

Serverless Data Pipeline

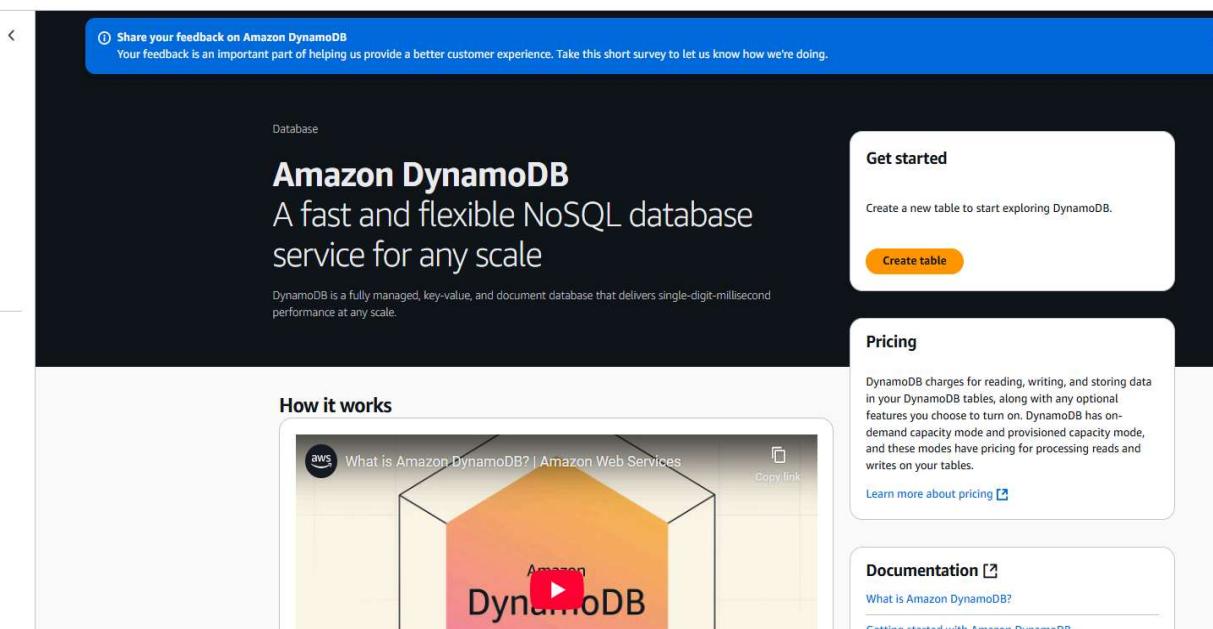
Overall Architecture

- Source: DynamoDB table (Orders)
- Transformation & Cataloging: Glue Crawlers & ETL Job
- Storage & Analytics: S3 (Parquet), Query via Athena
- Roles: One IAM role for DynamoDB crawling, another for S3 crawling and Glue ETL operational access

Step 1: Create DynamoDB Table

- Table Name: Orders
- Partition Key: OrderID (String)
- Sort Key: OrderDate (String)
- Capacity: On-Demand

Console Path: DynamoDB → Tables → Create Table



☰ DynamoDB > Tables > Create table

Create table

Table details Info

DynamoDB is a schemaless database that requires only a table name and a primary key when you create the table.

Table name
 This will be used to identify your table.
 Between 3 and 255 characters, containing only letters, numbers, underscores (_), hyphens (-), and periods (.)

Partition key
 The partition key is part of the table's primary key. It is a hash value that is used to retrieve items from your table and allocate data across hosts for scalability and availability.
 1 to 255 characters and case sensitive. String

Sort key - optional
 You can use a sort key as the second part of a table's primary key. The sort key allows you to sort or search among all items sharing the same partition key.
 1 to 255 characters and case sensitive. String

Table settings

Default settings
 The fastest way to create your table. You can modify most of these settings after your table has been created. To modify these settings now, choose 'Customize settings'.

Customize settings
 Use these advanced features to make DynamoDB work better for your needs.

Default table settings
 These are the default settings for your new table. You can change some of these settings after creating the table.

Setting	Value	Editable after creation
Setting	Value	Editable after creation

Tables (1) Info

Q Find tables Any tag key Any tag value

Name	Status	Partition key	Sort key	Indexes	Replication Regions	Deletion protection	Favorite	Read capacity mode	Write capacity mode	Total size	Table class
Orders	Active	OrderID (S)	OrderDate (S)	0	0	Off	☆	On-demand	On-demand	0 bytes	Standard

☰ DynamoDB > Tables > Orders

DynamoDB

- Dashboard
- Tables**
- Explore items
- PartiQL editor
- Backups
- Exports to S3
- Imports from S3
- Integrations
- Reserved capacity
- Settings

▼ DAX

- Clusters
- Subnet groups
- Parameter groups
- Events

Share your feedback on Amazon DynamoDB
 Your feedback is an important part of helping us provide a better customer experience. Take this short survey to let us know how we're doing.

