

TensorFlow Lite Extension to OpenNMT-tf Documentation

by Gerardo Cervantes and Stephen LaRocca

Approved for public release: distribution unlimited.

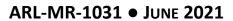
NOTICES

Disclaimers

The findings in this report are not to be construed as an official Department of the Army position unless so designated by other authorized documents.

Citation of manufacturer's or trade names does not constitute an official endorsement or approval of the use thereof.

Destroy this report when it is no longer needed. Do not return it to the originator.





TensorFlow Lite Extension to OpenNMT-tf Documentation

Gerardo Cervantes Advanced Resource Technologies, Inc.

Stephen LaRocca Computational and Information Sciences Directorate, DEVCOM Army Research Laboratory

Approved for public release: distribution unlimited.

| REPORT DOCUMENTATION | | | N PAGE | | Form Approved OMB No. 0704-0188 |
|--|---|--|---|---|--|
| data needed, and comple burden, to Department of Respondents should be a valid OMB control numb | ting and reviewing the collect f Defense, Washington Heado ware that notwithstanding an per. | ion information. Send commen uarters Services, Directorate fo | ts regarding this burden estin r Information Operations and son shall be subject to any pe | nate or any other aspect Reports (0704-0188) | structions, searching existing data sources, gathering and maintaining the et of this collection of information, including suggestions for reducing the , 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, nply with a collection of information if it does not display a currently |
| 1. REPORT DATE (| DD-MM-YYYY) | 2. REPORT TYPE | | | 3. DATES COVERED (From - To) |
| June 2021 | | Memorandum Re | port | | September 2020–March 2021 |
| 4. TITLE AND SUB | TITLE | | | | 5a. CONTRACT NUMBER |
| TensorFlow Li | te Extension to O | penNMT-tf Docum | nentation | | W911QX20D0012 |
| | | | | | 5b. GRANT NUMBER |
| | | | | | 5c. PROGRAM ELEMENT NUMBER |
| 6. AUTHOR(S) | | L D | | | 5d. PROJECT NUMBER |
| Gerardo Cerva | ntes and Stephen | LaRocca | | | 5e. TASK NUMBER |
| | | | | | 5f. WORK UNIT NUMBER |
| 7. PERFORMING C | RGANIZATION NAME | (S) AND ADDRESS(ES) | | | 8. PERFORMING ORGANIZATION REPORT NUMBER |
| DEVCOM Arr ATTN: FCDD | ny Research Labo -RLC-IB | oratory | | | ARL-MR-1031 |
| Adelphi, MD 20783-1138 | | | | | |
| | | (NAME(S) AND ADDRE | SS/ES) | | 10. SPONSOR/MONITOR'S ACRONYM(S) |
| 5. 51 01301110/1 | | | 33(E3) | | |
| | | | | | 11. SPONSOR/MONITOR'S REPORT NUMBER(S) |
| 12. DISTRIBUTION | AVAILABILITY STATE | MENT | | | |
| Approved for p | oublic release: dist | tribution unlimited. | | | |
| 13. SUPPLEMENT | | | | | |
| ORCID ID(s): | Cervantes, 0000-0 | 0002-4392-5017; L | aRocca, 0000-00 | 03-3341-5520 |) |
| including mach OpenNMT me use on Android Development O TensorFlow Li to the Lite vari occupying less | nine learning infer thods were achiev l platforms pendir Command Army F te and succeeded ant. Deployable o space. A pull req s how OpenNMT ttforms. | ence with artificial ed using TensorFlo ng completion of th Research Laborator in implementing a n Android devices, uest submitted to th | neural network r ow since 2018; ho e TensorFlow Lit y Shareable Com new method for c these converted ne OpenNMT was | nodels on An owever, most te library. The ponents proje converting Op models provid s approved an | ools for Neural Machine Translation (NMT) droid platforms. Rapid advances in of these advances were not deployable for e US Army Combat Capabilities of team closely tracked progress on benNMT models from standard TensorFlow de important gains in execution speed while id implemented in February 2021. This l types of models to TensorFlow Lite for use |
| | | ing, Android, open | development too | ls TensorFlo | wLite |
| 16. SECURITY CLAS | | ing, maioia, open | 17. LIMITATION OF | 18. NUMBER OF | 19a. NAME OF RESPONSIBLE PERSON |
| | | | ABSTRACT | PAGES | Stephen LaRocca |
| a. REPORT | b. ABSTRACT | c. THIS PAGE | UU | 18 | 19b. TELEPHONE NUMBER (Include area code) |
| Unclassified | Unclassified | Unclassified | I | | (301) 394-3198 |

| Standard Form 298 (Rev. | 8/98 |
|----------------------------|------|
| Prescribed by ANSI Std. Z. | |

