

# Lab 2: How Multilingual is Multilingual BERT?

## Multilingual NLP

Manon Fleury  
22410952

## 2. Fine-Tuning mBERT

```
1 !curl -O https://raw.githubusercontent.com/UniversalDependencies/UD_French
   -Sequoia/master/fr_sequoia-ud-test.conllu
2
3 import conllu
4 def load_conllu(filename):
5     for sentence in conllu.parse(open(filename, "rt", encoding="utf-8").read
      ()):
6         tokenized_words = [token["form"] for token in sentence]
7         gold_tags = [token["upos"] for token in sentence]
8         yield tokenized_words, gold_tags
9 corpus = list(load_conllu("fr_sequoia-ud-test.conllu")) # list of sets of
   two lists sentence and gold_tags
10
11 len(corpus) # 456 sentences
```

## 1 The Universal Dependencies Project

### 1. Why do we need “consistent annotation” for our experiment?

We need consistent annotation because it allows a fair comparison across languages. Indeed, with multilingual models, standardized grammatical categories ensure that the annotations (PoS, morphological features, syntactic dependencies) are labeled the same way, regardless of the language. This consistency is key for evaluating how well the model transfers knowledge from one language to another.

### 2. What do we use the yield keyword rather than a simple return (line 7) ? Why do we have to call the list constructor (line 9)?

The yield keyword produces a series of values over time, rather than computing them all at once and returning them in a list. This is particularly useful for processing large datasets, as it allows for more efficient memory usage since it doesn't load everything into memory at once, it generates data on demand. And if we only need a few sentences at a time for processing, yield allows to avoid loading the entire file.

The list constructor is called to convert the generator returned by load\_conllu into a list, to store all the parsed sentences and their PoS tags in memory for further analysis.

### 3. What is the distribution of labels in the test set of the Sequoia treebank?

```

1
2 from collections import Counter
3 all\_tags = [tag for \_, tags in corpus for tag in tags] # list of all the
4               tags (one after the other, so with repetition)
5 label\_distribution = Counter(all\_tags)
6 for label, count in label\_distribution.items():
7     print(f"{label}: {count}")

```

Output:

```

PRON: 410
VERB: 781
SCONJ: 106
DET: 1492
NOUN: 2161
ADJ: 638
PUNCT: 1084
ADV: 411
NUM: 243
ADP: 1634
PROP: 478
CCONJ: 221
AUX: 345
.: 310
X: 36
SYM: 4

```

#### 4. What are the multiword tokens in the test set of Sequoia? How are they annotated?

When we look at the Sequoia conllu, we see that the multiword token have a range of numbers (like 15-16) to indicate it covers multiple tokens and the next line explicits the multiple token with the first range number corresponding to the first token extracted from the multiple token, and so on.

For example:

```

15-16    au    -    -
15       à    ADP   -
16       le    DET   -

```

```

1 def extract_multiword_tokens(filename):
2     multiword_tokens = []
3     for sentence in conllu.parse(open(filename, "rt", encoding="utf-8").read
4                                   ()):
5         for token in sentence:
6             if isinstance(token["id"], tuple): # (like '15-16')
7                 token = token["form"].lower()
8                 if token not in multiword_tokens:
9                     multiword_tokens.append(token)
10    return multiword_tokens
11 print(extract_multiword_tokens("fr_sequoia-ud-test.conllu"))

```

Output:

```
['des', 'du', 'aux', 'au', 'desdites']
```

This is the list of the multiple tokens occurring in the conllu file, with no repetition.

## 5. Why did the UD project made the decision of splitting multiword tokens into several (grammatical) words? What is your opinion on this decision?

The Universal Dependencies (UD) project decided to split multiword tokens into separate grammatical words to maintain consistency across different languages and enhance the syntactic representation of each word's role. This approach simplifies tokenization, allowing for clearer syntactic dependencies and better handling of grammatical phenomena like agreement and case marking. By treating each component as an individual token, models can learn relationships more effectively, potentially improving performance in multilingual contexts. However, this decision may complicate the treatment of idiomatic expressions that might be better understood as single units.

