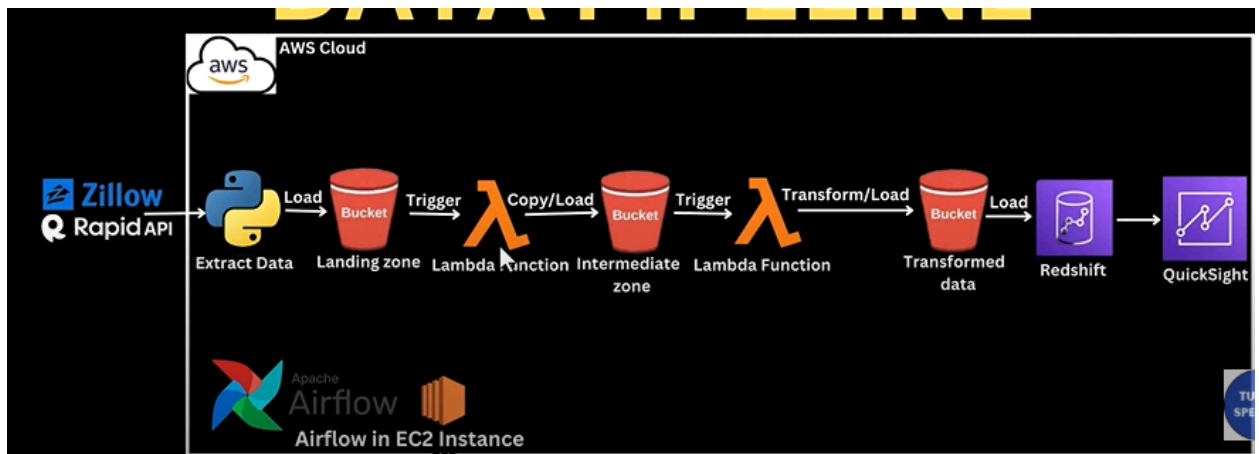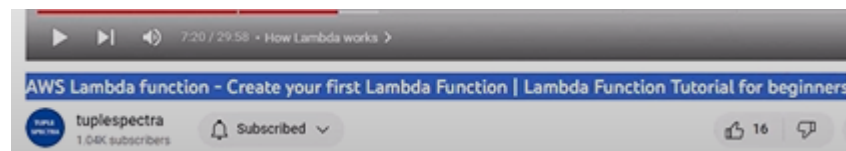Landing Zone - Intermediate zone



1.
2. Here we have extracted data from api and copies into initial s3 bucket
3. Our task is trigger a lambda function and load this data to intermediate zone
4. Now head to AWS and search for AWS lambda



5. Lambda tutorial :



6. Now we will create a lambda function



7. Next we will create a role for our lambda function …we will our s3 access to our lambda

role



8. We also need to give this policy to the role



9. So whatever is happening with our lambda function..we can watch them in cloudwatch

10. Next we give a name to our role



11. Next we will attach our role to the lambda function..and we create the lambda function
12. So now our landing_zone bucket needs to trigger the lambda..when ever there is data in



it
13. Here we can add trigger to our lambda function

14. Then we select source as s3



15. Next we have to choose the bucket and the event_type in the bucket

16. Here we can see .we've added trigger to our lambda function



17. Next we have coded our lambda function



18. Later to test the code..we have uploaded a sample file to pur s3 bucket..now as soon something lands in our s3 bucket..our lambda gets triggered

```python
import boto3
import json

s3_client = boto3.client('s3')

def lambda_handler(event, context):
    # TODO implement
    source_bucket = event['Records'][0]['s3']['bucket']['name']
    object_key = event['Records'][0]['s3']['object']['key']
    print(source_bucket)
    print(object_key)
```

19. So now we have to get our source_bucket name and object

20. To view them we go to monitor=>cloudwatch



21. Here we can see our bucket name and file name(that we have uploaded)



22. Basically what our code is doing is…getting the source_bucket and the files uploaded to source bucket

23. And copying this data to the target bucket(intermediate bucket)

```python
import boto3
import json

s3_client = boto3.client('s3')

def lambda_handler(event, context):
    # TODO implement
    source_bucket = event['Records'][0]['s3']['bucket']['name']
    object_key = event['Records'][0]['s3']['object']['key']

    target_bucket = 'copy-of-raw-json-bucket'
    copy_source = {'Bucket': source_bucket, 'Key': object_key}

    waiter = s3_client.get_waiter('object_exists')
    waiter.wait(Bucket=source_bucket, Key=object_key)
    s3_client.copy_object(Bucket=target_bucket, Key=object_key, CopySource=copy_source)
    return {
        'statusCode': 200,
        'body': json.dumps('Copy completed successfully')
    }
```

Explained : https://g.co/gemini/share/dbbb4f0dc7bc

24. But before copying we use "waiter" to confirm the entire file has been uploaded to source code and then we make copy

25. Now let us create the intermediate bucket



26. Here we have create a bucket



Successfully created bucket "copy-of-raw-json-bucket"
To upload files and folders, or to configure additional bucket settings choose **View details**.

