

# CSC Project

Name: Vineel Kumar

ID:2000031044

Title: Prevention of data leakage in various cloud computing platforms and detection of sensitive data(Data Analysis)

YOUTUBE LINK: [https://youtu.be/UBQ1We\\_j9DE](https://youtu.be/UBQ1We_j9DE)

## Introduction:

The company's Information security depends on employees by learning the rules through training and awareness-building sessions. And the data leakage is mainly happening via insider employees or other hacker which they try to steal or use it to their own or selling it This uncontrolled data leakage puts business in a vulnerable position. Once this data is no longer within the domain, then the company is at serious risk The main concept is water marking. Digital watermarking is the process of embedding information into a digital signal, such as an image, video, or audio file, in a way that is imperceptible to the human senses but can be detected and extracted by special software or hardware. Watermarking can serve various purposes, such as identifying the owner or distributor of the content, tracking its usage and distribution, protecting against copyright infringement, and detecting tampering or unauthorized modifications

## Services we will be using

Amazon S3: Amazon S3 is an object storage service that provides manufacturing scalability, data availability, security, and performance.

AWS IAM: This is nothing but identity and access management which enables us to manage access to AWS services and resources securely.

QuickSight: Amazon QuickSight is a scalable, serverless, embeddable, machine learning-powered business intelligence (BI) service built for the cloud.

AWS Glue: A serverless data integration service that makes it easy to discover, prepare, and combine data for analytics, machine learning, and application development.

AWS Lambda: Lambda is a computing service that allows programmers to run code without creating or managing servers.

## Reference link

<https://github.com/darshilparmar/dataengineering-youtube-analysis-project>

## Dataset Used

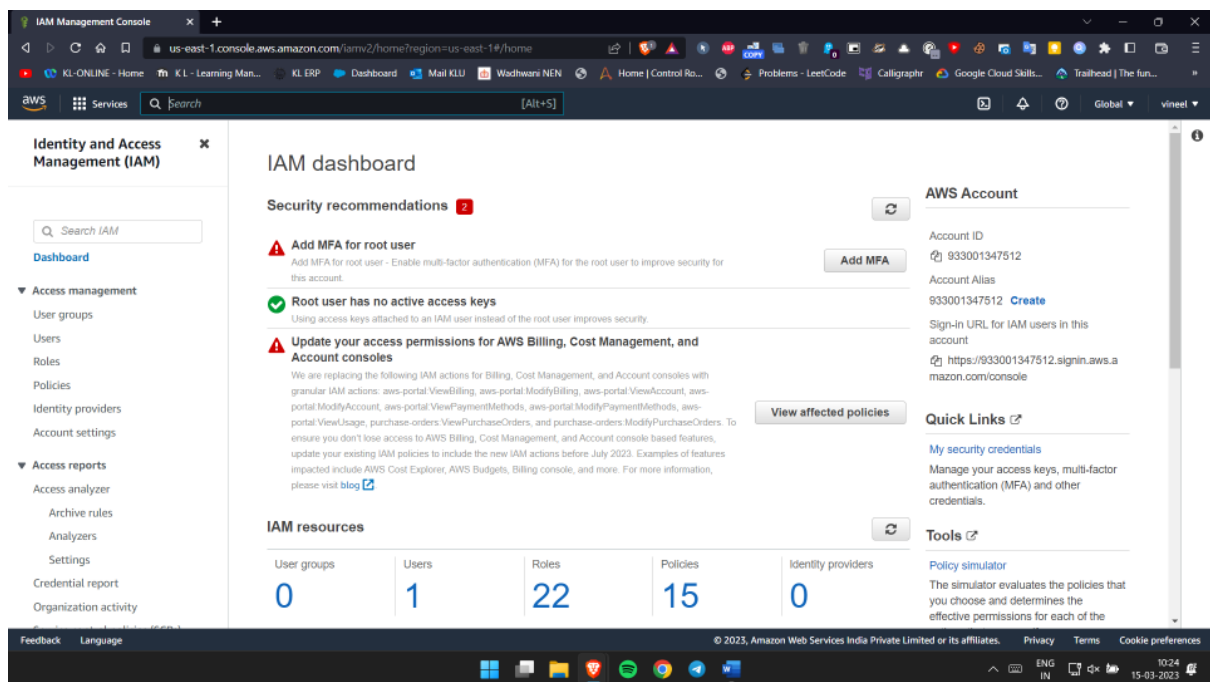
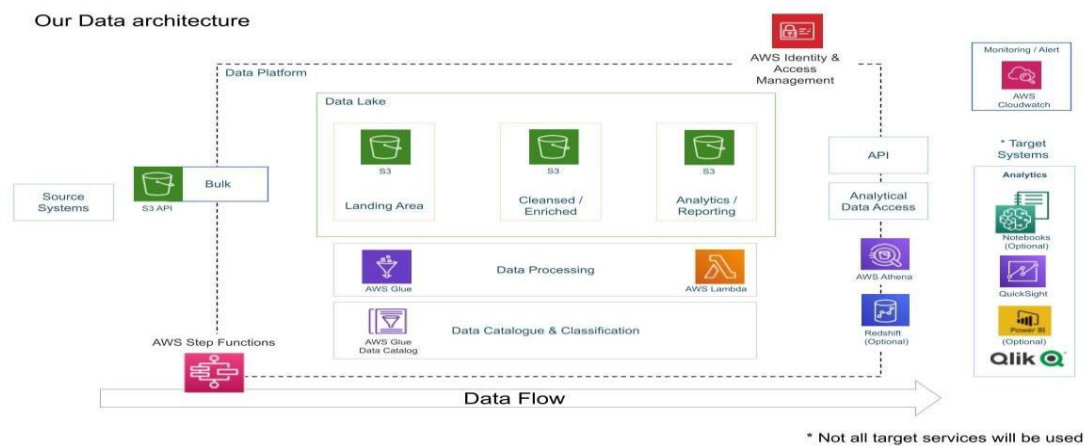
This Kaggle dataset contains statistics (CSV files) on daily popular YouTube videos over the course of many months. There are up to 200 trending videos published every day for many locations. The data for each region is in its own file. The video title, channel title, publication time, tags, views, likes and

