

## REALTIME AND BATCH PREDICTION SERVICE

The first 2 pages contain the project proposal, and the following pages go into the implementation details.

**Abstract:** The following 2 goals were implemented.

**GOAL #1:** Batch prediction. An automated workflow to generate batch predictions using S3, Lambda and EC2.

**GOAL #2:** Realtime prediction using API gateway, Lambda and EC2.

Prediction refers to generating summarized text (using ML model) for a given corpus/body of text. For this project we have used the Amazon Food Reviews dataset (<https://www.kaggle.com/snap/amazon-fine-food-reviews>), and we will be generating summarized reviews of this dataset. For example: If the review is 10 sentences long, the summarized review will be about 3 sentences long. The ML model will pick the key sentences. An example is shown below.

### Amazon Food Review

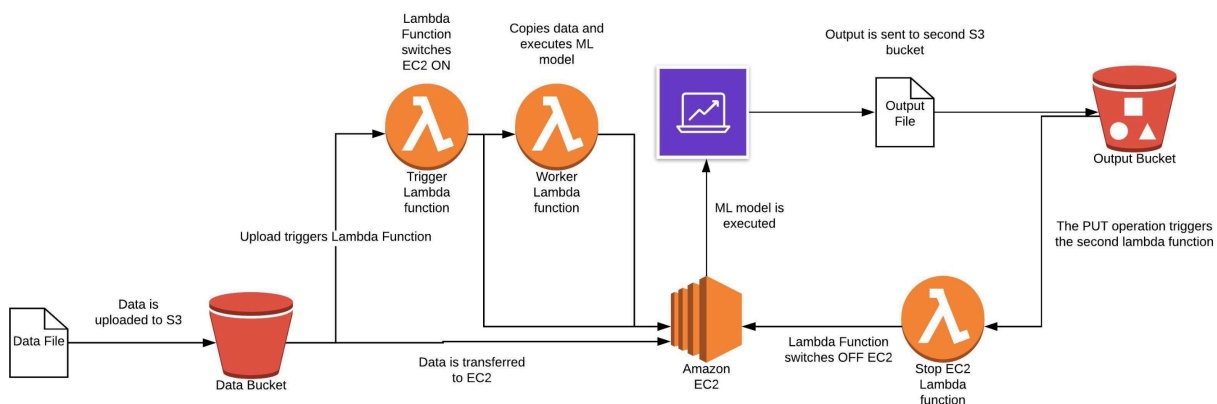
This taffy is so good. It is very soft and chewy. The flavors are amazing. I would definitely recommend you buying it. Very satisfying!!

Text Summarizer Model

### Summarized Review Output

"This taffy is so good. The flavors are amazing"

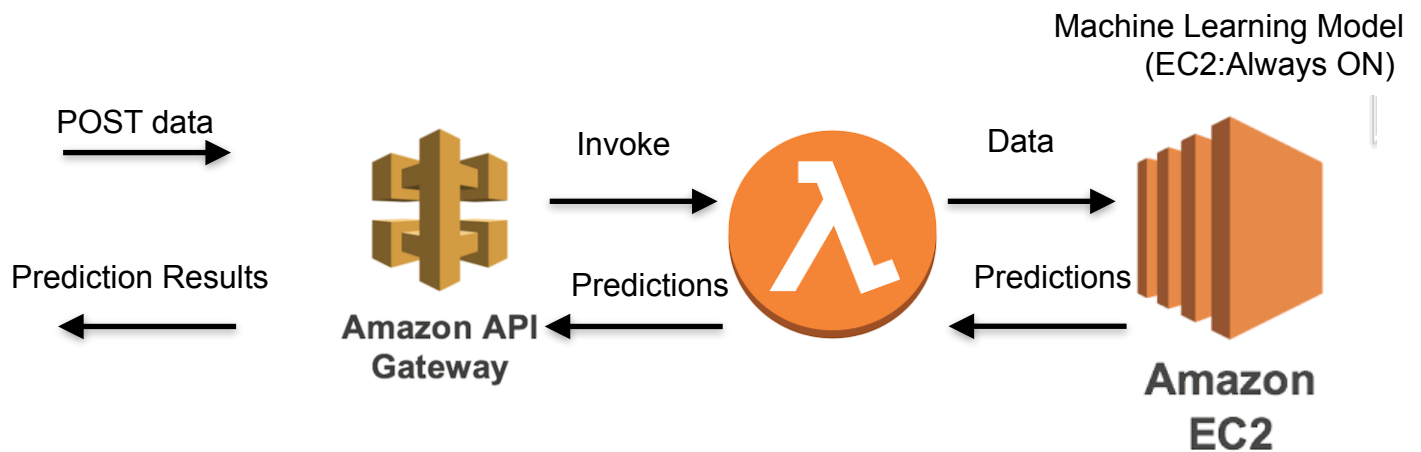
**GOAL #1("Compute as Glue"):** Visualization of automatic workflow. The below diagram has been borrowed from <https://towardsdatascience.com/automating-machine-learning-models-on-aws-bfa183fe4065>



