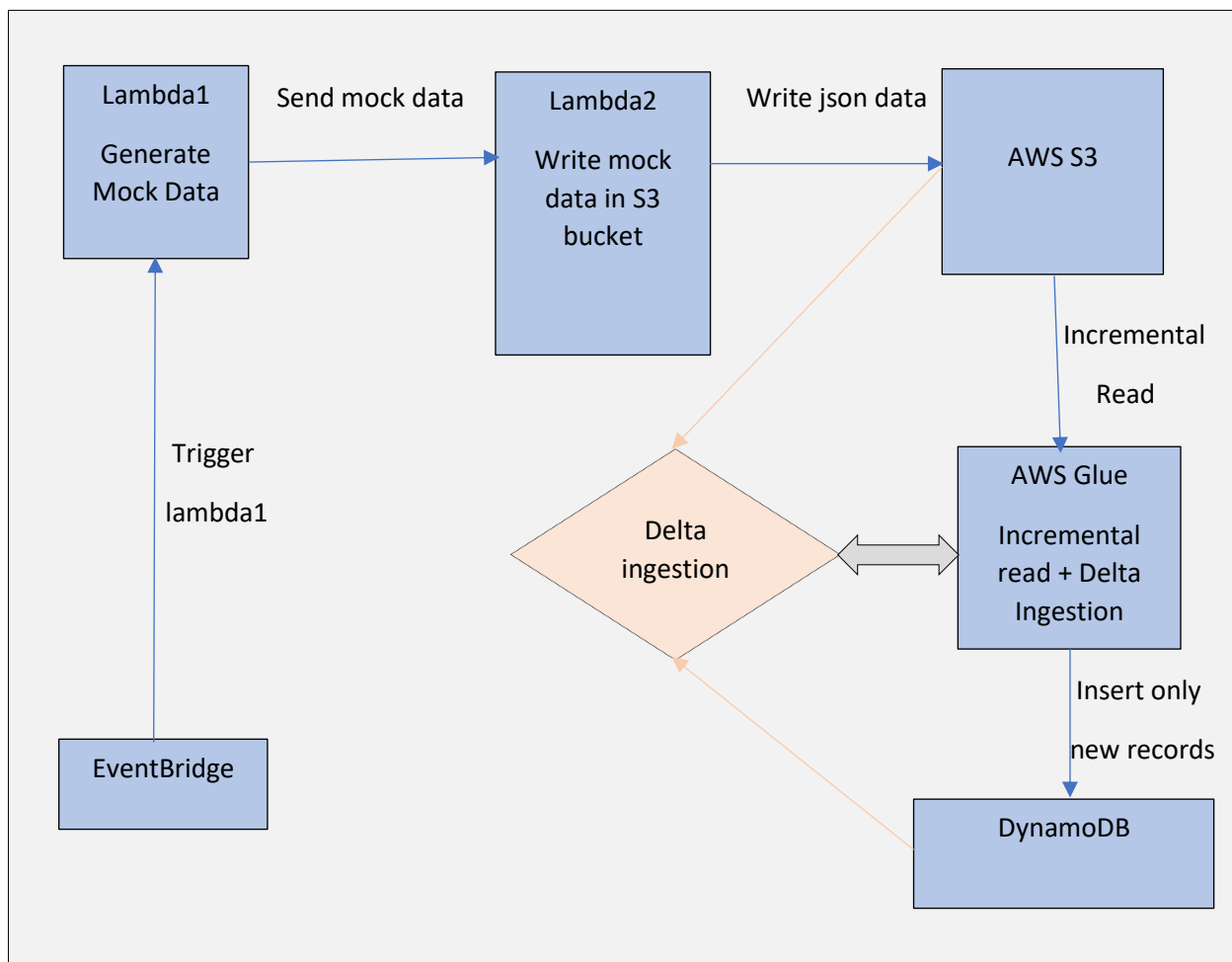# AWS PROJECT 1

# AWS BATCH DATA PIPELINE

_____

## Project Architecture:



## S3:

Create S3 bucket: employee-project1-data

# Create bucket Info

Buckets are containers for data stored in S3. Learn more 

## General configuration

Bucket name

employee-project1-data

Bucket name must be globally unique and must not contain spaces or uppercase letters. See rules for bucket naming 

AWS Region

Asia Pacific (Mumbai) ap-south-1

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.

Choose bucket

## Object Ownership Info

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

**○ ACLs disabled (recommended)**
All objects in this bucket are owned by this account. Access to this bucket and its objects is specified using only policies.

**○ ACLs enabled**
Objects in this bucket can be owned by other AWS accounts. Access to this bucket and its objects can be specified using ACLs.