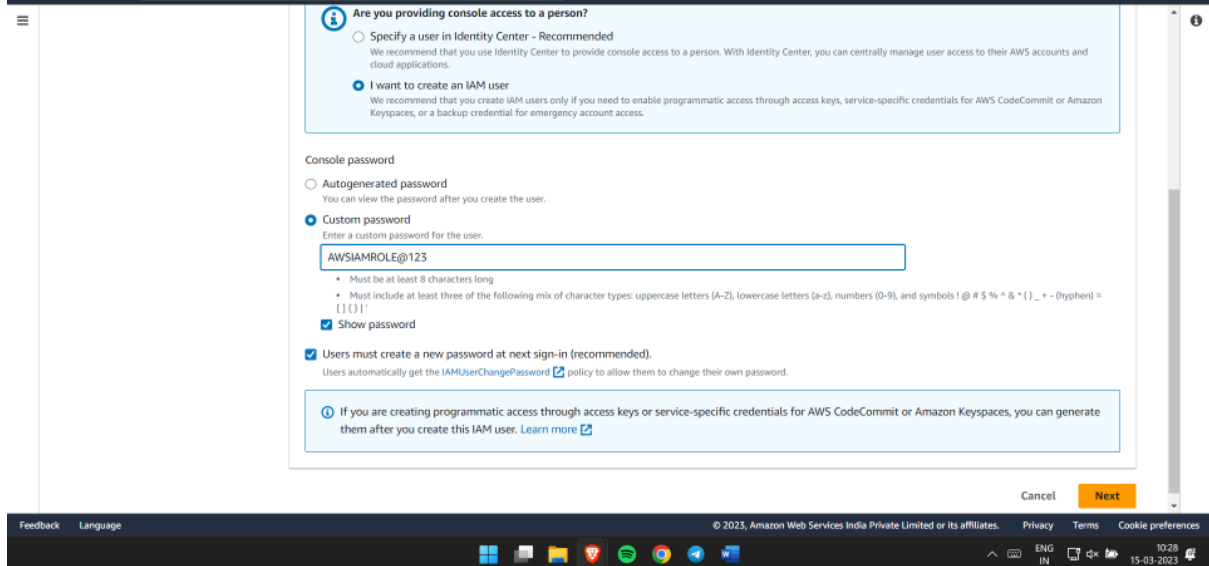
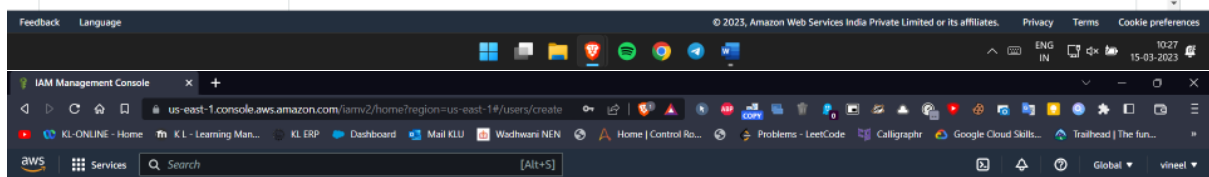
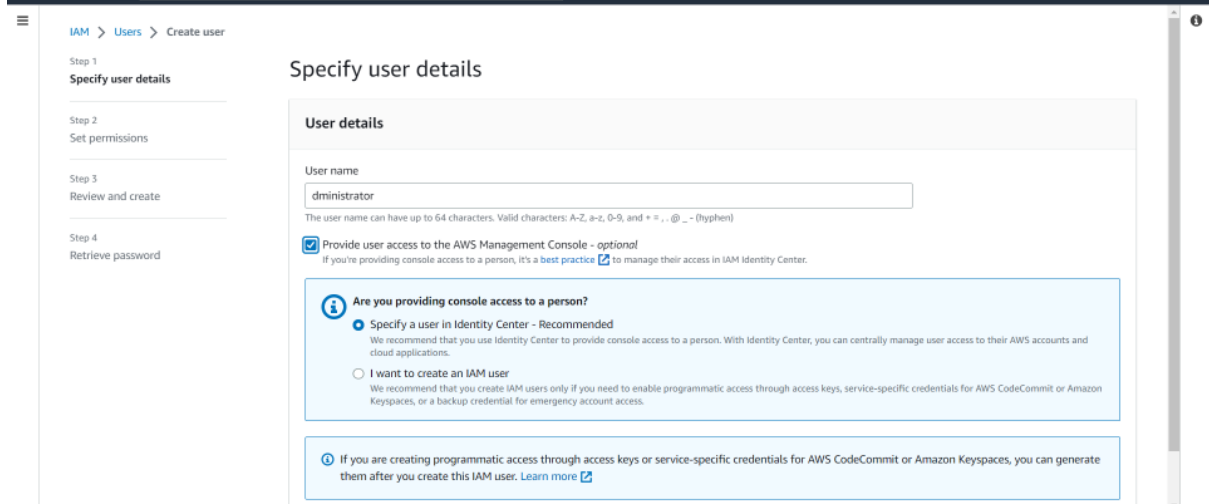
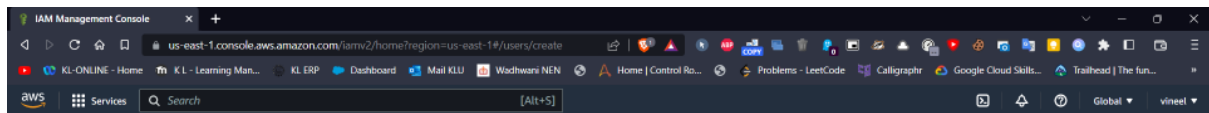
dislikes, description, and comment count are among the items included in the data. A category\_id field, which differs by area, is also included in the JSON file linked to the region.

