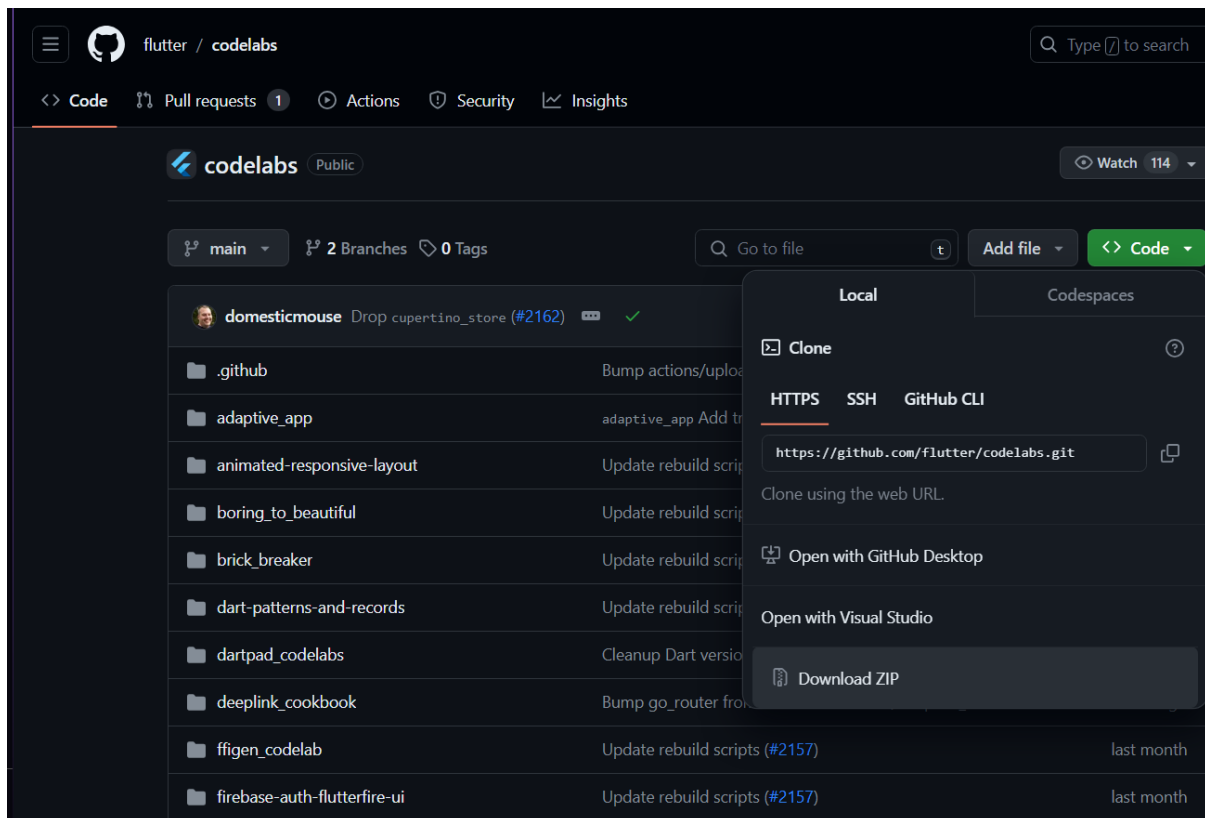**Lab session 6: Integrating Machine Learning Model in a Flutter application**

