

The screenshot shows the AWS Home page. On the left, there's a sidebar with links to Services, Features, Resources, Documentation, Knowledge articles, Marketplace, Blog posts, Events, and Tutorials. The main content area displays three service cards: S3 Scalable Storage in the Cloud, S3 Glacier Archive Storage in the Cloud, and AWS Snow Family Large Scale Data Transport. To the right, there's a dashboard section with a "Create application" button and a "Find applications" search bar.

The screenshot shows the Amazon S3 Buckets page. It features an "Account snapshot - updated every 24 hours" section with a link to "View Storage Lens dashboard". Below it, there are tabs for "General purpose buckets" (selected) and "Directory buckets". A table lists three general purpose buckets: "haritha-bucket", "haritha28", and "my-unique-bucket-1234566", each with its creation date and an "IAM Access Analyzer" link.

The screenshot shows the "Create bucket" page. It has sections for "General configuration" (AWS Region set to US East (N. Virginia) us-east-1, Bucket type set to General purpose), "Bucket name" (set to "glue-bucket28"), and "Copy settings from existing bucket - optional" (Choose bucket). The "Object Ownership" section at the bottom discusses controlling ownership of objects from other accounts and using ACLs.

Screenshot of the AWS IAM Roles page:

Identity and Access Management (IAM)

Roles (4) Info

An IAM role is an identity you can create that has specific permissions with credentials that are valid for short durations. Roles can be assumed by entities that you trust.

Role name	Trusted entities	Last activity
AWSServiceRoleForOrganizations	AWS Service: organizations (Service-Linked Role)	-
AWSServiceRoleForSSO	AWS Service: sso (Service-Linked Role)	15 hours ago
AWSServiceRoleForSupport	AWS Service: support (Service-Linked Role)	-
AWSServiceRoleForTrustedAdvisor	AWS Service: trustedadvisor (Service-Linked Role)	-

Roles Anywhere Info

Authenticate your non AWS workloads and securely provide access to AWS services.

- Access AWS from your non AWS workloads**: Operate your non AWS workloads using the same authentication and authorization strategy that you use within AWS.
- X.509 Standard**: Use your own existing PKI infrastructure or use AWS Certificate Manager Private Certificate Authority to authenticate identities.
- Temporary credentials**: Use temporary credentials with ease and benefit from the enhanced security they provide.

<https://us-east-1.console.aws.amazon.com/iam/home?region=us-east-1#/home>

Screenshot of the "Create role" wizard Step 2: Trusted entity type:

Step 2

- Add permissions
- Name, review, and create

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: 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

Glue

Choose a use case for the specified service.

Use case

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

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

Screenshot of the "Create role" wizard Step 2: Add permissions:

Step 1

- Select trusted entity

Step 2

- Add permissions
- Name, review, and create

Add permissions Info

Permissions policies (1/1032) Info

Choose one or more policies to attach to your new role.

Filter by Type

Policy name	Type	Description
AmazonDMSRedshiftS3Role	AWS managed	Provides access to manage S3 settings ...
AmazonS3FullAccess	AWS managed	Provides full access to all buckets via t...
AmazonS3ObjectLambdaExecutionRolePolicy	AWS managed	Provides AWS Lambda functions permit...
AmazonS3OutpostsFullAccess	AWS managed	Provides full access to Amazon S3 on ...
AmazonS3OutpostsReadOnlyAccess	AWS managed	Provides read only access to Amazon S...
AmazonS3ReadOnlyAccess	AWS managed	Provides read only access to all bucket...
AmazonS3TablesFullAccess	AWS managed	Provides full access to all S3 table bu...
AmazonS3TablesReadOnlyAccess	AWS managed	Provides read only access to all S3 tabl...
AWSBackupServiceRolePolicyForS3Backup	AWS managed	Policy containing permissions necessar...

[CloudShell](#) [Feedback](#)

Add permissions

Permissions policies (2/1032) Info

Choose one or more policies to attach to your new role.

Filter by Type All types 7 matches

Policy name	Type	Description
AmazonAPIGatewayPushToCloudWatchLogs	AWS managed	Allows API Gateway to push logs to us...
AmazonDMSCloudWatchLogsRole	AWS managed	Provides access to upload DMS replicat...
AWSSyncPushToCloudWatchLogs	AWS managed	Allows AppSync to push logs to user's ...
AWSOpsWorksCloudWatchLogs	AWS managed	Enables OpsWorks instances with the ...
CloudWatchLogsCrossAccountSharingConfig...	AWS managed	Provides capabilities to manage Obser...
CloudWatchLogsFullAccess	AWS managed	Provides full access to CloudWatch Logs
CloudWatchLogsReadOnlyAccess	AWS managed	Provides read only access to CloudWat...

Set permissions boundary - optional

Name, review, and create

Role details

Role name
Enter a meaningful name to identify this role.
glue-role2@

Description
Add a short explanation for this role.
Allows Glue to call AWS services on your behalf.

