

# ON-DEVICE MACHINE TRANSLATION WITH REAL-TIME ELABORATION AND QnA

**Aurosweta Mahapatra**  
206070948  
aurosweta99@ucla.edu

**Shrey Mathur**  
205928673  
shreyzmail@gmail.com

**Syam Sundar Kirubakaran**  
805905367  
syamk@ucla.edu

**Vaibhav Tiwari**  
105912719  
vaibhavtiwari@ucla.edu

## 1 Abstract

2 This project addresses language  
3 barriers and translation latency by  
4 developing an on-device machine  
5 translation system for English to  
6 Spanish. It employs RNN and  
7 Transformer models, evaluating  
8 translation quality using BLEU and  
9 Rouge metrics. The best model, the  
10 RNN, is converted to TensorFlow  
11 Lite for practical deployment. The  
12 system incorporates on-device text  
13 recognition from still images,  
14 allowing instant English-to-  
15 Spanish translations. Its unique  
16 feature is real-time elaboration,  
17 providing additional contextual  
18 information beyond translation.  
19 This enhances user experience and  
20 enables seamless communication  
21 across languages.

## 22 1. Introduction

### 23 1.1 Background

24 Language barriers and translation latency are  
25 persistent challenges in effective global  
26 communication. To address these issues, this  
27 project focuses on developing an on-device  
28 machine translation system that offers real-  
29 time elaboration and question-and-answer  
30 (QnA) capabilities. Our specific translation  
31 task involves converting English text to  
32 Spanish, employing recurrent neural network  
33 (RNN) and Transformer models. Evaluation  
34 of translation quality utilizes standard metrics  
35 like BLEU and Rouge.

36 To ensure practical implementation, we  
37 convert our best-performing model, the RNN,  
38 from TensorFlow to TensorFlow Lite,  
39 optimizing it for deployment on edge devices.  
40 We also integrate on-device text recognition

41 for still images, allowing instant English-to-  
42 Spanish translations.

43 The unique feature of our system lies in its real-  
44 time elaboration, going beyond basic  
45 translation. For example, if a user inputs a  
46 question such as "How to make an omelette,"  
47 our system not only provides the translation but  
48 also offers step-by-step instructions in Spanish.  
49 This anticipates users' potential future queries,  
50 enhancing the user experience.

51

### 52 1.2 Dataset

53 We have taken the dataset from the  
54 sentences\_detailed.csv file from tatoeba.org.  
55 ([http://tatoeba.org/files/downloads/sentences\\_](http://tatoeba.org/files/downloads/sentences_detailed.csv)  
56 [detailed.csv](http://tatoeba.org/files/downloads/sentences_detailed.csv)).

## 57 2 Models

58 This project delves into the exploration of two  
59 different models for our machine translation  
60 task: RNN and Transformer. The objective is to  
61 evaluate their performance and effectiveness in  
62 addressing these issues. In the following  
63 sections, we provide detailed explanations of  
64 each model and their respective  
65 methodologies.

66

### 67 2.1 RNN

68 A Recurrent Neural Network (RNN) is a type  
69 of neural network that is well-suited for  
70 processing sequential data. It's architecture  
71 allows it to retain information from previous  
72 steps and utilize it in subsequent steps,  
73 allowing it to capture temporal dependencies.  
74 Since RNNs can process sequential data and  
75 capture contextual dependencies between  
76 words, Neural machine translation systems  
77 are typically implemented with a Recurrent  
78 Neural Network (RNN) based encoder-  
79 decoder framework (Bahdanau et al., 2016).  
80 For our experiment, we decided to use a

81 bidirectional RNN encoder-decoder  
 82 architecture (Bahdanau et al., 2016; Yang et  
 83 al., 2017). In this setup, the source language  
 84 sentence is encoded using a bidirectional  
 85 RNN. The encoded representation of the  
 86 source sentence is then passed to a decoder  
 87 RNN, which generates the corresponding  
 88 translated sentence in the target language. We  
 89 have kept dimensionality of the embedded  
 90 layer as 256, batch size as 64 and used Adam  
 91 optimizer (Kingma and Ba, 2014). Due to  
 92 limited compute, we have used only 1 encoder  
 93 and decoder layer. Following (Srivastava,  
 94 2013), we used dropout after the RNN  
 95 decoder layer. By considering past and future  
 96 contexts of each word, bidirectional RNNs are  
 97 able to better capture contextual information,  
 98 leading to more accurate translations, as  
 99 compared to their vanilla counterpart. This  
 100 was the motivation behind using bidirectional  
 101 RNN.

