

The screenshot shows the AWS CloudWatch Metrics console. At the top, there are tabs for ChatGPT, Nagendra K | LinkedIn, AWS Skill Builder, Cloud Quest, Amazon Data Firehose, Terraform AWS Engineer, and a plus sign. The URL in the address bar is `us-east-1.console.aws.amazon.com/firehose/home?region=us-east-1#/home`. The main content area features a banner about Amazon Data Firehose supporting Apache Iceberg tables. Below the banner, the title "Amazon Data Firehose" is displayed, followed by a subtitle "Real-time streaming delivery for any data, at any scale, and at low-cost." A call-to-action button "Create Firehose stream" is visible. On the left, there's a sidebar titled "Analytics". On the right, a "Getting started" section provides instructions on creating a Firehose stream. The bottom of the page includes a "How it works" section and a "Pricing (United States (N. Virginia))" section.

The screenshot shows the AWS CloudWatch Metrics console. At the top, there are tabs for ChatGPT, Nagendra K | LinkedIn, AWS Skill Builder, Cloud Quest, Amazon Data Firehose, Terraform AWS Engineer, and a plus sign. The URL in the address bar is `us-east-1.console.aws.amazon.com/firehose/home?region=us-east-1#/create/put/s3`. The main content area shows the "Create Firehose stream" wizard. Step 1, "Choose source and destination", is selected. It shows "Source" set to "Direct PUT" and "Destination" set to "Amazon S3". Step 2, "Firehose stream name", is shown with the name "ClickStreamData" entered. Step 3, "Transform and convert records - optional", is shown with the note "Configure Amazon Data Firehose to transform and convert your raw data." The bottom of the page includes a "How it works" section and a "Pricing (United States (N. Virginia))" section.

Screenshot of the AWS Lambda function selection dialog:

**Choose an AWS Lambda function**

**AWS Lambda functions (4)**

Function name	Description	Runtime	Timeout
gbl_lab_monitoring		python3.11	5 minutes
AnalyticsDestinationFunction		python3.11	1 minute
LabStack-98ea6b70-bd2e-46-GbLabMonitoringgblproto-ePnvT6xj7z41		python3.11	1 minute
DataProcessingFunction		python3.11	1 minute

Minimum: 0 seconds, maximum: 900 seconds.

**Convert record format** | Info

Data in Apache Parquet or Apache ORC format is typically more efficient to query than JSON. Amazon Data Firehose can convert your JSON-formatted source records using a schema from a table defined in AWS Glue. For records that aren't in JSON format, create a Lambda function that converts them to JSON in the Transform source records with AWS Lambda section above.

Enable record format conversion

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences 34°F Sunny 08:56 29-11-2025

Screenshot of the AWS Lambda function configuration dialog:

**Transform and convert records - optional**

Configure Amazon Data Firehose to transform and convert your record data.

**Transform source records with AWS Lambda** | Info

Amazon Data Firehose can invoke an AWS Lambda function to transform, filter, decompress, convert and process your source data records. The specified AWS Lambda function can also be used to provide data partitioning keys for the incoming source data before its delivery to the specified destination.

Turn on data transformation

**AWS Lambda function**

arn:aws:lambda:us-east-1:734979449836:function:DataProcessingFunction:\$LATEST

Version or alias

\$LATEST

Browse  Create function

**Buffer size**

The AWS Lambda function has a 6 MB invocation payload quota. Your data can expand in size after it's processed by the AWS Lambda function. A smaller buffer size allows for more room should the data expand after processing.

1 MB

Minimum: 0.2 MB, maximum: 3 MB.

**Buffer interval**

The period of time during which Amazon Data Firehose buffers incoming data before invoking the AWS Lambda function. The AWS Lambda function is invoked once the value of the buffer size or the buffer interval is reached.

60 seconds

Minimum: 0 seconds, maximum: 900 seconds.

**Convert record format** | Info

Data in Apache Parquet or Apache ORC format is typically more efficient to query than JSON. Amazon Data Firehose can convert your JSON-formatted source records using a schema from a table defined in AWS Glue. For records that aren't in JSON format, create a Lambda function that converts them to JSON in the Transform source records with AWS Lambda section above.

Enable record format conversion

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences 34°F Sunny 08:57 29-11-2025

Screenshot of the AWS Cloud Console showing the creation of a new Firehose stream. The destination is set to an S3 bucket named 's3://databucket-3fc58d70'. The 'New line delimiter' option is disabled. The 'Dynamic partitioning' section is visible, and the 'S3 bucket prefix - optional' field contains 'processed\_data/'. The 'S3 bucket error output prefix - optional' field contains 'error/'. The status bar at the bottom shows a weather icon for 34°F and a timestamp of 09:02.

Screenshot of the AWS Cloud Console showing the successful creation of a ClickStreamData Firehose stream. The stream details include: Status: Active; Destination: Amazon S3; ARN: arn:aws:firehose:us-east-1:734979449836:deliverystream/ClickStreamData; Data transformation: Enabled; Dynamic partitioning: Not enabled; Creation time: November 29, 2025 at 09:19 EST; Error logs status: 0 Destination error logs. The 'Monitoring' tab is selected. The status bar at the bottom shows a weather icon for 36°F and a timestamp of 09:20.

