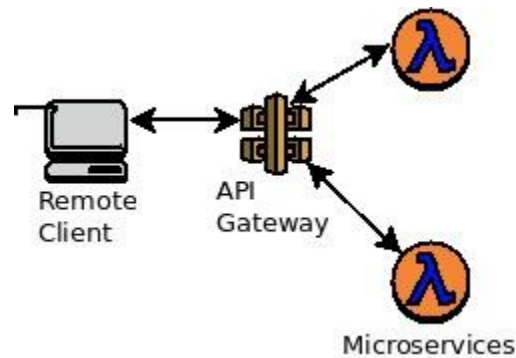


Tutorial 4 – Introduction to AWS Lambda with the Serverless Application Analytics Framework (SAAF)

Disclaimer: Subject to updates as corrections are found
Version 0.11
Scoring: 40 pts maximum

The purpose of this tutorial is to introduce creating Function-as-a-Service functions on the AWS Lambda FaaS platform, and then to create a simple two-service application where the application flow control is managed by the client:



This tutorial will focus on developing Lambda functions in Java using the Serverless Application Analytics Framework (SAAF). SAAF enables identification of the underlying cloud infrastructure used to host FaaS functions while supporting performance and resource utilization profiling of functions. SAAF helps identify infrastructure state to determine COLD vs. WARM function instances to help understand performance of the serverless Freeze-Thaw infrastructure lifecycle.

1. Download SAAF

To begin, using git, clone SAAF.

If you do not already have git installed, please do so.

Scroll down to find instructions on 'Install Git on Linux', and 'Debian/Ubuntu':

<https://github.com/git-guides/install-git>

For a full tutorial on the use of git, here is an old tutorial from TCSS 360:

http://faculty.washington.edu/wlloyd/courses/tcss360/assignments/TCSS360_w2017_Tutorial_1.pdf

If you prefer using a GUI-based tool, on Windows/Mac check out the GitHub Desktop:

<https://desktop.github.com/>

Once having access to a git client, clone the source repository:

```
git clone https://github.com/wlloyduw/SAAF.git
```

For tutorial #4, we will focus on using the SAAF provided AWS Lambda Java function template provided as a maven project. If you're familiar with Maven as a build environment, you can simply edit your Java Lambda function code using any text editor such as vi, emacs, pico/nano. However, working with an IDE tends to be easier, and many Java IDEs will open maven projects directly.

Next update your apt repository and local Ubuntu packages:

```
sudo apt update  
sudo apt upgrade
```

To install maven on Ubuntu:

```
sudo apt install maven
```

2. Build the SAAF Lambda function Hello World template

If you have a favorite Java IDE with maven support, feel free to try to open and work with the maven project directly. This project is confirmed to work in Apache NetBeans. Many students prefer using Microsoft Visual Studio Code with the "Extension Pack for Java". Other options include Eclipse, and IntelliJ.

Download Apache NetBeans IDE (NB) Installer for your platform:

<https://netbeans.apache.org/front/main/download/>

For Ubuntu 24.04, this will be "apache-netbeans-21-1_all.deb". Once downloaded, open a terminal in Ubuntu, navigate to the "Downloads" directory, and install the deb package:

```
# before installing, calculate and confirm authenticity using a SHA512 hash  
# compare with the SHA512 file downloaded from the website  
gpg --print-md SHA512 apache-netbeans_23-1_all.deb  
  
sudo dpkg -i apache-netbeans_23-1_all.deb
```

On Ubuntu, the netbeans snap package can be used as an installation alternative. (*this may take a while*)

For troubleshooting see: <https://snapcraft.io/install/netbeans/ubuntu>

```
sudo snap install netbeans
```

Alternatively, Microsoft Visual Studio Code can be used.

For Ubuntu 24.04 installation documentation see:

<https://linuxiac.com/how-to-install-vs-code-on-ubuntu-24-04-its/>

Once Visual Studio Code is installed, install the **Extension Pack for Java**:

<https://marketplace.visualstudio.com/items?itemName=vscjava.vscode-java-pack>

In addition, you may want (or need) to install additional plug-ins and extension packs to customize VS Code.

While the Netbeans installation may (or may not) automatically install the Java 21 Java Development Kit (JDK), it will be required to install the Java 21 JDK for Visual Studio Code:

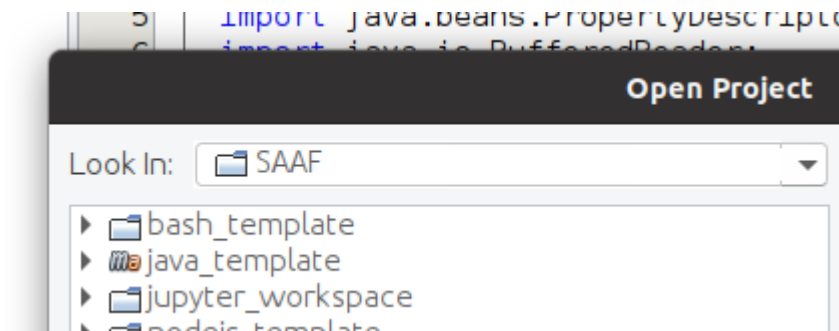
In Ubuntu, the Java 21 JDK can be installed with:

```
sudo apt install default-jdk
# After installing, verify the version of the java compiler:
javac -version
```

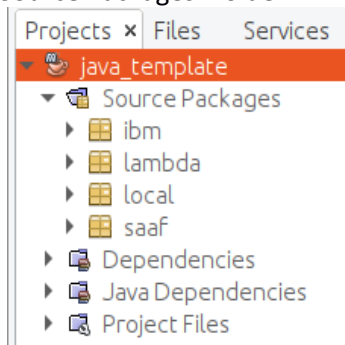
If working with Visual Studio Code, if the proper extensions are installed, you should be able to Open the Folder of the SAAF git project. Under the File menu, select “Open Folder” and navigate to the location where you clone the SAAF git hub repository, and select the SAAF/java_template directory. If Visual Studio Code is properly configured, it will scan your project, and recognize the Java project has a Maven build file. Then, in the lower left hand corner there should be a Maven menu listed. By opening the Maven menu and right-clicking on the java_template, there should be an option for “package”. If you select “package” it will build the Java JAR file that is required for deployment to AWS Lambda.

For Netbeans, once you’ve downloaded the IDE, you’ll be able to open the project directly without installing any plugins.

From Netbeans, Select “File | Open Project”, and navigate to where you have cloned the SAAF git project. Open the “SAAF” folder and then select “java_template”:



Then on the left-hand side, expand the “Source Packages” folder:



You’ll see a “lambda” package icon. Expand this.

This is where you’ll find the source code for a Hello World Lambda function that is provided as a starter template.

