# Let's dive into setting up AWS Glue with an ETL job, then querying the data via AWS Athena

Craig Godden-Payne
Feb 7 · 5 min read ★



Let's look at a simple end to end run through of using AWS Glue to transform data from a format, into a more queryable format, and then query it using AWS Athena. We will look at this through the console only, with more focus on how to automate this with terraform in the future post.

All the data in this post are from apache logs, which can be downloaded from Github. The data has been broken into 5 pieces, to simulate that the logs were uploaded at 5 different times.

You can download the splits here: log 1 log 2 log 3 log 4 log 5

·   ·   ·

## Uploading the data

In order to query the data in AWS, you will need to upload the data files into an S3 bucket, you can use the example files above, or just some random apache log files.

S3 bucket

. . .

## Setting up an AWS Glue Job

In the AWS console, search for Glue. Once it is open, navigate to the Databases tab. Create a new database, I created a database called craig-test. The database is used as a data catalog and stores information about schema information, not actual data.

**Data Driven Investor | Microsoft Having An 'Edge' Over Chrome**

A Brief History I was never a fan of browsers, well to be exact I was only a fan of one, Chrome. It has been my…

www.datadriveninvestor.com

Check the tables section, you will see the message that no tables have been defined in the data catalog. The next thing to do is to create a crawler, which will gather metadata about the data.



Data calatlog view

. . .

## Creating a crawler

Whilst remaining in glue, create a crawler.

Point the crawler to look at the S3 bucket, and set it to write to the database data catalogue that was created before.

Add a crawler

If you set the crawler to be on demand, you need to run it once you have finished creating the crawler.

Glue will pretty pretty quickly crawl and determine the format of the log file. For my case, it mentions the data looks like it is in apache format, and it detailed what the schema looks like.



Example of what the schema looks like

Schema

| | Column name | Data type | Partition key | Comment |
|---|---|---|---|---|
| 1 | clientip | string | | |
| 2 | ident | string | | |
| 3 | auth | string | | |
| 4 | timestamp | string | | |
| 5 | verb | string | | |
| 6 | request | string | | |
| 7 | httpversion | string | | |
| 8 | response | string | | |
| 9 | bytes | string | | |
| 10 | referrer | string | | |
| 11 | agent | string | | |

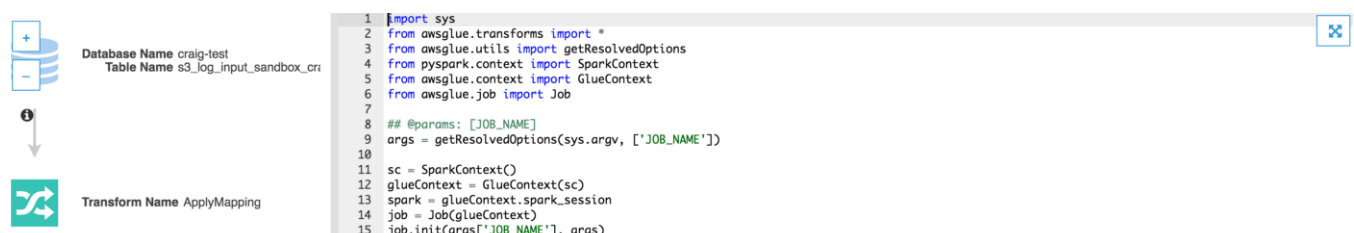More detail of the schema

. . .

# Creating the ETL job

The next stage is to create an ETL job, which will take the logs, and transform them into columnar data in parquet format, from apache log format.

It is much more efficient and cost effective to query in Athena using a format such as parquet rather than having to scan the whole data each time.

Created n Glue job, and add a data source, I used craig-test bucket. Set the destination bucket, I used craig-test-processed bucket.

Select that you want to have the destination in parquet format. On clicking next, the apache spark code will be auto generated, but easy enough to read and make changes to if it needs to.



```
Database Name craig-test
Table Name s3_log_input_sandbox_cra

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
8  ## @params: [JOB_NAME]
9  args = getResolvedOptions(sys.argv, ['JOB_NAME'])
10
11 sc = SparkContext()
12 glueContext = GlueContext(sc)
13 spark = glueContext.spark_session
14 job = Job(glueContext)
15 job.init(args['JOB_NAME'], args)

Transform Name ApplyMapping
```

Auto Generated Code



Auto mapped columns

Once you have completed creating a job, run the job. After a few minutes the job will be completed.

. . .

## Issues or Failures

On my first run, I forgot about permissions, but if you hit any issues, they are pretty easy to diagnose by looking at Cloudwatch.

| | | |
|---|---|---|
| ▸ | 12:45:18 | Completed 2.9 KiB/2.9 KiB (71.6 KiB/s) with 1 file(s) remaining |
| ▸ | 12:45:18 | download: s3://aws-glue-scripts-████████-eu-west-1/admin/craig-test- to ./script_2019-11-29-12-45-13.py |
| ▸ | 12:45:19 | SCRIPT_URL = /tmp/g-42f474d5ef6d04a9908ef3a429eefe1d0d922494-1503052698715096253/script_2019-11-29-12-45-13.py Skipping compilation step fo |
| ▸ | 12:45:19 | /usr/lib/spark/bin/spark-submit --conf spark.hadoop.yarn.resourcemanager.connect.max-wait.ms=60000 --conf spark.hadoop.fs.defaultFS=hdfs://ip-172-32- |

*No newer events found at the moment.* Retry.