The screenshot shows the AWS Glue Data Catalog homepage. The left sidebar includes links for AWS Glue, Data Catalog, Data Integration and ETL, and Legacy pages. The main content area features a large banner for AWS Glue, Serverless data integration, with a call-to-action button 'Get started'. Below the banner is a section titled 'What's new in Glue' listing three recent changes: support for audit context with Lake Formation, support for write operations with AWS Lake Formation fine-grained access controls, and the introduction of AWS Glue 5.1. To the right is a 'Benefits and features' section with tabs for 'AWS Glue Data Catalog' and 'Crawlers for data discovery'. A 'Pricing (US)' table shows rates for Jobs (\$0.44 per DPU-Hour) and Crawlers (\$0.44 per DPU-Hour). The bottom navigation bar includes CloudShell, Feedback, and a weather icon.

The screenshot shows the 'Databases' page under the AWS Glue section. The left sidebar lists AWS Glue, Data Catalog (Tables, Stream schema registries, Schemas, Connections, Crawlers, Classifiers, Catalog settings), Data Integration and ETL, and Legacy pages. The main content displays a table of databases, showing one entry: 'firehose-ingestion' with a 'Created on (UTC)' timestamp of November 29, 2025 at 14:26:27. The top navigation bar includes ChatGPT, LinkedIn, AWS Skill Builder, Cloud Quest, AWS Glue Console, Terraform AWS Engine, and other tabs. The bottom navigation bar includes CloudShell, Feedback, and a weather icon.

The screenshot shows the AWS Glue console with the 'Crawlers' page open. A new crawler is being created, and the first step, 'Set crawler properties', is selected. The left sidebar shows navigation options like 'AWS Glue', 'Data Catalog tables', and 'Data Catalog'. The main panel displays the 'Set crawler properties' configuration screen. The 'Name' field is populated with 'crawl\_processed\_data'. There are sections for 'Description - optional' and 'Tags - optional'. Buttons for 'Cancel' and 'Next' are at the bottom right.

The screenshot shows the 'Add data source' dialog box overlaid on the AWS Glue console. The 'Choose data sources' step is selected. The 'Data source' dropdown is set to 'S3'. Under 'Network connection - optional', there is a note about using a network connection for S3 targets. The 'Location of S3 data' section has 'In this account' selected. The 'S3 path' input field contains 's3://databucket-3fc58d70/processed'. The 'Subsequent crawler runs' section offers options to 'Crawl all sub-folders' or 'Crawl new sub-folders only'. Buttons for 'Cancel', 'Add an S3 data source', and 'Next' are visible at the bottom.

Screenshot of the AWS Glue console showing the "Choose data sources and classifiers" step of a crawler setup.

**AWS Glue** sidebar:

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

**Choose data sources and classifiers** wizard:

- Step 1: Set crawler properties
- Step 2: **Choose data sources and classifiers** (selected)
- Step 3: Configure security settings
- Step 4: Set output and scheduling
- Step 5: Review and create

**Data source configuration** section:

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://databucket-3fc58d70//process...	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.

Buttons: Cancel, Previous, Next

System status bar at the bottom:

- CloudShell Feedback
- 36°F Sunny
- Search
- © 2025, Amazon Web Services, Inc. or its affiliates.
- Privacy Terms Cookie preferences
- ENG IN 09:34 29-11-2025

Screenshot of the AWS Glue console showing the "Configure security settings" step of a crawler setup.

**AWS Glue** sidebar:

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

**Configure security settings** wizard:

- Step 1: Set crawler properties
- Step 2: Choose data sources and classifiers
- Step 3: **Configure security settings** (selected)
- Step 4: Set output and scheduling
- Step 5: Review and create

**IAM role** [Info](#)

Existing IAM role: AWSGlueServiceRole-lab-3fc58d70

Create new IAM role  Update chosen IAM role  View [Edit](#)

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.

Buttons: Cancel, Previous, Next

System status bar at the bottom:

- CloudShell Feedback
- 36°F Sunny
- Search
- © 2025, Amazon Web Services, Inc. or its affiliates.
- Privacy Terms Cookie preferences
- ENG IN 09:35 29-11-2025

Screenshot of the AWS Glue console showing the 'Add crawler' wizard at Step 4: Set output and scheduling.

The left sidebar shows the AWS Glue navigation menu with 'Set output and scheduling' selected. The main panel displays the 'Set output and scheduling' configuration screen.

**Output configuration** (Info): Target database is set to 'firehose-inbound'. Buttons for 'Clear selection' and 'Add database' are available.

**Table name prefix - optional**: A text input field with placeholder 'Type a prefix added to table names'.

**Maximum table threshold - optional**: A text input field with placeholder 'Type a number greater than 0'.

**Advanced options**: A section with a 'Frequency' dropdown set to 'On demand'.

**Crawler schedule**: A section for defining a cron-based schedule. It includes a 'Syntax' link and a 'Frequency' dropdown set to 'Custom (cron expression)'. Below it is a cron expression builder with fields for Minutes, Hours, Day of month, Month, Day of week, and Year, resulting in 'Every 5 minutes'.