## 6. According to UD tokenization guidelines, a token can contain spaces. Are there any token in the Sequoia corpora that contain spaces ? If so these spaces should be removed.

```
1 def tokens_space(filename):
2     tokens_with_space = []
3     for sentence in conllu.parse(open(filename, "rt", encoding="utf-8").
4         read()):
5         for token in sentence:
6             if " " in token["form"]:
7                 tokens_with_space.append(token["form"])
8     return tokens_with_space
9
10 tokens_with_space = tokens_space("fr_sequoia-ud-test.conllu")
11
12 if tokens_with_space:
13     for token in tokens_with_space:
14         print(token)
15
16 # Function to remove the spaces
17 def load_and_clean_conllu(filename):
18     for sentence in conllu.parse(open(filename, "rt", encoding="utf-8").read
19         ()):
20         tokenized_words = [token["form"].replace(" ", "") for token in
21             sentence]
22         gold_tags = [token["upos"] for token in sentence]
23         yield tokenized_words, gold_tags
24
25 cleaned_corpus = list(load_and_clean_conllu("fr_sequoia-ud-test.conllu"))
26
27 "print(token)" displayed numbers like "500 000", "3 862", that represent the tokens containing a space.
28 So with this function, these numbers become "500000".
```

## 2 mBERT tokenization

### 7. Why does mBERT use a subword tokenization?

mBERT uses subword tokenization to effectively handle out-of-vocabulary words by breaking them down into smaller, meaningful components, allowing the model to process even rare or infrequent terms. This approach also reduces vocabulary size while capturing the rich morphological structures found in various languages, enhancing generalization across tasks and languages.

8. How is the sentence “Pouvez-vous donner les mêmes garanties au sein de l’Union Européene” 4 tokenized according to the UD convention? How is it tokenized by mBERT tokenizer?

According to the UD convention, this sentence is tokenized as:

```
["Pouvez",  
 "-vous",  
 "donner",  
 "les",  
 "mêmes",  
 "garanties",  
 "au",  
 "sein",  
 "de",  
 "l",  
 "'",  
 "Union",  
 "Européenne".]
```

While mBERT tokenizer:

```
1 from transformers import AutoTokenizer  
2 tokenizer = AutoTokenizer.from_pretrained("bert-base-multilingual-cased")  
3 tokenizer.tokenize("Pouvez-vous donner les m mes garanties au sein de  
    l Union Europ enne")
```

Output:

```
['Po',  
 '##uve',  
 '##z',  
 '-',  
 'vous',  
 'donner',  
 'les',  
 'mêmes',  
 'gara',  
 '##nties',  
 'au',  
 'sein',  
 'de',  
 'l',  
 '[UNK]',  
 'Union',  
 'Euro',  
 '##pée',  
 '##nne']
```

We see that mBERT tokenizer splits infrequent words in smaller factors and the subword tokens are prefixed with ##.

9. Why is this difference in tokenization a problem for training a mBERT-based PoS tagger?

This mismatch means that the model may not directly learn the PoS tags that correspond to the tokens used in the training data, leading to inconsistencies in predictions.

Additionally, when the mBERT tokenizer generates [UNK] (unknown) tokens, it creates a gap in the training data for PoS tagging, that can degrade the model's ability to make correct predictions for those words.

