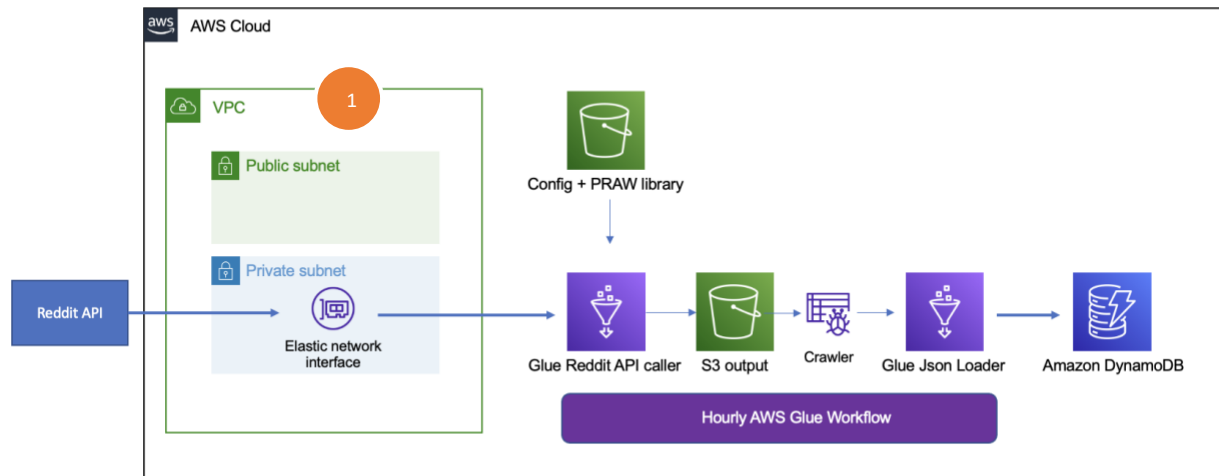# Architecture Diagram



Flow of the Architecture:

1.  We access the Reddit API via AWS VPC that has ENI in the Private subnet, so that any AWS resource that use this VPC would have access outside the public internet, but the public won't have access over our resources.
2.  We created an AWS Glue job, I named it as 'Reddit API Caller' that connects over the VPC and accesses the Reddit API.
3.  The problem is that PRAW library is not native in AWS Glue. So then this library is manually uploaded in an S3 bucket, along with a configuration file that contains all our access identities. The library and the configuration file are then used by the Reddit API Caller on runtime. This is designed this way so that all IDs and Accesses are not exposed on the scripts.
4.  Reddit API Caller then outputs the data into S3, in JSON format. This format is chosen because the fields we use from the response are semi-structured, so it will be problematic when handled in a traditional row-column (CSV) format.
5.  The JSON output is then *crawled* by an AWS Glue crawler which results to a Hive table that is accessible by any AWS resource.
6.  This hive table is then read by another AWS Glue job that loads the data into Dynamodb.
7.  The whole workflow is managed by AWS Glue Workflow. Each job is triggered whenever the task before it succeeded. This is scheduled to run every hour.

Factors in Designing the Architecture
- Serverless
  - I went for a serverless design. Putting this in EC2 server and bootstrapping it is one way to do it, but it would cap the performance acc to the size of the server. And once the server has an issue, the whole process stops. So designing it serverless-ly means microservice and scalability.
- Scalability
  - AWS Glue in PySpark is scalable. Taking advantage of the distributed parallel processing power and its elasticity is beneficial to this project i.e. we can size up or down the DPU of each job anytime we want.
- AWS Cloud Cost Optimization
  - Only paying for when you would run the job, and only upscaling when you need to, and archiving old data in S3, are all ways to make the project cost efficient.
- Semi-structured Data Compatibility
  - I chose JSON and Dynamodb because the response we get are semi-structured. We need the flexibility and ease of use in the future downstream processes so I opted to use Dynamodb.