Object Ownership

Bucket owner enforced

ⓘ **Upcoming permission changes to disable ACLs**
Starting in April 2023, to disable ACLs when creating buckets by using the S3 console, you will no longer need the s3:PutBucketOwnershipControls permission. Learn more 

## Block Public Access settings for this bucket

Public access is granted to buckets and objects through access control lists (ACLs), bucket policies, access point policies, or all. In order to

---

Successfully created bucket "employee-project1-data"
To upload files and folders, or to configure additional bucket settings choose **View details.**

View details

▶ Account snapshot

Storage lens provides visibility into storage usage and activity trends. Learn more 

View Storage Lens dashboard

**Buckets (1)** Info

Buckets are containers for data stored in S3. Learn more 

⟳  ⧉ Copy ARN   Empty   Delete   **Create bucket**

Q Find buckets by name

‹ 1 › ⚙

| Name | AWS Region | Access | Creation date |
|---|---|---|---|
| ○ employee-project1-data | Asia Pacific (Mumbai) ap-south-1 | Bucket and objects not public | February 9, 2023, 12:33:05 (UTC+05:30) |

---

Create folder inside bucket: inbox

# Create folder Info

Use folders to group objects in buckets. When you create a folder, S3 creates an object using the name that you specify followed by a slash (/). This object then appears as folder on the console. Learn more ⬚

ⓘ **Your bucket policy might block folder creation**
If your bucket policy prevents uploading objects without specific tags, metadata, or access control list (ACL) grantees, you will not be able to create a folder using this configuration. Instead, you can use the upload configuration to upload an empty folder and specify the appropriate settings.

## Folder

**Folder name**

inbox

Folder names can't contain "/". See rules for naming ⬚

## Server-side encryption

Server-side encryption protects data at rest.

ⓘ The following settings apply only to the new folder object and not to the objects contained within it.

**Encryption key type** Info

● Amazon S3-managed keys (SSE-S3)
○ AWS Key Management Service key (SSE-KMS)

Cancel | **Create folder**

---

⊘ Successfully created folder "**inbox**".
Operation successfully completed.

Amazon S3 > Buckets > employee-project1-data

# employee-project1-data Info

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

## 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 more ⬚

| C | ⧉ Copy S3 URI | ⧉ Copy URL | ⬇ Download | Open ⬚ | Delete | Actions ▼ | Create folder | **Upload** |

Q Find objects by prefix

| | Name | ▲ | Type | ▽ | Last modified | ▽ | |
|---|---|---|---|---|---|---|---|
| ☐ | 🗋 inbox/ | | Folder | | - | | |

**Lambda:**

# Create lambda function 1 to generate mock data: mock-data-generator

common use cases.

## Basic information

**Function name**
Enter a name that describes the purpose of your function.

mock-data-generator

Use only letters, numbers, hyphens, or underscores with no spaces.

**Runtime** Info
Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

Python 3.9

**Architecture** Info
Choose the instruction set architecture you want for your function code.

- ● x86_64
- ○ arm64

## Permissions Info

By default, Lambda will create an execution role with permissions to upload logs to Amazon CloudWatch Logs. You can customize this default role later when adding triggers.
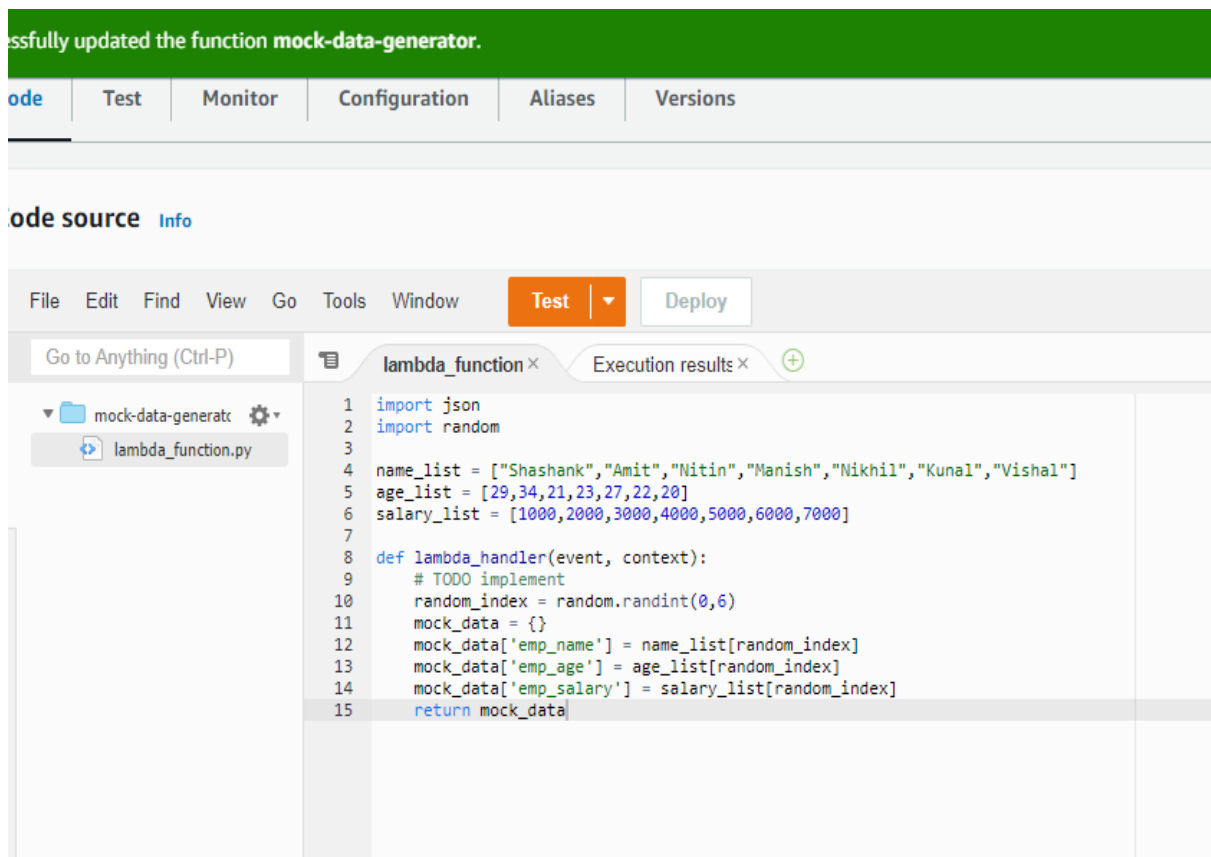
▼ Change default execution role

**Execution role**
Choose a role that defines the permissions of your function. To create a custom role, go to the IAM console ↗.

