# Advanced Data Engineering in Cloud

## ASSIGNMENT-2
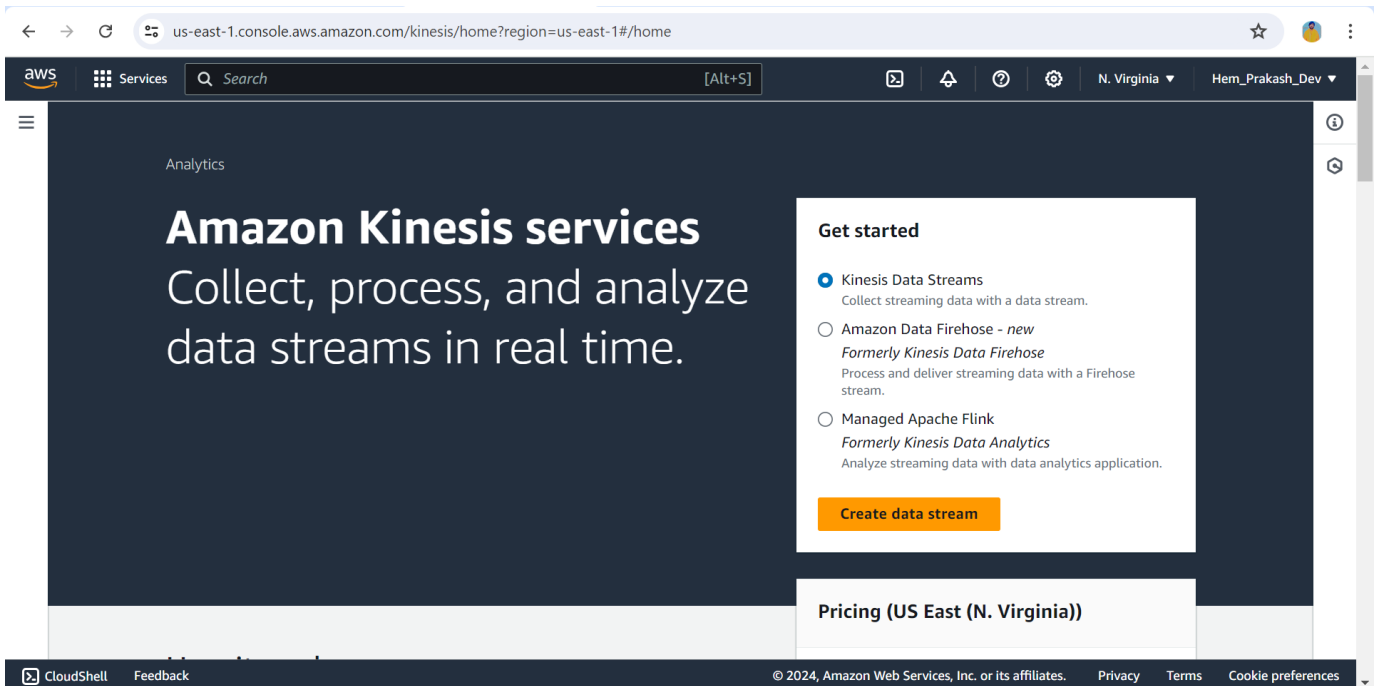
### Data Ingestion and Processing

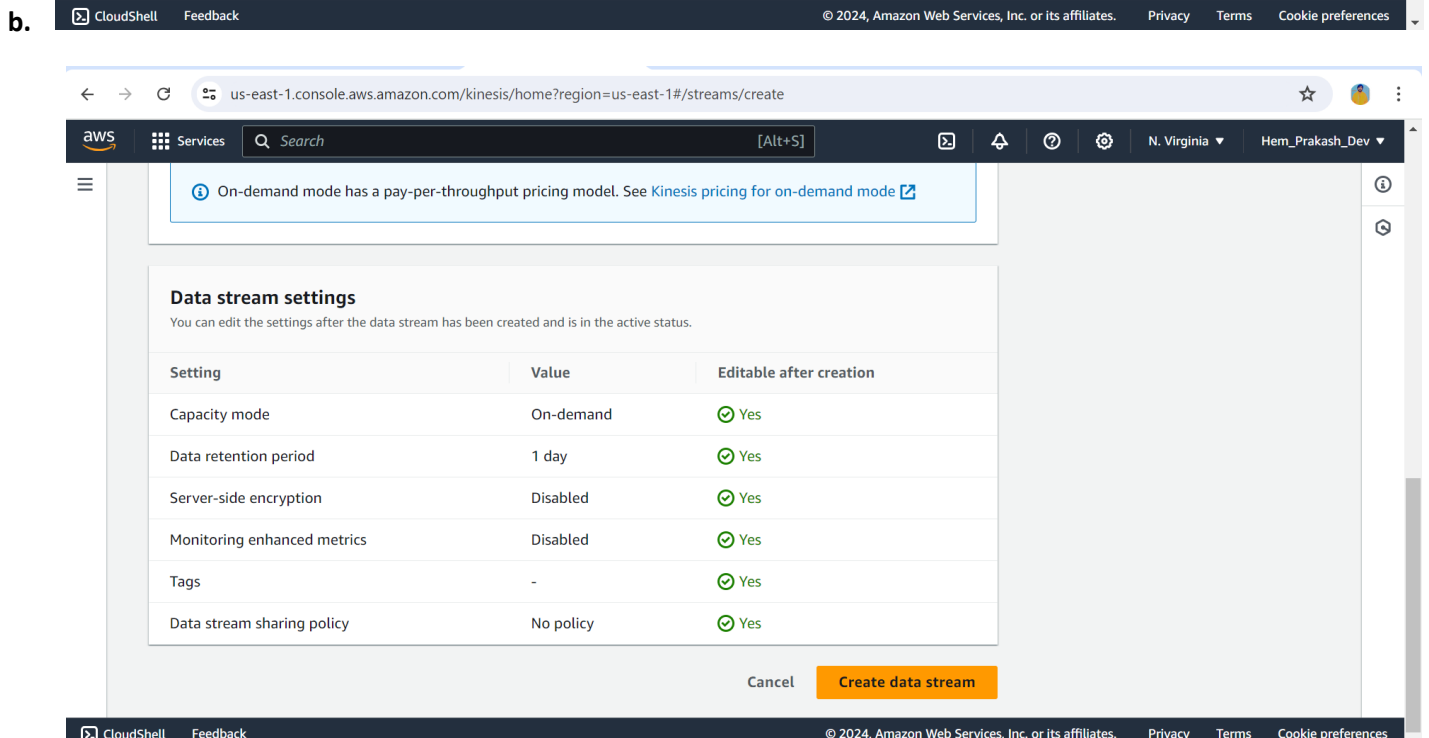**Hem Prakash Dev**                                    **Roll-G23AI1054**

1. **Implement the data ingestion mechanism using AWS Kinesis Data Streams or AWS Direct Connect to stream data from a source to Amazon S3.**

   a. **Create a Kinesis Data Stream**:

   

   b.

   

## b. stream Data to S3 using Kinesis Firehose:

aws | Services | Q Search | [Alt+S] | N. Virginia ▼ | Hem_Prakash_Dev ▼

Amazon Data Firehose > Firehose streams > Create Firehose stream

# Create Firehose stream  Info

▶ **Amazon Data Firehose: How it works**

## Choose source and destination

Specify the source and the destination for your Firehose stream. You cannot change the source and destination of your Firehose stream once it has been created.

**Source**  Info

Amazon Kinesis Data Streams  ▼

**Destination**  Info

Amazon S3  ▼

aws | Services | Q Search | [Alt+S] | N. Virginia ▼ | Hem_Prakash_Dev ▼

## Source settings

**Kinesis data stream**

arn:aws:kinesis:us-east-1:905418436402:stream/kinesis-db | Browse | Create ⬈

Format: arn:aws:kinesis:[Region]:[AccountId]:stream/[StreamName]

## Firehose stream name

**Firehose stream name**

kinesis-fh

Acceptable characters are uppercase and lowercase letters, numbers, underscores, hyphens, and periods.

## Transform and convert records - *optional*

Configure Amazon Data Firehose to transform and convert your record data.

**Transform source records with AWS Lambda**   Info

Amazon Data Firehose can invoke an AWS Lambda function to transform, filter, decompress, convert and process your source data records

aws | Services | Q Search | [Alt+S] | N. Virginia ▼ | Hem_Prakash_Dev ▼

## Destination settings  Info

Specify the destination settings for your Firehose stream.

**S3 bucket**

my-asw-s3-bucket | Browse | Create ⬈

Format: s3://bucket

**New line delimiter**

You can configure your Firehose stream to add a new line delimiter between records in objects that are delivered to Amazon S3.