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

In this documentation, we will learn how to run a text-classification inference from a Flutter app with **TensorFlow Serving** through **REST** and **gRPC**.
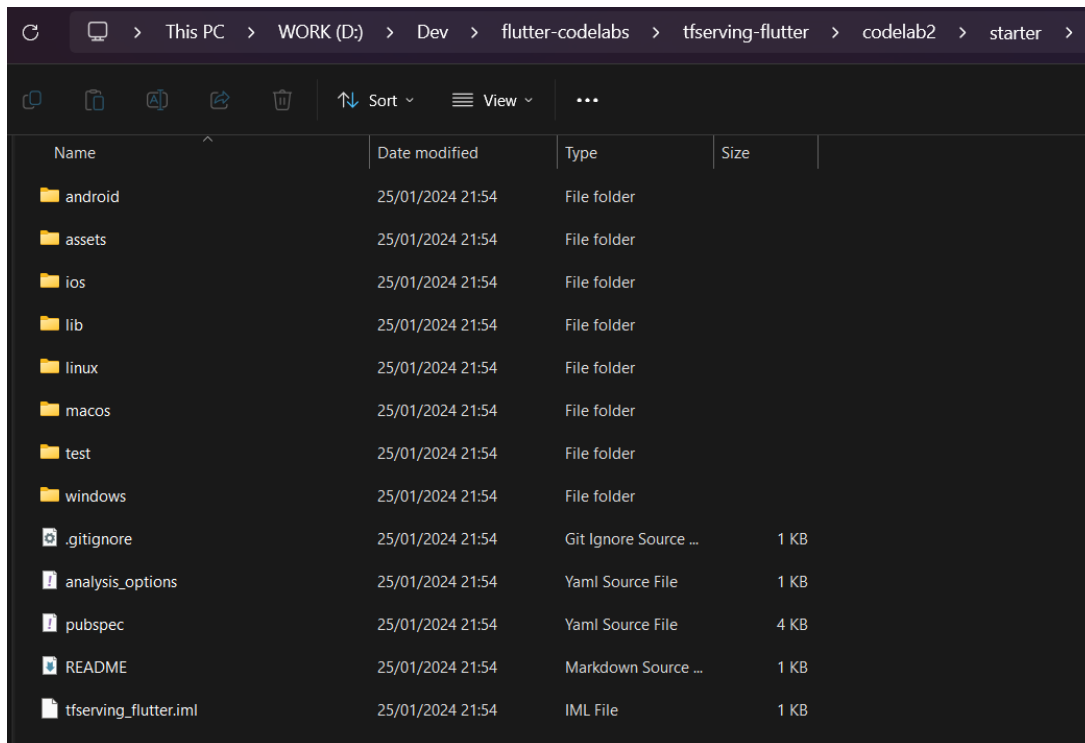
To download the code for this codelab:

1. Navigate to the GitHub repository for this codelab.
2. Click Code > Download zip to download all the code for this codelab.
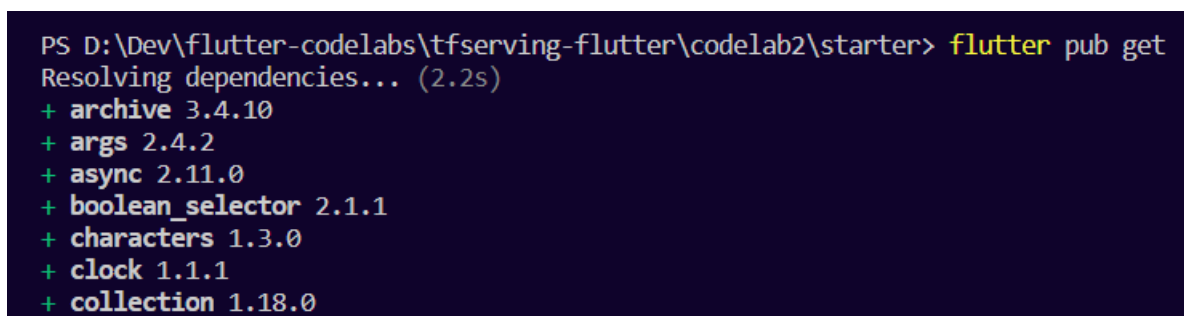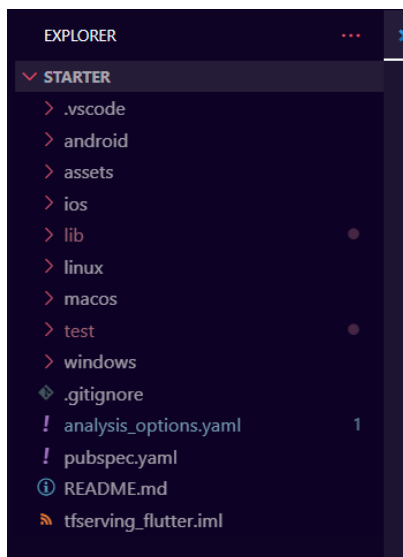