### 3 Reconciling the Two Tokenizations

10. Write a function that takes as parameter a sentence tokenized according to UD rules (i.e. a variable of type List[str] and its gold labels (also a variable of type List[str]) and implements the first principle explained above.

```

1     def update_ud_tokenization(tokens, labels):
2     '''tokens = List[str], sentence tokenized according to the UD rules
3     labels = List[str] of the labels corresponding to the tokens
4     Returns UD updated tokenization and UD updated labels'''
5     new_tokens = []
6     new_labels = []
7     i = 0
8     while i < len(tokens):
9         token = tokens[i]
10        label = labels[i]
11        if label == '_' and i+2 < len(tokens): # i+2: we must have 2 more
12            tokens after
13            multiword_token = token
14            concatene_label = labels[i+1] + '+' + labels[i+2]
15            new_tokens.append(token)
16            new_labels.append(concatene_label)
17            i += 2
18        else: # if the token is not a multiword token
19            new_tokens.append(token)
20            new_labels.append(label)
21            i += 1
22        return new_tokens, new_labels
23
24    #test
25    print(update_ud_tokenization(['aller', 'au', 'a', 'le', 'marche'], ['VERB', '_',
26        'ADP', 'DET', 'NOUN']))
27
28    ([ 'aller', 'au', 'marché'], ['VERB', 'ADP+DET', 'NOUN'])

```

11. Write a function that takes as input the sentences and labels created by “normalizing” UD sentences, apply the mBERT tokenizer and compute the corresponding label. You must ensure each subtoken (including padding symbols) has a label.

```

1     def labels_mBERT(sentences, labels):
2     '''sentences : list of list of words by UD tokenization
3     labels : list of list of labels by UD original labels'''
4     L_mBERT_labels = [] # to return
5     for sentence, label in zip(sentences, labels):
6         # HERE we process ONE sentence and the labels corresponding (the '
7         label' list)
8
9         # 1. Apply the mBERT tokenizer
10        info_sentence = tokenizer(sentence, is_split_into_words=True,

```

```

10         return_offsets_mapping=True,
11         padding=True,
12         truncation=True)
13     L_offset_mapping = info_sentence['offset_mapping'] # list of tuples of
14         offmap for this sentence, with two more tuple (0,0) one at the
15         beginning and one at the end
16     # We get rid of the first element and the last element of the list
17         ((0,0)), unrelevant for our processing
18     L_offset_mapping.pop(0)
19     L_offset_mapping.pop(-1)
20
21     # 2. Compute the corresponding label
22     for i,couple in enumerate(L_offset_mapping):
23         if couple[0] != 0: # <pad>
24             label.insert(i, '<pad>')
25     L_mBERT_labels.append(label) # we add the list of labels for one
26         sentence in the broader list labels containing all the labels
27     return L_mBERT_labels
28
29 # test
30 print(labels_mBERT([[ 'Revenons', '-en', 'aux', 'choses', 'essentielles', '
31     .'], [ 'Parcourir', "l'", 'itineraire', '!' ]], [[ 'VERB', 'PRON', 'ADP+DET',
32     'NOUN', 'ADJ', 'PUNCT'], [ 'VERB', 'DET', 'NOUN', 'PUNCT' ]]))
33
34 [[ 'VERB', '<pad>', '<pad>', 'PRON', '<pad>', 'ADP+DET', 'NOUN', 'ADJ', '<pad>', 'PUNCT'], [ 'VERB',
35     '<pad>', '<pad>', 'DET', '<pad>', 'NOUN', '<pad>', '<pad>', 'PUNCT']]