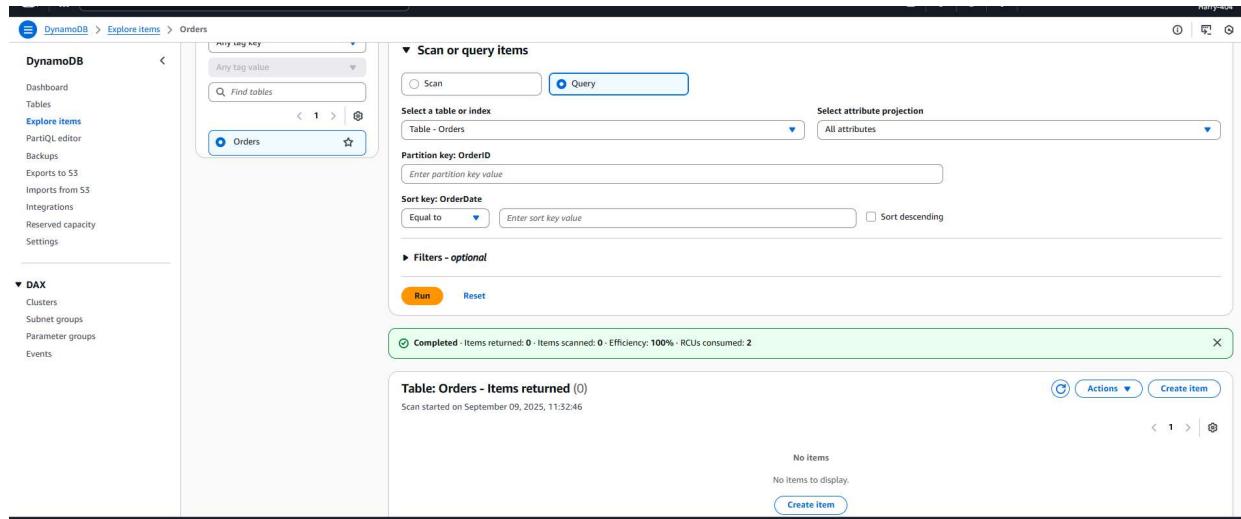
Orders

Settings Indexes Monitor Global tables Backups Exports and streams Permissions

Protect your DynamoDB table from accidental writes and deletes
 When you turn on point-in-time recovery (PITR), DynamoDB backs up your table data automatically so that you can restore to any given second in the preceding 1 to 35 days. Additional charges apply. [Learn more](#)

General information Info

Partition key OrderID (String)	Sort key OrderDate (String)	Capacity mode On-demand	Table status Active
Alarms 0 active alarms	Point-in-time recovery (PITR) <small>Info</small> Off	Item count 0	Table size 0 bytes
Average item size 0 bytes	Resource-based policy <small>Info</small> Not active	Get live item count	
Amazon Resource Name (ARN) arn:aws:dynamodb:ap-south-1:395938233352:table/Orders			
Additional info			
Read/write capacity <small>Info</small> The read/write capacity mode controls how you are charged for read and write throughput and how you manage capacity.			



Step 2: Insert Sample Data

Insert multiple orders manually via the DynamoDB Console → Orders → Create item.
Example (JSON view):

```
{  
  "OrderID": "O1001",  
  "OrderDate": "2025-09-05",  
  "Customer": "Alice",  
  "Amount": 250,  
  "Status": "Shipped"  
}
```

Repeat for other customers and statuses (like "Pending"), 10-15 records for better simulation.

DynamoDB > Explore Items: Orders > Create item

Share your feedback on Amazon DynamoDB
Your feedback is an important part of helping us provide a better customer experience. Take this short survey to let us know how we're doing.

Create item

You can add, remove, or edit the attributes of an item. You can nest attributes inside other attributes up to 32 levels deep. Learn more [\[?\]](#)

Attributes [View DynamoDB JSON](#)

```

1  {
2   "OrderID": {"S": "O1001"},
3   "OrderDate": {"S": "2025-09-05"},
4   "Customer": {"S": "Alice"},
5   "Amount": {"N": "250"},
6   "Status": {"S": "Shipped"}
7 }
8

```

[Copy](#)

DynamoDB > Explore Items > Orders

Scan or query items

Scan Query

Select a table or index Table: Orders

Select attribute projection All attributes

Partition key: OrderID Enter partition key value

Sort key: OrderDate Enter sort key value Sort descending

Filters - optional

Run Reset

Completed - Items returned: 0 · Items scanned: 0 · Efficiency: 100% · RCU consumed: 2

Table: Orders - Items returned (1)

Scan started on September 09, 2025, 11:52:46

OrderID (String)	OrderDate (String)	Amount	Customer	Status
O1001	2025-09-05	250	Alice	Shipped

Share your feedback on Amazon DynamoDB
Your feedback is an important part of helping us provide a better customer experience. Take this short survey to let us know how we're doing.

Create item

You can add, remove, or edit the attributes of an item. You can nest attributes inside other attributes up to 32 levels deep. Learn more [\[?\]](#)

Attributes [View DynamoDB JSON](#)

```

1 ↴ {
2   "OrderID": {"S": "O1002"},
3   "OrderDate": {"S": "2025-09-05"},
4   "Customer": {"S": "Bob"},
5   "Amount": {"N": "400"},
6   "Status": {"S": "Pending"}
7 }

```

[Copy](#)

Run Reset

Completed · Items returned: 0 · Items scanned: 0 · Efficiency: 100% · RCU consumed: 2

Table: Orders - Items returned (2)

Scan started on September 09, 2025, 11:32:46

<input type="checkbox"/> OrderID (String)	▼ OrderDate (String)	▼ Amount	▼ Customer	▼ Status
<input type="checkbox"/> O1002	2025-09-05	400	Bob	Pending
<input type="checkbox"/> O1001	2025-09-05	250	Alice	Shipped

Actions [\[?\]](#) Create item

Step 3: Glue Crawler 1 (OrdersCrawler): Catalog DynamoDB

- AWS Glue Console: Crawlers → Add crawler
- Name: OrdersCrawler
- Source: DynamoDB → Select "Orders" table

Crawler successfully starting
The following crawler is now starting: "OrdersCrawler"

OrdersCrawler

Crawler properties

Name: OrdersCrawler	IAM role: AWSGlueServiceRole-Orders	Database: ordersdb	State: READY
Description: -	Security configuration: -	Table prefix: -	

Crawler runs (3)

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
September 9, 2025 at 06:50:24	September 9, 2025 at 06:51:05	40 s	Failed	0.083	-
September 9, 2025 at 06:43:01	September 9, 2025 at 06:43:43	42 s	Failed	0.077	-
September 9, 2025 at 06:52:55	September 9, 2025 at 06:53:48	52 s	Completed	-	1 table change, 0 partition changes

IAM Role: Create (or select) a role with DynamoDB read access and Glue permissions

- Example: [AWSGlueServiceRole-Orders](#)
- Attach: [AmazonDynamoDBReadOnlyAccess](#), minimum [AWSGlueServiceRole](#)
- Output: Glue Data Catalog → Database: [ordersdb](#)

Run the crawler; it will infer the schema: *OrderID, OrderDate, Customer, Amount, Status*. IAM role trust and permissions are handled as per wizard prompts and AWS docs.

