

1. AWS Glue
2. AWS Redshift
3. AWS Athena
4. SQS & SNS
5. AWS Lambda
6. AWS EMR
7. CloudWatch & Monitoring

8. Apache Hadoop
9. Lambda function

Quiz :-

1. 7th June

2. 14th June (Post training Quiz)

AWS Glue

1. AWS Crawler

2. Glue ETL Script

3. AWS Glue Workflows

4. Glue concept

5. Read CSV file from S3 Bucket w/ Glue Catalogue

Glue Catalogue -

6. Insert data into Redshift w/o Glue Catalogue

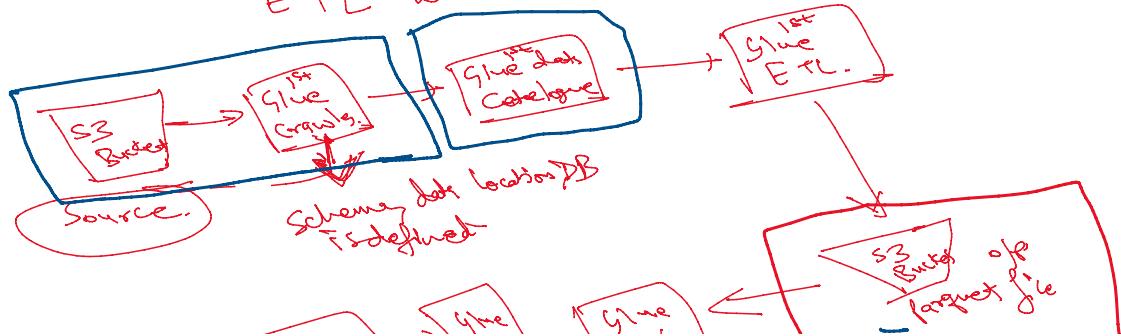
7. S3 + Glue + Athena

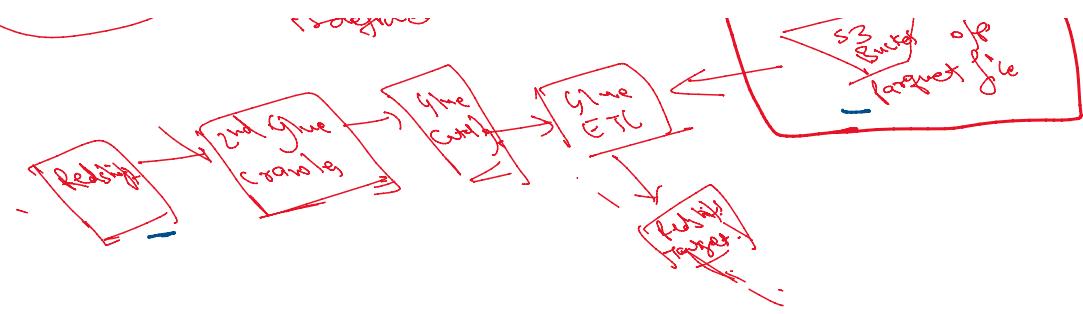
8. S3 + Glue + Redshift

9. By spark concept

10. Frame concept

AWS Glue provides console & API operation to setup & manage your ETL workload.





① Crawlers-

Connect with Source to Target - to store the Schema of the metadata information of the data.

② Glue Job (ETL)

Process the data python or PySpark-

Python - Small dataset - why? Not distributed processing. Numpy to perform.

PySpark - Large dataset - distributed processing → RAM
Speed very high
as compared to python

Glue
→
Large dataset
Processing → up to 1 hour
Serverless source

Lambda
→ small dataset
Processing → 15 min
Serverless service

Serverless & No
Cost

↓
Times of
execution

Steps to create the crawlers -

- ① Save data in S3 Bucket -
- ② Read data from S3 using Crawler
- ③ Create the crawler to reference S3 Bucket
- ④ Create IAM Role → S3 full access, AWS Lambda & CloudWatch Metrics
- ⑤ Column as partition column. Data will be visible in Athena
- ⑥ 1. ... moments update - append or modify files

~~1~~ Column as position column ---
~~2~~ Crawler increments update - append or modify files
to the existing catalogue files. It will not insert
the old data file again