RAKOTONDRAMANANA A.A. LAFATRA

INFO 3 – UNIVERSITE DES MASCAREIGNES

# TP session 6: Integrating Machine Learning Model in a Flutter application

## Task-1 TensorFlow Lite Tutorial for Flutter: Image Classification

In this task, we will delve into the utilization of TensorFlow Lite within the Flutter framework. This involves training a machine learning model using Teachable Machine and seamlessly integrating the results into a Flutter mobile app.

Getting Started: To initiate this process, we will utilize a provided starter project as our foundation. While the project currently allows users to select or drag-and-drop images into the app, it lacks image recognition functionality. In the upcoming sections, we will employ TensorFlow Lite to address this limitation.

Building a Model with Teachable Machine: The starter project already incorporates a pre-trained model (model_unquant.tflite) and classification labels stored in the labels.txt file. However, our focus will be on understanding and engaging in the training process itself.

Preparing the Dataset - Training the Model: Training is the pivotal stage where the computer acquires knowledge from data and formulates rules. Our exploration will involve learning how to train the model using Teachable Machine.

1. To commence, visit https://teachablemachine.withgoogle.com and click on "Get Started" to access the training tool.
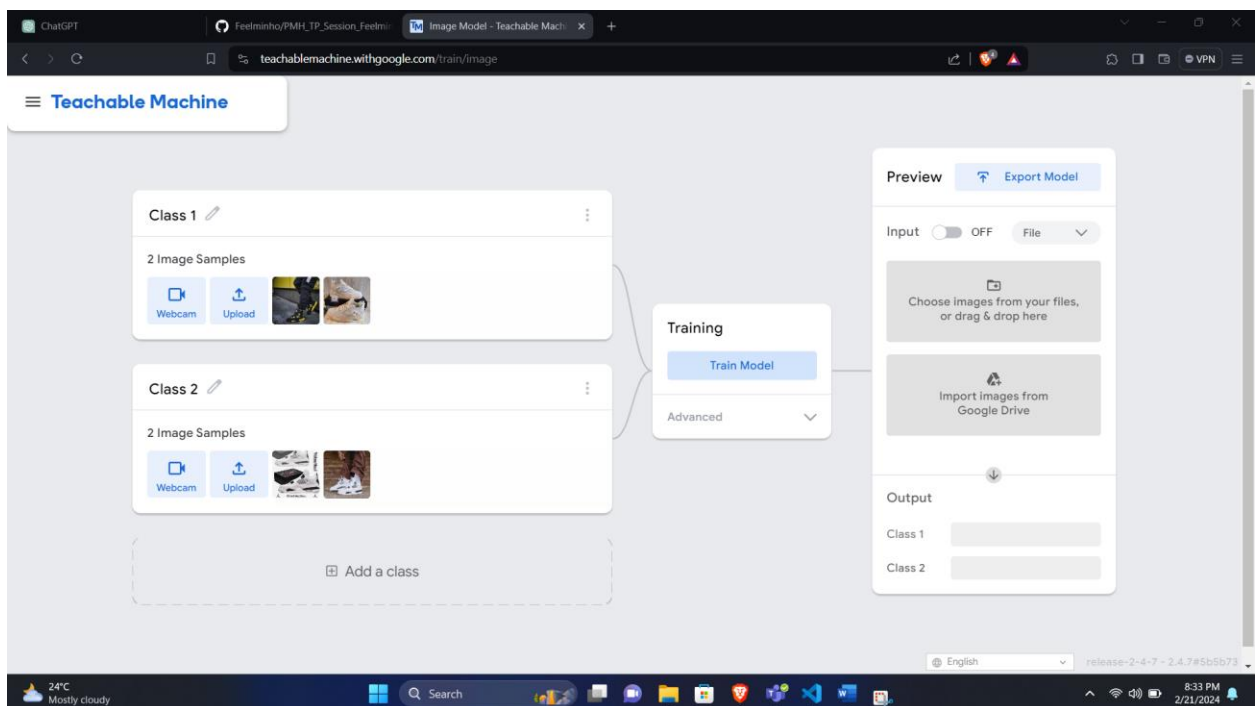
2. Then select Image Project. We Choose Standard Image Model, because we are not training a model to run on a microcontroller
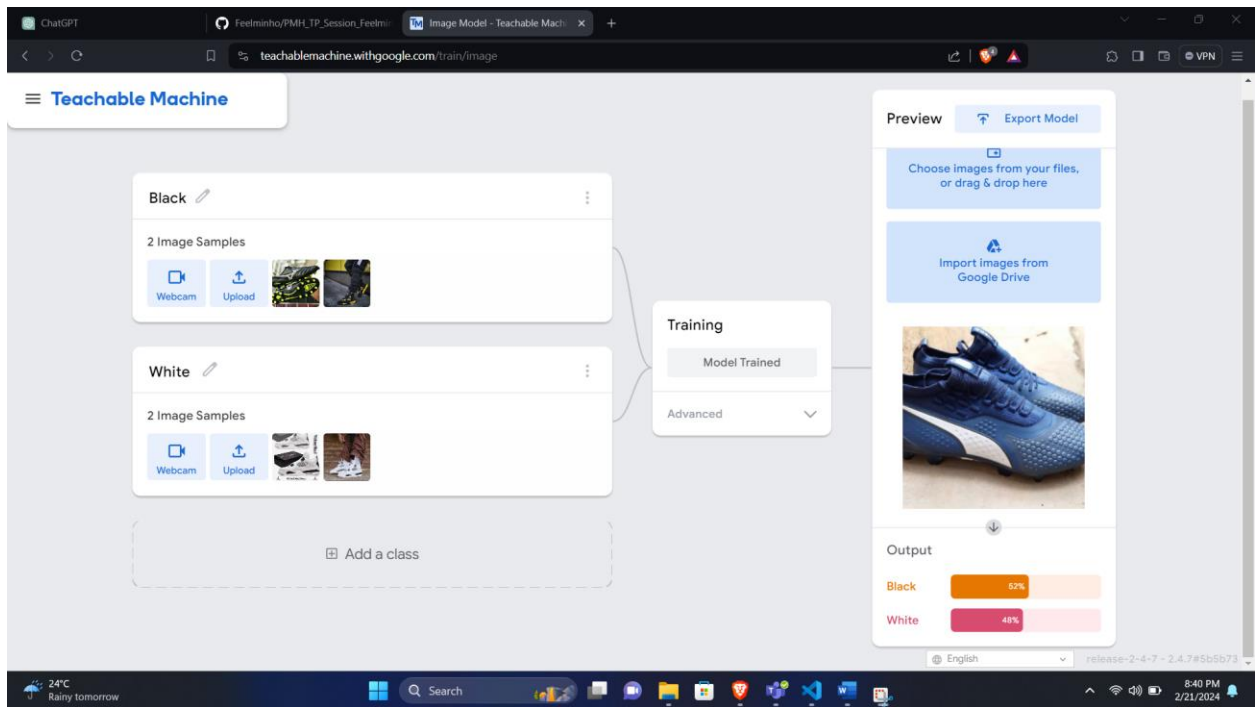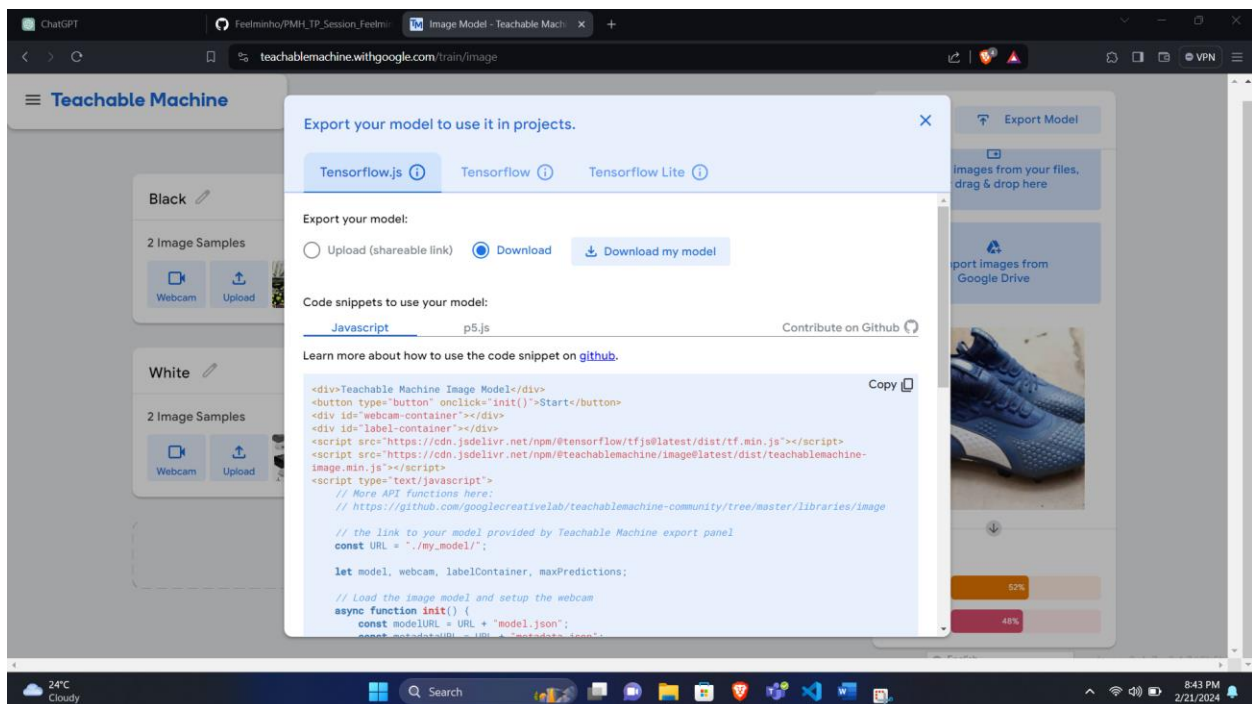
3. Next, add the training samples by clicking Upload under each class. Then, drag the folder of the appropriate plant type from the samples folder to the Choose images from your files … panel.