Step 1: Select trusted entities

Trust policy

```

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

```

Step 1: Select trusted entities

Trust policy

```

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

```

Step 2: Add permissions

Permissions policy summary

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

Step 3: Add tags

us-east-1.console.aws.amazon.com/s3/buckets/glue-bucket28?region=us-east-1&bucketType=general&tab=objects

Amazon S3 > Buckets > glue-bucket28

glue-bucket28 Info

Objects | Metadata | Properties | Permissions | Metrics | Management | Access Points

Objects (0)

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

No objects
You don't have any objects in this bucket.

[Upload](#)

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

us-east-1.console.aws.amazon.com/s3/buckets/glue-bucket28/object/create_folder?region=us-east-1&bucketType=general

Amazon S3 > Buckets > glue-bucket28 > Create folder

Create folder Info

Use folders to group objects in buckets. When you create a folder, S3 creates an object using the name that you specify followed by a slash (/). This object then appears as folder on the console. [Learn more](#)

Your bucket policy might block folder creation
If your bucket policy prevents uploading objects without specific tags, metadata, or access control list (ACL) grantees, you will not be able to create a folder using this configuration. Instead, you can use the [upload configuration](#) to upload an empty folder and specify the appropriate settings.

Folder

Folder name /

Folder names can't contain "/". See [rules for naming](#).

Server-side encryption Info
Server-side encryption protects data at rest.

The following encryption settings apply only to the folder object and not to sub-folder objects.

Server-side encryption

Don't specify an encryption key
The bucket settings for default encryption are used to encrypt the folder object when storing it in Amazon S3.

Specify an encryption key
The specified encryption key is used to encrypt the folder object before storing it in Amazon S3.

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

us-east-1.console.aws.amazon.com/s3/buckets/glue-bucket28

Amazon S3 > Buckets > glue-bucket28

Successfully created folder "landing_zone".

glue-bucket28 Info

Objects | Metadata | Properties | Permissions | Metrics | Management | Access Points

Objects (1)

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
landing_zone/	Folder	-	-	-

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Screenshot of the AWS S3 console showing the 'landing_zone/' folder in the 'glue-bucket28' bucket. The 'Objects (0)' section is displayed, with a search bar and various actions like Copy S3 URI, Copy URL, Download, Open, Delete, Actions, Create folder, and Upload. A message states 'No objects' and 'You don't have any objects in this folder.' A prominent 'Upload' button is at the bottom.

Screenshot of the AWS S3 console showing the 'Create folder' dialog for the 'landing_zone/' folder. It includes a note about bucket policy restrictions and a 'Your bucket policy might block folder creation' alert. The 'Folder' section has a 'Folder name' input field containing 'customer'. The 'Server-side encryption' section is present. A note states 'The following encryption settings apply only to the folder object and not to sub-folder objects.'

Screenshot of the AWS S3 console showing the 'Create folder' dialog for the 'landing_zone/' folder. It includes a note about bucket policy restrictions and a 'Your bucket policy might block folder creation' alert. The 'Folder' section has a 'Folder name' input field containing 'customer'. The 'Server-side encryption' section is present. A note states 'The following encryption settings apply only to the folder object and not to sub-folder objects.'

Screenshot of the AWS S3 console showing the 'landing_zone/' folder. The folder contains one object named 'customers/'. The 'Actions' menu is open, showing options like Copy S3 URI, Copy URL, Download, Open, Delete, Actions, Create folder, and Upload.

Name	Type
customers/	Folder

Screenshot of the AWS S3 console showing the 'customers/' folder. The folder is empty. The 'Actions' menu is open, showing options like Copy S3 URI, Copy URL, Download, Open, Delete, Actions, Create folder, and Upload.

No objects

You don't have any objects in this folder.

Screenshot of the AWS S3 console showing the 'Create folder' dialog for the 'customers/' folder. The 'Folder name' field is set to 'load_date=20230603'. A note states that bucket policies might block folder creation. The 'Server-side encryption' section is shown, with the 'Don't specify an encryption key' option selected. The 'Actions' menu is open, showing options like Copy S3 URI, Copy URL, Download, Open, Delete, Actions, Create folder, and Upload.

Your bucket policy might block folder creation

If your bucket policy prevents uploading objects without specific tags, metadata, or access control list (ACL) grantees, you will not be able to create a folder using this configuration. Instead, you can use the [upload configuration](#) to upload an empty folder and specify the appropriate settings.

Folder

Folder name

load_date=20230603 /

Folder names can't contain "/". See [rules for naming](#).

Server-side encryption [info](#)

Server-side encryption protects data at rest.

The following encryption settings apply only to the folder object and not to sub-folder objects.

Server-side encryption

Don't specify an encryption key

The bucket settings for default encryption are used to encrypt the folder object when storing it in Amazon S3.

Specify an encryption key

The specified encryption key is used to encrypt the folder object before storing it in Amazon S3.

aws | ⚙️ | Search [Alt+S] | United States (N. Virginia) | Haritha

Amazon S3 > Buckets > glue-bucket28 > landing_zone/ > customers/

Successfully created folder "load_date=20230603".

customers/ Copy S3 URI

Objects Properties

Objects (1)

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
load_date=20230603/	Folder	-	-	-

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Open Downloads Search Downloads

Organize New folder

Today

- customers.csv Microsoft Excel Comma Separate... 3.88 KB
- boto3-user_accessKeys.csv Microsoft Excel Comma Separate... 99 bytes

Yesterday

- VSCodeUserSetup-x64-1.97.2.exe Visual Studio Code Setup Microsoft Corporation
- Earlier this week

 - cli-user_accessKeys.csv Microsoft Excel Comma Separate... 99 bytes
 - Git-2.48.1-64-bit.exe Git Setup The Git Development Community
 - AWSCLIV2.msi
 - Key Note.docx

File name: customers.csv All Files (*.*) Open Cancel

No files or folders You have not chosen any files or folders to upload.

