

✓ HW 4 - POS Tagging with Hugging Face

In this exercise, you will create a part-of-speech (POS) tagging system for Thai text using NECTEC's ORCHID corpus. Instead of building your own deep learning architecture from scratch, you will leverage a pretrained tokenizer and a pretrained token classification model from Hugging Face.

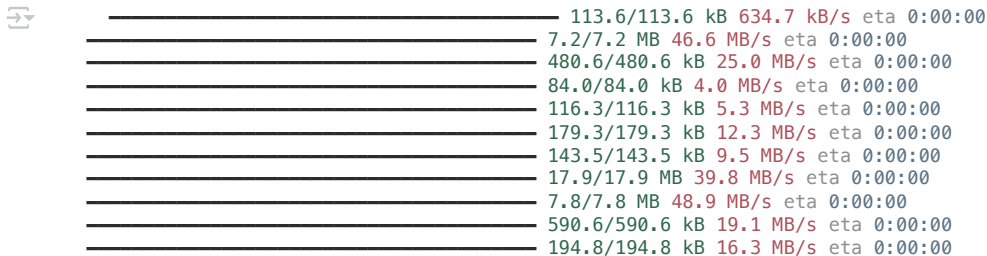
We have provided some starter code for data cleaning and preprocessing in this notebook, but feel free to modify those parts to suit your needs. You are welcome to use additional libraries (e.g., scikit-learn) as long as you incorporate the pretrained Hugging Face model. Specifically, you will need to:

1. Load a pretrained tokenizer and token classification model.
2. Fine-tune it on the ORCHID corpus for POS tagging.
3. Evaluate and report the performance of your model on the test data.

Don't forget to change hardware accelerator to GPU in runtime on Google Colab

✓ 1. Setup and Preprocessing

```
1 # Install transformers and thai2transformers
2 !pip install -q wandb
3 !pip install -q transformers==4.30.1 datasets evaluate thaixtransformers
4 !pip install -q emoji pythainlp sefr_cut tinydb sequeval sentencepiece pydantic jsonlines
5 !pip install -q peft==0.10.0
```



```
113.6/113.6 kB 634.7 kB/s eta 0:00:00
7.2/7.2 MB 46.6 MB/s eta 0:00:00
480.6/480.6 kB 25.0 MB/s eta 0:00:00
84.0/84.0 kB 4.0 MB/s eta 0:00:00
116.3/116.3 kB 5.3 MB/s eta 0:00:00
179.3/179.3 kB 12.3 MB/s eta 0:00:00
143.5/143.5 kB 9.5 MB/s eta 0:00:00
17.9/17.9 MB 39.8 MB/s eta 0:00:00
7.8/7.8 MB 48.9 MB/s eta 0:00:00
590.6/590.6 kB 19.1 MB/s eta 0:00:00
194.8/194.8 kB 16.3 MB/s eta 0:00:00
```

```
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour
sentence-transformers 3.3.1 requires transformers<5.0.0,>=4.41.0, but you have transformers 4.30.1 which is incompatible
torch 2.5.1+cu124 requires nvidia-cublas-cu12==12.4.5.8; platform_system == "Linux" and platform_machine == "x86_64", bu
torch 2.5.1+cu124 requires nvidia-cuda-cupti-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64", bu
torch 2.5.1+cu124 requires nvidia-cuda-nvrtc-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64", bu
torch 2.5.1+cu124 requires nvidia-cuda-runtime-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64", bu
torch 2.5.1+cu124 requires nvidia-cudnn-cu12==9.1.0.70; platform_system == "Linux" and platform_machine == "x86_64", but
torch 2.5.1+cu124 requires nvidia-cufft-cu12==11.2.1.3; platform_system == "Linux" and platform_machine == "x86_64", but
torch 2.5.1+cu124 requires nvidia-curand-cu12==10.3.5.147; platform_system == "Linux" and platform_machine == "x86_64",
torch 2.5.1+cu124 requires nvidia-cusolver-cu12==11.6.1.9; platform_system == "Linux" and platform_machine == "x86_64",
torch 2.5.1+cu124 requires nvidia-cuspars-cu12==12.3.1.170; platform_system == "Linux" and platform_machine == "x86_64",
torch 2.5.1+cu124 requires nvidia-nvjitlink-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64",
gcsfs 2024.10.0 requires fsspec==2024.10.0, but you have fsspec 2024.9.0 which is incompatible.
```

```
43.6/43.6 kB 1.5 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
8.7/8.7 MB 54.6 MB/s eta 0:00:00
118.3/118.3 kB 11.7 MB/s eta 0:00:00
1.3/1.3 MB 70.4 MB/s eta 0:00:00
Building wheel for sequeval (setup.py) ... done
199.1/199.1 kB 5.3 MB/s eta 0:00:00
363.4/363.4 MB 3.8 MB/s eta 0:00:00
13.8/13.8 MB 64.7 MB/s eta 0:00:00
24.6/24.6 MB 35.0 MB/s eta 0:00:00
883.7/883.7 kB 49.5 MB/s eta 0:00:00
664.8/664.8 MB 1.1 MB/s eta 0:00:00
211.5/211.5 MB 5.9 MB/s eta 0:00:00
56.3/56.3 MB 13.6 MB/s eta 0:00:00
127.9/127.9 MB 7.4 MB/s eta 0:00:00
207.5/207.5 MB 3.3 MB/s eta 0:00:00
21.1/21.1 MB 81.0 MB/s eta 0:00:00
```

