

Create Lark data mart checker app using AWS Lambda, AWS SAM

In this article, you will learn to create a data mart checker for country data marts (SG, MY) in Lark chat room (Data Pipeline Monitoring chat group) using AWS Lambda, [AWS Serverless Application Model \(SAM\) CLI](#), and Lark API.

Read more about AWS Lambda, [AWS Serverless Application Model \(SAM\) CLI](#) with [examples](#) attached.

Instruction

To build AWS Lambda Function, steps needed to cover:

1. Get data marts tables info from PostgreSQL
2. Get Lark Service API, Install AWS SAM CLI interface.
3. Define AWS Lambda function with the parameters we prepared, store the code file to Docker Image with the new build
4. Deploy the app with the newly built image to our server

Repositories:

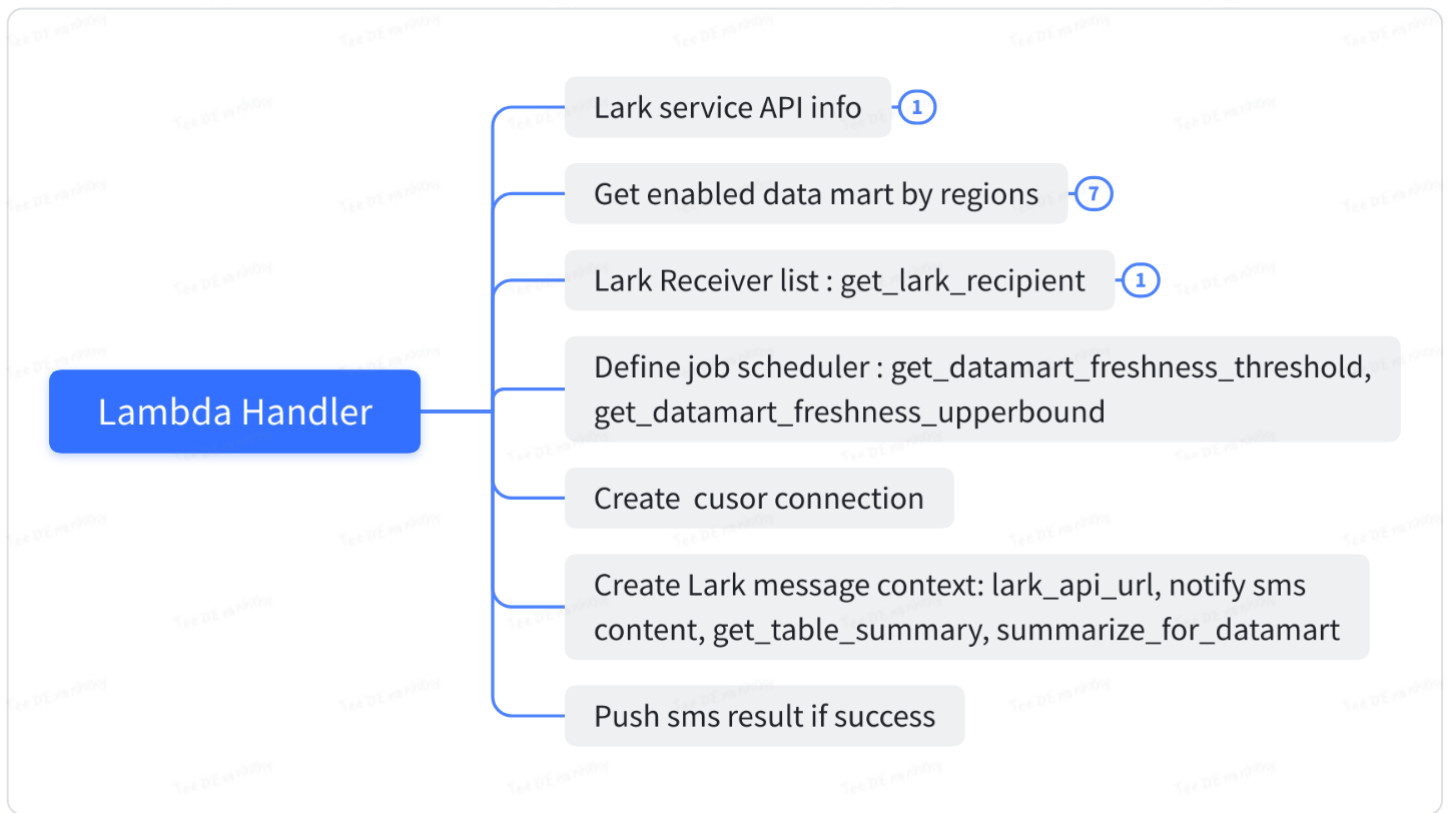
1. <https://github.com/TrustyCars/carro-datamart-monitoring>

