# Architecture OF ETL Pipeline Created Using AWS

# Spotify_api_data_extraction (Function)



```
1   import json
2   import os
3   import spotipy
4   from spotipy.oauth2 import SpotifyClientCredentials
5
6
7   import boto3 # package that connects you with the Amazon services
8   from datetime import datetime
9   def lambda_handler(event, context):
10
11      cilent_id=os.environ.get('client_id')
12      cilent_secret=os.environ.get('client_secret')
13
14      client_credentials_manager = SpotifyClientCredentials(client_id=cilent_id,client_secret=cilent_secret)
15      sp= spotipy.Spotify(client_credentials_manager = client_credentials_manager)
16      playlist_link= "https://open.spotify.com/playlist/37i9dQZEVXbNG2KDcFcKOF"
17      playlist_URI=playlist_link.split("/")[-1].split('?')[0]
18      spotify_data = sp.playlist_tracks(playlist_URI)
19      print(spotify_data)
20      cilent=boto3.client('s3')
21
22      filename="sotify_raw_"+ str(datetime.now()) + ".json"
23      cilent.put_object(
24          Bucket="spotify-etl-project-gurparteek",
25          Key="raw_data/to_processed/" + filename, #Path were you want to store your data
26          Body=json.dumps(spotify_data)
27          )#This will copvert entire data into JSON Dictionary Foramt
```

**This code with extract the data from Spotify using spotify(API)**

# Spotify_transformation_load_function (Function)

**Code source** Info

Upload from ▼

File  Edit  Find  View  Go  Tools  Window          Test  ▼     Deploy

Q  Go to Anything (Ctrl-P)

lambda_function ×   Environment Vari ×   Execution results ×   ⊕

▼ 📁 spotify_transformat ⚙▼
    📄 lambda_function.py

```python
1  import json
2  import boto3
3  from datetime import datetime
4  from io import StringIO # fOR THE CSV FILES
5  import pandas as pd
6
7  def album(data):
8   album_list = []
9   for row in data['items']:
10      album_id = row['track']['album']['id']
11      album_name = row['track']['album']['name']
12      album_release_date = row['track']['album']['release_date']
13      album_total_tracks = row['track']['album']['total_tracks']
14      album_url = row['track']['album']['external_urls']['spotify']
15      album_element = {'album_id':album_id,'name':album_name,'release_date':album_release_date,
16                      'total_tracks':album_total_tracks,'url':album_url}
17      album_list.append(album_element)
18   return album_list
19
20
21  def artist(data):
22   artist_list = []
23   for row in data['items']:
24      for key, value in row.items():
25          if key == "track":
26              for artist in value['artists']:
27                  artist_dict = {'artist_id':artist['id'], 'artist_name':artist['name'], 'external_url': artist['href']}
28                  artist_list.append(artist_dict)
29
30   return artist_list
31
32
33  def songs(data):
34   song_list = []
35   for row in data['items']:
36      song_id = row['track']['id']
37      song_name = row['track']['name']
```

1:1   Python   Spaces: 4 ⚙

The 3 functions namely **album,artist and songs** will extract data such as album_name, artist_name etc from the spotify data.

```python
def songs(data):
    song_list = []
    for row in data['items']:
        song_id = row['track']['id']
        song_name = row['track']['name']
        song_duration = row['track']['duration_ms']
        song_url = row['track']['external_urls']['spotify']
        song_popularity = row['track']['popularity']
        song_added = row['added_at']
        album_id = row['track']['album']['id']
        artist_id = row['track']['album']['artists'][0]['id']
        song_element = {'song_id':song_id,'song_name':song_name,'duration_ms':song_duration,'url':song_url,
                        'popularity':song_popularity,'song_added':song_added,'album_id':album_id,
                        'artist_id':artist_id
                       }
        song_list.append(song_element)

    return song_list



def lambda_handler(event, context):
    s3= boto3.client('s3') # Crearing the Object
    Bucket="spotify-etl-project-gurparteek"
    Key="raw_data/to_processed" # Path where the files are stored

    #This is a function
    #print(s3.list_objects(Bucket=Bucket,Prefix=Key)['Contents'])

    spotify_data=[]
    spotify_keys=[]

    # Here we are using Forloop because we have List and inside the list we
    # we have Mutiple Dictionaries
```

