## Inference Backend– Step-by-Step Implementation

- 1. We will start by uploading our model artifact to S3.
- Log in to your AWS account and navigate to the S3 homepage.
- On the top left corner, click + Create bucket .
- On the window that opens:
  - Bucket name: you name your bucket whatever you wish, but if you include "sagemaker" in the name it will save you a configuration step later on.
  - Region: Choose the region in which you wish to host your model.
  - Leave the rest as is.
  - On the bottom left corner, click Create

Name and region			
Name and region			
Bucket name ()			
sagemaker-name-you-want			
Region			
US West (Oregon)			~
Copy settings from an existing	ng bucket		
Copy settings from an existi	ng bucket		
Copy settings from an existing	ng bucket		Ŷ
Copy settings from an existing Select bucket (optional) 12 Bucket	ng bucket kets		v
Copy settings from an existiin	ng bucket kets		~
Copy settings from an existi Select bucket (optional)12 Buc	ng bucket	-	~
Copy settings from an existi Select bucket (optione)12 Buc	ng bucket		
Copy settings from an existi	ng bucket		

- Click on the name of the bucket you just created.
- On the top left corner, click + Create folder .
- Name the folder "Model" and hit Save
- Click on the folder that was just created.
- From the window that appeared, click Add files provided in the .zip file.

, and select the model.tar.gz

• Lastly make note of the Target path display at the top of the window.



- On the bottom right corner, click Upload
- **2.**Now that our model artifact has been uploaded to S3, we will turn it into an actual model.
- First open your AWS and navigate to the Amazon SageMaker console.
- On the left side of the page, select the Models link under the Inference tab.
- In the Model settings cell, configure:
  - Model name: Choose any name you wish, but make note of it as we will need it later.
  - IAM role: From the drop down menu select Create a new role.
    - This option will create a role with permissions described by the AmazonSageMakerFullAccess IAM policy.

• If the bucket in which you hosted the model artifact contains "sagemaker" in its name, the just click . Otherwise make sure Specific S3 buckets is selected from the radio menu, and enter your bucket's name in the text field.

Passing an IAM role gives Amazon SageMaker permission to perform actions in other AWS sen grant permissions described by the AmazonSageMakerFullAccess [2] IAM policy to the role y The IAM role you create will provide access to:	vices on your behalf. Creating a role here will rou create.
<ul> <li>S3 buckets you specify - optional</li> <li>Specific S3 buckets</li> </ul>	
your-buckets-name	
Comma delimited. ARNs, "*" and "/" are not supported.	
<ul> <li>Any S3 bucket Allow users that have access to your notebook instance access to any bucket and its of</li> </ul>	
	contents in your account.
○ None	contents in your account.
<ul> <li>None</li> <li>Any S3 bucket with "sagemaker" in the name</li> </ul>	contents in your account.
<ul> <li>None</li> <li>Any S3 bucket with "sagemaker" in the name</li> <li>Any S3 object with "sagemaker" in the name</li> </ul>	contents in your account.
<ul> <li>None</li> <li>Any S3 bucket with "sagemaker" in the name</li> <li>Any S3 object with "sagemaker" in the name</li> <li>Any S3 object with the tag "sagemaker" and value "true"</li> </ul>	contents in your account. See Object tagging [

- Click Create role .
- Now configure the Container Definition cell:
  - For the container input options select the first option on the list, Provide model artifacts and inference image location.
  - Now under the Provide model artifacts and inference image tab:
    - Location of inference code image: 763104351884.dkr.ecr.us-west-2.amazonaws.com/mxnet-inference:1.4.1-gpu-py2
    - Location of model artifacts: Enter the path to the model.tar.gz file you uploaded to your bucket. Make sure that the path ends with /model.tar.gz
- Leave everything else as is, and click Create model

3. Next we will host the model we created using Amazon SageMaker.

- On the left hand side of the page, select the Endpoints link under the Inference tab.
- On the top right corner click Create endpoint
- Configure the Cells as follows:
  - For the endpoint name, you may choose any name for the endpoint, but make sure you save it somewhere as we will need this name later to set up our lambda function.

Endpoint	
Endpoint name	
our application uses this name to access this endpoint.	
Pick a name for your endpoint (e.g PuenteEndpoint)	
Naximum of 63 alphanumeric characters. Can include hyphens (-), but not spaces. Must be unique w rour account in an AWS Region.	ithin

• In the next cell select Create a new endpoint configuration.



- For the new endpoint configuration:
  - Endpoint configuration name: Choose any name you wish.

