

CSCI 5409: Assignment 3: Part B:

Waleed R. Alhindi (B00919848)

GitLab A3 Part B Repository:

The A3-Part B code has been pushed to the GitLab repository's "A3" branch under the "Part_B" folder, which can be found at the following URL:

https://git.cs.dal.ca/alhindi/csci5410-summer-23-b00919848/-/tree/A3/Part_B

OR

Under the "A3" folder's "Part_B" sub-folder in the "main" branch, which can be found here:

https://git.cs.dal.ca/alhindi/csci5410-summer-23-b00919848/-/tree/main/A3/Part_B

Additionally, the professor and all TAs have been granted "Maintainer" access to this GitLab repository. This includes assigning "Maintainer" roles to the following GitLab accounts:

- @saurabh (Dr. Saurabh Dey)
- @mudgal (Ankush Mudgal)
- @bharatwaaj (Bharatwaaj Shankanarayanan)
- @rmacwan (Rahul Ashokkumar Macwan)

Operations Performed:

1. Creating S3 Buckets:

Our first step is to create the two AWS S3 [1] buckets: a "sampledata-b00919848" bucket, where the files from the "Tech" folder will be stored, and "tags-b00919848", where the extracted feature JSON of each file will be stored.

Note that although the S3 Buckets will be created manually, the files from the Tech folder will be uploaded to the "sampledata-b00919848" bucket using a script.

Thus, we can go ahead and create the "sampledata-b00919848" bucket as seen in Figure 1:

Amazon S3 > Buckets > Create bucket

Create bucket Info

Buckets are containers for data stored in S3. [Learn more](#)

General configuration

Bucket name

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

AWS Region

US East (N. Virginia) us-east-1

Fig 1. Creating “sampledata-b00919848” S3 Bucket [1]

Similarly, we create the “tags-b00919848” Bucket as seen in Figure 2:

Amazon S3 > Buckets > Create bucket

Create bucket Info

Buckets are containers for data stored in S3. [Learn more](#)

General configuration

Bucket name

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

AWS Region

US East (N. Virginia) us-east-1

Fig 2. Creating “tags-b00919848” S3 Bucket [1]

As seen in Figure 3, the two S3 Buckets have been successfully created:

Amazon S3

▼ Account snapshot Info

Last updated: Jul 10, 2023 by Storage Lens. Metrics are generated every 24 hours. [Learn more](#)

Total storage

860.0 B

Object count

21

Average object size

41.0 B

You can enable advanced metrics in the "default-account-dashboards" configuration.

Buckets (2) Info

Buckets are containers for data stored in S3. [Learn more](#)

↺ 1 ↻ ⚙

Name	AWS Region	Access	Creation date
<input type="radio"/> sampledata-b00919848	US East (N. Virginia) us-east-1	Bucket and objects not public	July 11, 2023, 10:47:37 (UTC-03:00)
<input type="radio"/> tags-b00919848	US East (N. Virginia) us-east-1	Bucket and objects not public	July 11, 2023, 10:05:41 (UTC-03:00)

Fig 3. S3 Buckets Successfully Created [1]

2. Creating DynamoDB Table:

In addition to the two S3 Buckets [1] created, we need to create a DynamoDB [2] table to store each extracted feature (named entity) and its number of occurrences. Thus, we can go ahead and create the “namedEntities” DyanmoDB [2] table as seen in Figure 4:

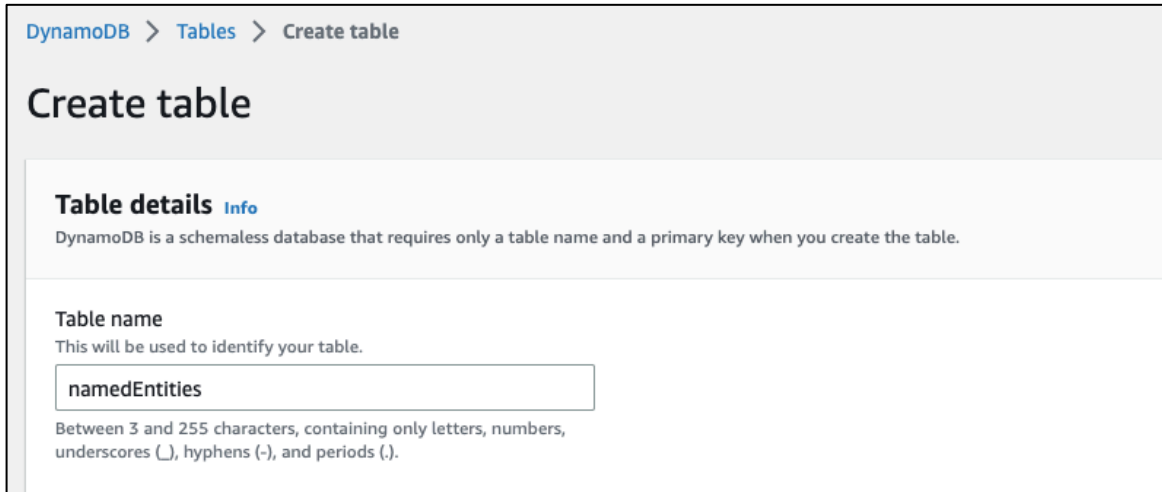


Fig 4. Creating “namedEntities” DynamoDB Table [2]

However, as seen in Figure 5, the “namedEntities” table is created with a partition key called “entity”, which refers to the named entity being stored. Therefore, each document will contain two fields: an “entity” partition key field that denotes the actual word being stored and a “tally” value that denotes how many times it has occurred. For example, the JSON {“USA”: 2, “Canada”: 1} will be stored as two documents: {“entity”: “USA”, “tally”: 2} and {“entity”: “Canada”, “tally”: 1}. Additionally, these documents will be updated as newer (applicable) JSONs are stored in the “tags-b00919848” S3 Bucket [1].

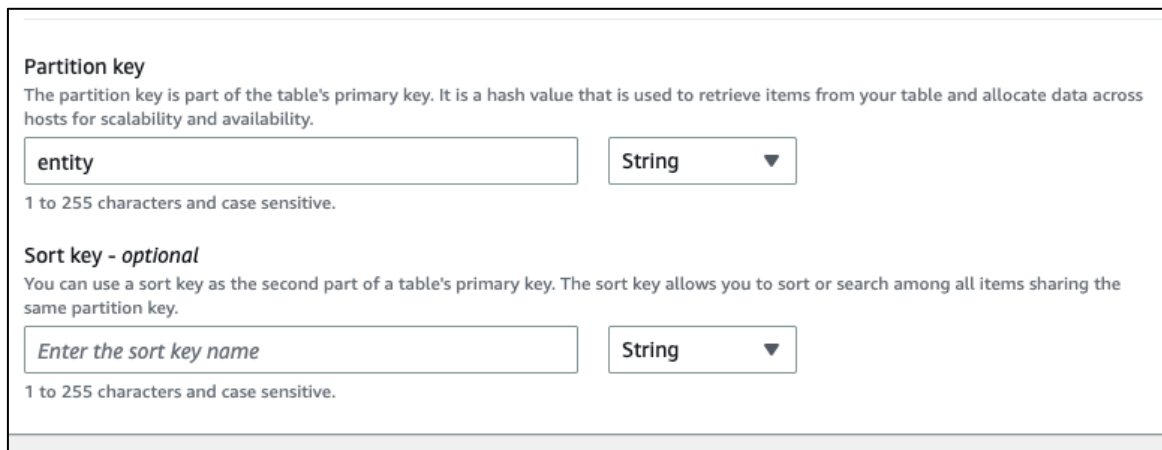


Fig 5. Defining “entity” Parititon Key in “namedEntities” Table [2]

Finally, as seen in Figure 6, the “namedEntities” DyanmoDB table [2] has been created successfully:

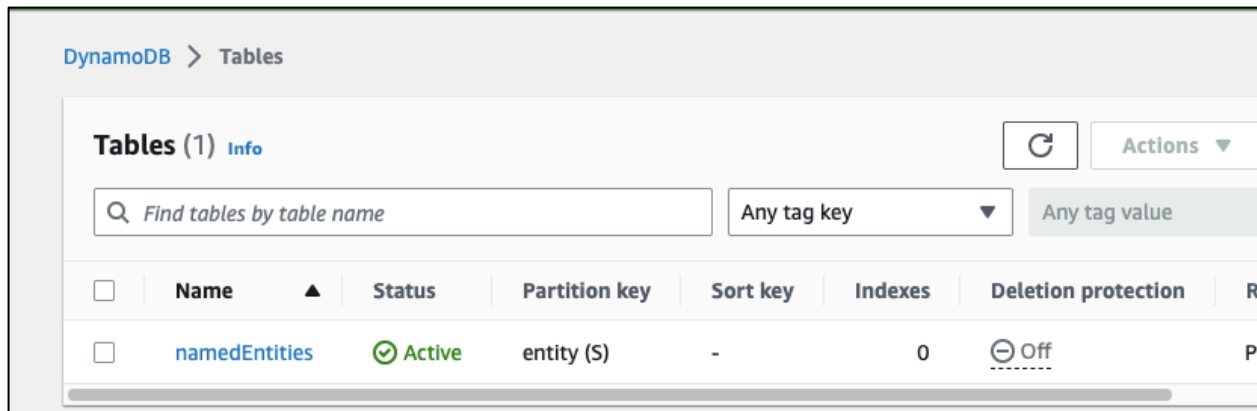


Fig 6. “namedEntities” DyanmoDB Table Successfully Created [2]

3. *Composing extractFeatures Lambda:*

Now that our S3 Buckets [1] and DynamoDB table [2] are set up, we can begin implementing the “extractFeatures” AWS Lambda [3] function. This function is meant to be triggered any time an object is added to our “sampledata-b00919848” S3 Bucket [1], and will extract all named entities in the newly added object (file from Tech folder) into a JSON string that it will then store in our “tags-b00919848” S3 Bucket [1].

Let us first create the “extractFeatures” function as seen below:

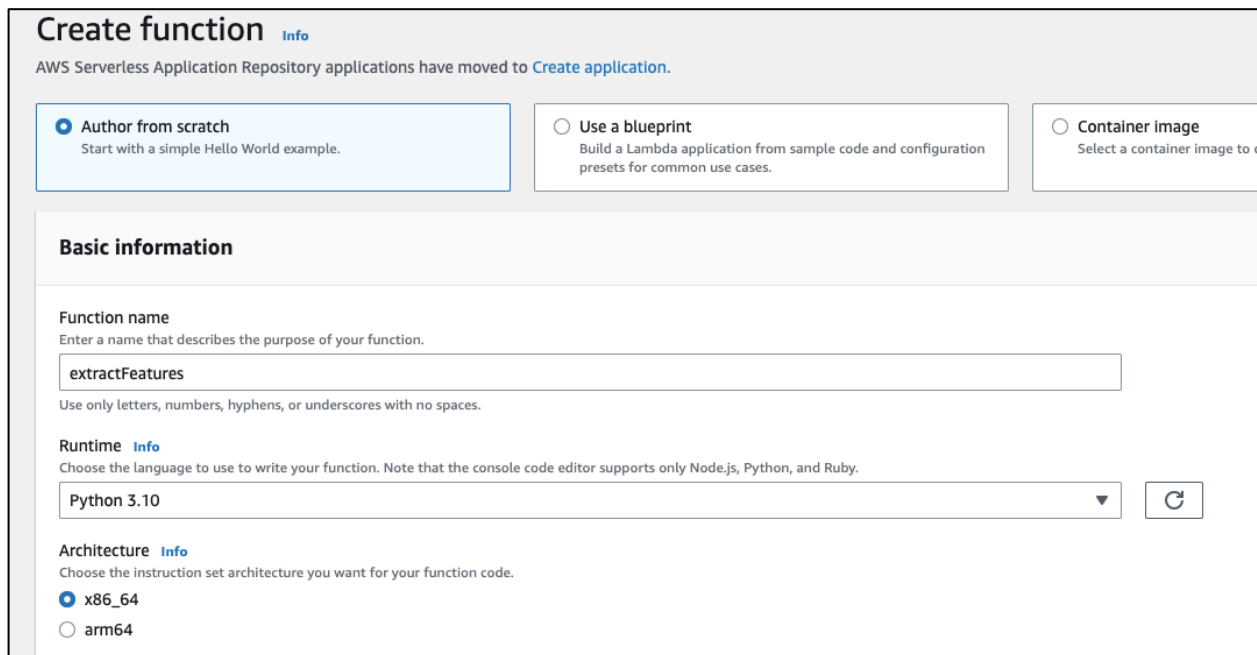



Fig 7. Creating “extractFeatures” Lambda Function [3]

After creating the function, we need to configure a trigger on to the “extractFeatures” that will be invoked whenever an object is pushed to our “sampledata-b00919848” S3 Bucket [1]:

Lambda > Add trigger

Add trigger

Trigger configuration [Info](#)

 **S3**
aws storage

Bucket
Please select the S3 bucket that serves as the event source. The bucket must be in the same region as the function.


Bucket must be in region us-east-1


Event types
Select the events that you want to have trigger the Lambda function. You can optionally set up a prefix or suffix for an event. However, for each bucket, individual events cannot have multiple configurations with overlapping prefixes or suffixes that could match the same object key.

Fig 8. Adding Trigger to “extractFeatures” Lambda [3]

extractFeatures

▼ **Function overview** [Info](#)

 **extractFeatures**

 **Layers** (0)


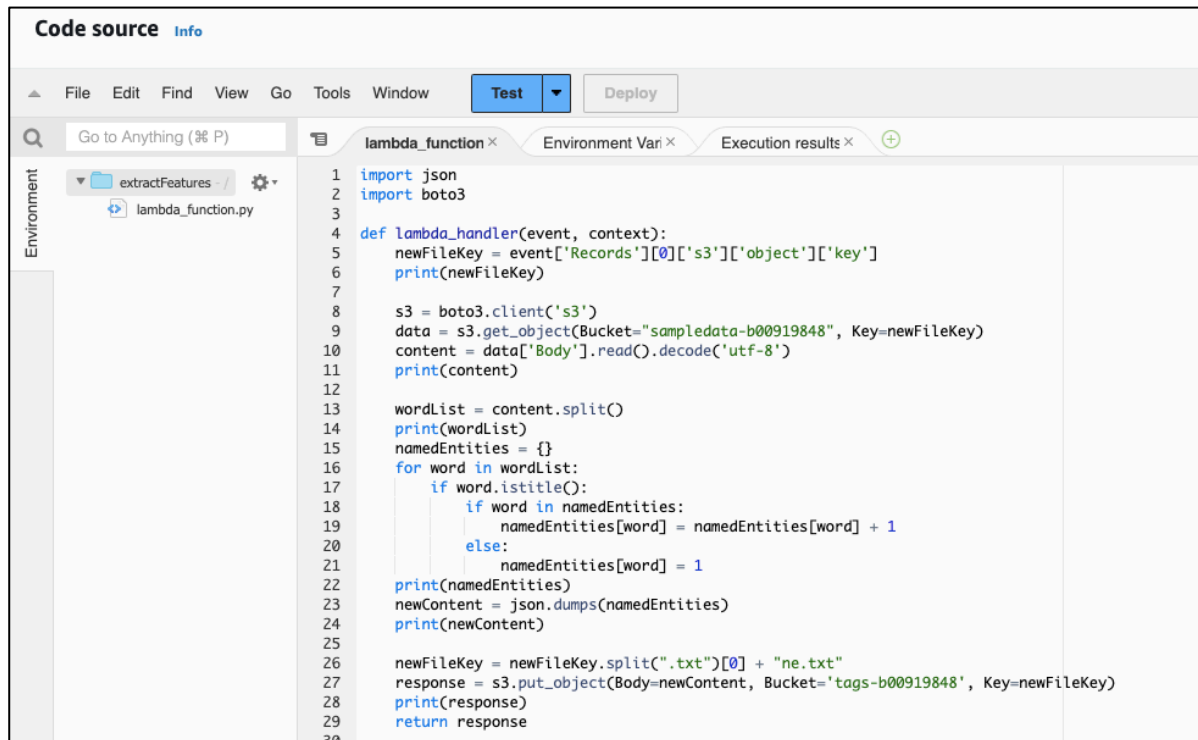
 **S3**

Fig 9. Trigger Successfully Added to “extractFeatures” Function [3]