3. Unzip the downloaded zip file to unpack a codelabs-main root folder with all the resources that you need.

For this codelab, we only need the files in the tfserving-flutter/codelab2 subdirectory in the repository, which contains the starter folder with the starter code that we will build upon for this codelab.
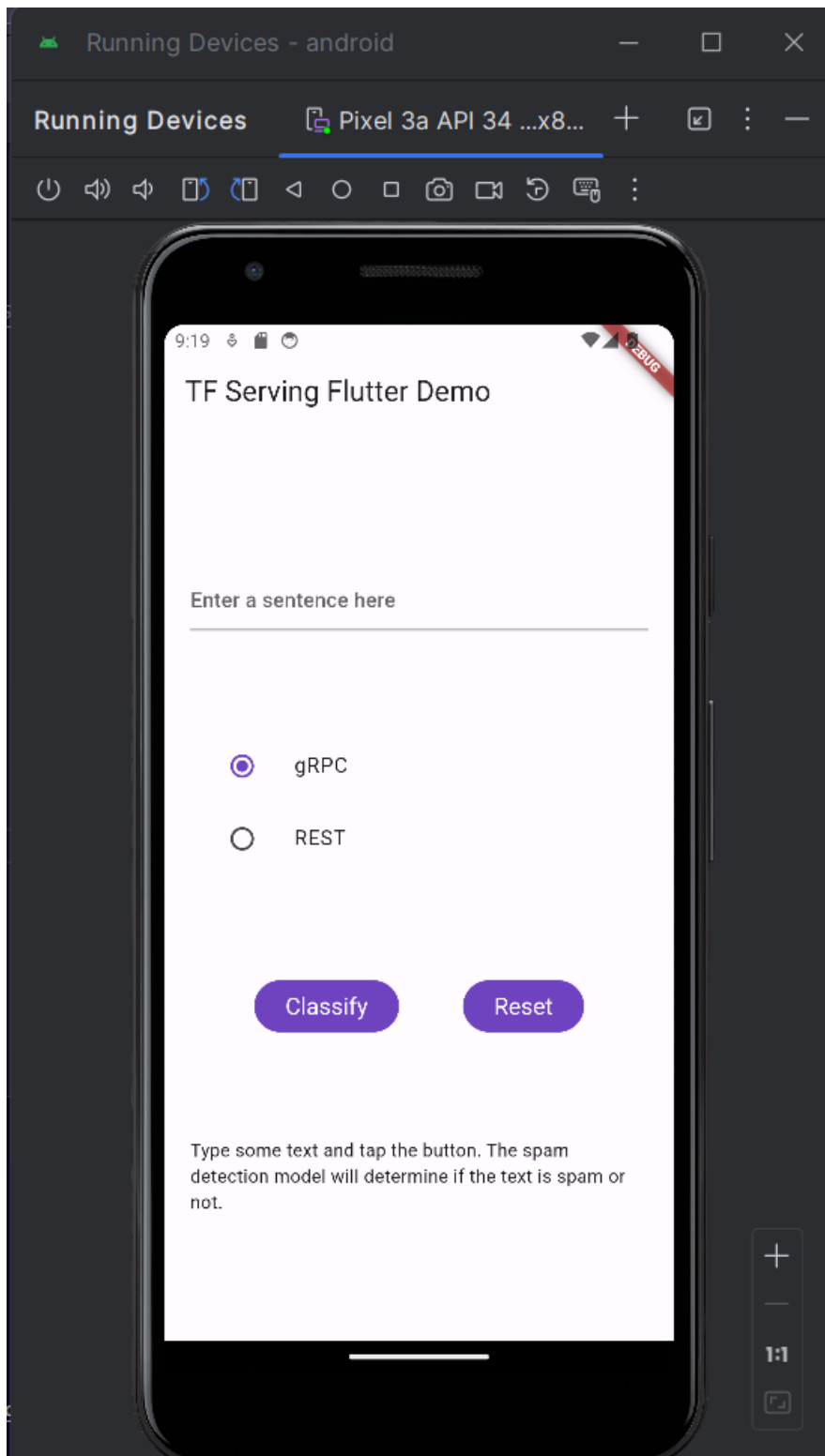
## Download the dependencies for the project