Cloudwatch

In my case, it was due to IAM permissions. I forgot to allow the glue role to have access to my output bucket.

| Permissions | Policy usage | Policy versions | Access Advisor |
|---|---|---|---|

| Policy summary | {} JSON | Edit policy |
|---|---|---|

```
1  {
2      "Version": "2012-10-17",
3      "Statement": [
4          {
5              "Effect": "Allow",
6              "Action": [
7                  "s3:*"
8              ],
9              "Resource": [
10                 "arn:aws:s3:::sandbox-craig-test/*",
11                 "arn:aws:s3:::sandbox-craig-test-processed/*"
12             ]
13         }
14     ]
15 }
```

Bucket policy

| History | Details | Script | Metrics |
|---|---|---|---|

View run metrics    Rewind job bookmark      Showing: 1 - 3 ‹ ›

| Run ID | Retry attempt | Run status | Error | Logs | Error logs | Glue version | Maximum capacity | Triggered by | Start time | End time | Execution time | Timeout Delay | Job run input |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ◯ jr_ca0e663061… | - | Succeeded | | Logs | | 1.0 | 10 | | 29 N… | 29 N… | 52 secs | 2880 mins | s3://aws-glue-t… |
| ◯ jr_28d1dc01ec… | - | Failed | ❗P… | Logs | Error logs | 1.0 | 10 | | 29 N… | 29 N… | 1 min | 2880 mins | s3://aws-glue-t… |

Job runs

. . .

# Rebuilding the schema

Check the s3 processed bucket, and you should find the parquet files are there ready for me to query using AWS Athena.



S3 output

You will need to rebuild the schema to include the new data sets. In my case I updated to include the transformed columnar parquet format.

This time, create a new crawler, and crawl the S3 bucket containing the processed files.

Add crawler

Once it has been ran, and the schema has been added to the data catalogue, it is time to move over to query.

.   .   .

## Query using AWS Athena

Once in Athena, you can query the data from the logs. Querying the parquet logs is very cost effective and quick. whereas querying the original data is expensive and slow.



AWS Athena

**Results**

| | clientip | ident | auth | timestamp | verb | request | httpversion | response | bytes | referrer |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 216.14.102.16 | | | | HEAD | /projects/xboxproxy | 1.1 | 301 | | |
| 2 | 89.170.74.95 | | | | HEAD | /projects/firefox-urledit/urledit-screencast.html | 1.1 | 200 | | |
| 3 | 212.48.66.64 | | | | HEAD | / | 1.0 | 200 | | |
| 4 | 216.14.102.16 | | | | HEAD | /blog/geekery/freebsd-ports-master-sites-sorting.html | 1.1 | 200 | | |
| 5 | 216.14.102.16 | | | | HEAD | /projects/xboxproxy | 1.1 | 301 | | |
| 6 | 89.170.74.95 | | | | HEAD | /projects/firefox-urledit/urledit-screencast.html | 1.1 | 200 | | |
| 7 | 212.48.66.64 | | | | HEAD | / | 1.0 | 200 | | |
| 8 | 216.14.102.16 | | | | HEAD | /files/fastest_sites/fastest_sites-20110317.py | 1.1 | 200 | | |
| 9 | 216.14.102.16 | | | | HEAD | /projects/fex/ | 1.1 | 200 | | |
| 10 | 37.115.186.244 | | | | POST | /blog/geekery/xvfb-firefox | 1.0 | 200 | 10975 | http://zoon |

Athena Query

The difference in this simple query, is with the parquet data. If you look at the amount of data processed, you can see how much more expensive it would be

| Query submitted time | Query | Encryption type ⓘ | State | Run time(s) | Data scanned | Action |
|---|---|---|---|---|---|---|
| 2019/11/29 14:01:59 UTC+0 | SELECT * FROM s3_log_input_sandbox_craig_test WHERE verb <> 'GET' Limit 10 | N/A | Succeeded | 2.83 | 2.26 MB | Download results |
| 2019/11/29 13:58:00 UTC+0 | SELECT * FROM parquet_sandbox_craig_test_processed WHERE verb <> 'GET' Limit 10 | N/A | Succeeded | 2.08 | 173.92 KB | Download results |

## Conclusion

Wherever possible, query data in a structured format

*Data Driven Investor*

## Gain Access to Expert Views

Email

First Name

Give me access!

☐ I agree to leave Medium.com and submit this information, which will be collected and used according to Upscribe's privacy policy.

**4.2K signups**

AWS      Glue      Athena      Python      Data Analytics

About   Help   Legal

Get the Medium app