**Steps:**

- 1) Data posted to input S3 will trigger a Lambda function ("Trigger Lambda") which will turn on the EC2 instance.
- 2) "Trigger Lambda" will invoke "Worker Lambda" that copies the data and code from S3 to EC2 instance, and directs the EC2 instance to execute the text summarizer ML model.
- 3) "EC2" instance executes the ML model and writes the model to output S3.
- 4) The output S3 triggers a Lambda function ("Stop Ec2 Lambda") to stop the EC2 instance.

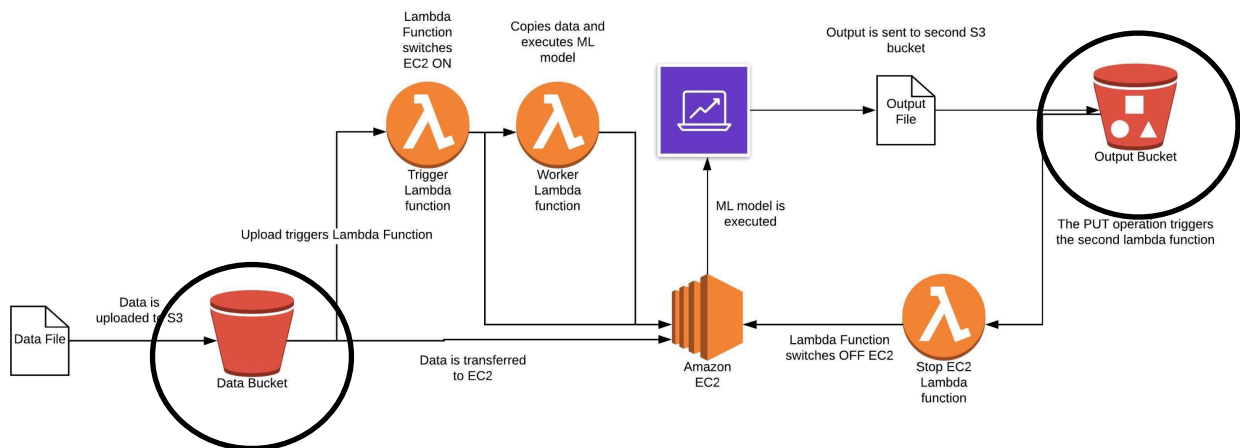
**Part 2 ("Compute as Backend"):** We will use the following building blocks to achieve the 2nd goal. This use case also enables the end users easy access to text summarizer service via the API Gateway.

**Steps:**

1. The API gateway invokes a lambda function when a "POST" request is sent on the configured resource.
2. The lambda function orchestrates the downloading of resources needed for the text summarizer model, executes the model, and collects the results, and passes the summarized text back to the api gateway.

## GOAL #1: IMPLEMENTATION DETAILS

### Generating batches of text summary



**S3 bucket:** Two S3 buckets were implemented namely cscie90-finalproject-inputdata and cscie90-finalproject-results. The S3 bucket cscie90-finalproject-inputdata was used to source the dataset, and the second S3 bucket cscie90-finalproject-results was used to write the summarized text result outputs.

Amazon Food Reviews file (<https://www.kaggle.com/snap/amazon-fine-food-reviews>) was uploaded to S3 bucket cscie90-finalproject-inputdata. This bucket was configured to trigger the lambda function on copy/put object event.