The bottom right shows the AWS navigation bar with account ID, region, and other links.

Screenshot of the AWS Glue console showing the 'Add crawler' wizard at Step 5: Review and create.

The left sidebar shows the AWS Glue navigation menu with 'Review and create' selected. The main panel displays the 'Review and create' configuration screen.

**Table name prefix - optional**: A text input field with placeholder 'Type a prefix added to table names'.

**Maximum table threshold - optional**: A text input field with placeholder 'Type a number greater than 0'.

**Advanced options**: A section with a 'Frequency' dropdown set to 'Custom (cron expression)'.

**Crawler schedule**: A section for defining a cron-based schedule. It includes a 'Syntax' link and a 'Frequency' dropdown set to 'Custom (cron expression)'. Below it is a cron expression builder with fields for Minutes, Hours, Day of month, Month, Day of week, and Year, resulting in 'Every 5 minutes'.

**Review and create**: A summary section showing the crawler details: Name 'crawler-test', Target database 'firehose-inbound', and Crawler schedule 'Every 5 minutes'. It also lists 'Tables' and 'Review and create'.

**Next**: A large orange 'Next' button at the bottom right.

The bottom right shows the AWS navigation bar with account ID, region, and other links.

Networking in AWS Cloud | Jaswanth Matta | LinkedIn | AWS Skill Builder | Cloud Quest | Crawlers - AWS Glue Conso | Terraform AWS Engineer | + | - | ○ | :| us-east-1.console.aws.amazon.com/glue/home?region=us-east-1#/v2/data-catalog/crawlers/add

aws Search [Alt+S] United States (N, Virginia) Account ID: 7349-7944-9836 AWSLabsUser-wDy3nDpvPjQPxr24FLGM98/98ea6b...

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**

Step 1  
Step 2  
Step 3  
Step 4  
Step 5  
**Review and create**

**Review and create**

**Step 1: Set crawler properties**

**Set crawler properties**

Name	Description	Tags
crawl_processed_data	-	-

**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://databucket-3fc58d70//processed...	Recrawl all

**Step 3: Configure security settings**

**Configure security settings**

IAM role	Security configuration	Lake Formation configuration
AWSGlueServiceRole-lab-3fc58d70	-	-

**Step 4: Set output and scheduling**

What's New L

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences ENG IN 09:38 29-11-2025

Networking in AWS Cloud | Jaswanth Matta | LinkedIn | AWS Skill Builder | Cloud Quest | Crawlers - AWS Glue Conso | Terraform AWS Engineer | + | - | ○ | :| us-east-1.console.aws.amazon.com/glue/home?region=us-east-1#/v2/data-catalog/crawlers/view/crawl\_processed\_data

aws Search [Alt+S] United States (N, Virginia) Account ID: 7349-7944-9836 AWSLabsUser-wDy3nDpvPjQPxr24FLGM98/98ea6b...

AWS Glue > Crawlers > crawl\_processed\_data

**crawl\_processed\_data**

**One crawler successfully created**  
The following crawler is now created: "crawl\_processed\_data"

Last updated (UTC) November 29, 2025 at 14:58:14

**Crawler properties**

Name	IAM role	Database	State
crawl_processed_data	AWSGlueServiceRole-lab-3fc58d70	firehose-ingestion	READY
Description	Security configuration	Lake Formation configuration	Table prefix
-	-	-	-

**Advanced settings**

**Crawler runs (0)**  
The list of crawler runs for this crawler.

Stop run View CloudWatch logs View run details Filter data Filter by a date and time range

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences ENG IN 09:38 29-11-2025

The screenshot shows the AWS EC2 Instances page. The left sidebar is collapsed, showing the EC2 navigation bar. The main area displays a table of instances. One instance is selected, labeled "Web UI" with the ID "i-013b6bc64f43c0a5b". The instance is shown as "Running" in the "Status check" column. Other columns include Instance ID, Instance state, Instance type, Status check, Alarm status, Availability Zone, and Public IPv4. Below the table, a detailed view for the selected instance is shown, with tabs for Details, Status and alarms, Monitoring, Security, Networking, Storage, and Tags. The "Details" tab is active, showing the Instance ID as "i-013b6bc64f43c0a5b", Public IPv4 address as "3.230.3.203 | open address", and Instance state as "Running".

The screenshot shows a web browser window with the URL "3.230.3.203/firehose". The page title is "Clickstream Publishing System". A message at the top reads: "Press the button below to send a sample set of data to your delivery stream. This will be used by your Amazon Kinesis Data Analytics application for schema discovery." Below this is a large blue button labeled "Start". At the bottom of the browser window, the Windows taskbar is visible, showing the date and time as "29-11-2025" and "09:42".

Clickstream Publishing System

Press the button below to send a sample set of data to your delivery stream. This will be used by your Amazon Kinesis Data Analytics application for schema discovery.