```

## 12. Write a function that encodes the labels into integers.

```

1 def encode_labels(labels):
2     '''labels: [list of [list of labels for one sentence]]'''
3
4     # GENERATE DICT OF UNIQUE LABEL TO UNIQUE INTEGER
5     unique_labels = set(label for sublist in labels for label in sublist)
6     label_dict = {label: (-100 if label == "<pad>" else idx + 1) for idx,
7         label in enumerate(unique_labels)}
8
9     # ENCORE LABELS INTO INTEGERS
10    list_integers_of_all_label_lists = []
11    for L_label in labels:
12        list_integers_of_one_label_list = []
13        for label in L_label:
14            list_integers_of_one_label_list.append(label_dict[label])
15        list_integers_of_all_label_lists.append(
16            list_integers_of_one_label_list)
17    return list_integers_of_all_label_lists, label_dict
18
19 # test
20 print(encode_labels([[ 'VERB', '<pad>', '<pad>', 'PRON', '<pad>', 'ADP+DET',
21     'NOUN', 'ADJ', '<pad>', 'PUNCT'],
22     [ 'VERB', '<pad>', '<pad>', 'DET', '<pad>', 'NOUN', '<pad>', '<pad>',
23     'PUNCT' ]]))
24
25 ([[6, -100, -100, 3, -100, 4, 2, 5, -100, 1], [6, -100, -100, 8, -100, 2, -100, -100, 1]],
26     'PUNCT': 1, 'NOUN': 2, 'PRON': 3, 'ADP+DET': 4, 'ADJ': 5, 'VERB': 6, '<pad>': -100, 'DET':

```

8)

## 4 Creating a Dataset

13. Using the `Dataset.from_list` method, write a method that creates a `Dataset` that encapsulates a corpus in the conllu.

```
1 from datasets import Dataset
2
3 def create_dataset(sentences, labels):
4     '''sentences : list of list of words by UD tokenization
5     labels : list of list of labels by encoded_labels'''
6     tokenized_sentences = tokenizer(sentences,
7                                     is_split_into_words=True,
8                                     return_offsets_mapping=False,
9                                     padding=True, truncation=True)
10    max_length = max(len(ids) for ids in tokenized_sentences['input_ids'])
11
12    padded_labels = []
13    for label in labels:
14        padded_label = label + [-100] * (max_length - len(label))
15        padded_labels.append(padded_label)
16
17    dataset_entries = []
18    for i in range(len(sentences)):
19        dataset_entries.append({
20            'input_ids': tokenized_sentences['input_ids'][i],
21            'attention_mask': tokenized_sentences['attention_mask'][i],
22            'labels': padded_labels[i]
23        })
24
25    dataset = Dataset.from_list(dataset_entries)
26    return dataset
```

14. Create three instances of `Dataset` : one for the train set, one for the dev set and the last one for the test set.

```
1 from datasets import Dataset
2 import random
3 import numpy as np
4
5 def split_dataset(dataset, train_ratio=0.8, dev_ratio=0.1, test_ratio=0.1)
6     :
7     total_samples = len(dataset)
8     train_size = int(train_ratio * total_samples)
9     dev_size = int(dev_ratio * total_samples)
10    test_size = total_samples - train_size - dev_size
11
12    indices = list(range(total_samples))
13    random.shuffle(indices)
14
15    train_indices = indices[:train_size]
```

```

15     dev_indices = indices[train_size:train_size + dev_size]
16     test_indices = indices[train_size + dev_size:]
17
18     train_set = dataset.select(train_indices)
19     dev_set = dataset.select(dev_indices)
20     test_set = dataset.select(test_indices)
21
22     return train_set, dev_set, test_set

```

Main program:

```

1 unique_labels_list = []
2 L_tokens = []
3 L_labels = []
4 for one_set in corpus: # corpus = list of sets of 2 lists (sentence,
    labels)
5     # one_set = ([list tokens], [list labels])
6     tokens = one_set[0]
7     labels = one_set[1]
8     updated_tokens, updated_labels = update_ud_tokenization(tokens, labels)
9     L_tokens.append(updated_tokens) # list of lists of sentence tokens # 456
    lists
10    L_labels.append(updated_labels) # list of lists of sentence labels # 456
    lists
11
12    aligned_labels = labels_mBERT(L_tokens, L_labels) # <pad>
13    encoded_labels, dico_label = encode_labels(aligned_labels) # integers
14    print('dico_label:', dico_label)
15
16    dataset = create_dataset(L_tokens, encoded_labels)
17
18    train_set, dev_set, test_set = split_dataset(dataset)
19    print(train_set)
20    print(dev_set)
21    print(test_set)
22
23    unique_labels = set(label for sublist in aligned_labels for label in
    sublist)
24    len(unique_labels) # 17 unique labels, useful later for the model