✓ Setup

1. Register [Wandb account](#) (and confirm your email)
2. wandb login and copy paste the API key when prompt

```
1 import wandb
```

```
1 wandb.login(key=r'648f0ebca50c7021eefe306ab62fcbf0029574da')
```

```
wandb: WARNING If you're specifying your api key in code, ensure this code is not shared publicly.
wandb: WARNING Consider setting the WANDB_API_KEY environment variable, or running `wandb login` from the command line.
wandb: Appending key for api.wandb.ai to your netrc file: /root/.netrc
wandb: Currently logged in as: jirayuwat12 (myfistteam) to https://api.wandb.ai. Use `wandb login --relogin` to force re
True
```

We encourage you to login to your Hugging Face account so you can upload and share your model with the community. When prompted, enter your token to login

```
1 %pip install -q ipywidgets
2 from huggingface_hub import notebook_login
3
4 notebook_login()
```

```
1.6/1.6 MB 18.0 MB/s eta 0:00:00
```

Download the dataset from Hugging Face

```
1 from datasets import load_dataset
2
3 orchid = load_dataset("Thichow/orchid_corpus")
```

```
/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens),
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
  warnings.warn(
README.md: 100% 79.0/79.0 [00:00<00:00, 5.69kB/s]
orchid_corpus.py: 100% 7.91k/7.91k [00:00<00:00, 525kB/s]
The repository for Thichow/orchid_corpus contains custom code which must be executed to correctly load the dataset. You
You can avoid this prompt in future by passing the argument `trust_remote_code=True`.

Do you wish to run the custom code? [y/N] y
Downloading data: 100% 5.24M/5.24M [00:00<00:00, 69.6MB/s]
Downloading data: 100% 1.36M/1.36M [00:00<00:00, 62.9MB/s]
Generating train split: 18500/0 [00:01<00:00, 12576.35 examples/s]
Generating test split: 4625/0 [00:00<00:00, 10888.97 examples/s]
```

```
1 orchid
```

```
DatasetDict({
  train: Dataset({
    features: ['id', 'label_tokens', 'pos_tags', 'sentence'],
    num_rows: 18500
  })
  test: Dataset({
    features: ['id', 'label_tokens', 'pos_tags', 'sentence'],
    num_rows: 4625
  })
})
```

```
1 orchid['train'][0]
```

```
{'id': '0',
 'label_tokens': ['การ', 'ประชุม', 'ทาง', 'วิชาการ', ' ', 'ครั้ง', 'ที่ 1'],
 'pos_tags': [21, 39, 26, 26, 37, 4, 18],
 'sentence': 'การประชุมทางวิชาการ ครั้งที่ 1'}
```

```
1 orchid['train'][0]['sentence']
```

```
'การประชุมทางวิชาการ ครั้งที่ 1'
```

```
1 ''.join(orchid['train'][0]['label_tokens'])
```

```
'การประชุมทางวิชาการ ครั้งที่ 1'
```

```
1 label_list = orchid["train"].features[f"pos_tags"].feature.names
2 print('total type of pos_tags :', len(label_list))
3 print(label_list)
```

```
total type of pos_tags : 47
['ADVI', 'ADVN', 'ADVP', 'ADVS', 'CFQC', 'CLTV', 'CMTR', 'CMTR@PUNC', 'CNIT', 'CVBL', 'DCNM', 'DDAC', 'DDAN', 'DDAQ', 'D
```

```

1 import numpy as np
2 import numpy.random
3 import torch
4
5 from tqdm.auto import tqdm
6 from functools import partial
7
8 #transformers
9 from transformers import (
10     CamembertTokenizer,
11     AutoTokenizer,
12     AutoModel,
13     AutoModelForMaskedLM,
14     AutoModelForSequenceClassification,
15     AutoModelForTokenClassification,
16     TrainingArguments,
17     Trainer,
18     pipeline,
19 )
20
21 #thaixtransformers
22 from thaixtransformers import Tokenizer
23 from thaixtransformers.preprocess import process_transformers