You will see four relevant class files:

- Hello.java** Provides an example implementation of a Java AWS Lambda Function. The `handleRequest()` method is called by Lambda as the entry point of your Lambda FaaS function. You'll see where in this function template your code should be inserted. `Hello.java` uses a `HashMap` for the request and response JSON objects. The incoming response JSON is "serialized" into a `hashmap` automatically by Lambda. The outgoing response JSON is created based on the `HashMap` that is returned.
- HelloPOJO.java** `HelloPOJO.java` is the same as `Hello.java` except that instead of using a `HashMap` for the Request (incoming) data, instead an explicitly defined Request class is defined with getter and setter methods to accept input from the user. The advantage with `HelloPOJO` is the Request object can perform post-processing on input parameters provided from the function caller before they are used. Post-processing includes operations such as formatting data or transforming values into another form before actual use in the FaaS function. User inputs to the FaaS function could trigger other behavior in the FaaS function automatically when the values are loaded.
- HelloMain.java** `HelloMain.java` is identical to `Hello.java` except that it also contains a public static void `main()` method to allow command line execution of the function package. This template is provided as an example. This allows Lambda functions to be first tested locally on the command line before deployment to Lambda. The local implementation could also be used to facilitate off-line unit testing of FaaS functions. As you develop your FaaS function, it will be necessary to continue to add to the implementation of the `main()` method to include required parameters for interacting with the function. The `main()` method creates a mock `Context()` object which fools the program into thinking it is running in context of AWS Lambda.
- Request.java** This class is a Plain Old Java Object (POJO). You'll want to define getter and setter methods and private variables to capture data sent from the client to the Lambda function. JSON that is sent to your Lambda function is automatically marshalled into this Java object for easy consumption at runtime.
- Response.java** (REMOVED) There is no longer a Response class POJO. This has been removed in favor of simply using a `HashMap`. A Response POJO could be implemented alternatively to add logic to getter and setter methods to perform data formatting, transformation, or validation operations.

In Netbeans, if you see exclamation marks on the source file icons, where the exclamation mark is on every file, this indicates that the Java Platform for the IDE is not properly configured. Close the project by clicking on "java_template" and selecting "Close".

Then, under Tools | Java Platforms, you will need to add an entry to point to JDK 21 on your system. JDK 21 is typically installed at: `/usr/lib/jvm/default-jvm`. Then under Tools | Options under the Java tab, configure your newly provided Java Platform under the “Java” tab. Verify settings under the “Java Shell” and “Maven” sub-tabs, and then save the settings. Then, reopen the project and check if the exclamation mark is resolved.

Additional versions of Java such as Java 11 or 17 can also be installed (not required). In Netbeans, you can toggle between different Java versions easily under Tools | Java Platforms.

On the Ubuntu command line, it is also possible to install multiple versions of Java. This enables working directly from the CLI to compile projects etc. For this, check the version of Java currently used. Make sure to match the version of functions to be deployed on AWS Lambda. To check which Java versions are installed, use the following command, but do not select a version, just press ENTER when prompted:

```
sudo update-alternatives --config java
```

In most cases, you will have just one version installed. It is possible to install multiple versions of Java, and switch between them on the command-line using the “sudo update-alternatives” command. **This sets the version of Java in the command-line environment.** This is different than the version of Java that is configured for the Netbeans project.

To inspect the version of Java used in your project in Netbeans, in the project explorer on the left-hand side, right-click on the project name, and select “**Properties**” at the bottom of the list. First, under the “**Build**” option, select “**Compile**”, and in the dialog box select the proper Java Platform, such as JDK 11.

For Fall 2024, since Ubuntu 24.04 LTS defaults to Java 21, this tutorial focuses on use of Java 21.

Now compile the project using maven from the IDE:

From the NetBeans IDE right click on the name of the project “java_template” in the left-hand list of Projects and click “**Clean and Build**”.

Now try compiling directly from the command line, under the “SAAF/java_template” directory:

```
cd {base directory where project was cloned}/SAAF/java_template/  
# Clean and remove old build artifacts  
mvn clean -f pom.xml
```

Then to build the project jar file:

```
# Rebuild the project jar file  
mvn verify -f pom.xml
```

3. Test Lambda function locally before deployment

From a terminal, navigate to:

```
cd {base directory where project was cloned}/SAAF/java_template/target
```

Execute your function from the command line to first test your Lambda function locally:

```
java -cp lambda_test-1.0-SNAPSHOT.jar lambda.HelloMain Susan
```

Output should be provided as follows:

```
cmd-line param name=Susan
function result:{cpuType=11th Gen Intel(R) Core(TM) i7-1165G7 @ 2.80GHz,
cpuNiceDelta=0, vmuptime=1729139076, cpuModel=140, linuxVersion=Linux lapetus
6.8.0-47-generic #47-Ubuntu SMP PREEMPT_DYNAMIC Fri Sep 27 21:40:26 UTC 2024
x86_64 x86_64 x86_64 GNU/Linux, cpuSoftIrqDelta=0, cpuUsrDelta=0, uid=d6775673-
4fce-4ee7-be21-920cf0f3d790, platform=Unknown Platform, contextSwitches=2783232,
cpuKrn=1613, cpuIdleDelta=0, cpuIowaitDelta=0, newcontainer=1, cpuNice=18,
startTime=1729142606442, lang=java, cpuUsr=4483, majorPageFaultsDelta=0,
freeMemory=217836, value=Hello Susan! This is from a response object!,
frameworkRuntime=70, contextSwitchesDelta=168, frameworkRuntimeDeltas=2,
vmcpusteal=0, cpuKrnDelta=0, cpuIdle=310260, runtime=84, message=Hello Susan! This
is a custom attribute added as output from SAAF!, version=0.5, cpuIrqDelta=0,
pageFaultsDelta=210, cpuIrq=0, totalMemory=4010468, cpuCores=1, cpuSoftIrq=447,
cpuIowait=2692, endTime=1729142606526, majorPageFaults=6603, vmcpustealDelta=0,
pageFaults=4187686, userRuntime=12}
```

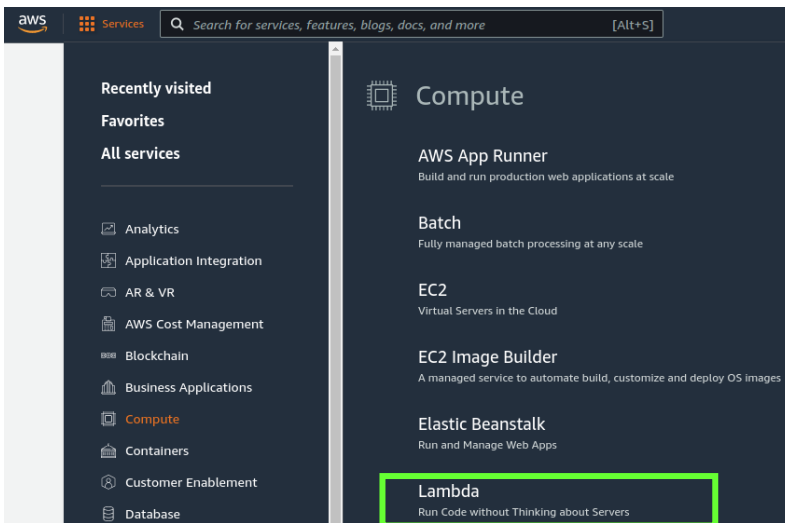
Whoa! That's a lot of output! The actual Lambda function output is highlighted.

Other values represent data collected by the SAAF profiling framework. Of course, since you're testing locally, this profiling data is for your local Linux environment, not the cloud.

4. Deploy the function to AWS Lambda

If the Lambda function has worked locally, the next step is to deploy to AWS Lambda.

Log into your AWS account, and under "services" locate "Lambda" by searching, or selecting "Compute":



Click the button to create a new Function:

[Create function](#)

Using the wizard, use the "Author from scratch" mode.

Next provide the following values:

Lambda > Functions > Create function

Create function Info

Choose one of the following options to create your function.

- Author from scratch**
Start with a simple Hello World example.
- Use a blueprint**
Build a Lambda application from sample code and configuration presets for common use cases.
- Container image**
Select a container image to deploy for your function.

Basic information

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

Function name must be 1 to 64 characters, must be unique to the Region, and can't include spaces. Valid characters are a-z, A-Z, 0-9, hyphens (-), and underscores (_).

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

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.