Welcome to AWS Glue

Get started by setting up your account and users, cataloging your data, and building ETL jobs to prepare data for analytics.

Prepare your account for AWS Glue

Admins: Grant access to AWS Glue and set a default IAM role.

[Set up roles and users](#)

Catalog and search for datasets

View your databases & tables and catalog data using Crawlers.

[Go to the Data Catalog](#)

Move and transform data

Transform data using a visual, notebook, or code interface.

[Author and edit ETL jobs](#)

Resources and tutorials

Getting started with AWS Glue: [Documentation](#) [AWS Training](#)

Glue In 5 Minutes Videos: [Authoring](#) [General Monitoring](#) [Orchestration](#)

[Using connectors and connections](#)

[AWS Glue Documentation home](#)

Examples: [AWS Glue blog posts](#) [AWS Glue on GitHub](#)

Data integration and management

Monitor & debug ETL jobs and track usage

[Go to job run monitoring](#)

Connect to your data stores

[Go to connections](#)

Orchestrate jobs to build data pipelines

[Go to workflows](#)

Manage and protect data

[Create a job to evaluate data quality](#) [Create a job to detect sensitive data](#)

What's new in Glue

- Amazon EMR on EC2 Adds Apache Spark native FGAC and AWS Glue Data Catalog Views Support Aug 29, 2025
- AWS Glue now supports Microsoft Dynamics 365 as a data source Jul 24, 2025
- AWS Glue Data Quality now supports Amazon 53 Tables and Iceberg Tables Jul 23, 2025

[View more what's new in Glue](#)

Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

Crawlers (0) info		Last updated (UTC)	Action	Run	Create crawler														
View and manage all available crawlers.		September 5, 2025 at 06:14:37																	
<input style="width: 100%;" type="text" value="Filter crawlers"/> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 10%;">Name</th> <th style="width: 10%;">State</th> <th style="width: 10%;">Schedule</th> <th style="width: 10%;">Last run</th> <th style="width: 10%;">Last run timestamp</th> <th style="width: 10%;">Log</th> <th style="width: 10%;">Table changes from last run</th> </tr> </thead> <tbody> <tr> <td colspan="7" style="text-align: center;">No resources</td> </tr> </tbody> </table>		Name	State	Schedule	Last run	Last run timestamp	Log	Table changes from last run	No resources										
Name	State	Schedule	Last run	Last run timestamp	Log	Table changes from last run													
No resources																			

Set crawler properties

Step 1 Set crawler properties

Crawler details

Name OrdersCrawler

Name can be up to 255 characters long. Some character set including control characters are prohibited.

Description - optional

Enter a description

Descriptions can be up to 2048 characters long.

Tags - optional

Use tags to organize and identify your resources.

[Cancel](#) [Next](#)

Choose data sources and classifiers

Data source configuration

Is your data already mapped to Glue tables?

Not yet Yes

Select one or more data sources to be crawled.

Data sources

The list of data sources to be crawled.

Type: **DynamoDB**

Table name

Orders

Scanning rate - optional

This field sets the percentage of DynamoDB table Read Capacity Units to be used by the crawler. If not specified, defaults to 0.5% for provisioned tables and 1/4 of maximum configured capacity for On-Demand tables.

Enter a value between 0.1 and 1.5.

Enable data sampling

Select to crawl a data sample only. If not selected, the entire table is crawled.

Add data source

Cancel **Add a DynamoDB data source**

AWS | AWS Glue | glue c | Add crawler | Account ID: 3959-1823-1352 | Asia Pacific (Mumbai) | Harry-404

Getting started

- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Data Catalog**
- Data Integration and ETL**
- Legacy pages

What's New | Documentation | AWS Marketplace

Enable compact mode | Enable new navigation

Choose data sources and classifiers

Data source configuration

Is your data already mapped to Glue tables?

Not yet Yes

Select one or more data sources to be crawled.

Data sources (1) info

The list of data sources to be crawled.

Type	Data source	Parameters
DynamoDB	Orders	-

Custom classifiers - optional

A classifier checks whether a given file is in a format the crawler can handle. If it is, the classifier creates a schema in the form of a StructType object that matches that data format.

Cancel **Previous** **Next**

AWS | AWS Glue | glue c | Add crawler | Account ID: 3959-1823-1352 | Asia Pacific (Mumbai) | Harry-404

Getting started

- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Data Catalog**
- Data Integration and ETL**
- Legacy pages

What's New | Documentation | AWS Marketplace

Enable compact mode | Enable new navigation

Configure security settings

IAM role info

Existing IAM role

Choose an IAM role Create new IAM role Update chosen IAM role

View

Only IAM roles created by the AWS Glue console and have the prefix "AWSGlueServiceRole" can be updated.

Security configuration - optional

Enable at-rest encryption with a security configuration.

Cancel **Previous** **Next**

and scheduling

create

▶ Security configuration - optional

Create new IAM role

Enter new IAM role

AWSGlueServiceRole-

Cancel Create

Step 2 Choose data sources and classifiers

Step 3 **Configure security settings**

Step 4 Set output and scheduling

Step 5 Review and create

IAM role info

Existing IAM role

Choose an IAM role

Create new IAM role

Update chosen IAM role

Only IAM roles created by the AWS Glue console and have the prefix "AWSGlueServiceRole-" can be updated.

▶ Security configuration - optional

Enable at-rest encryption with a security configuration.

Create new IAM role

Enter new IAM role

AWSGlueServiceRole-Orders

Cancel Create

Step 1 **Select trusted entity**

Step 2 Add permissions

Step 3 Name, review, and create

Select trusted entity

Trusted entity type

AWS service

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

AWS account

Allow entities in other AWS accounts belonging to you or a 3rd party to perform actions in this account.