Under cscie90-finalproject-inputdata properties, events were configured as below to call “**trigger\_handler**” lambda function when any object was written/uploaded to this bucket.

**Events**

☐ PUT

☐ POST

☐ COPY

☐ Multipart upload completed

☒ All object create events

☐ Object in RRS lost

☐ Permanently deleted

☐ Delete marker created

☐ All object delete events

☐ Restore initiated

☐ Restore completed

☐ Replication time missed threshold

☐ Replication time completed after threshold

☐ Replication time not tracked

☐ Replication time failed

**Prefix****Suffix****Send to****Lambda**☒ 1 Active notifications

Cancel

Save

MyResultsEventForPut

Events ⓘ

☒ PUT
 ☐ All object delete events

☐ POST
 ☐ Restore initiated

☒ COPY
 ☐ Restore completed

☐ Multipart upload completed
 ☐ Replication time missed threshold

☐ All object create events
 ☐ Replication time completed after threshold

☐ Object in RRS lost
 ☐ Replication time not tracked

☐ Permanently deleted
 ☐ Replication time failed

☐ Delete marker created

Prefix ⓘ

e.g. images/

Suffix ⓘ

e.g. .jpg

Send to ⓘ

Lambda Function

Lambda

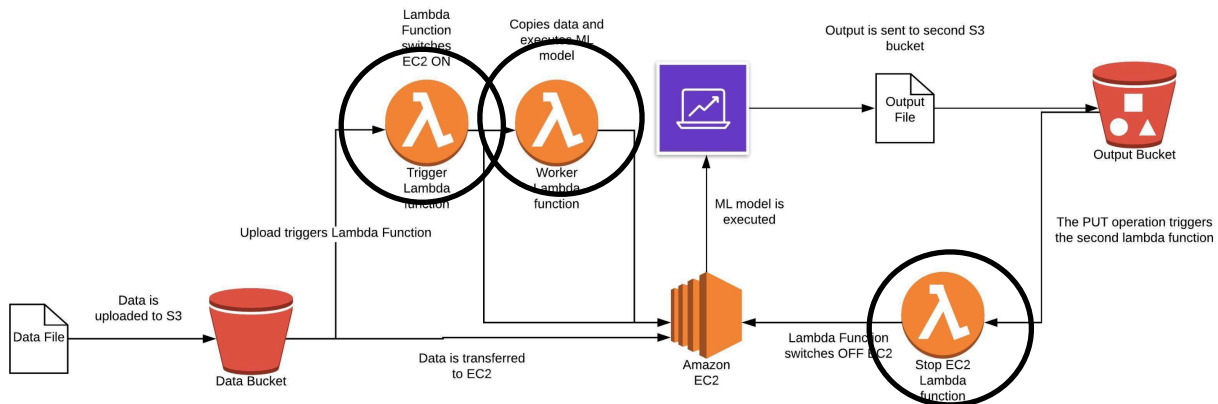
stopEC2

☒ 1 Active notifications
 

Cancel Save

Similarly the bucket cscie90-finalproject-results was configured to call **“stopEC2”** lambda function once results were posted/copied to this bucket.

## Lambda Functions



### Trigger\_Handler Lambda:

When dataset is uploaded to the input s3 bucket, the lambda function trigger\_handler is invoked.

This handler starts an EC2 instance, and invokes the lambda function worker\_function. Please note a policy was attached to this lambda function to allow access to EC2, and allow invocation of another lambda function. The source code and policy are shown in the below screenshots.

```
1 import json
2 import boto3
3
4 # Enter the region your instances are in. Include only the region without specifying Availability Zone; e.g.: 'us-east-1'
5 region = 'us-east-1'
6 # Enter your instances here: ex. ['X-XXXXXXX', 'X-XXXXXXX']
7 instances = ['X-XXXXXXX']
8
9 AWSAccessKeyId='secretid'
10 AWSSecretKey='secretkey'
11
12 def lambda_handler(event, context):
13     client = boto3.client('ec2', region_name=region, aws_access_key_id=AWSAccessKeyId, aws_secret_access_key=AWSSecretKey)
14     client.start_instances(instance_ids=instances)
15     print('started your instances: ' + str(instances))
16
17     #Get IP addresses of EC2 instances
18     #client = boto3.client('ec2')
19     instDict=client.describe_instances(
20         Filters=[{'Name':'tag:Environment','Values':['Dev']}]
21     )
22
23     hostList=[]
24
25     for r in instDict['Reservations']:
26         for inst in r['Instances']:
27             hostList.append(inst['PublicIpAddress'])
28
29     #Invoke worker function for each IP address
30     client = boto3.client('lambda')
31
32     for host in hostList:
33         print("Invoking worker_function on " + host)
34         invokeResponse=client.invoke(
35             FunctionName='worker_function',
36             InvocationType='Event',
37             LogType='Tail',
38             Payload='{"IP":"' + host + '"}'
39         )
40         print('response is' + str(invokeResponse))
41
42     return{
43         'message' : "Trigger function finished"
44     }
45
46
47
48
```

