



Serverless x GenAI BKK Workshop

Generative AI at scale: Serverless workflows for enterprise-ready apps

[Pre-requisite] Enable foundation model access in Amazon Bedrock

[Pre-requisite] Configuring the front-end application

► Playground

▼ Use cases

▼ Building a RAG pipeline

Testing the need for RAG

Building the RAG Pipeline

Running the Pipeline

Testing the RAG Inference

High level Code Walkthrough

Summary

► Document extraction and summarization

► Workshop Cleanup

▼ AWS account access

[Open AWS console \(us-west-2\)](#)

[Get AWS CLI credentials](#)

Exit event

[Event dashboard](#) > [Use cases](#) > [Building a RAG pipeline](#) > [Running the Pipeline](#)

© 2008 - 2025, Amazon Web Services, Inc. or its affiliates. All rights reserved.

[Privacy policy](#)

[Terms of use](#)

[Cookie preferences](#)

In this section, you are going to test the RAG pipeline created earlier. To demonstrate a very important feature of multiple files processing, we have introduced error code in a Lambda function. When you test the pipeline, one of the iterations will fail. You will fix the Lambda code and restart the pipeline from the point of failure rather than processing all the files all over again.

Running the workflow

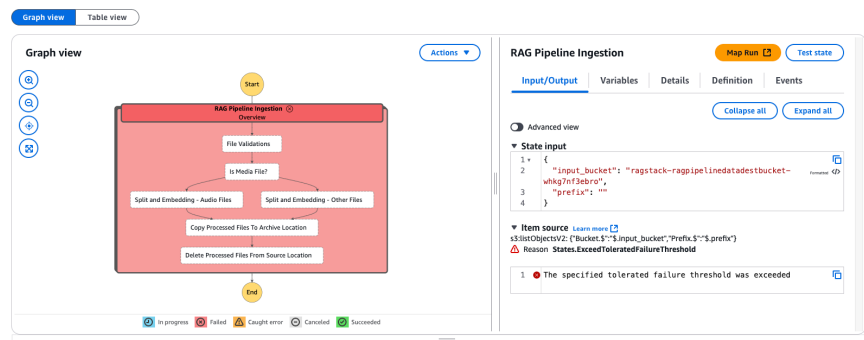
1. Open the [S3 console](#) in a different tab and search for `ragpipelinedatatestbucket` and note down the name.
2. Return to Step Functions console. Select **Execute** button
3. Paste the following input in the dialog that appears. Make sure to substitute the `{{INPUT_BUCKET}}` with the S3 bucket name in step 1. Select **Start execution**.

```
1 {
2   "input_bucket": "{{INPUT_BUCKET}}",
3   "prefix": ""
4 }
```

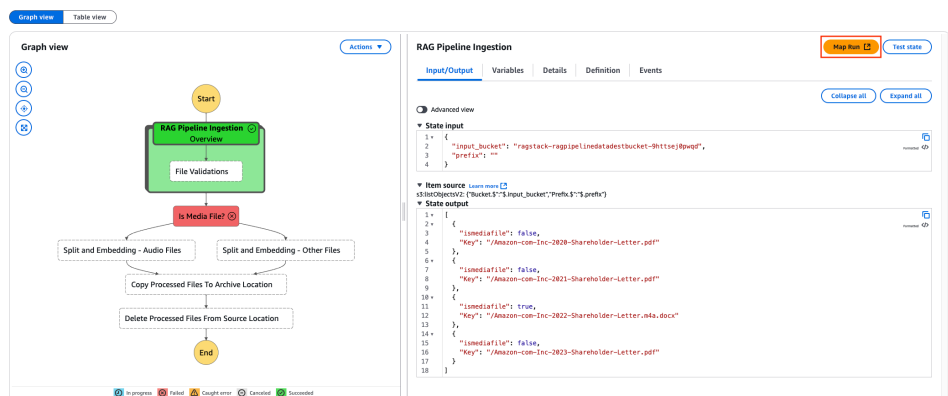


4. Select the map state in the graphical view.

Do not be alarmed by the deliberately introduced "States.ExceedToleratedFailureThreshold" error.

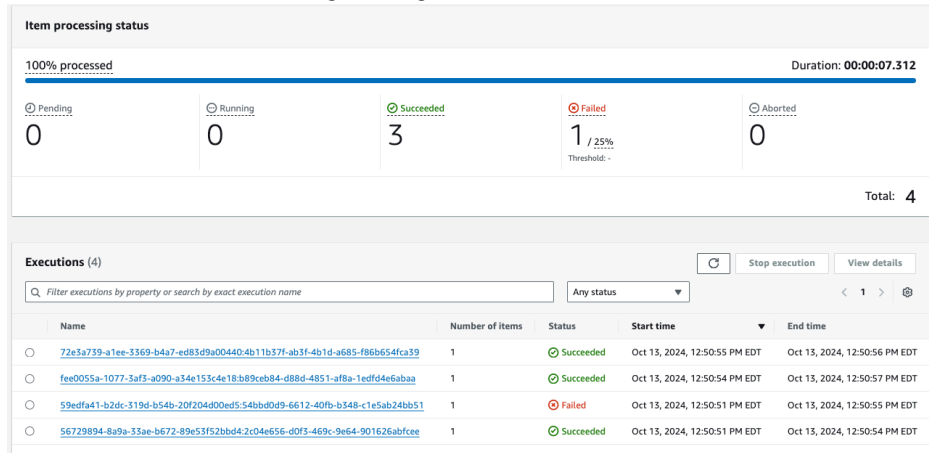


5. Select **Map Run** to view details of the Distributed Map execution.



5. Take sometime to review the details of the map run and the configurations.

6. Notice one of the iterations failing. You might see the results similar to below screenshot



Fixing the issue that caused the failure

The reason for failure is we deliberately introduced error in `rag-pipeline-ingestion-file-validation` lambda function. You are going to fix it and redrive the workflow.