Upload

Amazon S3 REST API. [Learn more](#)

Upload here, or choose Add files or Add folder.

Add files Add folder Remove

Destination Info

Destination s3://glue-bucket28/landing_zone/customers/load_date=20230603/

Destination details Bucket settings that impact new objects stored in the specified destination.

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

aws | ⚙️ | Search [Alt+S] | United States (N. Virginia) | Haritha

Upload: status Close

After you navigate away from this page, the following information is no longer available.

Summary

Succeeded	Failed
1 file, 3.9 KB (100.00%)	0 files, 0 B (0%)

Files and folders Configuration

Files and folders (1 total, 3.9 KB)

Find by name

Name	Folder	Type	Size	Status	Error
customers.csv	-	text/csv	3.9 KB	Succeeded	-

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

The screenshot shows the AWS S3 console interface. At the top, the navigation bar includes the AWS logo, a search bar, and a 'Search' button. The location bar shows the path: Amazon S3 > Buckets > glue-bucket28 > landing_zone/ > customers/ > load_date=20230603/. On the right side of the header, there are account information (United States (N. Virginia) - Haritha) and a sign-out button. Below the header, the main content area has a title 'load_date=20230603/'. Underneath it, there are tabs for 'Objects' (selected) and 'Properties'. A sub-header 'Objects (1)' is displayed. A message states: '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.' Below this is a search bar labeled 'Find objects by prefix' and a table listing the object 'customers.csv'. The table columns include Name, Type, Last modified, Size, and Storage class. The object 'customers.csv' is listed with a size of 3.9 KB and a storage class of Standard. At the bottom of the page, there are links for CloudShell, Feedback, and copyright information.

AWS Glue

AWS Glue is a serverless data integration service that allows users to prepare, transform, and analyze data. It is commonly used for ETL (Extract, Transform, Load) workflows.

Key Components of AWS Glue:

- **Glue Data Catalog** – A centralized metadata repository for managing datasets.
- **Glue Crawlers** – Automatically discover and catalog datasets.
- **Glue Jobs** – ETL scripts that process data.
- **Glue Triggers** – Automate job execution based on schedules or events.

AWS Glue ETL (Extract, Transform, Load):

AWS Glue ETL is used to process and transform large-scale data from various sources. It supports Python and Spark-based transformations.

The screenshot shows the AWS Glue console interface. At the top, the navigation bar includes the AWS logo, a search bar, and a 'Search' button. The location bar shows the path: AWS Glue > Databases. On the right side of the header, there are account information (United States (N. Virginia) - Haritha) and a sign-out button. Below the header, the main content area has a title 'Databases (0)'. A message states: 'A database is a set of associated table definitions, organized into a logical group.' Below this is a search bar labeled 'Filter databases' and a table listing databases. The table columns include Name, Description, Location URI, and Created on (UTC). The message 'No resources' is displayed. At the bottom of the page, there are links for CloudShell, Feedback, and copyright information.

Create a database

Create a database in the AWS Glue Data Catalog.

Database details

Name
project_db

Database name is required, in lowercase characters, and no longer than 255 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.

An S3 location is required for managed tables and Zero-ETL integrations.

Create database

Databases (1)

A database is a set of associated table definitions, organized into a logical group.

Name		Description	Location URI	Created on (UTC)
project_db				February 17, 2025 at 01:34:41

project_db

Database properties

Name	Description	Location	Created on (UTC)
project_db	-	-	February 17, 2025 at 01:34:41

Tables (0)

View and manage all available tables.

Name	Database	Location	Classification	Deprecated	View data	Data quality	Column stats...
No available tables							

AWS Glue > Tables > Add table

Announcing new optimization features for Apache Iceberg tables
Optimize storage for Apache Iceberg tables with automatic snapshot retention and orphan file deletion. Learn more.

Step 1: Set table properties

Step 2: Choose or define schema

Step 3: Review and create

Table details

Name: customers

If you plan to access the table from Amazon Athena, then the name should be under 256 characters and contain only lowercase letters (a-z), numbers (0-9), and underscore (_). For more information, see [Athena names](#).

Database: project_db

Description - optional: Enter a description

Descriptions can be up to 2048 characters long.

Table format

Data Catalog managed tables support data optimization for Apache Iceberg table type. Learn more.

Standard AWS Glue table (default)
Create a standard AWS Glue table.

Apache Iceberg table
Create a table in Apache Iceberg table format.

