

# Real-Time Data Pipeline with Kinesis Data Streams, Lambda, Firehose, S3, and Redshift for Data Validation and Ingestion





# Objective

- AWS Kinesis Data Streams: Capture live, real-time data from various sources.
- AWS Kinesis Data Firehose: Stream the captured data into Amazon S3 and use AWS Lambda for real-time validation and processing.
- AWS Lambda: Validate and process the incoming data before it is inserted into Amazon Redshift.
- Amazon S3: Temporarily store the data before it undergoes validation and transformation through Lambda.
- Amazon Redshift: Load and store validated and processed data for fast querying and analytics.



# Tools Used

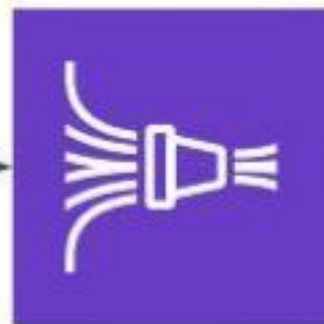
- AWS Kinesis Data Streams: Captures live, real-time data from various sources.
- AWS Kinesis Data Firehose: Streams the validated data to Amazon S3 and integrates with Lambda for processing.
- AWS Lambda: Validates and processes the incoming data before it is inserted into Redshift.
- Amazon S3: Serves as a temporary storage location for data before it is validated and ingested into Redshift.
- Amazon Redshift: A data warehouse where processed data is stored and made available for analytics.
- AWS SDK (boto3): Used within Lambda to interact with AWS services, including inserting data into Redshift.



AWS Cloud



Amazon Kinesis  
Data Streams



Amazon Kinesis  
Data Firehose



Lambda  
Function



Amazon Redshift





## Amazon Kinesis



Dashboard

### Data streams

Amazon Data Firehose [New](#)

Managed Apache Flink [New](#)

### ▼ Resources

CloudFormation templates

AWS Glue Schema Registry

## Data streams (1) [Info](#)

Process data in real time

Create a Firehose stream

Actions ▼

Create data stream

Find data streams

< 1 >



<input type="checkbox"/>	Name ▲	Status ▼	Capacity mode ▼	Provisioned shards ▼	Sharing policy ▼	Data retention period ▼	Encryption ▼
<input type="checkbox"/>	<a href="#">app_clickstream_events</a>	✓ Active	On-demand	-	No	1 day	Disabled



## Functions (1)

Last fetched 0 seconds ago



Actions ▼

Create function

🔍 Filter by attributes or search by keyword



1



Function name



Description



Package type



Runtime



Last modified



[clickstream\\_events\\_enrichment](#)

-

Zip

Python 3.10

[2 hours ago](#)





# clickstream\_events\_enrichment

Throttle

📄 Copy ARN

Actions ▼

## ▼ Function overview [Info](#)

Diagram

Template



clickstream\_events\_e  
nrichment



Layers

(0)

+ Add trigger

+ Add destination

Export to Infrastructure Composer

Download ▼

### Description

-

### Last modified

2 hours ago

### Function ARN

📄 arn:aws:lambda:us-east-1:418272783111:function:c  
lickstream\_events\_enrichment

### Function URL [Info](#)

-

[Code](#)

[Test](#)

[Monitor](#)

[Configuration](#)

[Aliases](#)

[Versions](#)



## Code source Info

Upload from



clickstream\_events\_enrichment



lambda\_function.py

lambda\_function.py

```

1  import json
2  import base64
3  import logging
4
5  # Configure logging
6  logging.basicConfig(level=logging.INFO)
7  logger = logging.getLogger()
8
9  # Helper function to determine if the browser indicates a mobile device
10 def is_mobile_browser(browser):
11     # List of mobile indicators based on typical mobile browser strings
12     mobile_indicators = ['Mobile', 'iPhone', 'Android', 'iPad']
13     return any(indicator in browser for indicator in mobile_indicators)
14
15 # Helper function to categorize traffic sources
16 def categorize_traffic_source(source):
17     if 'Google' in source or 'Bing' in source:
18         return 'Search Engine'
19     elif 'Facebook' in source or 'YouTube' in source:

```

```

17 def is_mobile_browser(browser):
18     # List of mobile indicators based on typical mobile browser strings
19     mobile_indicators = ['Mobile', 'iPhone', 'Android', 'iPad']
20     return any(indicator in browser for indicator in mobile_indicators)
21
22 # Helper function to categorize traffic sources
23 def categorize_traffic_source(source):
24     if 'Google' in source or 'Bing' in source:
25         return 'Search Engine'
26     elif 'Facebook' in source or 'YouTube' in source:

```

EXPLORER

CLICKSTREAM\_EVENTS\_ENRICHMENT

lambda\_function.py

DEPLOY

Deploy (Ctrl+Shift+U)

Test (Ctrl+Shift+I)

TEST EVENTS [NONE SELECTED]

+ Create new test event





lambda.py > categorize\_traffic\_source

```
1  import json
2  import base64
3  import logging
4
5  # Configure logging
6  logging.basicConfig(level=logging.INFO)
7  logger = logging.getLogger()
8
9  # Helper function to determine if the browser indicates a mobile device
10 def is_mobile_browser(browser):
11     # List of mobile indicators based on typical mobile browser strings
12     mobile_indicators = ['Mobile', 'iPhone', 'Android', 'iPad']
13     return any(indicator in browser for indicator in mobile_indicators)
14
15 # Helper function to categorize traffic sources
16 def categorize_traffic_source(source):
17     if 'Google' in source or 'Bing' in source:
18         return 'Search Engine'
19     elif 'Facebook' in source or 'YouTube' in source:
20         return 'Social Media'
21     elif 'Email' in source:
22         return 'Email'
23     else:
24         return 'Direct'
25
26 def lambda_handler(event, context):
27     output = []
28
29     # Process each record in the event batch
30     for record in event['records']:
31         try:
```

```
25
26 def lambda_handler(event, context):
27     output = []
28
29     # Process each record in the event batch
30     for record in event['records']:
31         try:
32             # Decode the incoming data from base64 and convert from JSON
33             payload = json.loads(base64.b64decode(record['data']))
34
35             # Extract fields to be included in the CSV
36             id = payload.get('id', '')
37             user_id = payload.get('user_id', '')
38             session_id = payload.get('session_id', '')
39             ip_address = payload.get('ip_address', '')
40             city = payload.get('city', '')
41             state = payload.get('state', '')
42             browser = payload.get('browser', '')
43             traffic_source = payload.get('traffic_source', '')
44             uri = payload.get('uri', '')
45             event_type = payload.get('event_type', '')
46
47             # Derive is_mobile and traffic_source_category
48             is_mobile = 'Yes' if is_mobile_browser(browser) else 'No'
49             traffic_source_category = categorize_traffic_source(traffic_source)
50
51             # Quote each field to ensure proper CSV formatting
52             data = f'"{id}","{user_id}","{session_id}","{ip_address}","{city}","{state}","{browser}","{is_mobile}","{traffic_source_category}","{uri}","{event_type}"'
53
54             # Encode the CSV string to base64
55             encoded_csv = base64.b64encode(data.encode('utf-8')).decode('utf-8')
```



```

26 def lambda_handler(event, context):
52     data = f'"{id}"","{user_id}"","{session_id}"","{ip_address}"","{city}"","{state}"","{browser}"","{is_mobile}"","{traffic"
53
54     # Encode the CSV string to base64
55     encoded_csv = base64.b64encode(data.encode('utf-8')).decode('utf-8')
56
57     # Prepare the output record with unique recordId and transformed data
58     output_record = {
59         'recordId': record['recordId'],
60         'result': 'Ok',
61         'data': encoded_csv
62     }
63     output.append(output_record)
64
65     logger.info(f"Processed record ID: {record['recordId']}")
66
67 except Exception as e:
68     logger.error(f"Error processing record {record['recordId']}: {str(e)}")
69     # Return the original record in case of failure
70     output_record = {
71         'recordId': record['recordId'],
72         'result': 'ProcessingFailed',
73         'data': record['data'] # Return original data
74     }
75     output.append(output_record)
76
77 # Return the transformed records to Firehose
78 logger.info(f"Total records processed: {len(output)}")
79 return {'records': output}