In the first place, we will open the code in vs code and get all the dependancies

## Run and explore the app

Once we have all the necessary dependancies, we will run the app to get an overview of all the features we can see in it and how it's displayed:



## Deploy a text-classification model with TensorFlow Serving

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

3

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: Downloading [>                                    ]    277.8kB/27.51MB
ce077e3fadc4: Downloading [====================================>]    1.996MB/2.65MB
806c774cb78b: Downloading [>                                    ]    1.602MB/115.3MB
c588a3276cac: Waiting
050d4101433f: Waiting
```

```
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:

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.

## Tokenize input sentence

Once the backend ready, we need to tokenize the input sentence to be able to send client requests to TensorFlow Serving.

If we inspect the input tensor of the model, we can see that it expects a list of 20 integer numbers instead of raw strings.

Tokenization is when we map the individual words typed in the app to a list of integers based on a vocabulary dictionary before sending them to the backend for classification.

1. In the lib/main.dart file, add this code to the predict() method to build the _vocabMap vocabulary dictionary.

```
lib > main.dart > _TFServingDemoState > predict
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:

```
lib > main.dart > _TFServingDemoState > predict
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        }
```

6

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

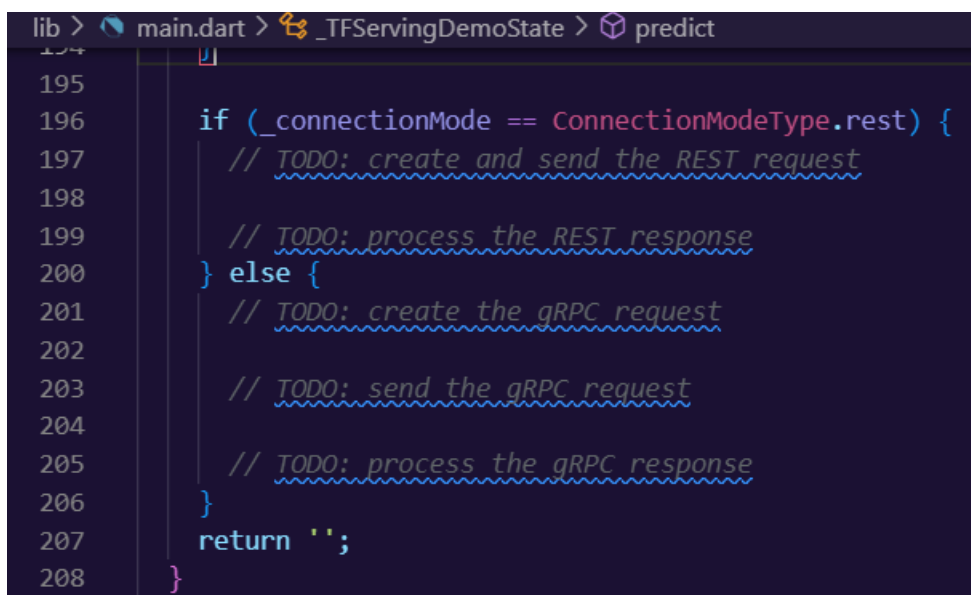There are two ways to send requests to TensorFlow Serving:

- 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:

```dart
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:

```dart
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.

## Connect the Flutter app with TensorFlow Serving through gRPC

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.



Add this code to create the gRPC request:

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);
```

