# Social Media Data Integration & Analysis

*MADE BY- Ashish Chamel*

*COURSE- Extract, Transform, and Load (ETL)*

*DATE OF SUBMISSON-05/04/2025*

## Step-1

1. Creation of a input bucket in s3 ie etl-twitter-blog-ashish



## Step-2

1. Creation of two folders inside etl-twitter-blog-ashish bucket ie

   - blog-data and etl-social-media



2. In blog-data folder add sample_blogs.csv

3. In etl-social-media folder add sample_tweets.csv

**etl-social-media/**

| Objects | Properties |
|---------|-----------|

**Objects** (1)　　　　　　　　　　　　　　　　⟳　Copy S3 URI　｜　Copy URL　｜　⬇ Download　｜　Open ⧉　｜

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ⧉ to get a list of all objects in your bucket. For others to access your objects, you'll nee

| Q Find objects by prefix |
|---|

| ☐ | Name ▲ | Type | Last modified ▽ | Size |
|---|---|---|---|---|
| ☐ | 📄 sample_tweets.csv | csv | April 4, 2025, 19:08:36 (UTC+05:30) | |

# *Step-3*

1. Creation of a new output bucket named etl-cep-output-ashish

⊘ **Successfully created bucket "etl-cep-output-ashish"**
To upload files and folders, or to configure additional bucket settings, choose **View details**.　　　　　　　　　View d

▶ **Account snapshot** - *updated every 24 hours* `All AWS Regions`　　　　　　　　　　　　　　　View Storage Lens
Storage lens provides visibility into storage usage and activity trends. Metrics don't include directory buckets. Learn more ⧉

| General purpose buckets | Directory buckets |
|---|---|

**General purpose buckets** (3) `Info` `All AWS Regions`　　　　　　　　　⟳　Copy ARN　｜　Empty　｜　Delete　｜　Cre
Buckets are containers for data stored in S3.

| Q Find buckets by name | | | | ⟨ |
|---|---|---|---|---|

| ○ | Name ▲ | AWS Region | IAM Access Analyzer ▽ | Creation date |
|---|---|---|---|---|
| ○ | demobuckedrishi-16-03-2025 | US East (N. Virginia) us-east-1 | View analyzer for us-east-1 | March 16, 2025, 21:29:09 (UTC+05:30) |
| ○ | etl-cep-output-ashish | US East (N. Virginia) us-east-1 | View analyzer for us-east-1 | April 4, 2025, 19:15:21 (UTC+05:30) |
| ○ | etl-twitter-blog-ashish | US East (N. Virginia) us-east-1 | View analyzer for us-east-1 | April 4, 2025, 19:04:58 (UTC+05:30) |

# *Step-4*

1. Navigating to AWS GLUE and creating a new database ie. Social-media-data

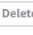**Databases** (1)　　　　　　　　　　　　　　　　Last updated (UTC)　⟳　Edit　｜　Delete
　　　　　　　　　　　　　　　　　　　　　　　　　April 4, 2025 at 13:47:41
A database is a set of associated table definitions, organized into a logical group.

| Q Filter databases |
|---|

| ☐ | Name ▲ | Description | Location URI ▽ | Created on (UTC) |
|---|---|---|---|---|
| ☐ | social-media-data | - | - | April 4, 2025 at 13:47:37 |

# *Step-5*

1. Classifier creation ie
   - Twitter_data
   - Blog_data

2. Creation of twitter_data classifier its settings are as follows
   - Classifier name as twitter_data
   - Classifier type and properties as CSV
   - CSV Serde – optional as None o Column delimiter as comma(,)
   - Quote symbol as Double-quote(")
   - Column headings as Has headings

## Create classifier Info

### Classifier details

**Classifier name**

| twitter_data |

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

### Classifier type and properties

**Classifier type**

| ○ Grok<br>Best for parsing unstructured text (e.g., application logs). | ○ XML<br>Extract data out of XML documents. |
| ○ JSON<br>Extract fields out of JSON files. | ● CSV<br>Filter and extract data out of CSV files. |

**CSV Serde - *optional***

| None ▼ |

Enter a CSV Serde option

**Column delimiter**

| Comma (,) ▼ |

Must be a single character. Use syntax like "001" or " " for special characters.