https://us-east-1.console.aws.amazon.com/lambda/home?region=us-east-1#/functions/spotify_transformation_load_function?tab=code

aws | Services | Q Search | [Alt+S] | N. Virginia ▼ | Gurparteek Gill ▼

## Code source  Info

Upload from ▼

File  Edit  Find  View  Go  Tools  Window  **Test** ▼  Deploy

Go to Anything (Ctrl-P)

lambda_function ×  Environment Var ×  Execution results ×  ⊕

spotify_transformat ⚙▼

lambda_function.py

```python
52
53
54
55   def lambda_handler(event, context):
56       s3= boto3.client('s3') # Creating the Object
57       Bucket="spotify-etl-project-gurparteek"
58       Key="raw_data/to_processed" # Path where the files are stored
59
60       #This is a function
61       #print(s3.list_objects(Bucket=Bucket,Prefix=Key)['Contents'])
62
63       spotify_data=[]
64       spotify_keys=[]
65
66       # Here we are using Forloop because we have List and inside the list we
67       # we have Mutiple Dictionaries
68       for file in s3.list_objects(Bucket=Bucket,Prefix=Key)['Contents']:
69           #print(file['Key'])
70           file_key=file['Key']
71           if file_key.split('.')[-1] == "json": # Means it will pick only Json files
72
73               # Inside the s3 object we have this function availbale
74               response =s3.get_object(Bucket = Bucket,Key = file_key)
75               content=response['Body']
76               jsonObject= json.loads(content.read())
77               spotify_data.append(jsonObject)
78               spotify_keys.append(file_key)
79               # print(jsonObject)
80
81       for data in spotify_data:
82           album_list= album(data)
83           artist_list= artist(data)
84           song_list= songs(data)
85
86           # print(album_list)
87           album_df = pd.DataFrame.from_dict(album_list)
88           album_df = album_df.drop_duplicates(subset=['album_id'])
89
```

1:1  Python  Spaces: 4 ⚙

**Now we will store the extracted data inside the (raw_data/to_processed) Bucket .**

⟨ → C ⌂ ⚷ https://us-east-1.console.aws.amazon.com/lambda/home?region=us-east-1#/functions/spotify_transformation_load_function?tab=code ☆ ⬙ | ≡ᵈ ◻ G ⋮

aws ▦ Services 🔍 Search [Alt+S] ⬚ ⌂ ? ⚙ N. Virginia ▼ Gurparteek Gill ▼

**Code source** Info

Upload from ▼

▲ File Edit Find View Go Tools Window Test ▼ Deploy ⤢ ⚙

🔍 Go to Anything (Ctrl-P) ▤ lambda_function ✕ | Environment Var ✕ | Execution results ✕ ⊕

```python
80
81     for data in spotify_data:
82         album_list= album(data)
83         artist_list= artist(data)
84         song_list= songs(data)
85
86         # print(album_list)
87     album_df = pd.DataFrame.from_dict(album_list)
88     album_df = album_df.drop_duplicates(subset=['album_id'])
89
90     artist_df = pd.DataFrame.from_dict(artist_list)
91     artist_df = artist_df.drop_duplicates(subset=['artist_id'])
92
93     song_df = pd.DataFrame.from_dict(song_list)
94
95     album_df['release_date'] = pd.to_datetime(album_df['release_date'])
96     song_df['song_added'] =  pd.to_datetime(song_df['song_added'])
97
98
99     songs_key= "transformed_data/songs_data/songs_transformed" + str(datetime.now()) + ".csv"
100    song_buffer=StringIO()# For string conversion
101    song_df.to_csv(song_buffer,index=False) # Will finally covert
102    song_content= song_buffer.getvalue()
103    s3.put_object(Bucket=Bucket,Key=songs_key,Body= song_content)
104
105    album_key= "transformed_data/album_data/album_transformed" + str(datetime.now()) + ".csv"
106    album_buffer=StringIO()# For string conversion
107    album_df.to_csv(album_buffer,index=False) # Will finally covert
108    album_content= album_buffer.getvalue()
109    s3.put_object(Bucket=Bucket,Key=album_key,Body= album_content)
110
111    artist_key= "transformed_data/artist_data/artist_transformed" + str(datetime.now()) + ".csv"
112    artist_buffer=StringIO()# For string conversion
113    artist_df.to_csv(artist_buffer,index=False) # Will finally covert # Index false because glue caller will not be able to access the entire schema
114    artist_content= artist_buffer.getvalue()
115    s3.put_object(Bucket=Bucket,Key=artist_key,Body= artist_content)
116
```