• Encryption key: Leave this on the default option of "No Custom Encryption".

models to a	nodels to Amaz deploy, and the	on SageMaker, relative traffic	, first create an weighting and	endpoint configu hardware require	ration. In the confi ments for each.	iguration, specify	which
Endpoint co	onfiguration na	me					
Choose a	name						
Maximum of Must be uniq	63 alphanumeric ue within your ac	characters. Can in count in an AWS F	nclude hyphens (-) Region.	, but not spaces.			
Encryption Encrypt your	<b>key - <i>optional</i></b> data. Choose an e	existing KMS key o	or enter a key's AR	IN.			
No Custo	m Encryption			•			
Production Model name	rvariants Training job	Variant name	Instance type	Elastic Inference	Initial instance count	Initial weight	Actions
	-		There are cu	rrently no resourc	:es.		

- Now click Add Model.
- From the menu that popped up, select the model we created in the last step and hit **Save**.
- Once you have selected your model, under the actions column you can click edit and choose the instance type you want to host your model on. For the original PoC we used the ml.t2.medium instance.
- Click Create endpoint configuration .
- Lastly you can use the Tags section to tag your endpoint. This will make it easier to group together the resources that belong to this PoC so you can tear down the application or keep track of its costs later on.
- Click Create endpoint

4. Now we need to create a lambda function to invoke the endpoint we just created.

- Navigate to the AWS Lambda console.
- Click Create function .

Author from scratch	0	Use a blueprint	Browse serverless app repository	
Start with a simple Hello World example.		Build a Lambda application from sample code and configuration presets for common use cases.	Deploy a sample Lambda application from the AWS Serverless Application Repository.	

- From the radio menu at the top of the page, select Author from scratch.
- Configure:
  - Function name: Any name you wish.
  - Runtime: Python 2.7
  - Expand the choose or create permissions tab.
    - From the drop down menu, ensure that Create a new role with basic Lambda permissions is selected.
  - Click Create function .
  - Scroll down to the Function code cell.
    - Highlight everything in the editor and replace it with the code provided in the .zip .

e entry type		Runtime		Handler Info	
it code inline	•	Python 2.7	•	lambda_function.lambda_handler	
File Edit Find View Go	Tools Window				K 2
v 📄 delete-me-please 🔅 v	lambda_function × 🕀				
Iambda_function.py	1 import json				
	3 def lambda_handler(event	;, context):			
	5 return {				
	6 'statusCode': 20 7 'body': json.dur	№, ıps('Hello from Lambda!')			
	8 }				

- Scroll to the top of the page, click **Save**.
- **5.** Lastly we will use API Gateway to create an API to be called from code.
- Navigate to the Amazon API Gateway console
- Choose + Create API .
- On the page that comes up, configure:
  - Protocol: REST
  - Create new API: New API
  - API Name: choose a name
  - Description: Optional
  - Endpoint Type: Regional
- Click Create API .
- From the drop down menu labeled actions, select Create Method.

s Re	esources Actions	/ Methods	
Documentation	/ RESOURCE AC	IONS	
documentation2	Create N	lethod	
1 -	Create R	esource	
Resources	Enable C	UK5	
Stages	EurcRes	dice Documentation	
Authorizers	APLACTIONS		
Gateway Beenonses	Deploy A	PI	
datonay nooponooo	Import A	PI	
Models	Delete A	PI	
Resource Policy			
Documentation			
Settinge			
Ootunga			
Puente			
age Plans			
Keys			
stom Domain Names			
ent Certificates			
C Links			
tings			
11190			

• A drop down menu will appear, from it select Post then click the checkmark button.



- The section to the right of Post will be populated with some options to configure:
  - Integration type: Lambda Function
  - Use Lambda Proxy Integration: Check
  - Lambda Region: Choose the region which you created the lambda function in.
  - Lambda Function: Start to type out the name of lambda function and a drop down menu should appear. Select the lambda function you created in the last section.
- Click Save
- A window will appear stating that "You are about to give API Gateway permission to invoke your lambda function".
- Choose okay.
- From the actions menu, choose deploy API.
- On the window that comes up, configure:
  - For deployment stage: choose "[New Stage]"
  - For Stage Name: choose a name (e.g. Dev)
  - For Stage description: Optional

- For Deployment description: Optional
- Click Deploy .
- Once the deployment is done, you can navigate to the dashboard link on the left hand side menu under the name of your API. This is the URL you will use to contact the endpoint.