# Multilingual NLP — Lab 3
# How Reliable is MT Evaluation?

Lina Conti

January 2023

## 1 Re-evaluating the role of Bleu in machine translation research

```python
from math import exp, log
from collections import Counter
from nltk.tokenize import word_tokenize
import nltk
nltk.download('punkt')


def brevity_penalty(r, h):
    if h > r:
        return 1.
    else:
        return exp( 1 - (r/h) )


def ngrams(sequence, n):
    if len(sequence) < n: # no ngrams
      return

    sequence = iter(sequence)
    history = []
    while n > 1:
        history.append(next(sequence))
        n -= 1
    for item in sequence:
        history.append(item)
        yield tuple(history)
        del history[0]


def ngram_matches(ref, hyp, n):
    bag_ref = Counter(ngrams(ref, n))
    bag_hyp = Counter(ngrams(hyp, n))

    common = sum((min(bag_hyp[ngram], bag_ref[ngram]) for ngram in bag_hyp))
    total = sum(bag_hyp.values())
    return common, total


def my_BLEU(references, hypotheses, max_n):
    r = 0 # number of tokens in the reference corpus
    h = 0 # number of tokens in the hypotheses corpus
    matched_ngrams = [0] * max_n  # number of matched ngrams for each n
    total_ngrams = [0] * max_n  # number of ngrams in hypotheses corpus for each n

    for ref, hyp in zip(references, hypotheses):
        hyp = word_tokenize(hyp)
        ref = word_tokenize(ref)

        r += len(ref)
        h += len(hyp)

        for n in range(max_n):
```

```
52              common, total = ngram_matches(ref, hyp, n+1)
53              matched_ngrams[n] += common
54              total_ngrams[n] += total
55
56      precisions = []
57      for n in range(max_n):
58          precisions.append(matched_ngrams[n]/total_ngrams[n])
59      log_precisions = [log(pn) for pn in precisions]
60
61      avg_precision = sum(log_precisions)/len(log_precisions)
62      bp = brevity_penalty(r, h)
63      bleu = bp * exp(avg_precision)
64      return {'bleu': round(bleu,4),
65              'precisions': [round(pn, 3) for pn in precisions],
66              'bp': round(bp,3),
67              'ratio': round(r/h,3),
68              'hyp_len': h,
69              'ref_len': r}
```

```
1  import re
2
3  def sgm_to_txt(sgm_file, txt_file):
4      with open(sgm_file, 'r') as f_in:
5          with open(txt_file, 'w') as f_out:
6              for line in f_in:
7                  if line.startswith("<seg id"):
8                      f_out.write(re.sub('(<seg id="\d*">)|(<\/seg>)', "", line))
9
10  sgm_to_txt("/content/test/newsdiscusstest2015-enfr-src.en.sgm", "
       newsdiscusstest2015-enfr-src.en")
11  sgm_to_txt("/content/test/newsdiscusstest2015-enfr-ref.fr.sgm", "
       newsdiscusstest2015-enfr-ref.fr")
```

```
1  from tqdm.notebook import tqdm
2
3  def generate_batches(src_sentences, tokenizer, batch_size, device):
4      bstart = 0
5      N = len(src_sentences)
6      while bstart < N:
7          bend = min(bstart+batch_size, N)
8          b_src_sentences = src_sentences[bstart:bend]
9          b_encoded_src = tokenizer(b_src_sentences, return_tensors="pt", padding=
       True, truncation=True, max_length=512)
10          yield b_encoded_src.to(device)
11          bstart += batch_size
12
13
14  def translate(src_file, out_file, tokenizer, model, batch_size, device):
15      with open(src_file, 'r') as f:
16          src_sentences = f.readlines()
17
18      nb_batches = (len(src_sentences)/batch_size)
19      for b_encoded_src in tqdm(generate_batches(src_sentences, tokenizer,
       batch_size, device), total=nb_batches):
20
21          print("torch.cuda.memory_reserved: %fGB"%(torch.cuda.memory_reserved(0)
       /1024/1024/1024))
22
23          try:  # MBart
24              generated_tokens = model.generate(**b_encoded_src,  max_new_tokens
       =512, forced_bos_token_id=tokenizer.lang_code_to_id["fr_XX"])
25          except AttributeError:  # MarianMT
26              generated_tokens = model.generate(**b_encoded_src,  max_new_tokens
       =512)
27
28          b_decoded_trg = tokenizer.batch_decode(generated_tokens,
       skip_special_tokens=True)
29
30          with open(out_file, 'a') as f:
31              for line in b_decoded_trg:
32                  f.write(f"{line}\n")
```

```
1  from transformers import MarianTokenizer, MarianMTModel
2
3  model_name = "Helsinki-NLP/opus-mt-en-fr"
4  marian_model = MarianMTModel.from_pretrained(model_name).to(DEVICE)
5  marian_tokenizer = MarianTokenizer.from_pretrained(model_name)
6
7  translate('newsdiscusstest2015-enfr-src.en',
8            'marian-newsdiscusstest2015-enfr-pred.fr',
9            marian_tokenizer,
10           marian_model,
11           batch_size=128,
12           device=DEVICE)
```

```
1  from transformers import MBartForConditionalGeneration, MBart50TokenizerFast
2
3  mbart_model = MBartForConditionalGeneration.from_pretrained("facebook/mbart-large
       -50-many-to-many-mmt").to(DEVICE)
4  mbart_tokenizer = MBart50TokenizerFast.from_pretrained("facebook/mbart-large-50-
       many-to-many-mmt")
5  mbart_tokenizer.src_lang = "en_XX"
6
7  translate('newsdiscusstest2015-enfr-src.en',
8            'mbart-newsdiscusstest2015-enfr-pred.fr',
9            mbart_tokenizer,
10           mbart_model,
11           batch_size=16,
12           device=DEVICE)
```

```
1  with open('newsdiscusstest2015-enfr-ref.fr', 'r') as f:
2      references = f.readlines()
3
4  with open('marian-newsdiscusstest2015-enfr-pred.fr', 'r') as f:
5      marian_hypotheses = f.readlines()
6
7  with open('mbart-newsdiscusstest2015-enfr-pred.fr', 'r') as f:
8      mbart_hypotheses = f.readlines()
9
10 print('MarianMT Bleu:')
11 print(my_BLEU(references, marian_hypotheses, 4))
12
13 print('\nMBart Bleu:')
14 print(my_BLEU(references, mbart_hypotheses, 4))
```