**Quote symbol**

| Double-quote (") ▼ |

Must be a single character and different than the column delimiter. Use syntax like "001" or " " for special characters.

**Column headings**

| Has headings ▼ |

| |

Enter a comma-delimited list.
Headings use the delimiter and quote symbol specified above.

**Processing options**

- ☐ Allow files with single column
- ☑ Disable whitespace trimming before identifying column values

3. Creation of a 2nd classifier ie blog_data and its settings are as follows-
   - Classifier name as blog_data
   - Classifier type and properties as CSV
   - CSV Serde – optional as None
   - Column delimiter as comma(,)
   - Quote symbol as Double-quote(")
   - Column headings as Has headings

## Create classifier Info

### Classifier details

**Classifier name**

```
blog_data
```

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

### Classifier type and properties

**Classifier type**

○ **Grok**
Best for parsing unstructured text (e.g., application logs).

○ **XML**
Extract data out of XML documents.

○ **JSON**
Extract fields out of JSON files.

● **CSV**
Filter and extract data out of CSV files.

**CSV Serde - optional**

```
None                                                    ▼
```
Enter a CSV Serde option

**Column delimiter**

```
Comma (,)                                               ▼
```
Must be a single character. Use syntax like "001" or " " for special characters.

**Quote symbol**

```
Double-quote (")                                        ▼
```
Must be a single character and different than the column delimiter. Use syntax like "001" or " " for special characters.

**Column headings**

```
Has headings                                            ▼
```
```

```
Enter a comma-delimited list.
Headings use the delimiter and quote symbol specified above.

**Processing options**

☐ Allow files with single column
☑ Disable whitespace trimming before identifying column values

Custom datatypes - optional | Info

### Classifiers

Classifiers are triggered during a crawl task. A classifier checks whether a given file is in a format the crawler can handle. If it is, the classifier creates a schema in the form of a StructType object that matches that data format.

Last updated (UTC)
April 4, 2025 at 13:52:41   ( Edit )  ( De

**Classifiers (2)** Info
View and manage all available classifiers.

```
🔍 Filter classifiers
```

| ☐ | Name | ▲ | Type | | Classification | ▽ | Last updated (UTC) |
|---|---|---|---|---|---|---|---|
| ☐ | blog_data | | CSV | | - | | April 4, 2025 at 13:52:40 |
| ☐ | twitter_data | | CSV | | - | | April 4, 2025 at 13:51:21 |

## *Step-6*

Navigated to the AWS console and configured an IAM role ie glue-role given administrator access.

**glue-role** Info                                                            ( Delete )
Allows Glue to call AWS services on your behalf.

### Summary                                                                    ( Edit )

**Creation date**                          **ARN**
April 04, 2025, 19:26 (UTC+05:30)          📋 arn:aws:iam::512253127681:role/glue-role

**Last activity**                          **Maximum session duration**
-                                          1 hour

| Permissions | Trust relationships | Tags | Last Accessed | Revoke sessions |
|---|---|---|---|---|

**Permissions policies (1)** Info            ( ⟳ ) ( Simulate 🗗 ) ( Remove ) ( Add permissions ▼ )
You can attach up to 10 managed policies.

Filter by Type
```
🔍 Search                                All types                    ▼
```
                                                                    ‹ 1 › ⚙

| ☐ | Policy name 🗗 | ▲ | Type | ▽ | Attached entities | ▽ |
|---|---|---|---|---|---|---|
| ☐ | ⊞ 🛡 AdministratorAccess | | AWS managed - job function | | 3 | |

## Step-7

1. In AWS GLUE creation of the 1st crawler as tweet-crawl where
   - IAM role is glue-role
   - Database is social-media data

2. Running the tweet-crawl successfully to get metadata



3. In AWS GLUE creation of a 2nd crawler as blog-crawler
   - IAM role is glue-role
   - Database is social-media data
   -
4. Running the blog-crawler successfully to get metadata



