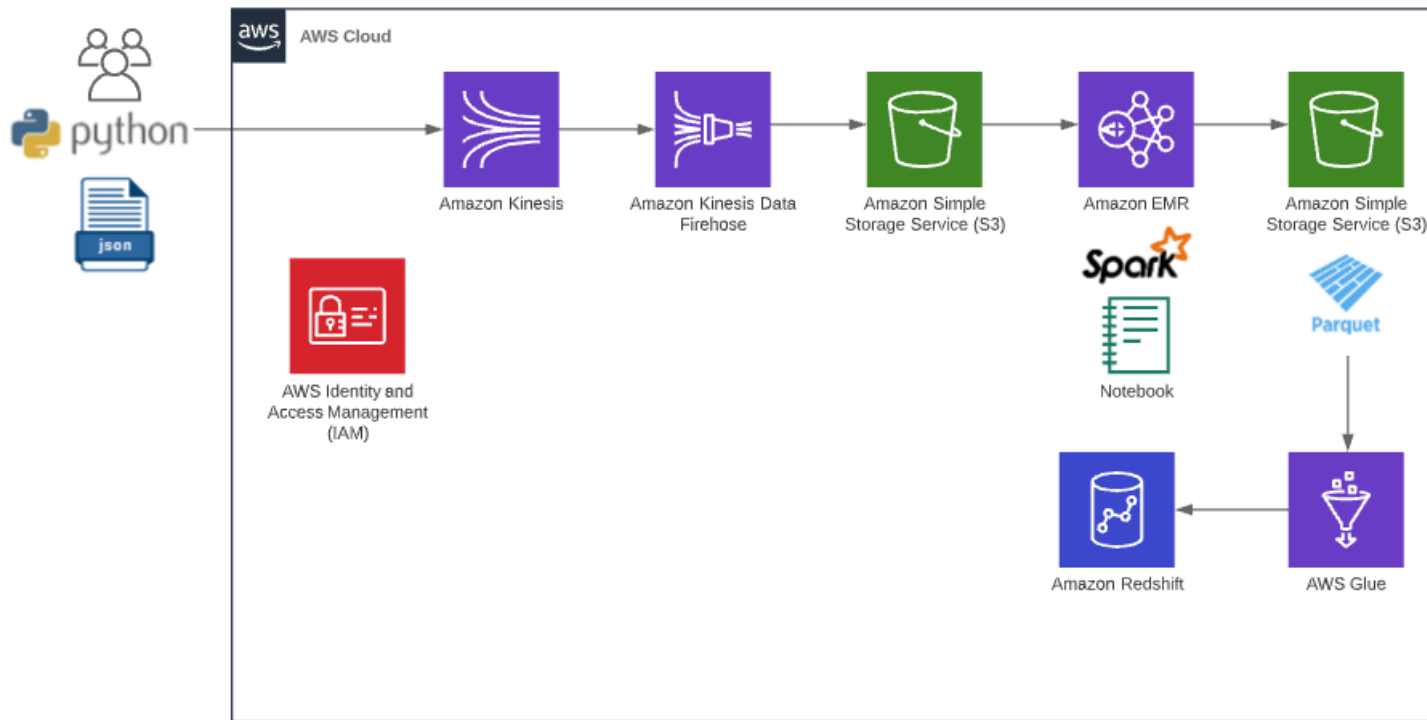


I spent a few days to familiarize myself a bit with AWS. I have been doing a lot of reading lately about it, but had only created a couple AWS VMs. I'm much more familiar with Azure, having worked with it for the past year or so as part of the UW Big Data Technology certificate program.

Here's the overall architecture I came up with, which I'll describe in more detail below:



I was especially interested in working with some streaming data. As with most cloud platforms there are multiple ways in AWS to ingest data streams, including Kinesis, Amazon Managed Streaming for Apache Kafka (Amazon MSK), or Kafka on an AWS VM (for this POC). I opted for Kinesis just to get experience with it. Perhaps I'll try MSK next, as I do have some experience using Kafka.

I first set up a Kinesis Data stream. Then I wrote small Python scripts to act as producer and consumer for Kinesis, generating some random data (see the notebooks). I was able to see the streaming data in the consumer script.