○ Not enabled
● Enabled

**Dynamic partitioning**   Info

Dynamic partitioning enables you to create targeted data sets by partitioning streaming S3 data based on partitioning keys. You can partition your source data with inline parsing and/or the specified AWS Lambda function. You can enable dynamic partitioning only when you create a new Firehose stream. You cannot enable dynamic partitioning for an existing Firehose stream. Enabling dynamic partitioning incurs additional costs per GiB of partitioned data. For more information, see Amazon Data Firehose pricing ⬈.

○ Not enabled
● Enabled

ⓘ You are choosing to enable dynamic partitioning for this Firehose stream. Once this Firehose stream is

▶ Example record
Use the example below to define the dynamic partitioning keys and key values.

Dynamic partitioning keys
Use the fields below to specify key names and JQ expressions to be used as dynamic partitioning keys. Amazon Data Firehose only supports expressions in JQ 1.6.

Key name                                JQ expression

senser-id                               sensor-id

Add dynamic partitioning key

Dynamic partitioning keys must be unique. You can add 49 more dynamic partitioning key values.

S3 bucket prefix
For dynamic partitioning, you must use the following expression format in your S3 bucket prefix: !{namespace:value}, where namespace can be either partitionKeyFromQuery or partitionKeyFromLambda, or both. If you are using inline parsing to create the partitioning keys for your source data, you must specify an S3 bucket prefix value that consists of expressions specified in the following format: "!{partitionKeyFromQuery:keyID}". If you are using an AWS Lambda function to create partitioning keys for your source data, you must specify an S3 bucket prefix value that consists of expressions specified in the following format: "!{partitionKeyFromLambda:keyID}".

!{partitionKeyFromQuery:senser-id}/

Apply dynamic partitioning keys

▼ Buffer hints, compression, file extension and encryption
The fields below are pre-populated with the recommended default values for S3. Pricing may vary depending on storage and request costs.

S3 buffer hints
Amazon Data Firehose buffers incoming records before delivering them to your S3 bucket. Record delivery is triggered once the value of either of the specified buffering hints is reached.

ⓘ For Firehose streams with dynamic partitioning enabled, we recommend a buffer size of 128 MiB for optimized processing of data. For more information, see Amazon Data Firehose pricing ↗.

Buffer size
The higher buffer size may be lower in cost with higher latency. The lower buffer size will be faster in delivery with higher cost and less latency.

128          MiB

Minimum: 64 MiB, maximum: 128 MiB. Recommended: 128 MiB.

Buffer interval
The higher interval allows more time to collect data and the size of data may be bigger. The lower interval sends the data more frequently and may be more advantageous when looking at shorter cycles of data activity.

60           seconds

Minimum: 0 seconds, maximum: 900 seconds. Recommended: 300 seconds.

# Kinesis firehose successfully created

**2. Develop and test the data processing pipeline using AWS Glue or Amazon EMR (Elastic MapReduce) with Apache Spark or Hadoop.**

**Data processing pipeline using Amazon EMR (Elastic MapReduce) with Apache Spark**

## EMR-Console



**Cluster creation**

Services  Search  [Alt+S]  N. Virginia ▾  Hem_Prakash_Dev ▾

⊘ Your cluster "my-emr-Cluster" has been successfully created.

Amazon EMR > EMR on EC2: Clusters > my-emr-Cluster

# my-emr-Cluster

Updated less than a minute ago  ↻  Terminate  Clone in AWS CLI  Clone

## ▼ Summary

### Cluster info

Cluster ID
j-1LNLGJ2F24G7O

Cluster configuration
Instance groups

Capacity
1 Primary | 2 Core | 0 Task

### Applications

Amazon EMR version
emr-7.1.0

Installed applications
Hadoop 3.3.6, Hive 3.1.3,
JupyterEnterpriseGateway 2.6.0, Livy
0.8.0, Spark 3.5.0

### Cluster management

Log destination in Amazon S3
aws-logs-905418436402-us-east-1/elasticmapreduce

Primary node public DNS
-

### Status and time

Status
⊙ Starting

Creation time
July 16, 2024, 21:46 (UTC+05:30)

Elapsed time
0 seconds

Properties | Bootstrap actions | Instances (Hardware) | Steps | Applications | Configurations | Monitoring | Events | Tags (1)