102 Since RNNs process data in a sequential  
 103 manner, they can be slow. Another drawback  
 104 of RNN is that it is not able to retain long term  
 105 dependencies. So, another popular choice can  
 106 be to use RNN with Long-Short Term  
 107 Memory (LSTM) for the machine translation  
 108 tasks (Jozefowicz et al., 2016; Lample et  
 109 al., 2018). However, our dataset didn't  
 110 comprise of very long sentences, and due to  
 111 limited computing power, we decided not to  
 112 use LSTM.

### 113 2.2 Transformer

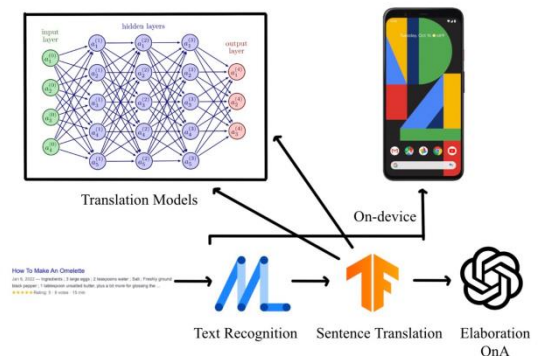
114 A transformer is a neural network architecture  
 115 that uses attention mechanism to process  
 116 sequential data. Unlike recurrent neural  
 117 networks (RNNs), transformers operate in  
 118 parallel, enabling more efficient computation.  
 119 We are choosing the Transformer model  
 120 (Vaswani et al., 2017) as it has shown to be  
 121 very effective for Machine Translation tasks  
 122 (Ott et al., 2018), including multilingual  
 123 machine translation tasks (Lakew et al., 2018;  
 124 Sachan and Neubig, 2018). This is due to their  
 125 attention mechanism, parallel processing, and  
 126 encoder-decoder architecture. The attention  
 127 mechanism allows them to capture

128 dependencies between words in a sentence,  
 129 facilitating the understanding of contextual  
 130 relationships. Parallel processing enables  
 131 efficient handling of long sequences, and the  
 132 encoder decoder architecture handles  
 133 variable-length input and output sequences.

134 For our experiment, we focused on the  
 135 Transformer in the "Base" configuration. We  
 136 refer the reader to Vaswani et al. (2017) to get  
 137 a better understanding of the model  
 138 architecture. Due to compute restrictions, we  
 139 have only used 1 encoder and 1 decoder block.  
 140 The encoder comprises of 8 attention heads, 1  
 141 attention layer, 2 add and normalization layers  
 142 and 2 fully connected layers with 2048 and  
 143 256 neurons. The decoder block comprises of  
 144 8 attention heads, 2 attention layers ( self and  
 145 cross attention). 3 add and normalization  
 146 layers and 2 fully connected layers with 2048  
 147 and 256 neurons. We have kept  
 148 dimensionality of the embedded layer as 256,  
 149 batch size as 64 and used Adam optimizer  
 150 (Kingma and Ba, 2014)

### 151 3 Workflow

152 The workflow of our project is mentioned in  
 153 this section with a detailed diagram in figure 1.



154  
 155

Figure 1: Project workflow

156 After training the selected model for an  
 157 adequate number of epochs, the TensorFlow  
 158 model is converted into a TensorFlow Lite  
 159 model using the following command:

```
tflite_convert --saved_model_dir=
/path/to/saved_model --output_file=
/path/to/output.tflite
```

160

161 The primary purpose of this conversion is to  
 162 leverage the advantages of TensorFlow Lite  
 163 models. Unlike TensorFlow models, these  
 164 lightweight models allow us to deploy high-  
 165 performance models on embedded systems or

166 mobile devices without significant degradation  
167 in performance or accuracy. An additional  
168 advantage is the availability of prebuilt and  
169 popular Android SDKs, such as MLKit, which  
170 natively support the loading of these lite  
171 models at runtime on any Android device. The  
172 converted model is then saved in the Android  
173 resource directory and loaded at runtime using  
174 the following code:

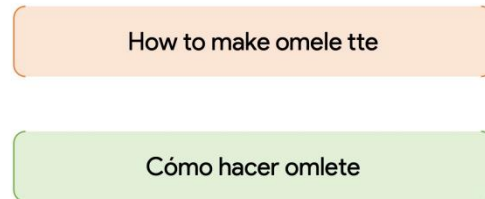
```
val localModel = LocalModel.Builder()  
    .setAssetFilePath(MODEL)  
    .build()
```

