count

trip\_distance

ratecodeid

store\_and\_fwd\_flag

pulocationid

dolocationid

employees.val.name	company	employees	index	employees.val.email	id
foo2	DummyCorp2	2	0	foo@company2.com	2
bar2	DummyCorp2	2	1	bar@company2.com	2
foo3	DummyCorp3	3	0	foo@company3.com	3
bar3	DummyCorp3	3	1	bar@company3.com	3
foo1	DummyCorp1	1	0	foo@company1.com	1
bar1	DummyCorp1	1	1	bar@company1.com	1

## Chapter 5: Designing Data Layouts

Visual

Script

Job details

Runs

Schedules

Source

Transform

Target

Undo

Redo

Remove

Data source - S3 bucket

S3 bucket

Transform - ApplyMapping

ApplyMapping

Data target - S3 bucket

S3 bucket

Script (Locked)

Info

```
10 glueContext = GlueContext(sc)
11 spark = glueContext.spark_session
12 job = Job(glueContext)
13 job.init(args["JOB_NAME"], args)
14
15 # Script generated for node S3 bucket
16 S3bucket_node1 = glueContext.create_dynamic_frame.from_options(
17     format_options={"multiline": False},
18     connection_type="s3",
19     format="json",
20     connection_options={"paths": ["s3://your-bucket-path"], "recurse": True},
21     transformation_ctx="S3bucket_node1",
22 )
23
24 # Script generated for node ApplyMapping
25 ApplyMapping_node2 = ApplyMapping.apply(
26     frame=S3bucket_node1,
27     mappings=[
28         ("product_name", "string", "product_name", "string"),
29         ("price", "string", "price", "string"),
30         ("customer_id", "int", "customer_id", "int"),
31         ("order_id", "string", "order_id", "string"),
32         ("datetime", "string", "datetime", "string"),
33         ("category", "string", "category", "string"),
34     ],
35     transformation_ctx="ApplyMapping_node2",
36 )
37
38 # Script generated for node S3 bucket
39 S3bucket_node3 = glueContext.write_dynamic_frame.from_options(
40     frame=ApplyMapping_node2,
41     connection_type="s3",
42     format="glueparquet",
43     connection_options={
44         "path": "s3://your-target-bucket-and-path/",
45         "partitionKeys": [],
46     },
47     format_options={"compression": "snappy"},
48     transformation_ctx="S3bucket_node3",
49 )
50
51 job.commit()
```

Node properties

Data target properties - S3

Output schema

Data pr

Format

Parquet

Compression Type

Snappy

S3 Target Location

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

s3://your-target-bucket-and-path/

Data Catalog update options

Info

☒ Do not update the Data Catalog

☐ Create a table in the Data Catalog and on subsequent runs, update the schema and a

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

Partition keys - optional

Add partition keys.

Partition (0)

category

Add a partition key

**Data management and security**

**Governance choice cannot be changed after table creation**

Governed tables support ACID (atomic, consistent, isolated, durable) transactions to guarantee data integrity with concurrent changes across multiple tables. Governed tables also come with automatic data layout optimizations by compacting small and delta files in order to maximize query performance.

☒ **Enable governed data access and management**

Check this option to enable transactional workloads and automatic data compaction.

☐ **Automatic compaction**

Enable this option to turn on automatic compaction on a governed table.



## Create lifecycle rule

### Lifecycle rule configuration

#### Lifecycle rule name

Up to 255 characters

#### Choose a rule scope

- ☒ Limit the scope of this rule using one or more filters
- ☐ Apply to all objects in the bucket

#### Filter type

You can filter objects by prefix, object tags, object size, or whatever combination suits your usecase.

##### Prefix

Add filter to limit the scope of this rule to a single prefix.

Don't include the bucket name in the prefix. Using certain characters in key names can cause problems with some applications and protocols. [Learn more](#)

#### Object tags

You can limit the scope of this rule to the key/value pairs added below.

Add tag

#### Object size

You can limit the scope of this rule to apply to objects based on their size. For example, you can filter out objects that might not be cost effective to transition to Glacier Flexible Retrieval (formerly Glacier) because of per-object fees.

- ☐ Specify minimum object size
- ☐ Specify maximum object size

### Lifecycle rule actions

Choose the actions you want this rule to perform. Per-request fees apply. [Learn more](#) or see [Amazon S3 pricing](#)

- ☐ Move current versions of objects between storage classes
- ☐ Move noncurrent versions of objects between storage classes
- ☐ Expire current versions of objects
- ☐ Permanently delete noncurrent versions of objects
- ☐ Delete expired object delete markers or incomplete multipart uploads
- These actions are not supported when filtering by object tags or object size.

# Chapter 6: Data Management

```
root
|-- id: long
|-- index: int
|-- entries.val.id: int
|-- entries.val.values.k1: string
|-- entries.val.values.k2: string
```

id	index	entries.val.id	entries.val.values.k1	entries.val.values.k2
1	0	1	aaa	bbb
1	1	2	ccc	ddd

```
root
|-- count: integer (nullable = true)
|-- entries: long (nullable = true)
|-- id: long (nullable = true)
|-- index: integer (nullable = true)
|-- entries.val.id: integer (nullable = true)
|-- entries.val.values.k1: string (nullable = true)
|-- entries.val.values.k2: string (nullable = true)
```

count	entries	id	index	entries.val.id	entries.val.values.k1	entries.val.values.k2
2	1	1	0	1	aaa	bbb
2	1	1	1	2	ccc	ddd

product_id	product_name	category	price
11	Introduction to C...	Ebooks	15
12	Best practices on...	Ebooks	25
21	Data Quest	Video games	30
22	Final Shooting	Video games	20

uid	customer_name	email	phone
A103	Barbara Gordon	gordon@example.com	117.835.2584
A042	Rebecca Thompson	thompson@example.net	001-469-964-3897x9041
A805	Rachel Gilbert	gilbert@example.com	001-510-198-4613x23986
A404	Tanya Fowler	tanya@example.net	(067)150-0263

product_id	product_by	purchased_at
21	A042	2022-03-30T01:30:00Z
22	A805	2022-04-01T02:00:00Z
11	A103	2022-04-21T11:40:00Z
12	A404	2022-04-28T08:20:00Z

ticket_id	purchased_by	purchased_at
216	A042	2022-03-30T01:30:00Z
217	A805	2022-04-01T02:00:00Z
218	A103	2022-04-21T11:40:00Z

customer_name	purchased_at
Rachel Gilbert	2022-04-01T02:00:00Z
Barbara Gordon	2022-04-21T11:40:00Z
Tanya Fowler	2022-04-28T08:20:00Z

product_name	category	price	customer_name	email	phone	purchased_at
Data Quest	Video games	30	Rebecca Thompson	thompson@example.net	001-469-964-3897x...	2022-03-30T01:30:00Z
Final Shooting	Video games	20	Rachel Gilbert	gilbert@example.com	001-510-198-4613x...	2022-04-01T02:00:00Z
Introduction to C...	Ebooks	15	Barbara Gordon	gordon@example.com	117.835.2584	2022-04-21T11:40:00Z
Best practices on...	Ebooks	25	Tanya Fowler	tanya@example.net	(067)150-0263	2022-04-28T08:20:00Z

product_name	category	price	customer_name	email	phone	purchased_at
Final Shooting	Video games	20	Rachel Gilbert	gilbert@example.com	001-510-198-4613x...	2022-04-01T02:00:00Z
Introduction to C...	Ebooks	15	Barbara Gordon	gordon@example.com	117.835.2584	2022-04-21T11:40:00Z
Best practices on...	Ebooks	25	Tanya Fowler	tanya@example.net	(067)150-0263	2022-04-28T08:20:00Z

product_name	category	price	customer_name	email	phone	purchased_at
Data Quest	Video games	30	Rebecca Thompson	thompson@example.net	***-***-***-****x...	2022-03-30T01:30:00Z
Final Shooting	Video games	20	Rachel Gilbert	gilbert@example.com	***-***-***-****x...	2022-04-01T02:00:00Z
Introduction to C...	Ebooks	15	Barbara Gordon	gordon@example.com	***.***.****	2022-04-21T11:40:00Z
Best practices on...	Ebooks	25	Tanya Fowler	tanya@example.net	(***)***-****	2022-04-28T08:20:00Z

product_name	category	price	customer_name	email	phone	purchased_at
Data Quest	Video games	30	Rebecca Thompson	97cf4c3dfff3a1245...	***-***-***-****x...	2022-03-30T01:30:00Z
Final Shooting	Video games	20	Rachel Gilbert	2857f8c8a7b8c1b7f...	***-***-***-****x...	2022-04-01T02:00:00Z
Introduction to C...	Ebooks	15	Barbara Gordon	b8ba2a41ce2d45a99...	***.***.****	2022-04-21T11:40:00Z
Best practices on...	Ebooks	25	Tanya Fowler	b822364443d400f56...	(***)***-****	2022-04-28T08:20:00Z

check	check_level	check_status	constraint	constraint_status	constraint_message
Review Check	Warning	Warning	SizeConstraint(Size=None)	Success	
Review Check	Warning	Warning	CompletenessConstraint(Completeness(id=None))	Success	
Review Check	Warning	Warning	UniquenessConstraint(Uniqueness(List(id),None))	Success	
Review Check	Warning	Warning	CompletenessConstraint(Completeness(productName=None))	Failure	Value: 0.8 does not meet the constraint requirement!
Review Check	Warning	Warning	ComplianceConstraint(Compliance(priority contained in high,medium,low,'priority' IS NULL OR 'priority' IN ('high','medium','low'),None))	Success	
Review Check	Warning	Warning	ComplianceConstraint(Compliance(numViews is non-negative,COALESCE(CAST(numViews AS DECIMAL(20,10)),0.0) >= 0,None))	Success	
Review Check	Warning	Warning	containsURL(description)	Failure	Value: 0.4 does not meet the constraint requirement!
Review Check	Warning	Warning	ApproxQuantileConstraint(ApproxQuantile(numViews,0.5,0.01,None))	Success	

customer_name	email	phone
Barbara Gordon	gordon@example.com	117.835.2584
Gordon, Barbara	gordon@example.com	117-835-2584
Rebecca Thompson	thompson@example.net	001-469-964-3897x9041
Rachel Gilbert	gilbert@example.com	001-510-198-4613x23986
Gilbert, R.	gilbert@example.com	
Tanya Fowler	tanya@example.net	(067)150-0263

customer_name	email	phone	match_id
Barbara Gordon	gordon@example.com	117.835.2584	1
Gordon, Barbara	gordon@example.com	117-835-2584	1
Rebecca Thompson	thompson@example.net	001-469-964-3897x9041	2
Rachel Gilbert	gilbert@example.com	001-510-198-4613x23986	3
Gilbert, R.	gilbert@example.com		3
Tanya Fowler	tanya@example.net	(067)150-0263	4

product_id	product_name	category	price
11	Introduction to Cloud	Ebooks	15
12	Best practices on data lakes	Ebooks	25
21	Data Quest	Video games	30
22	Final Shooting	Video games	20

uid	customer_name	email	phone
A103	Barbara Gordon	gordon@example.com	117.835.2584
A042	Rebecca Thompson	thompson@example.net	001-469-964-3897x9041
A805	Rachel Gilbert	gilbert@example.com	001-510-198-4613x23986
A404	Tanya Fowler	tanya@example.net	(067)150-0263