Cluster logs Info

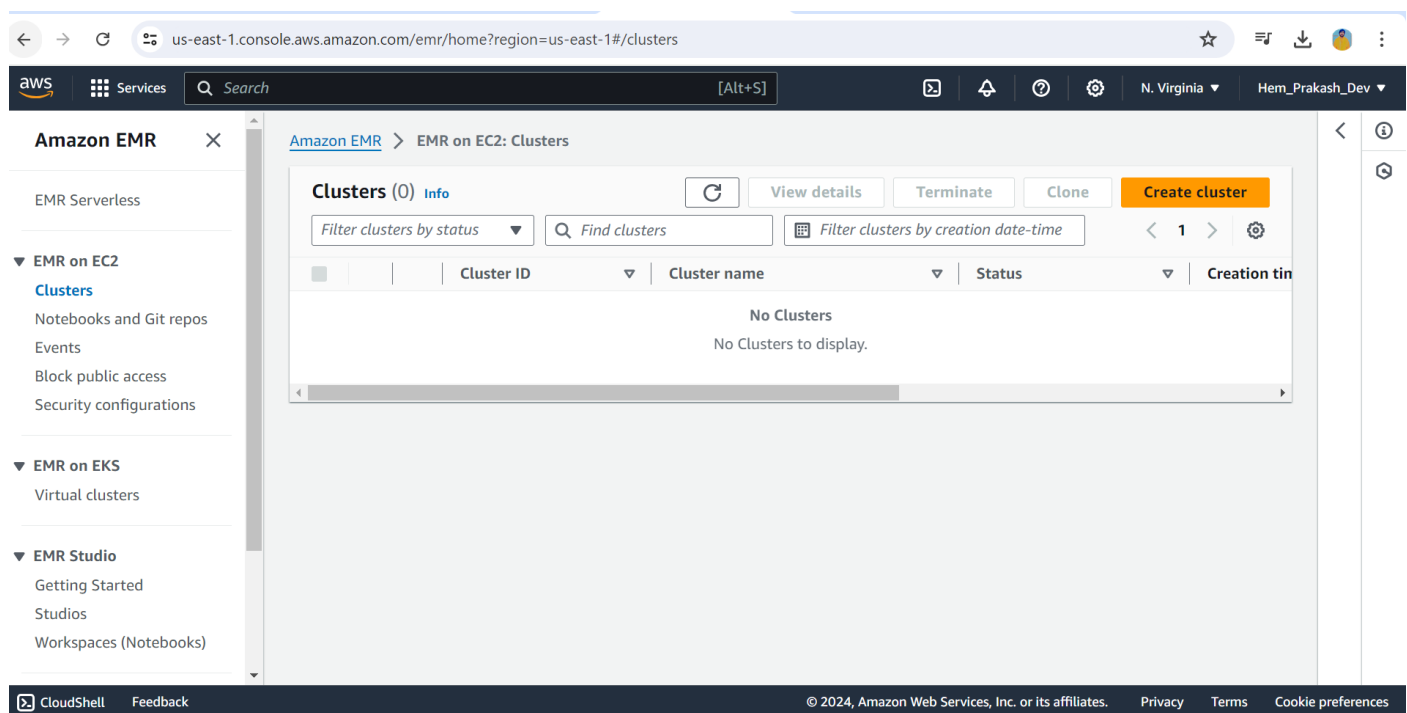Cluster termination and node replacement Info  Edit

CloudShell  Feedback  © 2024, Amazon Web Services, Inc. or its affiliates.  Privacy  Terms  Cookie preferences

---

### Operating system Info

Amazon Linux release
2023.5.20240708.0

### Cluster logs Info

Archive log files to Amazon S3
Turned on

Amazon S3 location
s3://aws-logs-905418436402-us-east-1/elasticmapreduce/ ↗

Encryption for logs
Turned off

### Cluster termination and node replacement Info

Edit

Termination option
Automatically terminate cluster after idle time

Idle time
1 hour

Termination protection
Off

Unhealthy node replacement
On

## Network and security Info

### Network

Virtual Private Cloud (VPC)
vpc-02d092a3c77cc5804 ↗

### Security configuration

Security configuration
None

### Permissions

Service role for Amazon EMR
AmazonEMR-ServiceRole-20240716T214601 ↗

CloudShell  Feedback  © 2024, Amazon Web Services, Inc. or its affiliates.  Privacy  Terms  Cookie preferences