4. After the training completes, test the model with other images. Use the images in the samples-test folder, like so:

5. Finally, export the model by clicking Export Model on the Preview panel. In the dialog, choose TensorFlow Lite. That's because the target platform is mobile.
6. Next, select Floating point conversion type for the best predictive performance. Then, click Download my model to convert and download the model.
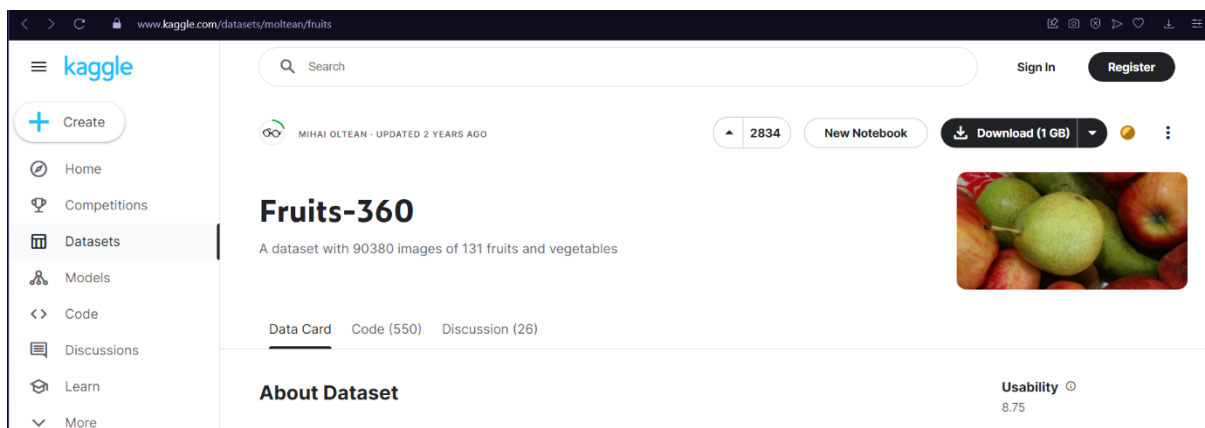
# Task-2 Create an image classification Flutter application:

In this Task, we will implement an image classification Flutter application that will recognize/identify fruits and vegetable images.
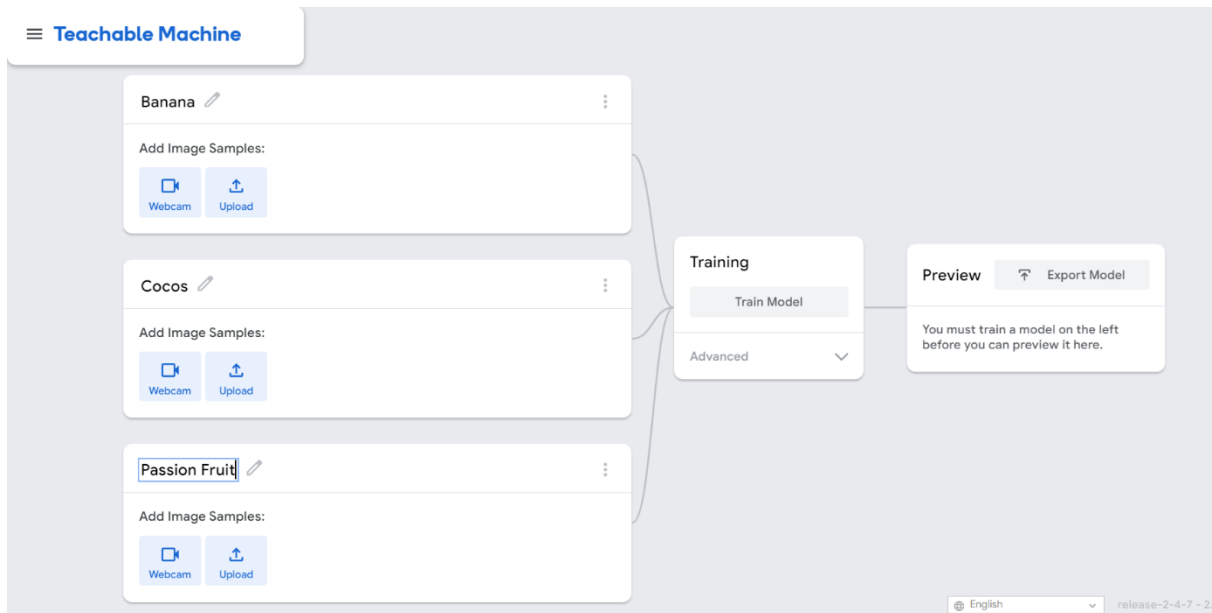
To create this application, we will build a model using Teachable Machine like we have just done in the previous task. This time we will use the Fruits-360 datasets from Kaggle to train the model.

