rp6578

**Data to Dashboards**

▼ **AWS account access**

Open AWS console (us-east-1)

**Get AWS CLI credentials**

Exit event

# Access and Validate processed data

Now, Lets test the flow and validate the data.

In the **Data Processing using Amazon Managed Apache Flink** module we initiated ingestion of processed data into **playerkillratio** data stream.

Amazon Data Firehose will deliver this processed data to Amazon S3 as destination under the prefix provided. The buffer is set to 30 seconds or 5 MB, whichever happens first. It may be 30 seconds before you see data in S3.

1. Navigate to the Firehose Console and click on the **gameanalysis-kds-firehose** Data Firehose delivery stream.



2. Click on the Monitoring tab and you should start to see some metrics from your Firehose stream.



3. Click on the Configuration tab. Scroll down to the Destination setting section. Then click on the S3 bucket link.

4. Give it a couple of minutes and you should see some data in your S3 bucket. If not, please wait a bit longer. Navigate into the game-processed-data folder and through the subfolders until you get to the data files in JSON format.



5. Select one of the file. Then under Actions, choose the Select Object Actions then Query with S3 Select.



6. Keep JSON as selected for Input settings. Choose JSON for Output Settings, then select Run SQL Query.

**Query with S3 Select** Info

Use Amazon S3 Select to retrieve a subset of data from an object using standard SQL queries. Pricing is based on the size of the input, query results, and data transferred. Learn more ↗ or see Amazon S3 pricing ↗

**Input settings**

**Path**
s3://game-processed-data-cxde33/year=2024/month=06/day=28/gameanalysis-kds-firehose-3-2024-06-28-15-48-36-506afc35-2321-4f42-a9be-76d1be237f83

**Size**
318.0 B

**Format**
○ CSV
● JSON
○ Apache Parquet

**JSON content type**
● Lines
  Each line in the input data contains a single JSON object
○ Document
  A single JSON object can span multiple lines in the input

**Compression**
● None
○ GZIP
○ BZIP2

**Output settings**

**Format**
○ CSV
● JSON

**SQL query**                              [Add SQL from templates]  [Run SQL query]

Amazon S3 Select supports only the SELECT SQL command. Using the S3 console, you can extract up to 40 MB of records from an object that is up to 128 MB in size. To work with larger files or more records, use the AWS CLI, AWS SDK, or Amazon S3 REST API. For more complex SQL queries, use Amazon Athena ↗

```
1  /* To create reference point for writing SQL queries, you can display the first 5 records of input data by running the following SQL query: SELECT * FROM
      s3object s LIMIT 5 */
2  SELECT * FROM s3object s LIMIT 5
```

SQL    Ln 1, Col 1    ⊗ Errors: 0  | ⚠ Warnings: 0                                                            ⚙

**Query results**                                                      [⊟ Download results]

Query results are not available after you choose **Close** or navigate away. Choose **Download results** to download a copy of the following query results.

**Status**
⊘ Successfully returned 3 records in 565 ms

Bytes returned: 329 B

```
1  {
2      "player_id": 5,
3      "kills": 731,
4      "death": 700,
5      "kill_ratio": 1.0442857142857143,
6      "window_end": "2024-06-28 15:48:30"
7  }
8  {
9      "player_id": 1,
10     "kills": 547,
11     "death": 555,
12     "kill_ratio": 0.9855855855855856,
13     "window_end": "2024-06-28 15:48:30"
14 }
15 {
16     "player_id": 3,
17     "kills": 172,
18     "death": 159,
19     "kill_ratio": 1.0817610062893082,
20     "window_end": "2024-06-28 15:49:00"
21 }
```

The processed game data can be observed in the S3 bucket, indicating that the Amazon Data Firehose service has successfully completed its task.

# AWS Glue Setup

1. Navigate to the AWS Glue Console ↗. Choose Crawlers in the navigation pane.
2. Select **Create Crawler**.



3. Provide name of the crawler as **gamedata-crawler** and select **Next** .

AWS Glue > Crawlers > Add crawler

Step 1
**Set crawler properties**

Step 2
Choose data sources and
classifiers

Step 3
Configure security settings

Step 4
Set output and scheduling

Step 5
Review and create

## Set crawler properties

### Crawler details  Info

Name

gamedata-crawler

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

Description - *optional*

To Crawl S3 bucket

Descriptions can be up to 2048 characters long.