**Test data set**

```
{ "Encrypted": false, "RecordId": "XRUwMFpRpk+EiyBxPvUil0UGGOrInY1f16kQJ0sy0/9tE9DQGdwQ7oPMsb3HvuMzehIGBGd48m70n0jpQjXo5WvrppjTVxBSR91jkPQ", "ResponseMetadata": { "HTTPHeaders": { "Content-Length": "257", "Content-Type": "application/x-amz-json-1.1", "Date": "Sat, 29 Nov 2025 14:44:51 GMT", "x-amz-id-2": "VHUUUT1RZ5j6nBA9BcaqLs8o+4Y/cfEvPAR70l3qCaf43TMY9pSgDhGgsBidtb8bsq+j7Mgpnis8czaFlN3jUfqduypfew", "x-amzn-requestid": "c94e7fcc-199e-c280-a8eb-76ef03f2aefb" }, "HTTPStatusCode": 200, "RequestId": "c94e7fcc-199e-c280-a8eb-76ef03f2aefb", "RequestType": "POST" } }
```

**Data set 1**

Start

37°F Sunny

ENG IN 09:44 29-11-2025

Clickstream Publishing System

Press the button below to send a sample set of data to your delivery stream. This will be used by your Amazon Kinesis Data Analytics application for schema discovery.

**Test data set**

```
{ "Encrypted": false, "RecordId": "zMB90gYKfrfxPdl+RUY728RvZ2tV1NjJeCz5yEkzoJvP28rB87lguiuc6pyFFUUrEm4qTE1i0h1ITmSp+ztdjXm//ePrngitJ1DFBFH", "ResponseMetadata": { "HTTPHeaders": { "Content-Length": "257", "Content-Type": "application/x-amz-json-1.1", "Date": "Sat, 29 Nov 2025 14:45:14 GMT", "x-amz-id-2": "sz7QJuvtL1ekGrolkxtWt5UsaplLoVhM5Gen9CSwggkFnyFhdHDJ/3Mj1QiTyrgChQcjMkYJaraR820ukS3d1uEakZ+E47UkT", "x-amzn-requestid": "c7ce5b83-3b59-b424-a66b-52b92135d85f" }, "HTTPStatusCode": 200, "RequestId": "c7ce5b83-3b59-b424-a66b-52b92135d85f", "RequestType": "POST" } }
```

**Data set 1**

Start

**Data set 2**

Start

37°F Sunny

ENG IN 09:45 29-11-2025

Introducing the AWS Serverless generative AI webinar series  
Join AWS Serverless experts for a hands-on workshops on Agentic AI, Intelligent Document Processing, and building real-time APIs for chatbots.

See more details

Last fetched 0 seconds ago

Actions Create function

Functions (4)

Function name	Description	Package type	Runtime	Last modified
gbl_lab_monitoring	-	Zip	Python 3.11	1 hour ago
AnalyticsDestinationFunction	-	Zip	Python 3.11	1 hour ago
LabStack-98ea6b70-bd2e-46-GblLabMonitoringgblproto-ePnvT6xj7z41	-	Zip	Python 3.11	1 hour ago
DataProcessingFunction	-	Zip	Python 3.11	1 hour ago

Code source Info

Open in Visual Studio Code Upload from

EXPLORER

DATAPROCESSINGFUNCTION

process\_data.py

Deploy (Ctrl+Shift+U) Test (Ctrl+Shift+I)

TEST EVENTS (NONE SELECTED)

```
process_data.py
process_data.py
def handler(event, context):
    # Clickstream data labels.
    labels = ['prev', 'curr', 'type', 'n']
    # Iterate list of input records.
    for record in event.get('records'):
        records = []
        # Get data from Kinesis record.
        data = base64.b64decode(record.get('data')).decode('utf-8')
        logger.info(data)
        # Split into lines.
        lines = data.split('\n')
        for line in lines:
            # Skip any empty lines.
            if line == "":
                continue
            # Process line.
```

The screenshot shows the AWS DynamoDB console with the 'Tables' section selected. A specific table named 'OutputTable' is being viewed. The left sidebar includes links for Dashboard, Tables, Explore items, PartiQL editor, Backups, Exports to S3, Imports from S3, Integrations, Reserved capacity, and Settings. Under 'DAX', there are links for Clusters, Subnet groups, Parameter groups, and Events.

**Table capacity**

Read capacity auto scaling	Write capacity auto scaling
Off	Off
Provisioned read capacity units	Provisioned write capacity units
5	5

**Estimated read/write capacity cost**

**Auto scaling activities (0)**

Last updated November 29, 2025, 09:49 (UTC-5:00) [Edit](#)

No auto scaling activities found

**Warm throughput** [Info](#)

Last updated November 29, 2025, 09:49 (UTC-5:00) [Edit](#)

Prepare your table for planned peak events, without impacting your application performance or availability. Learn more about Amazon DynamoDB pricing [Learn more](#)

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences 37°F Sunny ENG IN 09:49 29-11-2025

The screenshot shows the AWS S3 console with the 'Buckets' section selected. A bucket named 'databucket-3fc58d70' is being viewed. The left sidebar includes links for General purpose buckets (Directory buckets, Table buckets, Vector buckets), Access management and security (Access Points, Access Points for FSx, Access Grants, IAM Access Analyzer), and Storage management and insights (Storage Lens, Batch Operations). There are also links for Account and organization settings and AWS Marketplace for S3.

