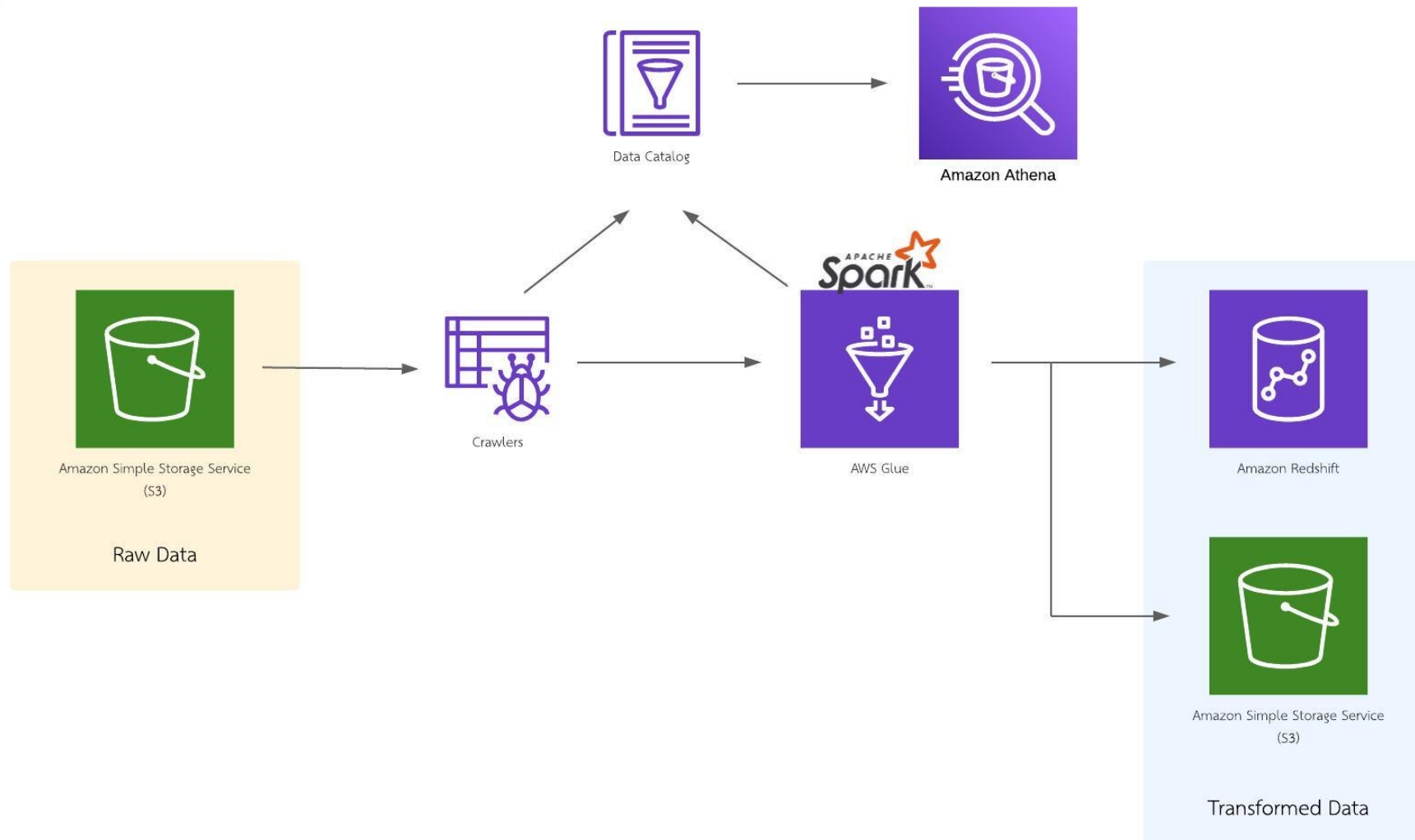


ETL Data Pipeline Development for Car-Sharing Data

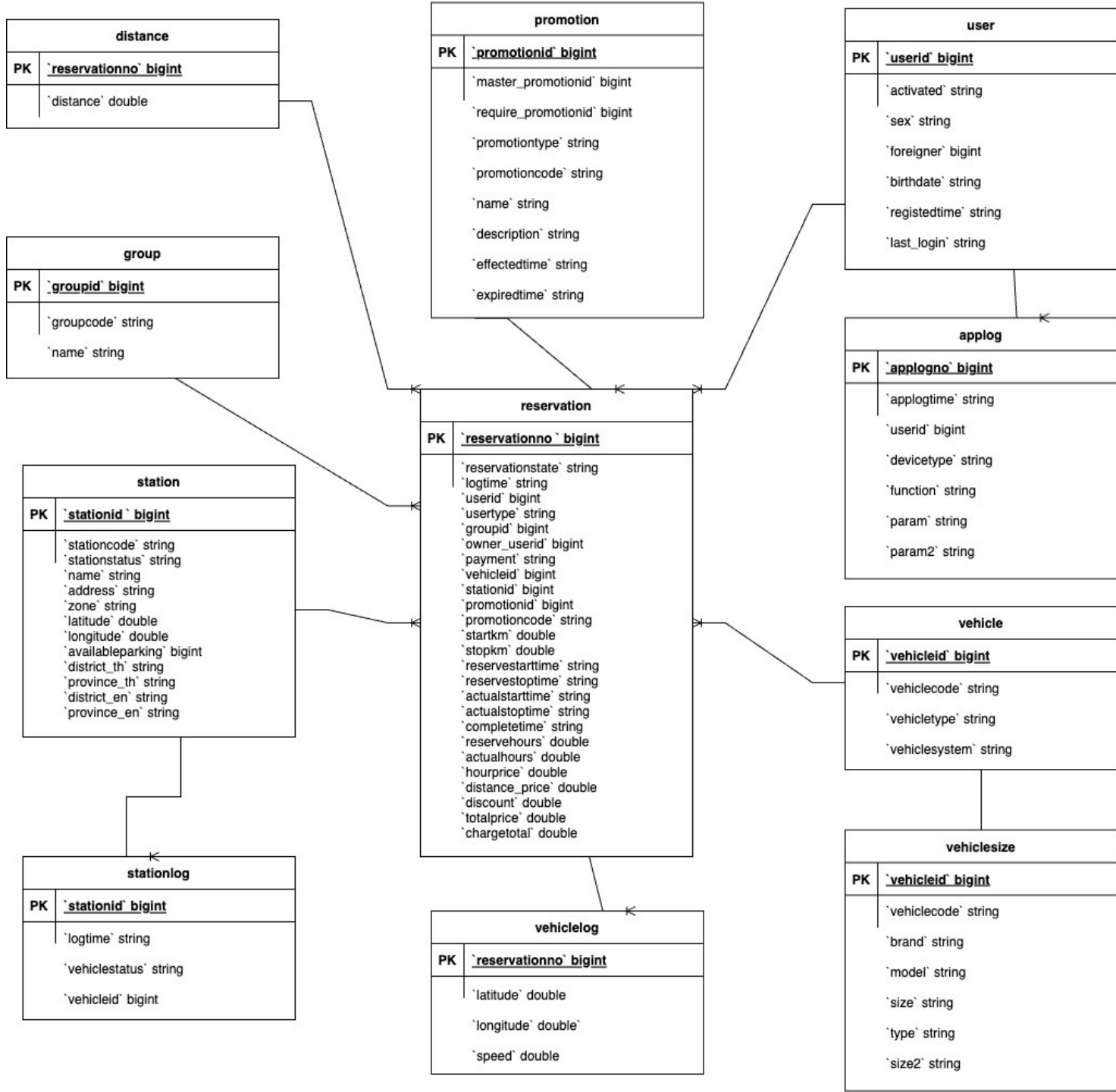
Designed and built a robust ETL data pipeline to transfer raw car-sharing data from the data lake (Amazon S3) to the data warehouse (Amazon Redshift) in the AWS cloud platform