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Tables > Add table

Choose S3 path

S3 buckets

Buckets (4)

Find bucket:

Name	Creation date
glue-bucket28	February 17, 2025 at 00:39:25
haritha-bucket	February 16, 2025 at 21:22:17
haritha28	February 16, 2025 at 18:39:52
my-unique-bucket-1234566	February 16, 2025 at 04:35:16

Data store

Select the type of source

S3

Data format

Classification

Choose the format of the data in your table.

Avro

CSV

JSON

XML

Parquet

Cancel Choose

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Tables > Add table

Standard AWS Glue table (default)
Create a standard AWS Glue table.

Apache Iceberg table
Create a table in Apache Iceberg table format.

Data store

Select the type of source

S3

Choose S3 path

S3 buckets > glue-bucket28 > landing_zone/

Objects (1/1)

Find object by prefix:

Key	Last modified	Size
customers/	-	-

Data format

Classification

Choose the format of the data in your table.

Avro

CSV

JSON

XML

Parquet

Cancel Choose

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Tables > Add table

AWS Glue

- Getting started
- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations [New](#)

Data Catalog

- Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings

Data Integration and ETL

Legacy pages

What's New [New](#)

Documentation [New](#)

CloudShell Feedback

Search [Alt+S]

Kinesis
Kafka

Data location is specified in
 my account
 another account

Include path
 s3://glue-bucket28/landing_zone/custom [View](#) [Browse S3](#)

Data format

Classification
Choose the format of the data in your table.
 Avro
 CSV
 JSON
 XML
 Parquet
 ORC

Delimiter
 Comma (,)

Cancel Next

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Tables > Add table

AWS Glue

- Getting started
- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations [New](#)

Data Catalog

- Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings

Data Integration and ETL

Legacy pages

What's New [New](#)

Documentation [New](#)

CloudShell Feedback

Search [Alt+S]

Step 2 Choose or define schema

Step 3 Review and create

Schema

Define or upload schema
Manually define schema

Choose from Glue Schema Registry
Select existing schema from your Glue Schema Registry.

Schema (0)
View and manage the table schema.

Filter schemas

#	Column name	Data type	Partition key	Comment
No table schema				

Partition indexes - optional (0)
AWS Glue partition indexes are an important configuration to reduce overall data transfers and processing, and reduce query processing time. A maximum of 3 partition indexes can be created on a given table.

Index name	Index keys	Index status
No indexes		

Provide indexes

Add indexes

Cancel Previous Next

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Databases > project_db

AWS Glue

- Getting started
- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations [New](#)

Data Catalog

- Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings

Data Integration and ETL

Legacy pages

What's New [New](#)

Documentation [New](#)

CloudShell Feedback

Search [Alt+S]

project_db

Last updated (UTC) February 17, 2025 at 01:45:21 [Edit](#) [Delete](#)

Database properties

Name	Description	Location	Created on (UTC)
project_db	-	-	February 17, 2025 at 01:34:41

Tables (0)
View and manage all available tables.

Filter tables

Name	Database	Location	Classification	Deprecated	View data	Data quality	Column statis...
No available tables							

Last updated (UTC) February 17, 2025 at 01:43:23 [Edit](#) [Delete](#) [Add tables using crawler](#) [Add table](#)

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Crawlers > Add crawler

Set crawler properties

Crawler details info

Name: customers
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**

AWS Glue

- Getting started
- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations New
- Data Catalog**
 - Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings
- Data Integration and ETL**
- Legacy pages**

What's New [?] Documentation [?]

CloudShell Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Crawlers > Add crawler

Choose data sources and classifiers

Data source configuration
Is your data already mapped to Glue tables?
 Not yet
Select one or more data sources to be crawled.
 Yes
Select existing tables from your Glue Data Catalog.

Data sources (0) Info
The list of data sources to be scanned by the crawler.

Type	Data source	Parameters
You don't have any data sources.		
Add a data source		

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 Glue

- Getting started
- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations New
- Data Catalog**
 - Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings
- Data Integration and ETL**
- Legacy pages**

What's New [?] Documentation [?]

CloudShell Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Crawlers > Add crawler

Add data source

Data source
Choose the source of data to be crawled.
 S3

Network connection - optional
Optionally include a Network connection to use with this S3 target. Note that each crawler is limited to one Network connection so any other S3 targets will also use the same connection (or none, if left blank).
 Clear selection Add new connection

Location of S3 data
 In this account
 In a different account
Browse for or enter an existing S3 path.
 s3://glue-bucket28/landing_zone/cust View [?] Browse S3

Subsequent crawler runs
This field is a global field that affects all S3 data sources.
 Crawl all sub-folders
Crawl all folders again with every subsequent crawl.
 Crawl new sub-folders only
Only Amazon S3 folders that were added since the last crawl will be crawled. If the schemas are compatible, new partitions will be added to existing tables.
 Crawl based on events
Rely on Amazon S3 events to control what folders to crawl.

Cancel **Add an S3 data source**

AWS Glue

- Getting started
- ETL jobs
- Visual ETL
- Notebooks
- Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations New
- Data Catalog**
 - Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings
- Data Integration and ETL**
- Legacy pages**

What's New [?] Documentation [?]

CloudShell Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Crawlers > Add crawler

Choose data sources and classifiers

Data source configuration
Is your data already mapped to Glue tables?