Building a Model with Teachable Machine

1. We will first download the dataset from Kaggle:



2. Next, we will create a new Image Project in Teachable machine and choose Standard Image Model when prompted.

3. Once in the training tool, add the classes and edit the labels of each class, as shown below:

4. Next, add the training samples by clicking Upload under each class. Then, drag the folder of the appropriate plant type from the samples folder to the Choose images from your files panel.

5. After you've added all the training samples, click Train Model to train the model:

6. After the training completes, test the model with a fruit image. Use the images in the test folder, like so:

7. Finally, export the model by clicking Export Model on the Preview panel. In the dialog, choose TensorFlow Lite. That's because the target platform is mobile.

8. After you have the model file converted_tflite.zip in hand, decompress it and copy labels.txt and model_unquant.tflite to the ./assets folder in the starter project.

Implement the flutter app

We will try to modify the project created in Task 1 to recognize the fruit instead of plant.

1. We will first replace the machine learning model – model_unquant.tflite and the classification labels — labels.txt in the project plant recognition by those we got from the Teachable Machine platform.

# Task-3: Create a Flutter app to classify texts:

This documentation guides you through executing a text classification inference from a Flutter application using TensorFlow Serving via both REST and gRPC protocols.

To obtain the code for this codelab, follow these steps:

1. Visit the GitHub repository dedicated to this codelab.

2. Click on the "Code" button and select "Download ZIP" to acquire all the necessary code files for this tutorial.



In this codelab, you will exclusively utilize the files located in the "tfserving-flutter/codelab2" subdirectory within the repository. This subdirectory comprises two main folders:

- **Starter Folder:** This contains the initial code that serves as the foundation for your work in this codelab.

- **Finished Folder:** Here, you can find the fully developed code representing the completed sample app.

3. Retrieve the project dependencies by following these steps:

- Open Visual Studio Code (VS Code).

- Navigate to "File" > "Open Folder."

- Choose the starter folder from the previously downloaded source code. Run this command: flutter pub get

4. Execute and navigate through the application:

Once launched, the app will appear on your Android Emulator or iOS Simulator. The user interface is intuitive, featuring a text field allowing users to input text. Users have the option to transmit data to the backend using either REST or gRPC. The backend employs a TensorFlow model to conduct text classification on the processed input, subsequently relaying the classification outcome to the client app. The app, in turn, updates the UI to reflect the result.



Deploy a text-classification model with TensorFlow Serving

Text classification is a very common machine learning task that classifies texts into predefined categories.

In this codelab, we will deploy a pretrained model from the Train a comment-spam detection model with TensorFlow Lite Model Maker codelab with TensorFlow Serving and call the backend from our Flutter frontend to classify the input text as **spam** or **not spam**.

**Note:** A pretrained SavedModel along with the vocabulary and label files is provided in the *tfserving-flutter/codelab1/mm_spam_savedmodel* folder.