5. Metadata inside the social-media-data database as follows:

## *Step-8*

1. Navigating to AWS GLUE in ETL jobs using Visual ETL for a job creation.
2. Importing AWS GLUE DATA CATALOG 1, settings are as follows
   - Database- social-media-data
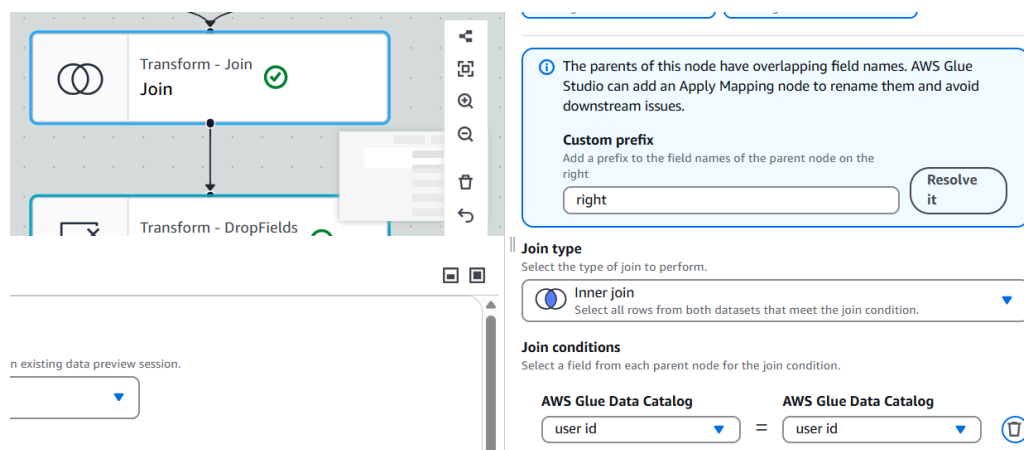   - Table- blog_data



3. Importing AWS GLUE DATA CATALOG 2, settings are as follows
   - Database- social-media-data
   - Table- etl_social_media



4. Importing Join from transforms and allotting parent nodes as AWS GLUE DATA CATALOG 1 and AWS GLUE DATA CATALOG 2 settings are as follows-
   - Custom prefix right
   - Join type – inner join
   - Join conditions- User id= User id

5. Importing Drop Fields from transforms where parent nodes are Join and its settings are as follows-
   - Drop fields as .userid



6. Importing Regex Extractor from Transforms choosing node parents as Drop Fields. Settings are as follows-
   - Column to extract from – tweet text
   - Regular expression #\(w+)
   - Extracted column as hashtags

7. Importing Aggregate from Transforms the settings are as follows-
   - Fields to group by as right_user id
   - Field to aggregate as tweet id and Aggregation function as count
   - In  Aggregate a column,
     Field to aggregate as timestamp and Aggregation function as min



8. Importing Amazon S3 from Data Target node parent as Aggregate.
   - S3 target location as – s3://etl-cep-output-ashish