- Not yet Select one or more data sources to be crawled.
- Yes Select existing tables from your Glue Data Catalog.

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

Type	Data source	Parameters
S3	s3://glue-bucket28/landing_zone/cu...	Recrawl all

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 Glue > Crawlers > Add crawler

Configure security settings

IAM role Info
Existing IAM role
gule-role28

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.

Lake Formation configuration - optional
Allow the crawler to use Lake Formation credentials for crawling the data source. [Learn more](#).

Use Lake Formation credentials for crawling S3 data source
Checking this box will allow the crawler to use Lake Formation credentials for crawling the data source. If the data source is registered in another account, you must provide the registered account ID. Otherwise, the crawler will crawl only those data sources associated to the account. Only applicable to S3, Glue Catalog, Iceberg, and Hudi data sources.

Security configuration - optional
Enable at-rest encryption with a security configuration.

Cancel Previous Next

AWS Glue > Crawlers > Add crawler

Set output and scheduling

Output configuration Info
Target database
project_db

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

Maximum table threshold - optional
This field sets the maximum number of tables the crawler is allowed to generate. In the event that this number is surpassed, the crawl will fail with an error. If not set, the crawler will automatically generate the number of tables depending on the data schema.
Type a number greater than 0

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

us-east-1.console.aws.amazon.com/glue/home?region=us-east-1#/v2/data-catalog/crawlers/add

AWS Glue > Crawlers > Add crawler

AWS Glue

- Getting started
- ETL jobs
 - Visual ETL
 - Notebooks
 - Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations [New](#)

Data Catalog

- Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings

Data Integration and ETL

Legacy pages

What's New [?](#) Documentation [?](#)

CloudShell Feedback

Review and create

Step 1: Set crawler properties

Set crawler properties

Name	Description	Tags
customers	-	-

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
S3	s3://glue-bucket28/landing_zone/cust...	Recrawl all

Step 3: Configure security settings

Configure security settings

IAM role	Security configuration	Lake Formation configuration
glue-role28	-	-

Step 4: Set output and scheduling

Set output and scheduling

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

aws Search [Alt+S]

AWS Glue > Crawlers > customers

AWS Glue

- Getting started
- ETL jobs
 - Visual ETL
 - Notebooks
 - Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations [New](#)

Data Catalog

- Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings

Data Integration and ETL

Legacy pages

What's New [?](#) Documentation [?](#)

CloudShell Feedback

customers

One crawler successfully created
The following crawler is now created: "customers"

Last updated (UTC) February 17, 2025 at 01:54:57

[Run crawler](#) [Edit](#) [Delete](#)

Crawler properties

Name	IAM role	Database	State
customers	glue-role28 ?	project_db	READY
Description	Security configuration	Lake Formation configuration	Table prefix
-	-	-	-

Maximum table threshold

Advanced settings

Crawler runs (0)

The list of crawler runs for this crawler.

Filter data	Filter by a date and time range	Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
-	-	-	-	-	-	-	-

You don't have any crawler runs.

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

aws Search [Alt+S]

AWS Glue > Crawlers

AWS Glue

- Getting started
- ETL jobs
 - Visual ETL
 - Notebooks
 - Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations [New](#)

Data Catalog

- Databases
 - Tables
 - Stream schema registries
 - Schemas
 - Connections
 - Crawlers
 - Classifiers
 - Catalog settings

Data Integration and ETL

Legacy pages

What's New [?](#) Documentation [?](#)

CloudShell Feedback

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 (1) [Info](#)

View and manage all available crawlers.

Filter crawlers	Name	State	Schedule	Last run	Last run times...	Log	Table changes fr...
-	customers	Stopping	-	-	-	-	-

Last updated (UTC) February 17, 2025 at 01:56:53

[Action](#) [Run](#) [Create crawler](#)

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS Glue > Crawlers > customers

customers

Last updated (UTC) February 17, 2025 at 01:54:57

Crawler properties

Name	IAM role glue-role28	Database	State
customers		project_db	READY
Description	Security configuration	Lake Formation configuration	Table prefix
-	-	-	-

Advanced settings

Crawler runs

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
February 17, 2025 at 01:55:29	February 17, 2025 at 01:56:31	01 min 02 s	Failed	-	-

Stop run View CloudWatch logs View run details

CloudShell Feedback

AWS Glue > Crawlers > customers

customers

Last updated (UTC) February 17, 2025 at 01:54:57

Crawler properties

Name	IAM role	Database	State
customers	glue-role28	project_db	READY
Description	Security configuration	Lake Formation configuration	Table prefix
-	-	-	-

Advanced settings

Crawler runs (2)

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
February 17, 2025 at 02:02:19	February 17, 2025 at 02:04:54	02 min 35 s	Completed	-	1 table change, 1 partition change
February 17, 2025 at 01:55:29	February 17, 2025 at 01:56:31	01 min 02 s	Failed	-	-

Stop run View CloudWatch logs View run details

CloudShell Feedback

AWS Glue > Tables

Tables

A table is the metadata definition that represents your data, including its schema. A table can be used as a source or target in a job definition.