Lambda will create an execution role named Hello-UWT-f2024-role-jasa0kvk, with permission to upload logs to Amazon CloudWatch Logs.

Function name: hello (choose a unique name of your choice)

Runtime: Java 21

Execution Role: "Create a new role with basic Lambda permissions"

(Additional policies and permissions can be added to this role if needed.)

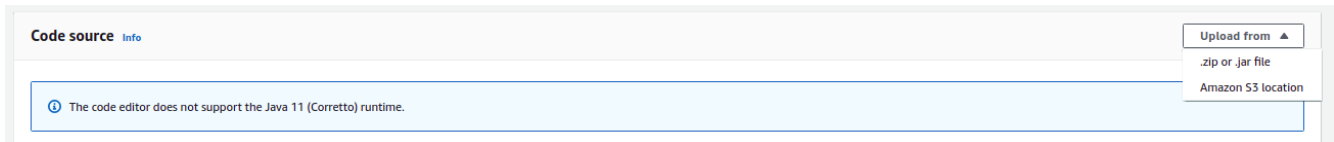
(Roles can be inspected under IAM | Roles in the AWS Management Console)

Once filling the form, click the button:

Create function

Note: if your Ubuntu Linux environment has limited memory (less than 5GB), it is possible this step will fail. If so, try increasing available memory.

Next, upload your compiled Java JAR file to AWS Lambda. Under “Code” and “Code Source”, select “Upload from” and “.zip or .jar file”.



Click the “Upload” button to navigate and locate your JAR file. The jar file is under the “target” directory. It should be called “lambda_test-1.0-SNAPSHOT.jar”. Once selecting the file press ‘Save’ to upload.

Next scroll down to “Runtime settings”. Click **Edit**, and in the dialog box change the “Handler” to:

```
lambda.Hello : :handleRequest
```

***** IF THE HANDLER IS NOT UPDATED, LAMBDA WILL NOT BE ABLE TO LOCATE THE ENTRY POINT TO YOUR CODE. THE LAMBDA FUNCTION WILL FAIL TO RUN *****

5. Create an API-Gateway REST URL

Next, in the AWS Management Console, navigate to the **API Gateway**.

This appears under the Network & Content Delivery services group, but using the search bar may be the fastest way:

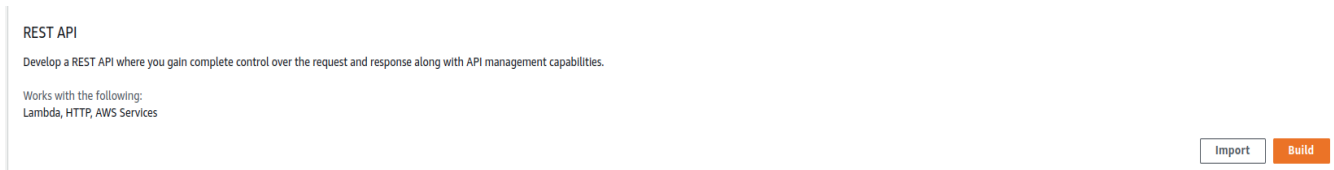


The very first time you visit the API Gateway console screen, there may be a “splash” screen. IF so, select the button:



If this is not your first visit, just click the ‘**Create API**’ button.

Next choose an API type. We will select “**REST API**” and click **Build**:



Click on the “Build” button. Next, specify these settings for the API details:

select: <*> NEW API

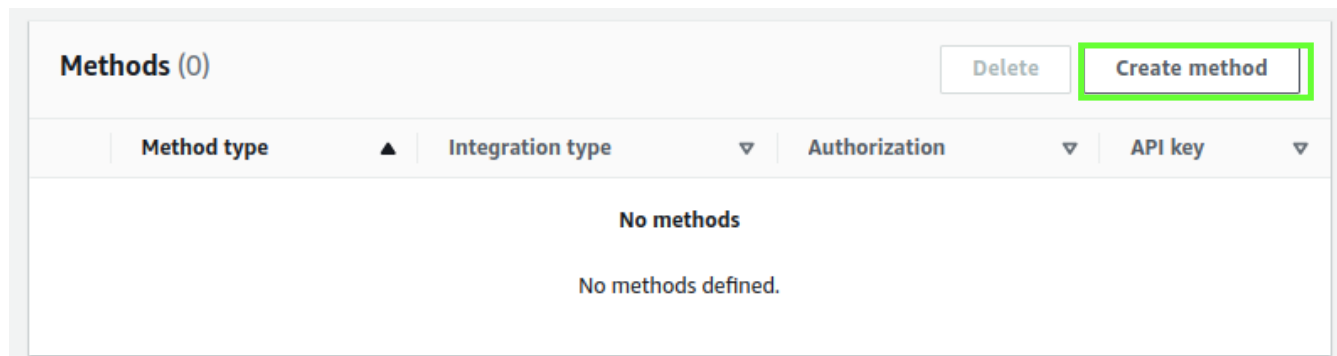
API Name: hello_uwt

Description: <can leave blank>

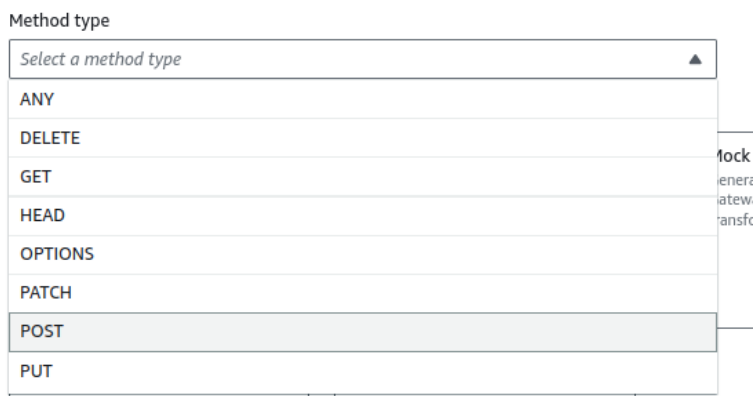
Endpoint Type: choose Regional

Press the “**Create API**” button.

Next, under Method(0), press the “**Create Method**” button:

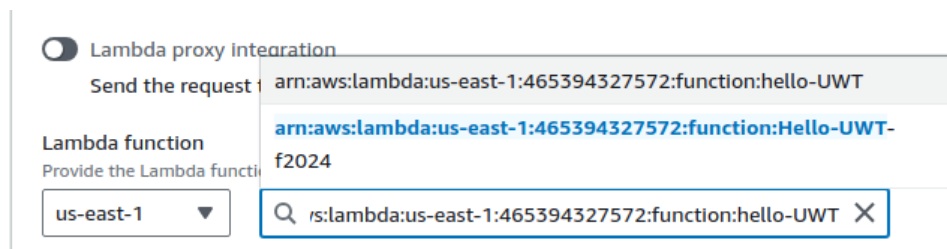


Select the ‘Post’ Method Type from the drop down list :



Next, confirm that “Lambda function” is the integration type.

Next, select the Lambda function. Search and find your function in the list:

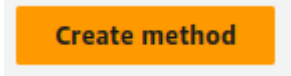


Confirm that the default timeout is 29000 milliseconds.

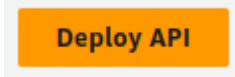
This is the maximum runtime for a synchronous Lambda function that is invoked using an HTTP REST endpoint provided using the API Gateway.

The API Gateway integration timeout for synchronous calls can be set between 50 and 29,000 milliseconds. Be sure to provide the maximum synchronous timeout “29000” milliseconds.