Next, the function’s source code was implemented as seen in Figure 10:



```
1 import json
2 import boto3
3
4 def lambda_handler(event, context):
5     newFileKey = event['Records'][0]['s3']['object']['key']
6     print(newFileKey)
7
8     s3 = boto3.client('s3')
9     data = s3.get_object(Bucket="sampledata-b00919848", Key=newFileKey)
10    content = data['Body'].read().decode('utf-8')
11    print(content)
12
13    wordList = content.split()
14    print(wordList)
15    namedEntities = {}
16    for word in wordList:
17        if word.istitle():
18            if word in namedEntities:
19                namedEntities[word] = namedEntities[word] + 1
20            else:
21                namedEntities[word] = 1
22    print(namedEntities)
23    newContent = json.dumps(namedEntities)
24    print(newContent)
25
26    newFileKey = newFileKey.split(".txt")[0] + "ne.txt"
27    response = s3.put_object(Body=newContent, Bucket='tags-b00919848', Key=newFileKey)
28    print(response)
29    return response
30
```

Fig 10. “extractFeature” Function’s Source Code [3]

The code depicted in Figure 10 is triggered whenever an object is pushed to the “sampledata-b00919848” bucket; supplying the function with an “event” that contains information about the object that was just pushed. It then retrieves the newly pushed object, scans through its data, and added each named entity (i.e., each word that begins with a capital letter) into a JSON. This JSON is then converted into an encodable string and stored in our “tags-b00919848” bucket.

Additionally, the pseudocode of the code depicted in Figure 10 is outlined below:

Get new object’s key from “event”

Create empty JSON

Retrieve object using key

Split object’s string data into an array of words

Iterate through each word in that array

If that word’s first letter is capital

If that word exists in our JSON

Increment its value by 1

Else

Add {word: 1} to the JSON

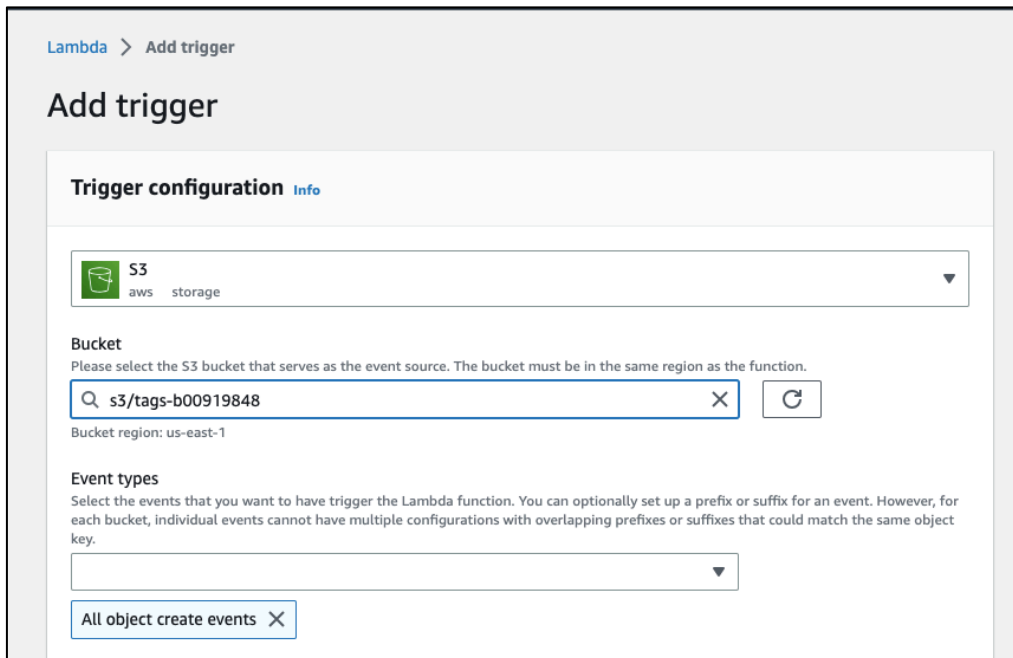
Convert JSON into string (dumps)

Upload JSON string to “tags-b00919848” bucket

4. Composing accessDB Lambda:

The “accessDB” Lambda [3] function is meant to be triggered once the “extractFeatures” Lambda [3] has uploaded an object to the “tags-b00919848” bucket. Once triggered, it will retrieve the named entities from the newly uploaded file and update the “namedEntities” DynamoDB [2] table accordingly. In other words, when a JSON string file is uploaded to the “tags-b00919848” bucket, this function will retrieve that file and use the named entities stored in it to update the “namedEntities” table accordingly.

Implementing the “accessDB” Lambda [3] function is very similar to our “extractFeatures” function above. We first create the “accessDB” function, then create a trigger for it. However, as seen in Figure 11, this time our trigger will be for the “tags-b00919848” bucket:



The screenshot shows the 'Add trigger' configuration page in the AWS Lambda console. The 'Trigger configuration' section is active, showing 'S3' as the provider. The bucket 's3/tags-b00919848' is selected, and the region is 'us-east-1'. Under 'Event types', 'All object create events' is selected.

Fig 11. Adding Trigger to “accessDB” Function [3]



Fig 12. Trigger Successfully Added to “accessDB” Function [3]

After which, we can implement the function’s source code, which can be seen in Figure 13 below:

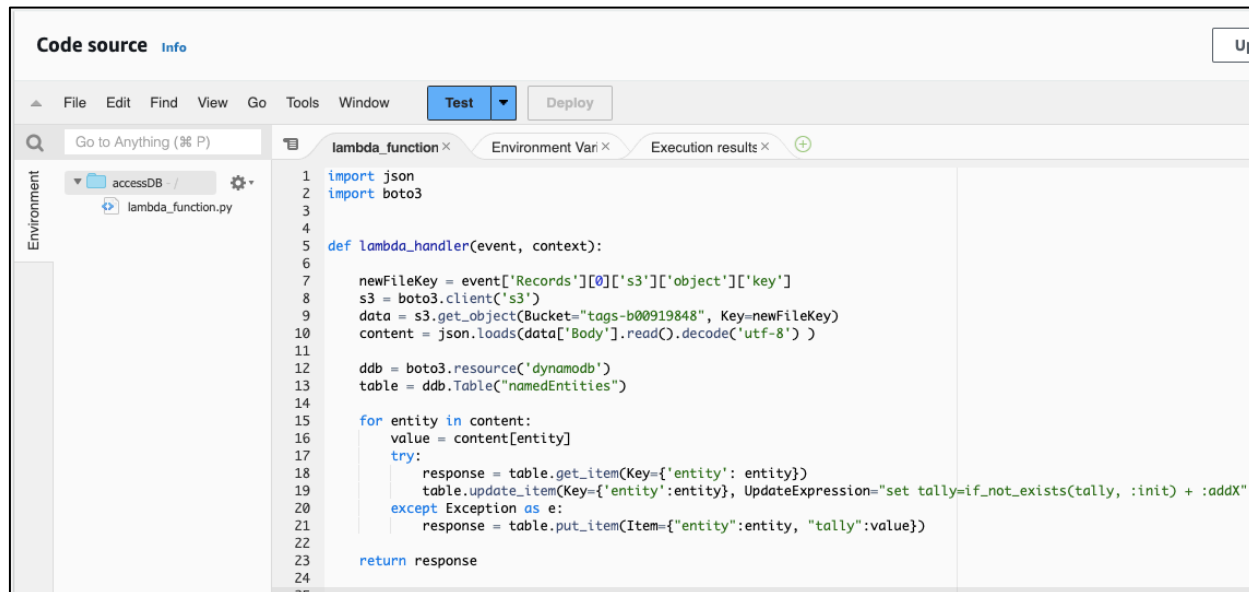


Fig 13. “accessDB” Function’s Source Code [3]

The code depicted above is invoked whenever a new object is added to the target S3 bucket and is passed information about the newly added object. Using that information, the object is retrieved. Then, the JSON string is converted back into a JSON object (loads). Then, for each key in the JSON, it either adds that named entity to our DynamoDB [3] table if it does not already exist or increments that entity’s tally in the table if it does exist.