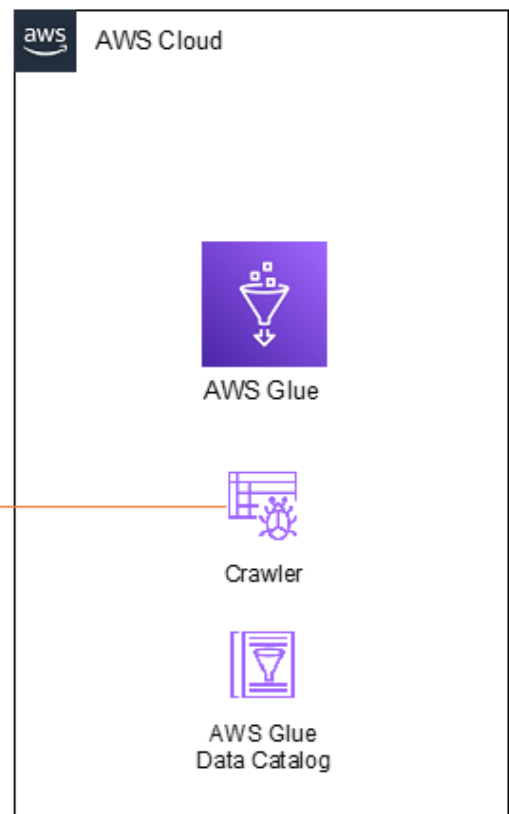
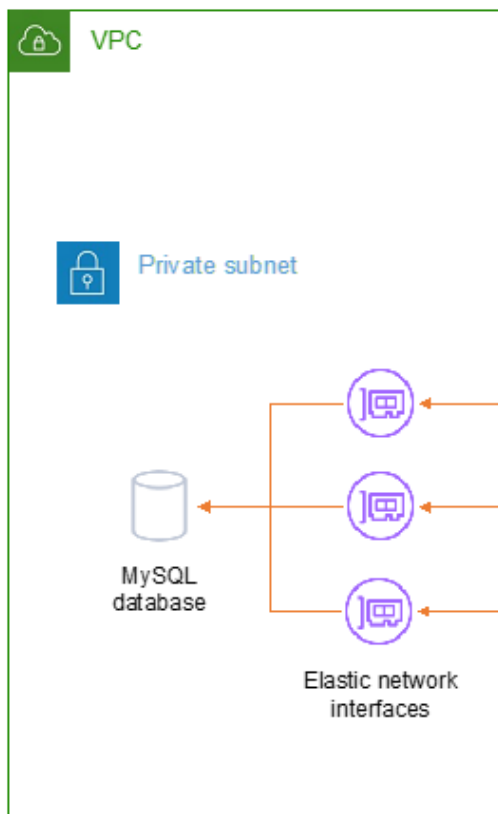
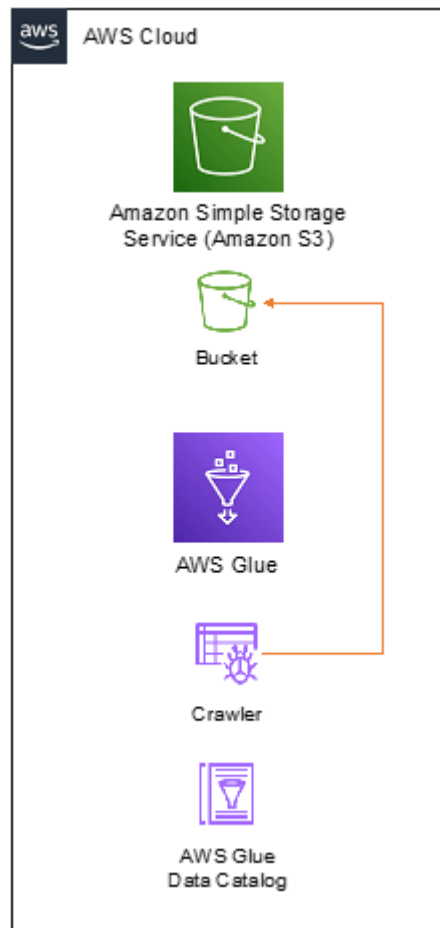


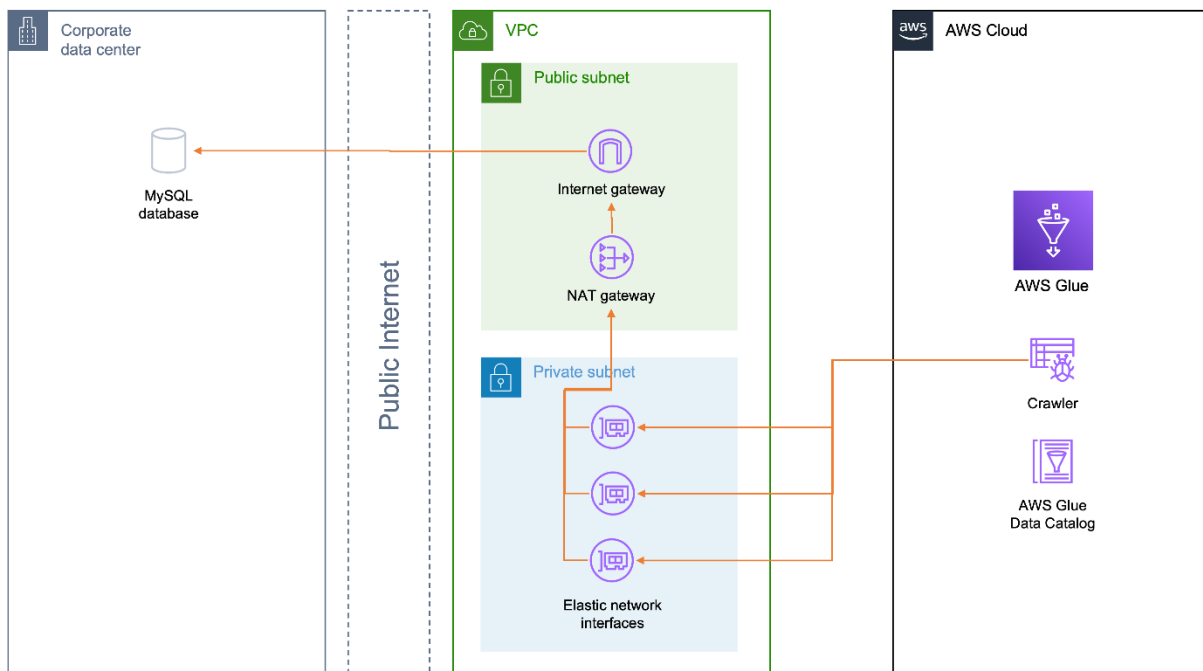
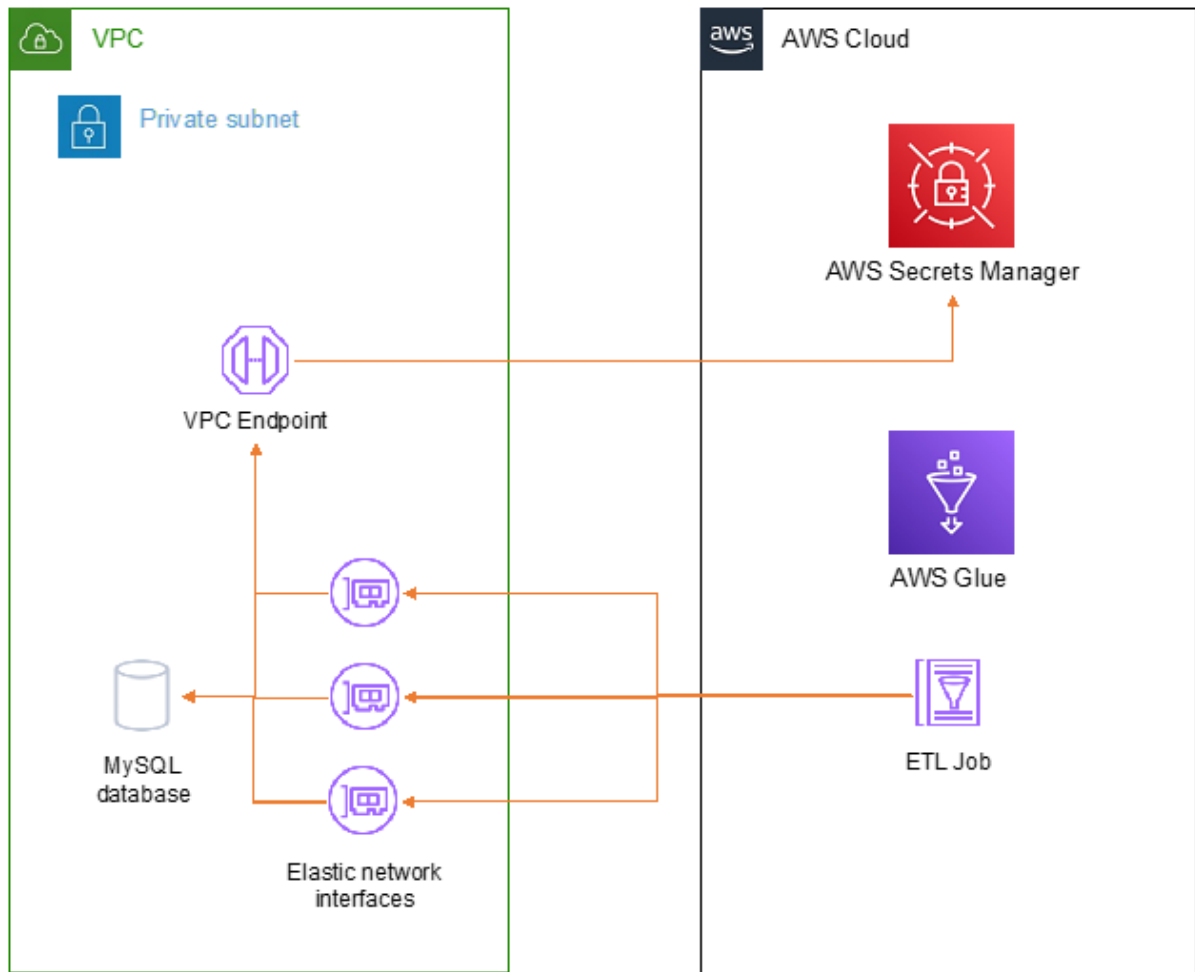
## Chapter 7: Metadata Management

No images...