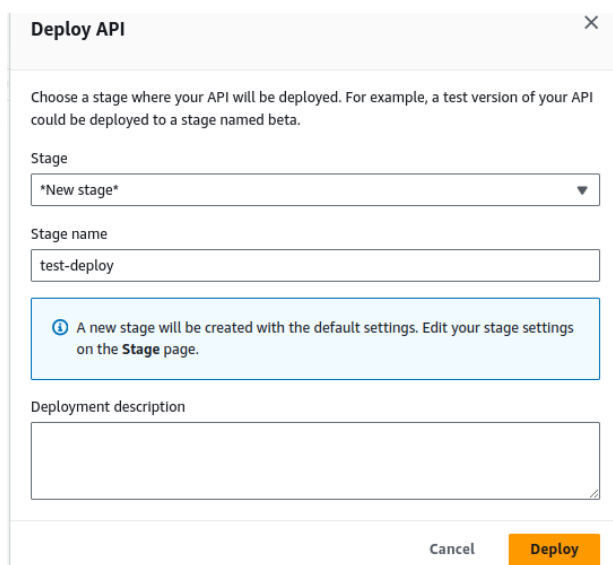
Now click ‘Create method’



Next, click on the “Deploy API” button:



Next complete the form. Select the “Stage”, as “New stage” and assign a Stage name. You can write something like ‘test-deploy’:

A dialog box titled "Deploy API" with a close button (X) in the top right corner. The dialog contains the following elements: a paragraph of text: "Choose a stage where your API will be deployed. For example, a test version of your API could be deployed to a stage named beta."; a "Stage" dropdown menu with the selected option being "*New stage*"; a "Stage name" text input field containing the text "test-deploy"; a light blue information box with a question mark icon and the text: "A new stage will be created with the default settings. Edit your stage settings on the Stage page."; a "Deployment description" text area which is currently empty; and at the bottom, two buttons: "Cancel" and "Deploy".

When complete, press the [Deploy] button.

The API-Gateway supports configuring RESTful Uniform Resource Locators (URLs) to associate with Lambda function backends so that clients can invoke the function just like the weather service in Tutorial #2. The API-Gateway is not limited to AWS Lambda functions. It can also point to other backend services hosted by AWS. When configuring a URL for a RESTful webservice, the AWS hosted URL acts as a proxy that routes incoming traffic to the configured backend, in this case the Hello Lambda function. This enables traffic to be routed using a URL through the API-Gateway, where the API-Gateway acts as an intermediary. Traditionally a proxy layer like the API Gateway can be used to introduce common logging, authentication, and/or features to RESTful backends.

Using the API-Gateway it is possible to host multiple versions of a function to support Agile/devops software development processes. An organization may want to maintain multiple LIVE versions of a function in various stages of development such as: (dev)elopment, test, staging, and (prod)uction

Once, deployed, locate and COPY the “InvokeURL” for your AWS Lambda function, from the “Stage details”:

Stage details [Info](#)

Stage name test-deploy	Rate Info 10000	Web ACL -
API cache ⊖ Inactive	Burst Info 5000	Client certificate -

Invoke URL
📄 [https://\[redacted\].amazonaws.com/test-deploy](https://[redacted].amazonaws.com/test-deploy)

COPY THE INVOKE URL TO THE CLIPBOARD:

Press the COPY ICON to the left of the URL.

Use this URL in the callservice.sh test script below.

6. Install package dependencies and configure your client to call Lambda

Return to the command prompt and create and navigate to a new directory

```
cd {base directory where project was cloned}/SAAF/java_template/  
mkdir test  
cd test
```

Using a text editor such as vi, pico, nano, vim, or gedit, create a file called “callservice.sh” as follows:

```
#!/bin/bash  
  
# JSON object to pass to Lambda Function  
json="{\"name\":\"Susan\u0020Smith\", \"param1\":1, \"param2\":2, \"param3\":3}"  
  
echo "Invoking Lambda function using API Gateway"  
time output=`curl -s -H "Content-Type: application/json" -X POST -d $json {INSERT  
API GATEWAY URL HERE}`  
echo ""  
  
echo ""  
echo "JSON RESULT:"  
echo $output | jq  
echo ""  
  
echo "Invoking Lambda function using AWS CLI (Boto3)"  
time output=`aws lambda invoke --invocation-type RequestResponse --cli-binary-  
format raw-in-base64-out --function-name {INSERT AWS FUNCTION NAME HERE} --region  
us-east-1 --payload $json /dev/stdout | head -n 1 | head -c -2 ; echo`  
  
echo ""
```

```
echo "JSON RESULT:"
echo $output | jq
echo ""
```

Replace {INSERT API GATEWAY URL HERE} with your URL.
Be sure to include the small quote mark at the end: `

This quote mark is next to the number 1 on US keyboards.

Next, locate the lines:

```
echo "Invoking Lambda function using AWS CLI"
time output=`aws lambda invoke --invocation-type RequestResponse --cli-binary-
format raw-in-base64-out --function-name {INSERT AWS FUNCTION NAME HERE} --region
us-east-1 --payload $json /dev/stdout | head -n 1 | head -c -2 ; echo`
```

Replace {INSERT AWS FUNCTION NAME HERE} with your Lambda function name "hello" (or whatever name you have used for your function – **do not include the parentheses**).

If using the AWS CLI major version 1.x, then the cli-binary-format binary must be removed.

Save the script and then provide execute permissions:

```
chmod u+x callservice.sh
```

Before running this script, it is necessary to install some packages.

You should have curl installed from tutorial #2. If not, please install it:

```
sudo apt install curl
```

Next, install the AWS command line interface (CLI) (***this should have been completed previously for Tutorial 0, but if not, do it now***):

```
sudo apt install awscli
```

Please refer to Tutorial 0 for detailed instructions for configuring the AWS CLI.

Note that running 'aws configure' in Tutorial 0 will create two hidden files at:

/home/ubuntu/.aws/config

/home/ubuntu/.aws/credentials

Use "ls -alt /home/ubuntu/.aws" to see them.

At any time, if needing to update the configuration, these files can be edited manually, or "aws configure" can be re-run. Amazon suggests changing the access key and secret access key every 90 days.

**NEVER UPLOAD YOUR ACCESS KEYS TO A GIT REPOSITORY.
AVOID HARD CODING THESE KEYS DIRECTLY IN SOURCE CODE WHERE FEASIBLE.**

Now install the “jq” package if you haven’t already from tutorial #2:

```
sudo apt install jq
```

7. Test your Lambda function using the API-Gateway and AWS CLI

It should now be possible to test your Lambda function using the `callservice.sh` script.

Run the script:

```
./callservice.sh
```

Output should be provided (abbreviated below):

```
Invoking Lambda function using API Gateway

real    0m0.129s
user    0m0.017s
sys     0m0.009s
\"

JSON RESULT:
{
  "cpuType": "Intel(R) Xeon(R) Processor @ 2.50GHz",
  "cpuNiceDelta": 0,
  "vmuptime": 1698054046,
  "cpuModel": "63",
  "linuxVersion": "Linux 169.254.39.177 5.10.186-200.751.amzn2.x86_64 #1 SMP Wed Aug 23
03:37:53 UTC 2023 x86_64 x86_64 x86_64 GNU/Linux",
  "cpuSoftIrqDelta": 0,
  "cpuUsrDelta": 0,
  "uuid": "14b99804-bdb2-403d-8f23-64a30d307674",
  "platform": "AWS Lambda",
  "contextSwitches": 27020,
  "cpuKrn": 90,
  "cpuIdleDelta": 3,
  "cpuIowaitDelta": 0,
  "newcontainer": 0,
  "cpuNice": 0,
  "startTime": 1698054523530,
  "lang": "java",
  "cpuUsr": 100,
  "majorPageFaultsDelta": 0,
  "freeMemory": "504404",
  "value": "Hello Susan Smith! This is from a response object!",
  .....
}
```

The script calls Lambda twice. The first instance uses the API gateway. As a synchronous call the curl connection is limited to 29 seconds.