Contents

| List | of Ta | bles | | v |
|------|-------|---------------|--------------------------------------|---|
| 1. | Intr | oductic | on | 1 |
| | 1.1 | Docum | nent Contents | 1 |
| | 1.2 | Purpos | se | 1 |
| | 1.3 | Distrib | pution | 1 |
| | 1.4 | Tools | | 1 |
| | 1.5 | Conve | rtible Models | 1 |
| 2. | Run | ning th | e Conversion | 2 |
| | 2.1 | Installi | ing | 2 |
| | 2.2 | Runnir | ng | 2 |
| | 2.3 | Data C | Configuration File | 3 |
| | 2.4 | Vocab | ulary Files and IDs Used for Running | 3 |
| | | 2.4.1 | | 3 |
| | | 2.4.2 | Description | 3 |
| | | 2.4.3 | Use | 3 |
| | | 2.4.4 | Example Code Snippet | 4 |
| | 2.5 | Runnir | ng on Android | 4 |
| | 2.6 | Quant | ization Options | 5 |
| | | 2.6.1 | Model Quantization Options | 5 |
| | | 2.6.2 | How | 5 |
| | | 2.6.3 | Advantages | 6 |
| | | 2.6.4 | Future | 6 |
| | 2.7 | Tests Created | | 6 |
| | 2.8 | Functi | on Converted Description | 6 |
| | | 2.8.1 | Parameter | 7 |
| | | 2.8.2 | Returns | 7 |
| 3. | Add | itional | Technical Details | 7 |
| | 3.1 | Code N | Modifications/Additions | 7 |
| | | 3.1.1 | Source | 7 |

| | 3.1.2 | Testing | 8 |
|-------------|---------|---|----|
| 3.2 | Convers | sion Process | 8 |
| | 3.2.1 | Initial Running Tests | 8 |
| | 3.2.2 | TensorFlow Conversion Methods | 9 |
| | 3.2.3 | Custom Inference Function | 9 |
| | 3.2.4 | Major Changes to Running a Normal Inference | 9 |
| Distributio | on List | | 11 |

List of Tables

| Table 1 | Automated conversion tests | 6 |
|---------|--------------------------------------|---|
| Table 2 | OpenNMT-tf source code modifications | 8 |
| Table 3 | Testing files modified | 8 |

1. Introduction

1.1 Document Contents

This report covers the new extension to the popular Machine Translation (MT) open-source package OpenNMT-tf. This extension allows models that were trained using this package to be mobile-compatible with TensorFlow Lite.

1.2 Purpose

The purpose of this report is to convert MT models trained with the OpenNMT-tf package to be mobile friendly with TensorFlow Lite. The conversion of these models to TensorFlow Lite improves the models' run time on mobile environments and allows the models to run using the onboard graphics processing unit (GPU) of the mobile device.

1.3 Distribution

The OpenNMT-tf open-source package can be found on GitHub, <u>https://github.com/OpenNMT/OpenNMT-tf</u>. This code is a pull request to the open-source library that was reviewed and accepted by the OpenNMT developers.

1.4 Tools

The following tools are needed in order to convert OpenNMT models to TensorFlow Lite:

- Python 3.5+
- Python pip
- OpenNMT-tf package, <u>https://pypi.org/project/OpenNMT-tf/</u>

1.5 Convertible Models

The following models from OpenNMT are convertible to TensorFlow Lite:

- NMTSmallV1
- NMTMediumV1
- NMTBigV1
- Luong Attention

2. Running the Conversion

2.1 Installing

To install and run the extension, first download Python (recommended version is Python 3.7). Make sure this download adds Python and pip onto your PATH environment variables so you can use the commands in the command line.

To install the extension in the command line, run this command to get the necessary Python packages:

```
pip install OpenNMT-tf
```

2.2 Running

Make sure to install the Python packages before trying to run the conversion.

