

# Lab5: Translating by Prompting a LLM On the Importance of Parallel Data

He HUANG 22107447

November 23, 2023

## 1 Introduction

## 2 Your work

We will use the challenge set available <sup>1</sup> to evaluate the capacity of a LLM to translate (either one of Mistral LLM or LLaMA) and compare their performance to mBART.

**1. What is the BLEU score achieved on the challenge set by an LLM ? by mBART ?**  
Mistral LLM :

To make it easier to clean up the generated text, here I set up the LLM model without much explanation.

```
1 # Mistral LLM
2 # Use the LLM prompt method and pipeline to generate the translation
3 pipe = pipeline("text-generation",
4                 model="HuggingFaceH4/zephyr-7b-alpha",
5                 torch_dtype=torch.bfloat16,
6                 device_map="auto")
7
8 def mistral_translate(dataset):
9     translated = []
10
11     for i in range(0, len(dataset)):
12         # take the source sentence and add the prompt
13         message = [{"role": "user", "content": "translate into French and do not
14         add explanations or notes: " + dataset['source'][i]}]
15         # tokenize the prompt
16         prompt = pipe.tokenizer.apply_chat_template(message,
17                                                     tokenize=False,
18                                                     add_generation_prompt=True)
19         # generate the translation
20         outputs = pipe(prompt,
21                       max_new_tokens=256,
22                       do_sample=True,
23                       temperature=0.7,
24                       top_k=50,
25                       top_p=0.95)
```

---

<sup>1</sup><https://aclanthology.org/D17-1263/>

```

26         translated.append(outputs[0]["generated_text"])
27     return translated
28
29 # extract the content translated by Mistral
30 mistral_en_fr = []
31 for sent in translated_mistral:
32     sent = re.split("\n", sent)
33     if '<|assistant|>' in sent:
34         ids = sent.index('<|assistant|>')
35         mistral_en_fr.append(sent[ids+1])
36
37 # compute BLEU score for Mistral
38 mistral_results = bleu.compute(predictions=mistral_en_fr, references=corpus["
    reference"])
39 print(f"Bleu score of Mistral model : {mistral_results}")

```

mBART :

```

1 # define the mBart model and tokenizer
2 from transformers import MBartForConditionalGeneration, MBart50TokenizerFast
3
4 mbart_name = "facebook/mbart-large-50-many-to-many-mmt"
5 mbart_model = MBartForConditionalGeneration.from_pretrained(mbart_name).to(
    device)
6 mbart_tokenizer = MBart50TokenizerFast.from_pretrained(mbart_name, src_lang="
    en_XX")
7
8 # translate English sentences into French
9 def mbart_translate(en_source):
10     # encode the french sentence
11     inputs = mbart_tokenizer(en_source["source"], return_tensors="pt", truncation=
        True, max_length=512, padding=True).to(device)
12     # generate
13     generated_tokens = mbart_model.generate(**inputs, forced_bos_token_id=
        mbart_tokenizer.lang_code_to_id["fr_XX"])
14     # decode the english sentence
15     eng_sent = mbart_tokenizer.batch_decode(generated_tokens, skip_special_tokens=
        True)
16     return {"mbart_en_fr": eng_sent}
17
18
19 # Compute BLEU score
20 import evaluate
21 bleu = evaluate.load("bleu")
22
23 mbart_results = bleu.compute(predictions=corpus["mbart_en_fr"], references=
    corpus["reference"])
24 print(f"Bleu score of mBART model : {mbart_results}")

```

Results :

As can be seen from the results, there is a difference between the bleu score of the LLM model and the mBART translation model, although it is not particularly large.

```

1 Output:
2 Bleu score of mBART model : {'bleu': 0.5017724846364539, 'precisions':
    [0.7444561774023232, 0.5578069129916567, 0.4322845417236662,
    0.35313001605136435], 'brevity_penalty': 1.0, 'length_ratio':
    1.0021164021164022, 'translation_length': 947, 'reference_length': 945}
3
4 Output:

```

```

5 Bleu score of Mistral model : {'bleu': 0.45006544291275874, 'precisions':
    [0.7083333333333334, 0.5048309178743962, 0.3861111111111111,
    0.3088235294117647], 'brevity_penalty': 0.9904306953846911, 'length_ratio':
    0.9904761904761905, 'translation_length': 936, 'reference_length': 945}

```

2. Compute for each major category of difficulty the proportion of sentences in which the difficulty (identified by square brackets) is correctly translated ?

And here are all the results :

	category_major	GNMT_correct
0	Lexico-Syntactic	56.097561
1	Morpho-Syntactic	72.413793
2	Syntactic	73.684211

Figure 1: Accuracy of GNMT model

	category_major	PBMT1_correct
0	Lexico-Syntactic	39.024390
1	Morpho-Syntactic	17.241379
2	Syntactic	28.947368

Figure 2: Accuracy of PBMT1 model

	category_major	PBMT2_correct
0	Lexico-Syntactic	43.902439
1	Morpho-Syntactic	17.241379
2	Syntactic	26.315789

Figure 3: Accuracy of PBMT2 model

	<code>category_major</code>	<code>NMT_correct</code>
<b>0</b>	Lexico-Syntactic	46.341463
<b>1</b>	Morpho-Syntactic	75.862069
<b>2</b>	Syntactic	34.210526

Figure 4: Accuracy of NMT model