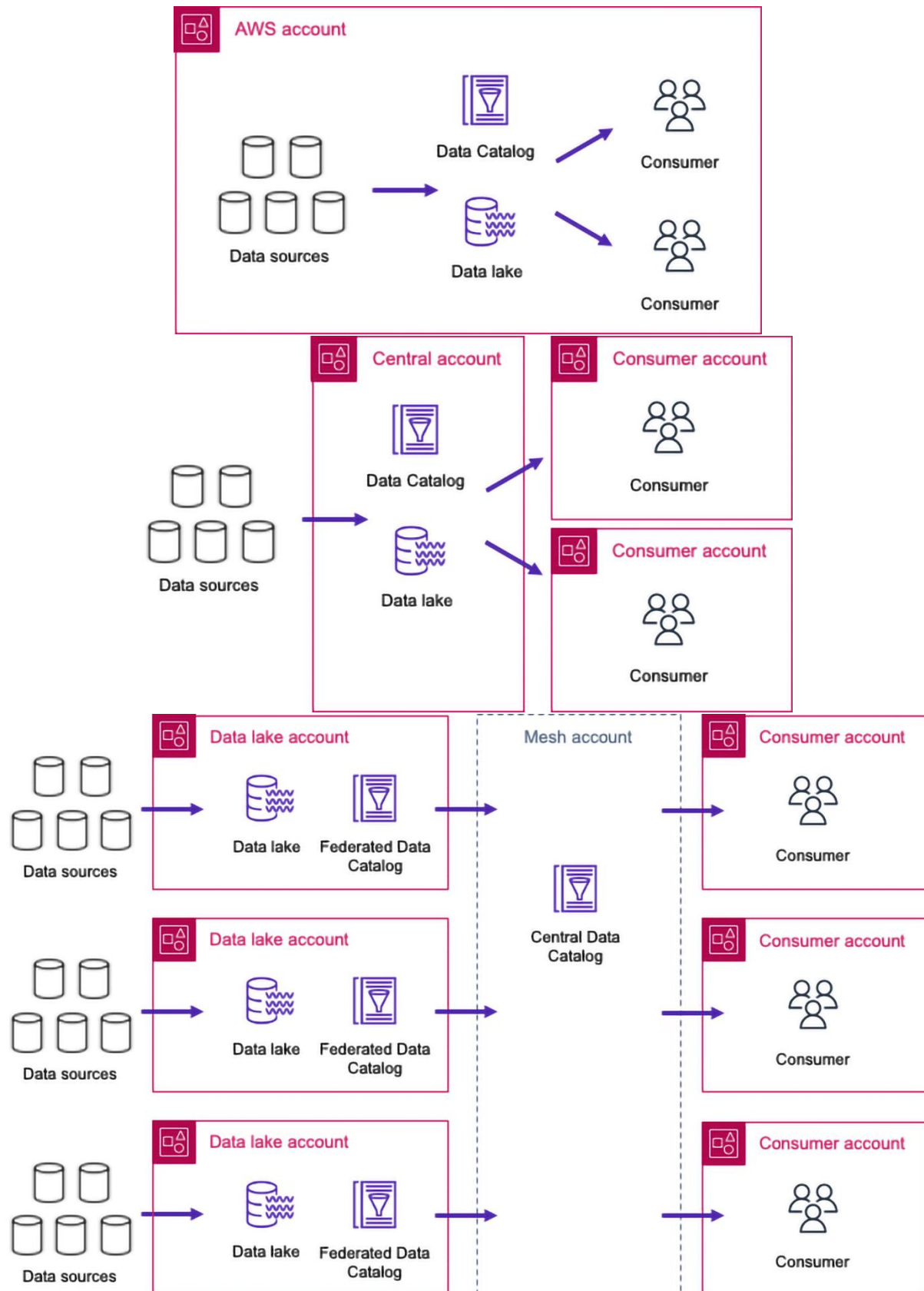
## Chapter 8: Data Security







## Chapter 9: Data Sharing



EditorRecent queriesSaved queriesSettings

Workgroupprimary

Data

Data Source

AwsDataCatalog

Database

simple\_datalake

Tables and views

Create

Filter tables and views

Tables (1)

<1>

pcs

:

Views (0)

<1>

Query 1

1 SELECT \* FROM simple\_datalake.pcs;

SQLLn 1, Col 34

Run againCancelSaveClearCreate

Completed

Time in queue: 0.177 secRun time: 0.484 secData scanned: 0.78 KB

Results (4)

CopyDownload results

Search rows

#	product_name	category	price	purchased_at	customer_name	email	phone
1	Introduction to Cloud	Ebooks	15	2022-04-21T11:40:00Z	Barbara Gordon	gordon@example.com	117.835.2584
2	Best practices on data lakes	Ebooks	25	2022-04-28T08:20:00Z	Tanya Fowler	tanya@example.net	(067)150-0263
3	Data Quest	Video games	30	2022-03-30T01:30:00Z	Rebecca Thompson	thompson@example.net	001-469-964-3897x9041
4	Final Shooting	Video games	20	2022-04-01T02:00:00Z	Rachel Gilbert	gilbert@example.com	001-510-198-4613x23986

EditorRecent queriesSaved queriesSettings

Workgroupprimary

Data

Data Source

producer\_catalog

Database

simple\_datalake

Tables and views

Create

Filter tables and views

Tables (1)

<1>

pcs

:

Views (0)

<1>

Query 1 × Query 2 ×

1 SELECT \* FROM simple\_datalake.pcs;

SQLLn 1, Col 34

Run againCancelSaveClearCreate

Completed

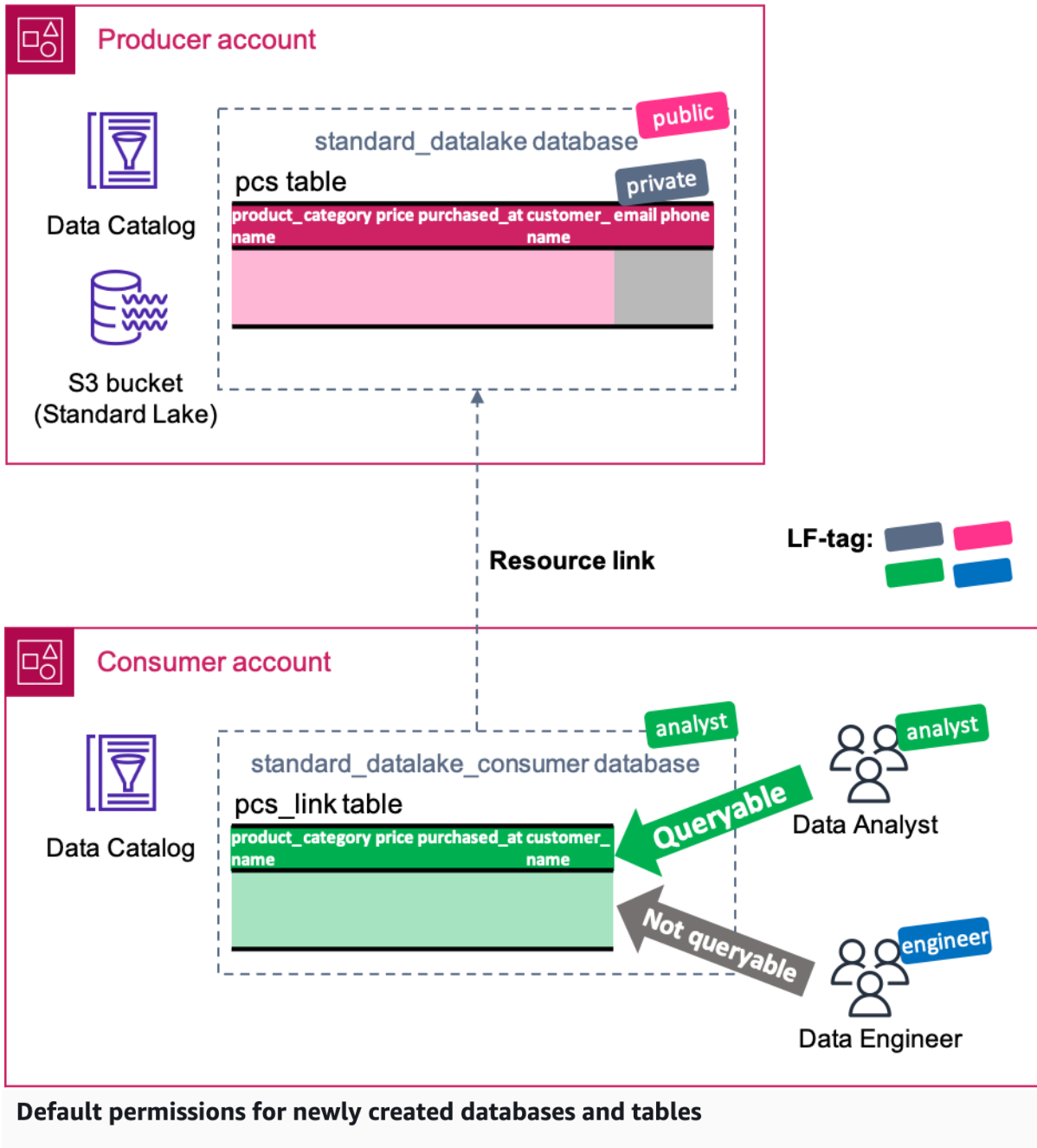
Time in queue: 0.515 secRun time: 0.892 secData scanned: 0.78 KB

Results (4)

CopyDownload results

Search rows

#	product_name	category	price	purchased_at	customer_name	email	phone
1	Introduction to Cloud	Ebooks	15	2022-04-21T11:40:00Z	Barbara Gordon	gordon@example.com	117.835.2584
2	Best practices on data lakes	Ebooks	25	2022-04-28T08:20:00Z	Tanya Fowler	tanya@example.net	(067)150-0263
3	Data Quest	Video games	30	2022-03-30T01:30:00Z	Rebecca Thompson	thompson@example.net	001-469-964-3897x9041
4	Final Shooting	Video games	20	2022-04-01T02:00:00Z	Rachel Gilbert	gilbert@example.com	001-510-198-4613x23986



These settings maintain existing AWS Glue Data Catalog behavior. You can still set individual permissions on databases and tables, which will take effect when you revoke the Super permission from IAMAllowedPrincipals. See [Changing Default Settings for Your Data Lake](#).

- ☐ Use only IAM access control for new databases
- ☐ Use only IAM access control for new tables in new databases

# Register location

## Amazon S3 location

Register an Amazon S3 path as the storage location for your data lake.

### Amazon S3 path

Choose an Amazon S3 path for your data lake.

*e.g.: s3://bucket/prefix/*

**Browse**

### Review location permissions - strongly recommended

Registering the selected location may result in your users gaining access to data already at that location. Before registering a location, we recommend that you review existing location permissions on resources in that location.

[Review location permissions](#)

### IAM role

To add or update data, Lake Formation needs read/write access to the chosen Amazon S3 path. Choose a role that you know has permission to do this, or choose the `AWSServiceRoleForLakeFormationDataAccess` service-linked role. When you register the first Amazon S3 path, the service-linked role and a new inline policy are created on your behalf. Lake Formation adds the first path to the inline policy and attaches it to the service-linked role. When you register subsequent paths, Lake Formation adds the path to the existing policy.

AWSServiceRoleForLakeFormationDataAccess ▼

 Do not select the service linked role if you plan to use EMR.

**Cancel**

**Register location**

## Add LF-Tag [Learn More](#)



LF-Tags have a key and one or more values that can be associated with data catalog resources. Tables automatically inherit from database LF-tags, and columns inherit from table LF-tags.

Example: Key = Confidentiality | Values = private, sensitive, public

### Key

Confidentiality

Key string must be less than 128 characters long, and cannot be changed once LF-tag is created.

### Values

Type a single value and select [Enter] or specify multiple values separated by commas.

**Add**

private ✕

public ✕

Enter up to 15 values; each value must be less than 256 characters long.

**Cancel**

**Add LF-tag**



Edit LF-Tags: standard\_datalake [Learn More](#)

LF-Tags

After they are associated with catalog resources, LF-Tags allow you to create scalable permissions.

Assigned keys

Q Confidentiality X

Values

public ▼

Remove

Assign new LF-Tag

You can add 49 more LF-tags.

Cancel

Save

Edit LF-Tags: review\_body [Learn More](#)

LF-Tags

After they are associated with catalog resources, LF-Tags allow you to create scalable permissions.

Inherited keys

Q Confidentiality

Values

private ▲

private

public (inherited)

Revert

Assign new LF-Tag

You can add 49 more LF-tags.

Cancel

Save

## Grant LF-tag permissions

Select the principals to grant permissions to, the LF-Tags to grant permissions on, and the specific set of permissions.

### Principals

☐ IAM users and roles  
Users or roles from this AWS account.

☐ SAML users and groups  
SAML users and group or QuickSight ARNs.

☒ External accounts  
AWS accounts or AWS organizations outside of this account.

#### AWS account or AWS organization

Enter one or more AWS account IDs or AWS organization IDs. Press Enter after each ID.

🔍 Choose AWS account ID or AWS organization ID

222222222222 ✕  
Account

### LF-Tags

#### LF-tag permission scope

Choose to grant permissions on all or a subset of LF-Tags.

##### Key

🔍 Confidentiality ✕

##### Values

Choose LF-tag values ▼

Remove

public ✕

Add LF-Tag

### Permissions

#### LF-tag permissions

Select the specific access permissions to grant.

☒ Describe

☐ Associate

#### Grantable permissions

Select the permissions that the grant recipient can grant to other principals.

☐ Describe

☐ Associate

Cancel

Grant

## Grant data permissions

### Principals

☐ **IAM users and roles**  
Users or roles from this AWS account.

☐ **SAML users and groups**  
SAML users and group or QuickSight ARNs.

☒ **External accounts**  
AWS accounts or AWS organizations outside of this account.

#### AWS account or AWS organization

Enter one or more AWS account IDs or AWS organization IDs. Press Enter after each ID.

222222222222 X  
Account

Granting data permissions to organizations is not supported when granting permissions by using LF-Tags.