```

🔄 The cache for model files in Transformers v4.22.0 has been updated. Migrating your old cache. This is a one-time only op
0/0 [00:00<?, ?it/s]

Next, we load a pretrained tokenizer from Hugging Face. In this work, we utilize WangchanBERTa, a Thai-specific pretrained model, as the tokenizer.

✓ Choose Pretrained Model

In this notebook, you can choose from 5 versions of WangchanBERTa, XLMR and mBERT to perform downstream tasks on Thai datasets. The datasets are:

- wangchanberta-base-att-spm-uncased (recommended) - Largest WangchanBERTa trained on 78.5GB of Assorted Thai Texts with subword tokenizer SentencePiece
- xlm-roberta-base - Facebook's [XLMR](#) trained on 100 languages
- bert-base-multilingual-cased - Google's [mBERT](#) trained on 104 languages
- wangchanberta-base-wiki-newmm - WangchanBERTa trained on Thai Wikipedia Dump with PyThaiNLP's word-level tokenizer newmm
- wangchanberta-base-wiki-syllable - WangchanBERTa trained on Thai Wikipedia Dump with PyThaiNLP's syllable-level tokenizer syllable
- wangchanberta-base-wiki-sefr - WangchanBERTa trained on Thai Wikipedia Dump with word-level tokenizer SEFR
- wangchanberta-base-wiki-spm - WangchanBERTa trained on Thai Wikipedia Dump with subword-level tokenizer SentencePiece

In the first part, we require you to select the wangchanberta-base-att-spm-uncased.

Learn more about using wangchanberta at [wangchanberta_getting_started_ai_reseach](#)

- **You need to set the transformers version to transformers==4.30.1.**

In the first part, we require you to select the wangchanberta-base-att-spm-uncased.

✓ Choose Pretrained Model

```

1 model_names = [
2     'airesearch/wangchanberta-base-att-spm-uncased',
3     'airesearch/wangchanberta-base-wiki-newmm',
4     'airesearch/wangchanberta-base-wiki-ssg',
5     'airesearch/wangchanberta-base-wiki-sefr',
6     'airesearch/wangchanberta-base-wiki-spm',
7 ]
8
9 #@title Choose Pretrained Model
10 model_name = "airesearch/wangchanberta-base-att-spm-uncased"
11
12 #create tokenizer
13 tokenizer = Tokenizer(model_name).from_pretrained(
14     f'{model_name}',

```

```

15         revision='main',
16         model_max_length=416,)
17

```

```

➔ /usr/local/lib/python3.11/dist-packages/huggingface_hub/file_download.py:795: FutureWarning: `resume_download` is deprecated
  warnings.warn(
sentencepiece.bpe.model: 100%          905k/905k [00:00<00:00, 4.76MB/s]
tokenizer_config.json: 100%           282/282 [00:00<00:00, 10.2kB/s]
config.json: 100%                     546/546 [00:00<00:00, 34.6kB/s]
The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may
The tokenizer class you load from this checkpoint is 'CamembertTokenizer'.
The class this function is called from is 'WangchanbertaTokenizer'.
The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may
The tokenizer class you load from this checkpoint is 'CamembertTokenizer'.
The class this function is called from is 'WangchanbertaTokenizer'.