With my implementation, MarianMT obtains a Bleu score of 0.3747 and MBart of 0.3231. It looks like MarianMT's model is better at translating from English to French then MBart's.

3. Let's take the sentence 'I think it 's called democracy .' The reference translation is 'Je crois que c'est ce qu'on appelle la démocratie .' and MBart's hypothesis is 'Je pense qu'il s'agit de la démocratie .'

Bigram mismatches are indicated with pipe symbols: 'Je | pense | qu'il | s'agit | de | la démocratie .' An example of permutation that does not lower the BLEU score is 'Je | la démocratie . | qu'il | s'agit | pense | de'.

This happens because the number of matching unigrams does not change with permutation, the number of matching ngrams obviously does not change when permutations happen on bigram mismatches and neither does the number of higher order ngrams, since the latter are composed of matching bigrams. For instance, in the example above, the trigram 'la démocratie .' will never be separated by such permutations. The length of the hypotheses will not change either, so all the building blocks of the BLEU score stay the same with these permutations.

4. Given a sentence with n words and b bigram mismatches, there are $b' = n - b + 1$ bigram matches, so $(n - b')! = (b + 1)!$ sentences can be generated with this principle (cf. Callison-Burch et al. 2006).

```
1  from math import factorial
2
3  def possible_permutations(references, hypotheses):
```

```
4
5      permutations = 0
6      for ref, hyp in zip(references, hypotheses):
7          hyp = word_tokenize(hyp)
8          ref = word_tokenize(ref)
9
10         n = len(hyp)
11         b, _ = ngram_matches(ref, hyp, 2)
12         permutations += factorial(n-b)
13
14     return f"{permutations:.2e}"
15
16 possible_permutations(references, mbart_hypotheses)
```

On the full WMT'15 test set, there are 9.90e+163 permutations for MarianMT's translations and 3.01e+213 for MBart's.