- ● Create a new role with basic Lambda permissions
- ○ Use an existing role
- ○ Create a new role from AWS policy templates

ⓘ Role creation might take a few minutes. Please do not delete the role or edit the trust or permissions policies in this role.

| ode | Test | Monitor | Configuration | Aliases | Versions |

## ode source  Info

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

Go to Anything (Ctrl-P)

lambda_function ×        Execution results ×    ⊕

▼ 📁 mock-data-generato ⚙▾
  ◇ lambda_function.py

```
1  import json
2  import random
3
4  name_list = ["Shashank","Amit","Nitin","Manish","Nikhil","Kunal","Vishal"]
5  age_list = [29,34,21,23,27,22,20]
6  salary_list = [1000,2000,3000,4000,5000,6000,7000]
7
8  def lambda_handler(event, context):
9      # TODO implement
10     random_index = random.randint(0,6)
11     mock_data = {}
12     mock_data['emp_name'] = name_list[random_index]
13     mock_data['emp_age'] = age_list[random_index]
14     mock_data['emp_salary'] = salary_list[random_index]
15     return mock_data
```

Create lambda function 2 to write data inside S3 bucket's folder "inbox": send-data-to-S3

## Create function  Info

AWS Serverless Application Repository applications have moved to Create application.

| Author from scratch ● | Use a blueprint ○ | Container ima |
|---|---|---|
| Start with a simple Hello World example. | Build a Lambda application from sample code and configuration presets for common use cases. | Select a containe |

**Basic information**

Function name
Enter a name that describes the purpose of your function.

> send-data-to-S3

Use only letters, numbers, hyphens, or underscores with no spaces.

Runtime  Info
Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

> Python 3.9                                                                 ▼

Architecture  Info
Choose the instruction set architecture you want for your function code.

● x86_64
○ arm64

Permissions  Info
By default, Lambda will create an execution role with permissions to upload logs to Amazon CloudWatch Logs. You can customize this default role later when adding triggers.

▼ **Change default execution role**

Execution role

# Add IAM role to allow lambda function to access S3 bucket:

Function URL  Info

-

| Code | Test | Monitor | Configuration | Aliases | Versions |

**General configuration**

Triggers

**Permissions**

Destinations

Function URL

Environment variables

Tags

VPC

Monitoring and operations tools

Concurrency

Asynchronous invocation

Code signing

Database proxies

File systems

## Execution role

Role name
send-data-to-S3-role-ygtp55bb ↗

## Resource summary

Vie

Amazon CloudWatch Logs
3 actions, 2 resources

To view the resources and actions that your function has permission to access, choose a service.

| By action | **By resource** |

| Resource | Actions |
| --- | --- |
| arn:aws:logs:ap-south-1:788659805261:* | Allow: logs:CreateLogGroup |
| arn:aws:logs:ap-south-1:788659805261:log-group:/aws/lambda/send-data-to-S3:* | Allow: logs:CreateLogStream<br>Allow: logs:PutLogEvents |

# Edit basic settings

## Basic settings  Info

Description - *optional*

Memory  **Info**
Your function is allocated CPU proportional to the memory configured.

128  MB

Set memory to between 128 MB and 10240 MB

### Ephemeral storage  **Info**
You can configure up to 10 GB of ephemeral storage (/tmp) for your function. View pricing

512  MB

Set ephemeral storage (/tmp) to between 512 MB and 10240 MB.

### Timeout

5  min   0  sec

### Execution role
Choose a role that defines the permissions of your function. To create a custom role, go to the IAM console .

○ **Use an existing role**
○ Create a new role from AWS policy templates

### Existing role
Choose an existing role that you've created to be used with this Lambda function. The role must have permission to upload logs to Amazon CloudWatch Logs.

service-role/send-data-to-S3-role-ygtp55bb  ▼

View the send-data-to-S3-role-ygtp55bb role  on the IAM console.

ℹ **Introducing the new IAM roles experience**
We've redesigned the IAM roles experience to make it easier to use. Let us know what you think.

IAM > Roles > send-data-to-S3-role-ygtp55bb

# send-data-to-S3-role-ygtp55bb

Delete

## Summary

Edit

| Creation date | ARN |
|---|---|
| February 09, 2023, 12:51 (UTC+05:30) | 🗐 arn:aws:iam::788659805261:role/service-role/send-data-to-S3-role-ygtp55bb |
| **Last activity** | **Maximum session duration** |
| None | 1 hour |

| **Permissions** | Trust relationships | Tags | Access Advisor | Revoke sessions |
|---|---|---|---|---|

### Permissions policies (1) Info
You can attach up to 10 managed policies.

⟳  Simulate  Remove  Add permissions ▲

Attach policies
Create inline policy

🔍 Filter policies by property or policy name and press enter.

| ☐ | Policy name 🗗 ▽ | Type ▽ | Description |
|---|---|---|---|
| ☐ | ⊞ AWSLambdaBasicExecutionRole-23b053c7-ee28-49ba-8478-bf2597345100 | Customer managed | |

### Permissions boundary - (not set) Info
Set a permissions boundary to control the maximum permissions this role can have. This is not a common setting but can be used to delegate permission management to others.