SAML 2.0 federation

Allow users federated with SAML 2.0 from a corporate directory to perform actions in this account.

Custom trust policy

Create a custom trust policy to enable others to perform actions in this account.

Use case

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

Service or use case

DynamoDB

Choose a use case for the specified service.

Use case

Amazon DynamoDB Accelerator (DAX) - DynamoDB access

Allows DAX to call DynamoDB on your behalf.

DynamoDB Accelerator (DAX) - Cluster management

Allows DAX to manage resources within clusters on your behalf.

DynamoDB - Global Tables

Allows DynamoDB to manage cross-region replication for Global Tables on your behalf.

Cancel Next

AWS Glue > Crawlers > Add crawler

Getting started

ETL jobs

Visual ETL

Notebooks

Job run monitoring

Data Catalog tables

Data connections

Workflows (orchestration)

Data Catalog

Data integration and ETL

Legacy pages

What's New

Documentation

AWS Marketplace

Enable compact mode

Error creating IAM Role "AWSGlueServiceRole-Orders"

We've encountered an error while attempting to create the IAM Role "AWSGlueServiceRole-Orders"

Message from service: "There were problems with the supplied replacements: [Missing single-value replacements required by the template for [tableName]]"

Configure security settings

IAM role info

Existing IAM role

Choose an IAM role

Create new IAM role

Update chosen IAM role

Only IAM roles created by the AWS Glue console and have the prefix "AWSGlueServiceRole-" can be updated.

Security configuration - optional

Enable at-rest encryption with a security configuration.

Cancel Previous Next

Step 1: Select trusted entity

Select trusted entity

Trusted entity type

- AWS service: Allow AWS services like EC2, Lambda, or others to perform actions in this account.
- AWS account: Allow entities in other AWS accounts belonging to you or a 3rd party to perform actions in this account.
- Web identity: Allows users federated by the specified external web identity provider to assume this role to perform actions in this account.
- SAML 2.0 federation: Allows users federated with SAML 2.0 from a corporate directory to perform actions in this account.
- Custom trust policy: Create a custom trust policy to enable others to perform actions in this account.

Use case

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

Service or use case

Glue

Choose a use case for the specified service.

Use case

Glue: Allows Glue to call AWS services on your behalf.

Step 2: Add permissions

Role details

Role name

AWSGlueServiceRole-Order1

Description

Allows Glue to call AWS services on your behalf.

Step 1: Select trusted entities

Trust policy

```

1 <[{"Version": "2012-10-17", "Statement": [{"Effect": "Allow", "Principal": {"Service": "glue.amazonaws.com"}, "Action": "sts:AssumeRole"}]}]>

```

Step 2: Add permissions

Permissions policy summary

Policy name	Type	Attached as
AmazonDynamoDBReadOnlyAccess	AWS managed	Permissions policy

Step 2: Add permissions

Permissions policy summary

Policy name	Type	Attached as
AmazonDynamoDBReadOnlyAccess	AWS managed	Permissions policy
AmazonS3FullAccess	AWS managed	Permissions policy

Step 3: Add tags

Add tags - optional

Tags are key-value pairs that you can add to AWS resources to help identify, organize, or search for resources.

No tags associated with the resource.

Add new tag

You can add up to 50 more tags.

Cancel **Previous** **Create role**

The screenshot shows two sequential steps in the AWS Glue console for creating a crawler:

Configure security settings (Step 3):

- IAM role**: AWSGlueServiceRole-Orders (Existing IAM role)
- Create new IAM role** and **Update chosen IAM role** buttons
- Security configuration - optional** section

Set output and scheduling (Step 4):

- Output configuration**:
 - Target database**: Choose a database (dropdown menu)
 - Table name prefix - optional**: Type a prefix added to table names
- Advanced options**
- Crawler schedule**:
 - You can define a time-based schedule for your crawlers and jobs in AWS Glue. The definition of these schedules uses the Unix-like cron syntax. Learn more.
 - Frequency**: On demand

Create Glue Database:

- In AWS Glue Console, go to Databases → Add database
- Database name: ordersdb (must be all lowercase)
- Click Create database

AWS Glue > Databases > Add database

Create a database
Create a database in the AWS Glue Data Catalog.

Database details

Name
ordersdb
Database name is required, in lowercase characters, and no longer than 256 characters.

Description - optional
Enter text
Descriptions can be up to 2048 characters long.

Database settings

Location - optional
Set the URI location for use by clients of the Data Catalog.

Create database

AWS Glue > Crawlers > Add crawler

Set output and scheduling

Output configuration
Target database: ordersdb
Clear selection Add database

Table name prefix - optional
Type a prefix added to table names

Crawler schedule
You can define a time-based schedule for your crawlers and jobs in AWS Glue. The definition of these schedules uses the Unix-like cron syntax. Learn more

Frequency
On demand

Cancel **Previous** **Next**

AWS Glue > Crawlers > Add crawler

Review and create

Step 1: Set crawler properties

Set crawler properties

Name	Description	Tags
OrdersCrawler	-	-

Step 2: Choose data sources and classifiers

Data sources (1) info
The list of data sources to be scanned by the crawler.

