

① → create s3 bucket → (CH-Project-s3-by-adi)

↳ enable public

↳ Acknowledge ✓

↳ Go to created bucket

↳ create subfolders

INPUT/

OUTPUT/

script/

temp/

↳ Go to input folder.

↳ create subfolders of

INPUT/Product/

↳ Go to INPUT Product Folders

↳ create subfolders

INPUT/Product/year=2021/

→ (Upload data into Partition folders after creation).

② Create Database

→ Go to Glue Service.

→ Go to Databases.

↳ Add database.

↳ Name → (CH-Project-db-by-adi)

↳ Create database.

③ Create IAM role ①

→ Go to IAM service.

→ Go to Access Management

↳ Roles

↳ Create role

→ PTO →

Cloud
Project
S3 data
integrated
with
Snowflake.

= AWS service under Trusted Only Type.

↳ use case → S3 -

↳ Next

↳ Permission Policies

↳ s3FullAccess.

↳ cloudWatchFullAccess.

↳ assumeServiceRole.

↳ Next.

↳ Role details

↳ Role name - (etl-project-~~iam~~-role-by-adi)

↳ Create role -

① Create Crawlers

→ Go to AWS Glue Service.

→ Go to Data Catalog

↳ Select Crawlers.

↳ Create crawler

↳ Name

(etl-project-crawlers-hatch-from-s3-inet
by-adi)

→ Next

↳ Choose Database & classifier

↳ ~~Not yet~~ ↳ Not yet

↳ Add a data source.

↳ Data source - S3

↳ Location of S3 data → this

↳ S3 Path : bucket/input/product

→ Add an S3 data source. ← (1)

→ Next

↳ Configure security settings.

↳ IAM role

↳ Select (ctt-project-lambda-role-by-adi)

→ Next

↳ Set up E scheduling. ②

↳ Target database.

ctt-project-blb-by-adi;

↳ Create schedule

↳ On-demand

↳ Review & Create

↳ Create crawlers.

→ Upload 2021 data file to S3 bucket under
Year=2021 folder.

→ Run crawlers now

- name (ctt-project-crawler-fetch-blb)
- S3 input (by-adi)

→ Go to tables under data catalog. and check
the table(s) and all

- ④ → Go to Athena service.
↳ check database (etl-project-hadoop-by-ad)
↳ check table (Product)
↳ compare the results by previous results.

⑤ Create ETL job

- Go to Step Functions service
↳ ETL jobs.
↳ visual ETL.
↳ under ~~general~~ job details tab
↳ name → etl-project-read-from-s3-job-by-ad
↳ IAM Role
↳ etl-project-iam-role1-by-ad
↳ type → spark
↳ Glue version → 3.0
↳ language - Python
↳ worker type - 16GB RAM (Select lower RAM)
↳ job timeout - 5 min.
↳ under Advanced Properties,
↳ script path
↳ bucket / script
↳ temporary path
↳ bucket / temp

↳ Under visual tab

↳ Select source

↳ S3

↳ ↳ Select source type

↳ Data catalog table

↳ database (etl-project-db-by-ad)

↳ Table (Product)

↳ Select Target

new columns in parquet
newyear
Count
Quantity

↳ S3

↳ Format (Parquet)

↳ S3 Target location

S3 bucket / output

↳ Under script tab

↳ edit script

↳ upload data from local drive.

→ Run job.

→ After success

→ Go to S3 bucket

bucket / output / newproduct /

→ Good.

⑤i) → Go to Give Service

↳ under 'Go Workflows (orchestration)'

↳ Add workflow

name (etl-project-workflow-by-adi)

↳ Create workflow

→ Open workflow.

↳ Go to workflow details.