**Example Policy Attached to the Trigger Lambda function:**

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "VisualEditor0",
      "Effect": "Allow",
      "Action": "ec2:Describe*",
      "Resource": "*"
    },
    {
      "Sid": "VisualEditor1",
      "Effect": "Allow",

```

```

    "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "lambda:InvokeFunction"
    ],
    "Resource": [
        "arn:aws:lambda:us-east-1:926962607868:function:worker_function"
    ]
}
]
}

```

### Worker Function Lambda:

This function creates a SSH tunnel into the EC2 instance, and executes various s3 commands, runs the text summarizer “kmeans.py” model, and copies the results to the output s3 results bucket (cscie90-finalproject-results). One very important aspect of this lambda function was the installation of paramiko package which provides methods to SSH into an EC2 instance. The steps outlined in this blog (<https://towardsdatascience.com/hosting-your-ml-model-on-aws-lambdas-api-gateway-part-1-9052e6b63b25>) were followed to install the paramiko package into the worker\_function lambda layer.

Administrative EC2 access policy was attached to this lambda function to execute the s3 and python code on the EC2 instance.

```

import json
import boto3
import paramiko

def lambda_handler(event, context):
    s3_client = boto3.client('s3')
    #Download private key file from secure S3 bucket
    s3_client.download_file('cscie90-finalproject-virginia','virginia.pem', '/tmp/virginia.pem')

    k = paramiko.RSAKey.from_private_key_file("/tmp/virginia.pem")
    c = paramiko.SSHClient()
    c.set_missing_host_key_policy(paramiko.AutoAddPolicy())

    host=event['IP']
    print ("Connecting to " + host)
    c.connect( hostname = host, username = "ec2-user", pkey = k )
    print ("Connected to " + host)

    commands = [
        "aws s3 cp s3://cscie90-finalproject-inputdata/Reviews.csv --region us-east-1 /home/ec2-user/Reviews.csv",
        "python3 /home/ec2-user/kmeans.py",
        "aws s3 cp /home/ec2-user/top_500_summary.csv s3://cscie90-finalproject-results/summary.csv",
        "rm /home/ec2-user/top_500_summary.csv"
    ]

    for command in commands:
        print ("Executing {}".format(command))
        stdin, stdout, stderr = c.exec_command(command)
        print (stdout.read())
        print (stderr.read())

