



DISTRIBUTED TRAINING ON EDGE

ENGR 295B, Spring 2020

Ahmed Hambaba, Farshid Marbouti

Shafie Mukhre, Pavani Somarouthu, Brian Hu

Table of Contents



- Problem Statements
- Motivations
- Technical Objective

Centralized Training

- Architecture
- Centralized training
- Model training
- Training on AWS
- API endpoint
- Testing endpoint
- Tensorflow training on AWS

Distributed Training

- Distributed Learning
- Options
- Tensorflow Federated
- Architecture
- FL Simulation Implementation
- Obstacles
- Discussion

PROBLEM STATEMENTS



1. **Privacy:** ML and DL utilizes lots user's data for model training. However, these data most often are personal and sensitive data such as images, health informations and so on.
2. **Decentralized data:** Most cosumer's data is decentralized in every edge devices, while the existing centralized training is necessary for complex model training, there are cases that this overkill and unnecessary.

MOTIVATIONS



1. To preserve the privacy of user data, and therefore increase the usage or implementation of ML and DL applications.

TECHNICAL OBJECTIVES



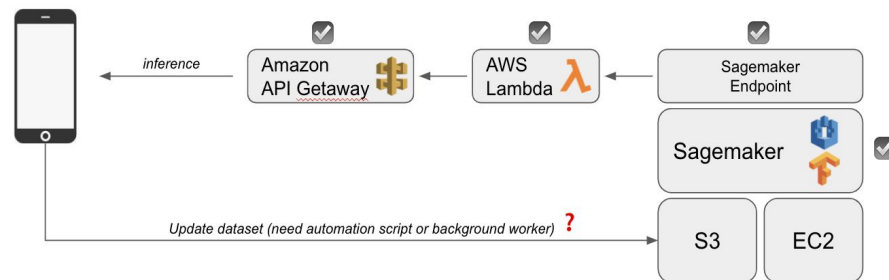
- The overall objective of the project is to create 2 ML Architectures:
 - a. The first architecture is one based on Centralized Cloud Training
 - b. The second is one that distributes the learning across IoT Devices

The rest of this presentation will go into each of these approaches in vast detail.