The second instance uses the AWS command line interface. This runtime is limited by the AWS Lambda function configuration. It can be set to a maximum of 15 minutes. The default is 15 seconds. **Both of these calls are performed synchronously to AWS Lambda.**

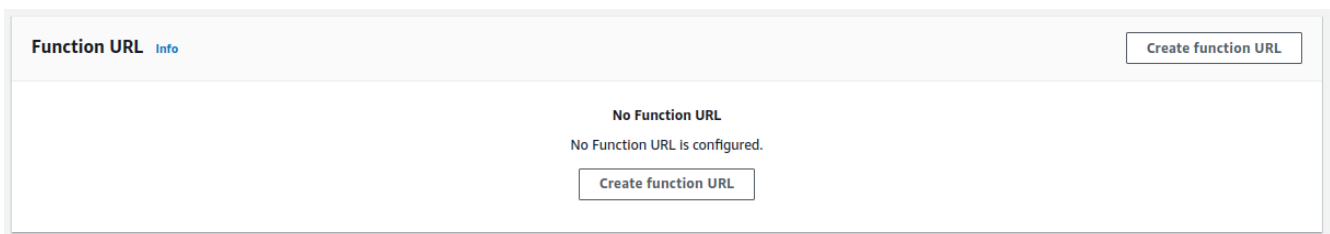
7B. (NEW) AWS Lambda Function URLs - BONUS

As of April 2022, AWS introduced Function URLs as an alternative to creating http REST endpoints using the API Gateway. This provides a faster way to create a REST URL that is associated with a Lambda function with one caveat (issue). For Function URLs, the Request object that is provided by curl is nested under the **body** tag. This means that the function source code must first read the body tag/value pair, and then convert the value of the body tag into a JSON object or Python dictionary. This adds an extra step, and makes code deployed using the API Gateway (or invoked using the AWS CLI) not directly compatible with code using Function URLs.

You can read about this feature in the April 2022 press release:

<https://aws.amazon.com/blogs/aws/announcing-aws-lambda-function-urls-built-in-https-endpoints-for-single-function-microservices/>

To create a Function URL for your Lambda function, from the AWS Lambda GUI, click on the **Configuration** tab. On the left-hand side select **Function URL**. Creating a Function URL is then a simple matter of clicking the **Create function URL** button.



For the 'Auth type' it will be easiest to use NONE. Function permissions can be restricted using advanced security capabilities by configuring specific IAM users and roles to have access. Additionally cross-origin resource sharing (CORS) can be used to restrict access further. Press **Save** to then create the Function URL.

The Lambda function's handler must be adapted to support Function URLs. The request object is now embedded in the body key-value pair of the input. Add the following code at the start of the `handleRequest()` method in **Hello.java** to support Function URLs after the "START FUNCTION IMPLEMENTATION" comment:

```
// Add imports at the top of the file (in netbeans can press ctrl-shift-I to auto-update imports)
import com.google.gson.Gson;
import com.google.gson.reflect.TypeToken;

// Get body which is nested as a parameter in the request when using a function URL
Object reqo = request.get("body");
if (reqo != null)
{
    String req = reqo.toString();
    // Create a hash map from the body
    HashMap<String, String> map = new Gson().fromJson(req, new TypeToken<HashMap<String,
String>>() { }.getType());
    // copy individual key value pairs from nested body to request obj
    request.put("name",map.get("name"));
}
```

A third invocation method can now be added to `call_service.sh`. Duplicate the http REST call in `call_service.sh` and update the URL to use the Function URL.

```

echo "Invoking Lambda function using Function URL"
time output=`curl -s -H "Content-Type: application/json" -X POST -d $json {INSERT FUNCTION URL
HERE}`
echo ""

echo ""
echo "JSON RESULT:"
echo $output | jq
echo ""

```

Try invoking your function now with all three calling conventions. Compare the function invocation times. Run the script ten times to make sure the function is warm. Inspect the “real” time for the calls. Compare the real times for the API Gateway Function URL, the Lambda Function URL, and the AWS CLI. ***** Which method reports the fastest function roundtrip time? Is there a clear winner? *****

Optional: Function Deployment from the Command Line and Use of Availability Zones

SAAF provides a command line tool that automates deploying and updating FaaS functions to different cloud providers. Here, we demonstrate the use for the hello function for AWS Lambda.

Navigate to:

```
cd {base directory where project was cloned}/SAAF/java_template/deploy
```

Backup the config.json script:

```
cp config.json config.json.bak
```

Now modify config.json to deploy your hello function:

```

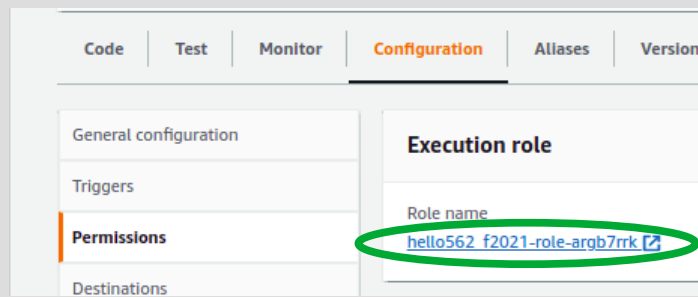
{
  "README": "See ./tools/README.md for help!",
  "functionName": "hello",
  "lambdaHandler": "lambda.Hello::handleRequest",
  "lambdaRoleARN": "arn:aws:iam::465394327572:role/service-role/simple_microservice_rolef19",
  "lambdaSubnets": "",
  "lambdaSecurityGroups": "",
  "lambdaEnvironment": "Variables={EXAMPLEVAR1=VAL1,EXAMPLEVAR2=VAL2}",
  "ibmHandler": "ibm.Hello",

  "test": {
    "name": "Bob"
  }
}

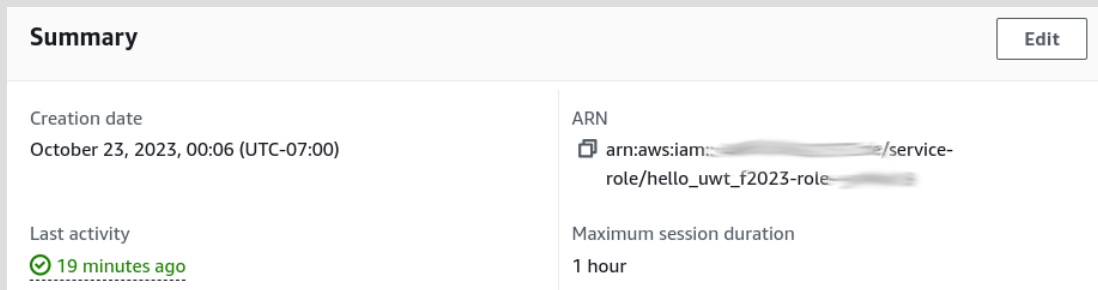
```

Function name: specify your function name ‘hello’.

LambdaRoleARN: This is the Amazon Resource Name (ARN) for the Lambda Role previously created for the Lambda function. The ARN can be found by editing the Lambda function configuration in the AWS management console web GUI. Under the **Configuration** tab, select **Permissions** on the left. Where it says **Role name**, click on the link and open the role under Identity Access Management.



This opens the role for editing in the IAM console.