175  
176 In our implementation, we utilize two models:  
177 MLKit's native OCR model for text  
178 recognition and a custom-built translation  
179 model developed by our team. The Optical  
180 Character Recognition model is imported as a  
181 vision dependency using  
182 "com.google.ml.vision.DEPENDENCIES".  
183 We process the bitmap image through the OCR  
184 model and utilize convenient methods such as  
185 "onSuccessListener" and "onFailureListener"  
186 to obtain the desired results. The recognized  
187 text is extracted from the onSuccess method  
188 and passed into the TFLite translation model,  
189 which is loaded using MLKit. It is worth  
190 mentioning that the OCR model is downloaded  
191 during the initial run of the app and cached for  
192 subsequent runs. However, this may lead to  
193 race conditions if users attempt to translate text  
194 before the OCR model finishes loading.  
195 Once the translation model is loaded, we set the  
196 source and target languages as translator  
197 options and invoke the on-device model to  
198 obtain the translated text:

```
val options = TranslatorOptions.Builder()  
    .setSourceLanguage(TranslateLanguage.ENGLISH)  
    .setTargetLanguage(TranslateLanguage.SPANISH)  
    .build()
```

199  
200 In this specific scenario, we set the source  
201 language to English and the target language to  
202 Spanish. The translated text is obtained  
203 through the "onSuccess" method, similar to the  
204 callback methods used for the OCR model. We  
205 display the recognized and translated text on  
206 the user interface as soon as we receive it from  
207 the model.  
208 To provide further elaboration of the translated  
209 text, we make a POST call to OpenAI's  
210 completion API  
211 (<https://api.openai.com/v1/chat/completions>).  
212 It is important to obtain a bearer token from

213 OpenAI's developer dashboard for  
214 authenticating the API calls.  
215 Figure 2 displays the translation result: Notice  
216 how the language translation model gives the  
217 correct result despite a word gap issue with  
218 the recognized text.



219  
220 *Figure 2: Translation Result*  
221

#### 222 4 Result and Comparison

223 Initially, we used Precision, Recall and F-1  
224 score to evaluate our models. However, we  
225 were not getting satisfactory results. This is  
226 expected because these metrics on their own  
227 are not able to capture nuanced nature of  
228 translation quality. Precision and recall are  
229 metrics commonly used in tasks like  
230 information retrieval or binary classification,  
231 where the goal is to classify instances into  
232 categories. However, in machine translation,  
233 the output is a sequence of words or phrases,  
234 and a direct classification evaluation is not  
235 appropriate. F-1 score combines precision and  
236 recall into a single metric, which can be useful  
237 in certain tasks. However, it would still not be  
238 able to capture the complexities and subtleties  
239 of translation quality. So, we decided to use  
240 BLEU (Bilingual Evaluation Understudy) and  
241 ROUGE (RecallOriented Understudy for  
242 Gisting Evaluation) metrics as these will  
243 consider factors such as n-gram overlap and  
244 semantic similarity providing a more  
245 comprehensive assessment of translation  
246 quality. As we can see from the below table  
247 Bidirectional RNN Encoder Decoder  
248 outperforms Transformer model the main  
249 reason behind it is that the transformer model  
250 has just one encoder-decoder block and has  
251 been trained for insufficient number of epochs  
252 due to compute constraints.

253 Following are the BLEU and ROUGE scores  
254 for English to Spanish translation task.

	BLEU	ROUGE - 1	ROUGE - 2	ROUGE-L
<b>Bidirectional RNN</b>	0.619	0.742	0.503	0.729
<b>Transformer</b>	0.571	0.681	0.409	0.67

255 *Table 1: Comparing evaluation metrics*

256

257 Furthermore, we also successfully ported it on-  
 258 device and tested the performance of our model  
 259 on android, the screenshots attached below  
 260 testify our successful implementation of the  
 261 project.  
 262



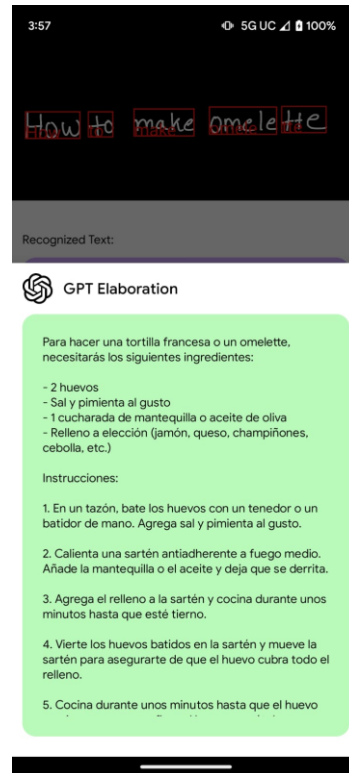
263

264 *Figure 3: Translation Result- On-device*

265

266 Figure 3 shows the result of processing and  
 267 extracting the text from the image using  
 268 Optical Character Recognition(OCR) model  
 269 and translation results using our custom model.  
 270 Both of them uses MLKit as SDK on android.  
 271

272 Also, Figure 4 below shows the results of the  
 273 remote API call made to OpenAI services  
 274 using the “gpt-3.5-turbo” model that gives the  
 275 context aware elaboration results.