▶ **Tags -** *optional*
  Use tags to organize and identify your resources.

Cancel    Next

4. Click on the add a data source and browse the s3 bucket 'gameanalysis-kds-firehose-*'. Select **Next**.

AWS Glue > Crawlers > **Add crawler**

Step 1
Set crawler properties

Step 2
**Choose data sources and
classifiers**

Step 3
Configure security settings

Step 4
Set output and scheduling

Step 5
Review and create

## 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                    Edit    Remove    Add a data source
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 ⟵

## 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

**S3 path**
Browse for or enter an existing S3 path.

| 🔍 s3://game-processed-data-xxxxxx        ✕ |   [ View ⬈ ]   |   [ Browse S3 ] |
|---|---|---|

All folders and files contained in the S3 path are crawled. For example, type s3://MyBucket/MyFolder/ to crawl all objects in MyFolder within MyBucket.

**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.

☐ Sample only a subset of files

☐ Exclude files matching pattern

                                                    Cancel     **Add an S3 data source**

5. Create a new IAM role and Select **Next**.

AWS Glue > Crawlers > Add crawler

**Configure security settings**

Step 1
Set crawler properties

Step 2
Choose data sources and classifiers

Step 3
**Configure security settings**

Step 4
Set output and scheduling

Step 5
Review and create

**IAM role**  Info

Existing IAM role

| AWSGlueServiceRole-s3crawl ▼ | ↻ | View ⧉ |

| 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**

6. Choose **gameanalytics** as as database and Select **Next**.

AWS Glue > Crawlers > Add crawler

**Set output and scheduling**

Step 1
Set crawler properties

Step 2
Choose data sources and classifiers

Step 3
Configure security settings

Step 4
**Set output and scheduling**

Step 5
Review and create

**Output configuration**  Info

Target database

| gameanalytics ▼ | ↻ |

| Clear selection | Add database ⧉ |

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 ▼ |

Cancel    Previous    **Next**

7. Review the crawler and click on **Create crawler**.

AWS Glue > Crawlers > Add crawler

**Review and create**

Step 1
Set crawler properties

Step 2
Choose data sources and classifiers

Step 3
Configure security settings

Step 4
Set output and scheduling

Step 5
**Review and create**

**Step 1: Set crawler properties**    Edit

Set crawler properties

| Name | Description | Tags |
|------|-------------|------|
| gamedata-crawler | - | - |

**Step 2: Choose data sources and classifiers**    Edit

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

| Type | Data source | Parameters |
|------|-------------|------------|
| S3 | s3://game-processed-data- | Recrawl all |

**Step 3: Configure security settings**    Edit

Configure security settings

| IAM role | Security configuration | Lake Formation configuration |
|----------|------------------------|------------------------------|
| AWSGlueServiceRole-s3crawl | - | - |

**Step 4: Set output and scheduling**    Edit

Set output and scheduling

| Database | Table prefix - *optional* | Maximum table threshold - *optional* | Schedule |
|----------|---------------------------|--------------------------------------|----------|
| gameanalytics | - | - | On demand |

Cancel    Previous    **Create crawler**

8. After successfully creating the crawler, run the crawler.

**One crawler successfully created**
The following crawler is now created: "gamedata-crawler"

AWS Glue > Crawlers > gamedata-crawler

## gamedata-crawler

Last updated (UTC)
June 25, 2024 at 05:12:02   Run crawler   Edit   Delete

**Crawler properties**

| Name | IAM role | Database | State |
|---|---|---|---|
| gamedata-crawler | AWSGlueServiceRole-s3crawl | gameanalytics | READY |

| Description | Security configuration | Lake Formation configuration | Table prefix |
|---|---|---|---|
| - | - | - | - |

Maximum table threshold
-

▶ Advanced settings

9. Check the status of the crawler.

Crawler runs   Schedule   Data sources   Classifiers   Tags

**Crawler runs (1)**
The list of crawler runs for this crawler.

Stop run   View CloudWatch logs   View run details

| | Start time (UTC) | End time (UTC) | Current/last duration | Status | DPU hours | Table changes |
|---|---|---|---|---|---|---|
| ○ | June 25, 2024 at 05:12:34 | June 25, 2024 at 05:13:57 | 01 min 23 s | ⊘ Completed | 0.052 | 1 table change, 1 partition change |