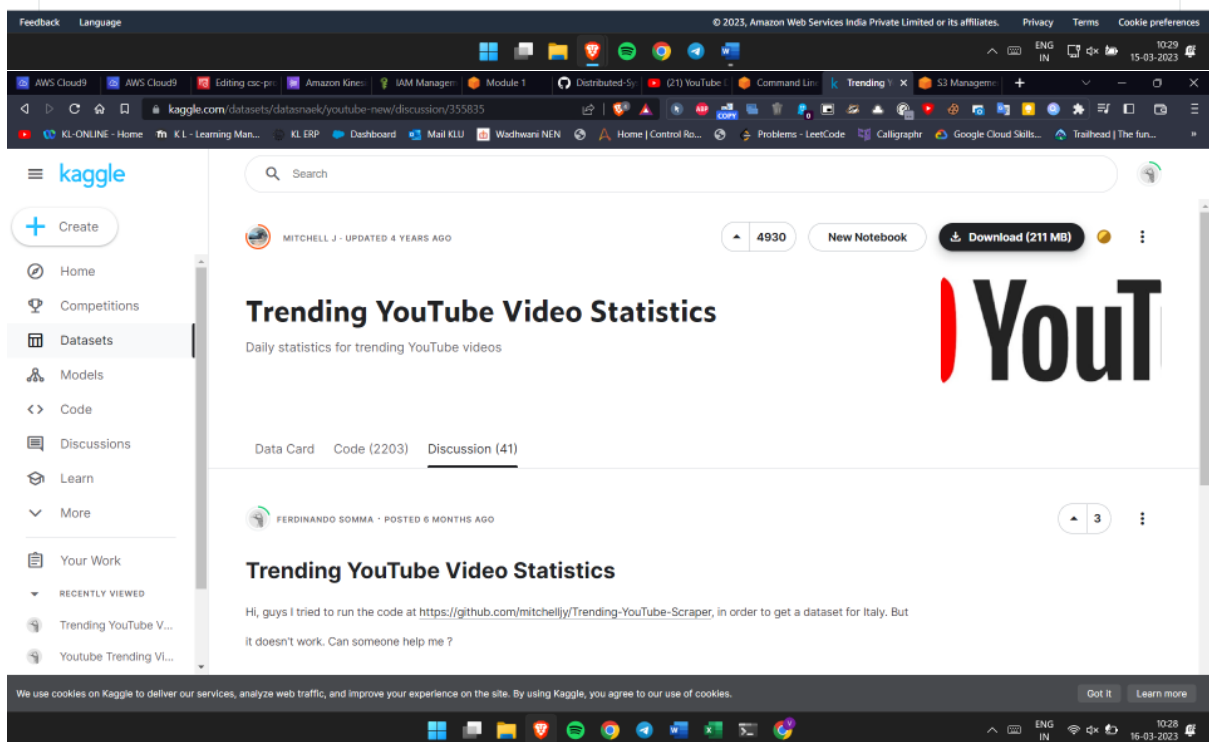
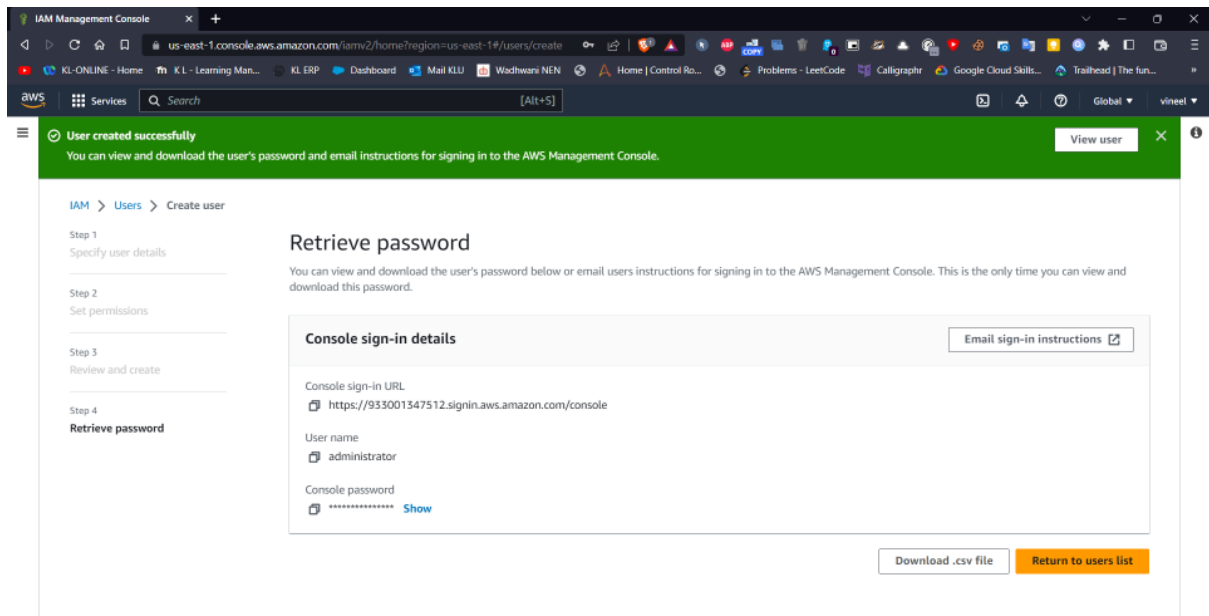
<https://www.kaggle.com/datasets/datasnaek/youtube-new>