### LF-Tags or catalog resources

☒ **Resources matched by LF-Tags (recommended)**  
Manage permissions indirectly for resources or data matched by a specific set of LF-Tags.

☐ **Named data catalog resources**  
Manager permissions for specific databases or tables, in addition to fine-grained data access.

Key

X

Values

Remove

public X

Add LF-Tag

### Database permissions

#### Database permissions

Choose specific access permissions to grant.

☐ Create table ☐ Alter ☐ Drop  
☒ Describe

#### Grantable permissions

Choose the permission that may be granted to others.

☐ Create table ☐ Alter ☐ Drop  
☒ Describe

☐ **Super**

This permission is the union of all the individual permissions to the left, and supersedes them.

☐ **Super**

This permission allows the principal to grant any of the permissions to the left, and supersedes those grantable permissions.

# Create resource link

## Table resource link details

Create a table resource link in the AWS Glue Data Catalog.

Resource link name

Name may contain letters (A-Z), numbers (0-9), hyphens (-), or underscores (\_), and must be less than 256 characters long.

Database

Resource link will be contained in this database.

Shared table

Enter or choose a shared table.

Shared table's database

Enter the database containing the shared table.

Shared table's owner ID

Enter the AWS account ID of the shared table owner.

[Cancel](#)[Create](#)

Editor

Recent queries

Saved queries

Settings

Workgroup primary

Data

Data Source

AwsDataCatalog

Database

standard\_datalake\_consumer

Tables and views

Create

Filter tables and views

Tables (1)

pcs\_link

product\_name string

category string

price int

purchased\_at string

customer\_name string

Views (0)

Query 1

1 SELECT \* FROM standard\_datalake\_consumer.pcs\_link

SQL Ln 1, Col 50

Run again Cancel Save Clear Create

Completed Time in queue: 0.12 sec Run time: 2.122 sec Data scanned: 0.78 KB

Results (4) Copy Download results

Search rows

#	product_name	category	price	purchased_at	customer_name
1	Introduction to Cloud	Ebooks	15	2022-04-21T11:40:00Z	Barbara Gordon
2	Best practices on data lakes	Ebooks	25	2022-04-28T08:20:00Z	Tanya Fowler
3	Data Quest	Video games	30	2022-03-30T01:30:00Z	Rebecca Thompson
4	Final Shooting	Video games	20	2022-04-01T02:00:00Z	Rachel Gilbert

Add LF-Tag

[Learn More](#)

✕

LF-Tags have a key and one or more values that can be associated with data catalog resources. Tables automatically inherit from database LF-tags, and columns inherit from table LF-tags.  
Example: Key = Confidentiality | Values = private, sensitive, public

Key

Division

Key string must be less than 128 characters long, and cannot be changed once LF-tag is created.

Values

Type a single value and select [Enter] or specify multiple values separated by commas.

Add

analyst ✕

engineer ✕

Enter up to 15 values; each value must be less than 256 characters long.

Cancel

Add LF-tag

Edit LF-Tags: standard\_datalake\_consumer

[Learn More](#)

✕

LF-Tags

After they are associated with catalog resources, LF-Tags allow you to create scalable permissions.

Assigned keys

🔍 Division ✕

Values

analyst ▼

Remove

Assign new LF-Tag

You can add 49 more LF-tags.

Cancel

Save

## Grant LF-tag permissions

Select the principals to grant permissions to, the LF-Tags to grant permissions on, and the specific set of permissions.

### Principals

☒ **IAM users and roles**  
Users or roles from this AWS account.

☐ **SAML users and groups**  
SAML users and group or QuickSight ARNs.

☐ **External accounts**  
AWS accounts or AWS organizations outside of this account.

**IAM users and roles**  
Add one or more IAM users or roles.

Choose IAM principals to add ▼

DataAnalyst X  
User

### LF-Tags

**LF-tag permission scope**  
Choose to grant permissions on all or a subset of LF-Tags.

Key

Division X

Values

Choose LF-tag values ▼

analyst X

Remove

Add LF-Tag

### Permissions

**LF-tag permissions**  
Select the specific access permissions to grant.

☒ Describe☐ Associate

**Grantable permissions**  
Select the permissions that the grant recipient can grant to other principals.

☐ Describe☐ Associate

Cancel

Grant

## Grant data permissions

### Principals

☒ **IAM users and roles**  
Users or roles from this AWS account.

☐ **SAML users and groups**  
SAML users and group or QuickSight ARNs.

☐ **External accounts**  
AWS accounts or AWS organizations outside of this account.

**IAM users and roles**  
Add one or more IAM users or roles.

Choose IAM principals to add

DataAnalyst X  
User

### LF-Tags or catalog resources

☒ **Resources matched by LF-Tags (recommended)**  
Manage permissions indirectly for resources or data matched by a specific set of LF-Tags.

☐ **Named data catalog resources**  
Manager permissions for specific databases or tables, in addition to fine-grained data access.

Key

Division X

Values

Choose LF-tag values

Remove

analyst X

Add LF-Tag

### Database permissions

**Database permissions**  
Choose specific access permissions to grant.

☐ Create table ☐ Alter ☐ Drop  
☒ Describe

☐ **Super**  
This permission is the union of all the individual permissions to the left, and supersedes them.

**Grantable permissions**  
Choose the permission that may be granted to others.

☐ Create table ☐ Alter ☐ Drop  
☐ Describe

☐ **Super**  
This permission allows the principal to grant any of the permissions to the left, and supersedes those grantable

Amazon Athena

Query editor

Editor

Recent queries

Saved queries

Settings

Workgroupprimary

Data

Data Source

AwsDataCatalog

Database

standard\_datalake\_consumer

Tables and views

Create

Filter tables and views

Tables (1)

pcs\_link

product\_namestring

categorystring

priceint

purchased\_atstring

customer\_namestring

Views (0)

Query 1

1SELECT \* FROM standard\_datalake\_consumer.pcs\_link

SQLLn 1, Col 50

Run againCancelSaveClearCreate

CompletedTime in queue: 0.12 secRun time: 2.122 secData scanned: 0.78 KB

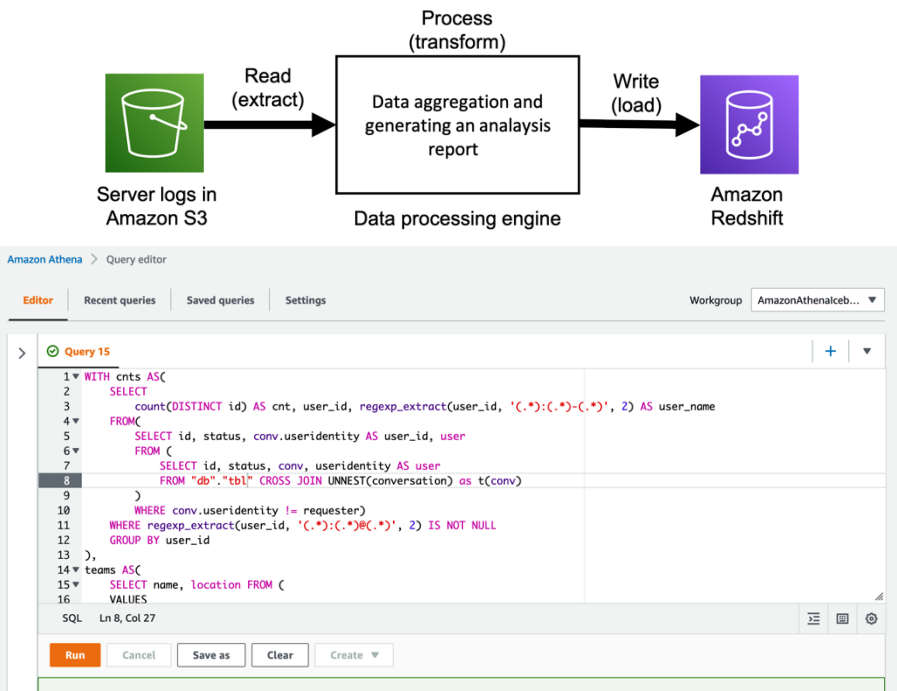
Results (4)CopyDownload results

Search rows

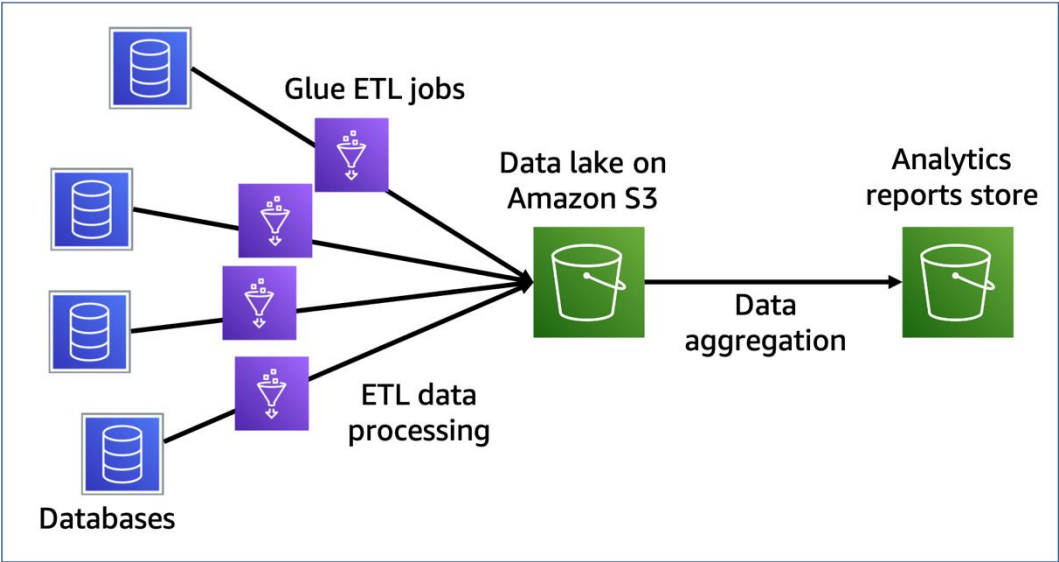
#	product_name	category	price	purchased_at	customer_name
1	Introduction to Cloud	Ebooks	15	2022-04-21T11:40:00Z	Barbara Gordon
2	Best practices on data lakes	Ebooks	25	2022-04-28T08:20:00Z	Tanya Fowler
3	Data Quest	Video games	30	2022-03-30T01:30:00Z	Rebecca Thompson
4	Final Shooting	Video games	20	2022-04-01T02:00:00Z	Rachel Gilbert

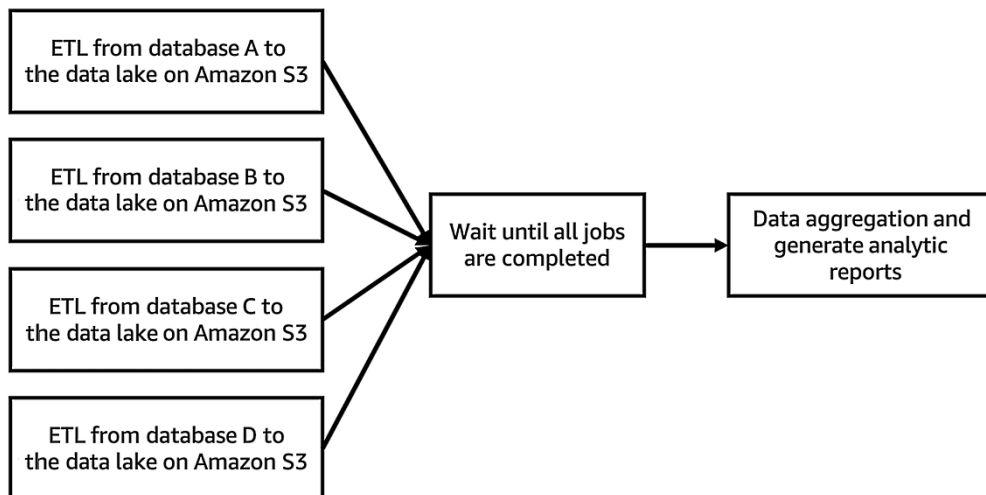


# Chapter 10: Data Pipeline Management



Allocated Resources per Worker (at the time of writing)	Standard	G.1X	G.2X
Memory	16 GB	16 GB	32 GB
vCPUs	4	4	8
Disk	50 GB	64 GB	128 GB





AWS Glue Studio > Jobs

### Jobs Info

**Create job Info**

☐ Visual with a source and target  
Start with a source, ApplyMapping transform, and target.

