# ETL GCP-Snowflakes

Then Tsze Yen

**Step 1:**

Create Cloud Storage Bucket
(to connect snowflakes, as my
snowflakes is in us central region,
then create bucket in us region)

# Step 2 :

Create a Cloud Storage Integration in Snowflake

# Step 3:

## Create Role

IAM & Admin 🔔

PAM [NEW]

Principal Access Boundary

Identity & Organization

Policy Troubleshooter

Policy Analyzer [NEW]

Organization Policies

Service Accounts

Workload Identity Federat...

Workforce Identity Federa...

Labels

Tags

Settings

Privacy & Security

Identity-Aware Proxy

Roles

Audit Logs

Manage Resources

← Edit Role

Custom roles let you group permissions and assign them to principals in your project or organization. You can manually select permissions or import permissions from another role. Learn more ↗

ID                    projects/airflow-workings/roles/CustomRole

Title *
snowflake-integration-role

26 / 100 characters

Description
Snowflake integration role

26 / 256 characters

Role launch stage
Alpha                                          ▼

➕ ADD PERMISSIONS

## 5 assigned permissions

Filter   Enter property name or value                    ❓   III

| Permission ↑ | Status |
|---|---|
| ☑ storage.buckets.get | Supported |
| ☑ storage.objects.create | Supported |
| ☑ storage.objects.delete | Supported |
| ☑ storage.objects.get | Supported |
| ☑ storage.objects.list | Supported |

# Step 4:

- Desc the integration to get the service account
- Assign role to Cloud Storage Service Account

# Step 5:

## Create Cloud Composer Environment & Add PYPI package

Step 6:
Add airflow-snowflake connection

# Step 6:

Upload DAG file

ONE



TWO

# Output of the Snowflake

Thank You