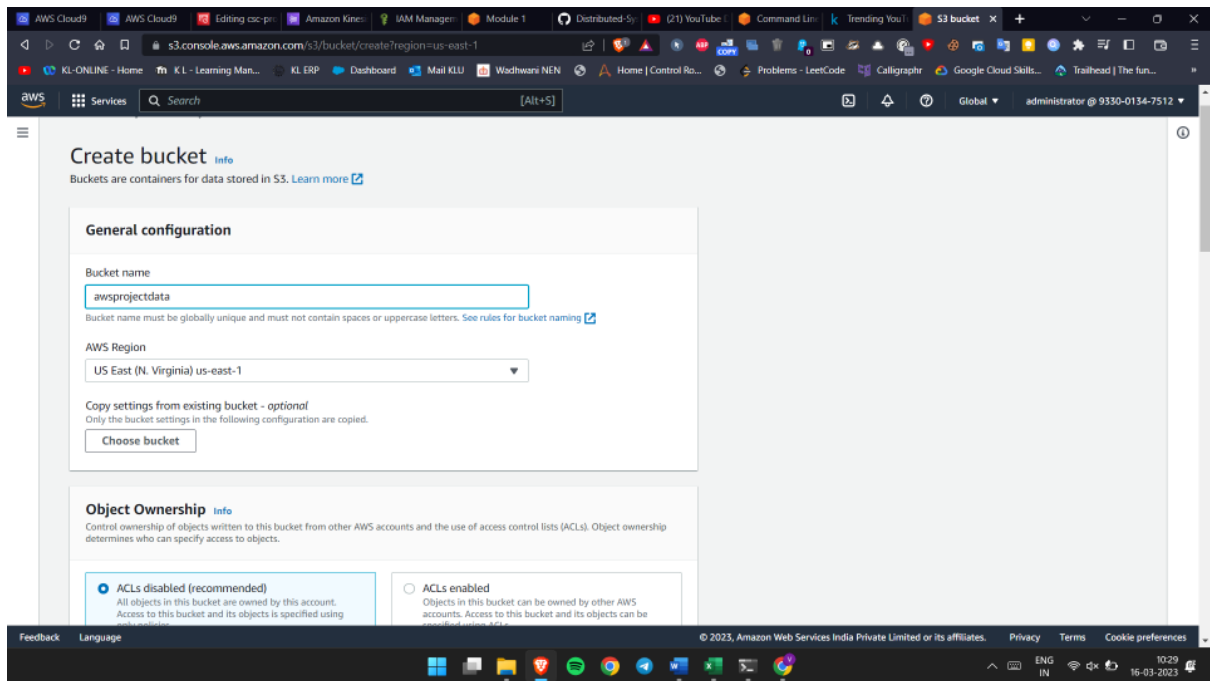
LinkedIn : [https://www.linkedin.com/posts/vineel-kumar-vukoti\\_youtube-data-analysis-activity-7052557651016126464-](https://www.linkedin.com/posts/vineel-kumar-vukoti_youtube-data-analysis-activity-7052557651016126464-7052557651016126464-)

ZwoO?utm\_source=li\_share&utm\_content=feedcontent&utm\_medium=g\_dt\_web&utm\_campaign=copy









Command Prompt - aws s3 c

```
'ls' is not recognized as an internal or external command,
operable program or batch file.

C:\Users\ASUS\Downloads\archive>aws s3 cp . s3://awsprojectdatast/youtube/raw_statistics_reference_data/ --recursive --e
xclude "*" --include "*.json"
upload: \JP_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/JP_category_id.json
upload: \RU_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/RU_category_id.json
upload: \CA_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/CA_category_id.json
upload: \WR_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/WR_category_id.json
upload: \FR_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/FR_category_id.json
upload: \DE_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/DE_category_id.json
upload: \GB_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/GB_category_id.json
upload: \US_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/US_category_id.json
upload: \MX_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/MX_category_id.json
upload: \IN_category_id.json to s3://awsprojectdatast/youtube/raw_statistics_reference_data/IN_category_id.json

C:\Users\ASUS\Downloads\archive>aws s3 cp CAvideos.csv s3://awsprojectdatast/youtube/raw_statistics/region=ca/
upload: .\CAvideos.csv to s3://awsprojectdatast/youtube/raw_statistics/region=ca/CAvideos.csv

C:\Users\ASUS\Downloads\archive>aws s3 cp DEvideos.csv s3://awsprojectdatast/youtube/raw_statistics/region=de/
upload: .\DEvideos.csv to s3://awsprojectdatast/youtube/raw_statistics/region=de/DEvideos.csv

C:\Users\ASUS\Downloads\archive>aws s3 cp FRvideos.csv s3://awsprojectdatast/youtube/raw_statistics/region=fr/
upload: .\FRvideos.csv to s3://awsprojectdatast/youtube/raw_statistics/region=fr/FRvideos.csv

C:\Users\ASUS\Downloads\archive>aws s3 cp GBvideos.csv s3://awsprojectdatast/youtube/raw_statistics/region=gb/
upload: .\GBvideos.csv to s3://awsprojectdatast/youtube/raw_statistics/region=gb/GBvideos.csv

C:\Users\ASUS\Downloads\archive>aws s3 cp INvideos.csv s3://awsprojectdatast/youtube/raw_statistics/region=in/
Completed 56.8 MiB/56.8 MiB (7.8 MiB/s) with 1 file(s) remaining
```

YouTube Data Analysis | END TO END DATA ENGINEERING PROJECT

Darshil Parmar 49.8K subscribers

78,859 views 28 Mar 2022 #dataengineer #project #darshil

DATA ENGINEERING | 13:46

BUILD THIS! Twitter Data Pipeline using Airflow for Beginners | Data...