Why did we introduce the error?

When processing large number of files, processing might fail due to data quality or exceptions in business logic. In such situations, restarting the processing of all the files can incur significant time loss and money. For example, suppose if processing failed at 1002 file when processing 1500 files, you do not want to reprocess all the 1001 files which were successfully processed. In such situations, you can use Step Functions redrive feature to restart from the point of failure.

1. Access the [Lambda console](#) and edit `rag-pipeline-ingestion-file-validation` function.
2. Search for `####`. This change is done to introduce error scenario.
3. Replace the entire return statement (from **line#15** to **line#21**) that included the searched string with the following code.

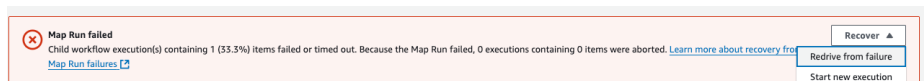
```
1      return {
2          'ismediafile': True,
3          'Key': '/' + filename
4      }
```



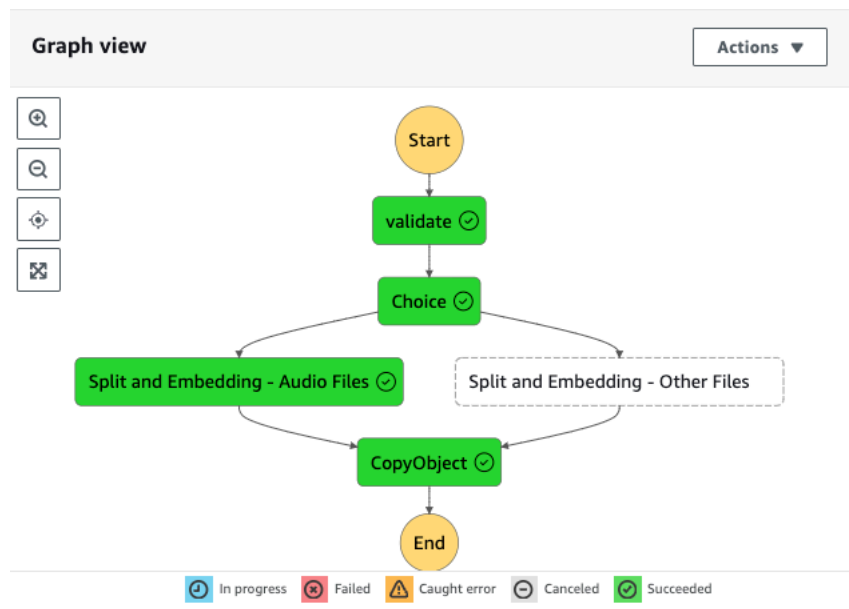
4. Make sure the indentations are correct.
5. Deploy the Lambda function

Restarting from the point of failure (Redrive)

1. Return to the Step Functions console.
2. Access the failed maprun
3. Restart the workflow by selecting the `recover` option from the failure banner. Select `Redrive from failure` from failure.



4. In the following popup, select the button `Redrive` execution.
5. Go to **Map Run**. Notice that just the one iteration that failed got restarted. This time the file is processed by the correct audio file processing step



Verifying the processed files

1. Navigate to [S3 console](#). Search for `ragpipelinearchivebucket` bucket.
2. Check the Archive folder. Once all the files are processed successfully those will be moved to Archive folder.

Amazon S3 > Buckets > ragstack-ragpipelinearchivebucket-jcwymzhinuq

ragstack-ragpipelinearchivebucket-jcwymzhinuq info

Objects Properties Permissions Metrics Management Access Points

Objects (4) info Find objects by prefix Copy S3 URI Copy URL Download Open Delete Actions ▾ Create folder Upload

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	Amazon-com-Inc-2020-Shareholder-Letter.pdf	pdf	October 9, 2024, 04:51:22 (UTC-05:00)	469.2 KB	Standard
<input type="checkbox"/>	Amazon-com-Inc-2021-Shareholder-Letter.pdf	pdf	October 9, 2024, 04:51:22 (UTC-05:00)	102.8 KB	Standard
<input type="checkbox"/>	Amazon-com-Inc-2022-Shareholder-Letter.m4a	m4a	October 9, 2024, 04:54:00 (UTC-05:00)	1.3 MB	Standard
<input type="checkbox"/>	Amazon-com-Inc-2023-Shareholder-Letter.pdf	pdf	October 9, 2024, 14:39:15 (UTC-05:00)	98.8 KB	Standard

Verifying Vector store

1. Navigate to the [Amazon OpenSearch Service Console](#)
2. Select **Collections** from left menu
3. Select **rag-pipeline-collection**
4. Navigate to indexes tab and review the indexes **rag-pipeline-knowledge-base**. You will see the number of indexed documents.

Amazon OpenSearch Service > Serverless Collections > rag-pipeline-collection

Start indexing your vector data. A vector index consists of vector embeddings alongside metadata that describes the data. Create vector index

rag-pipeline-collection info Delete collection OpenSearch Dashboards

Overview Monitor **Indexes** Tags

Indexes (1) Refresh Delete Create vector index

Indexing is the method by which search engines organize data for fast retrieval. [Learn more](#)

Index name	Total size (bytes)	Total document count	Total vector field count	Created date
rag-pipeline-knowledge-base	36.3mb	983	1	2024-10-10 12:51 UTC

✔ Congratulations! You have successfully built and tested the RAG pipeline. The data is ready in vector format in the vector database. In the next section, you will test if the chatbot is able to use the vector data to fetch you the right results

[Previous](#)[Next](#)