☐ Visual with a blank canvas  
Author using an interactive visual interface.

☒ Spark script editor  
Write or upload your own Spark code.

☐ Python Shell script editor  
Write or upload your own Python shell script.

☐ Jupyter Notebook - Preview  
Write your own code in a Jupyter Notebook for interactive development (Preview).

**Options Info**

☐ Create a new script with boilerplate code

☒ Upload and edit an existing script  
Choose a local file.

**File upload**

Limited to Python (\*.py, \*.py3) and Scala (\*.scala) files only.

ch10\_1\_example\_workflow\_gen\_report.py

2.09 KB

February 05, 2022

Run job

Actions

95 jobs

Script filename

ch10\_1\_example\_workflow\_gen\_report.py

Script path

S3 location of the script. Path must be in the form s3://bucket/prefix/path/. It must end with a slash (/) and not include any files.

s3://your-bucket/script/

×

View

Browse S3

☒ Job metrics Info

Enable the creation of CloudWatch metrics when this job runs.

☒ Continuous logging Info

Enable logs in CloudWatch.

☒ Spark UI Info

Enable using Spark UI for monitoring this job.

Spark UI logs path

s3://your-bucket/sparkeventlog/

×

View

Browse S3

Maximum concurrency

Sets the maximum number of concurrent runs that are allowed for this job. An error is returned when this threshold is reached.

1

Temporary path

Working directory. Path must be in the form s3://bucket/prefix/path/. It must end with a slash (/) and not include any files.

s3://your-bucket/temp/

×

View

Browse S3

Delay notification threshold (minutes) Info

## Workflows (5)

A workflow is an orchestration used to visualize and manage the relationship and execution of multiple triggers, jobs and crawlers.

Name	Last run	Last run status	Last modified
ch10_1_example_workflow_gen_report	-	-	Fri, 04 Feb 2022 18:43:14 GMT

Graph Details History

Legend: ● Start ◆ Trigger □ Job ■ Crawler ⚙ Incomplete ✖ Error ⌛ Deleting

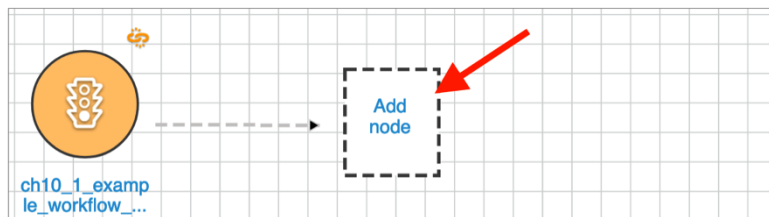
Remove Action

The workflow is empty

Add trigger

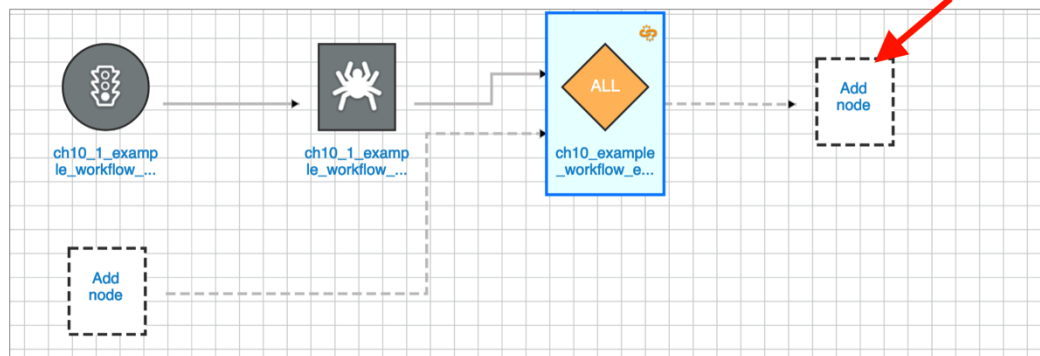
Graph Details History

Legend: ● Start ◆ Trigger □ Job ■ Crawler ⚙ Incomplete ✖ Error ⌛ Deleting



Graph Details History

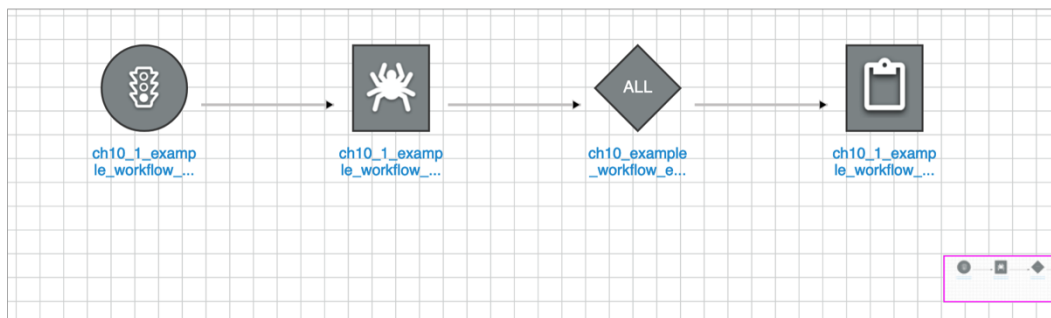
Legend: ● Start ◆ Trigger □ Job ■ Crawler ⚙ Incomplete ✖ Error ⌛ Deleting



Graph Details History

Legend: ● Start ◆ Trigger □ Job ■ Crawler ⚙ Incomplete ✖ Error ⌛ Deleting

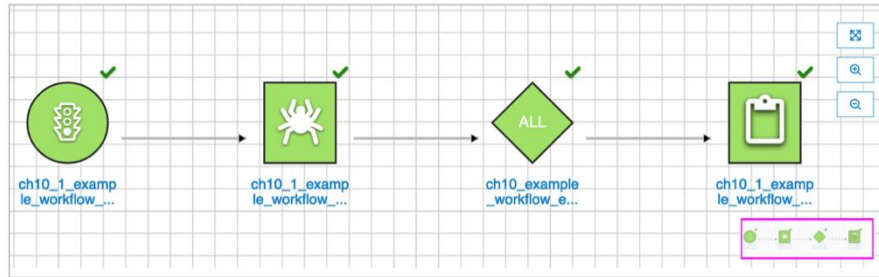
Remove



Graph

Select the graph nodes to resume and then choose Resume run.

**Legend:** ✔ Succeeded 🔄 Running ✖ Stopped ✖ Failed ✖ Timeout ✖ Error ⚠ Warning 🔄 Resume ■ Not started



category	total_sales	report_year
drink	906	2021
grocery	983	2021
kitchen	4535	2021
office	1316	2021

## Define state machine

### Design your workflow visually

Drag and drop your workflow together with Step Functions Workflow Studio. New

### Write your workflow in code

Author your workflow using Amazon States Language. You can generate code snippets to easily build out your workflow steps.

### Run a sample project

Deploy and run a fully functioning sample project in minutes using CloudFormation.

### Type

#### Standard

Durable, checkpointed workflows for machine learning, order fulfillment, IT/DevOps automation, ETL jobs, and other long-duration workflows.

#### Express

Event-driven workflows for streaming data processing, microservices orchestration, IoT data ingestion, mobile backends, and other short duration, high-event-rate workflows.

▶ Help me decide

### Definition

Define your workflow using [Amazon States Language](#). Test your data flow with the new [Data Flow Simulator](#).

Generate code snippet

Format JSON

```
1 {
2   "Comment": "A workflow to run Glue Crawler and ETL Job",
3   "StartAt": "StartCrawler",
4   "States": {
5     "StartCrawler": {
6       "Type": "Task",
7       "Parameters": {
8         "Name.$": "$.crawler_name"
9       },
10      "Resource": "arn:aws:states:::aws-sdk:glue:startCrawler",
11      "Next": "WaitForCrawlerRun",
12      "ResultPath": null
13    },
14    "WaitForCrawlerRun": {
15      "Type": "Wait",
16      "Seconds": 20,
17      "Next": "GetCrawler"
18    },
19    "GetCrawler": {
20      "Type": "Task",
21      "Next": "IsGlueCrawlerCompleted",
22      "Parameters": {
23        "Name.$": "$.crawler_name"
24      },
25      "Resource": "arn:aws:states:::aws-sdk:glue:getCrawler",
26      "ResultPath": "$.crawler",
27      "ResultSelector": {
28        "crawler_state.$": "$.Crawler.State"
29      }
30    }
31  }
32 }
```



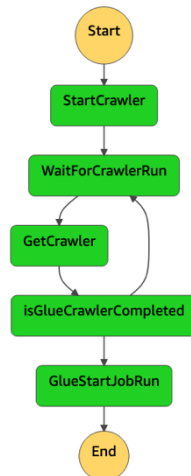
Name - *optional*

ffdb5e40-9dc9-0025-f619-63cc56083822

Input - *optional*

Enter input values for this execution in JSON format

```
1 {  
2   "crawler_name": "ch10_2_example_workflow",  
3   "etl_job_name": "ch10_2_example_workflow_gen_report",  
4   "etl_job_args": {  
5     "--datalake_location": "s3://<your-bucket-and-path>",  
6     "--database": "<your-database>",  
7     "--table": "example_workflow_sfn_sales",  
8     "--report_year": "2021"  
9   }  
10 }
```



In Progress Succeeded Failed Cancelled Caught Error

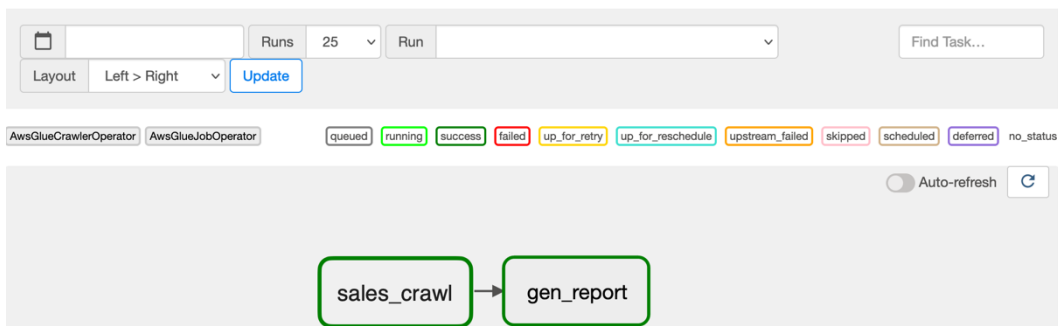
**DAG: ch10\_3\_example\_workflow\_mwaa**

**success** Schedule:

Next Run:

[Tree](#) [Graph](#) [Calendar](#) [Task Duration](#) [Task Tries](#) [Landing Times](#) [Gantt](#)

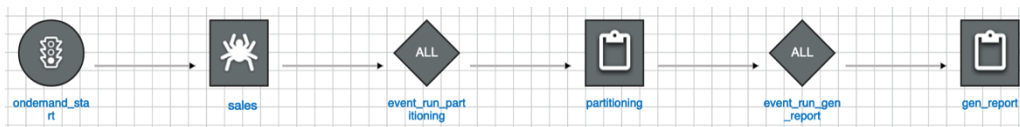
[Details](#) [Code](#)



```

Resources:
  GlueDatabase:
    Type: AWS::Glue::Database
    Properties:
      CatalogId: !Ref AWS::AccountId
      DatabaseInput:
        Name: 'glue_db'
  GlueTable:
    Type: 'AWS::Glue::Table'
    DependsOn:
      - GlueDatabase
    Properties:
      CatalogId: !Ref AWS::AccountId
      DatabaseName: 'glue_db'
      TableInput:
        Name: 'glue_table'

```



**Stack name**

Stack name

Stack name can include letters (A-Z and a-z), numbers (0-9), and dashes (-).

