**Guided Lab: Simplifying ETL Job Creation with AWS Glue Visual ETL**

**Description**

Welcome to the AWS Glue Visual ETL lab! This lab will walk you through creating and saving an ETL (Extract, Transform, Load) job using AWS Glue's visual interface with minimal coding effort. This lab is designed for beginners and experienced professionals to explore the advantages of visual data manipulation over traditional script-based methods and show you how straightforward it is to set up and configure an ETL process visually.

AWS Glue is a fully managed extract, transform, and load (ETL) service that makes it easy for customers to prepare and load their data for analytics. It provides a Python shell script that can be used to perform data transformations, including redacting sensitive data.

The visual ETL (Extract, Transform, Load) in AWS Glue refers to the graphical interface provided by AWS Glue Studio, which allows you to author and manage Glue jobs without writing code directly. AWS Glue Studio offers a drag-and-drop interface where you can define the flow of your data sources, transformations, and targets.

**Prerequisites**

This lab assumes you have experience creating AWS Glue Jobs using Python Shell and experience creating an Amazon S3 bucket and are familiar with its basic components.

If you find any gaps in your knowledge, consider taking the following labs:

- Creating an Amazon S3 bucket

- Introduction to AWS Glue Using Python Shell

**Objectives**

In this lab, you will:

- Simulate creating an ETL job using AWS Glue Visual ETL

- Understand the interface and tools provided by Glue Visual ETL.

- Apply redaction, and change schema transformations to a dataset.

**Subscribe to access AWS PlayCloud Labs**

**Lab Steps**

**Prepare your Environment**

**1. Log into the AWS Management Console**.

**2. Create an S3 Bucket**:

- Go to the S3 service in the AWS Console.

- Click **Create bucket**.

- Give your bucket a unique name.

- Leave the rest of the settings as default and click **Create bucket**

**3. Create a two folders**

- Navigate to your newly created bucket.

- Create two folders

  o **input**

  o **output**



**4. Download the CSV Data:**

https://media.tutorialsdojo.com/public/transactions.csv

*Take a look at the provided CSV file by opening the CSV file using Excel or any other spreadsheet application to visualize the data. When you are done, you can continue uploading the file.*

**5. Upload the file to the *input/* folder**



**Creating your ETL Job using Visual ETL**

1. Navigate to the **AWS Glue Console**.



2. You can either click on **ETL Jobs** on the left bar or click on the **Author and edit ETL jobs**

3. Click on **Visual ETL**



4. Then, you will be redirected to the Visual ETL interface

- If you notice in the **+ Add nodes,** there are tabs inside it:

-

o **Sources:**



- Sources are used to connect to the data sources, such as databases, data lakes, or other data stores. These allow you to read in data using various connectors provided by AWS Glue.

o **Transforms:**



- AWS Glue Studio offers a range of visual transforms, such as Concatenate, Split string, Pivot rows to columns, Lookup, and Derived column, which can be used to build more sophisticated data pipelines without writing code.

- **Targets:**



  - Targets are used to specify the destination for writing out the transformed data. These allow you to define the location and format for the output data, such as a database table, a file in an S3 bucket, or a data stream.

- **Popular:**



- The **Popular** tab is designed to help you quickly build ETL pipelines by providing easy access for commonly used nodes.

5. Now, let's **rename our job.**

- Click on the **Untitled job** in the upper left corner right beside the pencil icon.



- Name it uniquely



6. Next, click on the Job details tab

- Select the **PlayCloud-Sandbox** IAM role in the dropdown.