**Objects (1/1)**

Name	Type	Last modified	Size	Storage class
ClickStreamData-1-2025-11-29-14-43-58-937358d0-b4c6-45eb-b9aa-b8c9b08fcee	-	November 29, 2025, 09:48:59 (UTC-05:00)	9.5 KB	Standard

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences 37°F Sunny ENG IN 09:52 29-11-2025

File Edit Selection View Go Run Terminal Help ↶ ↷

Restricted Mode is intended for safe code browsing. Trust this window to enable all features. Manage Learn More

ClickStreamData-1-2025-11-29-14-43-58-937358d0-b4c6-45eb-b9aa-b8c9b08fcecc X vehicles\_function.py •

```
C: > Users > jaswa > Downloads > ClickStreamData-1-2025-11-29-14-43-58-937358d0-b4c6-45eb-b9aa-b8c9b08fcecc
1 other-search,"Doherty_(Floyd County, Texas)",external,17
2 other-search,schubring,external,14
3 other-search,bluejeans_Network,external,309
4 other-empty,BlueJeans_Network,external,175
5 other-empty,Kastell_Dunaszekcső,external,12
6 Shunt_(Medizin),Blalock-Taussig Anastomose,link,11
7 Fallot-Tetralogie,Blalock-Taussig Anastomose,link,17
8 other-search,Blalock-Taussig Anastomose,external,200
9 Rashkind-Manter,Blalock-Taussig Anastomose,link,13
10 other-empty,Blalock-Taussig Anastomose,external,83
11 Elisabeth_von_Matt,Francis_Triesnecker,link,20
12 other-search,Mussidan,external,20
13 Madeleine_Delibré,Mussidan,link,16
14 Exorcist_The_Exorcist_(Fernsehserie),link,13
15 other-internal,The_Exorcist_(Fernsehserie),external,20
16 other-search,The_Exorcist_(Fernsehserie),external,655
17 Brienne_Howey_The_Exorcist_(Fernsehserie),link,45
18 Brianna_Hildebrand_The_Exorcist_(Fernsehserie),link,20
19 Geena_Davis_The_Exorcist_(Fernsehserie),link,32
20 Der_Exorcist_(Roman)_The_Exorcist_(Fernsehserie),link,14
21 Ben_Daniels_The_Exorcist_(Fernsehserie),link,37
22 John.Cho_The_Exorcist_(Fernsehserie),link,21
23 other-empty,The_Exorcist_(Fernsehserie),external,173
24 other-external,Kronstädter_Schriftstellerprozess,external,10
25 other-search,Kronstädter_Schriftstellerprozess,external,17
26 Eginald_Schlattner_Kronstädter_Schriftstellerprozess,link,16
27 Hans_Bergel_Kronstädter_Schriftstellerprozess,link,14
28 other-empty,Kronstädter_Schriftstellerprozess,external,12
29 Newness,Danny_Huston,link,27
30 Gramma,Danny_Huston,link,15
31 All_I_See_Is_You,Danny_Huston,link,22
32 Number_23,Danny_Huston,link,20
33 John_Huston,Danny_Huston,link,49
34 American_Horror_Story_Freak_Show,Danny_Huston,link,130
35 Birth_(Film),Danny_Huston,link,13
36 American_Horror_Story_Coven,Danny_Huston,link,64
```

Ln 1, Col 1 Spaces: 4 UTF-8 LF { } Plain Text ⚡ Prettier

Restricted Mode ⚡ Launchpad ⚡ 0 ⚡ 0 ⚡ 100 ⚡ 37°F Sunny 09:53 29-11-2025

Networking in AWS C | LinkedIn Jaswanth Matta | LinkedIn AWS Skill Builder Cloud Quest Crawlers - AWS Glue | Terraform AWS Engine Clickstream Publishing | + -

aws Search [Alt+S] Account ID: 7349-7944-9836 AWSLabsUser-wOy3nDpPjQPxr24FL6M98/98ea5b...

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

CloudShell Feedback © 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences 37°F Sunny 09:56 29-11-2025

The screenshot shows the AWS Glue Crawler configuration page. The crawler is named 'crawl\_processed\_data'. It uses the 'AWSGlueServiceRole-lab-3fc58d70' IAM role. The target is 'firehose-inbound' and it is configured with 'Lake Formation configuration'. The status is 'RUNNING'. The 'Advanced settings' section is collapsed. Below the main configuration, there is a table titled 'Crawler runs (4)' showing four completed runs on November 29, 2025, at various times between 14:00 and 14:38 UTC. The table includes columns for Start time (UTC), End time (UTC), Current/last duration, Status, DPU hours, and Table changes.