    return
    {
        'message' : "Script execution completed. See Cloudwatch logs for complete output"
    }

```

## A sample CloudWatch Log showing execution of worker lambda function: The logs are a useful resource to debug lambda functions.

CloudWatch > Log Groups > /aws/lambda/worker\_function > 2019/12/02[\$LATEST]f1bacbd99bf14d5fb286c49321e99791

**Try CloudWatch Logs Insights**

CloudWatch Logs Insights allows you to search and analyze your logs using a new, purpose-built query language. Click [here](#) to experience it. If you want to learn more, read [the AWS blog](#) or visit [our documentation](#).

Expand all Row Text

Filter events all 2019-12-01 (21:52:06)

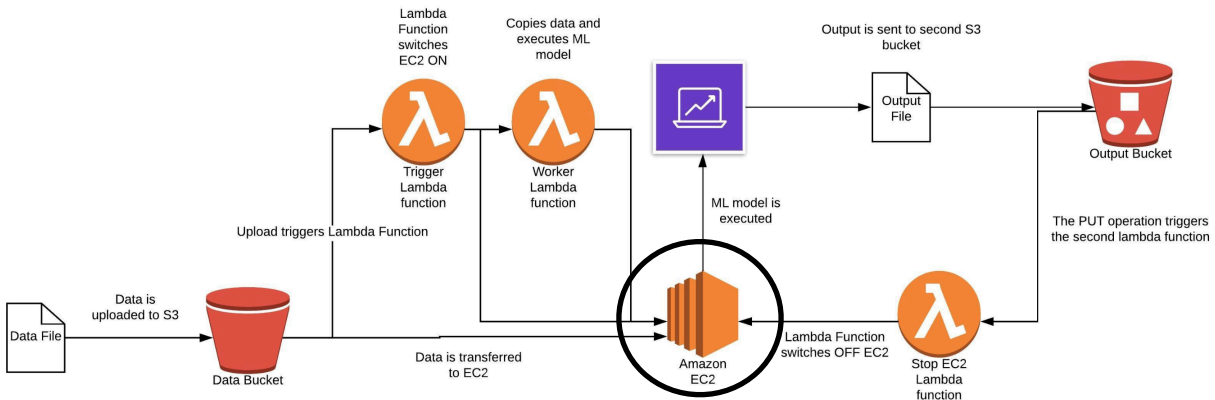
Time (UTC +00:00)	Message
2019-12-02	No older events found at the moment. <a href="#">Retry</a> .
21:48:32	START RequestId: 98a69e33-249a-4c33-994c-8a5b39269e6f Version: \$LATEST
21:48:34	Connecting to 54.146.185.40
21:48:34	Connected to 54.146.185.40
21:48:34	Executing aws s3 cp s3://cscie90-finalproject-inputdata/Reviews.csv --region us-east-1 /home/ec2-user/Reviews.csv
21:48:38	b'Completed 256.0 KiB/287.0 MiB (2.5 MiB/s) with 1 file(s) remainingvCompleted 512.0 KiB/287.0 MiB (4.8 MiB/s) with 1 file(s) remainingvCompleted 768.0 KiB/287.0 MiB (7.0 MiB/s) with 1 file(s) remainingvComple
21:48:38	b''
21:48:38	Executing python3 /home/ec2-user/kmeans.py
21:52:05	b' Id ... Text\n0 1 ... I have bought several of the Vitality canned d... \n1 2 ... Product arrived labeled as Jumbo Salted Peanut... \n2 3 ... This is a confection that has been around a fe... \n\n[3 rows x 10 columns]\nnum
21:52:05	b'Using TensorFlow backend.\n[nltk_data] Downloading package punkt to /home/ec2-user/nltk_data... \n[nltk_data] Package punkt is already up-to-date\n/home/ec2-user/kmeans.py:196: SettingWithCopyWarning:
21:52:05	Executing aws s3 cp /home/ec2-user/top_500_summary.csv s3://cscie90-finalproject-results/summary.csv
21:52:06	b'Completed 256.0 KiB/458.1 KiB (4.3 MiB/s) with 1 file(s) remainingvCompleted 458.1 KiB/458.1 KiB (3.4 MiB/s) with 1 file(s) remainingvupload: ./top_500_summary.csv to s3://cscie90-finalproject-results/summary
21:52:06	b''
21:52:06	Executing rm /home/ec2-user/top_500_summary.csv
21:52:06	b''
21:52:06	b''
21:52:06	END RequestId: 98a69e33-249a-4c33-994c-8a5b39269e6f
21:52:06	REPORT RequestId: 98a69e33-249a-4c33-994c-8a5b39269e6f Duration: 213912.98 ms Billed Duration: 214000 ms Memory Size: 128 MB Max Memory Used: 92 MB Init Duration: 370.70 ms
	No newer events found at the moment. <a href="#">Retry</a> .

## StopEC2 Lambda:

This function is invoked when results are posted to the output S3 bucket (cscie90-finalproject-results). Again, administrative EC2 access privileges were granted to this lambda function to stop the EC2 instance.