---

**Introducing the new IAM roles experience**
We've redesigned the IAM roles experience to make it easier to use. Let us know what you think.

IAM > Roles > send-data-to-S3-role-ygtp55bb > Add permissions

## Attach policy to send-data-to-S3-role-ygtp55bb

▶ **Current permissions policies** (1)

### Other permissions policies (Selected 1/811)

🔍 Filter policies by property or policy name and press enter.     1 match

"S3FULLACCESS" ✕    **Clear filters**

| ☑ | Policy name 🗗 | ▽ | T |
|---|---|---|---|
| ☑ | ⊞ 🎁 AmazonS3FullAccess | | A |

IAM > Roles > send-data-to-S3-role-ygtp55bb

# send-data-to-S3-role-ygtp55bb

Delete

## Summary

Edit

| | |
|---|---|
| **Creation date** | **ARN** |
| February 09, 2023, 12:51 (UTC+05:30) | ⧉ arn:aws:iam::788659805261:role/service-role/send-data-to-S3-role-ygtp55bb |
| **Last activity** | **Maximum session duration** |
| None | 1 hour |

**Permissions** | Trust relationships | Tags | Access Advisor | Revoke sessions

### Permissions policies (2) Info
You can attach up to 10 managed policies.

🔍 Filter policies by property or policy name and press enter.

🔄 | Simulate | Remove | Add permissions ▼

‹ 1 › ⚙

| ☐ | Policy name ⧉ | | Type | | Description |
|---|---|---|---|---|---|
| ☐ | ⊞ AWSLambdaBasicExecutionRole-23b053c7-ee28-49ba-8478-bf2597345100 | ▽ | Customer managed | ▽ | |
| ☐ | ⊞ 🛡 AmazonS3FullAccess | | AWS managed | | Provides full access to all buckets via the AWS Manage... |

### Permissions boundary - (not set) Info
Set a permissions boundary to control the maximum permissions this role can have. This is not a common setting but can be used to delegate

---

**✓ Successfully updated the function send-data-to-S3.**

**Code** | Test | Monitor | Configuration | Aliases | Versions

## Code source Info

△ | File | Edit | Find | View | Go | Tools | Window | **Test** ▼ | Deploy | Changes not deployed

🔍 Go to Anything (Ctrl-P) | 🗏 | lambda_function × | ⊕

Environment
▼ 📁 send-data-to-S3 ⚙▼
  ‹› lambda_function.py

```python
1   import json
2   import boto3
3   import json
4   import datetime
5
6   def lambda_handler(event, context):
7       print(event['responsePayload'])
8       employee_data = event['responsePayload']
9       BUCKET_NAME = "employee-project1-data"
10      current_epoch_time = datetime.datetime.now().timestamp()
11
12      print("Start Data Write in S3")
13      s3 = boto3.resource('s3')
14      s3object = s3.Object(BUCKET_NAME, f"inbox/{str(current_epoch_time)}_employee_data.json")
15      s3object.put(
16          Body=(bytes(json.dumps(employee_data).encode('UTF-8')))
17      )
18      print("Data Write Successfull in S3")
19
```

Epoch time is used here to give unique names to files generated by this lambda function.

An epoch timestamp is a way of representing a specific point in time using a number. This numerical value can be used for various date and time operations, making it easier to compare and manipulate different dates and times.

# Add destination to lambda function 1 "mock-data-generator"

"send-data-to-S3" will be the destination for "mock-data-generator"

Lambda > Functions > mock-data-generator > Add destination

## Add destination

### Destination configuration
Send invocation records to a destination when your function is invoked asynchronously, or if your function processes records from a stream.

**Source**
The type of invocation that maps to the destination.
- ● Asynchronous invocation
- ○ Stream invocation

**Condition**
The condition for using the destination.
- ○ On failure
- ● On success

**Destination type**
An SQS queue, SNS topic, Lambda function, or EventBridge event bus.

| Lambda function | ▼ |

ⓘ Your function's execution role doesn't have permission to send result to the destination. By clicking save we'll attempt to add permission to the role for you.

**Destination**

🔍 arn:aws:lambda:ap-south-1:788659805261:function:send-data-to-S3        ✕        ↻

Cancel        **Save**

Lambda > Functions > mock-data-generator

## mock-data-generator

Throttle | Copy ARN | Actions ▼

▼ Function overview  Info

mock-data-generator

⊗ Layers                                    (0)

+ Add trigger

λ AWS Lambda

+ Add destination

Description
-

Last modified
5 minutes ago

Function ARN
arn:aws:lambda:ap-south-1:788659805261:function:mock-data-generator

Function URL  Info
-

Code | Test | Monitor | Configuration | Aliases | Versions
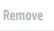
General configuration
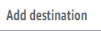
Triggers

Permissions

Destinations

Function URL

### Destinations  Info

Remove | Edit | Add destination

Q Find destinations

| Source | Stream | Condition | Destination |
|--------|--------|-----------|-------------|
| ○ Asynchronous invocation | - | On success | arn:aws:lambda:ap-south-1:788659805261:function:send-data-to-S3 |

## Eventbridge:

Create eventbridge rule: schedule mock-data-generator function to trigger.

