# Data Lake Proof of Concept – Detailed Guide

By Carlos Cruz Mejia

## Step by Step

### 1. Data Sources – MySQL

First of all, we need to set up our data source. We will proceed with a MySQL database with a standard configuration for this example. Make sure to select the free tier to keep costs as low as possible.



Then we must select the free tier template.

We keep all default configurations in the instance identifier and set our master password.

In the storage section, we set the minimum allowed to 20 GiB and make sure that the "Enable storage autoscaling" box is unchecked.



In the connectivity section, we use all default options and allow public access to make the process of connecting to the database easier.

We select the default VPC and leave everything else as it is.



We set the" Database authentication" with the password option. We leave everything else with the default configuration. We click on "create database" and wait for it to finish. Launching a database with our previously defined settings will take a few minutes.

When the database is ready, we must create a new parameter group to run Sakila scripts successfully.



In our new parameter group, we set **log_bin_trust_function_creators** to "1" and save the changes.

Then we modify our existing database and set the parameter group to our custom one.



We click on the "Modify DB" instance and check on "Apply immediately."

Then we proceed to reboot our instance to apply our latest change.



Once the reboot is over, we click on our database to get more details and the endpoints to set the connection. (Don't worry, this DB will no longer exist by the time you're reading this.)

With the information we got, we can establish a connection using MySQL Workbench, DBeaver, etc. We run the Sakila schema script first and then the Sakila data to finish setting up our data source.



We will have a database full of data, which we can use for our data lake.

## 2. Data Lake layers – S3

With the setup for our data source done, we can proceed to work on the layers. We will classify and store the data. We need to create three different layers for each stage of data.

## Create bucket Info

Buckets are containers for data stored in S3. Learn more ↗

### General configuration

**Bucket name**

data-lake-coal-layer

Bucket name must be 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**

We create another bucket for the resources we will need in the future. We will end up with a setup like this.

| | | | |
|---|---|---|---|
| ○ | data-lake-resources | US East (N. Virginia) us-east-1 | Bucket and objects not public |
| ○ | data-lake-pressure-layer | US East (N. Virginia) us-east-1 | Bucket and objects not public |
| ○ | data-lake-diamond-layer | US East (N. Virginia) us-east-1 | Bucket and objects not public |
| ○ | data-lake-coal-layer | US East (N. Virginia) us-east-1 | Bucket and objects not public |

## 3.  Credentials - Secret Manager

We set the name of our secret, and if our database is in RDS, we only need to select it and add the password. We can also store credentials for other databases.



In the end, our secret will look like this.

## 4.   Permissions – IAM

We set up the required IAM role for our lambda functions.

### Select trusted entity

**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
  Allow 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:

Choose a service to view use case ▼

Cancel    Next

### Name, review, and create

**Role details**

Role name
Enter a meaningful name to identify this role.

data-lake-lambda-etl

Maximum 64 characters. Use alphanumeric and '+=,.@-_' characters.

Description
Add a short explanation for this policy.

Allows Lambda functions to call AWS services on your behalf.

Maximum 1000 characters. Use alphanumeric and '+=,.@-_' characters.

Step 1: Select trusted entities                    Edit

```
1  {
2      "Version": "2012-10-17",
3      "Statement": [
4          {
5              "Effect": "Allow",
6              "Action": [
7                  "sts:AssumeRole"
8              ],
9              "Principal": {
10                 "Service": [
11                     "lambda.amazonaws.com"
12                 ]
13             }
14         }
15     ]
16 }
```

| | | | | |
|---|---|---|---|---|
| ☐ | ⊞ 🔷 SecretsManagerReadWrite | | AWS managed | Provides read/write access to AWS Secrets Mana... |
| ☐ | ⊞ 🔷 AmazonS3FullAccess | | AWS managed | Provides full access to all buckets via the AWS M... |
| ☐ | ⊞ 🔷 AmazonAthenaFullAccess | | AWS managed | Provide full access to Amazon Athena and scope... |
| ☐ | ⊞ 🔷 AWSGlueServiceRole | | AWS managed | Policy for AWS Glue service role which allows ac... |
| ☐ | ⊞ 🔷 AWSLambdaExecute | | AWS managed | Provides Put, Get access to S3 and full access to... |

## 5.   Setting up Glue

We create a database for our data source and each layer in the data lake.

### Add database

**Database name**

sakila_coal

▾ Description and location (optional)

**Location** ℹ

s3://data-lake-coal-layer/sakila_coal/

**Description**

Enter description...

---

| Add database | View tables | Action ▾ | | Showing: 1 - 4 ⟨ ⟩ ⟳ ❓ |
|---|---|---|---|---|
| ☐ **Name** | | **Description** | | |
| ☐ sakila_coal | | | | |
| ☐ sakila_diamond | | | | |
| ☐ sakila_pressure | | | | |

## 6.   Querying Data – Amazon Athena

If it is the first time we set up Athena, we need to finish the initial configuration.

Analytics

# Amazon Athena
## Start querying data instantly.

Amazon Athena is an interactive query service that makes it easy to analyze data in Amazon S3 and other federated data sources using standard SQL.

**Begin querying your data**

Discover the query editor and start querying right away.