```

Let's try using a pretrained tokenizer.

```

1 text = 'ศิลปะไม่เป็นเจ้านายใคร และไม่เป็นที่ซ้ำใคร'
2 print('text :', text)
3 tokens = []
4 for i in tokenizer([text], is_split_into_words=True)['input_ids']:
5     tokens.append(tokenizer.decode(i))
6 print('tokens :', tokens)

➔ text : ศิลปะไม่เป็นเจ้านายใคร และไม่เป็นที่ซ้ำใคร
  tokens : ['<s>', ' ', 'ศิลปะ', 'ไม่เป็น', 'เจ้านาย', 'ใคร', '<_>', 'และ', 'ไม่เป็น', 'ซ้ำ', 'ใคร', '</s>']

```

model:* wangchanberta-base-att-spm-uncased

First, we print examples of label tokens from our dataset for inspection.

```

1 example = orchid["train"][0]
2 for i in example :
3     print(i, ': ', example[i])

➔ id : 0
  label_tokens : ['การ', 'ประชุม', 'ทาง', 'วิชาการ', ' ', ' ', 'ครั้ง', 'ที่ 1']
  pos_tags : [21, 39, 26, 26, 37, 4, 18]
  sentence : การประชุมทางวิชาการ ครั้งที่ 1

```

Then, we use the sentence 'การประชุมทางวิชาการครั้งที่ 1' to be tokenized by the pretrained tokenizer model.

```

1 text = 'การประชุมทางวิชาการ ครั้งที่ 1'
2 tokenizer(text)

➔ {'input_ids': [5, 10, 882, 8222, 8, 10, 1014, 8, 10, 59, 6], 'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]}

```

These are already mapped into discrete values. We can uncover the original token text from the tokens by.

```

1 for i in tokenizer(text)['input_ids']:
2     print(tokenizer.convert_ids_to_tokens(i))

➔ <s>
  _
  การประชุม
  ทางวิชาการ
  <_>
  _
  ครั้งที่
  <_>
  _
  1
  </s>

```

Now let's look at another example.

```

1 example = orchid["train"][1899]
2 print('sentence :', example["sentence"])
3 tokenized_input = tokenizer([example["sentence"]], is_split_into_words=True)
4 tokens = tokenizer.convert_ids_to_tokens(tokenized_input["input_ids"])
5 print('tokens :', tokens)
6 print('label tokens :', example["label_tokens"])
7 print('label pos :', example["pos_tags"])

```

```

sentence : โดยพิจารณาจากพจนานุกรมภาษาคู่ (Bilingual transfer dictionary)
tokens : ['<s>', '_โดย', 'พิจารณาจาก', 'พจนานุกรม', 'ภาษา', 'คู่', '<_>', '(', '<unk>', 'i', 'ling', 'ual', '<_>', '_', 'tr',
label tokens : ['โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', ' ', '(', 'Bilingual transfer dictionary', ')']
label pos : [25, 39, 38, 26, 26, 5, 37, 37, 26, 37]

```

Notice how B becomes an <unk> token. This is because this is an uncased model, meaning it only handles small English characters.

✓ #TODO 0

Convert the dataset to lowercase.

```

1 # Create a lowercase dataset for uncased BERT
2 def lower_case_sentences(examples):
3     lower_cased_examples = examples
4     lower_cased_examples['sentence'] = lower_cased_examples['sentence'].lower()
5     lower_cased_examples['label_tokens'] = [t.lower() for t in lower_cased_examples['label_tokens']]
6
7     # fill code here to lower case the "sentence" and "label_tokens"
8
9     return lower_cased_examples

```

```

1 orchidl = orchid.map(lower_case_sentences)

```

```

Map: 100%                               18500/18500 [00:04<00:00, 6222.74 examples/s]
Map: 100%                               4625/4625 [00:00<00:00, 11289.08 examples/s]

```

```

1 orchidl

```

```

DatasetDict({
  train: Dataset({
    features: ['id', 'label_tokens', 'pos_tags', 'sentence'],
    num_rows: 18500
  })
  test: Dataset({
    features: ['id', 'label_tokens', 'pos_tags', 'sentence'],
    num_rows: 4625
  })
})

```

```

1 orchidl["train"][1899]

```

```

{'id': '1899',
 'label_tokens': ['โดย',
 'พิจารณา',
 'จาก',
 'พจนานุกรม',
 'ภาษา',
 'คู่',
 ' ',
 '(',
 'bilingual transfer dictionary',
 ')'],
 'pos_tags': [25, 39, 38, 26, 26, 5, 37, 37, 26, 37],
 'sentence': 'โดยพิจารณาจากพจนานุกรมภาษาคู่ (bilingual transfer dictionary)'}

```

Now let's examine the labels again.

```

1 example = orchidl["train"][1899]
2 print('sentence :', example["sentence"])
3 tokenized_input = tokenizer([example["sentence"]], is_split_into_words=True)
4 tokens = tokenizer.convert_ids_to_tokens(tokenized_input["input_ids"])
5 print('tokens :', tokens)
6 print('label tokens :', example["label_tokens"])
7 print('label pos :', example["pos_tags"])