Step 1
**Define rule detail**

Step 2
Define schedule

Step 3
Select target(s)

Step 4 - *optional*
Configure tags

Step 5
Review and create

# Define rule detail Info

## Rule detail

Name

mock-data-generator-trigger

Maximum of 64 characters consisting of numbers, lower/upper case letters, .,-,_.

Description - *optional*

Enter description

Event bus | Info
Select the event bus this rule applies to, either the default event bus or a custom or partner event bus.

default ▼

🔵 Enable the rule on the selected event bus

Rule type | Info

○ **Rule with an event pattern**
A rule that runs when an event matches the defined event pattern. EventBridge sends the event to the specified target.

🔵 **Schedule**
A rule that runs on a schedule

Cancel    **Next**

---

Step 1
Define rule detail

Step 2
**Define schedule**

Step 3
Select target(s)

Step 4 - *optional*
Configure tags

Step 5
Review and create

# Define schedule Info

## Schedule pattern

Schedule pattern
Choose the schedule type that best meets your needs.

○ A fine-grained schedule that runs at a specific time, such as 8:00 a.m. PST on the first Monday of every month.

🔵 A schedule that runs at a regular rate, such as every 10 minutes.

Rate expression | Info
Enter a value and the unit of time to run the schedule.

rate ( 3    Minutes ▼ )

Value    Unit, e.g. mins, hours...

Cancel    Previous    **Next**

Amazon EventBridge > Rules > Create rule

# Select target(s)

ℹ️ **Permissions**
Note: When using the EventBridge console, EventBridge will automatically configure the proper permissions for the selected targets. If you're using the AWS CLI, SDK, or CloudFormation, you'll need to configure the proper permissions.

## Target 1

**Target types**
Select an EventBridge event bus, EventBridge API destination (SaaS partner), or another AWS service as a target.

○ EventBridge event bus
○ EventBridge API destination
● AWS service

**Select a target**  Info
Select target(s) to invoke when an event matches your event pattern or when schedule is triggered (limit of 5 targets per rule)

| Lambda function ▼ |
| --- |

**Function**

| mock-data-generator ▼ | ⟳ |
| --- | --- |

▶ Configure version/alias

▶ Additional settings

| Add another target | | Cancel | Skip to Review and create | Previous | **Next** |

---

✓ Rule mock-data-generator-trigger was created successfully

Amazon EventBridge > Rules

# Rules

A rule watches for specific types of events. When a matching event occurs, the event is routed to the targets associated with the rule. A rule can be associated with one or more targets.

## Select event bus

**Event bus**
Select or enter event bus name

| default ▼ |
| --- |

**Rules** (1/1)  ⟳  Delete  Enable  Edit  CloudFormation Template ▼  **Create rule**

🔍 Find rules     Any status ▼                    ⟨ 1 … ⟩ ⚙

| ☐ | Name | ▲ | Status | ▽ | Type | ▽ | Description |
| --- | --- | --- | --- | --- | --- | --- | --- |
| ☐ | mock-data-generator-trigger | | ⊘ Enabled | | Scheduled Standard | | |

Enable the rule and check logs:



Check data inside S3:



-------**Mock data is generated and written successfully inside S3**-------

**Glue:**

Create crawler: employee-json-data-crawler

Add data source to crawl:

**Choose S3 path**

S3 buckets > employee-project1-data

**Objects (1/1)**

🔍 Find object by prefix

| | Key | ▽ | Last modified |
|---|---|---|---|
| 🔘 | 📄 inbox | | 2023-02-09T07:03:51.000Z |

🔵 Crawl all sub-folders
Crawl all folders again with every subsequent crawl.

## Add data source ✕

**Data source**
Choose the source of data to be crawled.

| S3 | ▼ |
|---|---|

**Network connection - *optional***
Optionally include a Network connection to use with this S3 target. Note that each crawler is limited to one Network connection so any other S3 targets will also use the same connection (or none, if left blank).

| | ▼ | | ↻ |
|---|---|---|---|

| Clear selection | | Add new connection ↗ |
|---|---|---|

**Location of S3 data**
🔵 In this account
⚪ In a different account

**S3 path**
Browse for or enter an existing S3 path.

| 🔍 s3://employee-project1-data/inbox | ✕ | View ↗ | Browse |
|---|---|---|---|

All folders and files contained in the S3 path are crawled. For example, type s3://MyBucket/MyFolder/ to crawl all objects in MyFolder within MyBucket.

**Subsequent crawler runs**
This field is a global field that affects all S3 data sources.
🔵 Crawl all sub-folders
Crawl all folders again with every subsequent crawl.
⚪ Crawl new sub-folders only
Only Amazon S3 folders that were added since the last crawl will be crawled. If the schemas are compatible, new partitions will be added to existing tables.
⚪ Crawl based on events
Rely on Amazon S3 events to control what folders to crawl.

☐ Sample only a subset of files

☐ Exclude files matching pattern

Cancel      **Add an S3 data source**

Add classifier: Crawler will crawl the files present in selected data source using this classifier.

# Create classifier Info

## Classifier details

### Classifier name

```
employee-json-data-parser
```

Name may contain letters (A-Z), numbers (0-9), hyphens (-), or underscores (_), and can be up to 255 characters long.

## Classifier type and properties

### Classifier type

○ **Grok**
Best for parsing unstructured text (e.g., application logs).