The pseudocode of the code depict in Figure 13 is outlined below:

Get new object’s key from “event”

Retrieve that object using that key

Load object’s data into a JSON

Iterate through each key (named entity) in the JSON

If the key already exists in the DyanmoDB table

Increment its corresponding tally based on the key’s value

Else

Create new item in the table to represent this key and set its tally to the key’s value

5. Implementing Script to Upload Files from Tech Folder to S3:

The last piece of our event-driven puzzle is the script that will upload each of the 401 files in the provided Tech folder to the “sampledata-b00919848” bucket with a 100 millisecond delay between each upload.

Figure 14 showcases the simple Python script used to upload each file in the Tech folder:

```
uploadFiles.py
1  import os
2  import json
3  import boto3
4  import time
5
6  s3 = boto3.resource('s3')
7
8  for filename in os.listdir('./tech'):
9      if filename == '.DS_store':
10         continue
11     print(filename)
12     response = s3.meta.client.upload_file(os.path.join(os.getcwd(), "tech/"+filename), 'sampledata-b00919848', filename)
13     time.sleep(0.1)
14
15
```

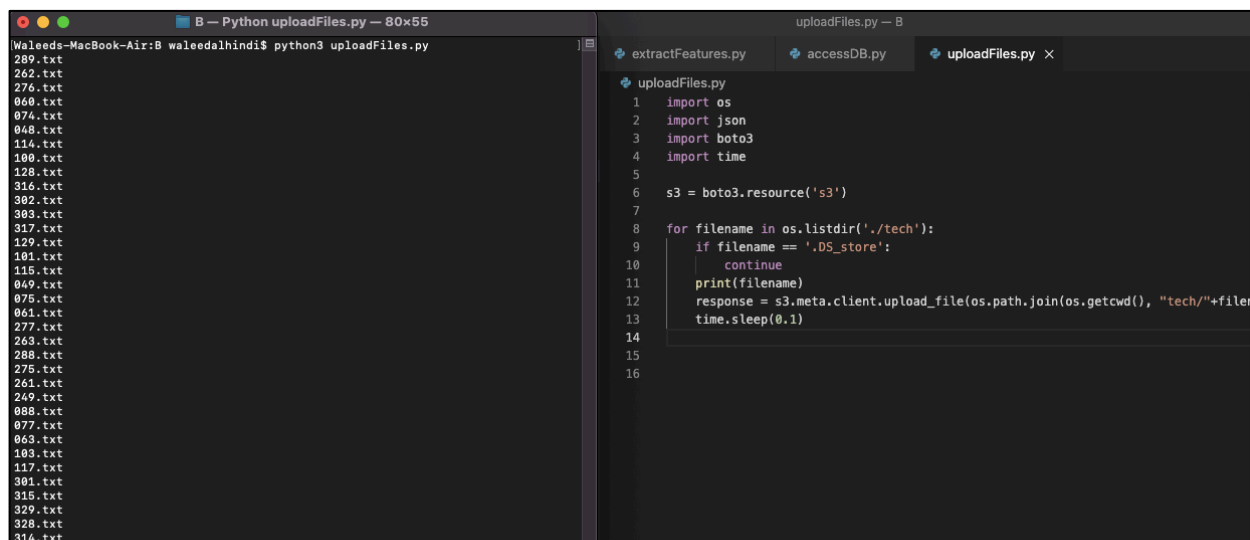
Fig 14. Python Script to Upload Tech Folder's Files to “sampledata-b00919848” Bucket

The code depicted above is a simple Python script that uses my locally stored AWS credentials (i.e., the ones stored in ~/.aws/credentials) and the boto3 library to upload each file in the provided Tech folder to our “sampledata-b00919848” bucket. Note, however, that after each upload, a 100 millisecond delay is awaited.

6. Testing and Verification:

Now that all we have configured our S3 Buckets, DynamoDB table, Lambda functions, and file upload script, all that is left to do is to verify and test its behavior. Note that running the file upload script should initiate the entire trigger pipeline. In other words, we shall execute the file upload script, then verify whether files are being correctly extracted and stored into the S3 buckets and DynamoDB table.

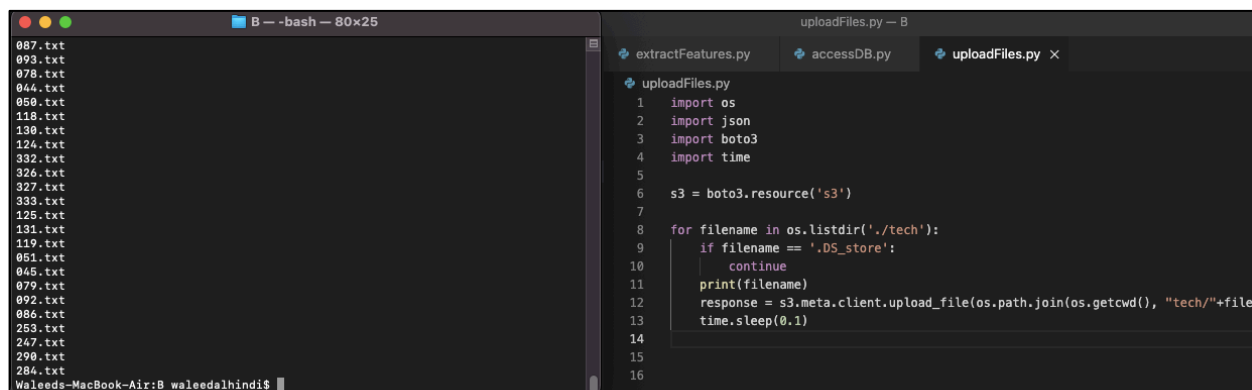
As such, we go ahead and *locally* execute the script as seen in Figure 15:



The screenshot shows a terminal window on the left and a code editor on the right. The terminal window, titled 'B — Python uploadFiles.py — 80x55', shows the command 'python3 uploadFiles.py' being executed, followed by a list of files in the 'tech' folder: 289.txt, 262.txt, 276.txt, 060.txt, 074.txt, 048.txt, 114.txt, 100.txt, 128.txt, 316.txt, 302.txt, 303.txt, 317.txt, 129.txt, 101.txt, 115.txt, 049.txt, 075.txt, 061.txt, 277.txt, 263.txt, 288.txt, 275.txt, 261.txt, 249.txt, 083.txt, 077.txt, 063.txt, 103.txt, 117.txt, 301.txt, 315.txt, 329.txt, 328.txt, and 314.txt. The code editor, titled 'uploadFiles.py — B', shows the same Python script as in Figure 14, with the file 'uploadFiles.py' selected in the tab bar.

Fig 15. Begin Executing UploadFiles.py Script

The script takes some time to execute due to the number of files, but once complete we can see that it has finished executing without any exceptions being raised; meaning the files should have been successfully uploaded.