Darshil Parmar 56K views · 1 month ago

us-east-1 console.aws.amazon.com/athena/home?region=us-east-1#/query-editor

Amazon Athena > Query editor

Editor Recent queries Saved queries Settings Workgroup primary

Data

Data source: AwsDataCatalog

Database: awsprojectdatabase

Tables and views: Create Filter tables and views

Tables (2)

- raw\_statistics
- raw\_statistics\_reference\_data

Views (0)

Query 1: X Query 2: X

```
1 SELECT * FROM "awsprojectdatabase"."raw_statistics" limit 10;
```

SQL Ln 1, Col 1

Run again Explain Cancel Clear Create

Query results Query stats

Reuse query results \*Athena engine version 3 only

© 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor...

Completed Time in queue: 107 ms Run time: 1.255 sec Data scanned: 2.79 MB

Results (10) [Copy](#) [Download results](#)

Search rows

#	video_id	trending_date	title
1	n1WpP7iowLc	17.14.11	"Eminem - Walk On Water (Audio) ft. Beyoncé"
2	OdBikQ4Mz1M	17.14.11	"PLUSH - Bad Unboxing Fan Mail"
3	5qpjK5DgCt4	17.14.11	"Racist Superman   Rudy Mancuso"
4	d380meDOWOM	17.14.11	"I Dare You: GOING BALD?!"
5	2Vv-BFVoq4g	17.14.11	"Ed Sheeran - Perfect (Official Music Video)"
6	OyIWz1XEyc	17.14.11	"Jake Paul Says Alissa Violet CHEATED with LOGAN PAUL! #DramaAlert Team 10 vs Martinez Twins!"
7	_uM5kFkhB8	17.14.11	"Vanoss Superhero School - New Students"
8	2ky565vSYSE	17.14.11	"WE WANT TO TALK ABOUT OUR MARRIAGE"
9	JzCsM1vtn78	17.14.11	"THE LOGAN MADE HISTORY. LOL. AGAIN."
10	43sm-QwLcx4	17.14.11	"Finally Sheldon is winning an argument about the existence of God"

Feedback Language © 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

us-east-1.console.aws.amazon.com/gluestudio/home?region=us-east-1#/jobs

AWS Glue Studio > Jobs

### Jobs [Info](#)

**Create job [Info](#)** [Create](#)

☒ **Visual with a source and target**  
Start with a source, ApplyMapping transform, and target.

☐ **Visual with a blank canvas**  
Author using an interactive visual interface.

☐ **Spark script editor**  
Write or upload your own Spark code.

☐ **Python Shell script editor**  
Write or upload your own Python shell script.

☐ **Jupyter Notebook**  
Write your own code in a Jupyter Notebook for interactive development.

☐ **Ray script editor [New](#)**  
Write your own code to run on Ray.

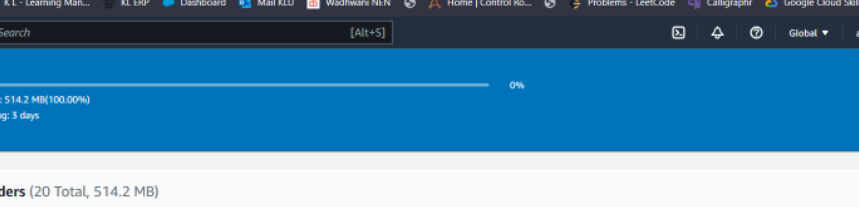
Source: **Amazon S3**  
JSON, CSV, or Parquet files stored in S3.

Target: **Amazon S3**  
S3 bucket by specifying a bucket path as the data target.

**Your jobs (0) [Info](#)** [Refresh](#) [Actions](#) [Run job](#)

Filter jobs

Job name	Type	Last modified	AWS Glue version
----------	------	---------------	------------------



**Uploading** 0%

Total remaining: 19 files: 514.2 MB(100.00%)  
 Estimated time remaining: 3 days  
 Transfer rate: 2.1 KB/s

**Files and folders (20 Total, 514.2 MB)**

Name	Folder	Type	Size	Status	Error
CA_category_id.json	-	application/json	7.7 KB	Pending	-
CAvideos.csv	-	text/csv	61.1 MB	Pending	-
DE_category_id.json	-	application/json	7.7 KB	Pending	-
DEvideos.csv	-	text/csv	60.1 MB	Pending	-
FR_category_id.json	-	application/json	7.7 KB	Pending	-
FRvideos.csv	-	text/csv	49.0 MB	Pending	-
GB_category_id.json	-	application/json	8.0 KB	Pending	-
GBvideos.csv	-	text/csv	50.7 MB	Pending	-
IN_category_id.json	-	application/json	8.0 KB	Pending	-
INvideos.csv	-	text/csv	56.8 MB	Pending	-



**AWS Glue** ×

Getting started  
ETL jobs  
Visual ETL  
Notebooks  
Job run monitoring  
Data Catalog tables  
Data connections  
Workflows (orchestration)

▼ **Data Catalog**  
Databases  
Tables  
Stream schema registries  
Schemas  
Connections  
**Crawlers**  
Classifiers  
Catalog settings

► Data Integration and ETL  
► Legacy pages

What's New

## Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

**Crawlers (0)** [Info](#) Last updated (UTC) March 16, 2023 at 05:09:24 Action Run Create crawler

View and manage all available crawlers.