**Explore the query editor**

We need to set up a location for the query results to be stored on. We will use this one:



We set the path for our primary workspace.

## 7.   Data Processing – Lambda

Then we create a new lambda function for each layer. We use the settings specified in the article.

**Basic information**

Function name
Enter a name that describes the purpose of your function.

sakila-coal-layer

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.

Python 3.8

Architecture  Info
Choose the instruction set architecture you want for your function code.
● x86_64
○ arm64

Permissions  Info
By default, Lambda will create an execution role with permissions to upload logs to Amazon CloudWatch Logs. You can customize this default role later when adding triggers.

▼ Change default execution role

Execution role
Choose a role that defines the permissions of your function. To create a custom role, go to the IAM console.
○ Create a new role with basic Lambda permissions
● Use an existing role
○ Create a new role from AWS policy templates

Existing role
Choose an existing role that you've created to be used with this Lambda function. The role must have permission to upload logs to Amazon CloudWatch Logs.

data-lake-lambda-etl

View the data-lake-lambda-etl role on the IAM console.

**Layers**  Info                                                                            Edit    Add a layer

| Merge order | Name | Layer version | Compatible runtimes | Compatible architectures | Version ARN |
|---|---|---|---|---|---|
| | | | There is no data to display. | | |

It is vital to add AWS Data wrangler with a layer provided by AWS.

Finally, we set a good amount of memory and maximum execution time so we can run our lambda functions. Remember to add the code from the repository and update it with the values from your data lake.

## Basic settings  Info

**Description - *optional***

[                                                    ]

**Memory  Info**
Your function is allocated CPU proportional to the memory configured.

[ 128                    ]  MB

Set memory to between 128 MB and 10240 MB

**Ephemeral storage  Info**
You can configure up to 10 GB of ephemeral storage (/tmp) for your function. View pricing ⬈

[ 512                    ]  MB

Set ephemeral storage (/tmp) to between 512 MB and 10240 MB.

**Timeout**

[ 5     ]  min  [ 0     ]  sec

**Execution role**
Choose a role that defines the permissions of your function. To create a custom role, go to the IAM console.

⦿ Use an existing role

◯ Create a new role from AWS policy templates

**Existing role**
Choose an existing role that you've created to be used with this Lambda function. The role must have permission to upload logs to Amazon CloudWatch Logs.

[ data-lake-lambda-etl                          ▼ ]   [ ⟳ ]

View the data-lake-lambda-etl role on the IAM console.

Cancel      **Save**

## 8. Business Intelligence – Quicksight

We need to sign up for Quicksight. We will make use of the 30-day trial for this project and set the following values;

Use your own buckets for this step.

## Select Amazon S3 buckets ✕

| **S3 Buckets Linked To QuickSight Account** | **S3 Buckets You Can Access Across AWS** |
|---|---|

Select the buckets that you want QuickSight to be able to access.

Selected buckets have read only permissions by default. However, you must give write permissions for Athena Workgroup feature.

☑ Select all

| S3 Bucket | Write permission for Athena Workgroup |
|---|---|
| ☐ | ☐ |
| ☐ | ☐ |
| ☑ data-lake-coal-layer | ☐ |
| ☑ data-lake-diamond-layer | ☐ |
| ☑ data-lake-melting-layer | ☐ |
| ☐ data-lake-resources | ☐ |
| ☐ | ☐ |

Cancel                                                         Finish

Check the following services. Make sure your account has enough permissions. This step creates an IAM role for Quicksight.

## QuickSight access to AWS services

Make your existing AWS data and users available in QuickSight. Learn more

**Allow access and autodiscovery for these resources**

- ☑ Amazon Redshift
- ☑ Amazon RDS
- ☑ IAM
- ☑ Amazon S3 (3 buckets selected)
  - Select S3 buckets
- ☑ Amazon Athena
  - Make sure you've chosen the right Amazon S3 buckets for QuickSight access
- ☐ Amazon S3 Storage Analytics
- ☐ AWS IoT Analytics
- ☐ Amazon OpenSearch Service
- ☐ Amazon SageMaker
- ☐ Amazon Timestream

**Finish**

---

**QuickSight**      DevCharles ⌄

Find analyses & more

- ★ Favorites
- 🕐 Recent
- 📊 Dashboards
- 📈 Analyses
- 🗄 Datasets
- 💬 Community  New

### Datasets

**New dataset**

| Name | | Owner | Last Modified ⌄ | |
|------|---|-------|-----------------|---|
| Web and Social Media Analytics | SPICE | Me | a few seconds ago | ... |
| Sales Pipeline | SPICE | Me | a few seconds ago | ... |
| People Overview | SPICE | Me | a few seconds ago | ... |
| Business Review | SPICE | Me | a few seconds ago | ... |

Now we will create a data source using Athena. That dataset will get information from our diamond layer.

We click on "Edit/Preview data" so we can make sure the data is loaded correctly.

## Choose your table

×

athena-primary

**Catalog: contain sets of databases.**

AwsDataCatalog ⌄

**Database: contain sets of tables.**

sakila_diamond ⌄

**Tables: contain the data you can visualize.**

◉ film_sales

Edit/Preview data    Use custom SQL    Select

Now we set a name for our dataset and click on "Save & visualize".

And with that, we are done. We can create visualizations for our data.