To run the conversion from the command line, navigate to the main directory of the code (where the README resides).

To convert an OpenNMT-tf model to TensorFlow Lite, you must first train a model. A small guide on training a model can be found here: https://opennmt.net/OpenNMT-tf/quickstart.html.

Run this command with your model-specific options to convert the model to TensorFlow Lite:

```
python opennmt/bin/main.py --model_type NMTSmallV1 --config
opensubs/data.yml export --output_dir
./opensubs/smallnmt/tflite --export_format tflite
```

Parameter descriptions are as follows:

- --*model_type* The model architecture you are trying to load to convert to TFLite; these can be found in opennmt/models/catalog.py.
- --config The data configuration file, written in YAML. The file also contains the path to the model. Parameters can be found here: https://opennmt.net/OpenNMT-tf/configuration.html
- *export* This tells the algorithm that you want to save the model as a TensorFlow Lite model.
- --output dir This is the directory where you want your model to be saved.
- --format This tells OpenNMT to export as a TensorFlow Lite model; TensorFlow Lite options are "tflite" and "tflite_float16".

The conversion should take about $1-2 \min$ to run.

2.3 Data Configuration File

Running anything on OpenNMT requires that you give OpenNMT a data configuration file of structure <u>YAML</u>. All the possible data configurations can be found here: <u>https://opennmt.net/OpenNMT-tf/configuration.html</u>

The data.yml file created when using the quick-start guide in model training is the data configuration file. This data configuration file is needed to run any operation in OpenNMT and is required to convert models to TensorFlow Lite.

One optional parameter added to OpenNMT with the code extension is

params:

tflite_output_size: 250

This parameter tells the algorithm what size array to make the output; if this optional parameter is not provided, then the default value used is 250.

2.4 Vocabulary Files and IDs Used for Running

2.4.1 Acquiring Vocabulary Files

During the process of training an OpenNMT-tf model, two vocabulary files will have been created. The location of these vocabulary files are specified in the data configuration file and are required to convert a model to TensorFlow Lite and to run a TensorFlow Lite–converted model on Android.

2.4.2 Description

The vocabulary files consist of a list of words of the language separated by new line characters, where each word in the language is given a unique number, commonly referred to as the word ID. To convert a word not in the vocabulary to an ID, you will have to mark it as an unknown word ID, which will be the biggest unique ID available plus 1.

2.4.3 Use

Running the TensorFlow Lite model requires an array of integers instead of text. To translate text using the model, you must convert the text to IDs using the vocabulary file of the language you are converting from. This is achieved by separating each word with white space. Then, you must look up the unique ID of each word in the vocabulary file. After getting the IDs of all of the words, you can store them in an array and pass the IDs to the model to get a translation. The returned translation from the model will be integer IDs convertible to text using the other vocabulary file.

2.4.4 Example Code Snippet

The following Android code snippet produces a HashMap that maps a word to the unique ID given in the file.

```
public HashBiMap<String, Integer> readVocab(InputStream file){
    try {
        BufferedReader brFile = new BufferedReader(new
InputStreamReader(file));
        HashBiMap<String, Integer> vocab =
HashBiMap.create();
        String wordRead = brFile.readLine();
        int index = 0;
        while(wordRead != null){
            vocab.put(wordRead, index);
            index += 1;
            wordRead = brFile.readLine();
        }
        return vocab;
    }
        catch(IOException e) {
            return null;
        }
      }
```

2.5 Running on Android

This section provides portions of code that will load and run the model on Android devices.

2.5.1 Gradle Dependencies

Adding these dependencies in the Android Studio Gradle file (build.gradle) will allow TensorFlow Lite models to be run on Android. The following lines provide the "nightly" version for TensorFlow, which is the beta version.

```
implementation 'org.tensorflow:tensorflow-lite:0.0.0-nightly'
implementation 'org.tensorflow:tensorflow-lite-select-tf-ops:0.0.0-
nightly'
implementation 'org.tensorflow:tensorflow-lite-support:0.0.0-nightly'
```