○ **XML**
Extract data out of XML documents.

● **JSON**
Extract fields out of JSON files.

○ **CSV**
Filter and extract data out of CSV files.

### JSON path

```
$.emp_name,$.emp_age,$emp_salary
```

The JSON path expression defines a JSON structure and is used to define a table schema.

Cancel    Create

---

## Classifiers

Classifiers are triggered during a crawl task. A classifier checks whether a given file is in a format the crawler can handle. If it is, the classifier creates a schema in the form of a StructType object that matches that data format.

### Classifiers (1/1) Info
View and manage all available classifiers.

Last updated (UTC)
February 9, 2023 at 07:43:16    Edit    Delete    **Add classifier**

Filter classifiers

< 1 >

| | Name | Type ▽ | Classification ▽ | Last updated (UTC) ▽ |
|---|---|---|---|---|
| ☑ | employee-json-data-parser | JSON | - | February 9, 2023 at 07:43:12 |

## Choose data sources and classifiers

### Data source configuration

Is your data already mapped to Glue tables?

- ● **Not yet**
  Select one or more data sources to be crawled.

- ○ **Yes**
  Select existing tables from your Glue Data Catalog.

### Data sources (1)  Info

The list of data sources to be scanned by the crawler.

[Edit]  [Remove]  [Add a data source]

| | Type | Data source | Parameters |
|---|---|---|---|
| ○ | S3 | s3://employee-project1-data/inbox | Recrawl all |

▼ **Custom classifiers - optional**

A classifier checks whether a given file is in a format the crawler can handle. If it is, the classifier creates a schema in the form of a StructType object that matches that data format.

**Custom classifiers**   Info

Select one or more classifiers to use with this crawler.

| Choose one or more classifiers | ▲ |
|---|---|

☑ employee-json-data-parser

[Clear selection]  [Add new classifier ↗]

Cancel    [Previous]    [Next]

---

Only IAM roles created by the AWS Glue console and have the prefix "AWSGlueServiceRole-" can be updated.

Lake F

Allow the

Use L

Check
only t

## Create new IAM role                    ✕

Enter new IAM role

AWSGlueServiceRole-employee-json-parser

Cancel    [Create]

▶ Security configuration - optional

Enable at-rest encryption with a security configuration.

vler

## Set output and scheduling

### Output configuration  Info

Target database

Choose a database

Clear selection    Add database ⧉

⚠ Target database is required

Table name prefix - *optional*

Type a prefix added to table names

Maximum table threshold - *optional*

This field sets the maximum number of tables the crawler is allowed to generate. In the event that this number is surpassed, the crawl will fail with an error. If not set, the crawler will automatically generate the number of tables depending on the data schema.

Type a number greater than 0

▶ Advanced options

### Crawler schedule

You can define a time-based schedule for your crawlers and jobs in AWS Glue. The definition of these schedules uses the Unix-like cron ⧉ syntax. Learn more ⧉.

Frequency

On demand

Cancel    Previous    Next

Create database: crawler will create table inside this database to store metadata.

AWS Glue › Databases › Add new database

## Create a database

Create a database in the AWS Glue Data Catalog.

### Database details

Name

employee_data_db

Database name is required, in lowercase characters, and no longer than 255 characters.

Location - *optional*

Set the URI location for use by clients of the Data Catalog.

Description - *optional*

Enter text

Descriptions can be up to 2048 characters long.

Give prefix to the table name that crawler will create for storing metadata -> this prefix will help in identifying the required table.

wler

## Set output and scheduling

### Output configuration Info

Target database

employee_data_db ▼

[ Clear selection ]  [ Add database 🔗 ]

Table name prefix - *optional*

employee_data_json

Maximum table threshold - *optional*

This field sets the maximum number of tables the crawler is allowed to generate. In the event that this number is surpassed, the crawl will fail with an error. If not set, the crawler will automatically generate the number of tables depending on the data schema.

Type a number greater than 0

▶ Advanced options

### Crawler schedule

You can define a time-based schedule for your crawlers and jobs in AWS Glue. The definition of these schedules uses the Unix-like cron 🔗 syntax. Learn more 🔗.

Frequency

On demand ▼

[ Cancel ]  [ Previous ]

---

AWS Glue > Crawlers > employee-json-data-crawler

# employee-json-data-crawler

Last updated (UTC)
February 9, 2023 at 07:47:41

[ C ]  [ Run craw ]

### Crawler properties

| Name | IAM role | Database | State |
|---|---|---|---|
| employee-json-data-crawler | AWSGlueServiceRole-employee-json-parser 🔗 | employee_data_db | READY |

| Description | Security configuration | Lake Formation configuration | Table prefix |
|---|---|---|---|
| - | - | - | employee_data_json |

Maximum table threshold
-

▶ Advanced settings

**Crawler runs** | Schedule | Data sources | Classifiers | Tags

### Crawler runs (0)
The list of crawler runs for this crawler.

[ C ]  [ Stop run ]  [ View CloudWatch logs ]

🔍 Filter data

🗓 Filter by a date and time range

| Start time (UTC) ▲ | End time (UTC) ▽ | Current/last duration ▽ | Status ▽ | DPU hours ▽ | Table changes |
|---|---|---|---|---|---|

You don't have any crawler runs.