1:1   Python   Spaces: 4 ⚙

**Now we will transform the JSON data into pandas Dataframe format.**

## Code source  Info

Upload from ▾

File   Edit   Find   View   Go   Tools   Window        Test ▾      Deploy

Go to Anything (Ctrl-P)

▾ 📁 spotify_transformat ⚙▾
   ⟨⟩ lambda_function.py

```python
      album_key= transformed_data/album_data/album_transformed + str(datetime.now()) + ".csv"
106   album_buffer=StringIO()# For string conversion
107   album_df.to_csv(album_buffer,index=False) # Will finally covert
108   album_content= album_buffer.getvalue()
109   s3.put_object(Bucket=Bucket,Key=album_key,Body= album_content)
110
111   artist_key= "transformed_data/artist_data/artist_transformed" + str(datetime.now()) + ".csv"
112   artist_buffer=StringIO()# For string conversion
113   artist_df.to_csv(artist_buffer,index=False) # Will finally covert # Index false because glue caller will not be able to access the entire schema
114   artist_content= artist_buffer.getvalue()
115   s3.put_object(Bucket=Bucket,Key=artist_key,Body= artist_content)
116
117
118
119   s3_resource = boto3.resource('s3')
120   for key in spotify_keys:
121       copy_source ={
122           'Bucket':Bucket,
123           'Key':key
124                               #Targated Bucket
125       }
126       s3_resource.meta.client.copy(copy_source,Bucket,'raw_data/processed/' + key.split("/")[-1])
127       s3_resource.Object(Bucket,key).delete()
128
```

**This will store the transformed data into (raw_data/processed) Bucket.**

1:1   Python   Spaces: 4 ⚙

# s3 Bucket



Amazon S3

https://s3.console.aws.amazon.com/s3/buckets/spotify-etl-project-gurparteek?region=us-east-1&tab=objects

aws    Services    Q Search    [Alt+S]      Global ▼    Gurparteek Gill ▼

**Amazon S3** ✕

**Buckets**

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

▼ **Storage Lens**

Dashboards

AWS Organizations settings

Feature spotlight   7

▶ AWS Marketplace for S3

Amazon S3 > Buckets > spotify-etl-project-gurparteek

# spotify-etl-project-gurparteek   Info

| Objects | Properties | Permissions | Metrics | Management | Access Points |
|---------|-----------|-------------|---------|-----------|---------------|

## Objects (3)

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 ↗

| | | | | | |
|---|---|---|---|---|---|
| ↻ | Copy S3 URI | Copy URL | Download | Open ↗ | Delete | Actions ▼ | Create folder | **Upload** |

Find objects by prefix      ‹ 1 › ⚙

| ☐ | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|--------|--------|-----------------|--------|-----------------|
| ☐ | 📄 raw_data/ | Folder | - | - | - |
| ☐ | 📄 transformed_data/ | Folder | - | - | - |
| ☐ | 📄 Unsaved/ | Folder | - | - | - |