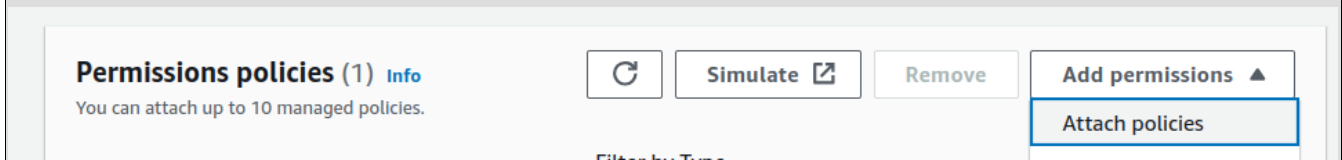
At the top of the Role Summary you'll see the **Role ARN** name. Click on the **COPY** icon on the LEFT to copy the ARN name to the clipboard. Paste this into your `config.json` file for the ARN.

The other attributes of note include **lambdaSubnets** and **lambdaSecurityGroups**.

A subnet specifies a virtual network within a Virtual Private Cloud (VPC). Selecting a subnet allows the function to be deployed to a specific Availability Zone within an AWS Region. An availability zone is a specific physical data center. These facilities are miles apart and considered physically separate locations.

The motivation to locate a Lambda function in an availability zone is to co-locate the function with other cloud resources that share the VPC. This way virtual machines (clients) and Lambda functions (backends) can be deployed in the same data center (aka physical site). This co-location provides the lowest possible network latency as all network traffic is local. The network communication between resources does not have to leave the ~ physical building.

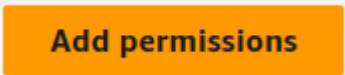
<OPTIONAL - VPC Setup – will be reviewed again in Tutorial #6> To create a Lambda function in a VPC, the Execution Role must be modified to include the `AWSLambdaVPCLambdaAccessExecutionRole` policy. This policy can be added when copying the ARN name to setup `config.json`. In the permission policies frame, drop down the "Add permissions" list and select "Attach policies" :



Then search for “VPC” policies and select the policy by finding the policy and checkmarking it: **AWSLambdaVPCAccessExecutionRole**.



Click the button to add the permissions:



This will attach the policy to your Lambda Execution role.

This is required to deploy a Lambda function to a VPC.

In the AWS Lambda function GUI, explore these options under **VPC**. On the Left select “**VPC**” and press the ‘**Edit**’ button. To deploy a function to a VPC, first select the “Default VPC”. This enables the “**Subnets**” drop-down list. By default AWS has provided a subnet for each availability zone in the Region. These subnet IDs are what is added to `config.json` to deploy the Lambda function to a specific availability zone.

You may ignore the error message that says: “*We recommend that you choose at least 2 subnets for Lambda to run functions in high availability mode.*”

High availability is a great feature for production deployment.

For development, experimentation, performance testing, and research however, it may be preferable for reproducing results to execute function code from the availability zone (aka datacenter) every day. Executing the function in random availability zones (what happens without a VPC) may increase the likelihood of receiving different hardware (e.g. CPUs) and the latency to the availability zones may vary. Use of multiple zones could increase function runtime variation.

Explore the GUI to write down subnet IDs and security group IDs for `config.json`:

VPC

ⓘ When you connect a function to a VPC in your account, it does not have access to the Internet unless your VPC provides access. To give your function access to the Internet, route outbound traffic to a NAT gateway in a public subnet. [Learn more](#)

VPC Info
Choose a VPC for your function to access.

vpc-b1a8e1d6 (172.31.0.0/16)

Allow IPv6 traffic for dual-stack subnets
You can allow outbound IPv6 traffic to subnets that have both IPv4 and IPv6 CIDR blocks.

Subnets
Select the VPC subnets for Lambda to use to set up your VPC configuration.

Choose subnets

subnet-f261d8bb (172.31.32.0/20) us-west-2a X

⚠ We recommend that you choose at least 2 subnets for Lambda to run your functions in high availability mode.

Once you've configured `config.json` it's very easy to recompile and deploy your Lambda function using the command line. Simply run the script `publish.sh` with the arguments below. The last argument is the desired function memory size. The 0s are for deployment to other FaaS platforms which are not used in this tutorial: Google Cloud Functions, IBM Cloud Functions, and Azure Functions.

```
# Deploy to AWS with 3GBs:
./publish.sh 1 0 0 0 3008
```

Additional documentation on the deploy tool can be found here:
https://github.com/wlloyduw/SAAF/tree/master/java_template/deploy

8. TO DO: Parallel Client Testing of AWS Lambda

SAAF provides the "FaaS Runner" Python-based client tool for orchestrating multi-threaded concurrent client tests against FaaS function end points. "FaaS Runner" allows the function end points to be defined in JSON objects, and for the repeatable experiments to be defined as JSON objects.

Before starting, install dependencies for FaaS Runner:

```
sudo apt install python3 python3-pip
sudo apt install python3-requests python3-boto3 python3-botocore
#OLD: pip3 install requests boto3 botocore
```

Verify after installation that your AWS CLI is still working:

```
aws --version
```

If the CLI, stops working, reinstall the following:

```
pip3 install awscli==1.19.53
pip3 install botocore==1.22.0
```

For detailed instructions on the FaaS Runner, please refer to the GitHub repository mark down documentation page:

FaaS Runner Documentation:

<https://github.com/wlloydw/SAAF/tree/master/test>

There also exists a Bash script for performing multi-threaded concurrent tests that is available on request from the instructor called 'partest'.

To tryout the FaaS Runner, navigate to the "test" directory:

```
cd {base directory where project was cloned}/SAAF/test
```

First, create a function JSON file under the SAAF/test/functions directory that describes your AWS Lambda function.

```
cd functions
cp exampleFunction.json hello.json
```

Edit the file hello.json function file to specifically describe your Lambda function:

```
{
  "function": "hello",
  "platform": "AWS Lambda",
  "source": "../java_template",
  "endpoint": ""
}
```

Function is the name of your AWS Lambda function.

Platform describes the FaaS platform where the function is deployed.

Source points to the source directory tree of the function.

Endpoint is used to specify a API Gateway URL.

If endpoint (URL) is left blank, the function can be invoked if the callWithCLI is set to true in the experiment file described below.

Next, create an experiment JSON file to describe your experiment again using the example template provided:

```
$ cd ..
$ cd experiments/
$ cp exampleExperiment.json hello.json
```

Next edit the hello.json experiment file to specifically describe your desired experiment using the hello function:

```
{
  "callWithCLI": true,
  "memorySettings": [0],
  "payloads": [
    { "name": "Bob" },
    { "name": "Joe" },
    { "name": "Steve" }
  ],
  "runs": 50,
  "threads": 50,
  "iterations": 3,
  "sleepTime": 5,
  "randomSeed": 42,

  "outputGroups": ["uuid", "cpuType", "vmuptime", "newcontainer", "endpoint", "containerID", "vmID",
  "zAll", "zTenancy[vmID]", "zTenancy[vmID[iteration]]"],
```

```

"outputRawOfGroup": ["zTenancy[vmID[iteration]]", "zTenancy[vmID]", "cpuType"],
"showAsList": ["vmuptime", "cpuType", "endpoint", "containerID", "vmID", "vmID[iteration]"],
"showAsSum": ["newcontainer"],
"ignoreFromAll": ["zAll", "lang", "version", "linuxVersion", "platform", "hostname"],
"ignoreFromGroups": ["1_run_id", "2_thread_id", "cpuModel", "cpuIdle", "cpuIowait", "cpuIrq",
"cpuKrn", "cpuNice", "cpuSoftIrq", "cpuUsr", "finalCalc"],
"ignoreByGroup": {
"containerID": ["containerID"],
"cpuType": ["cpuType"],
"vmID": ["vmID"],
"zTenancy[vmID]": ["cpuType"],
"zTenancy[vmID[iteration]]": ["cpuType"]
},

"invalidators": {},
"removeDuplicateContainers": false,
"openCSV": true,
"combineSheets": false,
"warmupBuffer": 1
}

```

A detailed description of experiment configuration parameters is included on the GitHub page. Please modify the following:

Runs: This is the total number of function calls. **Keep this set to 50.** (this is the default concurrency limit)

Threads: This is the total number of threads used to invoke the **Runs**. **Keep this set this to 50.** Keeping a 1 : 1 ratio between runs and threads ensures that each run will be performed by the client in parallel using a dedicated thread.