Start TensorFlow Serving

- In a terminal, start TensorFlow Serving with Docker, but replace the *PATH/TO/SAVEDMODEL* placeholder with the absolute path of the mm_spam_savedmodel folder on the computer.

```
PS D:\Dev\flutter-codelabs\tfserving-flutter\codelab2\starter> docker pull tensorflow/serving
Using default tag: latest
latest: Pulling from tensorflow/serving
96d54c3075c9: Pull complete
ce077e3fadc4: Pull complete
806c774cb78b: Pull complete
c588a3276cac: Pull complete
050d4101433f: Pull complete
Digest: sha256:fdc296e313fa4454173c5728fceda38f5d18cdb44c71a9f279ce61bc5818335e
```

```
PS D:\Dev\flutter-codelabs\tfserving-flutter\codelab2\starter> docker run -it --rm -p 8500:8500 -p 8501:8501 -v "D:\Dev\flutter-codelabs\tfserving-flutter\codelab1\m
m_spam_savedmodel:/models/spam-detection" -e MODEL_NAME=spam-detection tensorflow/serving
2024-02-02 18:41:13.470563: I external/org_tensorflow/tensorflow/core/util/port.cc:111] oneDNN custom operations are on. You may see slightly different numerical res
ults due to floating-point round-off errors from different computation orders. To turn them off, set the environment variable `TF_ENABLE_ONEDNN_OPTS=0`.
2024-02-02 18:41:13.472897: I tensorflow_serving/model_servers/server.cc:74] Building single TensorFlow model file config:  model_name: spam-detection model_base_pat
h: /models/spam-detection
2024-02-02 18:41:13.473230: I tensorflow_serving/model_servers/server_core.cc:467] Adding/updating models.
2024-02-02 18:41:13.473284: I tensorflow_serving/model_servers/server_core.cc:596]  (Re-)adding model: spam-detection
2024-02-02 18:41:13.630192: I tensorflow_serving/core/basic_manager.cc:739] Successfully reserved resources to load servable {name: spam-detection version: 123}
2024-02-02 18:41:13.630231: I tensorflow_serving/core/loader_harness.cc:66] Approving load for servable version {name: spam-detection version: 123}
2024-02-02 18:41:13.630238: I tensorflow_serving/core/loader_harness.cc:74] Loading servable version {name: spam-detection version: 123}
2024-02-02 18:41:13.632782: I external/org_tensorflow/tensorflow/cc/saved_model/reader.cc:83] Reading SavedModel from: /models/spam-detection/123
2024-02-02 18:41:13.646410: I external/org_tensorflow/tensorflow/cc/saved_model/reader.cc:51] Reading meta graph with tags { serve }
2024-02-02 18:41:13.646453: I external/org_tensorflow/tensorflow/cc/saved_model/reader.cc:146] Reading SavedModel debug info (if present) from: /models/spam-detectio
n/123
2024-02-02 18:41:13.647297: I external/org_tensorflow/tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU ins
tructions in performance-critical operations.
To enable the following instructions: AVX2 AVX_VNNI FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
2024-02-02 18:41:13.665364: I external/org_tensorflow/tensorflow/compiler/mlir/mlir_graph_optimization_pass.cc:382] MLIR V1 optimization pass is not enabled
2024-02-02 18:41:13.666096: I external/org_tensorflow/tensorflow/cc/saved_model/loader.cc:233] Restoring SavedModel bundle.
2024-02-02 18:41:13.734929: I external/org_tensorflow/tensorflow/cc/saved_model/loader.cc:217] Running initialization op on SavedModel bundle at path: /models/spam-d
etection/123
2024-02-02 18:41:13.740667: I external/org_tensorflow/tensorflow/cc/saved_model/loader.cc:316] SavedModel load for tags { serve }; Status: success: OK. Took 107852 m
icroseconds.
2024-02-02 18:41:13.741774: I tensorflow_serving/servables/tensorflow/saved_model_warmup_util.cc:80] No warmup data file found at /models/spam-detection/123/assets.e
```

Docker automatically downloads the TensorFlow Serving image first, which takes a minute. Afterward, TensorFlow Serving should start. The log should look like this code snippet:

```
PROBLEMS  7   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS

2024-02-02 18:41:13.630238: I tensorflow_serving/core/loader_harness.cc:74] Loading se
2024-02-02 18:41:13.632782: I external/org_tensorflow/tensorflow/cc/saved_model/reader
2024-02-02 18:41:13.646410: I external/org_tensorflow/tensorflow/cc/saved_model/reader
2024-02-02 18:41:13.646453: I external/org_tensorflow/tensorflow/cc/saved_model/reader
n/123
2024-02-02 18:41:13.647297: I external/org_tensorflow/tensorflow/core/platform/cpu_fea
tructions in performance-critical operations.
To enable the following instructions: AVX2 AVX_VNNI FMA, in other operations, rebuild
2024-02-02 18:41:13.665364: I external/org_tensorflow/tensorflow/compiler/mlir/mlir_gr
2024-02-02 18:41:13.666096: I external/org_tensorflow/tensorflow/cc/saved_model/loader
2024-02-02 18:41:13.734929: I external/org_tensorflow/tensorflow/cc/saved_model/loader
etection/123
2024-02-02 18:41:13.740667: I external/org_tensorflow/tensorflow/cc/saved_model/loader
icroseconds.
2024-02-02 18:41:13.741774: I tensorflow_serving/servables/tensorflow/saved_model_warm
xtra/tf_serving_warmup_requests
2024-02-02 18:41:13.788668: I tensorflow_serving/core/loader_harness.cc:95] Successful
2024-02-02 18:41:13.790450: I tensorflow_serving/model_servers/server_core.cc:488] Fin
2024-02-02 18:41:13.790537: I tensorflow_serving/model_servers/server.cc:118] Using In
2024-02-02 18:41:13.790556: I tensorflow_serving/model_servers/server.cc:383] Profiler
2024-02-02 18:41:13.791912: I tensorflow_serving/model_servers/server.cc:409] Running
[warn] getaddrinfo: address family for nodename not supported
2024-02-02 18:41:13.795643: I tensorflow_serving/model_servers/server.cc:430] Exporting
[evhttp_server.cc : 245] NET_LOG: Entering the event loop ...
```