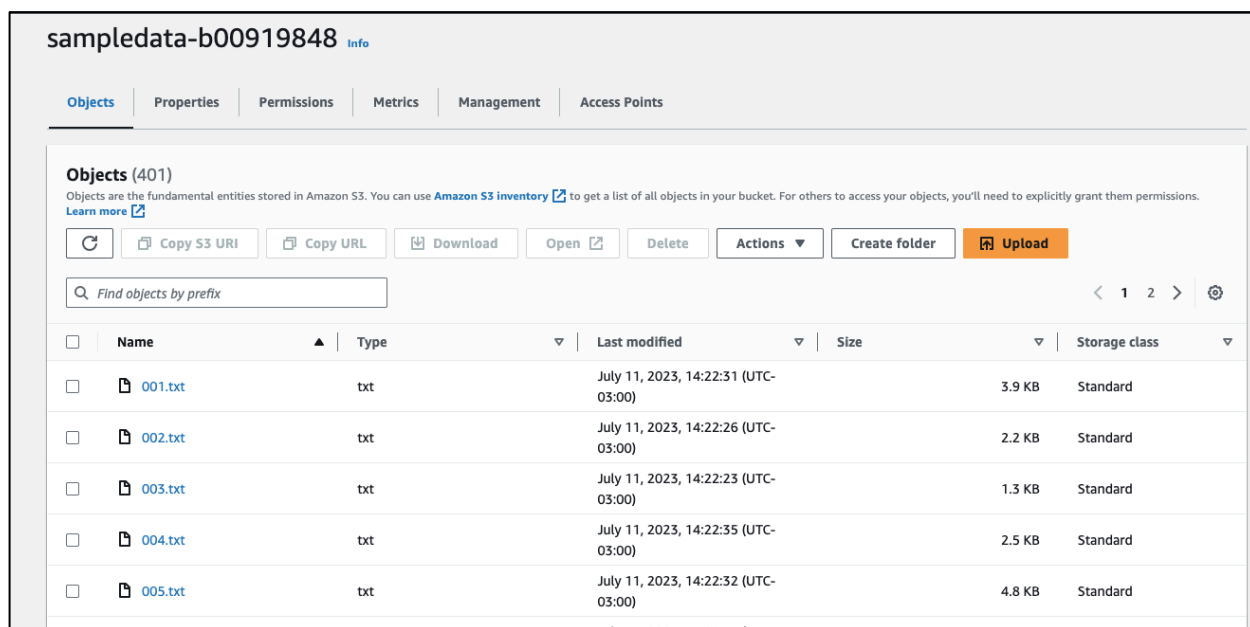
The image shows a terminal window on the left and a code editor on the right. The terminal window displays a list of files (087.txt to 284.txt) and a prompt for the user 'waleedalhindi'. The code editor shows the contents of the 'uploadFiles.py' script, which imports 'os', 'json', 'boto3', and 'time', initializes an S3 resource, and loops through files in the './tech' directory to upload them to an S3 bucket named 'sampledata-b00919848'.

```
087.txt
093.txt
078.txt
044.txt
050.txt
118.txt
130.txt
124.txt
332.txt
326.txt
327.txt
333.txt
125.txt
131.txt
119.txt
051.txt
045.txt
079.txt
092.txt
086.txt
253.txt
247.txt
290.txt
284.txt
Waleeds-MacBook-Air:~ waleedalhindi$

uploadFiles.py
1 import os
2 import json
3 import boto3
4 import time
5
6 s3 = boto3.resource('s3')
7
8 for filename in os.listdir('./tech'):
9     if filename == '.DS_store':
10         continue
11     print(filename)
12     response = s3.meta.client.upload_file(os.path.join(os.getcwd(), "tech/"+filename), "sampledata-b00919848", filename)
13     time.sleep(0.1)
14
15
16
```

Fig 16. UploadFiles.py Script’s Execution Completed

While Figures 15 and 16 suggest that all files have been uploaded to the bucket, we still need to verify that the files were indeed uploaded by navigating to our “sampledata-b00919848” bucket. As expected, each of the 401 files were uploaded to the bucket as seen in Figure 17:



The image shows the Amazon S3 console interface for the bucket 'sampledata-b00919848'. The 'Objects' tab is selected, showing a list of 401 objects. The table displays columns for Name, Type, Last modified, Size, and Storage class. The first five objects are listed as 001.txt, 002.txt, 003.txt, 004.txt, and 005.txt, all with a 'Standard' storage class.

Name	Type	Last modified	Size	Storage class
001.txt	txt	July 11, 2023, 14:22:31 (UTC-03:00)	3.9 KB	Standard
002.txt	txt	July 11, 2023, 14:22:26 (UTC-03:00)	2.2 KB	Standard
003.txt	txt	July 11, 2023, 14:22:23 (UTC-03:00)	1.3 KB	Standard
004.txt	txt	July 11, 2023, 14:22:35 (UTC-03:00)	2.5 KB	Standard
005.txt	txt	July 11, 2023, 14:22:32 (UTC-03:00)	4.8 KB	Standard

Fig 17. Tech Files Successfully Uploaded to “sampledata-b00919848” Bucket [1]

Furthermore, let us verify that the contents of the files are correct as well by accessing one of the files then downloading and verifying its contents as seen in Figures 18 and 19:

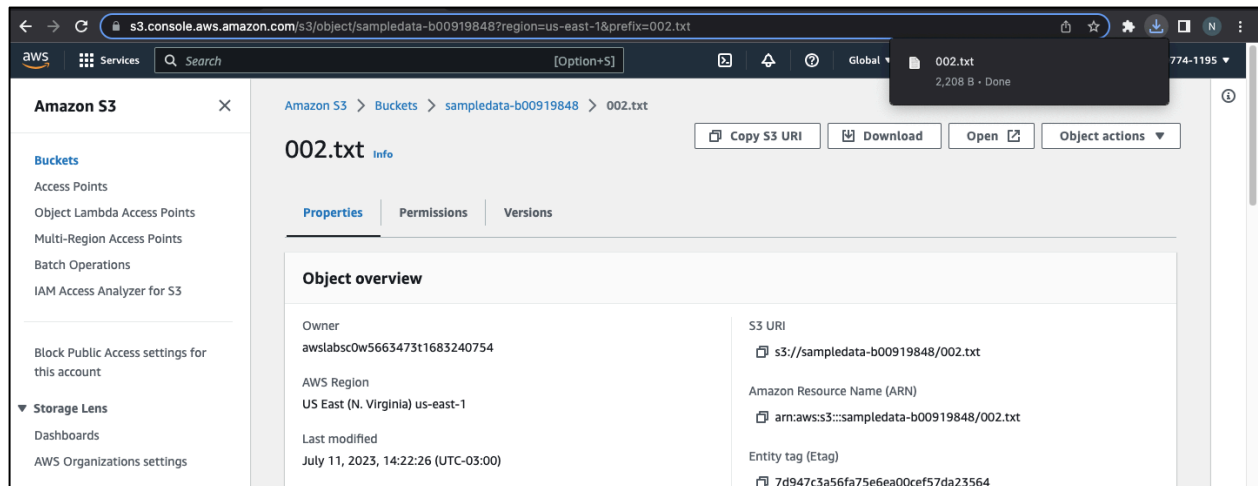


Fig 18. Downloading a File to Verify its Contents [1]