- Scroll down and look for **Requested number of workers.** Type 2 for number of workers. (optional step)
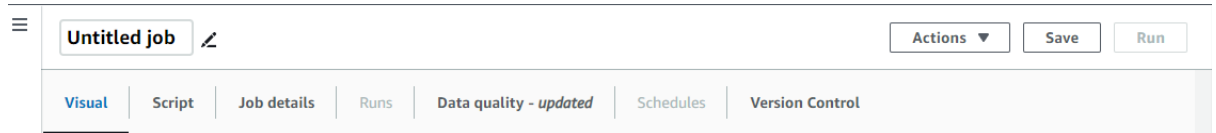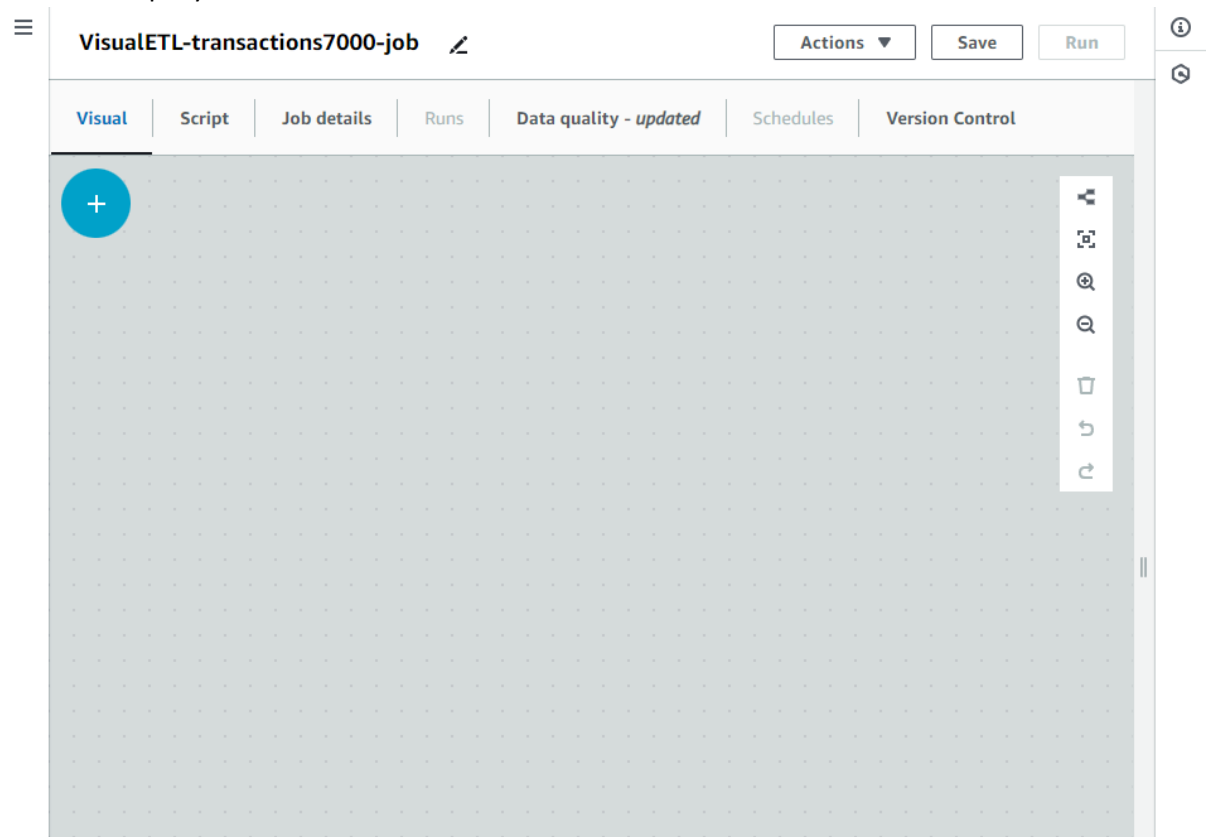  *NOTE: The **Number of Workers** determines the **DPU usage** of your job, depending on your worker type. The higher your DPU usage, the higher the bill you will pay. So be mindful of your DPU usage when doing AWS GLUE experiments in your AWS Account.*



7. Return to the previous tab (Visual) and click on the **blue circle with the+ sign** to add nodes.



8. Let's add the **Source node** first. In the **sources,** click on **Amazon S3.**

*Notice that after clicking on the source **Amazon S3**, a node appeared in our interface*

9. Click on the **Amazon S3 node** that just appeared.



- Paste the directory of your object or browse and select the bucket where you uploaded the transactions.csv file previously.

s3://{your-bucket-name}/input

*DO NOT FORGET TO CHANGE THE PLACEHOLDER, {bucket_name}*
*DO NOT FORGET TO CHANGE THE PLACEHOLDER, {bucket_name}*
*DO NOT FORGET TO CHANGE THE PLACEHOLDER, {bucket_name}*

**VisualETL-transaction700-job**  ✎

⚠ Job has not been saved      Actions ▾      Save      Run

**Data source properties - S3**

Name

Amazon S3

S3 source type   Info

◉ S3 location
Choose a file or folder in an S3 bucket.

○ Data Catalog table

S3 URL

🔍 s3://visualetl-transaction7000/input            ✕      View ⬀      Browse S3

☑ Recursive
Read files in all subdirectories.

- For the **Data Format**, select **CSV.**



Leave the rest as default.

*If you encounter an error similar to the example below, you can safely ignore it. You will still be able to complete this lab even if a 'failed to infer schema' error occurs.*

Failed to infer schema from s3://sample1234567890/input/transactions (1).csv
{ "message": "Unauthorized access for account id                , "code": "AccessDeniedException", "[Message]": "See error.Message for details.", "time": "2024-10-03T11:43:03.987Z", "requestId": "db221a3f-fe5f-4a6
1-8ca1-8c54c94d9c31", "statusCode": 403, "retryable": false }

- If you wait for a while, you will see the information inside the CSV in your bucket in the Data preview.