```

Output:

```

dico_label: {'ADV': 1, 'ADP': 2, 'PUNCT': 3, 'CCONJ': 4, 'X': 5, 'NUM': 6, 'NOUN': 7,
'PRON': 8, 'ADP+DET': 9, 'SYM': 10, 'SCONJ': 11, 'PROPN': 12, 'ADJ': 13, 'VERB': 14,
'<pad>': -100, 'DET': 16, 'AUX': 17}

```

```

Dataset({
  features: ['input_ids', 'attention_mask', 'labels'],
  num_rows: 364
})
Dataset({
  features: ['input_ids', 'attention_mask', 'labels'],
  num_rows: 45
})
Dataset({

```



```

        features: ['input_ids', 'attention_mask', 'labels'],
        num_rows: 47
    })

```

## 5 Fine-Tuning mBERT

15. How can you modify the code of Figure 4 to report the PoS tagging accuracy during optimization. Why is this information important?

```

1  from evaluate import load
2  import numpy as np
3  from transformers import AutoModelForTokenClassification,
   TrainingArguments, Trainer
4
5  base_model_checkpoint = "bert-base-multilingual-cased"
6  model = AutoModelForTokenClassification.from_pretrained(
   base_model_checkpoint, num_labels=len(list(unique_labels))+1)
7
8  metric = load("accuracy")
9
10 def compute_metrics(pred):
11     logits, labels = pred.predictions, pred.label_ids
12     predictions = np.argmax(logits, axis=-1)
13     # Mask of the -100 labels to ignore the padding labels
14     true_labels = labels[labels != -100]
15     true_predictions = predictions[labels != -100]
16     accuracy = metric.compute(predictions=true_predictions, references=
   true_labels)
17     return {"accuracy": accuracy["accuracy"]}
18
19
20 training_args = TrainingArguments(
21     output_dir="./mBERT_fine_tuning",
22     eval_strategy="epoch", # evaluate at each epoch
23     learning_rate=2e-5,
24     per_device_train_batch_size=16,
25     per_device_eval_batch_size=16,
26     num_train_epochs=3,
27     weight_decay=0.01,
28     logging_dir='./logs',
29     logging_steps=10,
30     report_to="none" #Deactivate the integration with 'wandb'
31 )
32
33 trainer = Trainer(
34     model=model,
35     args=training_args,
36     train_dataset= train_set,
37     eval_dataset= dev_set,
38     compute_metrics=compute_metrics
39 )
40
41 trainer.train()

```

Epoch	Training Loss	Validation Loss	Accuracy
1	2.204400	1.821981	0.506024
2	1.598800	1.264914	0.653614
3	1.254000	1.050898	0.707831

Figure 1: Fine-tuning results for French

## 6 Evaluating the multilingual capacity of mBERT

16. Choose 5 languages from the UD project. For each language train a PoS tagger using the code of the previous section and test it on the 5 languages you have chosen.

```

1 corpus_breton = list(load_conllu("br_keb-ud-test.conllu"))
2 corpus_chinese = list(load_conllu("zh_hk-ud-test.conllu"))
3 corpus_russian = list(load_conllu("ru_pud-ud-test.conllu"))
4 corpus_french = corpus
5
6 def preprocess_corpus(corpus):
7     unique_labels_list = []
8     L_tokens = []
9     L_labels = []
10    for one_set in corpus:
11        tokens = one_set[0]
12        labels = one_set[1]
13        updated_tokens, updated_labels = update_ud_tokenization(tokens, labels)
14        L_tokens.append(updated_tokens)
15        L_labels.append(updated_labels)
16
17    aligned_labels = labels_mBERT(L_tokens, L_labels)
18    encoded_labels, dico_label = encode_labels(aligned_labels)
19
20    dataset = create_dataset(L_tokens, encoded_labels)
21    train_set, dev_set, test_set = split_dataset(dataset)
22
23    unique_LABELS = set(label for sublist in aligned_labels for label in
24                        sublist)
25
26    return train_set, dev_set, test_set, unique_LABELS
27
28 from evaluate import load
29 import numpy as np
30 from transformers import AutoModelForTokenClassification,
31     TrainingArguments, Trainer
32
33 metric = load("accuracy")
34
35 def compute_metrics(pred):
36     logits, labels = pred.predictions, pred.label_ids