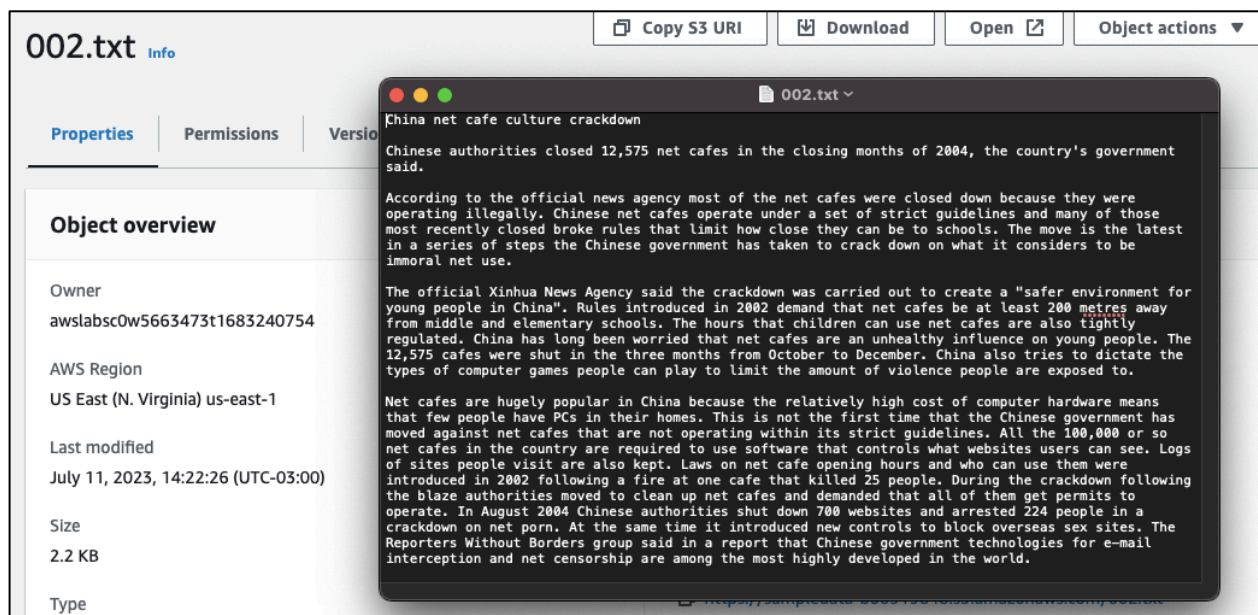


Fig 19. Verifying Downloaded File's Contents

Now that we have verified that the UploadFiles.py script has indeed uploaded all the Tech folder's files to the correct S3 bucket, we need to confirm that "extractFeature" Lambda was triggered on each upload to the "sampledata-b00919848" bucket. To verify that it has been triggered successfully, we navigate to the function's logs in Cloud Watch [4], as seen in Figure 20, where we can see a bunch of logs have been created:

Log streams

Metric filters

Subscription filters

Contributor Insights

Tags

Data protection

Log streams (62)

Delete

Create log stream

Search all log streams

Filter log streams or try prefix search

Exact match

Show expired

Info

1

<div><div></div></div> Log stream	Last event time
<div><div></div><div>2023/07/11/[\$LATEST]e6f6e19d7bff403ab64210c1604eb4e0</div></div>	2023-07-11 14:11:15 (UTC-03:00)
<div><div></div><div>2023/07/11/[\$LATEST]855947eea2a14975a704f15d5f55c2d9</div></div>	2023-07-11 14:11:04 (UTC-03:00)
<div><div></div><div>2023/07/11/[\$LATEST]049070a2e1fd45289208151994514d44</div></div>	2023-07-11 14:02:07 (UTC-03:00)
<div><div></div><div>2023/07/11/[\$LATEST]a1ccd183933b4d03a09482c4da3af9ce</div></div>	2023-07-11 14:02:07 (UTC-03:00)
<div><div></div><div>2023/07/11/[\$LATEST]ab2badc5bf0848509bf43a1f9e8b5f95</div></div>	2023-07-11 14:02:06 (UTC-03:00)
<div><div></div><div>2023/07/11/[\$LATEST]955f2bcd6c34380958c01baa5993727</div></div>	2023-07-11 14:02:05 (UTC-03:00)
<div><div></div><div>2023/07/11/[\$LATEST]8481cc764ed64774a1b9b8d2b91a6cf9</div></div>	2023-07-11 14:02:02 (UTC-03:00)
<div><div></div><div>2023/07/11/[\$LATEST]384a6acba4e043a8dbf35f46dea8933</div></div>	2023-07-11 14:02:02 (UTC-03:00)

Fig 20. “extractFeatures” Lambda’s Execution Logs in Cloud Watch [4]

Now, let us navigate inside one of these log streams to verify that “extractFeatures” is executing properly:

▶	2023-07-11T14:12:07.754-03:00	{ "Macrovision.": 1, "The": 4, "Macrovision": 4, "Some": 1, "News": 1, "Linux.": 1, "Hollywood": 1, "...
▶	2023-07-11T14:12:07.992-03:00	{ "ResponseMetadata": { "RequestId": "J4YESZJM8JDKGFG8", "HostId": "K5Ke/8LXkPhWy08KAoSWdddkbiWfNmHyAN1...
▶	2023-07-11T14:12:07.994-03:00	END RequestId: d93e587c-dba5-4d73-83ff-96d01d11f691
▶	2023-07-11T14:12:07.994-03:00	REPORT RequestId: d93e587c-dba5-4d73-83ff-96d01d11f691 Duration: 523.34 ms Billed Duration: 524 ms Me...
▶	2023-07-11T14:12:08.956-03:00	START RequestId: 1a73ab35-1c4f-4662-8462-5f279d93900d Version: \$LATEST
▶	2023-07-11T14:12:08.957-03:00	131.txt
▶	2023-07-11T14:12:09.251-03:00	Seamen sail into biometric future
▶	2023-07-11T14:12:09.251-03:00	The luxury cruise liner Crystal Harmony, currently in the Gulf of Mexico, is the unlikely setting for...
▶	2023-07-11T14:12:09.251-03:00	As holidaymakers enjoy balmy breezes, their ship's crew is testing prototype versions of the world's ...
▶	2023-07-11T14:12:09.251-03:00	Authenti-corp, the US technology consultancy, has been working with the ILO on its technical specific...
▶	2023-07-11T14:12:09.251-03:00	"If you're issued a seafarer's ID in your country, you want to be sure that when the ship lands in a ...
▶	2023-07-11T14:12:09.251-03:00	['Seamen', 'sail', 'into', 'biometric', 'future', 'The', 'luxury', 'cruise', 'liner', 'Crystal', 'Har...
▶	2023-07-11T14:12:09.251-03:00	{ "Seamen": 1, "The": 3, "Crystal": 3, "Harmony": 2, "Gulf": 1, "Mexico": 1, "As": 1, "Along": 1, "S...
▶	2023-07-11T14:12:09.251-03:00	{ "Seamen": 1, "The": 3, "Crystal": 3, "Harmony": 2, "Gulf": 1, "Mexico": 1, "As": 1, "Along": 1, "S...
▶	2023-07-11T14:12:09.489-03:00	{ "ResponseMetadata": { "RequestId": "QKG6FKBV0RH0X5GD", "HostId": "w/9GzcyfAffaQM/r0YxoPKWGq5pUeiRQmZv...
▶	2023-07-11T14:12:09.514-03:00	END RequestId: 1a73ab35-1c4f-4662-8462-5f279d93900d
▶	2023-07-11T14:12:09.514-03:00	REPORT RequestId: 1a73ab35-1c4f-4662-8462-5f279d93900d Duration: 558.09 ms Billed Duration: 559 ms Me...