**Parameters**

Parameters are defined in your template and allow you to input custom values when you create or update a stack.

**DataLakeLocation**

The combination of S3 bucket name and path that stores the analytic sales data and a sales report. This location must end with a slash (/) and not include any files.

**DatabaseName**

Database name for the table of the sales data.

**GlueCrawlerRoleArn**

IAM Role ARN for the Glue Crawler.

**GlueJobRoleArn**

IAM Role ARN for the Glue ETL jobs.

**GlueJobScriptLocation**

The combination of S3 bucket name and path that locates ch10\_4\_example\_cf\_partitioning.py and ch10\_4\_example\_cf\_gen\_report.py. The combination of this location and each script name is specified in each Glue job as its script location. This location must end with a slash (/) and not include any files.

**ReportYear**

The year when you want to aggregate the dataset and generate a report.

**SalesDataLocation**

The combination of S3 bucket name and path that stores sales-data.json that you downloaded from GitHub repository. This location must end with a slash (/) and not include any files.

**TableName**

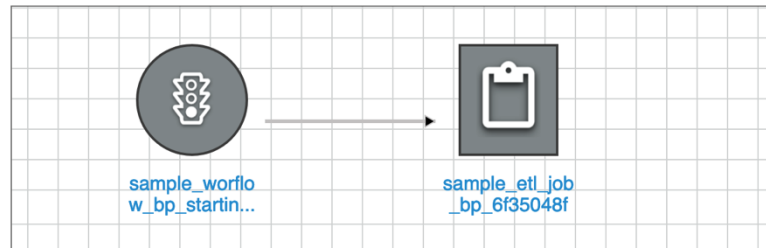
Table name for the table of sales data. If you use a custom table name, set the table name here.

Cancel
Previous
Next

	Name	Last run	Last run status
	sample_workflow_bp	-	-

Graph Details History

**Legend:** ● Start ◆ Trigger 📄 Job 🏠 Crawler ⚡ Incomplete ❌ Error ⌚ Deleting



## Create a workflow from ch10\_5\_example\_bp

AWS Glue will run the blueprint to create a workflow.

### WorkflowName

Name for the workflow.

### ScriptLocation

Specify the S3 path to store your glue job scripts.



### SalesDataLocation

Specify the S3 path to store the sales-data.json.



### DataLakeLocation

Specify the S3 path to store your sales data.



### GlueCrawlerRoleName

Choose an IAM role for Glue Crawler.

### GlueJobRoleName

Choose an IAM role for Glue ETL Job.

### DatabaseName

Specify a database name for the table of sales data.

### ReportYear

Specify the year when you want to aggregate the dataset and generate a report.

### IAM role

Role assumed by AWS Glue with permission to create workflows and their AWS resources. For more information, see [Create an IAM Role for AWS Glue](#).



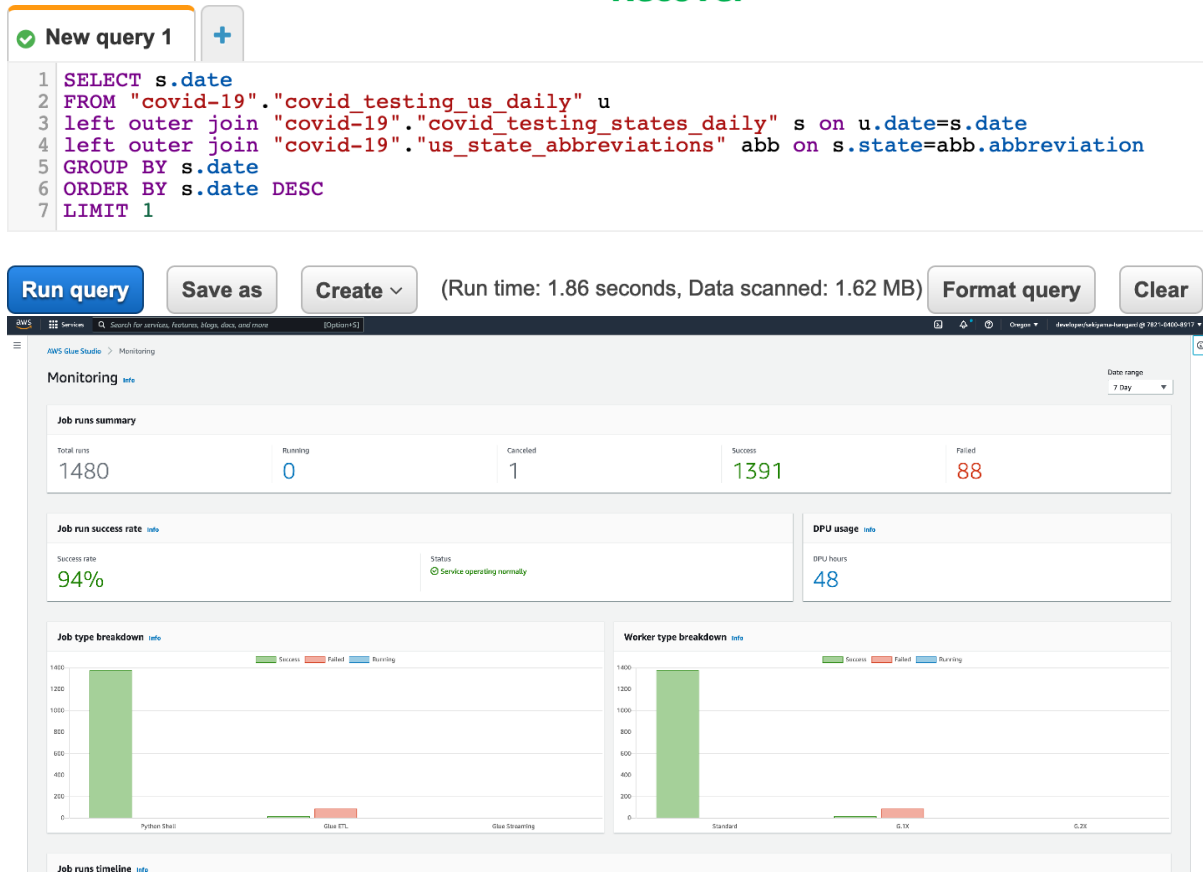
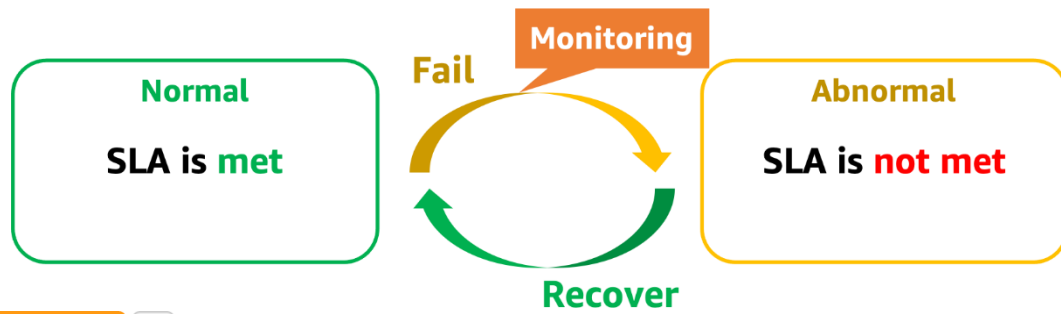
## Blueprint Runs (4)

Actions




	Run Id	State	Workflow name	Started on	Completed on
	bpr_b26bd73939210b7a3dfd1a1044...	SUCCEEDED	ch10_5_example_bp		

# Chapter 11: Monitoring



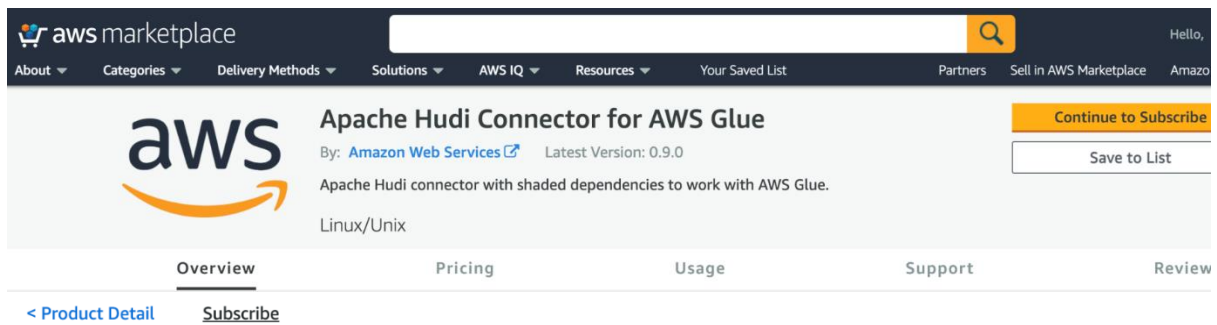




## Chapter 12: Tuning, Debugging, and Troubleshooting

No images...

## Chapter 13: Data Analysis



aws marketplace

About Categories Delivery Methods Solutions AWS IQ Resources Your Saved List Partners Sell in AWS Marketplace Amazon Web Services

aws

Apache Hudi Connector for AWS Glue

By: Amazon Web Services Latest Version: 0.9.0

Apache Hudi connector with shaded dependencies to work with AWS Glue.

Linux/Unix

Continue to Subscribe

Save to List

Overview Pricing Usage Support Review

< Product Detail Subscribe

### Subscribe to this software

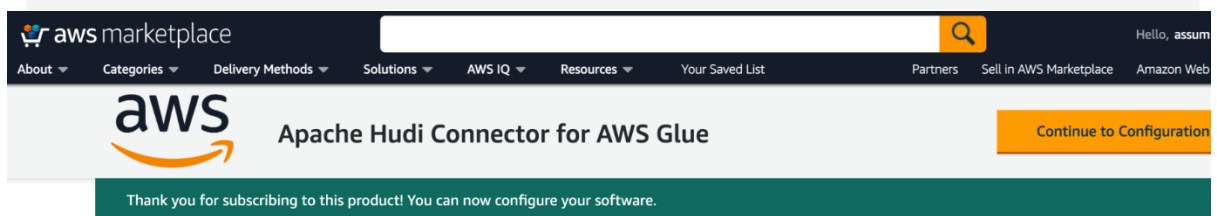
To create a subscription, review the pricing information and accept the terms for this software.

#### Terms and Conditions

##### Amazon Web Services Offer

By subscribing to this software, you agree to the pricing terms and the seller's [End User License Agreement \(EULA\)](#). You also agree and acknowledge that AWS may share information about this transaction (including your payment terms) with the respective seller, reseller or underlying provider, as applicable, in accordance with the [AWS Privacy Notice](#). AWS will issue invoices and collect payments from you on behalf of the seller through your AWS account. Your use of AWS services is subject to the [AWS Customer Agreement](#) or other agreement with AWS governing your use of such services.

Accept Terms



aws marketplace

About Categories Delivery Methods Solutions AWS IQ Resources Your Saved List Partners Sell in AWS Marketplace Amazon Web Services

aws

Apache Hudi Connector for AWS Glue

Continue to Configuration

Thank you for subscribing to this product! You can now configure your software.

#### Terms and Conditions

##### Amazon Web Services Offer

You have subscribed to this software and agreed that your use of this software is subject to the pricing terms and the seller's [End User License Agreement \(EULA\)](#). You agreed that AWS may share information about this transaction (including your payment terms) with the respective seller, reseller or underlying provider, as applicable, in accordance with the [AWS Privacy Notice](#). AWS will issue invoices and collect payments from you on behalf of the seller through your AWS account. Your use of AWS services remains subject to the [AWS Customer Agreement](#) or other agreement with AWS governing your use of such services.