This project contains source code and supporting files for a serverless application that you can deploy with the SAM CLI. It includes the following files and folders.

- src - Code for the application's Lambda function.
- events - Invocation events that you can use to invoke the function.
- tests - Unit tests for the application code.
- template.yaml - A template that defines the application's AWS resources.

The app builder process

The whole app builder pipeline on this project is as follows, which is to get information we need to push into AWS Lambda Handler function. The app code can be found at </src/app.py>



The resources we will need:

- carro-datamart-monitoring SAM template in [yaml](#) file (template that defines the application's AWS resources)
- [gitignore](#)
- ARN set up and [events](#) file (Invocation events that you can use to invoke the function)
- Besides, you can setup AWS SAM CLI and clone [tests](#) folder to play around with.

Requirements.txt

These packages are not supported in AWS Lambda, Requirement.txt file will help the system know what to install packages you desire to run your code.

If you test in your local IDE, make sure you have requirement packages installed and updated.

```
1 psycpg2-binary==2.9.3
2 requests==2.28.1
```

Code explain

@ Class EnvHandler : A class to handle/parse env vars.

Lark Service API info

```

1 def get_lark_service_api_info(self) -> dict:
2     api_info = dict()
3     attr_host = "LARK_SERVICE_HOST"
4     attr_port = "LARK_SERVICE_PORT"
5     attr_proto = "LARK_SERVICE_PROTO"
6     attr_endpoint = "LARK_SERVICE_ENDPOINT"
7     attr_api_key = "LARK_SERVICE_API_KEY"
8
9     api_info["host"] = os.environ.get(attr_host)
10    api_info["port"] = os.environ.get(attr_port)
11    api_info["proto"] = os.environ.get(attr_proto)
12    api_info["endpoint"] = os.environ.get(attr_endpoint)
13    api_info["api_key"] = os.environ.get(attr_api_key)
14
15    setattr(self, attr_host, api_info.get("host"))
16    setattr(self, attr_port, api_info.get("port"))
17    setattr(self, attr_proto, api_info.get("proto"))
18    setattr(self, attr_endpoint, api_info.get("endpoint"))
19    setattr(self, attr_api_key, api_info.get("api_key"))
20
21    return api_info

```

Get Lark Receivers

```

1 def get_lark_recipient(self) -> dict:
2     notify_info = dict()
3     available_recipient_keys = [
4         "NOTIFY_TO_LARK_CHAT_ID",
5         "NOTIFY_TO_LARK_USER_ID",
6         "NOTIFY_TO_LARK_OPEN_ID",
7         "NOTIFY_TO_LARK_EMAIL",
8     ]
9
10    for each in available_recipient_keys:
11        tmp_recipient = os.environ.get(each)
12        if tmp_recipient:
13            tmp_key_name = each.replace("NOTIFY_TO_LARK_", "").lower()
14            notify_info[tmp_key_name] = tmp_recipient
15            break
16
17    return notify_info

```

Get enabled data mart regions, data mart connections, table names info (for 5 countries: SG, MY, ID, TW, TH). A data mart would have this info as an example:

- 1 - DATAMART_SG_HOST
- 2 - DATAMART_SG_PORT
- 3 - DATAMART_SG_USER
- 4 - DATAMART_SG_PASSWORD
- 5 - DATAMART_SG_DB
- 6 - DATAMART_SG_SCHEMA

```
1 def get_tables(self) -> list:
2     tables = []
3     try:
4         tables = json.loads(os.environ.get("DATAMART_TABLES"))
5     except Exception as exc:
6         raise exc
7     return tables
8
9 def get_enabled_datamart_regions(self) -> list:
10    regions = []
11    try:
12        regions = json.loads(os.environ.get("ENABLED_DATAMART_REGIONS"))
13    except Exception as exc:
14        raise exc
15    return regions
16
17 def get_datamart_connection_info(self, country_code: str) -> dict:
18    conn_info = dict()
19    upper_cased = country_code.upper()
20    attr_host = f"DATAMART_{upper_cased}_HOST"
21    attr_port = f"DATAMART_{upper_cased}_PORT"
22    attr_user = f"DATAMART_{upper_cased}_USER"
23    attr_password = f"DATAMART_{upper_cased}_PASSWORD"
24    attr_dbname = f"DATAMART_{upper_cased}_DB"
25    attr_schemaname = f"DATAMART_{upper_cased}_SCHEMA"
26
27    conn_info["host"] = os.environ.get(attr_host)
28    conn_info["port"] = int(os.environ.get(attr_port, "5432"))
29    conn_info["user"] = os.environ.get(attr_user)
30    conn_info["password"] = os.environ.get(attr_password)
31    conn_info["dbname"] = os.environ.get(attr_dbname)
32    schema_name = os.environ.get(attr_schemaname)
33    conn_info["options"] = f"-c search_path={schema_name}"
34
```

```

35         setattr(self, attr_host, conn_info.get("host"))
36         setattr(self, attr_port, conn_info.get("port"))
37         setattr(self, attr_user, conn_info.get("user"))
38         setattr(self, attr_password, conn_info.get("password"))
39         setattr(self, attr_dbname, conn_info.get("dbname"))
40         setattr(self, attr_schemaname, schema_name)
41
42         return conn_info

```

As our data marts servers get fresh everyday every 3 hours to run through all ETL pipelines from raw data source, we should get info of their refresh, put into a range of [threshold, upperbound] = [0,24]

```

1 def get_datamart_freshness_threshold(self) -> int:
2     return int(os.environ.get("DATAMART_FRESHNESS_THRESHOLD") or 0)
3
4 def get_datamart_freshness_upperbound(self) -> int:
5     return int(os.environ.get("DATAMART_FRESHNESS_UPPERBOUND") or 24)

```

Create connection to server using **psycopg2**

```

1 import psycopg2
2 import requests
3 from psycopg2 import DatabaseError
4 from psycopg2.extras import RealDictCursor
5
6 def create_cursor(db_conn: dict) -> Tuple:
7     cursor_options = dict(cursor_factory=RealDictCursor)
8     conn = psycopg2.connect(**db_conn)
9     cursor = conn.cursor(**cursor_options)
10    return (conn, cursor)

```

@Class Lark Extension : Compile message output context

- Create lark API url, where we parse message result into
- Get Lark API and url
- Define notify message : message, title, chat_id, open_id, user_id, email, message_type, title_template (optional), lines (optional) to provide to Lark, in json format, or else raise error

```

1 def __init__(self, api_info: dict):