CENTRALIZED TRAINING

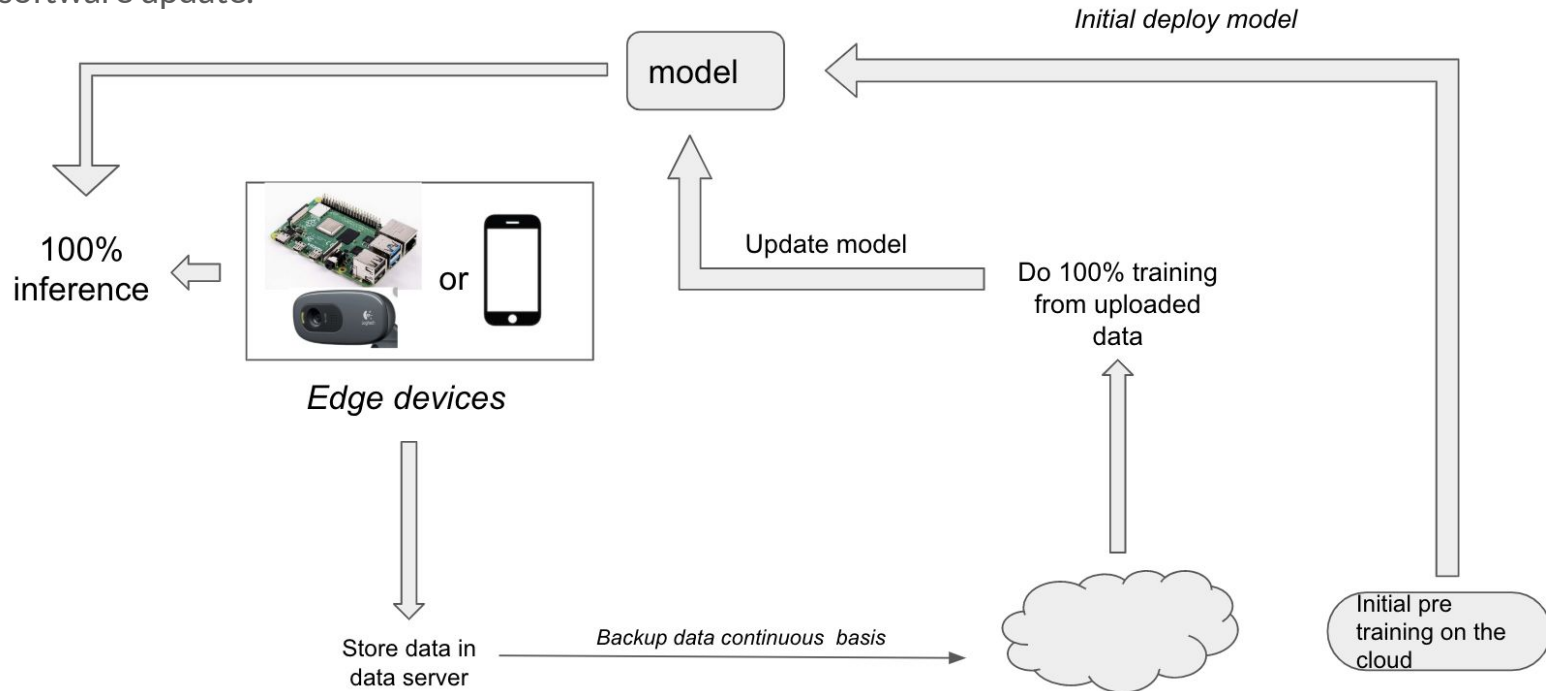
Centralized training architecture

- Centralized Training is the approach that most modern ML architectures and IoT applications use.
- The model is hosted and trained on a cloud.
- The application downloads a frozen version of the model and uses it to make inferences. The application can then send data to the cloud and download an updated retrained model.



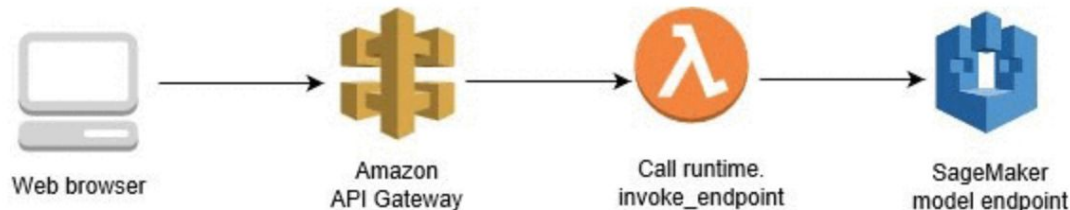
Model training

ML training is done exclusively on the cloud and the model is updated on regular basis via over-the-air software update.



CENTRALIZED TRAINING ON AWS

- Amazon Sagemaker: For access to jupyter notebook instances to build and train models.
- Amazon EC2 (Amazon Elastic Compute Cloud): Model Training is using accelerated computing provided by EC2
- Amazon S3 (Amazon Simple Storage Service): All files and datasets are download or uploaded to S3 storage
- AWS Lambda: To create a lambda function to handle the REST API request in JSON format
- Amazon API Gateway: As an intermediate communication between edge devices client and web server itself