Tables (1)

Last updated (UTC) February 17, 2025 at 02:09:20

Name	Database	Location	Classification	Deprecated	View data	Data quality	Column stats...
customers	project_db	s3://glue-bucket28/	CSV	-	Table data	View data quality	View statistics

Add tables using crawler Add table

View and manage all available tables.

Filter tables

CloudShell Feedback

AWS Glue > Databases > project_db

project_db

Last updated (UTC) February 17, 2025 at 02:10:49 Edit Delete

Database properties

Name	Description	Location	Created on (UTC)
project_db	-	-	February 17, 2025 at 01:34:41

Tables (1)

Last updated (UTC) February 17, 2025 at 02:10:50 Delete Add tables using crawler Add table

Name	Database	Location	Classification	Deprecated	View data	Data quality	Column stats...
customers	project_db	s3://glue-bucket28/	CSV	-	Table data	View data quality	View statistics

CloudShell Feedback

AWS Glue > Tables > customers

customers

Last updated (UTC) February 17, 2025 at 02:11:59 Version 1 (Current version) Actions

Table overview Data quality - new

Table details

Name	customers	Classification	Deprecated
Database	project_db	CSV	-
Description	-	Location	s3://glue-bucket28/landing_zone/customers/
Last updated	February 17, 2025 at 02:04:54	Connection	-

Advanced properties

Schema (6)

Edit schema as JSON Edit schema

#	Column name	Data type	Partition key	Comment
1	-	-	-	-

CloudShell Feedback

us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#query-editor

Amazon Athena > Query editor

Catalog None Database project_db

Tables and views Create

Tables (1) < 1 >

customerid	bigint
gender	string
age	bigint
annual income (k\$)	bigint
spending score (1-100)	bigint
load_date	string (Partitioned)

Views (0) < 1 >

SQL Ln 1, Col 1

Run Query Preview Table Generate table DDL Load partitions Insert Insert into editor Manage Delete table Generate statistics - new View properties

Reuse query results up to 60 minutes ago

Copy Download results CSV

No results Run a query to view results

CloudShell Feedback

aws | Search [Alt+S]

Amazon S3 > Buckets > glue-bucket28 > Create folder

Create folder Info

Use folders to group objects in buckets. When you create a folder, S3 creates an object using the name that you specify followed by a slash (/). This object then appears as folder on the console. [Learn more](#)

Your bucket policy might block folder creation
If your bucket policy prevents uploading objects without specific tags, metadata, or access control list (ACL) grantees, you will not be able to create a folder using this configuration. Instead, you can use the [upload configuration](#) to upload an empty folder and specify the appropriate settings.

Folder

Folder name /
Folder names can't contain "/". [See rules for naming](#)

Server-side encryption Info

Server-side encryption protects data at rest.

The following encryption settings apply only to the folder object and not to sub-folder objects.

Server-side encryption

Don't specify an encryption key
The bucket settings for default encryption are used to encrypt the folder object when storing it in Amazon S3.

Specify an encryption key

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

aws | Search [Alt+S]

Amazon Athena > Query editor

Settings successfully updated.

Editor Recent queries Saved queries **Settings**

Workgroup primary

Query result and encryption settings

Query result location s3://glue-bucket28/query_results/

Encrypt query results -

Expected bucket owner -

Assign bucket owner full control over query results Turned off [Manage](#)

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

aws | Search [Alt+S]

Amazon Athena > Query editor

Tables and views [Create](#)

Filter tables and views

Tables (1) < 1 >

customers Partitioned

Views (0) < 1 >

SQL Ln 1, Col 1

Run again Explain Cancel Clear Create

Reuse query results up to 60 minutes ago

Query results [Query stats](#)

Completed Time in queue: 65 ms Run time: 1.155 sec Data scanned: 3.89 KB

Results (10)

Search rows

#	customerid	gender	age	annual income (k\$)	spending score (1-100)	load_date
1	1	Male	19	15	39	20230603
2	2	Male	21	15	81	20230603
3	3	Female	20	16	6	20230603
4	4	Female	23	16	77	20230603
5	5	Female	31	17	40	20230603
6	6	Female	22	17	76	20230603

Copy Download results CSV

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Screenshot of the AWS Glue Connectors page:

The left sidebar shows navigation links for Notebooks, Data Catalog tables, Data connections (selected), Workflows (orchestration), Zero-ETL integrations, Data Catalog, Databases, Tables, Stream schema registries, Schemas, Connections, Crawlers, Classifiers, Catalog settings, Data Integration and ETL, and Legacy pages.

The main content area has two sections: Marketplace connectors and Custom connectors.

- Marketplace connectors:** Subtitle: "Subscribe to connectors from AWS partners to expand your data sources." Includes a "Go to AWS Marketplace" button.
- Custom connectors:** Subtitle: "Provide your own connector to expand your data sources. Creating custom connectors". Includes a "Create custom connector" button.

Connectors (0) Info: You can manage your connectors or use them to create connections. Includes a search bar and a table header: Name, Status, Type, Last modified. Below it says "No connectors" and "Create custom connector".

Connections (0) Info: You can manage your connections or use a connection in a job. Includes a search bar and a table header: Name, Status, Type, Last modified, Version. Below it says "No connections" and "Create custom connector".