Start time (UTC)	End time (UTC)	Current/last duration	Status	DPU hours	Table changes
November 29, 2025 at 14...	-	14 s	Running	-	-
November 29, 2025 at 14...	November 29, 2025 at 14...	38 s	Completed	0.107	-
November 29, 2025 at 14...	November 29, 2025 at 14...	01 min	Completed	0.039	-
November 29, 2025 at 14...	November 29, 2025 at 14...	37 s	Completed	0.109	-

The screenshot shows the Amazon Athena landing page. The page features a large 'Analytics' header and the 'Amazon Athena' logo with the tagline 'Start querying data instantly.' Below this, a paragraph explains that Athena is an interactive query service for analyzing data in Amazon S3 and other federated data sources using standard SQL. To the right, there is a 'Get started' section with two options: 'Query your data in Amazon SageMaker Unified Studio' (selected) and 'Query your data in Athena console'. A button labeled 'Open Amazon SageMaker Unified Studio' is present. Another section titled 'Pricing' provides details about costs: \$5.00 per TB for SQL queries on TB scanned, \$0.30 per DPU hour for SQL queries on Provisioned Capacity, and \$0.35 per DPU hour for Apache Spark calculations. The page also indicates the region is US East (N. Virginia). The bottom navigation bar includes links for CloudShell, Feedback, and the AWS logo.

The screenshot shows the Amazon Athena Query Editor interface. A prominent message at the top states: "Athena now supports typeahead code suggestions to speed up SQL query development. Typehead suggestions are turned on by default. You can change this setting in query editor preferences." On the left, the "Data" sidebar shows a "Data source" set to "AwsDataCatalog" and a "Database" set to "firehose-ingestion". Under "Tables and views", there is one table named "processed\_data" which is partitioned. The main area is titled "Query 1" and contains a single row of SQL code: "SQL Ln 1, Col 1". Below the SQL area are buttons for "Run", "Explain", "Cancel", "Clear", and "Create". To the right of these buttons is a checkbox for "Reuse query results up to 60 minutes ago". At the bottom, tabs for "Query results" and "Query stats" are visible, along with "Copy" and "Download results CSV" buttons.

This screenshot shows the "Manage query result location and encryption" dialog box overlaid on the main Amazon Athena Query Editor interface. The dialog has two main sections: "Location of query result - optional" and "Expected bucket owner - optional". In the "Location of query result" section, a text input field shows "s3://databucket-55610880" with "View" and "Browse S3" buttons. In the "Expected bucket owner" section, there is a text input field for "Enter AWS account ID". Below these sections are two checkboxes: "Assign bucket owner full control over query results" and "Encrypt query results". At the bottom of the dialog are "Cancel" and "Save" buttons. The background of the main interface shows the "Query settings" tab selected, with various configuration options like "Query result encryption" and "Execution controls". The status bar at the bottom indicates "CloudShell Feedback" and the date "02-12-2025".

The screenshot shows the AWS Cloud console with multiple tabs open. The active tab is 'Query editor | Athena | us-east-1'. The interface includes a sidebar for 'Database' (selected) and 'Tables and views'. Under 'Tables', there is one entry: 'processed\_data' (Partitioned). The main area contains an SQL editor with the query: `SELECT * FROM "firehose-inbound"."processed_data" limit 10;`. Below the editor is a 'Query results' section with a 'Results' tab showing the message 'No results'.

This screenshot shows the same setup as the first one, but the history bar at the top indicates a previous query was run. The SQL editor now displays the completed query: `1 | SELECT * FROM "firehose-inbound"."processed_data" limit 10;`. The 'Query results' section shows the status 'Completed' and provides performance metrics: Time in queue: 132 ms, Run time: 1.465 sec, Data scanned: 12.71 KB.

The screenshot shows the AWS Cloud Console with multiple tabs open. The active tab is 'Query editor | Athena | us-east-1'. The URL is <https://us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor/history/c005b584-d10c-43c4-b347-1c08f0d8eccf>. The account ID is 7349-7944-9836.

The interface includes a search bar, navigation buttons, and a toolbar with icons for copy, download, and refresh. The left sidebar shows the schema with partitions and views. The main area displays the results of a query with 10 rows. The columns are labeled #, col0, col1, col2, col3, partition\_0, partition\_1, partition\_2, and partition\_3. The data includes various search terms and their corresponding details like source type, count, and timestamp.

#	col0	col1	col2	col3	partition_0	partition_1	partition_2	partition_3
1	other-search	"Dougherty_(Floyd_County	_Texas)"		2025	12	02	20
2	other-search	Schubring	external	14	2025	12	02	20
3	other-search	BlueJeans_Network	external	309	2025	12	02	20
4	other-empty	BlueJeans_Network	external	75	2025	12	02	20
5	other-empty	Kastell_Dunaszekcső	external	12	2025	12	02	20
6	Shunt_(Medizin)	Blalock-Taussig-Anastomose	link	11	2025	12	02	20
7	Fallot-Tetralogie	Blalock-Taussig-Anastomose	link	17	2025	12	02	20
8	other-search	Blalock-Taussig-Anastomose	external	200	2025	12	02	20
9	Rashkind-Manöver	Blalock-Taussig-Anastomose	link	13	2025	12	02	20
10	other-empty	Blalock-Taussig-Anastomose	external	83	2025	12	02	20

The screenshot shows the AWS Cloud Console with multiple tabs open. The active tab is 'Query editor | Athena | us-east-1'. The URL is <https://us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor/history/2671b2b3-8993-4a7b-ad5e-4cb0e6c18c60>. The account ID is 7349-7944-9836.