```
1 import json
2 import boto3
3
4 # Enter the region your instances are in. Include only the region without specifying Availability Zone; e.g.; 'us-east-1'
5 region = 'us-east-1'
6 # Enter your instances here: ex. ['X-XXXXXXX', 'X-XXXXXXX']
7 instances = ['X-XXXXXXX']
8
9 AWSAccessKeyId='secretId'
10 AWSSecretKey='secretKey'
11
12 def lambda_handler(event, context):
13
14     client = boto3.client('ec2', region_name=region, aws_access_key_id=AWSAccessKeyId, aws_secret_access_key=AWSSecretKey)
15     client.stop_instances(InstanceIds=instances)
16     print ('stopped your instances: ' + str(instances))
17
18     return {
19         'statusCode': 200,
20         'body': json.dumps('Successfully stopped EC2 from stopEC2 Lambda')}
21
22
```

## Text Summarizer Model (kmeans.py):



The model has been taken from the GitHub [https://github.com/Preetikasri/CognitiveComputing-TextSummarizer/blob/master/Text\\_summarizer\\_unsupervised.ipynb](https://github.com/Preetikasri/CognitiveComputing-TextSummarizer/blob/master/Text_summarizer_unsupervised.ipynb).

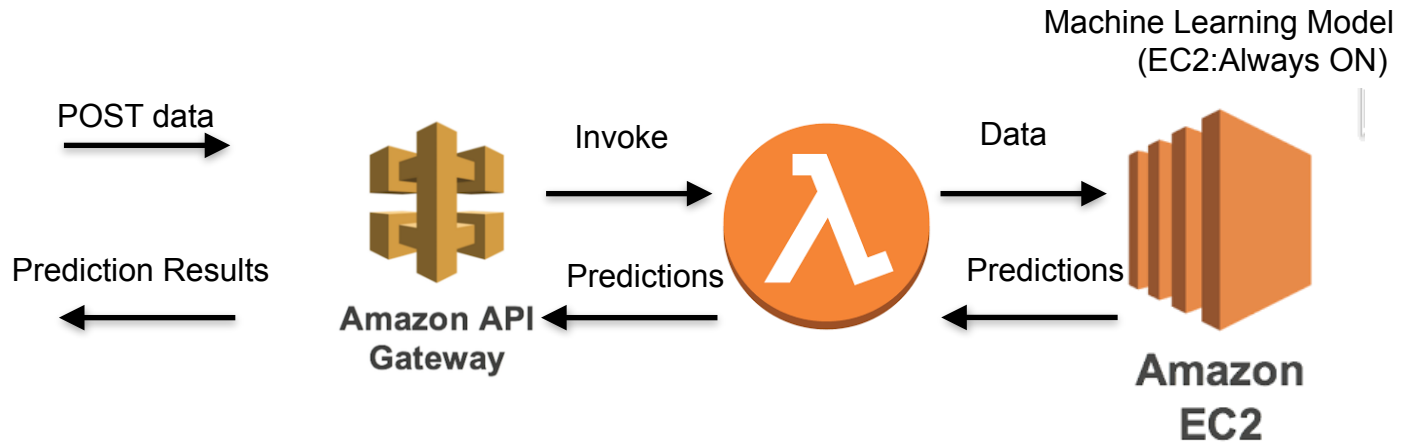
The source code for this model is in “kmeans.py” and is available at the GitHub link identified in the resources section towards the end of the document.

The model loads the glove embeddings for the tokens, and implements a kmeans algorithm. It essentially finds the cluster centroid and picks the sentences closest to the centroid. Essentially the centroid is the gist of the individual amazon food reviews. A sample of summarized text reviews are available in the file top\_500\_summary.csv on GitHub.

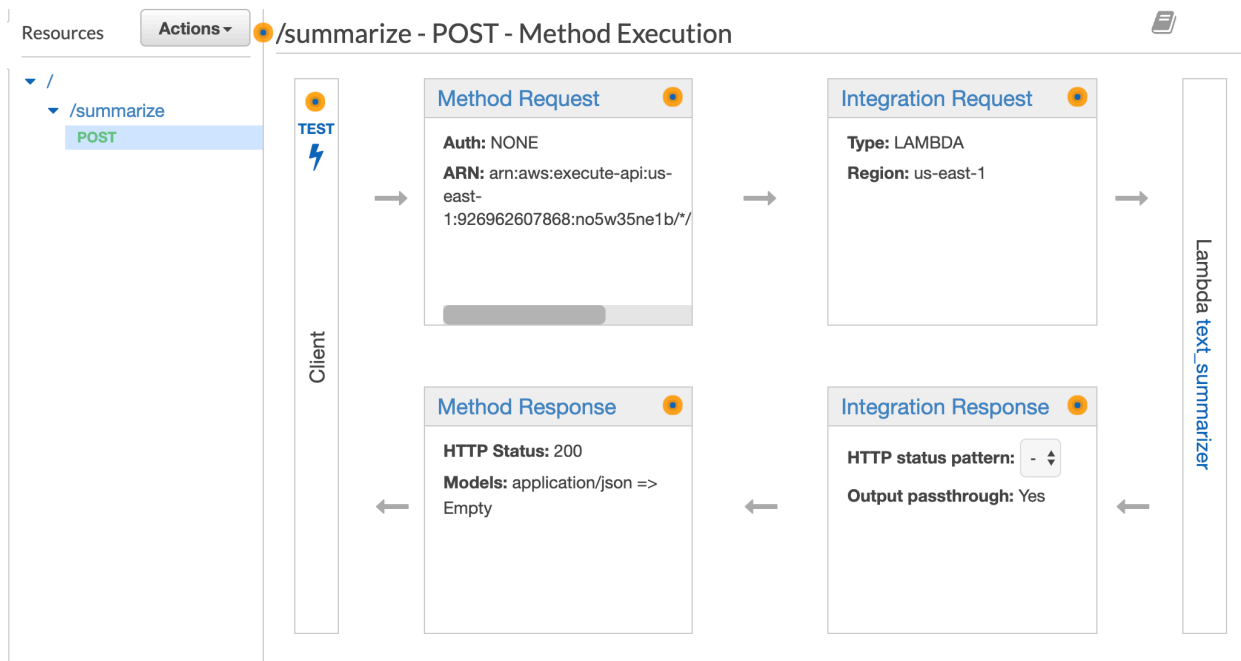


## GOAL #2 IMPLEMENTATION DETAILS:

Real time text summarization (predictions) via REST API call.