Edge application

ie: web app



Take an image

Convert image to
base64

API request

Plot detection on
image

*API
response*

Amazon
API Gateway



AWS
Lambda



base64

.jpg

Upload to S3
storage

Download to
/tmp storage

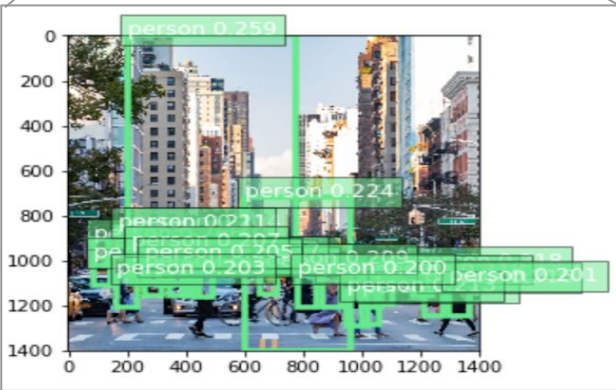
.jpg

bytearray

*Invoke
Model endpoint*

prediction

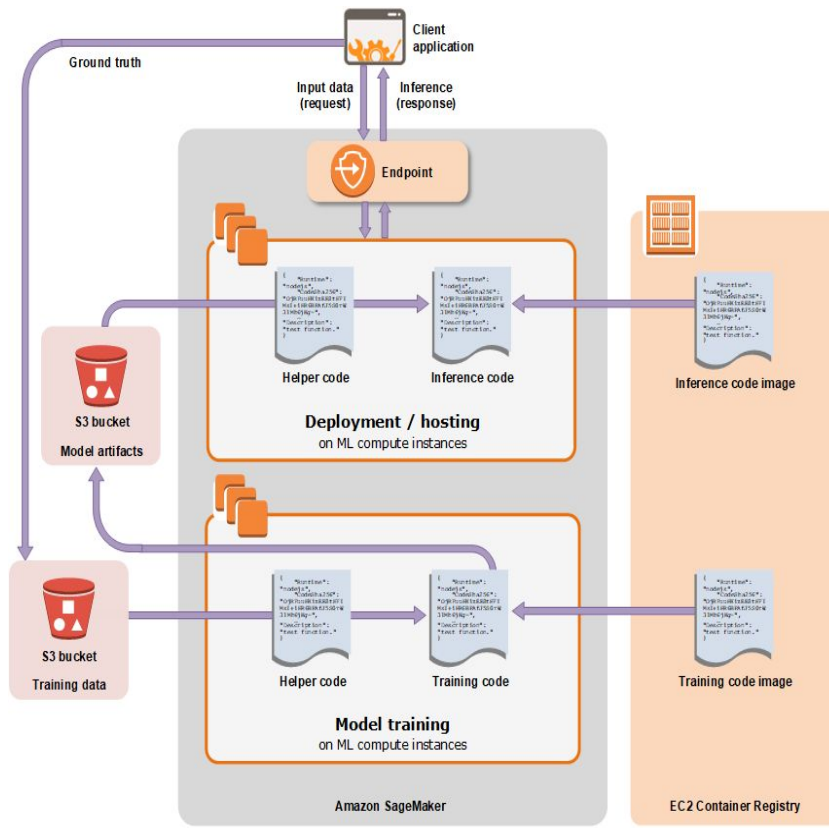
Edge device



Example of prediction

Tensorflow Training on AWS

- Sagemaker's own Object Detection Algorithm outputs image data in JSON format.
- It is difficult to deploy that output into applications so we had to do the exercise all over again in Tensorflow.
- Sagemaker has its own interface to train and deploy user-provided tensorflow code.
- Training code was taken from Tensorflow Object Detection API and adapted to Sagemaker.
- Sagemaker outputs a frozen pb model which can then be deployed to make predictions.

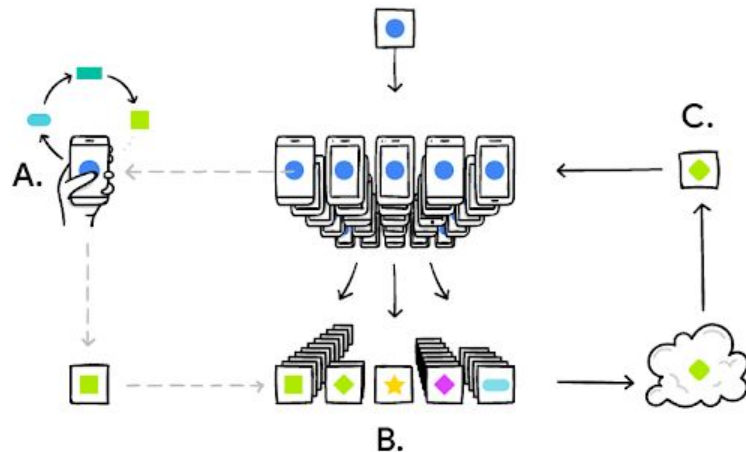


DISTRIBUTED TRAINING

Distributed training options

Distributed Training is training a ML or DL model of decentralized data on multiple nodes.