[ Run crawler ]

---

# Run crawler:

AWS Glue > Crawlers

# 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** (1/1) Info
View and manage all available crawlers.

Last updated (UTC)
February 9, 2023 at 07:51:53

Action ▼    Run    Cr

Q Filter crawlers

| | Name | ▽ | State | ▽ | Schedule | | Last run | ▽ | Last run timestamp | ▽ | Log | | Table chang |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ✓ | employee-json-data-crawler | | ⊘ Ready | | | | ⊘ Succeeded | | February 9, 2023 at 07:49:53 | | View log ↗ | | 1 created |

## This is the table created by crawler:

AWS Glue > Tables > employee_data_jsoninbox

# employee_data_jsoninbox

Last updated (UTC)
February 9, 2023 at 07:50:36

Version 0 (Current ve

**Table details**    **Advanced properties**

| Name | Description | Database | Classification |
|---|---|---|---|
| employee_data_jsoninbox | - | employee_data_db | json |

| Location | Connection | Deprecated | Last updated |
|---|---|---|---|
| s3://employee-project1-data/inbox/ | - | - | February 9, 2023 at 07:50:36 |

| Input format | Output format | Serde serialization lib | |
|---|---|---|---|
| org.apache.hadoop.mapred.TextInputFormat | org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat | org.openx.data.jsonserde.JsonSerDe | |

**Schema**    **Partitions**    **Indexes**

**Schema** (3)
View and manage the table schema.

Edit schema as JS

Q Filter schemas

| # | ▽ | Column name | ▽ | Data type | ▽ | Partition key | ▽ | Comment |
|---|---|---|---|---|---|---|---|---|
| 1 | | emp_name | | string | | - | | - |
| 2 | | emp_age | | int | | - | | - |
| 3 | | emp_salary | | int | | - | | - |

## Glue job:

## Select spark script

# Jobs Info

## Create job Info

| ○ Visual with a source and target | ○ Visual with a blank canvas | ● Spark script editor |
| --- | --- | --- |
| Start with a source, ApplyMapping transform, and target. | Author using an interactive visual interface. | Write or upload your own Spark code. |

| ○ Python Shell script editor | ○ Jupyter Notebook |
| --- | --- |
| Write or upload your own Python shell script. | Write your own code in a Jupyter Notebook for interactive development. |

## Options Info

● Create a new script with boilerplate code
○ Upload and edit an existing script
  Choose a local file.

---

## glue_to_dynamodb_ingestion  ☑

**Script**  |  **Job details**  |  Runs  |  Schedules  |  **Version Control**

### Script  Info

```
1   import sys
2   from awsglue.transforms import *
3   from awsglue.utils import getResolvedOptions
4   from pyspark.context import SparkContext
5   from awsglue.context import GlueContext
6   from awsglue.dynamicframe import DynamicFrame
7   from awsglue.job import Job
8
9   ## @params: [JOB_NAME]
10  args = getResolvedOptions(sys.argv, ['JOB_NAME'])
11
12  sc = SparkContext()
13  glueContext = GlueContext(sc)
14  spark = glueContext.spark_session
15  job = Job(glueContext)
16  job.init(args['JOB_NAME'], args)
17
18  empDf = glueContext.create_dynamic_frame.from_catalog(
19          database="employee_data_db",
20          table_name="employee_data_jsoninbox",
21          transformation_ctx = "s3_employee_new_json"
22          )
23
24  incrementalEmpDf = empDf.toDF()
25  print(incrementalEmpDf.count())
26
27▼ if incrementalEmpDf.count() == 0:
28      print("No New records were received, Do not ingest anything into DynamoDb")
```

Enable bookmark: for incremental read

## glue_to_dynamodb_ingestion ☑

**Script** | **Job details** | Runs | Schedules | **Version Control**

AWSGlueServiceRole-Glue ▼ ⟳

**Type**
The type of ETL job. This is set automatically based on the types of data sources you have selected.

Spark ▼

**Glue version**  Info

Glue 3.0 - Supports spark 3.1, Scala 2, Python 3 ▼

**Language**

Python 3 ▼

**Worker type**
Set the type of predefined worker that is allowed when a job runs.

G 1X
(4vCPU and 16GB RAM) ▼

☐ **Automatically scale the number of workers**
AWS Glue will optimize costs and resource usage by dynamically scaling the number of workers up and down throughout the job run. Requires Glue 3.0 or later.

**Requested number of workers**
The number of workers you want AWS Glue to allocate to this job.

2

☑ **Generate job insights**
AWS Glue will analyze your job runs and provide insights on how to optimize your jobs and the reasons for job failures.

**Job bookmark**  Info
Specifies how AWS Glue processes job bookmark when the job runs. It can remember previously processed data (Enable), update state information (Pause), or ignore state information (Disable)

# IAM role for glue job:

ce to make it easier to use. **Let us know what you think**.

## Select trusted entity  Info

**Trusted entity type**

| ⦿ AWS service | ○ AWS account | ○ Web identity |
|---|---|---|
| Allow AWS services like EC2, Lambda, or others to perform actions in this account. | Allow entities in other AWS accounts belonging to you or a 3rd party to perform actions in this account. | Allows users federated by the specified external web identity provider to assume this role to perform actions in this account. |

| ○ SAML 2.0 federation | ○ Custom trust policy |
|---|---|
| Allow users federated with SAML 2.0 from a corporate directory to perform actions in this account. | Create a custom trust policy to enable others to perform actions in this account. |

**Use case**
Allow an AWS service like EC2, Lambda, or others to perform actions in this account.

Common use cases

○ EC2
Allows EC2 instances to call AWS services on your behalf.