```

```

36     predictions = np.argmax(logits, axis=-1)
37     true_labels = labels[labels != -100]
38     true_predictions = predictions[labels != -100]
39     accuracy = metric.compute(predictions=true_predictions, references=
40         true_labels)
41     return {"accuracy": accuracy["accuracy"]}
42
43 def evaluate_on_all_languages(trainer, test_datasets, languages):
44     """Evaluate the fine-tuned model on all chosen languages.
45     test_sets (dict): Dictionary of test datasets keyed by language code.
46     languages (list): List of languages.
47     Returns: dictionary of evaluation accuracies keyed by language code."""
48
49     scores = {}
50     for lang in languages:
51         print(f"Evaluating on language: {lang}")
52         metrics = trainer.evaluate(test_datasets[lang])
53         scores[lang] = metrics.get("eval_accuracy", 0.0)
54         print(f"{lang}: {scores[lang]}")
55     return scores
56
57 def finetune(train_set, dev_set, unique_LABELS, language, test_datasets):
58     base_model_checkpoint = "bert-base-multilingual-cased"
59     model = AutoModelForTokenClassification.from_pretrained(
60         base_model_checkpoint, num_labels=len(list(unique_LABELS))+1)
61
62     training_args = TrainingArguments(
63         output_dir=f"./mBERT_fine_tuning_{language}",
64         eval_strategy="epoch",
65         learning_rate=2e-5,
66         per_device_train_batch_size=16,
67         per_device_eval_batch_size=16,
68         num_train_epochs=3,
69         weight_decay=0.01,
70         logging_dir=f"./logs_{language}",
71         logging_steps=10,
72         report_to="none"
73     )
74
75     trainer = Trainer(
76         model=model,
77         args=training_args,
78         train_dataset=train_set,
79         eval_dataset=dev_set,
80         compute_metrics=compute_metrics
81     )
82
83     trainer.train()
84
85     # Evaluate the model on all languages
86     scores = evaluate_on_all_languages(trainer, test_datasets, list(corpora.
87         keys()))
88     return scores

```

```

87
88
89 # Main program
90
91 corpora = {'breton': corpus_breton,
92           'chinese': corpus_chinese,
93           'russian': corpus_russian,
94           'french': corpus_french}
95 languages = list(corpora.keys())
96 evaluation_results = {}
97
98 for language, corpus in corpora.items():
99     train_set, dev_set, test_set, unique_LABELS = preprocess_corpus(corpus
100 )
101     test_datasets = {lang: preprocess_corpus(corpus)[2] for lang, corpus
102                     in corpora.items()}
103
104     # Fine-tune the model and evaluate
105     scores = finetune(train_set, dev_set, unique_LABELS, language,
106                       test_datasets)
107
108     # Store the scores
109     evaluation_results[language] = scores