2.5.2 Loading the Model

The following code is used to load the model and assumes you have stored the saved model to the assets folder. The code will look into the assets folder to find the model you specified. The *modelPath* variable should be set to the location of

the model. The *NUM_LITE_THREADS* should be set to an integer that tells it how many threads to use when running the model. Tested values are 1 and 4.

```
AssetManager assetManager =
this.context.getResources().getAssets();
    ByteBuffer buffer = loadModelFile(assetManager,
modelPath);
if (buffer == null) {
    Log.e("nmt-tf", "Could not load model");
    return;
    }
    Interpreter.Options opt = new Interpreter.Options();
    opt.setNumThreads(NUM_LITE_THREADS);
    tflite = new Interpreter(buffer, opt);
```

2.5.3 Running the Model

The following code will run the model with the given input. Since the function takes in IDs, you should read and load the vocabulary files onto the application. You should also have functions to convert sentences to an integer array of the respective IDs as well as a function to convert an integer array of IDs back to a sentence.

```
int[] input_ids = sentenceToIds("Hello World");
int[] output_ids = new int[250];
tflite.run(input_ids, output_ids);
String translatedSentence = IdsToSentence(output_ids);
System.out.println("Translated Sentence: " +
translatedSentence);
```

output_ids is initialized to the array. After *tflite.run* is called, it will modify this variable to have the result of the model inference. *output_ids* should be the same size as *tflite_output_size*.

2.6 Quantization Options

Quantizing the model can be useful for reducing the model size and increasing the run time.

2.6.1 Model Quantization Options

- Normal model
- 16-bit float quantized

2.6.2 How

The model can easily be quantized by modifying the argument value "format" when exporting the model to TensorFlow Lite.

2.6.3 Advantages

The model size is reduced by half, and according to TensorFlow, the run time should be faster.

2.6.4 Future

For an additional speedup, it may be worth quantizing differently. The following options require a representative dataset to be able to smartly know how to quantize the variables to 8 bits:

- Dynamic range quantization
 - Weights to 8-bit integers; outputs are still 32 bit
 - Four times smaller, two to three times speedup
- Full integer quantization
 - Weights to 8-bit integers
 - Four times smaller, three times plus speedup

These are great options to consider if the models need to run faster.

2.7 Tests Created

The tests listed in Table 1 check if the TensorFlow Lite function produces the same results as the original inference function and if the models could be converted.

| Test name | Model | Test description |
|----------------------|-----------------------------|---|
| testTFLiteOutput | Luong Attention NMTBigV1 | Checks that the model prediction is the same when outputting with the TensorFlow Lite conversion. |
| testTFLiteOutputFile | Luong Attention NMTBigV1 | Checks that the code is able to convert and output the model into a file. |

Table 1Automated conversion tests

The tests can be found and ran from the following file: opennmt/tests/tflite_test.py.

2.8 Function Converted Description

The function being converted is named *infer_tflite* and can be found in the file opennmt/models/sequencetosequence.py. This function runs when you predict using the TensorFlow Lite model that was converted; it runs on any model that uses SequenceToSequence.

2.8.1 Parameter

The function that takes in one parameter is named "ids".

To run IDs in Python, you give the function a 32-bit integer Tensor, which can be any size as long as it is 1-D. When running on Android, you will have to feed in a 1-D integer array; this array can be any size. The parameter is an integer array and not a String because the function takes in an ID for each word in the sentence. Refer to Section 2.4, Vocabulary Files and IDs Used for Running, to learn how to convert the sentence you want translated back to IDs.

2.8.2 Returns

The function returns a 1-D integer array. The size of the array is by default 250 unless a size was specified in the data configuration file. Section 2.3, Data Configuration File, covers how to specify a different output size.

In Python, the TensorFlow Lite function will return a 32-bit integer Tensor. Running on Android will return a 32-bit integer array. The function does not return a String; it will return IDs, which will have to be converted back to a sentence to find the translation of the sentence. Refer to Section 2.4, Vocabulary Files and IDs Used for Running, to go from IDs back to a sentence.

3. Additional Technical Details

3.1 Code Modifications/Additions

This section is a brief overview of why some files were modified or added, which files changed, and the reason for the modification.

Source file changes are the changes that were made to add the TensorFlow Lite support. Testing files are files that were modified or added to include automated unit tests to check if the TensorFlow Lite conversion works.