After starting the TensorFlow Serving docker image, we can visit *http://localhost:8501/v1/models/spam-detection/metadata* to inspect the details of the input and output tensors.

```
localhost:8501/v1/models/spam-detection/metadata

{
"model_spec":{
 "name": "spam-detection",
 "signature_name": "",
 "version": "123"
}
,
"metadata": {"signature_def": {
 "signature_def": {
  "serving_default": {
   "inputs": {
    "input_3": {
     "dtype": "DT_INT32",
     "tensor_shape": {
      "dim": [
       {
        "size": "-1",
        "name": ""
       },
       {
        "size": "20",
        "name": ""
       }
      ],
      "unknown_rank": false
     },
     "name": "serving_default_input_3:0"
    }
   },
   "outputs": {
    "dense_5": {
     "dtype": "DT_FLOAT",
     "tensor_shape": {
      "dim": [
       {
        "size": "-1",
        "name": ""
       },
       {
        "size": "2",
        "name": ""
       }
      ],
      "unknown_rank": false
     },
     "name": "StatefulPartitionedCall:0"
    }
   },
```

## Tokenize input sentence

After preparing the backend, the next step is to tokenize the input sentence, enabling the transmission of client requests to TensorFlow Serving. By examining the input tensor of the model, it becomes apparent that it anticipates a list of 20 integer numbers instead of raw strings.

Tokenization involves mapping individual words entered in the app to a list of integers based on a predefined vocabulary dictionary before forwarding them to the backend for classification.

To implement this, add the following code to the **predict()** method in the **lib/main.dart** file to construct the **_vocabMap** vocabulary dictionary.

```dart
164        // Build _vocabMap if empty.
165        if (_vocabMap.isEmpty) {
166          final vocabFileString = await rootBundle.loadString(vocabFile);
167          final lines = vocabFileString.split('\n');
168          for (final l in lines) {
169            if (l != "") {
170              var wordAndIndex = l.split(' ');
171              (_vocabMap)[wordAndIndex[0]] = int.parse(wordAndIndex[1]);
172            }
173          }
174        }
```

Immediately after the previous code snippet, add this code to implement tokenization:

```dart
174        }
175
176        // Tokenize the input sentence.
177        final inputWords = _inputSentenceController.text
178            .toLowerCase()
179            .replaceAll(RegExp('[^a-z ]'), '')
180            .split(' ');
181        // Initialize with padding token.
182        _tokenIndices = List.filled(maxSentenceLength, 0);
183        var i = 0;
184        for (final w in inputWords) {
185          if ((_vocabMap).containsKey(w)) {
186            _tokenIndices[i] = (_vocabMap)[w]!;
187            i++;
188          }
189
190          // Truncate the string if longer than maxSentenceLength.
191          if (i >= maxSentenceLength - 1) {
192            break;
193          }
194        }
```

This code lowercases the sentence string, removes non-alphabet characters, and maps the words to 20 integer indices based on the vocabulary table.

Connect the Flutter app with TensorFlow Serving through REST

There are two ways to send requests to TensorFlow Serving:

- REST

- gRPC

Send requests and receive responses through REST

There are three simple steps to send requests and receive responses through REST:

1. Create the REST request.

2. Send the REST request to TensorFlow Serving.

3. Extract the predicted result from the REST response and render the UI.

*Create and send the REST request to TensorFlow Serving*

1. Right now, the predict() function doesn't send the REST request to TensorFlow Serving. You need to implement the REST branch to create a REST request:

```
lib > ◎ main.dart > ⑂ _TFServingDemoState > ⬡ predict
 194         |
 195
 196         if (_connectionMode == ConnectionModeType.rest) {
 197           // TODO: create and send the REST request
 198
 199           // TODO: process the REST response
 200         } else {
 201           // TODO: create the gRPC request
 202
 203           // TODO: send the gRPC request
 204
 205           // TODO: process the gRPC response
 206         }
 207         return '';
 208       }
```