10. **Applying Detect Sensitive Data Transformations**

In this lab, we are gonna use Detect Sensitive Data Transformations first. This transformation scans the dataset for sensitive data like personal identifiable information (PII), then replace it or add a column with information on what was found and where.

- Click on the **+ sign** again (**Add Node**) and navigate to the **Transforms tab**

- Select **Detect Sensitive Data**



- Click on the **Detect Sensitive Data Transforms node**

Follow the following Configurations

- o Node parents: *Amazon S3*

- **Detect sensitive data:** Select **Find columns that contain sensitive data.**

**Detect sensitive data** Info

Scan the dataset for sensitive data like personal identifiable information (PII), then replace it or add a column with information on what was found and where.

○ Find sensitive data in each row
Scan the entire data set, and act on each cell individually.

● Find columns that contain sensitive data
Scan a sample of the dataset to quickly find columns that represent sensitive information.

**Sample portion**
The percentage of rows to sample out of the entire data set.

| 100 | % |

Between 0 and 100.

**Detection threshold**
To consider a field as containing PII, set the minimum percentage of detected rows out of the sampled rows.

| 10 | % |

Between 0 and 100.

- **Types of sensitive information to detect:** Select **Select specific patterns**

**Types of sensitive information to detect**
Select the types of sensitive information you would like to detect. For example, email or credit card number.

○ Include all available types (256)
This will select all types available at job authoring time.

○ Select categories
This will dynamically include all patterns in categories you select.

● Select specific patterns
Only patterns you explicitly select will be detected.

- **Select patterns: Search for Credit Card and select it**

Types of sensitive information to detect

Select the types of sensitive information you would like to detect. For example, email or credit card number.

○ Include all available types (256)

This will select all types available at job authoring time.

○ Select categories

This will dynamically include all patterns in categories you select.

● Select specific patterns

Only patterns you explicitly select will be detected.

Selected patterns

| Q Credit Card | ✕ | Browse |
| --- | --- | --- |

Use: "Credit Card"

**Credit Card**
Universal

- Scroll down and look for:
  **Select global action (required):** Select **REDACT. Redact detected text**

Select global action (required)

Choose an action to take on detected entities.

○ DETECT. Output detection results

Output a frame with detection information for each column in the data.

● REDACT. Redact detected text

Replace detected entity with a string you choose.

○ SHA256_HASH. Apply cryptographic hash.

Apply a SHA-256 cryptographic hash function to the input string.

Complete settings for global action

If unset, will be the default values.

Redaction Text:   ******

- In the **Data Preview,** you will see that the CardNumber Column has been redacted.



11. **Applying Change Schema Transformations**
**Next, let's do a simple Change Schema transformation.**

- Click on the **+ sign** again (**Add Node**) and navigate to the **Transforms tab**

- Select **Change Schema**



- Click on the **Change Schema Transforms node**
  Follow the following Configurations

  o  Node parents:  **Detect Sensitive Data**

- o **Change Schema (Apply mapping)**:
  - ▪ Under the Source key, look for **Amount** and change it to **TransactionAmount** under the Target Key



- You will see the changes made in the **Data preview**



12. Lastly, lets **Add the Target Node** for our ETL Job.

- Click on the **+ sign** again (**Add Node**) and navigate to the **Targets tab**

- Select **Amazon S3**



- Click on the Amazon S3 **Target node**
  Follow the following Configurations

- o Node parents: **Change Schema**

**Data target properties - S3** 1

Name

Amazon S3

Node parents

Choose which nodes will provide inputs for this one.

Choose one or more parent node ▲

Q Filter parent nodes

☐ **Data sources**

☐ Amazon S3
S3 - DataSource

■ **Transforms**

☐ Detect Sensitive Data
PIIDetection - Transform

☑ Change Schema
ApplyMapping - Transform

☐ **Unclassified nodes**

- o Format: **Parquet**
  *Parquet is an open-source, column-oriented data file format that is designed for efficient data storage and retrieval.*

- o *compression Type: **GZIP***

- o **S3 Target Location:**

s3://{your-bucket-name}/output/

*DO NOT FORGET TO CHANGE THE PLACEHOLDER, {bucket_name}*
*DO NOT FORGET TO CHANGE THE PLACEHOLDER, {bucket_name}*

*DO NOT FORGET TO CHANGE THE PLACEHOLDER, {bucket_name}*

Format

Parquet ▼

ⓘ After you save your job, it will ✕
use Glue Studio's optimized
Parquet writer.

Compression Type

GZIP ▼

S3 Target Location
Choose an S3 location in the format
s3://bucket/prefix/object/ with a trailing slash (/).