- Compared distributed training technique:
 - Large Batch Training
 - Federated Learning
- Compared multiple options to do federated learning
 - Tensorflow Federated
 - Tensorflow.js
 - Trains on AWS



TENSORFLOW FEDERATED



TensorFlow Federated (TFF) is an open-source framework for machine learning and other computations on decentralized data.

TFF's interface is organized in two layers

- Federated Learning (FL) API - High level interface
 - Models
 - Federated Computation Builders
 - Datasets
- Federated Code(FC) API - Low level interface

Tensorflow Federated Architecture

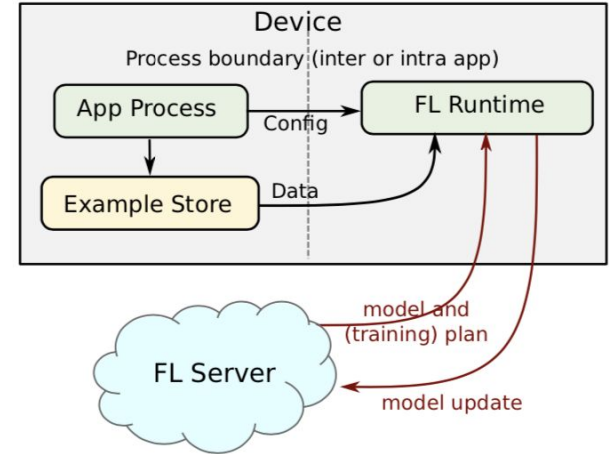
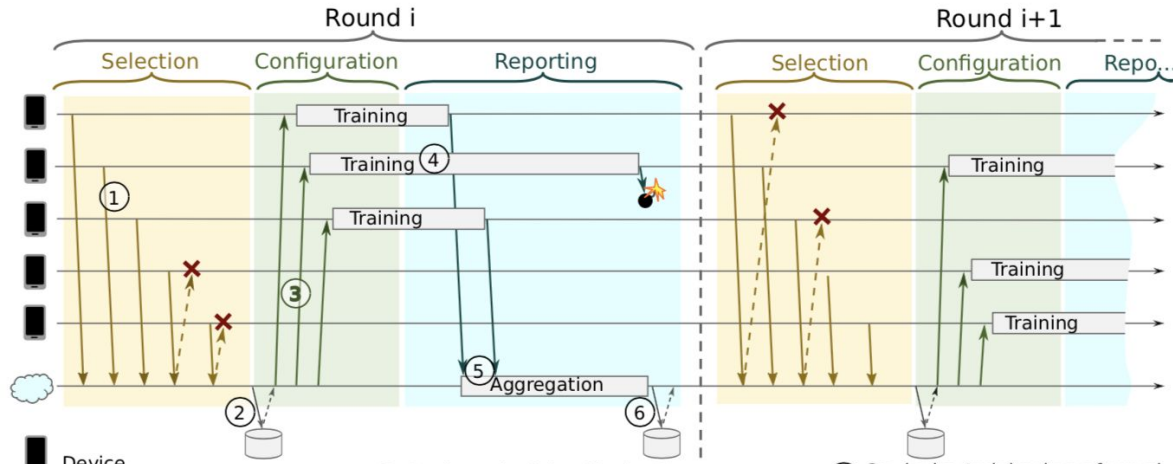
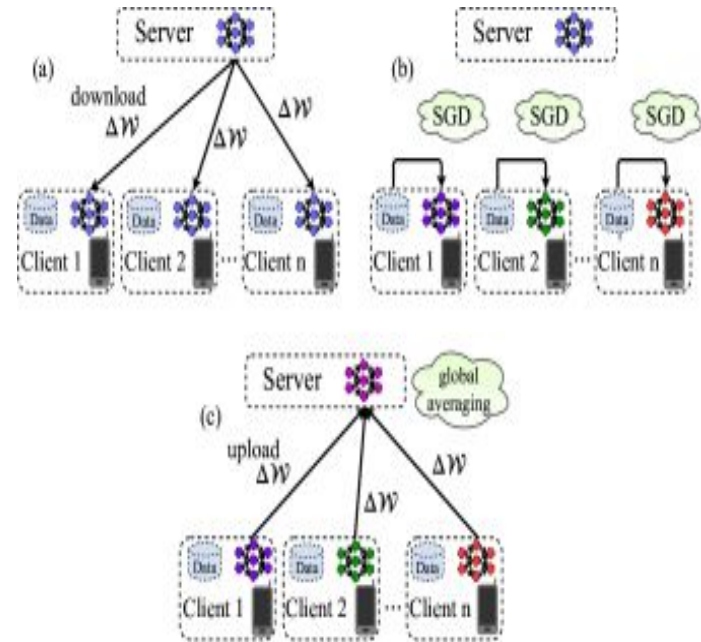