Fig 21. Sample “extractFeatures” Log Stream [4]

Figure 21 depicts a snippet of one of the log streams created by the “extractFeatures” invocation. The portion highlighted in the red box corresponds to a single Lambda invocation triggered by an upload to the “sampledata-b00919848” bucket. However, the most important part of the

highlighted section is the portion depicted in Figure 22 which logs the response the Lambda function got from uploading a file to the “tags-b00919848” bucket:



Fig 22. Expanded Response Log Line from Figure 21 [4]

Additionally, we must verify the contents of the “tags-b00919848” bucket to ensure that features were correctly extracted into a JSON string. So, we navigate to the “tags-b00919848” bucket as seen in Figure 23:

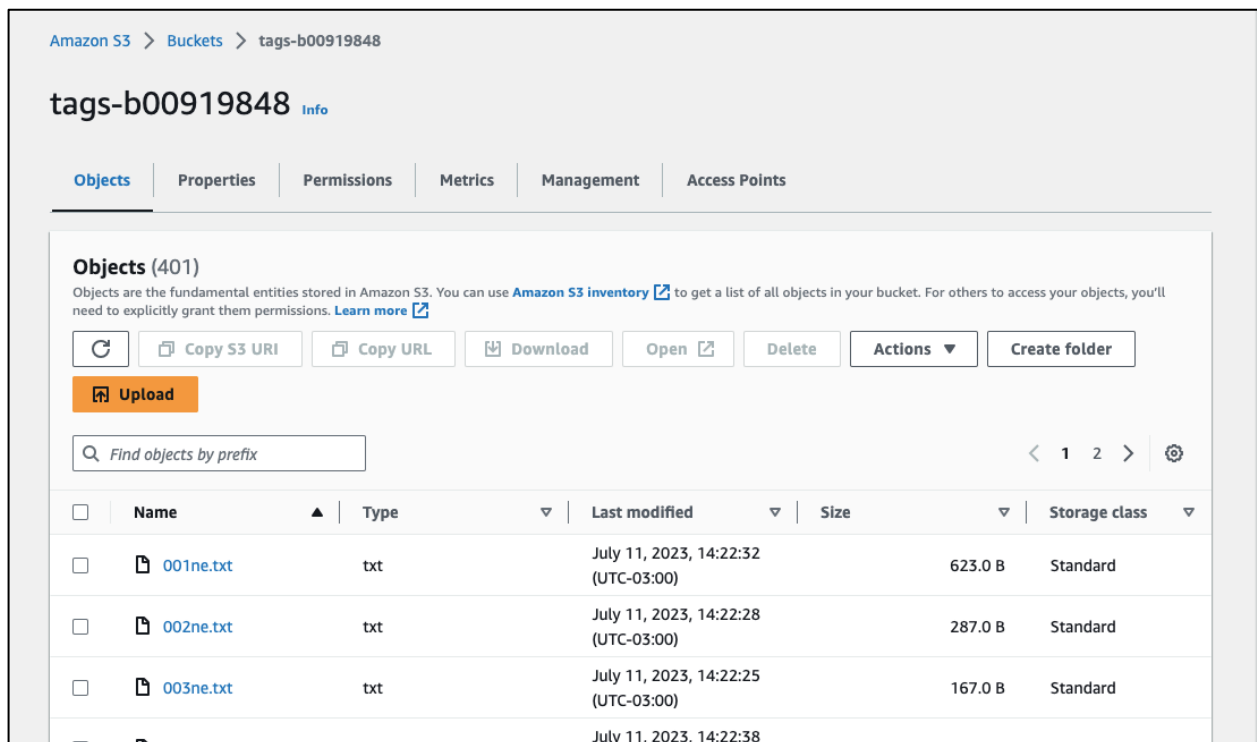


Fig 23. Extracted Named Entities Files Successfully Uploaded to “tags-b00919848” [1]

To verify the contents of the uploaded files, a sample file was downloaded and its contents verified as seen in Figure 24:

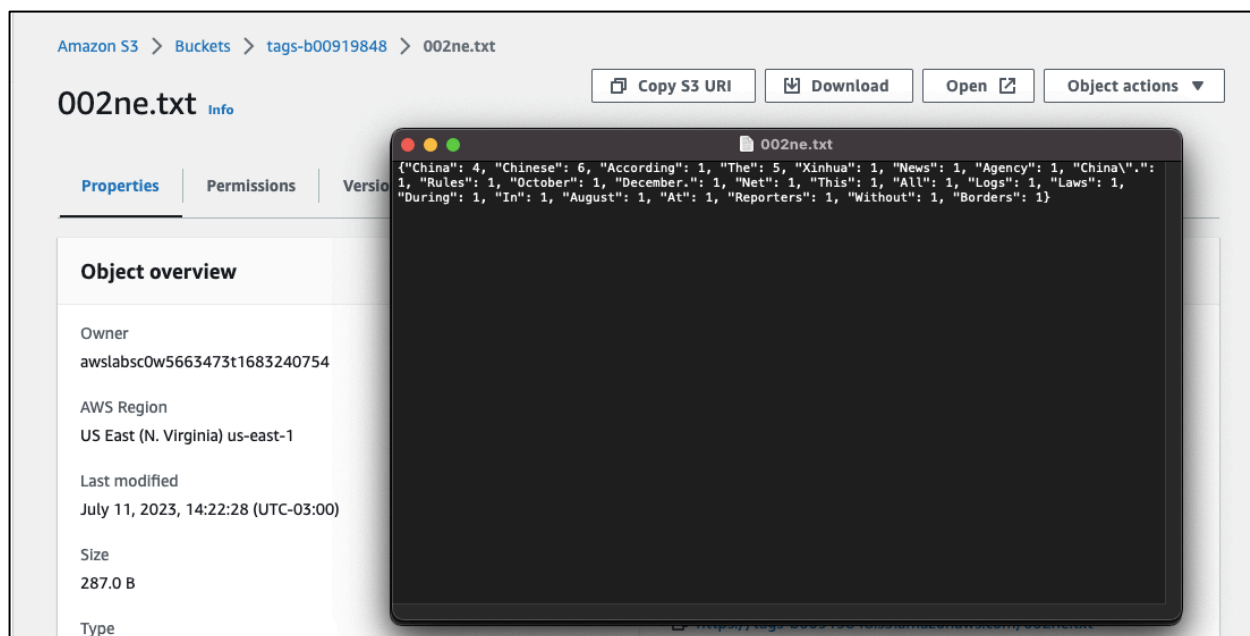


Fig 24. Verifying Sample Named Entity File in “tags-b00919848” [1]

Indeed, the named entity files have been successfully uploaded to the target “tags-b00919848” bucket via triggering the “extractFeatures” function by uploading objects to the “sampledata-b00919848” bucket.

Next, we must verify that our second Lambda function, “accessDB”, has been triggered correctly as “extractFeatures” uploads each file to the “tags-b00919848” bucket; wherein “accessDB” then updates the “namedEntities” DynamoDB table accordingly. So, we once again navigate to Cloud Watch [4]. However, this time we check the “accessDB” logs, where we can see many logs have been created:

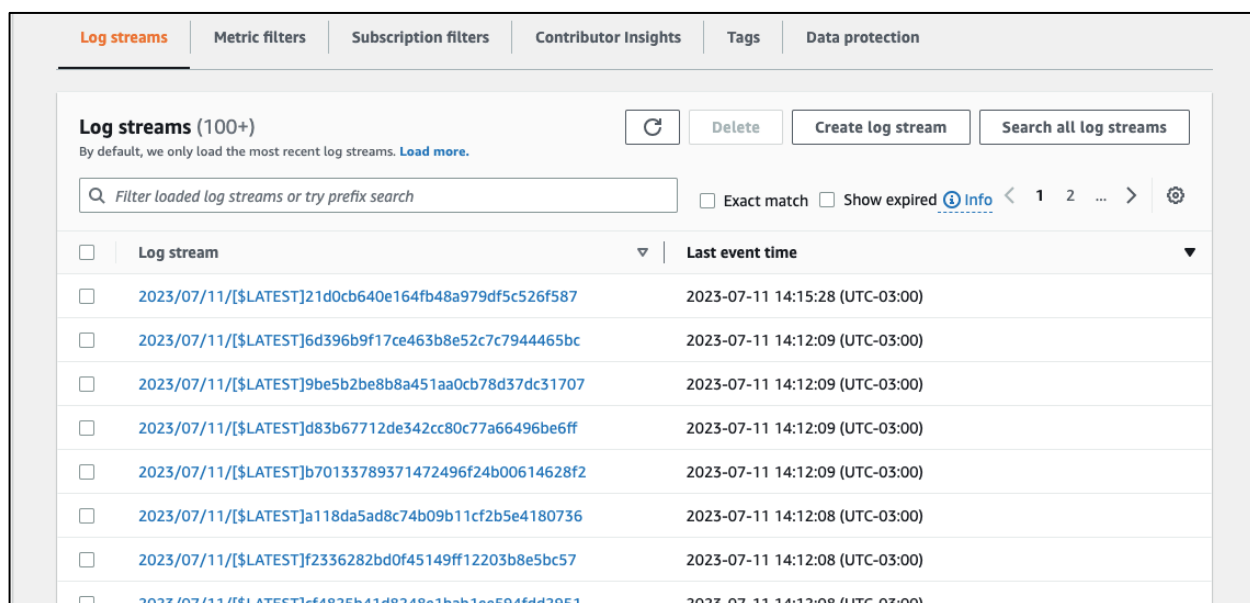


Fig 25. “accessDB” Cloud Watch Logs [4]