↳ Graph

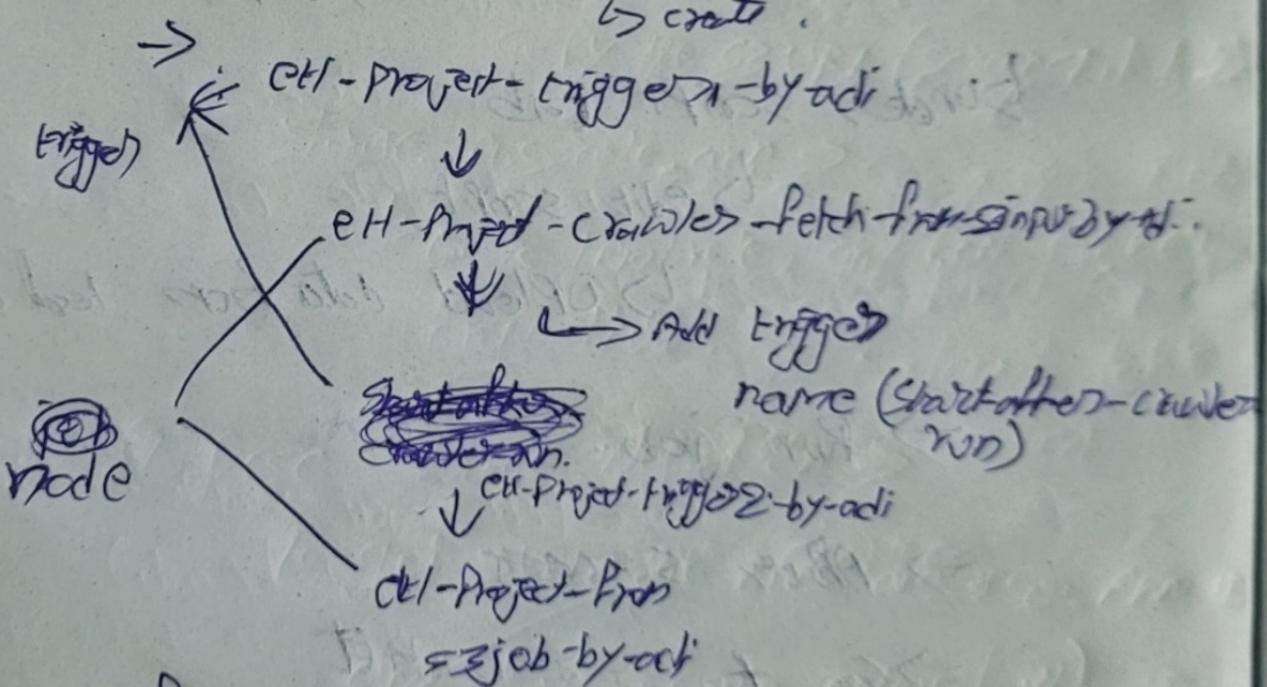
↳ Select Add trigger

↳ Add new

↳ name

(etl-project-trigger-by-adi)

↳ create.



All Good So Far

⑥ Snowflake works

① Create warehouse

↳ etl-project-snowwhby-adi;

② Create Database.

↳ etl-project-snowdb-by-adi;

③ Create Schema

↳ etl-project-snow^{db}schema-by-adi_Prod;

④ Create File Format

↳ etl-project-Parquet-Format-by-adi

⑤ Create IAM Role

→ Go to IAM service

→ Go to Access management
roles.

↳ Create role.

- ↳ Trusted entity type

↳ AWS account

↳ An AWS ~~Acc~~count

↳ This account

↳ Options

external ID

↳ For now give 0000

↳ Next

- ↳ Add Permissions

↳ AWS Full Access

↳ AmazonS3FullAccess

↳ Next

-> name, review, and create

↳ Role Name

↳ CH-Project-iam-role2-by-adi

↳ Create role

⑥ Create storage integration

⑦ Show

↳ name (CH-Project-Storage-intg-by-adi)

Name →

Type → External Stage

Enabled → True

Storage_Provider → S3

Storage_allowed_location ⇒

↳ Go to S3 output folder and

copy the S3 URL (Screenshot below)

Storage_AWS_ARN →

↳ Go to IAM role, copy ARN number
From there.

⑦ Desc storage integration →

To get → Storage-AWS_ARN & S3-ARN ID

Keep this ID under

1DM Trust edit Policy script

(Screenshot there)

and → Storage-AWS-external-ID) ID.

Keep this ID under

1DM must edit Policy script
(Screenshot There).

⑧ Create stage

Name → ETI-Project-stage-by-add

URL →

Go to S3 output bucket and copy URL
It's S3 URL (secret taken).

Storage-integration →

Name given to storage-integration

(stage-A,
ETI-Project-storage-by-add).

File-format = name given to fileformat

(ETI-Project-Request-Format-by-add)

⑨ Check file loaded or not

List ~~stage name~~ >

@ stagename.

if-way

⑩ Create target table.

⑪ Copy stage data to target table.

⑫ Verify results.

⑬ Create external table

For that do below.

① Create table first (Query available)

Name → Ext-Customer-Product.

② Show external tables.

↳ copy notification-channel value.

③ Go to S3 ~~output object bucket~~

bucket/output/object

bucket

↳ Properties

↳ event notifications.

↳ Create-event-notifications

↳ name -> ct1-Project-dataarch1
by-add

↳ Object creation

↳ All object creation events

↳ Object removal

↳ All object removal events

↳ Destination

↳ SNS output

↳ specifying sns queue

↳ enter sns queue ARN

↳ Paste the ARN copied
from external table
show command

→ save changes.

④ Go to S3 & upload data to S3/output
manually or through above job

⑤ Wait for some time query select.