A message box at the top left announces: "Athena now supports typeahead code suggestions to speed up SQL query development. Typeahead suggestions are turned on by default. You can change this setting in query editor preferences." There is a "Edit preferences" button next to it.

The interface includes a search bar, navigation buttons, and a toolbar with icons for run, explain, cancel, clear, and create. The left sidebar shows the data source (AwsDataCatalog), database (firehose-ingestion), and tables (processed\_data). The main area shows two queries: "Query 1" and "Query 2". Query 1 contains the following SQL:

```
1 select col2 as source, count(*) as count from "firehose-ingestion"."processed_data"
2 where col2 in ('external', 'link')
3 group by "col2"
```

The status bar at the bottom shows the date and time as 02-12-2025 16:01.

Screenshot of the AWS Athena Query Editor interface.

**Tables (1)**

- processed\_data (Partitioned)
  - col0 string
  - col1 string
  - col2 string
  - col3 bigint
  - partition\_0 string (Partitioned)
  - partition\_1 string (Partitioned)
  - partition\_2 string (Partitioned)
  - partition\_3 string (Partitioned)

**Views (0)**

**SQL** Ln 3, Col 16

Run again Explain Cancel Clear Create Reuse query results up to 60 minutes ago

**Query results** | Query stats

Completed Time in queue: 112 ms Run time: 1.035 sec Data scanned: 18.86 KB

Copy Download results CSV

**Results (2)**

#	source	count
1	external	172
2	link	268

Screenshot of the CloudShell interface.

CloudShell Feedback 36°F Rain and snow 02-12-2025

Screenshot of the AWS Athena Query Editor interface showing three JSON documents from a firehose endpoint.

**Data set 3**

```
{
  "x-amz-id-2": "EMKc73Dl3gI8NFeao8Cj1QYzfbqDXXJa0vcyX2TrCYe5mdeiwfSousCDwYZhUbbufQsFT1saH9c7SOZBQKmV3AztBVBGt9",
  "x-amzn-requestid": "d19f36f9-d3c0-f469-b036-60a9a7736d0c"
},  

  "HTTPstatusCode": 200,  

  "RequestId": "d19f36f9-d3c0-f469-b036-60a9a7736d0c",
  "RetryAttempts": 0
}
```

**Data set 4**

```
{
  "Encrypted": false,
  "RecordId": "UPA3gH8PPYOx1wKCPJkdwSHU1M0m01Z/NiPSX72EYf3k1a1jkfttrP0I300VKMczQNoSwyj9wsCEXuE6Ta90Lf5FujfOxIrsh6qTM9bfZf",
  "ResponseMetadata": {
    "HTTPHeaders": {
      "content-length": "257",
      "content-type": "application/x-amz-json-1.1",
      "date": "Tue, 02 Dec 2025 21:03:12 GMT",
      "x-amz-id-2": "WAZqKtA2Cyd4vDyEztomad87avAub1sQU/8ZHaIN+m75X8rauwmXbKiJGoIHyxmDap0fPADx2LvoWfq4jityIraCuAMKV",
      "x-amzn-requestid": "f57900c2-6a09-b7aa-94d0-5692f7d39960"
    },
    "HTTPstatusCode": 200,
    "RequestId": "f57900c2-6a09-b7aa-94d0-5692f7d39960",
    ...
  }
}
```

**Data set 5**

```
{
  "Encrypted": false,
  "RecordId": "AXe6r1BDpbQis1qlqBf1T90yKEbbyoU40X/XYS0GrLqR+KlwZnZpj1s1i9007u/h89NU1K/xVcgOJTFj7fnCRJfimXjE903Qy+E+fZWC000f",
  "ResponseMetadata": {
    "HTTPHeaders": {
      "content-length": "257",
      "content-type": "application/x-amz-json-1.1",
      "date": "Tue, 02 Dec 2025 21:03:12 GMT",
      ...
    }
  }
}
```

The screenshot shows the AWS Athena Query Editor interface. On the left, there's a sidebar with 'Tables (1)' containing 'processed\_data' (Partitioned) with columns col0, col1, col2, col3, partition\_0, partition\_1, partition\_2, and partition\_3. Below it is 'Views (0)'. The main area has a SQL editor with the query 'SELECT \* FROM processed\_data'. Below the editor is a 'Query results' section showing a table with two rows: 'external' and 'link'. The table has columns '#', 'source', and 'count'. The count column shows values 172 and 268 respectively.



The screenshot shows the 'Edit transform and convert records' page for Amazon Data Firehose. It includes sections for 'AWS Lambda function' (set to arnaws:lambda:us-east-1:734979449836:function:AnalyticsDestinationFunction:\$LATEST), 'Version or alias' (\$LATEST), and 'Create function'. Other settings include 'Buffer size' (1 MB), 'Buffer interval' (60 seconds), and 'Convert record format' (unchecked). The page also includes a note about record format conversion using AWS Glue.