```

```

sentence : โดยพิจารณาจากพจนานุกรมภาษาคู่ (bilingual transfer dictionary)
tokens : ['<s>', '_โดย', 'พิจารณาจาก', 'พจนานุกรม', 'ภาษา', 'คู่', '<_>', '(', 'bi', 'ling', 'ual', '<_>', '_', 'trans', 'fi',
label tokens : ['โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', ' ', '(', 'bilingual transfer dictionary', ')']
label pos : [25, 39, 38, 26, 26, 5, 37, 37, 26, 37]

```

```

1 example = orchidl["train"][0]
2 print('sentence :', example["sentence"])
3 tokenized_input = tokenizer([example["sentence"]], is_split_into_words=True)
4 tokens = tokenizer.convert_ids_to_tokens(tokenized_input["input_ids"])
5 print('tokens :', tokens)

```

```
6 print('label tokens :', example["label_tokens"])
7 print('label pos :', example["pos_tags"])
```

```
→ sentence : การประชุมทางวิชาการ ครั้งที่ 1
tokens : ['<s>', ' ', 'การประชุม', 'ทางวิชาการ', '<_>', ' ', 'ครั้งที่', '<_>', ' ', '1', '</s>']
label tokens : ['การ', 'ประชุม', 'ทาง', 'วิชาการ', ' ', ' ', 'ครั้งที่', 'ที่ 1']
label pos : [21, 39, 26, 26, 37, 4, 18]
```

In the example above, tokens refer to those tokenized using the pretrained tokenizer, while label tokens refer to tokens tokenized from our dataset.

Do you see something?

Yes, the tokens from the two tokenizers do not match.

- sentence : การประชุมทางวิชาการ ครั้งที่ 1

-
- tokens : ['<s>', ' ', 'การประชุม', 'ทางวิชาการ', '<_>', ' ', 'ครั้งที่', '<_>', ' ', '1', '</s>']
-

- label tokens : ['การ', 'ประชุม', 'ทาง', 'วิชาการ', ' ', ' ', 'ครั้งที่', 'ที่ 1']
- label pos : [21, 39, 26, 26, 37, 4, 18]

You can see that in our label tokens, 'การ' has a POS tag of 21, and 'ประชุม' has a POS tag of 39. However, when we tokenize the sentence using WangchanBERTa, we get the token 'การประชุม'. What POS tag should we assign to this new token?

What should we do ?

Based on this example, we found that the tokens from the WangchanBERTa do not directly align with our label tokens. This means we cannot directly use the label POS tags. Therefore, we need to reassign POS tags to the tokens produced by WangchanBERTa tokenization. The method we will use is majority voting:

- If a token from the WangchanBERTa matches a label token exactly, we will directly assign the POS tag from the label POS.
- If the token generated overlaps or combines multiple label tokens, we assign the POS tag based on the number of characters in each token: If the token contains the most characters from any label token, we assign the POS tag from that label token.

Example :

```
# "การประชุม" (9 chars) is formed from "การ" (3 chars) + "ประชุม" (6 chars).
# "การ" has a POS tag of 21,
# and "ประชุม" has a POS tag of 39.
# Therefore, the POS tag for "การประชุม" is 39,
# as "การประชุม" is derived more from the "ประชุม" part than from the "การ" part.
```

```
# "ทางวิชาการ" (10 chars) is formed from 'ทาง' (3 chars) + 'วิชาการ' (7 chars)
# "ทาง" has a POS tag of 26,
# and "วิชาการ" has a POS tag of 2.
# Therefore, the POS tag for "ทางวิชาการ" is 2,
# as "ทางวิชาการ" is derived more from the "ทาง" part than from the "วิชาการ" part.
```

✓ #TODO 1

****Warning:** Please be careful of <unk>, an unknown word token.**

****Warning:** Please be careful of "ํ", the 'am' vowel. WangchanBERTa's internal preprocessing replaces all "ํ" to "ั" and "ำ" to "ั"