```

```

2         self.api_info = api_info
3         self.api_url = self.create_lark_api_url()
4         self.headers = {
5             "Authorization": "Bearer {}".format(self.api_info.get("api_key"))
6         }
7
8     def create_lark_api_url(self) -> str:
9         return "{proto}://{host}:{port}/{endpoint}".format_map(self.api_info)
10
11     def notify(
12         self,
13         message: str,
14         title: str,
15         chat_id: str = None,
16         open_id: str = None,
17         user_id: str = None,
18         email: str = None,
19         message_type: str = "interactive",
20         title_template: str = "orange",
21         lines: int = 100,
22     ) -> Any:
23         if not (chat_id or open_id or user_id or email):
24             raise ValueError(
25                 "Please provide either chat_id or open_id or"
26                 " user_id or email."
27             )
28
29         data = {
30             "message_type": message_type,
31             "message": message,
32             "title": title,
33             "title_template": title_template,
34             "lines": lines,
35         }
36         if chat_id:
37             data["chat_id"] = chat_id
38         elif open_id:
39             data["open_id"] = open_id
40         elif user_id:
41             data["user_id"] = user_id
42         elif email:
43             data["email"] = email
44
45         resp = requests.post(self.api_url, json=data, headers=self.headers)
46         if resp.ok:
47             return resp.json()

```

Get table data summary

We want to check data freshness every day by 2 columns : [table_name]_created_at and [table_name]_updated_at.

We use functions that we already defined:

- Class EnvHandler
- get_datamart_connection_info
- create_cursor

The logic query applied here, for example:

```
1 SELECT 'activities' AS table_name,
2     COUNT(*) AS total,
3     MAX(activity_created_at) AS max_created_at,
4     MAX(activity_updated_at) AS max_updated_at,
5     now() - MAX(activity_created_at) AS freshness_interval
6 FROM activities
7 # table_name = {schema_name}.{table_name}
```

```
1 def get_table_summary(
2     table_name: str, region_name: str, pipe_conn: Any
3 ) -> dict:
4     row = dict()
5     env_handler = EnvHandler()
6     db_conn = env_handler.get_datamart_connection_info(region_name)
7     conn, cursor = create_cursor(db_conn)
8     info_query = (
9         "SELECT column_name FROM information_schema.columns"
10        " WHERE table_catalog = %s"
11        " AND table_schema = %s"
12        " AND table_name = %s"
13        " AND (column_name LIKE '%_created_at'"
14        " OR column_name LIKE '%_updated_at')"
15    )
16    try:
17        db_name = getattr(env_handler, f"DATAMART_{region_name.upper()}_DB")
18        schema_name = getattr(
19            env_handler, f"DATAMART_{region_name.upper()}_SCHEMA"
20        )
21        cursor.execute(info_query, (db_name, schema_name, table_name))
22        info_rows = cursor.fetchall()
```

```

23     # SELECT 'activities' AS table_name,
24     #     COUNT(*) AS total,
25     #     MAX(activity_created_at) AS max_created_at,
26     #     MAX(activity_updated_at) AS max_updated_at,
27     #     now() - MAX(activity_created_at) as freshness_interval
28     # FROM activities
29     col_mapped = {}
30     query = f"SELECT '{table_name}' AS table_name", COUNT(*) AS total"
31     for each_info in info_rows:
32         col_name = each_info.get("column_name")
33         if col_name.endswith("created_at"):
34             query += f", MAX({col_name}) AS max_created_at"
35             query += f", FLOOR(EXTRACT(EPOCH FROM now() - MAX({col_name}))/3
36         elif col_name.endswith("updated_at"):
37             query += f", MAX({col_name}) AS max_updated_at"
38         query += f" FROM {schema_name}.{table_name};"
39         cursor.execute(query)
40         row = cursor.fetchone()
41     except DatabaseError as de:
42         print(f"DatabaseError: {de}")
43     except Exception as ge:
44         print(f"GeneralError: {ge}")
45     import traceback
46
47     traceback.print_exc()
48     finally:
49         conn.close()
50
51     pipe_conn.send(row)
52     pipe_conn.close()

```

Get data mart conditions summary

Functions to be used:

- Class EnvHandler
- get_tables
- get_table_summary

```

1 def summarize_for_datamart(region_name: str) -> list:
2     result = []
3     env_handler = EnvHandler()
4     tables = env_handler.get_tables()
5
6     result = []