```

## Redshift query editor v2

Create ▼

Load data



Filter resources



&gt; Serverless: default-workgr... ⓘ ⋮

+ Untitled 1 x

Untitled 2 x

Untitled 3 x

Untitled 4 x

Untitled 5 x

Untitled 6 x



Run



Limit 100



Explain



Isolated session ⓘ

Serverless: de... ▼

clickstream\_db ▼



Schedule



```
1 create database clickstream_db;
2
3 CREATE TABLE prod_schema.events (
4     id BIGINT,
5     user_id BIGINT,
6     session_id CHAR(36),
7     ip_address VARCHAR(50),
8     city VARCHAR(255),
9     state VARCHAR(255),
10    browser VARCHAR(50),
11    is_mobile VARCHAR(3),
12    traffic_source VARCHAR(50),
13    traffic_source_category VARCHAR(50),
14    uri VARCHAR(255),
15    event_type VARCHAR(50)
16 );
17
18 SELECT * from prod_schema.events;
```



Result 1

Export ▼



Chart



Row 3, Col 14, Chr





General purpose buckets | Directory buckets

General purpose buckets (5) Info All AWS Regions

Refresh, Copy ARN, Empty, Delete, Create bucket

Buckets are containers for data stored in S3.

Find buckets by name

	Name ▲	AWS Region ▼	IAM Access Analyzer	Creation date ▼
<input type="radio"/>	<a href="#">aws-glue-assets-418272783111-us-east-1</a>	US East (N. Virginia) us-east-1	<a href="#">View analyzer for us-east-1</a>	December 8, 2024, 13:06:32 (UTC+05:30)
<input type="radio"/>	<a href="#">aws-logs-418272783111-us-east-1</a>	US East (N. Virginia) us-east-1	<a href="#">View analyzer for us-east-1</a>	December 10, 2024, 08:46:40 (UTC+05:30)
<input type="radio"/>	<a href="#">gd-airflow</a>	US East (N. Virginia) us-east-1	<a href="#">View analyzer for us-east-1</a>	December 20, 2024, 11:10:44 (UTC+05:30)
<input type="radio"/>	<a href="#">gd-aws-de-labs</a>	US East (N. Virginia) us-east-1	<a href="#">View analyzer for us-east-1</a>	December 20, 2024, 22:48:51 (UTC+05:30)
<input type="radio"/>	<a href="#">testing65987</a>	US East (N. Virginia) us-east-1	<a href="#">View analyzer for us-east-1</a>	December 8, 2024, 13:01:11 (UTC+05:30)





## Amazon Data Firehose



Firehose streams

### ▼ Resources

What's new

Developer guide

API reference

## Firehose streams (1) [Info](#)



Delete

Create Firehose stream

You can create a Firehose stream to set up a source, destination, and optional transformation for your streaming data delivery.

Find Firehose streams

< 1 >

	Name ▲	Status ▼	Creation... ▼	Source ▼	Data tra... ▼	Destinat... ▼	Destinat... ▼
	<a href="#">KDS-RED-b...</a>	Active	January 05...	<a href="#">app_clickst...</a>	<a href="#">clickstream...</a>	Amazon Re...	<a href="#">default-w...</a>





## Amazon Data Firehose



Firehose streams

### ▼ Resources

What's new

Developer guide

API reference

## Transform records

Edit

Configure Amazon Data Firehose to transform your record data.



### AWS Lambda function timeout

The current timeout of the specified AWS Lambda function is 3 seconds. To reduce the risk of the AWS Lambda function timing out before data transformation is complete, increase the timeout to 1 minute or longer in the Advanced settings section of your AWS Lambda configuration. [Go to AWS Lambda configuration.](#)

Transform source records with AWS  
Lambda

#### Info

On

Buffer size  
0.2 MiB

Buffer interval  
20 seconds

Lambda function  
[clickstream\\_events\\_enrichment](#)

Lambda function version  
\$LATEST

Description  
-

Runtime  
python3.10

Timeout  
3 seconds

## Destination settings [Info](#)

Edit

Specify the destination settings for your Firehose stream.



## Amazon Data Firehose



Firehose streams

### ▼ Resources

What's new [↗](#)

Developer guide [↗](#)

API reference [↗](#)

## Destination settings [Info](#)

Edit

Specify the destination settings for your Firehose stream.

**i** Ensure that your Amazon Redshift Serverless workgroup is publicly accessible and allows inbound access from this Amazon Data Firehose IP address: **52.70.63.192/27**. For more information, see [VPC Access to an Amazon Redshift Serverless workgroup](#) [↗](#).  
If you specify Amazon Redshift as your Firehose stream destination, once your Firehose stream is created, you cannot update the specified Amazon Redshift destination type.