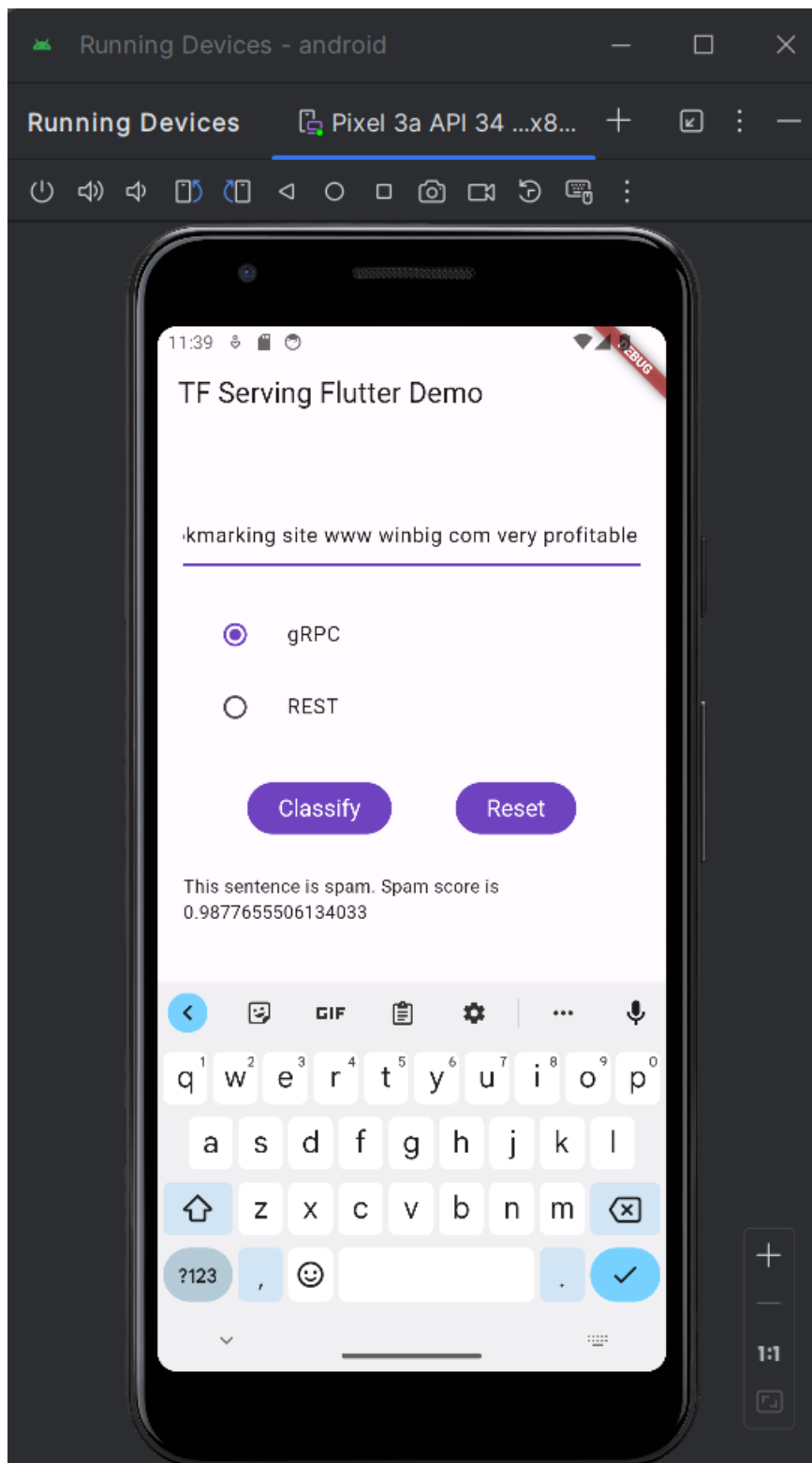## Process the gRPC response from TensorFlow Serving

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

```
lib > ● main.dart > ⟨⟩ _TFServingDemoState > ⊘ predict
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**.

## Conclusion

This lab helped us discover new features of TensorFlow. It made the text-classification to a quite simple process.