```

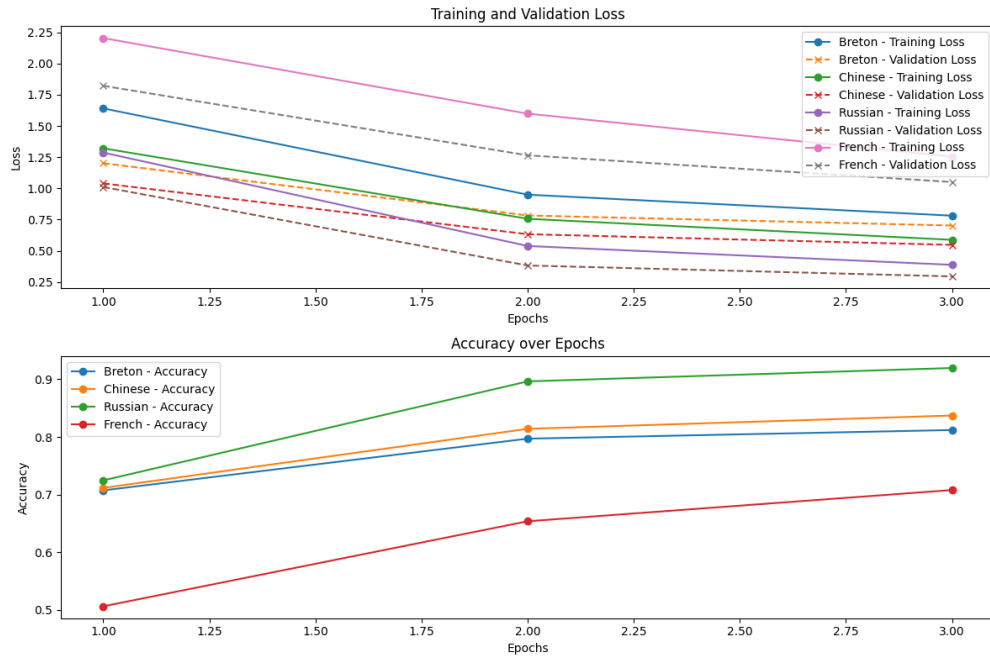


Figure 2: Finetuning results for multiple languages (Breton, Chinese, Russian, French)

These graphs show that the fine-tuning globally enhanced the performances: the training and validation loss decreases for all languages and the accuracy increases, with the Russian having better results than the others, suggesting potential challenges in fine-tuning.

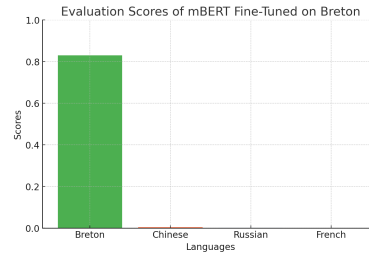
### Results for Breton:

- Result of finetuning:

Epoch	Training Loss	Validation Loss	Accuracy
1	1.641100	1.201294	0.707265
2	0.949500	0.782655	0.797009
3	0.781400	0.702319	0.811966

- Result of evaluation on all languages (scores):

breton: 0.83  
chinese: 0.0047  
russian: 0.0010  
french: 0.00085



The high score for Breton (0.83) indicates that the model fine-tuned on this language performs well for its own data. However, the dramatically lower scores for Chinese (0.0047), Russian (0.0010), and French (0.00085) reveal that this fine-tuning did not generalize well to other languages. This could indicate that mBERT struggles to transfer representations learned from a language like Breton.

## 17. Can you use the same hyperparameters for the different languages?

Using the same hyperparameters across all languages may not yield optimal results due to the variability observed in performance between languages. For instance, Breton achieved relatively high accuracy, while Chinese, Russian, and French had much lower scores when evaluated on the Breton-fine-tuned model. This suggests that adjusting hyperparameters, such as learning rate or training epochs, could improve results by better accommodating each language's unique characteristics.

## 18. What can you conclude?

The results indicate that mBERT's multilingual capabilities vary significantly across languages. While fine-tuning generally improves accuracy, some languages perform better than others, suggesting that mBERT's representations are more effective for certain languages. This highlights both the strengths and limitations of using a single multilingual model across diverse linguistic contexts, and suggests that language-specific adjustments may be necessary for optimal performance.