aws | Services | Q Search | [Alt+S] | Ohio ▼ | Hem_Prakash_Dev ▼

| | All UDP ▼ | UDP | 0 - 6553! | Cus... ▼ | Q | | Delete |

sg-05ccc38efd66854e3 ✕

| sgr-090ced6ff35884ba0 | All UDP ▼ | UDP | 0 - 6553! | Cus... ▼ | Q | | Delete |

sg-0346b33f6c921cbd4 ✕

| sgr-0061ffbbb43463292 | All TCP ▼ | TCP | 0 - 6553! | Cus... ▼ | Q | | Delete |

sg-0346b33f6c921cbd4 ✕

| sgr-0429ce2e0204e42c7 | SSH ▼ | TCP | 22 | Cus... ▼ | Q | | Delete |

0.0.0.0/0 ✕

| sgr-04e0d77e331e0ffd1 | All ICMP - IPv4 ▼ | ICMP | All | Cus... ▼ | Q | | Delete |

sg-05ccc38efd66854e ✕

CloudShell | Feedback | © 2024, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

---

aws | Services | Q Search | [Alt+S] | Ohio ▼ | Hem_Prakash_Dev ▼

- EC2 Dashboard  ✕
- EC2 Global View
- Events

▼ Instances
- Instances
- Instance Types
- Launch Templates
- Spot Requests
- Savings Plans
- Reserved Instances
- Dedicated Hosts
- Capacity Reservations

▼ Images
- AMIs
- AMI Catalog

▼ Elastic Block Store

✓ Inbound security group rules successfully modified on security group (sg-0346b33f6c921cbd4 | ElasticMapReduce-master)    ✕
▶ Details

EC2 > Security Groups > sg-0346b33f6c921cbd4 - ElasticMapReduce-master

# sg-0346b33f6c921cbd4 - ElasticMapReduce-master     Actions ▼

## Details

| Security group name | Security group ID | Description | VPC ID |
|---|---|---|---|
| ElasticMapReduce-master | sg-0346b33f6c921cbd4 | Master group for Elastic MapReduce created on 2024-07-09T19:49:55.783Z | vpc-08f2aeaae04f0a9eb ↗ |

| Owner | Inbound rules count | Outbound rules count | |
|---|---|---|---|
| 905418436402 | 8 Permission entries | 1 Permission entry | |

**Inbound rules**    Outbound rules    Tags

CloudShell | Feedback | © 2024, Amazon Web Services, Inc. or its affiliates. | Privacy | Terms | Cookie preferences

## Instance group settings Info

Edit cluster scaling option

| Cluster scaling option | Core |
|---|---|
| Manually set cluster size | Name and Maximum core nodes in the cluster |
| | Core  2 instances |

### Instance groups (2) Info

With the instance groups configuration, each node type consists of the same instance type and the same purchasing option for instances: On-Demand or Spot.

Resize instance group     Add task instance group

Q Find resource

< 1 >

| Node type ▲ | Name | ID | Status | Instances | Instance type | Purchasing option ▼ | Current |
|---|---|---|---|---|---|---|---|
| ○ Primary | Primary | ig-3RWQ96JY9LTFL | ⊘ Running | 1 | m5.xlarge | On-Demand | $0.192/ |
| ○ Core | Core | ig-IF3HLNNDJ95 | ⊘ Running | 2 | m5.xlarge | On-Demand | $0.192/ |

```python
from pyspark.sql import SparkSession


# Create SparkSession
spark = SparkSession.builder.getOrCreate()

# Specify the input file path
input_file = 's3://myemrbucket13/inputfolder/product_data.csv'

# Read CSV into a DataFrame
df = spark.read.option("header", "true").csv(input_file)


print(df.show())




# Write DataFrame as Parquet to the output folder
df.write.option("header", "true").mode("overwrite").parquet("s3://myemrbucket13/outputfolder")

# Stop the SparkSession
spark.stop()
```

## 3. Apply data transformation and cleansing techniques to prepare the data for aggregation and analysis.

Applying data transformation and cleansing techniques to prepare the data for aggregation and analysis.