```



```

7     processes = []
8     parent_connections = []
9
10    for table_name in tables:
11        parent_conn, child_conn = Pipe()
12        parent_connections.append(parent_conn)
13
14        process = Process(
15            target=get_table_summary,
16            args=(
17                table_name,
18                region_name,
19                child_conn,
20            ),
21        )
22        processes.append(process)
23
24    for each_process in processes:
25        each_process.start()
26
27    for each_process in processes:
28        each_process.join()
29
30    for parent_connection in parent_connections:
31        result.append(parent_connection.recv())
32
33    return

```

Environment set up

The Serverless Application Model Command Line Interface (SAM CLI) is an extension of the AWS CLI that adds functionality for building and testing Lambda applications. It uses Docker to run your functions in an Amazon Linux environment that matches Lambda. It can also emulate your application's build environment and API.

To use the SAM CLI, you need the following tools.

- SAM CLI - [Install the SAM CLI](#)
- Python 3 installed
- Docker - [Install Docker community edition](#)

To build and deploy your application for the first time, run the following in your shell:

```

1 sam build --use-container
2 sam deploy --guided

```

Moving forward, to build our application with the `sam build --use-container` command.

```
1 carro-db-monitoring$ sam build --use-container
```

The SAM CLI installs dependencies defined in `hello_world/requirements.txt`, creates a deployment package, and saves it in the `.aws-sam/build` folder.

Test a single function by invoking it directly with a test event. An event is a JSON document that represents the input that the function receives from the event source. Test events are included in the `events` folder in this project.

Run functions locally and invoke them with the `sam local invoke` command.

```
1 carro-db-monitoring$ sam local invoke CarroDatamartMonitoring --event events/ev
```

The SAM CLI can also emulate your application's API. Use the `sam local start-api` to run the API locally on port 3000.

```
1 carro-db-monitoring$ sam local start-api
2 carro-db-monitoring$ curl http://localhost:3000/
```

The SAM CLI reads the application template to determine the API's routes and the functions that they invoke. The `Events` property on each function's definition includes the route and method for each path.

```
1 Events:
2   CarroDatamartMonitoring:
3     Type: Api
4     Properties:
5       Path: /hello
6       Method: get
```

Fetch, tail, and filter Lambda function logs

To simplify troubleshooting, SAM CLI has a command called `sam logs`. `sam logs` lets you fetch logs generated by your deployed Lambda function from the command line. In addition to

printing the logs on the terminal, this command has several nifty features to help you quickly find the bug.

NOTE : This command works for all AWS Lambda functions; not just the ones you deploy using SAM.

```
1 carro-db-monitoring$ sam logs -n CarroDatamartMonitoring --stack-name carro-db-m
```

You can find more information and examples about filtering Lambda function logs in the [SAM CLI Documentation](#).

Tests

To understand more about AWS SAM, you can run tests folder. Tests are defined in the `tests` folder in this project. Use PIP to install the test dependencies and run tests.

```
1 carro-db-monitoring$ pip install -r tests/requirements.txt --user
2 # unit test
3 carro-db-monitoring$ python -m pytest tests/unit -v
4 # integration test, requiring deploying the stack first.
5 # Create the env variable AWS_SAM_STACK_NAME with the name of the stack we are t
6 carro-db-monitoring$ AWS_SAM_STACK_NAME=<stack-name> python -m pytest tests/inte
```

Cleanup

To delete the sample application that you created, use the AWS CLI. Assuming you used your project name for the stack name, you can run the following:

```
1 aws cloudformation delete-stack --stack-name carro-db-monitoring
```

Resources

See the [AWS SAM developer guide](#) for an introduction to SAM specification, the SAM CLI, and serverless application concepts.

Next, you can use AWS Serverless Application Repository to deploy ready to use Apps that go beyond hello world samples and learn how authors developed their applications: [AWS Serverless Application Repository main page](#)