	Name	State	Schedule	Last run	Last run times...	Log	Table changes fr...
No resources							
No resources to display.							

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

10:39 16-03-2023

**Select trusted entity**

Step 2  
Add permissions

Step 3  
Name, review, and create

### Trusted entity type

☒ **AWS service**  
Allow AWS services like EC2, Lambda, or others to perform actions in this account.

☐ **AWS account**  
Allow entities in other AWS accounts belonging to you or a 3rd party to perform actions in this account.

☐ **Web identity**  
Allows users federated by the specified external web identity provider to assume this role to perform actions in this account.

☐ **SAML 2.0 federation**  
Allow users federated with SAML 2.0 from a corporate directory to perform actions in this account.

☐ **Custom trust policy**  
Create a custom trust policy to enable others to perform actions in this account.

### Use case

Allow an AWS service like EC2, Lambda, or others to perform actions in this account.

**Common use cases**

☐ **EC2**  
Allows EC2 instances to call AWS services on your behalf.

☐ **Lambda**  
Allows Lambda functions to call AWS services on your behalf.

**Use cases for other AWS services:**

☒ **Glue**  
Allows Glue to call AWS services on your behalf.

us-east-1.console.aws.amazon.com/iamv2/home?region=us-east-1#/roles/create...

Step 2  
Add permissions

Step 3  
Name, review, and create

Permissions policies (Selected 1/832) Info

Choose one or more policies to attach to your new role.

Filter policies by property or policy name and press enter. 11 matches

\*s3\* Clear filters

	Policy name	Type	Description
<input type="checkbox"/>	AWSLambdaS3Execu...	Custom...	
<input type="checkbox"/>	AWSLambdaS3Execu...	Custom...	
<input type="checkbox"/>	AmazonDMSRedshi...	AWS m...	Provides access to manage S3 settings for Redshift endpoints for DMS.
<input checked="" type="checkbox"/>	AmazonS3FullAccess	AWS m...	Provides full access to all buckets via the AWS Management Console.
<input type="checkbox"/>	QuickSightAccessF...	AWS m...	Policy used by QuickSight team to access customer data produced by S3 Storage Management Analy...
<input type="checkbox"/>	AmazonS3ReadOnl...	AWS m...	Provides read only access to all buckets via the AWS Management Console.
<input type="checkbox"/>	AmazonS3Outposts...	AWS m...	Provides full access to Amazon S3 on Outposts via the AWS Management Console.
<input type="checkbox"/>	AWSBackupService...	AWS m...	Policy containing permissions necessary for AWS Backup to backup data in any S3 bucket. This inclu...
<input type="checkbox"/>	AWSBackupService...	AWS m...	Policy containing permissions necessary for AWS Backup to restore a S3 backup to a bucket. This inc...
<input type="checkbox"/>	AmazonS3ObjectLa...	AWS m...	Provides AWS Lambda functions permissions to interact with Amazon S3 Object Lambda. Also grants ...
<input type="checkbox"/>	AmazonS3Outposts...	AWS m...	Provides read only access to Amazon S3 on Outposts via the AWS Management Console.

us-east-1.console.aws.amazon.com/iamv2/home?region=us-east-1#/roles/details...

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

ENG IN 10:51 16-03-2023

<input type="checkbox"/>	AWSGlueServiceNotebookRole	AWS managed	Policy for AWS Glue service role which allows customer to manage notebook server
<input checked="" type="checkbox"/>	AWSGlueServiceRole	AWS managed	Policy for AWS Glue service role which allows access to related services including EC...
<input type="checkbox"/>	AWSGlueConsoleSageMakerNotebookFullAccess	AWS managed	Provides full access to AWS Glue via the AWS Management Console and access to s...
<input type="checkbox"/>	AWSGlueConsoleFullAccess	AWS managed	Provides full access to AWS Glue via the AWS Management Console
<input type="checkbox"/>	AwsGlueDataBrewFullAccessPolicy	AWS managed	Provides full access to AWS Glue DataBrew via the AWS Management Console. Also ...
<input type="checkbox"/>	AWSGlueSchemaRegistryReadOnlyAccess	AWS managed	Provides readonly access to the AWS Glue Schema Registry Service
<input type="checkbox"/>	AwsGlueSessionUserRestrictedNotebookPolicy	AWS managed	Provides permissions that allows users to create and use only the notebook sessions t...
<input type="checkbox"/>	AmazonSageMakerServiceCatalogProductsGlueServiceRolePolicy	AWS managed	Service role policy used by the AWS Glue within the AWS ServiceCatalog provisioned ...
<input type="checkbox"/>	AwsGlueSessionUserRestrictedNotebookServiceRole	AWS managed	Provides full access to all AWS Glue resources except for sessions. Allows users to cr...
<input type="checkbox"/>	AWSGlueSchemaRegistryFullAccess	AWS managed	Provides full access to the AWS Glue Schema Registry Service
<input type="checkbox"/>	AWSGlueDataBrewServiceRole	AWS managed	This policy grants permission to glue to perform action on user's glue data catalog, this...
<input type="checkbox"/>	AwsGlueSessionUserRestrictedPolicy	AWS managed	Provides permissions that allows users to create and use only the interactive sessions ...
<input type="checkbox"/>	AwsGlueSessionUserRestrictedServiceRole	AWS managed	Provides full access to all AWS Glue resources except for sessions. Allows users to cr...