**API Gateway:** The API Gateway was configured on summarize resource, with action “POST”. The API Gateway calls the text\_summarizer lambda function when the end-user sends a POST message with the text that needs to be summarized.



## Lambda Function:

The API Gateway invokes the text\_summarizer lambda function. This function is very similar in functionality to the worker lambda function described earlier. It creates a SSH tunnel with the always ON EC2 instance, and takes the “text” message from API Gateway and runs the text summarizer model (“kmeans.py”) on the text, and returns the summarized text back to the API Gateway. The screenshot for the text summarizer lambda code is below. It has similar administrative privileges policy assigned to it to execute code on the EC2 instance.

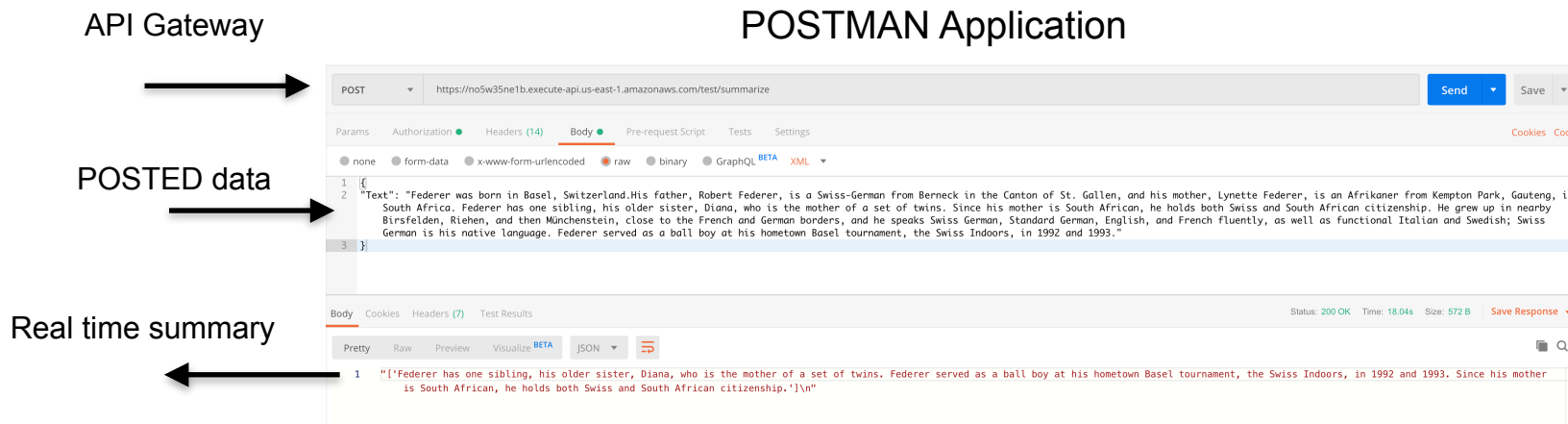
```
1 import json
2 import boto3
3 import paramiko
4
5 # Enter the region your instances are in. Include only the region without specifying Availability Zone; e.g.: 'us-east-1'
6 #region = 'us-east-1'
7 # Enter your instances here: ex. ['X-XXXXXXX', 'X-XXXXXXX']
8 #instances = ['i-0e220d9daf7348e8']
9
10 def lambda_handler(event, context):
11
12     print(event)
13     #Get IP addresses of EC2 instances
14     client = boto3.client('ec2')
15     instDict = client.describe_instances(
16         filters=[{'Name': 'tag:Text-Summarizer', 'Values': ['FinalProject']}])
17
18     hostList=[]
19     for r in instDict['Reservations']:
20         for inst in r['Instances']:
21             hostList.append(inst['PublicIpAddress'])
22
23     print(hostList)
24     s3_client = boto3.client('s3')
25     #Download private key file from secure S3 bucket
26     s3_client.download_file('cscie98-finalproject-virginia', 'virginia.pem', '/tmp/virginia.pem')
27
28     k = paramiko.RSAKey.from_private_key_file("/tmp/virginia.pem")
29
30     c = paramiko.SSHClient()
31     c.set_missing_host_key_policy(paramiko.AutoAddPolicy())
32
33     host=hostList[0]
34     #host='18.212.209.128'
35     #host='172.31.16.89'
36     print ("Connecting to " + host)
37     conn = c.connect( hostname = host, username = "ec2-user", pkey = k, timeout=10 )
38     if conn is None:
39         print ("Successfully Authenticated")