We then verify execution by checking a sample log stream as seen in Figure 26, where we can see that the request has been services with no exceptions being raised:

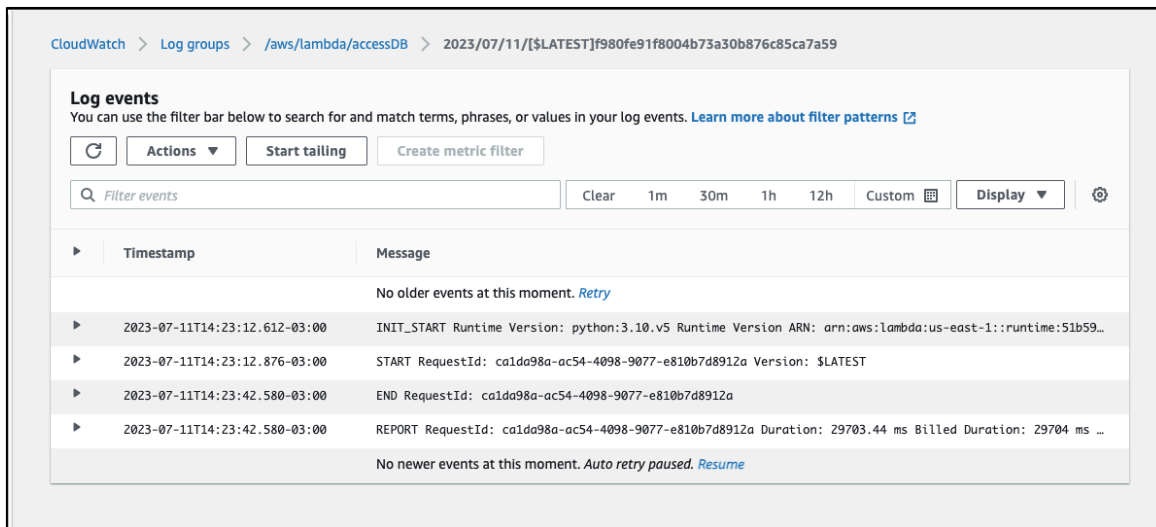


Fig 26. Sample “accessDB” Log Stream [4]

Figures 25 and 26 suggest that our “namedEntities” DynamoDB table should have been populated correctly since the many log streams correspond to each invocation triggered by “extractFeatures” uploading to “tags-b00919848”. However, to verify this fact, we navigate to our “namedEntities” DynamoDB table, depicted in Figure 27, where we can see that 4,671 named entities and their tallies have been successfully stored in the table:

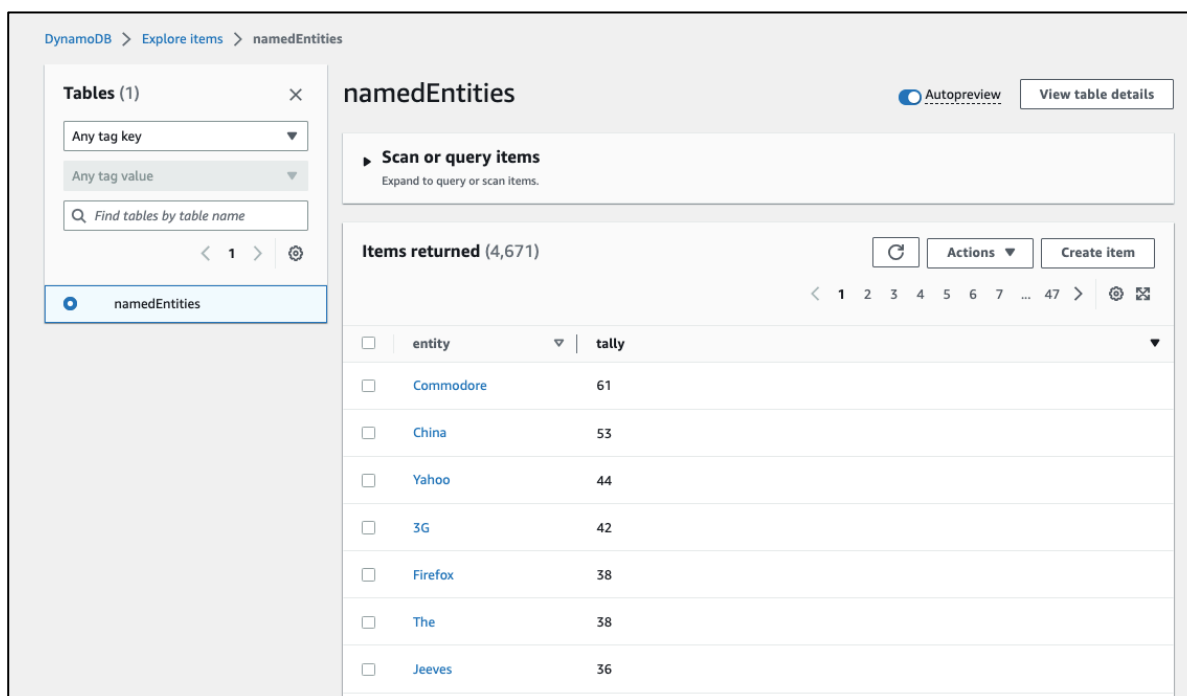


Fig 27. Verifying “accessDB” Correctly Stored Named Entities into “namedEntities” DynamoDB Table (Page 1)

The screenshot shows the AWS DynamoDB console interface for the 'namedEntities' table. On the left, there's a sidebar with 'Tables (1)' and a search bar. The main area shows the 'namedEntities' table with a 'Scan or query items' button. Below this, a table lists items with columns 'entity' and 'tally'. The items are: Media (13), Spam (13), Finnish (13), Warner (12), John (12), August (12), and Nielsen (12). The table indicates 4,671 items returned.

entity	tally
Media	13
Spam	13
Finnish	13
Warner	12
John	12
August	12
Nielsen	12

Fig 29. Verifying “accessDB” Correctly Stored Named Entities into “namedEntities” DynamoDB Table (Page 2)

Thus, we have verified that executing *only* the uploadFiles.py script correctly triggers “extractFeatures”, which extracts the named entities of newly uploaded files and uploads it to “tags-b00919848”; triggering “accessDB” to extract the named entities from the newly uploaded file and update our “namedEntities” table accordingly.

References:

- [1] Amazon Web Services Inc., “Amazon S3,” *Amazon Web Services Inc.* [Online], Available: <https://aws.amazon.com/s3/> [Accessed: July 11, 2023].
- [2] Amazon Web Services Inc., “Amazon DynamoDB,” *Amazon Web Services Inc.* [Online], Available: <https://aws.amazon.com/pm/dynamodb/> [Accessed: July 11, 2023].
- [3] Amazon Web Services Inc., “AWS Lambda,” *Amazon Web Services Inc.* [Online], Available: <https://aws.amazon.com/lambda/> [Accessed: July 11, 2023].
- [4] Amazon Web Services Inc., “Amazon CloudWatch,” *Amazon Web Services Inc.* [Online], Available: <https://aws.amazon.com/cloudwatch/> [Accessed: July 11, 2023].