Figure 2: Device Architecture

Federated Learning simulation implementation

Federated Learning Architecture

- a) Clients synchronize with the server.
- b) Clients compute a weight update independently based on their local data.
- c) Clients upload their local weight updates to the server, where they are averaged to produce the new master model.



Federated Averaging Algorithm



Federated Averaging algorithm, combines local stochastic gradient descent (SGD) on each client with a server that performs model averaging.

Client (On Device) training

$$w^1 = \text{model}(x, y, b, e, \eta)$$

Where, w^1 is new model trained on device

Server training

$$g = g + (n^k * w^{1k} / M)$$

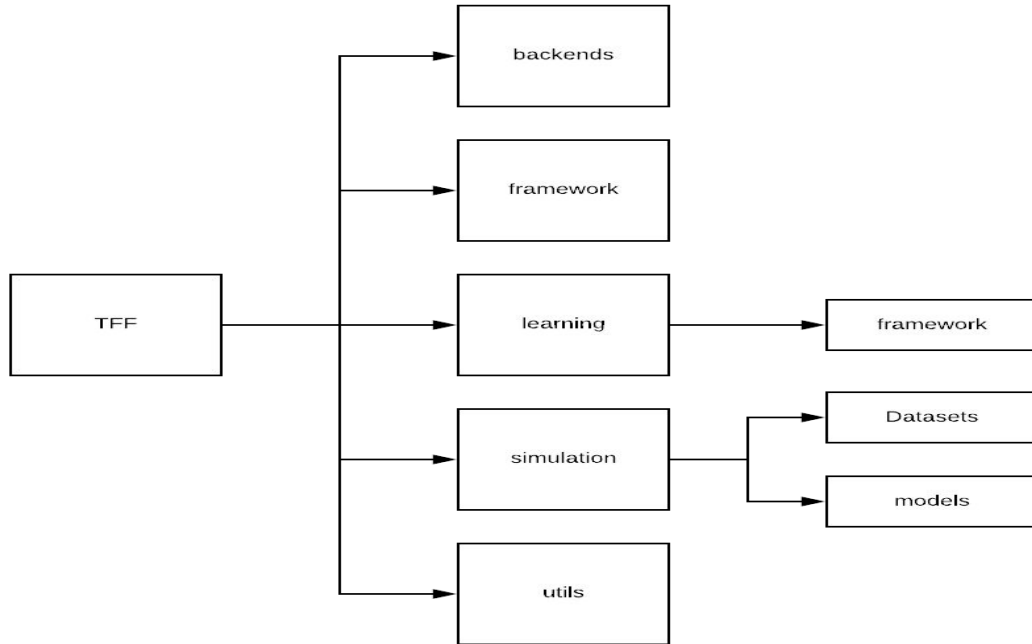
Where, g is the Stochastic Gradient Descent