Footer: © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Screenshot of the AWS Glue Create connection page:

The left sidebar shows the same navigation links as the previous screenshot.

The main content area shows a step-by-step process: Step 2 (Configure connection), Step 3 (Set properties), and Step 4 (Review and create). To the right is a "Data sources (72)" section with a search bar and a table header: Grid, List, Name, Status, Type, Last modified. It lists several data sources:

- Adobe Analytics - new**: Adobe Analytics is a business intelligence solution for combining and analyzing data from any digital point in the customer journey. Includes a "Learn more" link.
- Amazon Aurora**: Connect to Amazon Aurora. Includes a "Learn more" link.
- Amazon DocumentDB - new**: Connect to Amazon DocumentDB. Includes a "Learn more" link.
- Amazon OpenSearch Service - new**: Connect to Amazon OpenSearch Service. Includes a "Learn more" link.
- Amazon Redshift**: Connect to Amazon Redshift. Includes a "Learn more" link.
- Apache Kafka**: Connect to streaming data in Apache Kafka. Includes a "Learn more" link.
- Asana - new**
- Azure Cosmos - new**

Footer: © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Screenshot of the AWS Glue Studio Jobs page:

The left sidebar shows the same navigation links as the previous screenshots.

The main content area includes:

- AWS Glue Studio Info:** Subtitle: "Create job". Options: Visual ETL (selected), Notebook, Script editor.
- Example jobs Info:** Subtitle: "Create example job".
- Your jobs (0) Info:** Subtitle: "Your jobs". Includes a search bar and a table header: Job name, Type, Created by, Last modified, AWS Glue version. Below it says "No jobs" and "You have not created a job yet." Includes a "Create job from a blank graph" button.

Footer: © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Screenshot of AWS Glue Studio job configuration:

Untitled job

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

Choose an S3 bucket

S3 buckets > glue-bucket28

Objects (1/3)

Key

- landing_zone/ (unchecked)
- query_results/ (unchecked)
- transform_zone/ (checked)

Last updated: February 16, 2025, 20:43:39 (UTC-6)

To start a data preview session, choose an IAM role for this job. Changing the role will end an existing data preview session.

Create IAM role. ▾ Additional Settings

Cancel Choose

Add a partition key

CloudShell Feedback

us-east-1.console.aws.amazon.com/gluestudio/home?region=us-east-1#/editor/blank/details

Customers_ETL_dropField

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

Basic properties Info

Name: Customers_ETL_dropField

Description - optional:

Descriptions can be up to 2048 characters long.

IAM Role: glue-role28

Type: Spark

Glue version: Info

Glue 5.0 - Supports spark 3.5, Scala 2, Python 3

CloudShell Feedback

Successfully updated job

Successfully updated job Customers_ETL_dropField. To run the job choose the Run Job button.

Customers_ETL_dropField

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

Last modified on 2/16/2025, 8:46:42 PM

Basic properties Info

Name: Customers_ETL_dropField

Description - optional:

Descriptions can be up to 2048 characters long.

IAM Role: glue-role28

Type: Spark

CloudShell Feedback

Customers_ETL_dropField

Last modified on 2/16/2025, 8:46:42 PM [Actions](#) [Save](#) [Run](#)

Visual [Script](#) Job details Runs Data quality Schedules Version Control

Script (Locked) Info

```

1 import sys
2 from awsglue.transforms import *
3 from awsglue.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from awsglue.context import glueContext
6 from awsglue.job import Job
7 from awsglued.transforms import EvaluateDataQuality
8
9 args = getResolvedOptions(sys.argv, ['JOB_NAME'])
10 sc = SparkContext()
11 glueContext = glueContext(sc)
12 spark = glueContext.spark_session
13 job = Job(glueContext)
14 job.init(args['JOB_NAME'], args)
15
16 # Default ruleset used by all target nodes with data quality enabled
17 DEFAULT_DATA_QUALITY_RULESET = """
18 Rules = [
19
20     ]
21 """
22
23 # Script generated for node Amazon S3
24 AmazonS3_node1739759665838 = glueContext.create_dynamic_frame.from_catalog(database="project_db", table_name="customers", transformation_ctx="AmazonS3_node1739759665838")
25
26 # Script generated for node Drop Fields
27 DropFields_node1739759710768 = DropFields.apply(frame=AmazonS3_node1739759665838, paths=["gender"], transformation_ctx="DropFields_node1739759710768")
28
29 # Script generated for node Amazon S3
30 EvaluateDataQuality().process_rows(frame=DropFields_node1739759710768, ruleset=DEFAULT_DATA_QUALITY_RULESET, publishing_options={"dataQualityEvaluationContext": "EvaluateDataQuality_node_AmazonS3_node1739759868999": glueContext.write_dynamic_frame.from_options(frame=DropFields_node1739759710768, connection_type="s3", format="glueparquet", connection_options={"path": "s3://bucket-name/path"}}, transformation_ctx="EvaluateDataQuality_node_AmazonS3_node1739759868999")
31
32 job.commit()
33

```

[Download script](#) [Edit script](#)

Customers_ETL_dropField

Last modified on 2/16/2025, 8:46:42 PM [Actions](#) [Save](#) [Run](#)