Product	Effective date	Expiration date	Action
Apache Hudi Connector for AWS Glue	6/12/2022	N/A	<a href="#">Show Details</a>

aws marketplace

Hello, assume

About Categories Delivery Methods Solutions AWS IQ Resources Your Saved List Partners Sell in AWS Marketplace Amazon Web

aws

Apache Hudi Connector for AWS Glue

Continue to Launch

## Configure this software

Choose a fulfillment option and software version to launch this software.

Fulfillment option

Glue 3.0

Supported services [Learn more](#)

- Amazon ECS
- Amazon EKS

Software version

0.9.0 (Feb 16, 2022)

Fulfillment option description

Hudi connector with shaded dependencies to work with AWS Glue. - This delivery option is built with [hudi] (<https://github.com/apache/hudi>) 0.9.0. - This delivery option is compatible with AWS Glue 3.0.

aws

Apache Hudi Connector for AWS Glue

[< Product Detail](#) [Subscribe](#) [Configure](#) [Launch](#)

## Launch this software

Review the launch configuration details and follow the instructions to launch this so

Configuration details

Fulfillment option

Glue 3.0

Software version

0.9.0

Supported services

[Amazon ECS](#) [Amazon EKS](#)

Usage instructions

### Usage Instructions for 0.9.0 ✕

Please subscribe to the product from AWS Marketplace and [Activate the Glue connector from AWS Glue Studio](#).

## Create connection [Info](#)

### Connection properties [Info](#)

#### Connector



Apache Hudi Connector 0.9.0 for AWS Glue 3.0  
Connect to Apache Hudi tables from AWS Glue

#### Name

Enter a unique name for your connection.

Names can contain letters (A-Z), numbers (0-9), hyphens (-), or underscores (\_), and must be less than 256 characters long.

#### Description - optional

Descriptions can be up to 2048 characters long.

### Connection access [Info](#)

#### AWS Secret - optional [Info](#)

Choose a secret from [AWS Secrets Manager](#). [AWS secrets](#) eliminate hardcoding sensitive information.



#### ► Network options - optional

If your AWS Glue job needs to run on [Amazon Elastic Compute Cloud](#) (EC2) instances in a virtual private cloud (VPC) subnet, you must provide additional VPC-specific configuration information.

Cancel

Activate connector only

Create connection and activate connector



#### Refine results

##### Categories

Infrastructure Software (4)

Professional Services (1)

##### ▼ Delivery methods

☐ Container Image (3)

☐ Professional Services (1)

Elasticsearch Connector for AWS Glue (4 results) showing 1 - 4

< 1 > ⚙

Sort By: Relevance ▼



#### Elasticsearch Connector for AWS Glue

By Amazon Web Services | Ver 7.13.4-2

The Elasticsearch Connector for AWS Glue helps you read from and write to Elasticsearch using Apache Spark. By using this connector, you can focus on mining meaningful business insights from your data instead of writing and maintaining the connecting logic. For more details, please refer to Glue...

### AWS Glue Studio

**Jobs**

- Monitoring
- Connectors
- Sensitive data detection
- What's new

**Glue console**

- Glue catalog
- Crawlers
- Security configurations

Python Shell script editor  
Write or upload your own Python shell script.
Jupyter Notebook  
Write your own code in a Jupyter Notebook for interactive development.

**Source**  
 Amazon S3  
JSON, CSV, or Parquet files stored in S3.

→

**Target**  
 Amazon S3  
S3 bucket by specifying a bucket path as the data target.

**Your jobs (12)** Info

<input checked="" type="checkbox"/>	Job name	Type	Last modified	AWS Glue version
<input checked="" type="checkbox"/>	01 - Seed data job for Data Analysis Chapter	Glue ETL	6/12/2022, 2:12:00 AM	3.0

Python Shell
Glue ETL
Glue Streaming
Standard
G.1X
G.2X

**Job runs timeline** Info

■ Success 
 ■ Failed 
 ■ Running 
 ■ Canceled

**Job runs (1)** Info

Job name	Type	Start time	End time	Run status	Run time	Capacity	Worker type	DPU hours
<input type="radio"/> 01 - Seed data job for Data Analysis Chapter	Glue ETL	06/12/2022 02:29:46	-	<span style="color: blue;">●</span> Running	-	2	Standard	

### AWS Athena

**Query editor**

- Workgroups
- Data sources

**Jobs**

- Workflows New

Powered by Step Functions

☐ Enable compact mode

**Data**

Data source: AwsDataCatalog

Database: chapter-data-analysis-glue-database

Tables and views: Create

- Tables (1)
  - employees
- Views (0)

```
SELECT * FROM "chapter-data-analysis-glue-database"."employees" order by emp_no;
```

SQL Ln 1, Col 80

Run again Cancel Save Clear Create

Completed Time in queue: 121 ms Run time: 491 ms Data scanned: 0.71 KB

**Results (7)**

#	emp_no	name	department	city	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Tokyo	55000

## Database

dev

The database name must be 1-64 characters. Valid characters are lowercase alphanumeric characters.

## User name

## Authentication

- ☐ Federated user
- ☐ Temporary credentials
- ☒ Database user name and password

## Password

.....


☐ Show password

Provide a Password for the database that you are connecting to. The query editor v2 creates a secret on your behalf stored in AWS Secrets Manager. This secret contains credentials to connect to your database. [Learn more](#)


Cancel


Create connection


Cluster Serverless (admin) Database dev

**Redshift query editor v2**

Database

Queries

Notebooks  
(Preview)

Charts

+ Create

Load data

Filter resources

chapterdataanalysisredshift...

dev

chapter\_data\_analysi...

Tables1

employees

public

+ Create

Load data

Filter resources

chapterdataanalysisredshift...

dev

chapter\_data\_analysi...

Tables

employees

public

sample\_data\_dev

Untitled 1

Run

Limit 100

Explain

1 SELECT \* FROM "dev"."chapter\_data\_analysis\_schema"."employees" order by emp\_no;

Result 1 (7)

emp_no	name	department	city	salary
1	Adam	IT	SFO	50000
2	Susan	Sales	NY	60000
3	Jeff	Finance	Tokyo	55000
4	Bill	Manufacturing	New Delhi	70000
5	Joe	IT	Chicago	45000
6	Steve	Finance	NY	60000
7	Mike	IT	SFO	60000

Job runs (3) Info

Filter job runs

	Job name	Type	Start time	End time	Run status	Run time	Capacity	Worker type	DPU hours
<input type="radio"/>	02 - Hudi Init load for Data Analysis Chapter	Glue ETL	06/12/2022 18:16:04	-	Running	-	2	Standard	
<input type="radio"/>	01 - Seed data job for Data Analysis Chapter	Glue ETL	06/12/2022 17:13:34	06/12/2022 17:14:40	Succeeded	1 minute	2	Standard	0.03

Query 1

1 SELECT emp\_no, name, department, city, salary FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees\_cow" order by emp\_no;

SQL Ln 1, Col 97

Run again

Cancel

Save

Clear

Create

Completed

Time in queue: 115 ms

Run time: 1.079 sec

Data scanned:

Results (7)

Copy

Download results

Search rows

#	emp_no	name	department	city	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Tokyo	55000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

Job runs (4) Info

Filter job runs

	Job name	Type	Start time	End time	Run status	Run time	Capacity	Worker type	DPU hours
<input type="radio"/>	03 - Hudi Incremental load for Data Analysis Chapter	Glue ETL	06/12/2022 18:25:27	-	Running	-	2	Standard	
<input type="radio"/>	02 - Hudi Init load for Data Analysis Chapter	Glue ETL	06/12/2022 18:16:04	06/12/2022 18:18:16	Succeeded	2 minutes	2	Standard	0.06
<input type="radio"/>	01 - Seed data job for Data Analysis Chapter	Glue ETL	06/12/2022 17:13:34	06/12/2022 17:14:40	Succeeded	1 minute	2	Standard	0.03



Query 1

+

1

SELECT emp\_no, name, department, city, salary FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees\_cow" order by emp\_no;

SQL

Ln 1, Col 139

≡

📄

Run again

Cancel

Save ▼

Clear

Create ▼

Completed

Time in queue: 148 ms

Run time: 1.061 sec

Data scanned: 1.2

Results (7)

Copy

Download result

🔍 Search rows

< 1 >

# ▼	emp_no ▼	name ▼	department ▼	city ▼	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Cincinnati	75000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

Libraries [Info](#)

Python library path

/tmp/delta-core\_2.12-1.0.0.jar

Dependent JARs path

Referenced files path

Job parameters [Info](#)

Key

Q

--DELTALAKE\_CONN

X

Q

deltalake-connector

X

Remove

Q

--TARGET\_BUCKET

X

Q

scd-targets3bucket-i

X

Remove

Add new parameter

You can add 48 more parameters.

Tags

Key

Q

Project

X

Q

HandsonSeriesWithA

X

Remove

✔ Query 1

```
1 SELECT * FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees_deltalake" order by emp_no;
```

SQL Ln 1, Col 108

Run again

Cancel

Save ▼

Clear

Create ▼

✔ Completed

Time in queue: 110 ms

Run time: 1.265 sec

Data sca

Results (7)

Copy

Downl

Search rows

# ▼	emp_no ▼	name ▼	department ▼	city ▼	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Tokyo	55000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

✔ Query 1

```
1 SELECT * FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees_deltalake" order by emp_no;
```

SQL Ln 1, Col 108

Run again

Cancel

Save ▼

Clear

Create ▼

✔ Completed

Time in queue: 156 ms

Run time: 1.531 sec

Data sca

Results (7)

Copy

Downl

Search rows

# ▼	emp_no ▼	name ▼	department ▼	city ▼	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Cincinnati	70000
4	4	Bill	Manufacturing	New Delhi	70000

# Register location

## Amazon S3 location

Register an Amazon S3 path as the storage location for your data lake.

### Amazon S3 path

Choose an Amazon S3 path for your data lake.

s3://scd-targets3bucket-i7xhawru6kwi/employees\_governed\_table/

Browse

### Review location permissions - strongly recommended

Registering the selected location may result in your users gaining access to data already at that location. Before registering a location, we recommend that you review existing location permissions on resources in that location.

Review location permissions

### IAM role

To add or update data, Lake Formation needs read/write access to the chosen Amazon S3 path. Choose a role that you know has permission to do this, or choose the **AWSServiceRoleForLakeFormationDataAccess** service-linked role. When you register the first Amazon S3 path, the service-linked role and a new inline policy are created on your behalf. Lake Formation adds the first path to the inline policy and attaches it to the service-linked role. When you register subsequent paths, Lake Formation adds the path to the existing policy.

AWSServiceRoleForLakeFormationDataAccess

 Do not select the service linked role if you plan to use EMR.

Cancel

Register location

## Data locations (1)

Choose a storage location for which to review, grant or revoke user permissions.

e.g.: s3://bucket/prefix/

Browse

	Principal	Principal type	Resource	Own
<input type="radio"/>	HandsonSeriesWithAWSGlueJobRole	IAM role	s3://scd-targets3bucket-i7xhawru6kwi/employees_governed_table	-

### Query 1

```
1 SELECT * FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees_governed_table" order by emp_no;
```

SQL Ln 1, Col 95

Run again

Cancel

Save ▼

Clear

Create ▼

Completed

Time in queue: 113 ms

Run time: 5.062 sec

Da

### Results (7)

Copy

D

Search rows

# ▼	emp_no ▼	name ▼	department ▼	city ▼	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Tokyo	55000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

## Set up your connection's properties.

For more information, see [Working with Connections](#).

### Connection name

chapter-data-analysis-msk-connection

### Connection type

Kafka

☒ Amazon managed streaming for Apache Kafka (MSK)

☐ Customer managed Apache Kafka

### Select MSK cluster

msk-source-cluster-scd

### Kafka bootstrap server URLs

b-1.msksourceclusterscd.buosan.c10.kafka.us-west-2.amazonaws.com:9094,b-2.msksourceclusters

Enter a comma-separated list of bootstrap server URLs. Include the port number. Example: b-1.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094, b-2.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094, b-3.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094

☒ Require SSL connection

When selected, connection fails if unable to connect over SSL.