Type	Data source	Parameters
DynamoDB	Orders	-

Step 3: Configure security settings

Configure security settings

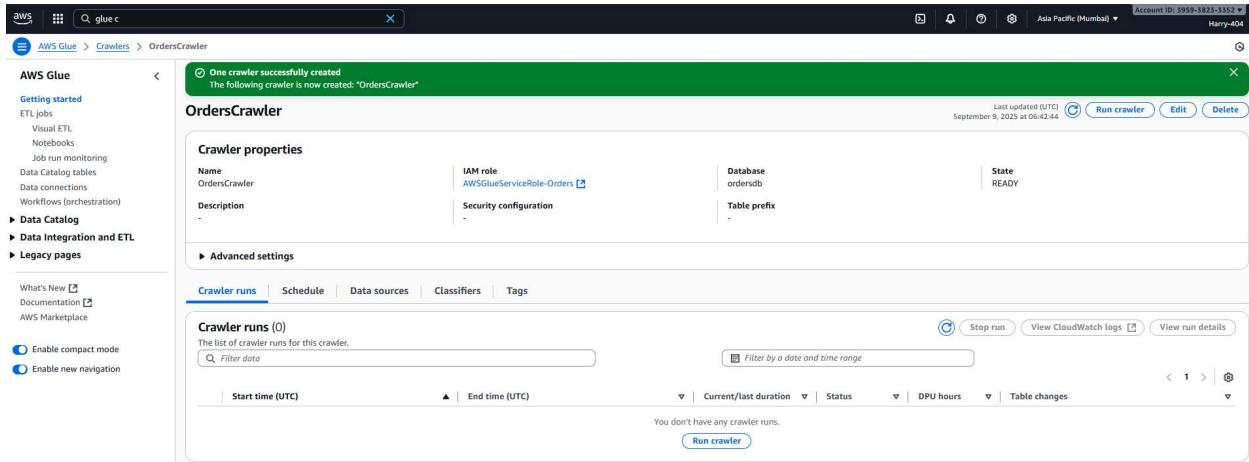
IAM role	Security configuration	Lake Formation configuration
AWSGlueServiceRole-Orders	-	-

Step 4: Set output and scheduling

Set output and scheduling

Database	Table prefix - optional	Schedule
ordersdb	-	On demand

Create crawler



Step 4: Glue Crawler 2 (OrdersS3Crawler): Catalog Cleaned S3 Data

- Purpose: After Glue ETL job writes transformed data to S3, this crawler catalogs the S3 Parquet files.
- Name: OrdersS3Crawler
- Source: S3 → Bucket or prefix (e.g., `s3://my-orders-analytics/shipped/`)
- IAM Role: Create/select role (e.g., `s3buck`) with S3 read, Glue, and DynamoDB read-only permissions
 - Attach: `AmazonS3FullAccess`, `AmazonDynamoDBReadOnlyAccess`, `AWSGlueServiceRole`
 - Inline Permission: `DynamodbRead`, `S3pol`
- Output: Glue Data Catalog → Database: `ordersdb` (or a separate database if preferred)

Run the crawler after ETL job has created Parquet files in S3 to update the schema in Glue Catalog as per S3 data structure.

Screenshot of the AWS IAM Role configuration for 'S3buck'.

Summary

Creation date: September 10, 2025, 14:34 (UTC+05:30)

Last activity: -

ARN: arn:aws:iam::395938233352:role/S3buck

Maximum session duration: 1 hour

Permissions (5):

- AmazonDynamoDBReadOnlyAccess (AWS managed)
- AmazonS3FullAccess (AWS managed)
- AWSGlueServiceRole (AWS managed)
- DynamodbRead (Customer inline)
- S3pol (Customer inline)

Trust relationships

```

1 < [{}]
2   "Version": "2012-10-17",
3   "Statement": [
4     {
5       "Effect": "Allow",
6       "Principal": {
7         "Service": [
8           "s3.amazonaws.com",
9           "glue.amazonaws.com"
10        ]
11      },
12      "Action": "sts:AssumeRole"
13    }
14  ]
15 ]
  
```

The screenshot shows the AWS Glue Crawler configuration page. The crawler is named 'OrdersS3Crawler' and is currently in a 'READY' state. It uses the 'ordersdb' database and the 'S3buck' IAM role. The crawler has run once, completed on September 10, 2025, at 07:33:57, with 52 seconds of DPU hours and 1 table change, 0 partition changes.

Step 5: Glue ETL Job (Visual Editor/Studio, PySpark Script)

- Source: Glue Data Catalog Table: ordersdb.orders (from DynamoDB, cataloged by OrdersCrawler)
- Target: S3 → `s3://my-orders-analytics/shipped/`
 - Format: Parquet (compression: Snappy, ideal for analytics)
- Transformation: Filter for orders where Status = "Shipped"
- Role: Use the same IAM role as for S3 Crawler, with access to S3, Glue, and DynamoDB.
- Method: Either visual mapping or PySpark code as below:

Sample Glue PySpark Script

python

```
import sys

from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
```

```
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

# Read DynamoDB (via Glue Catalog)
datasource = glueContext.create_dynamic_frame.from_catalog(
    database="ordersdb",
    table_name="orders"
)

# Filter only shipped orders
filtered = Filter.apply(datasource, lambda x: x["status"] == "Shipped")

# Write to S3 as Parquet
glueContext.write_dynamic_frame.from_options(
    frame=filtered,
    connection_type="s3",
    connection_options={"path": "s3://my-orders-analytics/shipped/"},
    format="parquet"
)
```

```
job.commit()
```

- In Glue Studio Visual Mode: Use drag-and-drop to set up the same source, filter (status = shipped), S3 as Parquet target. Map columns as shown in your screenshots.

The image consists of three vertically stacked screenshots of the AWS Glue Studio interface, illustrating the process of creating a data pipeline.

Screenshot 1: AWS Glue Studio Home

This screenshot shows the AWS Glue Studio home page. On the left, a sidebar navigation includes sections for AWS Glue (Getting started, Visual ETL, Notebooks, Job run monitoring, Data Catalog tables, Data connections, Workflows (orchestration)), Data Catalog (Data Integration and ETL, Legacy pages), and a general AWS section (What's New, Documentation, AWS Marketplace). On the right, the main area is titled "AWS Glue Studio" and "Create job". It features three options: "Visual ETL" (selected), "Notebook", and "Script editor". Below this is a "Example jobs" section and a "Your jobs (0)" section with a "Create job from a blank graph" button.