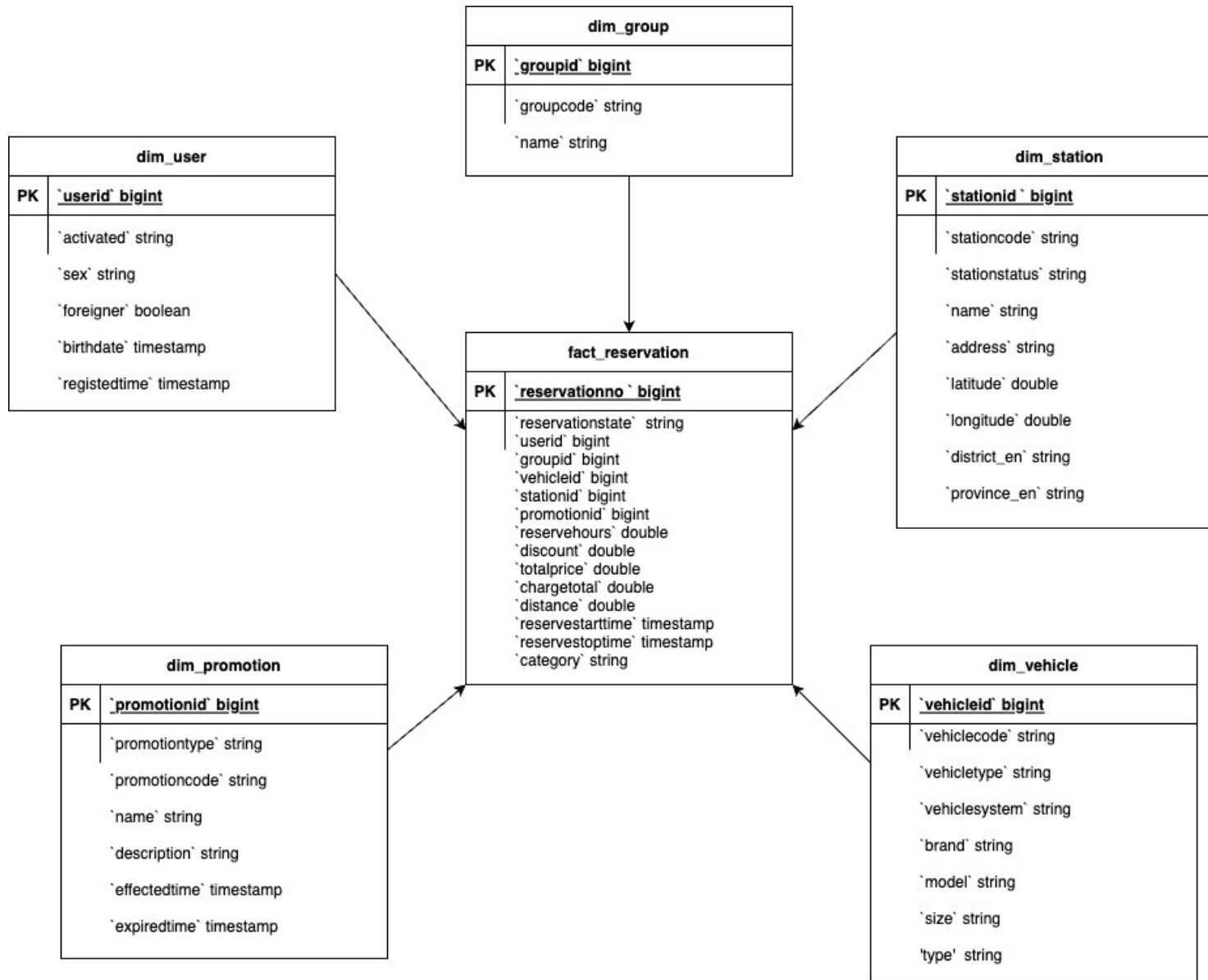


<https://github.com/bmbenz/ETL-Data-Pipeline-Development-for-Car-Sharing-Data.git>



ER Diagram





Star Schema

1. Create a bucket and upload data

Amazon S3 > Buckets > car-sharing-bucket > data/

data/

Copy S3 URI

ObjectsProperties

Objects (10)

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 permissions. [Learn more](#)

Refresh

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Find objects by prefix

< 1 > ⚙

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	applog/	Folder	-	-	-
<input type="checkbox"/>	distance/	Folder	-	-	-
<input type="checkbox"/>	group/	Folder	-	-	-
<input type="checkbox"/>	promotion/	Folder	-	-	-
<input type="checkbox"/>	reservation/	Folder	-	-	-
<input type="checkbox"/>	station/	Folder	-	-	-
<input type="checkbox"/>	user/	Folder	-	-	-
<input type="checkbox"/>	veehiclelog/	Folder	-	-	-
<input type="checkbox"/>	vehicle/	Folder	-	-	-
<input type="checkbox"/>	vehiclesize/	Folder	-	-	-

CloudShellFeedbackLanguage

© 2023, Amazon Web Services, Inc. or its affiliates. PrivacyTermsCookie preferences

2. Create data catalog using crawler

AWS Glue > Crawlers

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 (8) Info
View and manage all available crawlers.