2. Add this code to the REST branch:

```
195          J
196        if (_connectionMode == ConnectionModeType.rest) {
197          //Create the REST request.
198          final response = await http.post(
199            Uri.parse('http://' +
200                server +
201                ':' +
202                restPort.toString() +
203                '/v1/models/' +
204                modelName +
205                ':predict'),
206            body: jsonEncode(<String, List<List<int>>>{
207              'instances': [_tokenIndices],
208            }),
209          );
```

*Process the REST response from TensorFlow Serving*

- Add this code right after the previous code snippet to handle the REST response:

```
209          );
210
211          // Process the REST response.
212          if (response.statusCode == 200) {
213            Map<String, dynamic> result = jsonDecode(response.body) as Map<String, dynamic>;
214            if ((result['predictions']![0][1] as num) >= classificationThreshold.toDouble()) {
215              return 'This sentence is spam. Spam score is ' +
216                  result['predictions']![0][1].toString();
217            }
218            return 'This sentence is not spam. Spam score is ' +
219                result['predictions']![0][1].toString();
220          } else {
221            throw Exception('Error response');
222          }
```

The postprocessing code extracts the probability that the input sentence is a spam message from the response and displays the classification result in the UI.
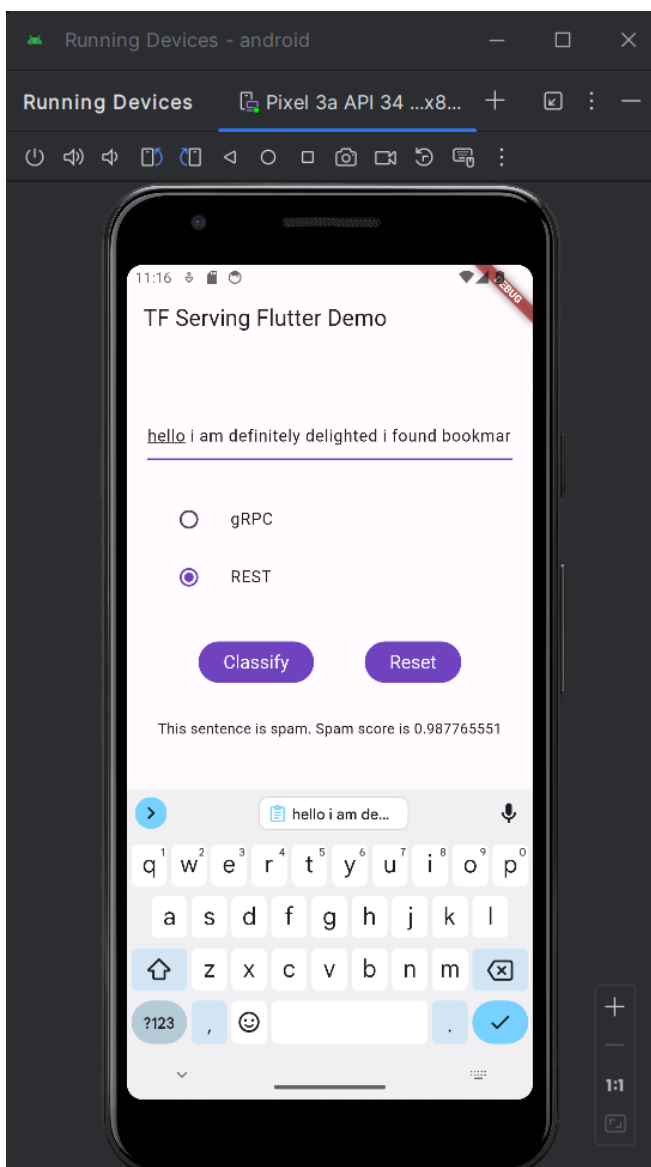
**Run it**