### Amazon Redshift destination

Serverless workgroup  
[default-workgroup](#) [↗](#)

Database  
clickstream\_db

Columns

-

Destination type  
Serverless workgroup

Table  
prod\_schema.events

User name  
admin

### Amazon Redshift Serverless COPY command [Info](#)

COPY command options  
CSV IGNOREHEADER 1

COPY command

COPY prod\_schema.events FROM 's3://prod-csv-data-1-ba-1-manifests-1-CREDENTIALS'





## Amazon Data Firehose



Firehose streams

### ▼ Resources

What's new

Developer guide

API reference

## Amazon Redshift Serverless COPY command [Info](#)

COPY command options

CSV IGNOREHEADER 1

COPY command

```
COPY prod_schema.events FROM 's3://gd-aws-de-labs/<manifest>' CREDENTIALS  
'aws_iam_role=arn:aws:iam::<aws-account-id>:role/<role-name>' MANIFEST CSV  
IGNOREHEADER 1 ;
```

Copy

Retry duration

3600 seconds

## Intermediate S3 destination

S3 bucket

[gd-aws-de-labs](#)

S3 bucket prefix

-

## Buffer hints

Buffer size

5 MiB

## Compression and encryption

Compression for data records


Not enabled




## Amazon Data Firehose <

Firehose streams


### ▼ Resources

[What's new](#) 

[Developer guide](#) 

[API reference](#) 

S3 bucket

[gd-aws-de-labs](#) 

S3 bucket prefix

-

### Buffer hints

Buffer size

5 MiB

Buffer interval

300 seconds

### Compression and encryption

Compression for data records

Not enabled

Encryption for data records

Not enabled

### Backup settings [Info](#)

[Edit](#)

Enabling source record backup ensures that source records can be recovered if record processing transformation does not produce the desired results.

Source record backup in Amazon S3

Not enabled

### Server-side encryption (SSE) [Info](#)

[Edit](#)



You can use AWS Key Management Service (KMS) to create and manage keys and to control the use of encryption across a wide range of AWS services in your applications.



# gd-aws-de-labs Info

- Objects
- Metadata - Preview
- Properties
- Permissions
- Metrics
- Management
- Access Points

Objects (5) Info

Copy S3 URI

Copy URL

Download

Open

Delete

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

< 1 > Settings

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	2025/	Folder	-	-	-
<input type="checkbox"/>	glue-scripts/	Folder	-	-	-
<input type="checkbox"/>	manifests/	Folder	-	-	-
<input type="checkbox"/>	spotify_data/	Folder	-	-	-
<input type="checkbox"/>	streams/	Folder	-	-	-



Editor

Queries

Notebooks

Charts

History

Scheduled queries

Redshift query editor v2

+

Create

▼

⬆

Load data

⏪

🔍

Filter resources

↻

>

🔗

Serverless: default-workgr...

ⓘ

⋮

+

Untitled 1

×

Untitled 2

×

Untitled 3

×

Untitled 4

×

Untitled 5

×

Untitled 6

×

▶

Run

■

🔵

Limit 100

⬜

Explain

🔵

Isolated session

ⓘ

Serverless: de...

▼

clickstream\_db

▼

📅

Schedule

💾

🔄

⋮

11

is\_mobile VARCHAR(3),

12

traffic\_source VARCHAR(50),

13

traffic\_source\_category VARCHAR(50),

14

uri VARCHAR(255),

15

event\_type VARCHAR(50)

16

);

17

18

SELECT \* FROM prod\_schema.events LIMIT 10;

Row 18, Col 1, Chr 416

📊

Result 1 (10)

Export

▼

⬜

Chart

🔄

⌵

<input type="checkbox"/>	id	user_id	session_id	ip_address	city
<input type="checkbox"/>	1822041	48474	54d4fa70-4baa-4b3f-a5b2...	215.107.45.172	Tianjin
<input type="checkbox"/>	1481573	73714	ffbed0dc-8f35-47b4-a229-...	51.228.47.204	Quanz
<input type="checkbox"/>	870743	28796	5f4f29ec-7aeb-4020-b5c3...	36.6.247.40	Caieira
<input type="checkbox"/>	239679	86653	3c82c565-91dc-4fa0-b90...	144.178.157.62	Tokyo
<input type="checkbox"/>	956291	52902	d1a0168b-bea2-4165-b73...	219.60.17.51	East P
<input type="checkbox"/>	1613231	45421	9056a8f8-e41a-4978-8f44...	158.108.8.144	Min

Query ID 2443318

Elapsed time: 76 ms

Total rows: 10