Tables - AWS Glue Con | spotify-etl-project-gu | Functions - Lambda | spotify_api_data_extra | spotify_transformatio | Query editor | Athena | Gucci Chick - YouT

https://s3.console.aws.amazon.com/s3/buckets/spotify-etl-project-gurparteek?region=us-east-1&prefix=raw_data/&showversions=false

aws | Services | Q Search [Alt+S] | Global ▼ | Gurparteek Gill ▼

**Amazon S3** ✕

Buckets

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

▼ Storage Lens

Dashboards

AWS Organizations settings

Feature spotlight 7

▶ AWS Marketplace for S3

Amazon S3 > Buckets > spotify-etl-project-gurparteek > raw_data/

Copy S3 URI

# raw_data/

**Objects** | Properties

**Objects** (2)

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 ↗

| ↻ | Copy S3 URI | Copy URL | Download | Open ↗ | Delete | Actions ▼ | Create folder | Upload |

Find objects by prefix

< 1 > ⚙

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 🗋 processed/ | Folder | - | - | - |
| ☐ | 🗋 to_processed/ | Folder | - | - | - |

https://s3.console.aws.amazon.com/s3/buckets/spotify-etl-project-gurparteek?region=us-east-1&prefix=transformed_data/&showversions=false

aws   ::: Services   🔍 Search   [Alt+S]   Global ▾   Gurparteek Gill ▾

**Amazon S3**   ✕

**Buckets**

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

▼ Storage Lens

Dashboards

AWS Organizations settings

Feature spotlight  7

▶ AWS Marketplace for S3

Amazon S3 > Buckets > spotify-etl-project-gurparteek > transformed_data/

# transformed_data/

[ ⧉ Copy S3 URI ]

**Objects** | Properties

## Objects (3)

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 ↗

| ↻ | Copy S3 URI | Copy URL | Download | Open ↗ | Delete | Actions ▼ | Create folder | ⬆ Upload |

🔍 Find objects by prefix                                          < 1 >   ⚙

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 🗀 album_data/ | Folder | - | - | - |
| ☐ | 🗀 artist_data/ | Folder | - | - | - |
| ☐ | 🗀 songs_data/ | Folder | - | - | - |

Crawler

# 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 (3) Info

View and manage all available crawlers.

Last updated (UTC)
October 27, 2023 at 13:57:27

Filter crawlers

| | Name ▽ | State ▽ | Schedule | Last run ▽ | Last run timesta... ▽ | Log | Table changes fro... |
|---|---|---|---|---|---|---|---|
| ☐ | spotify_album_craw... | ⊘ Ready | | | ⊘ Succeeded | October 27, 2023 a... | View log ↗ | 1 created |
| ☐ | spotify_artists_craw... | ⊘ Ready | | | ⊘ Succeeded | October 27, 2023 a... | View log ↗ | 1 created |
| ☐ | spotify_songs_crawler | ⊘ Ready | | | ⊘ Succeeded | October 27, 2023 a... | View log ↗ | 1 created |

**Here we have created 3 crawlers(album,artists,songs)**

**⊘ Crawler successfully starting**
The following crawler is now starting: "spotify_album_crawler"

AWS Glue

- Getting started
- ETL jobs
  - Visual ETL
  - Notebooks
  - Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)

▼ Data Catalog
- Databases
  - Tables
- Stream schema registries
  - Schemas
- Connections
- **Crawlers**
  - Classifiers
- Catalog settings

▶ Data Integration and ETL

▶ Legacy pages

What's New ↗
Documentation ↗
AWS Marketplace

**This is the preview of artists table.**

Views (0)   ‹ 1 ›

Query results    Query stats

⊘ Completed            Time in queue: 190 ms    Run time: 455 ms    Data scanned: 8.97 KB