276

277

278 *Figure 4: Elaboration Result- On-device*

278

## 279 5 Conclusion

280 In conclusion, this project successfully  
 281 developed a machine translation system that  
 282 operates on user devices and offers real-time  
 283 elaboration and Q&A capabilities. The RNN  
 284 model demonstrated superior performance  
 285 compared to the Transformer model, achieving  
 286 notable BLEU (0.617) and Rouge scores  
 287 across different metrics. The utilization of  
 288 TensorFlow Lite enabled efficient deployment  
 289 on edge devices, while on-device text  
 290 recognition facilitated instant translations from  
 291 English to Spanish. The incorporation of real-  
 292 time elaboration enhanced the user experience,  
 293 and the integration of OpenAI's completion  
 294 API provided additional information. Overall,  
 295 this project significantly contributes to  
 296 overcoming language barriers, with potential  
 297 avenues for further exploration and  
 298 optimization of models in the future.  
 299

## 300 6.Acknowledgments

301 We would like to thank Prof. Kai Wei Chan for  
 302 encouraging us to try new cutting edge projects  
 303 and helping us throughout the coursework. We  
 304 also like to thank the course TAs Tanmay,  
 305 Elaine and Masoud for helping us in every  
 306 stage of the project.

## 307 References

308 1. D. Bahdanau, K. Cho, and Y. Bengio.  
309 Neural Machine Translation by Jointly  
310 Learning to Align and Translate. 2016. arXiv  
311 preprint arXiv:1409.0473 [cs.CL].

312 2. Rafal Jozefowicz, Oriol Vinyals, Mike  
313 Schuster, Noam Shazeer, and Yonghui Wu.  
314 2016. Exploring the Limits of Language  
315 Modeling. In Proceedings of the 2016  
316 Conference on Empirical Methods in Natural  
317 Language Processing (EMNLP).

318 3. D. P. Kingma and J. Ba. "Adam: A method  
319 for stochastic optimization." 2014. arXiv  
320 preprint arXiv:1412.6980.

321 4. Surafel Melaku Lakew, Mauro Cettolo, and  
322 Marcello Federico. 2018. [A Comparison of  
323 Transformer and Recurrent Neural Networks  
324 on Multilingual Neural Machine Translation.](#)  
325 In *Proceedings of the 27th International  
326 Conference on Computational Linguistics*,  
327 pages 641–652, Santa Fe, New Mexico, USA.  
328 Association for Computational Linguistics.

329 5. Guillaume Lample, Myle Ott, Alexis  
330 Conneau, Ludovic Denoyer, and  
331 Marc'Aurelio Ranzato. 2018. Phrase-Based &  
332 Neural Unsupervised Machine Translation. In  
333 Proceedings of the 2018 Conference on  
334 Empirical Methods in Natural Language  
335 Processing (EMNLP).

336 6. Myle Ott, Sergey Edunov, David Grangier,  
337 and Michael Auli. 2018. [Scaling Neural  
338 Machine Translation.](#) In *Proceedings of the  
339 Third Conference on Machine Translation:  
340 Research Papers*, pages 1–9, Brussels,  
341 Belgium. Association for Computational  
342 Linguistics.

343 7. Devendra Sachan and Graham Neubig.  
344 2018. [Parameter Sharing Methods for  
345 Multilingual Self-Attentional Translation  
346 Models.](#) In *Proceedings of the Third  
347 Conference on Machine Translation:  
348 Research Papers*, pages 261–271, Brussels,

349 Belgium. Association for Computational  
350 Linguistics.

351 8. Nitish, Srivastava. 2013. Improving neural  
352 networks with dropout. PhD thesis,  
353 University of Toronto.

354 9. Zhen Yang, Wei Chen, Feng Wang, and Bo  
355 Xu. 2017. Improving Neural Machine  
356 Translation with Conditional Sequence  
357 Generative Adversarial Nets. In *Proceedings  
358 of the 2018 Conference of the North American  
359 Chapter of the Association for Computational  
360 Linguistics: Human Language Technologies,  
361 Volume 1 (Long Papers)*, pages 1346–1355,  
362 New Orleans, Louisiana. Association for  
363 Computational Linguistics.

364 10. A. Vaswani, N. Shazeer, N. Parmar, J.  
365 Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser,  
366 and I. Polosukhin. "Attention Is All You  
367 Need." 2017. arXiv preprint  
368 arXiv:1706.03762 [cs.CL].