Screenshot of the AWS DynamoDB console showing the 'Explore items' interface for the 'OutputTable' table.

The left sidebar shows the navigation menu:

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

**DAX**

- Clusters
- Subnet groups
- Parameter groups
- Events

The main area displays the 'OutputTable' configuration:

- Tables (1)**: Shows 1 table named 'OutputTable'.
- Scan or query items**: Set to 'Scan'.
- Select a table or index**: Table - OutputTable, Select attribute projection: All attributes.
- Filters - optional**: No filters applied.
- Run** button.

The results table shows 50 items returned:

timestamp (String)	value
2025-12-02 21:10:47...	Tom_Reece Georg_VI_(Vereinigtes_Königreich) link 112
2025-12-02 21:11:45...	Sandringham_House Georg_VI_(Vereinigtes_Königreich) link 147
2025-12-02 21:11:43...	other-search Georg_VI_(Vereinigtes_Königreich) external 27645
2025-12-02 21:11:45...	Schutzenraum Sonnenbergtunnel link 13
2025-12-02 21:11:45...	other-internal Sonnenbergtunnel external 17
2025-12-02 21:11:46...	Großer_Preis_von_Großbritannien_1950 Georg_VI_(Vereinigtes_Königreich) link 12
2025-12-02 21:10:45...	Karl_Obermayr Kehraus_(Film) link 30
2025-12-02 21:10:48...	Admiral_Graf_Spee Georg_VI_(Vereinigtes_Königreich) link 51
2025-12-02 21:10:54...	other-search Jonas_Myrin external 34
2025-12-02 21:11:44...	Psiloritis Kamares-Höhle link 13
2025-12-02 21:10:48...	other-empty Takeshi_Obara external 36

Screenshot of the AWS DynamoDB console showing the 'Explore items' interface for the 'OutputTable' table.

The left sidebar shows the navigation menu:

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

**DAX**

- Clusters
- Subnet groups
- Parameter groups
- Events

The main area displays the 'OutputTable' configuration:

- Tables (1)**: Shows 1 table named 'OutputTable'.
- Scan or query items**: Set to 'Scan'.
- Select a table or index**: Table - OutputTable, Select attribute projection: All attributes.
- Filters - optional**: No filters applied.
- Run** button.

The results table shows 50 items returned:

timestamp (String)	value
2025-12-02 21:10:47...	Tom_Reece Georg_VI_(Vereinigtes_Königreich) link 112
2025-12-02 21:11:45...	Sandringham_House Georg_VI_(Vereinigtes_Königreich) link 147
2025-12-02 21:11:43...	other-search Georg_VI_(Vereinigtes_Königreich) external 27645
2025-12-02 21:11:45...	Schutzenraum Sonnenbergtunnel link 13
2025-12-02 21:11:45...	other-internal Sonnenbergtunnel external 17
2025-12-02 21:11:46...	Großer_Preis_von_Großbritannien_1950 Georg_VI_(Vereinigtes_Königreich) link 12
2025-12-02 21:10:45...	Karl_Obermayr Kehraus_(Film) link 30
2025-12-02 21:10:48...	Admiral_Graf_Spee Georg_VI_(Vereinigtes_Königreich) link 51
2025-12-02 21:10:54...	other-search Jonas_Myrin external 34
2025-12-02 21:11:44...	Psiloritis Kamares-Höhle link 13
2025-12-02 21:10:48...	other-empty Takeshi_Obara external 36

The screenshot shows the AWS DynamoDB console with the 'Explore items' feature open for the 'OutputTable'. The left sidebar has 'Explore items' selected under the 'Tables' section. The main area displays a table titled 'Table: OutputTable - Items returned (50)' with a scan started on December 02, 2025, at 16:17:15. The table lists 50 items, each with a timestamp and a value. The values are links to external sources. The interface includes standard AWS navigation and search tools.

	timestamp (String)	value
1	2025-12-02 21:10:45...	other-search Kehraus_(Film) external 279
2	2025-12-02 21:10:46...	Dieter_Hildebrandt Kehraus_(Film) link 14
3	2025-12-02 21:11:43...	Eid_Gekreuzte_Finger link 29
4	2025-12-02 21:11:44...	other-empty Gottlob_Adolf_Ernst_von_Nostitz_und_Jänkendorf external 13
5	2025-12-02 21:11:45...	other-external Tikun_Olam external 11
6	2025-12-02 21:14:43...	Big_Hit_Music_RM_(Rapper) link 23
7	2025-12-02 21:10:48...	other-search Takeshi_Obata external 226
8	2025-12-02 21:10:48...	Angers Liste_der_Monuments_historiques_in_Angers link 21
9	2025-12-02 21:16:44...	Harry_Duke_of_Sussex_Georg_VI_(Vereinigtes_Königreich) link 57
10	2025-12-02 21:10:48...	Platinum_End Takeshi_Obata link 20
11	2025-12-02 21:10:50...	MantegnaAndrea_Mantegna link 71
12	2025-12-02 21:11:45...	Liste_der_Herrscher_namens_Georg_Georg_VI_(Vereinigtes_Königreich) link 10