Assigning the label -100 to the special tokens [<s>] and [</s>] and [_] so they're ignored by the PyTorch loss function (see [CrossEntropyLoss: ignore_index](#))

```
1 def majority_vote_pos(examples):
2
3     #####
4     # TO DO: Since the tokens from the output of the pretrained tokenizer
5     # do not match the tokens in the label tokens of the dataset,
6     # the task is to create a function to determine the POS tags of the tokens generated by the pretrained tokenizer.
7     # This should be done by referencing the POS tags in the label tokens. If a token partially overlaps with others,
8     # the POS tag from the segment with the greater number of characters should be assigned.
9     #
10    # Example :
11    # "การประชุม" (9 chars) is formed from "การ" (3 chars) + "ประชุม" (6 chars).
12    # "การ" has a POS tag of 21,
13    # and "ประชุม" has a POS tag of 39.
```

```

14 # Therefore, the POS tag for "การประชุม" is 39,
15 # as "การประชุม" is derived more from the "ประชุม" part than from the "การ" part.
16 #
17 # 'ทางวิชาการ' (10 chars) is formed from 'ทาง' (3 chars) + 'วิชาการ' (7 chars)
18 # "ทาง" has a POS tag of 26,
19 # and "วิชาการ" has a POS tag of 2.
20 # Therefore, the POS tag for "ทางวิชาการ" is 2,
21 # as "ทางวิชาการ" is derived more from the "ทาง" part than from the "วิชาการ" part.
22
23 # tokenize word by pretrained tokenizer
24 tokenized_inputs = tokenizer([examples["sentence"]], is_split_into_words=True)
25
26 # FILL CODE HERE
27 label_tokens = examples["label_tokens"]
28 pos_tags = examples["pos_tags"]
29 new_pos_result = []
30
31 new_tokens = tokenizer.convert_ids_to_tokens(tokenized_inputs["input_ids"])
32 # print(new_tokens, "\n", label_tokens, "\n", pos_tags)
33
34 label_idx = 0
35 i = 0
36 for t in new_tokens:
37     if t in ["<s>", "</s>", "_"]:
38         new_pos_result.append(-100)
39         continue
40     # if t == "<_>":
41     #     new_pos_result.append(37)
42     #     continue
43
44     buffer = ""
45     weights = {}
46     t = t.replace("'", "")
47     t = t.replace("<_>", " ")
48     if t[0] == "_":
49         t = t[1:]
50     while label_tokens[label_idx][i] != t[0]:
51         i += 1
52     if i == len(label_tokens[label_idx]):
53         label_idx += 1
54         i = 0
55
56     while buffer != t:
57         buffer += label_tokens[label_idx][i]
58         # print(t, buffer, t==buffer)
59         if pos_tags[label_idx] not in weights:
60             weights[pos_tags[label_idx]] = 0
61         weights[pos_tags[label_idx]] += 1
62         i += 1
63     if i == len(label_tokens[label_idx]):
64         label_idx += 1
65         i = 0
66     # print(weights)
67     max_key = max(weights, key=weights.get)
68     new_pos_result.append(max_key)
69
70 tokenized_inputs['tokens'] = new_tokens
71 tokenized_inputs['labels'] = new_pos_result
72
73 return tokenized_inputs
74 #####

```

```
1 tokenized_orchid = orchidl.map(majority_vote_pos)
```



Map: 100%

18500/18500 [00:22<00:00, 2113.48 examples/s]

Map: 100%

4625/4625 [00:02<00:00, 2126.82 examples/s]

```
1 tokenized_orchid
```



```

DatasetDict({
  train: Dataset({
    features: ['id', 'label_tokens', 'pos_tags', 'sentence', 'input_ids', 'attention_mask', 'tokens', 'labels'],
    num_rows: 18500
  })
  test: Dataset({
    features: ['id', 'label_tokens', 'pos_tags', 'sentence', 'input_ids', 'attention_mask', 'tokens', 'labels'],
    num_rows: 4625
  })
})

```

```
1 tokenized_orchid['train'][0]
```

```
{'id': '0',  
'label_tokens': ['การ', 'ประชุม', 'ทาง', 'วิชาการ', ' ', ' ', 'ครั้ง', 'ที่ 1'],  
'pos_tags': [21, 39, 26, 26, 37, 4, 18],  
'sentence': 'การประชุมทางวิชาการ ครั้งที่ 1',  
'input_ids': [5, 10, 882, 8222, 8, 10, 1014, 8, 10, 59, 6],  
'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1],  
'tokens': ['<s>',  
'_',  
'การประชุม',  
'ทางวิชาการ',  
'<_>',  
'_',  
'ครั้งที่',  
'<_>',  
'_',  
'1',  
'</s>'],  
'labels': [-100, -100, 39, 26, 37, -100, 4, 18, -100, 18, -100]}
```

```
1 example = tokenized_orchid["train"][0]  
2 for i in example :  
3     print(i, ":", example[i])
```

```
{'id': 0  
 'label_tokens': ['การ', 'ประชุม', 'ทาง', 'วิชาการ', ' ', ' ', 'ครั้ง', 'ที่ 1']  
 'pos_tags': [21, 39, 26, 26, 37, 4, 18]  
 'sentence': 'การประชุมทางวิชาการ ครั้งที่ 1'  
 'input_ids': [5, 10, 882, 8222, 8, 10, 1014, 8, 10, 59, 6]  
 'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]  
 'tokens': ['<s>', ' ', 'การประชุม', 'ทางวิชาการ', '<_>', ' ', 'ครั้งที่', '<_>', ' ', '1', '</s>']  
 'labels': [-100, -100, 39, 26, 37, -100, 4, 18, -100, 18, -100]}
```