Screenshot 2: Visual Mode - Data Sources

This screenshot shows the "Visual" mode of AWS Glue Studio. The top navigation bar includes "Actions", "Save", and "Run". The main area is titled "Untitled job" and shows a "Sources" tab. A search bar at the top of the list says "search sources, transforms and targets". The list includes "AWS Glue Data Catalog", "Amazon S3", "Amazon Kinesis", "Apache Kafka", "Relational DB", "Amazon Redshift", "MySQL", "PostgreSQL", and "Oracle SQL". A "Manage Connections" button is at the bottom.

Screenshot 3: Visual Mode - Data Preview

This screenshot shows the "Visual" mode continuing. The top navigation bar includes "Actions", "Save", and "Run". The main area is titled "Untitled job" and shows a "Data preview" tab. It displays a data source named "Data source - Data Catalog" with the sub-name "AWS Glue Data Catalog". The "Output schema" tab is also visible. On the right, "Data source properties - Data Catalog" are being configured, including "Name" (set to "AWS Glue Data Catalog"), "Database" (set to "ordersdb"), and "Table" (set to "orders"). A "Partition predicate - optional" field is also present. The bottom of the screen shows "Additional Settings" and a "Start session" button.

job 0

Visual | Script | **Job details** | Runs | Data quality | Schedules | Version Control

Basic properties [Info](#)

Name
job

Description - optional
Descriptions can be up to 2048 characters long.

IAM Role
Role assumed by the job with permission to access your data stores. Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job.

AWSGlueServiceRole-Orders [Actions](#) [Save](#) [Run](#)

Type
The type of ETL job. This is set automatically based on the types of data sources you have selected.

Spark

Glue version [Info](#)
Glue 5.0 - Supports spark 3.5, Scala 2, Python 3

Language
Python 3

Worker type
Set the type of predefined worker that is allowed when a job runs.

G 1v

job 0

Job has not been saved [Actions](#) [Save](#) [Run](#)

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Data source properties - Data Catalog

Name AWS Glue Data Catalog

Database
Choose a database.
ordersdb

Use runtime parameters

Table
orders

Use runtime parameters

Partition predicate - optional
Enter a boolean expression supported by Spark SQL, using only partition columns.
partition_predicate

Data preview [Info](#) **READY** [Actions](#) [End session](#) [Previewing 5 of 5 fields](#)

Output schema

Data preview (2) [Info](#) **READY** [Actions](#)

Filter sample dataset

orderid	amount	customer	orderdate	status
O1001	250	Alice	2025-09-05	Shipped
O1002	400	Bob	2025-09-05	Pending

aws Search [Alt+S] Account ID: 5959-3823-3352 Asia Pacific (Mumbai) Harry-404

job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Data source - Data Catalog AWS Glue Data Catalog

Data preview | Output schema

Schema AVAILABLE

Key	Data type
orderid	string
amount	long
customer	string
orderdate	string
status	string

Infer schema from session

Data source properties - Data Catalog

Name: AWS Glue Data Catalog

Database: Choose a database. ordersdb

▶ Use runtime parameters

Table: orders

▶ Use runtime parameters

Partition predicate - optional

Partition predicate syntax for Spark SQL is year >= year(date_sub(current_date, 7)) AND month >= month(date_sub(current_date, 7)) AND day >= day(date_sub(current_date, 7)).

aws Search [Alt+S] Account ID: 5959-3823-3352 Asia Pacific (Mumbai) Harry-404

job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Data source - Data Catalog AWS Glue Data Catalog

Data target - S3 bucket Amazon S3

Data preview | Output schema

Schema AVAILABLE

Key	Data type	Partition
orderid	string	-
amount	long	-
customer	string	-
orderdate	string	-
status	string	-

Infer schema from session

Data target properties - S3

Name: Amazon S3

Node parents: Choose which nodes will provide inputs for this one. Choose one or more parent node AWS Glue Data Catalog Catalog - DataSource

Format: 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://bucket/prefix/object/ View | Browse S3

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 add new partitions

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

File partitioning: Info

After you save your job, it will be saved into autogenerated or in a predesignated number of partitions.

Autogenerate files (Recommended)

Multiple file output

Partition keys - optional

aws Search [Alt+S] Account ID: 5959-3823-3352 Asia Pacific (Mumbai) Harry-404

job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Data source - Data Catalog AWS Glue Data Catalog

Data target - S3 bucket Amazon S3

Data preview | Output schema

Schema AVAILABLE

Key	Data type	Partition
orderid	string	-
amount	long	-
customer	string	-
orderdate	string	-
status	string	-

Infer schema from session

Data target properties - S3

Name: Amazon S3

Node parents: Choose which nodes will provide inputs for this one. Choose one or more parent node AWS Glue Data Catalog Catalog - DataSource

Format: 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://my-orders-analytics/shipped/ View | Browse S3

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 add new partitions

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

File partitioning: Info

After you save your job, it will be saved into autogenerated or in a predesignated number of partitions.

Autogenerate files (Recommended)

Multiple file output

job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Add nodes

Search sources, transforms and targets

Sources | Transforms | Targets | Popular

Change Schema
Change field names, data types and drop fields. Formerly known as Apply Mapping.

Join
Combine records from two datasets based on a set of conditions.

SQL Query
Use a SQL query to transform data.

Join

Detect Sensitive Data
Detect PII and other sensitive information.

Evaluate Data Quality
Evaluate the quality and completeness of your data.

Aggregate
Apply functions like sum or average to fields in the dataset.

Custom Transform
Write custom code to transform data.

Drop Duplicates
Remove duplicate records from your dataset.

Drop Fields
Remove the selected fields from your data.

Drop Null Fields
Remove columns that have only empty/null values.

Filter
Filter data based on conditions.

Add Transforms

Data target properties

Name
Amazon S3

Node parents
Choose which nodes will provide inputs for this one.
Choose one or more parent nodes

AWS Glue Data Catalog
Catalog - DataSource

Format
Parquet

Compression Type
Snappy

S3 Target Location
Choose an S3 location for the output
s3://my-orders-analy

Data Catalog update option
Choose how you want to update an S3 backed source.