Iterations: This is number of times the experiment will be repeated. **Set this to 1.**

openCSV: If your platform has a spreadsheet application that will automatically open CSV files, then specify true, otherwise specify false. (Linux or MAC only)

CombineSheets: When set to true, this will combine multiple **iterations** into one spreadsheet. Since we are only performing 1 iteration, set this to **false**.

To obtain 50 distinct execution environments on AWS Lambda (think sandboxes), on remote network connections It is necessary to add a sleep call in the function so that the client computer can concurrently invoke 50 functions to run in parallel. Without adding a sleep function, AWS Lambda is so fast that many of the functions will complete preventing the client computer from successfully invoking 50 functions in parallel. Instead some function environments (instances) will be reused. When function executions do not overlap in time existing function environments (think sandboxes) will be reused resulting in (**newcontainer=0**). When all function invocations overlap this forces AWS Lambda to create and run 50 distinct sandboxes at the same time. This create resource contention in the public cloud because the function instances will compete for resources across a set of cloud servers. Given that HelloWorld is not a computationally complex function, overlapping calls requires adding the sleep statement to extend the duration of execution to force an overlap. These functions, however, don't compete for resources because they just "sleep".

Try adding a sleep statement to force the cloud provider to create 50 distinct execution environments (i.e. sandboxes) for running your HelloWorld function at the same time. Success will be indicated by obtaining 50 functions with newcontainer=1. After sandboxes are created, they are reused on subsequent calls, so they report newcontainer=0. Function instances (e.g. sandboxes) are deprovisioned randomly by AWS Lambda starting approximately 5 minutes after the last function call. In experiments, previously deprovisioning 100 sandboxes has been shown to take from 10 to 40 minutes as the sandboxes are slowly retired by AWS.

Add a sleep function to overlap the execution of your functions on AWS Lambda to obtain 50 new containers:

```
// Sleep for 10 seconds
try
{
    Thread.sleep(10000);
}
catch (InterruptedException ie)
{
    System.out.println("Interruption occurred while sleeping...");
}
```

Now try the FaaS Runner python tool.

Before trying the tool, be sure to close any spreadsheets that may be open in Microsoft Excel or Open/LibreOffice Calc from previous SAAF experiment runs.

```
# navigate back to the test directory
cd {base directory where project was cloned}/SAAF/test

# Requires python3
python3 faas_runner.py -f functions/hello.json -e experiments/hello.json
```

Note that faas_runner will automatically use the default AWS region that has been configured using the ‘aws configure’ command. If for some reason you need to change your default region, please rerun ‘aws configure’.

If your platform has a spreadsheet or tool configured to automatically open CSV files, then the CSV file will automatically open once it is created if openCSV is “true” in the experiments/hello.json file. When opening the file, be sure that only the comma (“,”) is used as a field/column delimiter. In Ubuntu 24.04, libreoffice can be used. If you have not installed libreoffice, it can be installed as follows:

```
sudo apt update
sudo apt install libreoffice
```

CSV report files produced by FaaS Runner are saved under the “history” directory. They can be opened using a spreadsheet after FaaSRunner has completed and exited.

Increasing AWS Lambda Concurrency

By default, AWS Lambda now only allow 10 concurrent functions. This means that that maximum number of unique function instances indicated by “newcontainer=1” in the output, will be 10. **It is recommended that you ** increase your AWS Lambda account concurrency to at least 100 **, if not more, and then repeat the experiment before submitting results.**

In the AWS Management Console, in the search bar, search for “Service Quotas”. Then in the Service Quota user interface, on the right, select/search for “AWS Lambda”. The press the “View Quotas” button. Click on the “Concurrent executions” parameter. Click on the button in the upper-right corner of the screen that says “Request increase at account level”. For the “Increase quota value” specify 100. Optionally, you can specify more than 100, but not more than 1,000. It specifying more than 100, it is possible that the request could take

24 hours, and/or be rejected. If your request to increase the account quota is rejected, it is suggested that you inform the professor using discord, email, or mention this in the class.

>>> FOR SUBMISSION <<<

After increasing the concurrent executions to 100 or more, repeat the experiment, and then examine the CSV report output using a spreadsheet application to determine the following.

Capture answers to the questions below in a PDF file and upload the PDF file to Canvas.

Include your Name, Function Name, AWS Region, VPC (+ Availability Zone), or no VPC

0. Did you add `Thread.sleep(10000)` ? Yes / No
1. Report the total number of “Successful Runs”
2. Report the total number of unique container IDs
3. Report the total number of unique VM IDs
4. Report the number of runs with `newcontainer=0` (these are recycled runtime environments)
5. Report the number of runs with `newcontainer=1` (these are newly created runtime environments)
6. The `zAll` row aggregates performance results for all tests. Looking at this row, report the:
 - **avg_runtime** for your function calls (measured on the server side in milliseconds)
 - **avg_roundTripTime** for your function calls (measured from the client side in milliseconds)
 - **avg_cpuidleDelta** for your function calls (units are in centiseconds)
 - **avg_latency** for your function calls (in milliseconds)

The difference between the `avg_roundTripTime` and the `avg_runtime` should be the `avg_latency`.

`cpuidle` time is measured in centiseconds. Multiply this by 10 to obtain milliseconds.

Linux CPU time accounting is provided in SAAF to report the state of the processor when executing Lambda functions. The wall clock (or watch time) can be derived by adding up the available CPU metric deltas and dividing by the number of CPU cores (2 for AWS Lambda @ 3GB RAM) to obtain an estimate of the wall clock time (function runtime).

Once adding “`Thread.sleep(10000)`” to your hello function check the delta value for CPU IDLE time. By including `Thread.sleep(10000)` this value should be close to 10,000. Sleep essentially makes the CPU idle for most of the duration of the function’s runtime.

Difference Between AWS Lambda VPC and NO VPC function deployments:

Lambda functions that run in a Virtual Private Cloud (VPC) suffer from additional cold start overhead because when function instances are first called, there is a higher initialization cost to setup the VPC network connection for the function compared to standard non-VPC Lambda functions.

AWS Lambda Abstracts CPU type through the use of Micro VMs:

When executing Lambda functions, for each concurrent client request arriving at the same time, Lambda creates distinct virtual infrastructure known as “function instances”. One function instance is created for each user request received in parallel (at the same time). Function instances can be reused on subsequent calls. After 5 minutes of inactivity, function instances are gradually deleted. When function instances are continually used, they can stay alive for up to ~4 hours. Amazon will automatically replace function instances to continually refresh infrastructure at random. Function instances are implemented using micro-VMs which are a lighter-weight form of a full virtual machine. AWS has created the “Firecracker” MicroVM specifically for

serverless (FaaS and CaaS) workloads. MicroVMs provide better isolation from a resource accounting point of view. Using these micro-VMs, however, has led to further abstraction of the underlying hardware. For example the CPU type on Firecracker VMs are simply identified as: **Intel(R) Xeon(R) Processor @ 2.50GHz**. No model number is specified. This may be a virtual CPU designation provided by Firecracker which is based on KVM.