This is the result after we realigned the POS based on the majority vote.

- label_tokens: ['การ', 'ประชุม', 'ทาง', 'วิชาการ', ' ', ' ', 'ครั้ง', 'ที่ 1']
- pos_tags: [21, 39, 26, 26, 37, 4, 18]
- tokens: ['<s>', ' ', 'การประชุม', 'ทางวิชาการ', '<_>', ' ', 'ครั้งที่', '<_>', ' ', '1', '</s>']
- labels: [-100, -100, 39, 26, 37, -100, 4, 18, -100, 18, -100]

```
['<s>', ' ', '</s>'] : -100
```

Check :

"การประชุม" (9 chars) is formed from "การ" (3 chars) + "ประชุม" (6 chars).

"การ" has a POS tag of 21,

and "ประชุม" has a POS tag of 39.

Therefore, the POS tag for "การประชุม" is 39,

as "การประชุม" is derived more from the "ประชุม" part than from the "การ" part.

```
1 # hard test case  
2 example = tokenized_orchid["train"][1899]  
3 for i in example :  
4     print(i, ":", example[i])
```

```
{'id': 1899  
 'label_tokens': ['โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', ' ', '(', 'bilingual transfer dictionary', ')']  
 'pos_tags': [25, 39, 38, 26, 26, 5, 37, 37, 26, 37]  
 'sentence': 'โดยพิจารณาจากพจนานุกรมภาษาคู่ (bilingual transfer dictionary)'  
 'input_ids': [5, 489, 15617, 19737, 958, 493, 8, 1241, 4906, 11608, 12177, 8, 10, 11392, 9806, 8, 10, 2951, 15779, 8001,  
 'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]  
 'tokens': ['<s>', 'โดย', 'พิจารณาจาก', 'พจนานุกรม', 'ภาษา', 'คู่', '<_>', '(', 'bi', 'ling', 'ual', '<_>', ' ', 'trans', 'fer', '<_>',  
 'labels': [-100, 25, 39, 26, 26, 5, 37, 37, 26, 26, 26, 26, -100, 26, 26, 26, -100, 26, 26, 26, 37, -100]}
```

Expected output

```
id : 1899  
label_tokens : ['โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', ' ', '(', 'bilingual transfer dictionary', ')']  
pos_tags : [25, 39, 38, 26, 26, 5, 37, 37, 26, 37]  
sentence : โดยพิจารณาจากพจนานุกรมภาษาคู่ (bilingual transfer dictionary)  
input_ids : [5, 489, 15617, 19737, 958, 493, 8, 1241, 4906, 11608, 12177, 8, 10, 11392, 9806, 8, 10, 2951, 15779, 8001, 29, 6]  
attention_mask : [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]  
tokens : ['<s>', 'โดย', 'พิจารณาจาก', 'พจนานุกรม', 'ภาษา', 'คู่', '<_>', '(', 'bi', 'ling', 'ual', '<_>', ' ', 'trans', 'fer', '<_>',  
labels : [-100, 25, 39, 26, 26, 5, 37, 37, 26, 26, 26, 26, -100, 26, 26, 26, -100, 26, 26, 26, 37, -100]
```


✓ Train and Evaluate model

We will create a batch of examples using [DataCollatorWithPadding](#).

Data collators are objects that will form a batch by using a list of dataset elements as input. These elements are of the same type as the elements of `train_dataset` or `eval_dataset`.

`DataCollatorWithPadding` will help us pad the sentences to the longest length in a batch during collation, instead of padding the whole dataset to the maximum length. This allows for efficient computation during each batch.