# Set up access to your data store.

For more information, see [Working with Connections](#).

## VPC

Choose the VPC name that contains your data store.

vpc-00a9f717ef84e5925

## Subnet

Choose the subnet within your VPC.

subnet-0361ac6939d4d2bd1

## Security groups

Choose one or more security groups that allow access to the data store in your VPC. AWS Glue associates these security groups to the ENI attached to your subnet. To allow AWS Glue components to communicate and also prevent access from other networks, at least one chosen security group must specify a self-referencing inbound rule for all TCP ports.

<input type="checkbox"/>	Group ID	Group name
<input checked="" type="checkbox"/>	sg-07c9e5d071c2a40c6	chapter-data-analysis-sg

Query 1

```
1 SELECT emp_no,name,department,city,salary FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees_cow_streaming" order by emp_no;
```

SQL Ln 1, Col 1

Run again

Cancel

Save

Clear

Create

Completed

Time in queue: 141 ms

Run time: 1.12 sec

Data scanned: 0.58 KB

Results (7)

Copy

Download results

Search rows

#	emp_no	name	department	city	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Tokyo	55000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

Query 1

```
1 SELECT emp_no,name,department,city,salary FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees_deltastreamer" order by emp_no;
```

SQL Ln 1, Col 66

Run again

Cancel

Save

Clear

Create

Completed

Time in queue: 163 ms

Run time: 926 ms

Data scanned: 0.56 KB

Results (7)

Copy

Download results

Search rows

< 1 > ⚙

#	emp_no	name	department	city	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Tokyo	55000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

Query 1

```
1 SELECT emp_no,name,department,city,salary FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees_deltastreamer" order by emp_no;
```

SQL Ln 1, Col 15

Run again

Cancel

Save

Clear

Create

Completed

Time in queue: 191 ms

Run time: 1.34 sec

Data scanned: 0.56 KB

Results (7)

Copy

Download results

Search rows

< 1 > ⚙

#	emp_no	name	department	city	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Cincinnati	70000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

Query 1

```
1 SELECT emp_no,name,department,city,salary FROM "AwsDataCatalog"."chapter-data-analysis-glue-database"."employees_cow_streaming" order by emp_no;
```

SQL Ln 1, Col 127

Run again

Cancel

Save

Clear

Create

Completed

Time in queue: 129 ms

Run time: 1.064 sec

Data scanned: 0.58 KB

Results (7)

Copy

Download results

Search rows

< 1 > ⚙

#	emp_no	name	department	city	salary
1	1	Adam	IT	SFO	50000
2	2	Susan	Sales	NY	60000
3	3	Jeff	Finance	Cincinnati	70000
4	4	Bill	Manufacturing	New Delhi	70000
5	5	Joe	IT	Chicago	45000
6	6	Steve	Finance	NY	60000
7	7	Mike	IT	SFO	60000

# opensearch-connection

Edit

Delete

Create job

## Connection details [Info](#)

Connector type

MARKETPLACE

Subnet

-

Connector ECR URL

https://709825985650.dkr.ecr.us-east-1.amazonaws.com/amazon-web-services/glue/elasticsearch:7.13.4-glue3.0-2

Require SSL connection

-

Security groups

-

Description

-



# Edit connection

## Connection properties [Info](#)

### Connector



Elasticsearch Connector 7.13.4 for AWS Glue 3.0  
Connect to Elasticsearch from AWS Glue



### Name (Read-only)

opensearch-connection

### Description - *optional*

Descriptions can be up to 2048 characters long.

## Connection access [Info](#)

### AWS Secret - *optional* [Info](#)

Choose a secret from [AWS Secrets Manager](#). [AWS secrets](#) eliminate hardcoding sensitive information.

ChapterDataAnalysisOSSecret



## Select your tenant

Tenants are useful for safely sharing your work with other OpenSearch Dashboards users. You can switch your tenant anytime by clicking the user avatar on top right.

☐ Global

The global tenant is shared between every OpenSearch Dashboards user.

☒ Private

The private tenant is exclusive to each user and can't be shared. You might use the private tenant for exploratory work.

☐ Choose from custom

Cancel

Confirm



vpc-chapter-data-analysis-6rt

# OpenSearch Dashboards



Query Workbench



Home

## Recently viewed



No recently viewed items



## OpenSearch Dashboards



Overview

Discover

Dashboard

Visualize

## OpenSearch Plugins



Query Workbench



## Query editor

```
1 select * from employees order by emp_no;
```

Run

Clear

Explain

## Results

Output employees

## employees (7)

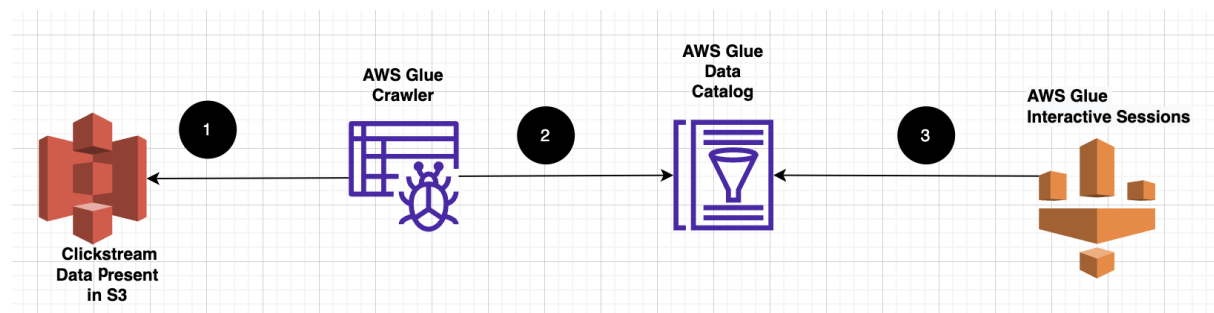
name	department	salary	city	emp_no
Adam	IT	50000	SFO	1
Susan	Sales	60000	NY	2
Jeff	Finance	55000	Tokyo	3
Bill	Manufacturing	70000	New Delhi	4

## Chapter 14: Machine Learning Integration

label	book_id	book_title	authors
	12367126	1-Safe Algorithms for Symmetric Site Configurations	John Hayes, Richard B. Bauchman
	08272651	1-Safe Algorithms for Symmetric Site Configurations	John Hayes, Richard B. Bauchman
	71616223	2003 SIGMOD Innovations Award Speech	Martha Smith
	12637181	2Q: A Low Overhead High-Performance Buffer Management	Elena Garcia
	72521341	2Q: A Low Overhead High-Performance Buffer Management	Elena Garcia

label	book_id	book_title	authors
0	12367126	1-Safe Algorithms for Symmetric Site Configurations	John Hayes, Richard B. Bauchman
0	08272651	1-Safe Algorithms for Symmetric Site Configurations	John Hayes, Richard B. Bauchman
1	71616223	2003 SIGMOD Innovations Award Speech	Martha Smith
2	12637181	2Q: A Low Overhead High-Performance Buffer Management	Elena Garcia
2	72521341	2Q: A Low Overhead High-Performance Buffer Management	Elena Garcia

# Chapter 15: Architecting Data Lakes for Real-World Scenarios and Edge Cases



```
%%time
%%sql
select count(*) from serverless_glue.tbl_without_index_clkstreamdata where customer = 2 and visityearmonth = 199812
```

CPU times: user 52.6 ms, sys: 19.1 ms, total: 71.7 ms  
Wall time: 29.7 s

Type: ☐ Table ☐ Pie

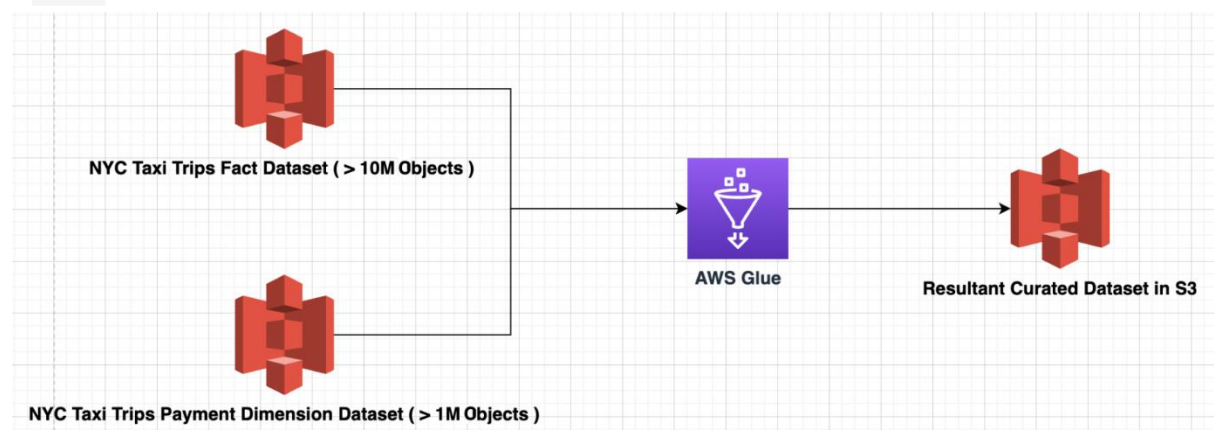
```
count(1)
7342094
```

```
%%time
%%sql
select count(*) from serverless_glue.tbl_with_index_clkstreamdata where customer = 2 and visityearmonth = 199812
```

CPU times: user 42.2 ms, sys: 2.21 ms, total: 44.4 ms  
Wall time: 9.43 s

Type: ☐ Table ☐ Pie

```
count(1)
7342094
```

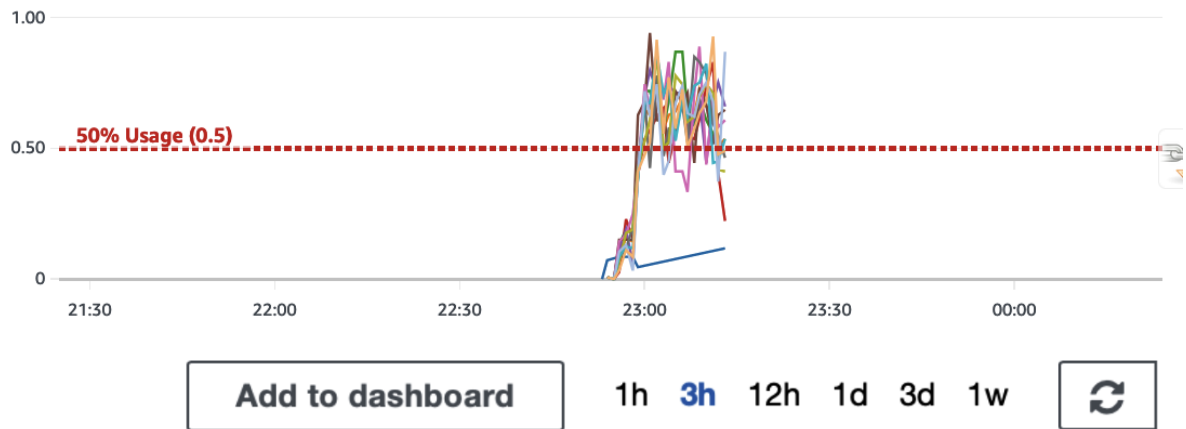


Memory Profile: Driver and Executors

No unit



## Memory Profile: Driver and Executors



## Memory Profile: Driver and Executors



## Edit table details



UPDATED\_BY\_CRAWLER

crawl\_aid



CrawlerSchemaSerializerVer:

1.0



recordCount

20726160



averageRecordSize

841



CrawlerSchemaDeserializerV

1.0



compressionType

none



classification

parquet



typeOfData

file



groupFiles

inPartition



groupSize

1048576



recurse

true



Apply



## S3 Eventual Consistency Model



If I try to access an object right after writing it, S3 might not show it until the write fully propagates