**Results** (10)                    Copy    Download results

🔍 Search rows                                    ‹ 1 ›  ⚙

| # ▽ | artist_id ▽ | artist_name ▽ | external_url ▽ |
|---|---|---|---|
| 1 | 4q3ewBCX7sLwd24euuV69X | Bad Bunny | https://api.spotify.com/v1/artists/4q3ewBCX7sLwd24euuV69X |
| 2 | 6HaGTQPmzraVmaVxvz6EUc | Jung Kook | https://api.spotify.com/v1/artists/6HaGTQPmzraVmaVxvz6EUc |
| 3 | 3MdXrJWsbVzdn6fe5JYkSQ | Latto | https://api.spotify.com/v1/artists/3MdXrJWsbVzdn6fe5JYkSQ |
| 4 | 5cj0lLjcoR7YOSnhnX0Po5 | Doja Cat | https://api.spotify.com/v1/artists/5cj0lLjcoR7YOSnhnX0Po5 |
| 5 | 2LRoIwlKmHjgvigdNGBHNo | Feid | https://api.spotify.com/v1/artists/2LRoIwlKmHjgvigdNGBHNo |
| 6 | 45dkTj5sMRSjrmBSBeiHym | Tate McRae | https://api.spotify.com/v1/artists/45dkTj5sMRSjrmBSBeiHym |
| 7 | 0jbo7KFNMilkfBR6ih0yhm | iñigo quintero | https://api.spotify.com/v1/artists/0jbo7KFNMilkfBR6ih0yhm |
| 8 | 3qsKSpcV3ncke3hw52JSMB | Young Miko | https://api.spotify.com/v1/artists/3qsKSpcV3ncke3hw52JSMB |
| 9 | 06HL4z0CvFAxyc27GXpf02 | Taylor Swift | https://api.spotify.com/v1/artists/06HL4z0CvFAxyc27GXpf02 |
| 10 | 7uMDnSZyUYNBPLhPMNuaM2 | Kenya Grace | https://api.spotify.com/v1/artists/7uMDnSZyUYNBPLhPMNuaM2 |

Databases - AWS Glue    spotify-etl-project-gu    Functions - Lambda    spotify_api_data_extra    spotify_transformatio    Query editor | Athena    Diljit Dosanjh: LO

https://us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor/history/20644263-73ef-4585-9312-40a1b3fd014c

aws    Services    Search    [Alt+S]      N. Virginia ▾    Gurparteek Gill ▾

Query results    Query stats

⊘ Completed      Time in queue: 307 ms    Run time: 879 ms    Data scanned: 8.54 KB

## Results (10)

Copy    Download results

Search rows      ‹ 1 › ⚙

**This is the preview of songs table** ←

| # | song_id | song_name | duration_ms | url |
|---|---------|-----------|-------------|-----|
| 1 | 4MjDJD8cW7iVeWInc2Bdyj | MONACO | 267194 | https://open.spotify.com/track/4MjDJD8cW7iVeWInc2Bdyj |
| 2 | 7x9aauaA9cu6tyfpHnqDLo | Seven (feat. Latto) (Explicit Ver.) | 184400 | https://open.spotify.com/track/7x9aauaA9cu6tyfpHnqDLo |
| 3 | 56y1jOTK0XSvJzVv9vHQBK | Paint The Town Red | 230480 | https://open.spotify.com/track/56y1jOTK0XSvJzVv9vHQBK |
| 4 | 7iQXYTyuG13aoeHxGG28Nh | PERRO NEGRO | 162767 | https://open.spotify.com/track/7iQXYTyuG13aoeHxGG28N |
| 5 | 3rUGC1vUpkDG9CZFHMur1t | greedy | 131872 | https://open.spotify.com/track/3rUGC1vUpkDG9CZFHMur1 |
| 6 | 2HafqoJbgXdtjwCOvNEF14 | Si No Estás | 184061 | https://open.spotify.com/track/2HafqoJbgXdtjwCOvNEF14 |
| 7 | 3nNmRE0DxHC6ZaKkrpUumS | FINA | 216327 | https://open.spotify.com/track/3nNmRE0DxHC6ZaKkrpUu |
| 8 | 1BxfuPKGuaTgP7aM0Bbdwr | Cruel Summer | 178426 | https://open.spotify.com/track/1BxfuPKGuaTgP7aM0Bbdw |
| 9 | 5mjYQaktjmjcMKcUIcqz4s | Strangers | 172964 | https://open.spotify.com/track/5mjYQaktjmjcMKcUIcqz4s |
| 10 | 01qFKNWq73UfEslI0GvumE | 3D (feat. Jack Harlow) | 201812 | https://open.spotify.com/track/01qFKNWq73UfEslI0Gvum |