Last updated (UTC)
August 1, 2023 at 10:07:43

↻

Action ▾

Run

Create crawler

Q Filter crawlers

< 1 > ⚙

<input type="checkbox"/>	Name ▾	State ▾	Schedule	Last run ▾	Last run tim... ▾	Log	Table changes ...
<input type="checkbox"/>	distance_crawler	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created
<input type="checkbox"/>	group_crawler	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created
<input type="checkbox"/>	promotion_craw...	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created
<input type="checkbox"/>	reservation_cra...	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created
<input type="checkbox"/>	station_crawler	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created
<input type="checkbox"/>	user_crawler	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created
<input type="checkbox"/>	vehicle_crawler	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created
<input type="checkbox"/>	vehiclesize_craw...	✔ Ready		✔ Succeeded	July 28, 2023 at...	View log	1 created

Permissions policies (10) Info
You can attach up to 10 managed policies.

↻

Simulate

Remove

Add permissions ▾

Q Filter policies by property or policy name and press enter.

< 1 > ⚙

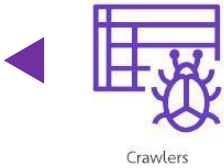
<input type="checkbox"/>	Policy name ▾	Type ▾	Description
<input type="checkbox"/>	AwsGluePassRolePolicy	Customer managed	
<input type="checkbox"/>	PowerUserAccess	AWS managed - job function	Provides full access t
<input type="checkbox"/>	AmazonRedshiftFullAccess	AWS managed	Provides full access t
<input type="checkbox"/>	AmazonS3FullAccess	AWS managed	Provides full access t
<input type="checkbox"/>	AWSGlueServiceRole	AWS managed	Policy for AWS Glue
<input type="checkbox"/>	AWSGlueConsoleFullAccess	AWS managed	Provides full access t
<input type="checkbox"/>	AWSGlueServiceNotebookRole	AWS managed	Policy for AWS Glue
<input type="checkbox"/>	AmazonRedshiftDataFullAccess	AWS managed	This policy provides f
<input type="checkbox"/>	AwsGlueSessionUserRestrictedNotebookPolicy	AWS managed	Provides permissions
<input type="checkbox"/>	AwsGlueSessionUserRestrictedNotebookServiceRole	AWS managed	Provides full access t



IAM
Roles



Data
Catalog



Crawler

AWS Glue > Databases > carsharing

carsharing

Last updated (UTC)
August 1, 2023 at 10:07:53



Edit

Delete

Database properties

Name	Description	Location	Created on (UTC)
carsharing	-	-	July 28, 2023 at 11:59:07

Tables (8)

View and manage all available tables.