🔍 s3://visualetl-transaction7000/outpu ✕
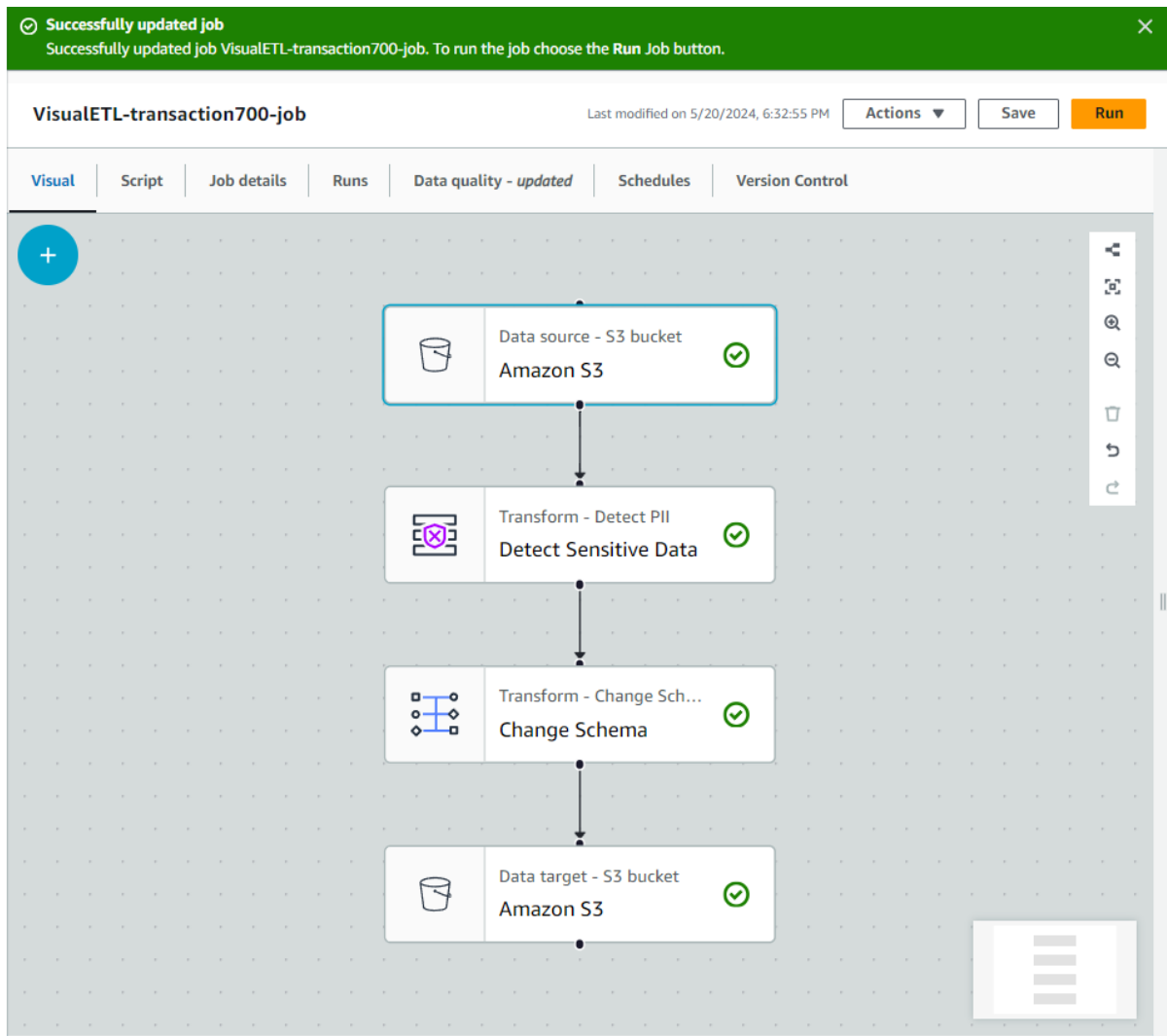
View ⬈ | Browse S3

- 
    o Leave the rest as Default

13. Click on **Save** in the upper right corner of the interface to save your first Visual ETL job.

*NOTE: After saving the Job, you won't be able to run it. This lab simulates creating a visual ETL job and is intended for demonstration purposes only.*

That's it congratulations on completing this lab! Through this hands-on experience, you've set up data sources, applied sensitive data transformations, and adjusted data schemas using AWS Glue's visual interface. The skills you've honed today are invaluable for efficiently managing and manipulating large datasets in real-world scenarios.

This lab has demonstrated how AWS Glue's Visual ETL can streamline your data transformation workflows, making the process more intuitive and less error-prone compared to hard coding your ETL job. As you continue to explore AWS Glue, leverage these advantages to enhance your data integration and analytics capabilities further.

One last thing! It is a good practice to clean up the resources created during this lab. Not only will it make you a better professional, but you will also become a more organized person. Happy learning!