Do not update the Data

Create a table in the Data Catalog

Create a table in the Data Store

File partitioning | Info
Choose if you want the output to be partitioned.

Autogenerate files (Rec)

Multiple file output

Last modified on 9/9/2025, 12:39:33 PM Actions Save Run

job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Transform

Name: Filter

Node parents: Choose one or more parent node
AWS Glue Data Catalog

Filter | Info
Build a new output by selecting records from the input data that satisfy a specified predicate function.

Global AND
All filter conditions will be applied as a global "AND."

Global OR
All filter conditions will be applied as a global "OR."

Filter condition | Info
Specify your filter condition by choosing the key, operator, and entering a value.

Key: status Operation: matches Value: shipped

Show full path

Add condition

Last modified on 9/9/2025, 12:39:33 PM Actions Save Run

Schema **Info** **AVAILABLE**

Key	Data type	Partition
orderid	string	-
amount	long	-
customer	string	-
orderdate	string	-
status	string	-

Data preview | Output schema

job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Transform

Name: Filter

Node parents: Choose one or more parent node
AWS Glue Data Catalog

Filter | Info
Build a new output by selecting records from the input data that satisfy a specified predicate function.

Global AND
All filter conditions will be applied as a global "AND."

Global OR
All filter conditions will be applied as a global "OR."

Filter condition | Info
Specify your filter condition by choosing the key, operator, and entering a value.

Key: status Operation: matches Value: shipped

Show full path

Add condition

Last modified on 9/9/2025, 12:39:33 PM Actions Save Run

Schema **Info** **AVAILABLE**

Key	Data type
-----	-----------

Data preview | Output schema

Amazon S3 > Buckets > [aws-glue-assets-395938233352-ap-south-1](#) > Edit bucket policy

Amazon S3

General purpose buckets

- Directory buckets
- Table buckets
- Vector buckets
- Access Grants
- Access Points (General Purpose Buckets, FSx file systems)
- Access Points (Directory Buckets)
- Object Lambda Access Points
- Multi-Region Access Points
- Batch Operations
- IAM Access Analyzer for S3

Block Public Access settings for this account.

Storage Lens

- Dashboards
- Storage Lens groups
- AWS Organizations settings

Feature spotlight 11

AWS Marketplace for S3

Edit bucket policy [Info](#)

Bucket policy

The bucket policy, written in JSON, provides access to the objects stored in the bucket. Bucket policies don't apply to objects owned by other accounts. [Learn more](#)

Bucket ARN

[arn:aws:s3:::aws-glue-assets-395938233352-ap-south-1](#)

Policy

```

1▼ {
2  "Version": "2012-10-17",
3  "Statement": [
4    {
5      "Effect": "Allow",
6      "Principal": {
7        "AWS": "arn:aws:iam::395938233352:role/AWSGlueServiceRole-Orders"
8      },
9      "Action": [
10        "s3:PutObject",
11        "s3:GetObject",
12        "s3>ListBucket"
13      ],
14      "Resource": [
15        "arn:aws:s3:::my-orders-analytics",
16        "arn:aws:s3:::my-orders-analytics/*"
17      ]
18    }
19  ]
20 }
21

```

job

Visual | Script | Job details | Runs | Data quality | Schedules | Version Control

Visual

Data target properties - S3

Name: my-orders-analytics

Node parents: Choose which nodes will provide inputs for this one. [Choose one or more parent node](#)

AWS Glue Data Catalog Catalog - DataSource

Format: Parquet

Compression Type: Snappy

S3 Target Location: Choose an S3 location in the format s3://bucket/prefix/object/ with a trailing slash (/). [s3://aws-glue-assets-395938233352-ap-south-1/sh](#) [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

File partitioning: [Info](#)

Choose if you want the output to be saved into autogenerated or in a predesignated number of partitions.

- Autogenerate files (Recommended)
- Multiple file output

Partition keys - optional

Data preview **Output schema**

Target node not supported

You have selected a data target node which is not supported for data preview. Please select another type of node instead.

Amazon S3 > Buckets > aws-glue-assets-395938233352-ap-south-1

aws-glue-assets-395938233352-ap-south-1 Info

Objects (3) Actions Create folder Upload

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](#)

Find objects by prefix:

Name	Type	Last modified	Size	Storage class
scripts/	Folder	-	-	-
shipped/	Folder	-	-	-
sparkHistoryLogs/	Folder	-	-	-

Block Public Access settings for this account

Storage Lens

Feature spotlight: 11

AWS Marketplace for S3

aws | athena | athena

Amazon S3 > Buckets > aws-glue-assets-395938233352-ap-south-1

Amazon S3

General purpose buckets

Block public access (bucket settings)

Block all public access Off

Individual Block Public Access settings for this bucket

Bucket policy

The bucket policy, written in JSON, provides access to the objects stored in the bucket. Bucket policies don't apply to objects owned by other accounts. [Learn more](#)

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "AWS": "arn:aws:iam::395938233352:role/AWSGlueServiceRole-Orders"
      },
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3>ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::aws-glue-assets-395938233352-ap-south-1",
        "arn:aws:s3:::aws-glue-assets-395938233352-ap-south-1/*"
      ]
    }
  ]
}
```

Amazon Q Free Tier

Troubleshooting AWS Glue Job Issues

```

  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "AWS": "arn:aws:iam::395938233352:role/AWSGlueServiceRole-Orders"
      },
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3>ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::your-target-bucket",
        "arn:aws:s3:::your-target-bucket/*"
      ]
    }
  ]
}

2. Verify the S3 Path:
  

- Check line 32 in your Glue script
- Ensure the S3 path is correct and the bucket exists
- Verify you're writing to the correct region (ap-south-1)