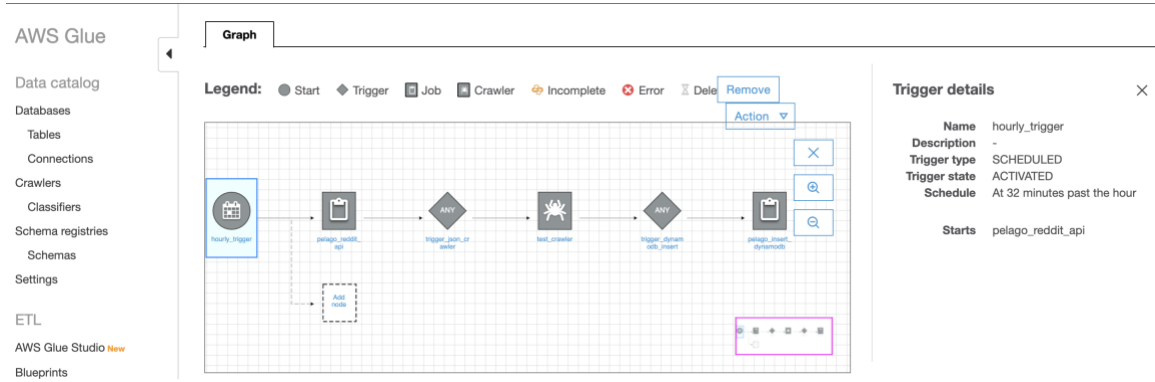GITHUB REPO: https://github.com/jergadi/coding-challenge

REDDIT dev account:

developed applications



create application

Please read the API usage guidelines before creating your application. After creating, you will be required to register for production API use.

AWS RESOURCES:

1. AWS Glue Workflow – hourly scheduled.



2. S3 bucket containing the JSON output records



3. Dynamodb containing the inserted Records

### 4. AWS Glue job Reddit API Caller settings

| pelago_reddit_api | | | Spark | python | s3://aws-glue-scripts-709... | 22 October 2021 7:56 PM UTC+8 | Disable |

History **Details** Script Metrics

| | | | | |
|---|---|---|---|---|
| Name | pelago_reddit_api | | Python lib path | - |
| IAM role | GlueRdsS3Role | | Jar lib path | - |
| Type | Spark | | Other lib path | - |
| Glue version | 2.0 | | Job parameters | --additional-python-modules    s3://aws-glue-kgalife-test/pelago/wheelhouse/praw-7.4.0-py3-none-any.whl |
| Python version | 2.0 | | Non-overrideable Job parameters | - |
| ETL language | python | | | |
| Script location | s3://aws-glue-scripts-709581768673-sa-east-1/root/pelago_reddit_api | | Connections | outside_connection |
| Temporary directory | s3://aws-glue-temporary-709581768673-sa-east-1/root | | Maximum capacity | 2 |

### 5. Aws Glue job JSON loader settings

| ☑ pelago_insert_dynamodb | | | Spark | python | s3://aws-glue-scripts-709... | 22 October 2021 8:12 PM UTC+8 | Disable |
| ☐ pelago_reddit_api | | | Spark | python | s3://aws-glue-scripts-709... | 22 October 2021 7:56 PM UTC+8 | Disable |

History **Details** Script Metrics

| | | | | |
|---|---|---|---|---|
| Name | pelago_insert_dynamodb | | Python lib path | - |
| IAM role | GlueRdsS3Role | | Jar lib path | - |
| Type | Spark | | Other lib path | - |
| Glue version | 2.0 | | Job parameters | - |
| Python version | 3 | | Non-overrideable Job parameters | - |
| ETL language | python | | | |
| Script location | s3://aws-glue-scripts-709581768673-sa-east-1/root/pelago_insert_dynamodb | | Connections | - |
| Temporary directory | s3://aws-glue-temporary-709581768673-sa-east-1/root | | Maximum capacity | 2 |
| Job bookmark | Disable | | Worker type | Standard |

### 6. S3 Bucket containing the config file and PRAW library

Amazon S3 > aws-glue-kgalife-test > pelago/

## pelago/

**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

| | Name | Type | Last modified | Size |
|---|---|---|---|---|
| ☐ | config.json | json | October 22, 2021, 14:01:10 (UTC+08:00) | 20 |
| ☐ | output/ | Folder | - | |
| ☐ | wheelhouse/ | Folder | - | |

Amazon S3 > aws-glue-kgalife-test > pelago/ > wheelhouse/

## wheelhouse/

**Objects** Properties

**Objects (1)**

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 m

| | Name | Type | Last modified |
|---|---|---|---|
| ☐ | praw-7.4.0-py3-none-any.whl | whl | October 22, 2021, 19:33:30 (UTC+08:00) |