27. Now we need to send data from end_to_end bucket to raw bucket



    ○    copy-of-raw-json-bucket
    ○    endtoendyoutube-ym-bucket

28. Lets us trigger our DAG and check what happend

29. Using our DAG we again loaded data to our initial bucket(end2end) and the lambda got triggered and copied the data to raw_bucket



30. Now we need to trigger another lambda function that gets triggered when the data is in intermediate zone and transforms and load to another bucket



Intermediate - Transformed

1. Now we need to write another lambda function..and while creating that function..we will attach the role which we created earlier

2. Next we will add trigger ..next we'll choose trigger source as s3

**Bucket**

Please select the S3 bucket that serves as the event source. The bucket must be in the same region as the function.

| 🔍 | | C |

cellpaintingimages-yml

copy-of-raw-json-bucket

**Event types**

Select the events that you want to have trigger the Lambda function. You can optionally set up a prefix or suffix for an event. However, for each bucket, individual events cannot have multiple configurations with overlapping prefixes or suffixes that could match the same object key.

| | ▼ |

All object create events ✕

3. Now Lets focus on the lambda code

```python
import boto3
import json
import pandas as pd

s3_client = boto3.client('s3')

def lambda_handler(event, context):
    # TODO implement
    source_bucket = event['Records'][0]['s3']['bucket']['name']
    object_key = event['Records'][0]['s3']['object']['key']
    print(object_key)
    print(source_bucket)
    target_bucket = 'cleaned-data-zone-csv-bucket'
    target_file_name = object_key[:-5]
    print(target_file_name)

    waiter = s3_client.get_waiter('object_exists')
    waiter.wait(Bucket=source_bucket, Key=object_key)

    response = s3_client.get_object(Bucket=source_bucket, Key=object_key)
    print(response)
    data = response['Body']
    print(data)
    data = response['Body'].read().decode('utf-8')
    print(data)
    data = json.loads(data)
    print(data)
    f = []
```
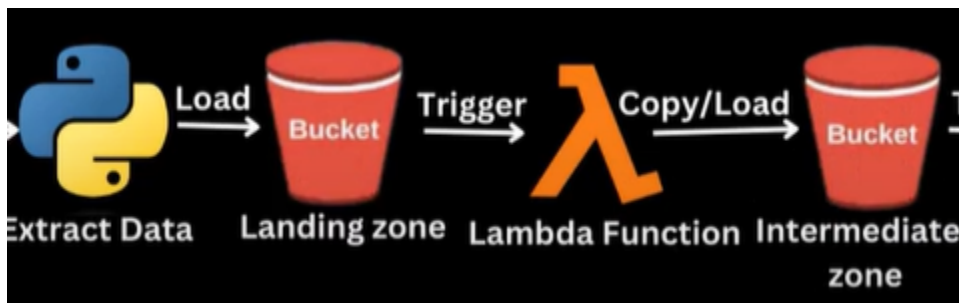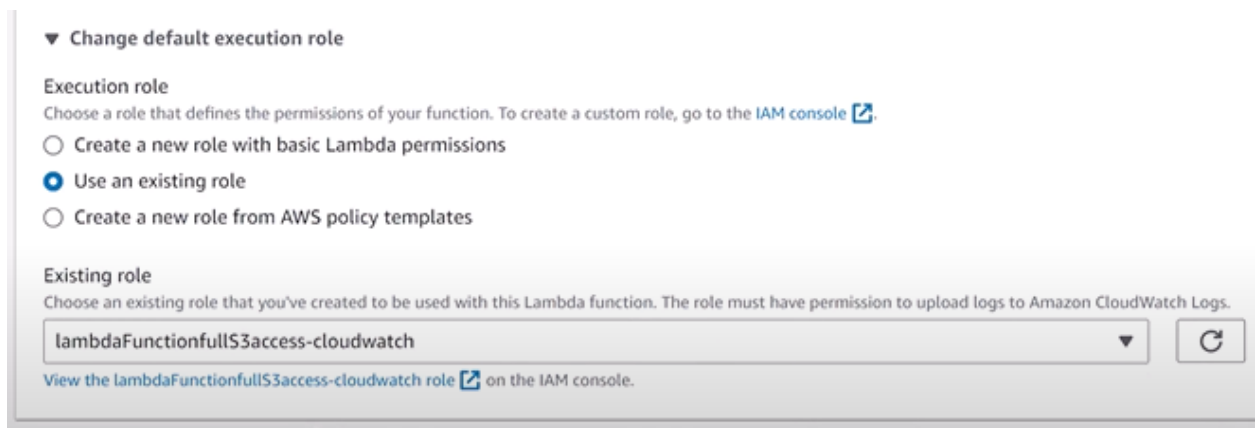
4. Next we'll create a bucket for cleaned_data(target_bucket)
5. Then we'll retrieve data from response['Body'] ..and then we need to decode it to 'utf-8'
6. To see the data…we used json.loads(data)