40     else:
41         print("Connection Not Successful")
42         print ("Connected to " + host)
43
44     input = event['Text']
45     input_str = "\"" + input + "\""
46     command_str = f'python3 kmeans.py {input_str}'
47     print(command_str)
48     commands = [
49         command_str
50     ]
51
52     for command in commands:
53         print ("Executing {}".format(command))
54         stdin, stdout, stderr = c.exec_command(command, timeout=120)
55         result = stdout.read().decode('ascii')
56         #print (stderr.read())
57
58     return result
59
60
```

## Text Summarizer Model:

The “kmeans.py” model is run on the EC2 instance. It takes the “text” argument passed to it from the text\_summarizer lambda function, and prints the summarized text to its console. The lambda function grabs the console output and passes it back to the API Gateway.

## Results:

Here is a screenshot of text summarizer model using POSTMAN. The sample text was taken from Roger Federer's wikipedia and sent via API Gateway POST on the summarize resource. The response in red is the summarized text from the model. The round trip time is approximately 30seconds.



## Resources:

- Project Report (PDF document): <https://github.com/amitrgupta27/cscie90-finalproject>
- Code: <https://github.com/amitrgupta27/cscie90-finalproject>
- Youtube URL (Project Overview): [https://youtu.be/aQPjzQAn\\_Ys](https://youtu.be/aQPjzQAn_Ys)
- Youtube URL (Deep Dive): <https://www.youtube.com/watch?v=s8wfpDI0Okk>
- Contact: [amitrgupta27@gmail.com](mailto:amitrgupta27@gmail.com)
- References:
  - <https://towardsdatascience.com/hosting-your-ml-model-on-aws-lambdas-api-gateway-part-1-9052e6b63b25>
  - <https://towardsdatascience.com/hosting-your-ml-model-on-aws-lambdas-api-gateway-part-2-23517609522b>
  - <https://towardsdatascience.com/automating-machine-learning-models-on-aws-bfa183fe4065>
  - <https://towardsdatascience.com/data-scientists-guide-to-summarization-fc0db952e363>