3.1.1 Source

The files listed in Table 2 were modified to implement the TensorFlow Lite conversion option.

| File | Reason |
|--|--|
| opennmt/utils/exporters.py | Adds "tflite" as a method to export, making models mobile friendly. This adds the class TFLiteExporter to handle exporting. The model produces an error if it is not compatible with TensorFlow Lite. Otherwise, the model is loaded and converted. |
| opennmt/inputters/text_inputter.py | Adds the function <i>tflite_call</i> that gets the word embeddings in a way that is safe for TensorFlow Lite. The fix and problem description can be found in <u>TensorFlow issue</u> <u>#42410</u> . |
| opennmt/models/sequence_to_sequence.py | Creates the custom inference function, which is the concrete function that will be converted. A parameter is added to <i>_dynamic_decode</i> to tell the function that it is running in a TensorFlow Lite safe way, and if running with TensorFlow Lite, then it returns IDs. |
| opennmt/decoders/decoder.py | Modifies the <i>dynamic_decode</i> function to use the modified word embedding function. |
| opennmt/utils/decoding.py | Modifies the decoding. The original function uses <i>TensorArrays</i> , which are not compatible with TensorFlow Lite unless they are specified as static and are pre-allocated. |

Table 2 OpenNMT-tf source code modifications

3.1.2 Testing

The files listed in Table 3 were modified or added to create the tests.

Table 3Testing files modified

| File | Reason |
|------------------------------|---|
| opennmt/tests/tflite_test.py | Contains four test cases created to check if the models NMTBigV1 and Luong Attention produce the same output when running as a TensorFlow Lite model and to check if they were able to convert and save the model. |

3.2 Conversion Process

The conversion runs through the command line. Running details are provided in Section 2.2.

3.2.1 Initial Running Tests

After arguments are parsed, the function *save_tflite* is called, which loads the model structure and restores the tensor weights of the model. To test for normal operation, it runs on translating the small example "Hello World". This is to ensure that the

model runs and not to verify that the output is correct. The model then tries to run the custom function created in this TensorFlow Lite modification, which runs the model in a way that is TensorFlow Lite safe; this will be explained later in the report.

3.2.2 TensorFlow Conversion Methods

There are three different ways to convert a model to TensorFlow Lite.

- From a TensorFlow Saved model
 - Converting from a Saved model was problematic because the word embeddings were not being saved onto the TensorFlow Lite model being converted. This resulted in an error.
- From a Keras model
 - We cannot convert using this method because our model is not a Keras model.
- From a concrete function
 - This method worked the best and was used to convert the model to TensorFlow Lite. TensorFlow did much of the work during conversion, including automatically specifying which variables to save. To convert from a concrete function, a custom inference function was created. After ensuring it was compatible with TensorFlow Lite operations, the function was converted.

3.2.3 Custom Inference Function

The custom inference function *infer_tflite* was used to convert the model to TensorFlow Lite. The function inputs and outputs are described in Section 2.8, Function Converted Description. To feed in a sentence to the algorithm, convert the sentence to IDs using the same vocabulary used to train the model.

3.2.4 Major Changes to Running a Normal Inference

Modification to word embedding: The function *tflite_call* was added to get the required word embeddings. This function is called twice: once during encoding and again during decoding. One problem we encountered is also described in <u>TensorFlow issue #42410</u>. The workaround described on the GitHub issue was implemented to fix errors that occurred while trying to make the TensorFlow Lite conversion work.

Different decoding: For decoding, we have to assume the batch size is 1. There is an object from the TensorFlow library named *TensorArray* that is not fully supported on TensorFlow. The object is only supported when the dynamic size flag is set to flag, and the size is specified when creating the *TensorArray*. During decoding, the *TensorArray* created is replaced with a static size *TensorArray*. We assume that the batch size is 1, which is a safe assumption since we are running on Android.

Inputs and outputs: This inference function only takes in an integer array as IDs and outputs IDs as integers, as explained in Section 2.8, Function Converted Description. Since the return of this function is only IDs, this function stops early and only returns the IDs when decoding. The regular inference takes in a dictionary that includes IDs and other required variables.

| 1 | DEFENSE TECHNICAL |
|-------|-------------------|
| (PDF) | INFORMATION CTR |
| | DTIC OCA |

| 1 | DEVCOM ARL |
|---|------------|
|---|------------|

| (PDF) | FCDD RLD DCI |
|-------|--------------|
| | TECH LIB |

| 1 | DEVCOM ARL |
|-------|-------------|
| (PDF) | FCDD RLC IB |
| | S LAROCCA |