7. Data will look like this

```
{
    "results": [
        {
            "bathrooms": 2.0,
            "bedrooms": 3.0,
            "city": "Houston",
            "country": "USA",
            "currency": "USD",
            "daysOnZillow": 0,
            "homeStatus": "FOR_SALE",
            "homeStatusForHDP": "FOR_SALE",
            "homeType": "SINGLE_FAMILY",
            "imgSrc": "https://photos.zillowstatic.com/fp/eaa31be6b0b2267323624c193367be4c-p_e.jpg",
            "isFeatured": false,
            "isNonOwnerOccupied": true,
            "isPreforeclosureAuction": false,
            "isPremierBuilder": false,
            "isShowcaseListing": false,
            "isUnmappable": false,
            "isZillowOwned": false,
            "latitude": 29.885355,
            "listing_sub_type": {
                "is_FSBA": true
            },
            "livingArea": 1060.0,
            "longitude": -95.61802,
            "lotAreaUnit": "sqft",
            "lotAreaValue": 5998.212,
            "price": 125000.0,
            "priceForHDP": 125000.0,
            "rentZestimate": 1566,
            "shouldHighlight": false,
            "state": "TX",
            "streetAddress": "13723 Smokey Trail Dr",
            "taxAssessedValue": 158004.0,
            "zestimate": 125900,
            "zipcode": "77041",
            "zpid": 28282152
        },
        {
            "bathrooms": 3.0,
            "bedrooms": 4.0,
            "city": "Houston",
            "country": "USA",
            "currency": "USD",
            "daysOnZillow": 0,
            "homeStatus": "FOR_SALE",
            "homeStatusForHDP": "FOR_SALE",
            "homeType": "SINGLE_FAMILY",
```

8. Basically we took some columns and converted it to df...then again we converted this data to CSV

```python
f.append(t)
df = pd.DataFrame(f)
# Select specific columns
selected_columns = ['bathrooms', 'bedrooms', 'city', 'homeStatus',
                'homeType','livingArea','price', 'rentZestimate','zipcode']
df = df[selected_columns]
print(df)

# Convert DataFrame to CSV format
csv_data = df.to_csv(index=False)

# Upload CSV to S3
bucket_name = target_bucket
object_key = f"{target_file_name}.csv"
s3_client.put_object(Bucket=bucket_name, Key=object_key, Body=csv_data)


return {
    'statusCode': 200,
    'body': json.dumps('CSV conversion and S3 upload completed successfully')
}
```

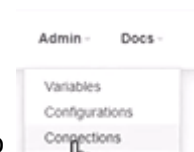9. Now we'll create an another task in our airflow

```python
is_file_in_s3_available = S3KeySensor(
    task_id='tsk_is_file_in_s3_available',
    bucket_key='{{ti.xcom_pull("tsk_extract_zillow_data_var")[1]}}',
    bucket_name=s3_bucket,
    aws_conn_id='aws_s3_conn',
    wildcard_match=False,  # Set this to True if you want to use wildcards in the prefix
    timeout=120,  # Optional: Timeout for the sensor (in seconds)
    poke_interval=5,  # Optional: Time interval between S3 checks (in seconds)
)
```

10. First we need to import this key sensor

```python
from airflow.providers.amazon.aws.sensors.s3 import S3KeySensor
```

The `S3KeySensor` is a sensor provided by Airflow's AWS provider that allows you to pause a DAG execution until a specific file (key) becomes available in an S3 bucket. It essentially waits for a file to be uploaded before proceeding with downstream tasks in your Airflow DAG.

11. More on s3keysensor : https://g.co/gemini/share/04c9d4eb3193

12. And also we need to give connection to our aws..so we go to

Admin     Docs

Variables
Configurations
Connections