To read more about the Firecracker MicroVM, see:

<https://firecracker-microvm.github.io/>

The FaaS Runner will store experiment results as CSV files under the history directory.

On some platforms, these filenames may automatically increment so they don't overwrite each other. On other platforms, it may be necessary to make a copy to preserve the files between runs.

Here is an example of making a copy:

```
cd history
cp "hello - hello - 0 - 0.csv" tcss462-562_ex1.csv
```

9. TO DO: Two-Function Serverless Application: Caesar Cipher

To complete tutorial #4, use the resources provided to construct a two-function serverless application that implements a Caesar Cipher. The Caesar cipher shifts an ASCII string forward to encode the message, and shifts the string backwards to decode.

LLM USE IS OKAY: For writing the Caesar cipher code, it is fine to use ChatGPT, GitHub copilot, or another LLM. Google search can also provide code for a Java implementation.

To get started, create a new directory under /home/ubuntu

Then clone the SAAF repository twice to have two separate empty Lambdas.

Alternatively, a single project can be used where there are separate encode and decode class files. The function handler can be adjusted to point to the specific class and/or method that serves as the Lambda function entry point to your Java code.

```
$ cd ~
:~$ mkdir tcss562
:~$ cd tcss562
:~/tcss562$ mkdir encode
:~/tcss562$ mkdir decode
:~/tcss562$ cd encode
:~/tcss562/encode$ git clone https://github.com/wlloyduw/SAAF.git
Cloning into 'SAAF'.....
:~/tcss562/encode$ cd ..
:~/tcss562$ cd decode
:~/tcss562/decode$ git clone https://github.com/wlloyduw/SAAF.git
Cloning into 'SAAF'.....
```

Next, implement two lambda functions.

One called “Encode”, and another “Decode” that implement the simple Caesar cipher.

In the SAAF template, the verbosity level can be adjusted to provide less output.

To explore verbosity levels offered by SAAF, try adjusting the number of metrics that are returned by replacing the line of code:

```
inspector.inspectAll();
```

with one of the following or simply remove inspectAll() altogether:

inspectCPU()	reports all CPU metrics
inspectContainer()	reports all Container-level metrics (e.g. metrics from the runtime environment)
inspectLinux()	reports the version of the Linux kernel hosting the function.
InspectMemory()	reports memory metrics.
InspectPlatform()	reports platform metrics.

At the bottom, the following line of code can be commented out or replaced:

```
inspector.inspectAllDeltas();
```

Less verbose options include:

inspectCPUDelta()	Reports only CPU metric changes
inspectMemoryDelta()	Reports only memory metric utilization changes

Detailed information about metrics collection by SAAF is described here:

https://github.com/wlloydw/SAAF/tree/master/java_template

For the Caesar Cipher, pass a message as JSON to your “encode” function as follows:

```
{
  "msg": "ServerlessComputingWithFaaS",
  "shift": 22
}
```

The encode function should shift the letters of an ASCII string forward to disguise the contents as shown in the example JSON below (SAAF metrics mostly removed):

```
{
  "msg": "OanranhaooykilqpejcsSepdBwwO",
  "uuid": "036c9df1-4a1d-4993-bb69-f9fd0ab29816",
  "vmuptime": 1539943078,
  "newcontainer": 0
  . . . output from SAAF truncated for brevity..
}
```

The second service, decrypt, should shift the letters back to decode the contents as shown in the JSON output:

```
{
  "msg": "ServerlessComputingWithFaaS",
  "uuid": "f047b513-e611-4cac-8370-713fb2771db4",
  "vmuptime": 1539943078,
```



```
"newcontainer": 0
. . . output from SAAF truncated for brevity...
}
```

Notice that the two services have different uids (container IDs) but the same vmuptime (VM/host ID). On AWS Lambda + VPC this behavior could occur if two functions share the same VMs. Note: *This behavior is no longer observable as AWS Lambda now uses the Firecracker MicroVM for hosting function which abstracts this information about shared hosts from users.*

Both services should accept two inputs:

integer	shift	number of characters to shift
String	msg	ASCII text message

The Internet has many examples of implementing the Caesar cipher in Java. Use of ChatGPT, etc. is also permitted.

<https://stackoverflow.com/questions/21412148/simple-caesar-cipher-in-java>

You'll notice that SAAF provides a lot of attributes in the JSON output. This verbosity may be optionally reduced to simplify the output. Instead of calling **inspectAll()** the code can be reworked to call a few functions that will then only provide a subset of the information. For example, this set would offer fewer attributes while retaining some helpful metrics:

```
inspector.inspectCPUDelta();
inspector.inspectContainer();
inspector.inspectPlatform();
```

Once implementing and deploying the two-function Caesar cipher Lambda application, modify the call_service.sh script and create a "cipher_client.sh" BASH script to serve as the client to test your two-function app.

Cipher_client.sh should create a JSON object to pass to the encode service. The output should be captured, parsed with jq, and sent to the decode service.

The result should be a simple pair of services for applying and removing the cipher. The Cipher_client.sh bash script acts as the client program that instruments the flow control of the two-function cipher application. Deploy all functions to operate synchronously just like the hello example service. Host functions in your account to support testing.

Use API gateway endpoints and curl to implement Cipher_client.sh. Do not use the AWS CLI to invoke Lambda functions. This will allow your two-function application to be tested using the Cipher_client.sh script that is submitted on Canvas.

IMPORTANT: DO NOT DELETE THE API GATEWAY ENDPOINTS from your account *until you receive a grade* on Tutorial 4. There is no cost for having an API gateway defined in your account. Charges are only incurred when the API Gateway endpoint is used. The API Gateway endpoint is required to test your submissions. After you receive a grade, it is okay to delete the API Gateway endpoints and Lambda functions. It is suggested to leave these items in the account for several weeks to allow time for grading.

SUBMISSION

Tutorial #4 should be completely **individually**. Everyone should have the experience of creating and working with serverless functions on AWS Lambda. Files will be submitted online using Canvas.

When possible please create and submit a Linux tar.gz file to capture all of your project's source files. From the command line, navigate to the SAAF directory for your encode/decode project. You may combine functions into a single project (*by modifying the function handler when deploying the Lambdas- recommended*) or submit separate tar.gz files for a separate encode and decode project to Canvas.

To create the tar.gz archive file, from the SAAF directory, use the command:

```
tar czf encode.tar.gz .
```

Once having the archive, the contents can be inspected as follows:

```
tar ztf encode.tar.gz | less
```

Use the 'f' key to go forward, 'b' key to go backward, and 'q' key to quit

For the submission, submit a working bash client script (**Cipher_client.sh**) that invokes both functions.

Be sure to include in the Canvas submission the tar.gz file that includes all source code for your Lambda functions. Alternatively a zip file can be submitted.

In addition, include a PDF file including answers to questions for #8.

Scoring

- | | |
|-----------|--|
| 20 points | Providing a PDF file answering questions using output from the FaaS Runner for #8. |
| 20 points | Providing a working Cipher_client.sh that instruments the two-functions Lambda app using REST URLs from the API Gateway. |

Change History

Version	Date	Change
0.1	10/19/2024	Original Version
0.11	10/28/2024	Added clarification on cli-binary-format binary parameter for AWS CLI.