Cancel Add permissions

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

ENG IN 10:53 16-03-2023

us-east-1.console.aws.amazon.com/glue/home?region=us-east-1#/v2/data-catal...

**AWS Glue**

Getting started  
ETL jobs  
Visual ETL  
Notebooks  
Job run monitoring  
Data Catalog tables  
Data connections  
Workflows (orchestration)

► Data Catalog  
► Data Integration and ETL  
► Legacy pages

What's New  
Documentation  
AWS Marketplace

Enable compact mode  
Enable new navigation

**Create a database**  
Create a database in the AWS Glue Data Catalog.

**Database details**

Name  
awsprojectdatabase  
Database name is required, in lowercase characters, and no longer than 255 characters.

Location - optional  
Set the URI location for use by clients of the Data Catalog.

Description - optional  
Enter text  
Descriptions can be up to 2048 characters long.

Cancel Create database

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

us-east-1.console.aws.amazon.com/glue/home?region=us-east-1#/v2/data-catal...

**AWS Glue**

Getting started  
ETL jobs  
Visual ETL  
Notebooks  
Job run monitoring  
Data Catalog tables  
Data connections  
Workflows (orchestration)

▼ Data Catalog  
Databases  
Tables  
Stream schema registries  
Schemas  
Connections  
Crawlers  
Classifiers  
Catalog settings

► Data Integration and ETL  
► Legacy pages

What's New

**One crawler successfully created**  
The following crawler is now created: "awsdatacrawler"

**awsdatacrawler**  
Last updated (UTC)  
March 16, 2023 at 05:27:41

Run crawler Edit Delete

**Crawler properties**

Name awsdatacrawler	IAM role cscprojectrole	Database awsprojectdatabase	State READY
Description -	Security configuration -	Lake Formation configuration -	Table prefix -
Maximum table threshold -			

► Advanced settings

Crawler runs Schedule Data sources Classifiers Tags

**Crawler runs (0)**  
The list of crawler runs for this crawler.

Stop run View CloudWatch logs View run details

Amazon Athena

Before you run your first query, you need to set up a query result location in Amazon S3. [Edit settings](#)

**Query editor**

Query editor [New](#)

Notebook editor [New](#)

Notebook explorer [New](#)

**Jobs**

Workflows  
Powered by Step Functions

**Administration**

Workgroups

Data sources

☐ Turn on compact mode

**Data**

Data source: AwsDataCatalog

Database: awsprojectdatabase

Tables and views [Create](#)

**Tables (20)**

- ca\_category\_id\_json
- cavideos\_csv
- de\_category\_id\_json
- devideos\_csv
- fr\_category\_id\_json
- frvideos\_csv

**Query 1**

```
1 SELECT * FROM "AwsDataCatalog"."awsprojectdatabase" limit 10;
```

SQL Ln 1, Col 62

[Run](#) [Explain](#) [Cancel](#) [Clear](#) [Create](#)

☐ Reuse query results  
\*Athena engine version 3 only

**Query results** [Query stats](#)

**Results** [Copy](#) [Download results](#)

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. [Privacy](#) [Terms](#) [Cookie preferences](#)

Lambda > Functions > Create function

## Create function [Info](#)

AWS Serverless Application Repository applications have moved to [Create application](#).

☒ **Author from scratch**  
Start with a simple Hello World example.

☐ **Use a blueprint**  
Build a Lambda application from sample code and configuration presets for common use cases.

☐ **Container image**  
Select a container image to deploy for your function.

**Basic information**

**Function name**  
Enter a name that describes the purpose of your function.  
  
Use only letters, numbers, hyphens, or underscores with no spaces.