10. Now, you can see the table and schema information under the data catalog section.

AWS Glue

AWS Glue > Tables > game_processed_data_cxde33

### game_processed_data_

Last updated (UTC)
June 25, 2024 at 05:32:28   Version 1 (Current version) ▾   Actions ▾

Table overview   Data quality New

Table details   Advanced properties

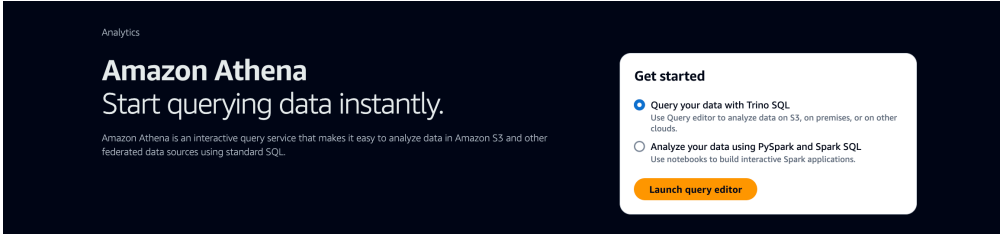| Name | Description | Database | Classification |
|---|---|---|---|
| game_processed_data_cxde33 | - | gameanalytics | JSON |

| Location | Connection | Deprecated | Last updated |
|---|---|---|---|
| s3://game-processed-data- | - | - | June 25, 2024 at 05:13:56 |

| Input format | Output format | Serde serialization lib |
|---|---|---|
| org.apache.hadoop.mapred.TextInputFormat | org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat | org.openx.data.jsonserde.JsonSerDe |

Schema   Partitions   Indexes   Column statistics - new

**Schema (7)**
View and manage the table schema.

Edit schema as JSON   Edit schema

| # | Column name | Data type | Partition key | Comment |
|---|---|---|---|---|
| 1 | player_id | int | - | - |
| 2 | kill_ratio | double | - | - |
| 3 | window_end | string | - | - |
| 4 | year | string | Partition (0) | - |
| 5 | month | string | Partition (1) | - |
| 6 | day | string | Partition (2) | - |
| 7 | hour | string | Partition (3) | - |

# Amazon Athena Setup

1. Navigate to the Athena Console ↗ and launch query editor.

Analytics

## Amazon Athena
### Start querying data instantly.

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

**Get started**

◉ Query your data with Trino SQL
Use Query editor to analyze data on S3, on premises, or on other clouds.

○ Analyze your data using PySpark and Spark SQL
Use notebooks to build interactive Spark applications.

Launch query editor

2. Click on edit setting to setup query result location in the S3 bucket.

Amazon Athena > Query editor

Editor   Recent queries   Saved queries   Settings     Workgroup   primary ▾

ⓘ Before you run your first query, you need to set up a query result location in Amazon S3.    Edit settings

ⓘ **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.    Edit preferences ✕

**Data**   Query 1
Data source
AwsDataCatalog
Database
gameanalytics

3. Choose S3 bucket athena-result-bucket-* and **Save** the setting.

**Choose S3 data set**                                                                                          ×

S3 buckets

**Bucket** (96)                                                                                                  ⟳

🔍 athena-result-bucket|                                    ✕    1 match                                    ‹  1  ›

| | Name | ▽ | Creation date | ▽ |
|---|---|---|---|---|
| ○ | athena-result-bucket- ▭▭▭ | | 2024-06-25T00:46:13.000-05:00 | |

                                                                                    Cancel        Choose

4. Choose the **AwsDataCatalog** as the data source, the **gameanalytics** database, and the table name according to your Glue table name. Then, run the following query based on your table name.

```
SELECT * from game_processed_data_xxxxxx;
```



# Conclusion

In this module, you gained hands-on experience managing and delivering processed real-time data to required destinations and making it available for decision making. With this let's move on to next module and visualizing this processed data.

Previous          Next