○ Lambda
Allows Lambda functions to call AWS services on your behalf.

Use cases for other AWS services:

Glue ▼

⦿ Glue
Allows Glue to call AWS services on your behalf.

Step 2: Add permissions

Permissions policy summary

| Policy name ⧉ | Type | Attached as |
|---|---|---|
| AWSGlueServiceRole-employee-json-parser-EZCRC-s3Policy | Customer managed | Permissions policy |
| AWSGlueConsoleFullAccess | AWS managed | Permissions policy |
| AWSGlueServiceRole | AWS managed | Permissions policy |
| AmazonDynamoDBFullAccess | AWS managed | Permissions policy |
| AmazonS3FullAccess | AWS managed | Permissions policy |

Tags

**Add tags** - *optional*   Info
Tags are key-value pairs that you can add to AWS resources to help identify, organize, or search for resources.

No tags associated with the resource.

[ Add tag ]

You can add up to 50 more tags.

Cancel     Prev

---

# glue_to_dynamodb_ingestion  ✎

**Script**  |  **Job details**  |  Runs  |  Schedules  |  **Version Control**

## Basic properties  Info

### Name

[ glue_to_dynamodb_ingestion ]

### Description - *optional*

[                                                                      ]

### IAM Role
Role assumed by the job with permission to access your data stores. Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job.

[ AWSGlueServiceRole-Glue                         ▼ ]   [ ↻ ]

### Type
The type of ETL job. This is set automatically based on the types of data sources you have selected.

[ Spark                                              ▼ ]

### Glue version   Info

[ Glue 3.0 - Supports spark 3.1, Scala 2, Python 3    ▼ ]

### Language

[ Python 3                                            ▼ ]

### Worker type

**glue_to_dynamodb_ingestion** ☑

Script | **Job details** | Runs | Schedules | Version Control

Libraries Info
Python library path

Dependent JARs path

Referenced files path

Job parameters Info
No parameters associated with the resource.

Add new parameter

You can add 50 more parameters.

Tags

Key
🔍 --JOB_NAME ✕

Value - *optional*
🔍 test ✕    Remove
Custom tag value

Add new tag

You can add 49 more tags.

## DynamoDB:

Create DynamoDB table:

DynamoDB > Tables > Create table

# Create table

## Table details  Info

DynamoDB is a schemaless database that requires only a table name and a primary key when you create the table.

**Table name**
This will be used to identify your table.

employee_data

Between 3 and 255 characters, containing only letters, numbers, underscores (_), hyphens (-), and periods (.).

**Partition key**
The partition key is part of the table's primary key. It is a hash value that is used to retrieve items from your table and allocate data across hosts for scalability and availability.

emp_name                                String ▼

1 to 255 characters and case sensitive.

**Sort key - optional**
You can use a sort key as the second part of a table's primary key. The sort key allows you to sort or search among all items sharing the same partition key.

emp_salary                              Number ▼

1 to 255 characters and case sensitive.

## Table settings

◉ Default settings          ○ Customize settings

# Run Glue Job:

**glue_to_dynamodb_ingestion**                                    Last modified on 2/9/2023, 1:35:49 PM    Actions ▼

Script    Job details    Runs    Schedules    Version Control

**Recent job runs** Info (1)

🔍 Filter job runs by property

**February 09, 2023 1:39:40 PM**                                                                                    Rewind jo

| Job name | Id | Run status | Glue version |
|---|---|---|---|
| glue_to_dynamodb_ingestion | jr_7db35e710e4ae4d26997bf2bc460723144484ee996314843c5 ed5aa513068722 | ⊘ Succeeded | 3.0 |

| Retry attempt number | Start time | End time | Start-up time |
|---|---|---|---|
| Initial run | February 09, 2023 1:39:40 PM | February 09, 2023 1:41:04 PM | 8 seconds |

| Execution time | Last modified on | Trigger name | Security configuration |
|---|---|---|---|
| 1 minute 15 seconds | February 09, 2023 1:41:04 PM | - | - |

| Timeout | Max capacity | Number of workers | Worker type |
|---|---|---|---|
| 2880 minutes | 2 DPUs | 2 | G.1X |

| Execution class | Log group name | Cloudwatch logs | Performance and debugging recommendations |
|---|---|---|---|
| STANDARD | /aws-glue/jobs | ○ All logs ↗<br>○ Output logs ↗<br>○ Error logs ↗ | ○ View in CloudWatch ↗ |

▶ **Input arguments (9)**
Arguments used when this job run was executed.

## Check the data in DynamoDB:

**DynamoDB** ✕

Dashboard
Tables
    Update settings
    Explore items
**PartiQL editor** New
Backups
Exports to S3
Imports from S3 New
Reserved capacity
Settings New

▼ **DAX**
Clusters
Subnet groups
Parameter groups
Events

Tables (1)

🔍 Find tables

‹ 1 › ⚙

▶ employee_data ...

```
1 select * from employee_data;
```

**Run**  **Clear**

**Table view**  **JSON view**

⊘ Completed
Started on 2/9/2023, 1:41:34 PM
Elapsed time 479ms

**Items returned** (2)

🔍 Find items

| emp_salary ▽ | emp_age ▽ | emp_name |
|---|---|---|
| 7000 | 20 | Vishal |
| 2000 | 34 | Amit |