1. Type flutter run in a new terminal. Make sure docker is running on the same time.

2. Enter some text and then select **REST > Classify**.

**Note:** Depending on your training procedure and the input sentence, the model may predict a different spam score. For example, "hello i am definitely delighted i found bookmarking site www winbig com very profitable" is an example sentence that generally receives a high spam score.
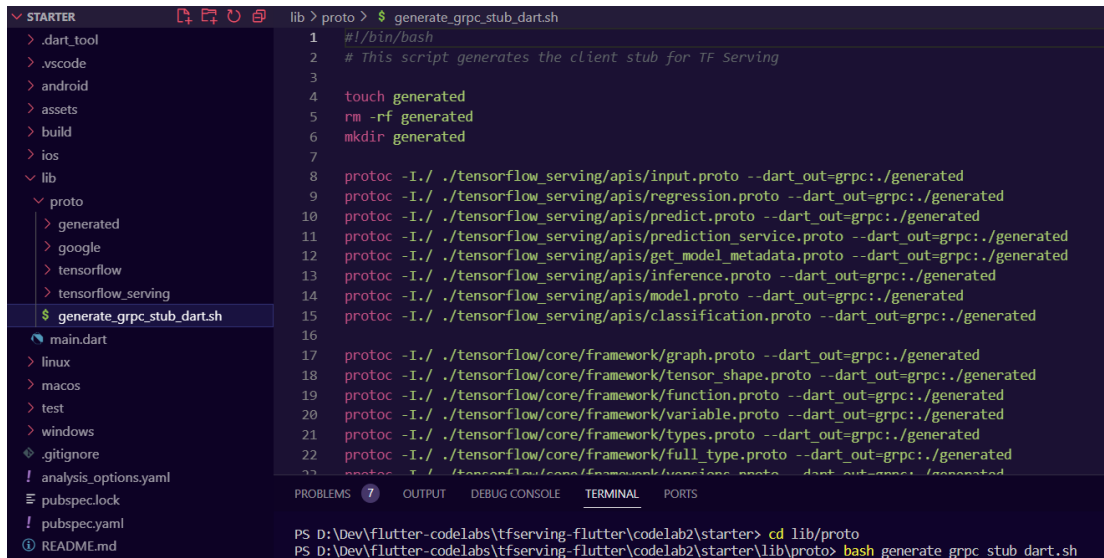
<u>Connect the Flutter app with TensorFlow Serving through gRPC</u>

In addition to REST, TensorFlow Serving also supports gRPC.

gRPC is a modern, open source, high-performance Remote Procedure Call (RPC) framework that can run in any environment. It can efficiently connect services in, and across, data centers with pluggable support for load balancing, tracing, health checking, and authentication. It's been observed that gRPC is more performant than REST in practice.

- In the terminal, navigate to the starter/lib/proto/ folder and generate the stub:

*bash generate_grpc_stub_dart.sh*



Create the gRPC request

Similar to the REST request, you create the gRPC request in the gRPC branch.

```
lib > main.dart > _TFServingDemoState > predict
222              } else {
223                  throw Exception('Error response');
224              }
225          } else {
226              // TODO: create the gRPC request
227
228              // TODO: send the gRPC request
229
230              // TODO: process the gRPC response
231          }
232          return '';
233      }
234  }
```

Add this code to create the gRPC request:

```
lib > main.dart > _TFServingDemoState > predict
223              throw Exception('Error response');
224              }
225          } else {
226              //Create the gRPC request.
227              final channel = ClientChannel(_server,
228                  port: grpcPort,
229                  options:
230                      const ChannelOptions(credentials: ChannelCredentials.insecure()));
231              _stub = PredictionServiceClient(channel,
232                  options: CallOptions(timeout: const Duration(seconds: 10))); // Predict
233
234              ModelSpec modelSpec = ModelSpec(
235                  name: 'spam-detection',
236                  signatureName: 'serving_default',
237              );
238
239              TensorShapeProto_Dim batchDim = TensorShapeProto_Dim(size: Int64(1));
240              TensorShapeProto_Dim inputDim =
241                  TensorShapeProto_Dim(size: Int64(maxSentenceLength));
242              TensorShapeProto inputTensorShape =
243                  TensorShapeProto(dim: [batchDim, inputDim]);
244              TensorProto inputTensor = TensorProto(
245                  dtype: DataType.DT_INT32
```

The input and output tensor names could differ from model to model, even if the model architectures are the same. Make sure to update them if you train your own model.

Send the gRPC request to TensorFlow Serving

Add this code after the previous code snippet to send the gRPC request to TensorFlow Serving:

```
lib > main.dart > _TFServingDemoState > predict
255              // Send the gRPC request.
256              PredictResponse response = await _stub.predict(request);
```

Process the gRPC response from TensorFlow Serving

Add this code after the previous code snippet to implement the callback functions to handle the response:

```dart
lib > main.dart > _TFServingDemoState > predict
257
258        // Process the response.
259        if (response.outputs.containsKey(outputTensorName)) {
260          if (response.outputs[outputTensorName]!.floatVal[1] >
261              classificationThreshold) {
262            return 'This sentence is spam. Spam score is ' +
263                response.outputs[outputTensorName]!.floatVal[1].toString();
264          } else {
265            return 'This sentence is not spam. Spam score is ' +
266                response.outputs[outputTensorName]!.floatVal[1].toString();
267          }
268        } else {
269          throw Exception('Error response');
270        }
```

Now the postprocessing code extracts the classification result from the response and displays it in the UI.

**Run it**

1.  Type flutter run in a new terminal. Make sure docker is running on the same time.

2.  Enter some text and then select **gRPC > Classify**.