Last updated (UTC)
August 1, 2023 at 10:07:53

↻

Delete

Data quality New

Add tables using crawler

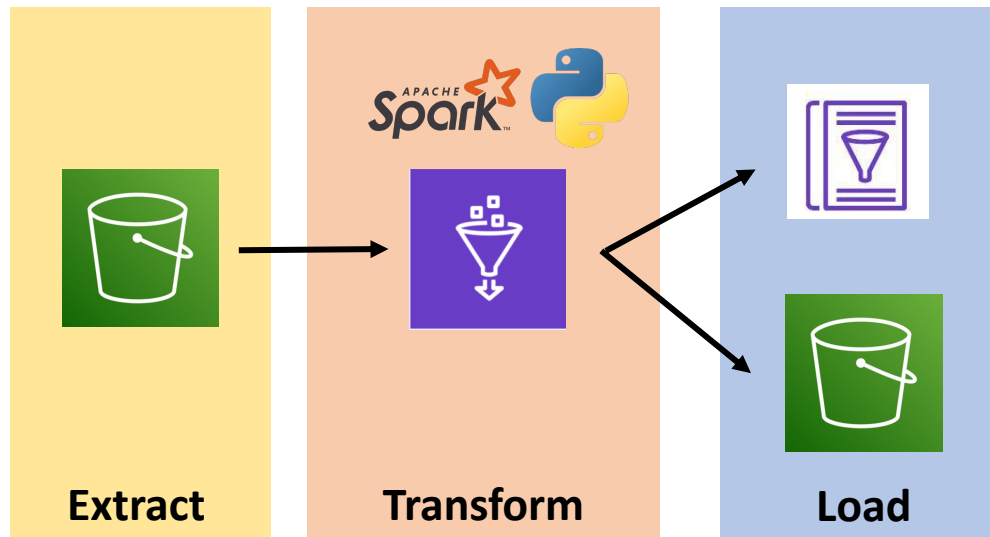
Add table

Q Filter tables

< 1 > ⚙

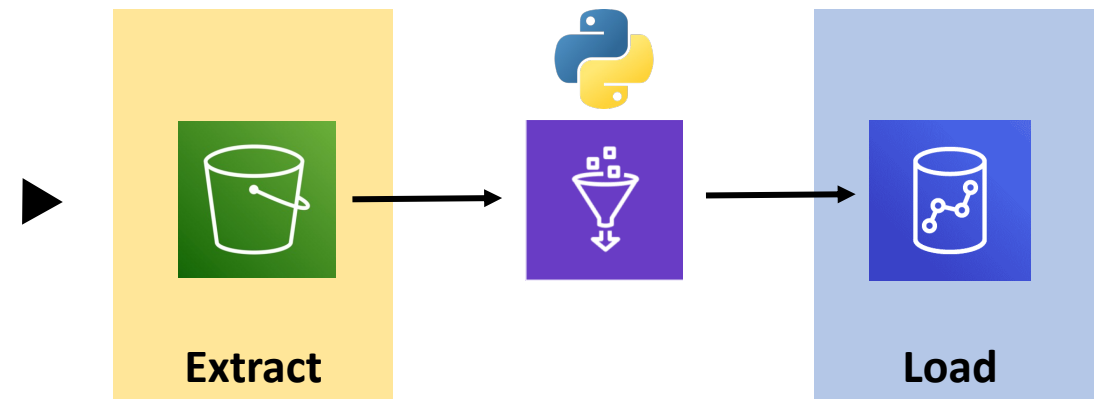
<input type="checkbox"/>	Name ▲	Database ▾	Location ▾	Classification ▾	Deprecated ▾	View data
<input type="checkbox"/>	carsharing_distance	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data
<input type="checkbox"/>	carsharing_group	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data
<input type="checkbox"/>	carsharing_promotion	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data
<input type="checkbox"/>	carsharing_reservation	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data
<input type="checkbox"/>	carsharing_station	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data
<input type="checkbox"/>	carsharing_user	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data
<input type="checkbox"/>	carsharing_vehicle	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data
<input type="checkbox"/>	carsharing_vehiclesize	carsharing	s3://car-sharing-bucket-f	CSV	-	Table data

3. ETL raw data using AWS Glue



◀ **ETL_S3_Athena:** Python Script for Extracting data from Amazon S3, Transforming data on Glue Notebook Studio using PySpark, Loading data to Amazon S3, and creating tables on the AWS Glue Catalog.

EL_Redshift: Python Script for Extracting data from Glue Catalog, and Loading data to Amazon Redshift with Redshift connector.

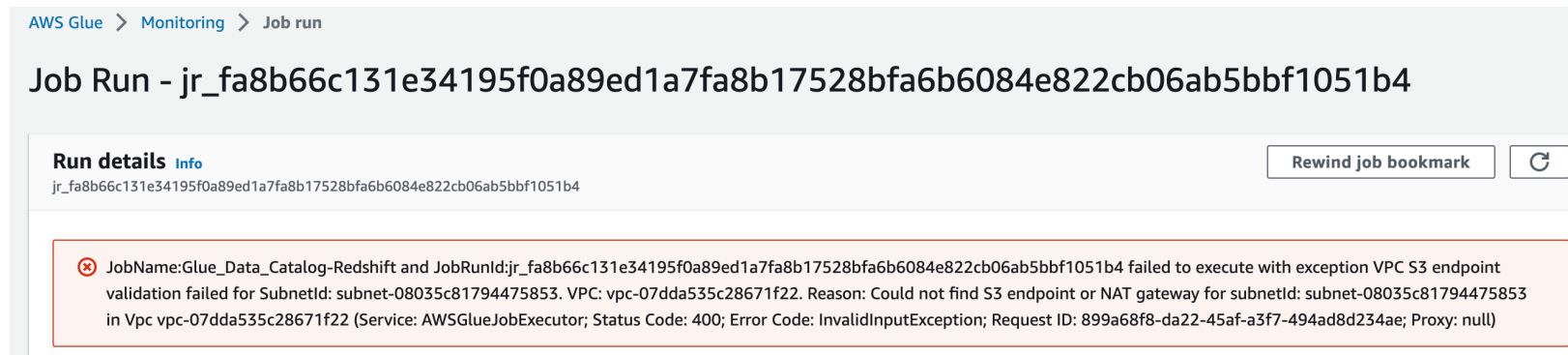


Problem

Can not load data to Redshift

Code with many script but it's not work, and finally...

1. Create Glue Connector (Connector type: Redshift) but...



2. Solved by create a new VPC Endpoint

<https://repost.aws/knowledge-center/glue-s3-endpoint-validation-failed>