3. Test with a Different S3 Location:
  

- Try writing to a test bucket you control to isolate the issue

```

This will help me provide more specific guidance for your situation.

Ask me anything about AWS

Amazon S3 > Buckets > aws-glue-assets-395938233352-ap-south-1 > sparkHistoryLogs/

Objects (20)

Name	Type	Last modified	Size	Storage class
jr_3b53cc3f41e8752f53197c08b5a6fdb2950aa5427854a92385dc59454913a126	-	September 9, 2025, 14:06:27 (UTC+05:30)	484.4 KB	Standard
jr_3b53cc3f41e8752f53197c08b5a6fdb2950aa5427854a92385dc59454913a126.inprogress	inprogress	September 9, 2025, 14:06:25 (UTC+05:30)	465.3 KB	Standard
jr_5e9b0ed9980d082bd22fb33e3653192fb41a2ae26d4c8c41a0e7be5d578224	-	September 9, 2025, 12:56:06 (UTC+05:30)	646.0 KB	Standard
jr_5e9b0ed9980d082bd22fb33e3653192fb41a2ae26d4c8c41a0e7be5d578224.inprogress	inprogress	September 9, 2025, 12:56:02 (UTC+05:30)	465.1 KB	Standard
jr_7f5ee0ed75f1d2ec38072cdd997d3726f7278c3c7730c80b3c2f0659eb43fff	-	September 9, 2025, 13:21:04 (UTC+05:30)	655.3 KB	Standard
jr_7f5ee0ed75f1d2ec38072cdd997d3726f7278c3c7730c80b3c2f0659eb43fff.inprogress	inprogress	September 9, 2025, 13:21:00 (UTC+05:30)	441.7 KB	Standard
jr_98cd77154c6500334c75e593bdb6369ae0d0de413c96c530e50f1502	-	September 9, 2025, 13:02:25 (UTC+05:30)	642.3 KB	Standard

Amazon Athena
Start querying data instantly.

Amazon Athena is an interactive query service that makes it easy to analyze data in Amazon S3 and other federated data sources using standard SQL.

Get started

- Query your data with Trino SQL. Use Query editor to analyze data on S3, on-premises, or in Amazon Redshift.
- Analyze your data using PySpark and Spark SQL. Use notebooks to build interactive Spark applications.

Launch query editor

Pricing

Region: Asia Pacific (Mumbai)

SQL queries per TB scanned	\$5.00 per TB
SQL queries on Provisioned Capacity	\$0.28 per DPU hour
Apache Spark calculations	\$0.45 per DPU hour

Getting started

Add the data source and identify the table and column names. Then you can begin querying your data right away. [Learn more](#)

Step 6: Athena Setup and Queries

1. Set S3 query output location:

Example: `s3://athena-query-results-54/`

2. In Athena Console:

- Data source: `AwsDataCatalog`
- Database: `ordersdb`
- Table: `shipped` (or whatever table is created by the S3 crawler)

Example SQL Queries

a) Top Customers by Spend

sql

```
SELECT Customer, SUM(Amount) AS TotalSpent
FROM shipped
```

```

WHERE Status = 'Shipped'
GROUP BY Customer
ORDER BY TotalSpent DESC;

```

b) Monthly Revenue

sql

```

SELECT OrderDate, SUM(Amount) AS MonthlyRevenue
FROM shipped
WHERE Status = 'Shipped'
GROUP BY OrderDate;

```

The screenshot shows the AWS Athena Query Editor. On the left, there's a sidebar for 'Data' with 'Data source' set to 'AwsDataCatalog', 'Catalog' set to 'None', and 'Database' set to 'ordersdb'. The main area has a 'Query 1' tab with the following SQL code:

```

1 SELECT Customer, SUM(Amount) AS TotalSpent
2 FROM shipped
3 WHERE Status = 'Shipped'
4 GROUP BY Customer
5 ORDER BY TotalSpent DESC;
6

```

Below the code, the 'SQL' tab shows 'Ln 2, Col 13'. At the bottom of the editor, there are buttons for 'Run again', 'Explain', 'Cancel', 'Clear', and 'Create'. The 'Query results' tab is selected, showing a table with one row:

Customer	TotalSpent
Alice	500

At the bottom right, there are buttons for 'Copy' and 'Download results CSV'. Above the results table, there are performance metrics: 'Time in queue: 60 ms', 'Run time: 710 ms', and 'Data scanned: 0.43 KB'.

Step 7: QuickSight Visualization

- Connect QuickSight to Athena (using the same database and tables)
- Example dashboards:
 - Bar: Customer vs. TotalSpent
 - Line: Monthly revenue trend

Step 8: Cleanup

1. Delete DynamoDB Table

- Go to DynamoDB Console.
- Select the Orders table.
- Click Delete table.
- Confirm deletion.

2. Delete AWS Glue Components

a) Delete Glue Jobs

- Go to AWS Glue Console → Jobs.
- Select all the jobs you created (related to this lab).
- Click Delete.

b) Delete Glue Crawlers

- Go to AWS Glue Console → Crawlers.
- Select all crawlers related to this lab (e.g., OrdersCrawler).
- Click Delete.

c) Delete Glue Databases and Tables (Optional)

- Go to AWS Glue Console → Databases.
- Select the databases (e.g., ordersdb).
- Delete related tables.
- Delete the database if desired.

3. Empty and Delete S3 Buckets

a) `my-orders-analytics` Bucket

- Navigate to S3 Console.
- Open the bucket.
- Select all folders/files (e.g., `shipped/`).
- Click Delete to empty bucket.
- Then delete the bucket.

b) Athena Query Result Bucket (e.g., `athena-query-results`)

- Similarly, empty all query result files.
- Delete the bucket.

c) `aws-glue-assets-*` Bucket

- Do not delete this bucket as it's managed by AWS for Glue internal use.

4. Delete AWS Athena Resources

Drop Athena Tables

- In Athena Console, connect to your database (ordersdb).
- Run `DROP TABLE` commands to remove created tables, e.g.,
- `sql`
- `DROP TABLE IF EXISTS orders;`