Visual [Script](#) Job details Runs Data quality Schedules Version Control

Script (Locked) Info

```

16 # Default ruleset used by all target nodes with data quality enabled
17 DEFAULT_DATA_QUALITY_RULESET = """
18 Rules = [
19     columnCount > 0
20 ]
21 """
22
23 # Script generated for node Amazon S3
24 AmazonS3_node1739759665838 = glueContext.create_dynamic_frame.from_catalog(database="project_db", table_name="customers", transformation_ctx="AmazonS3_node1739759665838")
25
26 # Script generated for node Drop Fields
27 DropFields_node1739759710768 = DropFields.apply(frame=AmazonS3_node1739759665838, paths=["gender"], transformation_ctx="DropFields_node1739759710768")
28
29 # Script generated for node Amazon S3
30 EvaluateDataQuality().process_rows(frame=DropFields_node1739759710768, ruleset=DEFAULT_DATA_QUALITY_RULESET, publishing_options={"dataQualityEvaluationContext": "EvaluateDataQuality_node_AmazonS3_node1739759868999": glueContext.write_dynamic_frame.from_options(frame=DropFields_node1739759710768, connection_type="s3", format="glueparquet", connection_options={"path": "s3://bucket-name/path"}}, transformation_ctx="EvaluateDataQuality_node_AmazonS3_node1739759868999")
31
32 job.commit()
33

```

[Download script](#) [Edit script](#)

Customers_ETL_dropField

Last modified on 2/16/2025, 8:46:42 PM [Actions](#) [Save](#) [Run](#)

Visual [Script](#) Job details [Runs](#) Data quality Schedules Version Control

Job runs (1/1) Info

Last updated (UTC) February 17, 2025 at 02:50:34 [View details](#) [Stop job run](#) [Troubleshoot with AI](#)

Run status	Retries	Start time (Local)	End time (Local)	Duration	Capacity (DPUs)	Worker type	Glue version
Succeeded	0	02/16/2025 20:48:40	02/16/2025 20:50:20	1 m 30 s	10 DPUs	G.1X	5.0

[Table View](#) [Card View](#)

Run details [Input arguments \(10\)](#) [Continuous logs](#) [Run insights](#) [Metrics](#) [Troubleshooting analysis - preview](#) [Spark UI](#)

Job name	Start time (Local)	Glue version	Last modified on (Local)
Customers_ETL_dropField	02/16/2025 20:48:40	5.0	02/16/2025 20:50:20
Id	End time (Local)	Worker type	Log group name

[CloudShell](#) [Feedback](#)

AWS Glue > Triggers

Triggers

A trigger starts a job when it fires.

Triggers (0)

View and manage all available triggers.

Last updated (UTC) February 17, 2025 at 02:52:02 Action Add trigger

Name	Status	Type	Parameters	Targets
No resources No resources to display.				

CloudShell Feedback

AWS Glue > Triggers > Add trigger

Set trigger properties

Step 1: Set trigger properties (selected)

Step 2: Choose jobs or crawlers to activate

Step 3: Review and create

Trigger details

Name: customer_ETL_trigger

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.

Trigger type

On demand: Fire the trigger immediately when started. (Selected)

Schedule: Fire the trigger on a timer.

Job or crawler event: Fire the trigger when job or crawler events match your watched list.

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 ↗

CloudShell Feedback

AWS Glue > Triggers > Add trigger

Choose jobs or crawlers to activate

Step 1: Set trigger properties

Step 2: Choose jobs or crawlers to activate (selected)

Step 3: Review

Resources to trigger (0)

Add target

Resource type: Job

The type of resource for this trigger target.

Job name: Customers_ETL_dropField

Job to start when this trigger fires

► Parameters passed down to job "Customers_ETL_dropField" when started - optional

Cancel Add

CloudShell Feedback

AWS Glue > Triggers > Add trigger

Review and create

Step 1: Set trigger properties

Trigger details

Name	Description	Tags	Schedule
customer_ETL_trigger	-	-	At 1 minutes past the hour

Step 2: Choose jobs or crawlers to activate

Resources to trigger (1)

List of resources to start once the trigger activates

Type	Name	Parameters
Job	Customers_ETL_dropField	-

Enable trigger on creation

Create

CloudShell Feedback

AWS Glue > Triggers > customer_ETL_trigger

customer_ETL_trigger

Last updated (UTC) February 17, 2025 at 0:25:09 **Edit trigger**

Trigger properties

Name	Description
customer_ETL_trigger	-

Status Created

Associated workflow -

Schedule At 1 minutes past the hour

Target resources

Resources to trigger (1)

List of resources to start once the trigger activates

Type	Name	Parameters
Job	Customers_ETL_dropField	-

CloudShell Feedback

AWS Glue > Triggers

Triggers

A trigger starts a job when it fires.

Last updated (UTC) February 17, 2025 at 02:55:45 **Action** **Add trigger**

Triggers (1)

View and manage all available triggers.

Name	Status	Type	Parameters	Targets
customer_ETL_trigger	Created	Scheduled	At 1 minutes past the hour	1 job: Customers_ETL_dropField

CloudShell Feedback