5. This questions the use of BLEU as an evaluation metric because it is unlikely that that many variations of each translation hypothesis are all really of equal quality. This would mean that BLEU cannot correlate well with human judgement for all hypotheses and that a higher BLEU score does not necessarily mean better translations (since a high score can also be achieve by permutated gibberish).

6.
```
1 from sacrebleu.metrics import BLEU
2 bleu = BLEU()
3
4 print('MarianMT sacrebleu:')
5 print(bleu.corpus_score(marian_hypotheses, [references]))
6
7 print('\nMBart sacrebleu:')
8 print(bleu.corpus_score(mbart_hypotheses, [references]))
```

With sacrebleu, MarianMT gets a score of 38.36 and MBart of 32.58. With my implementation the scores were respectively 37.47 and 32.31. The scores are not too different, however, they are not the same. For MarianMT especially, the difference is of almost one BLEU point. This means that one should be careful when comparing BLEU scores for different models that were computed by different teams for different papers, for example. If the implementation of BLEU is not the same the results may not be comparable.

7. Raw text:
```
1 print('MarianMT Bleu:')
2 print(my_BLEU(references, marian_hypotheses, 4))
3
4 print('\nMBart Bleu:')
5 print(my_BLEU(references, mbart_hypotheses, 4))
6
7 print('\nMarianMT sacrebleu:')
8 print(bleu.corpus_score(marian_hypotheses, [references]))
9
10 print('\nMBart sacrebleu:')
11 print(bleu.corpus_score(mbart_hypotheses, [references]))
```

Output:
```
1 MarianMT Bleu:
2 {'bleu': 0.3747, 'precisions': [0.655, 0.44, 0.316, 0.234], 'bp': 0.981, 'ratio':
      1.019, 'hyp_len': 28047, 'ref_len': 28589}
3
4 MBart Bleu:
5 {'bleu': 0.3231, 'precisions': [0.606, 0.381, 0.261, 0.182], 'bp': 0.999, 'ratio':
      1.001, 'hyp_len': 28550, 'ref_len': 28589}
6
7 MarianMT sacrebleu:
8 BLEU = 38.36 65.6/44.2/31.8/23.5 (BP = 1.000 ratio = 1.005 hyp_len = 28115 ref_len
      = 27975)
9
10 MBart sacrebleu:
11 BLEU = 32.58 60.7/38.4/26.3/18.4 (BP = 1.000 ratio = 1.021 hyp_len = 28575 ref_len
      = 27975)
```

Subword tokenization:

```
1  def subword_tokenization(tokenizer, in_file, out_file):
2      with open(in_file, 'r') as f_in:
3          with open(out_file, 'w') as f_out:
4              for line in f_in:
5                  f_out.write(' '.join(tokenizer.tokenize(line.strip())) + '\n')
6
7  subword_tokenization(mbart_tokenizer, 'marian-newsdiscusstest2015-enfr-pred.fr', '
       subtok-marian-newsdiscusstest2015-enfr-pred.fr')
8  subword_tokenization(mbart_tokenizer, 'mbart-newsdiscusstest2015-enfr-pred.fr', '
       subtok-mbart-newsdiscusstest2015-enfr-pred.fr')
9
10 with open('subtok-marian-newsdiscusstest2015-enfr-pred.fr', 'r') as f:
11     marian_hypotheses = f.readlines()
12
13 with open('subtok-mbart-newsdiscusstest2015-enfr-pred.fr', 'r') as f:
14     mbart_hypotheses = f.readlines()
15
16 print('MarianMT subword Bleu:')
17 print(my_BLEU(references, marian_hypotheses, 4))
18
19 print('\nMBart subword Bleu:')
20 print(my_BLEU(references, mbart_hypotheses, 4))
21
22 print('\nMarianMT subword sacrebleu:')
23 print(bleu.corpus_score(marian_hypotheses, [references]))
24
25 print('\nMBart subword sacrebleu:')
26 print(bleu.corpus_score(mbart_hypotheses, [references]))
```