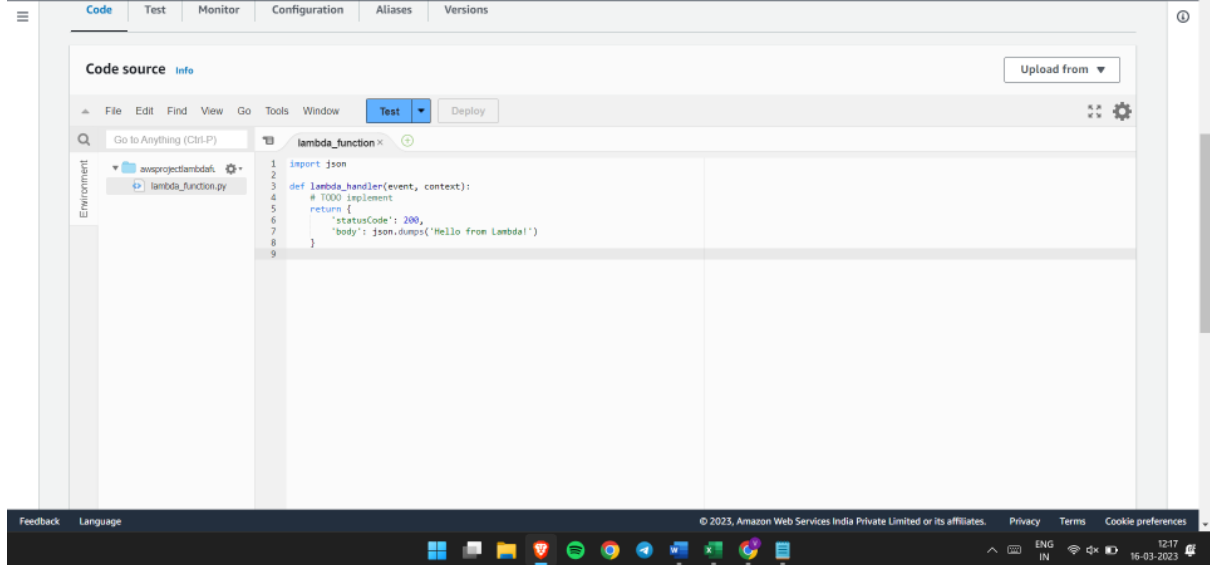
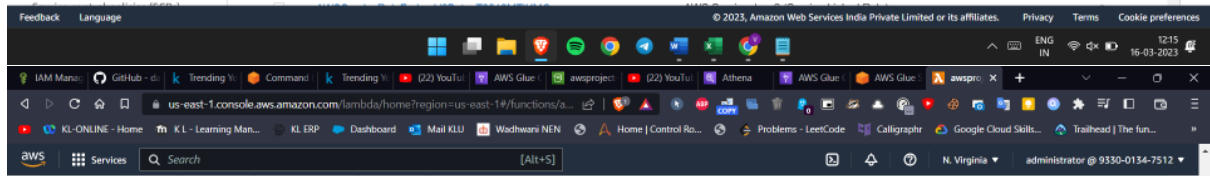
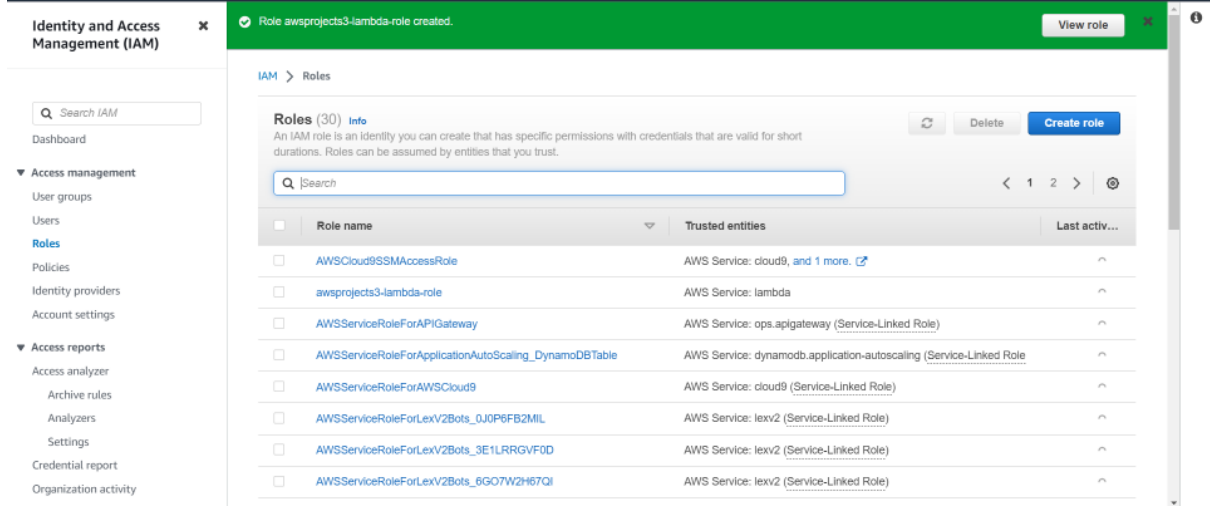
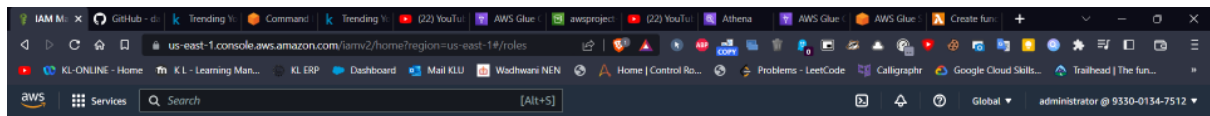
**Runtime [Info](#)**  
Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

**Architecture [Info](#)**  
Choose the instruction set architecture you want for your function code.  
☒ x86\_64  
☐ arm64

**Permissions [Info](#)**

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. [Privacy](#) [Terms](#) [Cookie preferences](#)



us-east-1.console.aws.amazon.com/s3/bucket/create?region=us-east-1

Services Search [Alt+S]

Amazon S3 > Buckets > Create bucket

Create bucket [Info](#)

Buckets are containers for data stored in S3. [Learn more](#)

General configuration

Bucket name

aws-cleaned-project-layer

Bucket name must be globally unique and must not contain spaces or uppercase letters. See rules for bucket naming

AWS Region

US East (N. Virginia) us-east-1

Copy settings from existing bucket - optional

Only the bucket settings in the following configuration are copied.

Choose bucket

Object Ownership [Info](#)

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

☒ ACLs disabled (recommended)

☐ ACLs enabled

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

12:23 16-03-2023

us-east-1.console.aws.amazon.com/lambda/home?region=us-east-1#/functions/...

Services Search [Alt+S]

Lambda > Functions > awsprojectlambdafunction > Edit environment variables

Edit environment variables

Environment variables

You can define environment variables as key-value pairs that are accessible from your function code. These are useful to store configuration settings without the need to change function code. [Learn more](#)

Key	Value	
glue_catalog_db_name	db_youtube_cleaned	Remove
glue_catalog_table_name	cleaned_statistics_reference_data	Remove
s3_cleaned_layer	s3://aws-cleaned-project-layer	Remove
write_data_operation	append	Remove

Add environment variable

Encryption configuration

Feedback Language

© 2023, Amazon Web Services India Private Limited or its affiliates. Privacy Terms Cookie preferences

12:23 16-03-2023