I then set up a Kinesis Firehose delivery stream to move the data from the data stream to an S3 bucket, where I was able to see the records:

petebchamp-kinesis-firehose-2-2020-07-30-01-19-38-942... [Latest version](#) ▼

SQL editor [Sample SQL expressions](#) [Learn more](#)

1 `select * from s3object s limit 5`

Copy

Run SQL

Result

```
{
  "createdate": "2020-07-29 18:19:39.550996",
  "firstname": "Cory",
  "lastname": "Parks",
  "gender": "M",
  "birthdate": "2013-02-01"
}
{
  "createdate": "2020-07-29 18:19:42.714322",
  "firstname": "David",
  "lastname": "Cooper",
  "gender": "F",
  "birthdate": "2001-09-02"
}
{
  "createdate": "2020-07-29 18:19:45.785644",
  "firstname": "Kayla",
  "lastname": "Chavez",
  "gender": "F",
  "birthdate": "1994-10-29"
}
{
  "createdate": "2020-07-29 18:19:48.844165",
  "firstname": "James",
  "lastname": "Johnson",
  "gender": "F",
  "birthdate": "1929-07-13"
}
{
  "createdate": "2020-07-29 18:19:51.924311",
  "firstname": "Adam",
  "lastname": "Garrett",
  "gender": "M",
  "birthdate": "1998-01-07"
}
```

I wanted to do some manipulation of the data, and looked into using AWS Lambda, but opted at this time to use Spark in AWS EMR. I created two Jupyter Spark notebooks, both of which check for and eliminate duplicate records and then write the data back to S3 as parquet files. The first one just deals with a static/batch set of data, as I was working through the Spark code. I then modified that one to read and write streaming data. Both notebooks are included here. I created several iterations of EMR clusters just to go through all the options, and needed to include Hive to be able to use Spark SQL code, after having that error the first time.

Here are the parquet files in S3:

Amazon S3 > petebchamp-s3 > parquet\_files

## petebchamp-s3

Overview

🔍 Type a prefix and press Enter to search. Press ESC to clear.

[Upload](#) [+ Create folder](#) [Download](#) [Actions](#) [Versions](#) [Hide](#) [Show](#)

<input type="checkbox"/> Name ▼	Last modified ▼
<input type="checkbox"/> _SUCCESS	Jul 31, 2020 8:32:36 PM GMT-0700
<input type="checkbox"/> part-00000-6bde7470-e200-491d-921d-f7e24bd7c412-c000.snappy.parquet	Jul 31, 2020 8:32:36 PM GMT-0700
<input type="checkbox"/> part-00000-9e7871a9-b4d4-43d4-97fc-3dc5c942adc1-c000.snappy.parquet	Jul 29, 2020 8:24:41 PM GMT-0700

I wanted to write the data directly from Spark to Redshift, and looked into installing the Redshift JDBC, but it seemed very involved, with installations of the driver and the correct version of Java components. Instead, I just decided to use AWS Glue to copy the data to a Redshift table I created.

```
create table person (  
    create_date TIMESTAMP ,  
    first_name varchar(20),  
    last_name varchar(50),  
    gender char(1),  
    birth_date date  
);
```

I did have to set up a new IAM role and associate it to the Glue for the job to complete successfully. Also, I first created the table schema in Glue manually, and then tried having a Crawler discover the schema based on the parquet files, and it came up with the same schema I had created.

If I were doing anything more elaborate with Redshift I would definitely want to use a client tool like SQL Workbench or pgAdmin, although it is convenient to be able to query right from the console.

Query editor

petebchamp-redshift-cluster

test

Query 1

Query 2

Query 3

+

1

select createdate, firstname, lastname, gender, birthdate

2

from person

Rows returned (50)

Export

Search rows

< 1 2 3 4 5 > ⚙

createdate	firstname	lastname	gender	birthdate
2020-07-29 18:20:40.989000	Benjamin	Wang	F	1943-09-19
2020-07-29 18:20:31.779000	Brandon	Sandoval	M	1954-12-22
2020-07-29 18:20:53.232000	Katherine	Kirby	F	1985-05-19
2020-07-29 18:20:59.365000	Caitlin	Moore	F	1928-07-03
2020-07-29 18:19:42.714000	David	Cooper	F	2001-09-02
2020-07-29 18:19:54.989000	Mark	Adkins	M	2004-07-30
2020-07-29 18:21:02.433000	Preston	May	M	1953-10-14