Output:

```
1  MarianMT subword Bleu:
2  {'bleu': 0.0029, 'precisions': [0.089, 0.006, 0.001, 0.0], 'bp': 1.0, 'ratio':
       0.727, 'hyp_len': 39315, 'ref_len': 28589}
3
4  MBart subword Bleu:
5  {'bleu': 0.003, 'precisions': [0.087, 0.006, 0.001, 0.0], 'bp': 1.0, 'ratio':
       0.712, 'hyp_len': 40143, 'ref_len': 28589}
6
7  MarianMT subword sacrebleu:
8  BLEU = 0.73 8.6/0.8/0.3/0.1 (BP = 1.000 ratio = 1.410 hyp_len = 39447 ref_len =
       27975)
9
10 MBart subword sacrebleu:
11 BLEU = 0.74 8.3/0.8/0.3/0.1 (BP = 1.000 ratio = 1.437 hyp_len = 40212 ref_len =
       27975)
```

Character tokenization:

```
1  def char_tokenization(in_file, out_file):
2      with open(in_file, 'r') as f_in:
3          with open(out_file, 'w') as f_out:
4              for line in f_in:
5                  f_out.write(' '.join(line))
6
7  char_tokenization('marian-newsdiscusstest2015-enfr-pred.fr', 'chartok-marian-
       newsdiscusstest2015-enfr-pred.fr')
8  char_tokenization('mbart-newsdiscusstest2015-enfr-pred.fr', 'chartok-mbart-
       newsdiscusstest2015-enfr-pred.fr')
9
10 with open('chartok-marian-newsdiscusstest2015-enfr-pred.fr', 'r') as f:
11     marian_hypotheses = f.readlines()
12
13 with open('chartok-mbart-newsdiscusstest2015-enfr-pred.fr', 'r') as f:
14     mbart_hypotheses = f.readlines()
15
16 print('MarianMT character Bleu:')
17 print(my_BLEU(references, marian_hypotheses, 4))
18
19 print('\nMBart character Bleu:')
```

```
20  print(my_BLEU(references, mbart_hypotheses, 4))
21
22  print('\nMarianMT character sacrebleu:')
23  print(bleu.corpus_score(marian_hypotheses, [references]))
24
25  print('\nMBart character sacrebleu:')
26  print(bleu.corpus_score(mbart_hypotheses, [references]))
```

Output:

```
1   MarianMT character Bleu:
2   {'bleu': 0.0016, 'precisions': [0.03, 0.002, 0.001, 0.0], 'bp': 1.0, 'ratio':
        0.228, 'hyp_len': 125250, 'ref_len': 28589}
3
4   MBart character Bleu:
5   {'bleu': 0.0016, 'precisions': [0.03, 0.002, 0.001, 0.0], 'bp': 1.0, 'ratio':
        0.223, 'hyp_len': 128008, 'ref_len': 28589}
6
7   MarianMT character sacrebleu:
8   BLEU = 0.24 2.8/0.3/0.1/0.0 (BP = 1.000 ratio = 4.477 hyp_len = 125250 ref_len =
        27975)
9
10  MBart character sacrebleu:
11  BLEU = 0.24 2.8/0.3/0.1/0.0 (BP = 1.000 ratio = 4.576 hyp_len = 128008 ref_len =
        27975)
```

The scores for tokenized text are oddly low (lower than 1.) I find this hard to explain, but given that it is the same both with sacrebleu and my implementation, I don't think it is a problem with how I am using either of those. The tokenized text files also look normal. Maybe it is due to the fact that there are a lot more tokens this way, so mistakes in the translations are amplified, but I am not sure.

Anyway, the moral is that BLEU scores change according to the implementation of BLEU **and** according to the tokenization of the texts, so one should be wary of any difference in those regards when comparing scores.