- **Transformation**: Filter, select, join, and aggregate data as needed.
- **Cleansing**: Handle missing values, remove duplicates, and standardize formats.

```
1   from pyspark.sql import SparkSession
2   from pyspark.sql.functions import col, trim
3
4   # Initialize a SparkSession
5   spark = SparkSession.builder.appName("DataCleaningExample").getOrCreate()
6
7   # Sample data
8 ▾ data = [
9       (1, "  Alice  ", None),
10      (2, "Bob", "value2"),
11      (3, "Alice", "value3"),
12      (1, "  Alice  ", "value1"),
13      (4, None, "value4")
14  ]
15
16  # Column names
17  columns = ["id", "name", "column_name"]
18
19  # Create DataFrame
20  df = spark.createDataFrame(data, columns)
21
22  # Remove duplicates
23  df_cleaned = df.dropDuplicates()
24
25  # Handle missing values
26  df_cleaned = df_cleaned.na.fill({'column_name': 'default_value', 'name': 'unknown'})
27
28  # Standardize formats
29  df_cleaned = df_cleaned.withColumn('trimmed_column', trim(col('name')))
30
31  # Show the result
32  df_cleaned.show()
33
34  # Stop the SparkSession
35  spark.stop()
```

**Output of the sample taken**

```
+---+--------+-----------+--------------+
| id|    name|column_name|trimmed_column|
+---+--------+-----------+--------------+
|  1|  Alice |     value1|         Alice|
|  1|  Alice |       null|         Alice|
|  3|   Alice|     value3|         Alice|
|  2|     Bob|     value2|           Bob|
|  4| unknown|     value4|       unknown|
+---+--------+-----------+--------------+
```

4. **Implement data partitioning and indexing strategies to optimize query performance.**
   Implementing data partitioning and indexing strategies to optimize query performance.
   **Partitioning**: Partition data based on commonly queried fields.

```
1   from pyspark.sql import SparkSession
2   from pyspark.sql.functions import col, trim
3
4   # Initialize a SparkSession with Hadoop AWS package
5   spark = SparkSession.builder \
6       .appName("DataCleaningExample") \
7       .config("spark.hadoop.fs.s3a.impl", "org.apache.hadoop.fs.s3a.S3AFileSystem") \
8       .config("spark.hadoop.fs.s3a.aws.credentials.provider", "com.amazonaws.auth.DefaultAW
9       .getOrCreate()
10
11  # Sample data
12  data = [
13      (1, "  Alice  ", None, "2024-01-01"),
14      (2, "Bob", "value2", "2024-01-02"),
15      (3, "Alice", "value3", "2024-01-01"),
16      (1, "  Alice  ", "value1", "2024-01-03"),
17      (4, None, "value4", "2024-01-02")
18  ]
19
20  # Column names
21  columns = ["id", "name", "column_name", "partition_column"]
22
23  # Create DataFrame
24  df = spark.createDataFrame(data, columns)
25
26  # Remove duplicates
27  df_cleaned = df.dropDuplicates()
28
29  # Handle missing values
30  df_cleaned = df_cleaned.na.fill({'column_name': 'default_value', 'name': 'unknown'})
31
32  # Standardize formats
33  df_cleaned = df_cleaned.withColumn('trimmed_column', trim(col('name')))
34
35  # Show the result
36  df_cleaned.show()
37
38  # Write to S3 in JSON format, partitioned by 'partition_column'
39  df_cleaned.write.partitionBy('partition_column').json('s3://my-first-bucket')
40
41  # Stop the SparkSession
42  spark.stop()
+---+---------+------------+----------------+--------------+
| id|     name| column_name|partition_column|trimmed_column|
+---+---------+------------+----------------+--------------+
|  1|    Alice|      value1|      2024-01-03|         Alice|
|  1|    Alice|        null|      2024-01-01|         Alice|
|  3|    Alice|      value3|      2024-01-01|         Alice|
|  2|      Bob|      value2|      2024-01-02|           Bob|
|  4|  unknown|      value4|      2024-01-02|       unknown|
+---+---------+------------+----------------+--------------+
```

## 5. Update the GitHub repository with the code and configuration files for data ingestion and processing.

Updating the GitHub repository with the code and configuration files for data ingestion and processing.

1. **Initializing a Git Repository**:

   **git init**

2. **Adding and Commit Code**:

   **git add .**

   **git commit -m "Initial commit with data ingestion and processing scripts"**

3. **Pushing to GitHub:**

```
git remote add origin https://github.com/Hem-Prakash-Dev-Bharadwaj/data-engineering
git push -u origin main
```