n^k is the number of data points

w^{1k} is the weights gathered from device k

M is the sum of the number of data points

TFF Library

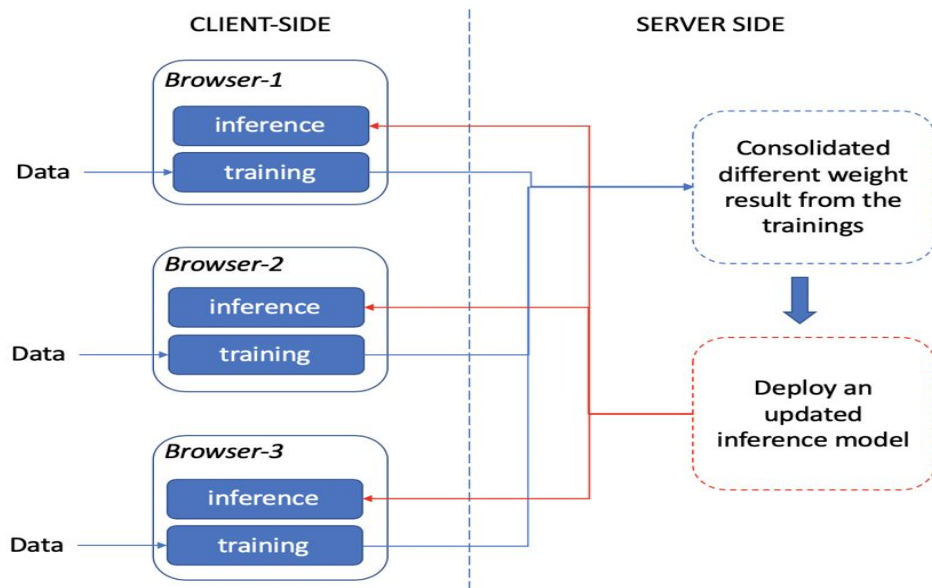


Obstacles we are facing with TFF



- Tensorflow federated is the only framework available for federated learning that have good documentations, but it is not for production.
- Tensorflow federated simulation is limited only to LEAF dataset which is designed or orchestrated for federated learning.

Tensorflow.js



DEMO

Colab interface showing the execution of a TensorFlow object detection model. The browser address bar displays the URL: `colab.research.google.com/github/AlaaSenjab/-Tutorial-Tensorflow_Object_Detection_API_On_Custom_Dataset/blob/master/weapon_detection_BLIpynb#scrollTo...`.

The file explorer on the left shows a folder named `sample_data`.

The main code editor displays the following log output:

```
2020-05-26 05:36:39.211145: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:983] successful NUMA node read from SysFS had negative value (-1), t
2020-05-26 05:36:39.211787: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1767] Adding visible gpu devices: 0
2020-05-26 05:36:39.212160: I tensorflow/core/platform/cpu_feature_guard.cc:142] Your CPU supports instructions that this TensorFlow binary was not com
2020-05-26 05:36:39.221382: I tensorflow/core/platform/profile_utils/cpu_utils.cc:94] CPU Frequency: 2000165000 Hz
2020-05-26 05:36:39.221569: I tensorflow/compiler/xla/service/service.cc:168] XLA service 0x1422f40 initialized for platform Host (this does not guar
2020-05-26 05:36:39.221597: I tensorflow/compiler/xla/service/service.cc:176] StreamExecutor device (0): Host, Default Version
2020-05-26 05:36:39.330167: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:983] successful NUMA node read from SysFS had negative value (-1), t
2020-05-26 05:36:39.330860: I tensorflow/compiler/xla/service/service.cc:168] XLA service 0x1423640 initialized for platform CUDA (this does not guar
2020-05-26 05:36:39.330892: I tensorflow/compiler/xla/service/service.cc:176] StreamExecutor device (0): Tesla T4, Compute Capability 7.5
2020-05-26 05:36:39.331077: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:983] successful NUMA node read from SysFS had negative value (-1), t
2020-05-26 05:36:39.331695: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1639] Found device 0 with properties:
name: Tesla T4 major: 7 minor: 5 memoryClockRate(GHz): 1.59
pciBusID: 0000:00:04.0
2020-05-26 05:36:39.331825: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcudart.so.10.1
2020-05-26 05:36:39.331873: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcublas.so.10
2020-05-26 05:36:39.331898: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcufft.so.10
2020-05-26 05:36:39.331924: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcuband.so.10
2020-05-26 05:36:39.331945: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcusolver.so.10
2020-05-26 05:36:39.331965: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcusparse.so.10
2020-05-26 05:36:39.331985: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcudnn.so.7
2020-05-26 05:36:39.332095: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:983] successful NUMA node read from SysFS had negative value (-1), t
2020-05-26 05:36:39.332658: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:983] successful NUMA node read from SysFS had negative value (-1), t
2020-05-26 05:36:39.333155: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1767] Adding visible gpu devices: 0
2020-05-26 05:36:39.333213: I tensorflow/stream_executor/platform/default/dso_loader.cc:44] Successfully opened dynamic library libcudart.so.10.1
2020-05-26 05:36:39.334410: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1180] Device interconnect StreamExecutor with strength 1 edge matrix:
2020-05-26 05:36:39.334436: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1186] 0
2020-05-26 05:36:39.334446: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1199] 0: N
2020-05-26 05:36:39.334558: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:983] successful NUMA node read from SysFS had negative value (-1), t
2020-05-26 05:36:39.335123: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:983] successful NUMA node read from SysFS had negative value (-1), t
2020-05-26 05:36:39.335615: W tensorflow/core/common_runtime/gpu/gpu_bfc_allocator.cc:39] Overriding allow_growth setting because the TF_FORCE_GPU_ALLOC
2020-05-26 05:36:39.335655: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1325] Created TensorFlow device (/job:localhost/replica:0/task:0/device:
INFO:tensorflow:Restoring parameters from training/model.ckpt-1140
```