9. Final outlook after the job creation
   - Job name-   etl-cep-output

etl-cep-output                                                          Last modified on 4/4/2025, 8:14:04 PM   [ Actions ▼ ]  [ Save ]  [ Run ]
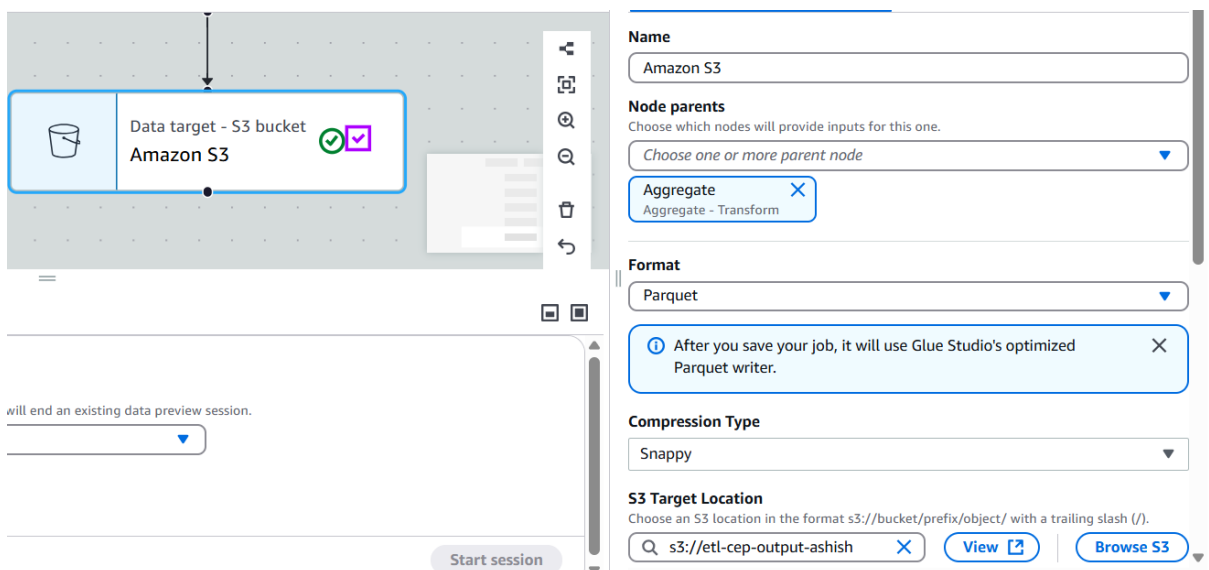
10. Job ran successfully

≡                                                                                                                                              ⓘ   ⊘

etl-cep-output                                                          Last modified on 4/4/2025, 8:40:53 PM   [ Actions ▼ ]  [ Save ]  [ Run ]

| Visual | Script | Job details | Runs | Data quality | Schedules | Version Control |

**Job runs** (1/1) Info                      Last updated (UTC)  ↻  [ View details ]  [ Stop job run ]  [ ✦ Troubleshoot with AI ]      [ Table View ] [ Card View ]
                                             April 4, 2025 at 15:12:31

🔍 Filter job runs by property                                                                                                       ⟨ 1 ⟩   ⚙

| | Run status ▽ | Retries ▽ | Start time (Local) ▽ | End time (Local) ▽ | Duration ▽ | Capacity (DPUs) ▽ | Worker type ▽ | Glue version ▽ |
|---|---|---|---|---|---|---|---|---|
| ● | ✓ Succeeded | 0 | 04/04/2025 20:40:57 | 04/04/2025 20:42:20 | 1 m 17 s | 10 DPUs | G.1X | 5.0 |

## *Step-9*

11. The output folder now has a file that is generated from the ETL job I performed

**etl-cep-output-ashish** Info

| Objects | Metadata | Properties | Permissions | Metrics | Management | Access Points |
|---|---|---|---|---|---|---|

**Objects (1)**    ⟳  Copy S3 URI   Copy URL   ⬇ Download   Open ⧉   Delete   Actions ▾   Create

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ⧉ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them perm

🔍 Find objects by prefix

| ☐ | Name ▲ | Type | Last modified | Size | ▼ | Storage class |
|---|---|---|---|---|---|---|
| ☐ | 📄 run-1743779521057-part-block-0-r-00003-snappy.parquet | parquet | April 4, 2025, 20:42:09 (UTC+05:30) | 769.0 B | | Standard |

12. Running the SQL query to obtain the final result

Amazon S3 > ... > run-1743779521057-part-block-0-r-00003-snappy.parquet > Query with S... ⓘ  ⤵  ⊘

## 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 1 record in 1187 ms
Bytes returned: 26 B

**Raw**    Formatted

```
1  123,1,2024-05-20T08:00:00
2  |
```

Raw   **Formatted**

⟨ 1 ⟩

| 123 | 1 | 2024-05-20T08:00:00 |
|---|---|---|

| User Id | Blog ID | Timestamp |
|---|---|---|
| 123 | 1 | 2024-05-20T08:00:00 |