The code editor also shows the following Python code snippet:

```
[ ] #downloads the frozen model that is needed for inference
files.download(output_directory + '/frozen_inference_graph.pb')
```

The bottom status bar indicates that 36.98 GB is available on the disk.

Conclusion



1. We have implemented Centralized Cloud Training architecture using Tensorflow framework and MS COCO dataset on AWS.
2. We have implemented a simulation of Distributed Training architecture using Federated Learning using Tensorflow Federated framework and LEAF's FEMNIST dataset on Google Colab notebook.
3. We have created federated learning environment using flask as the node server and tfjs to train the model.

Appendix

An example of Lambda function python script

```
lambda_function × (+)
1 import os
2 import io
3 import boto3
4 import json
5 import csv
6 import base64
7
8 s3 = boto3.client('s3')
9
10 ENDPOINT_NAME = os.environ['ENDPOINT_NAME']
11 runtime= boto3.client('runtime.sagemaker')
12
13 def lambda_handler(event, context):
14     print(event)
15     # if event['httpMethod'] == 'POST' :
16     data = json.loads(json.dumps(event))
17     name = data['name']
18     image = data['file']
19     image = image[image.find(",")+1:]
20     dec = base64.b64decode(image + "===")
21     s3.put_object(Bucket='objectdataset', Key=name, Body=dec)
22     s3.download_file('objectdataset', name, '/tmp/'+name)
23
24     file_name = '/tmp/' + name
25     with open(file_name, 'rb') as image:
26         f = image.read()
27         b = bytearray(f)
28     payload = b
29     response = runtime.invoke_endpoint(EndpointName=ENDPOINT_NAME,
30                                     ContentType='image/jpeg',
31                                     Body=payload)
32     print(response)
33     detections = json.loads(response['Body'].read().decode())
34     print(detections)
35     return detections
36
```

```
Execution Result × (+)
```

Testing API endpoint with postman

[illegible]

```
[{"prediction": [0.0, 0.29635363817214966, 0.7147926092147827, 0.761025071144104, 0.7601768970489502, 0.8612469434738159], [0.0, 0.2664850950241089, 0.060184113681316376, 0.6464931964874268, 0.09928541630506516, 0.7313581705093384], [0.0, 0.259016215801239, 0.1430124044418335, 0.0, 0.556414451408386, 0.7332626581192017], [0.0, 0.25492000579833984, 0.281993567943573, 0.7036433219909668, 0.31999415159225464, 0.788191556930542], [0.0, 0.24118885397910672, 0.2485395073890686, 0.7308893203735352, 0.2905610203742981, 0.8252425193786621], [0.0, 0.23789450526237488, 0.7131987810134888, 0.8191956877708435, 0.7555108070373535, 0.919554358690432], [0.0, 0.2293343949947357, 0.3080902397632599, 0.7273744940457751, 0.35195204615592957, 0.83082145454249939], [0.0, 0.2278672456741333, 0.21998275816440582, 0.7114567160606384, 0.257265083133773804, 0.7918701769214481], [0.0, 0.22411756217479706,
```