- `DataCollatorForTokenClassification`: `padding` (`bool`, `str` or `PaddingStrategy`, optional, defaults to `True`)
- `True` or `'longest'` (default): Pad to the longest sequence in the batch (or no padding if only a single sequence is provided).

```
1 from transformers import DataCollatorForTokenClassification
2
3 data_collator = DataCollatorForTokenClassification(tokenizer=tokenizer)
```

For evaluating your model's performance. You can quickly load a evaluation method with the [Evaluate](#) library. For this task, load the [sequeval](#) framework (see the Evaluate [quick tour](#) to learn more about how to load and compute a metric). `Sequeval` actually produces several scores: precision, recall, F1, and accuracy.

```
1 import evaluate
2
3 sequeval = evaluate.load("sequeval")
```



Downloading builder script: 100%

6.34k/6.34k [00:00<00:00, 295kB/s]

Huggingface requires us to write a `compute_metrics()` function. This will be invoked when huggingface evaluates a model.

Note that we ignore to evaluate on -100 labels.

```
1 import numpy as np
2 import warnings
3
4
5 def compute_metrics(p):
6     predictions, labels = p
7     predictions = np.argmax(predictions, axis=2)
8
9     true_predictions = [
10         [label_list[p] for (p, l) in zip(prediction, label) if l != -100]
11         for prediction, label in zip(predictions, labels)
12     ]
13     true_labels = [
14         [label_list[l] for (p, l) in zip(prediction, label) if l != -100]
15         for prediction, label in zip(predictions, labels)
16     ]
17
18     with warnings.catch_warnings():
19         warnings.filterwarnings("ignore")
20         results = sequeval.compute(predictions=true_predictions, references=true_labels)
21     return {
22         "precision": results["overall_precision"],
23         "recall": results["overall_recall"],
24         "f1": results["overall_f1"],
25         "accuracy": results["overall_accuracy"],
26     }
```

The total number of labels in our POS tag set.

```
1 id2label = {
2     0: 'ADVI',
3     1: 'ADVN',
4     2: 'ADVP',
5     3: 'ADVS',
6     4: 'CFQC',
7     5: 'CLTV',
8     6: 'CMTR',
9     7: 'CMTR@PUNC',
10    8: 'CNIT',
11    9: 'CVBL',
12   10: 'DCNM',
```

```

13     11: 'DDAC',
14     12: 'DDAN',
15     13: 'DDAQ',
16     14: 'DDBQ',
17     15: 'DIAC',
18     16: 'DIAQ',
19     17: 'DIBQ',
20     18: 'DONM',
21     19: 'EAFF',
22     20: 'EITT',
23     21: 'FIXN',
24     22: 'FIXV',
25     23: 'JCMP',
26     24: 'JCRG',
27     25: 'JSBR',
28     26: 'NCMN',
29     27: 'NCNM',
30     28: 'NEG',
31     29: 'NLBL',
32     30: 'NONM',
33     31: 'NPRP',
34     32: 'NTTL',
35     33: 'PDMN',
36     34: 'PNTR',
37     35: 'PPRS',
38     36: 'PREL',
39     37: 'PUNC',
40     38: 'RPRE',
41     39: 'VACT',
42     40: 'VATT',
43     41: 'VSTA',
44     42: 'XVAE',
45     43: 'XVAM',
46     44: 'XVBB',
47     45: 'XVBM',
48     46: 'XVMM',
49     # 47: 'O'
50 }
51 label2id = {}
52 for k, v in id2label.items() :
53     label2id[v] = k
54
55 label2id

```

```

↔ {'ADVI': 0,
   'ADVN': 1,
   'ADVP': 2,
   'ADVS': 3,
   'CFQC': 4,
   'CLTV': 5,
   'CMTR': 6,
   'CMTR@PUNC': 7,
   'CNIT': 8,
   'CVBL': 9,
   'DCNM': 10,
   'DDAC': 11,
   'DDAN': 12,
   'DDAQ': 13,
   'DDBQ': 14,
   'DIAC': 15,
   'DIAQ': 16,
   'DIBQ': 17,
   'DONM': 18,
   'EAFF': 19,
   'EITT': 20,
   'FIXN': 21,
   'FIXV': 22,
   'JCMP': 23,
   'JCRG': 24,
   'JSBR': 25,
   'NCMN': 26,
   'NCNM': 27,
   'NEG': 28,
   'NLBL': 29,
   'NONM': 30,
   'NPRP': 31,
   'NTTL': 32,
   'PDMN': 33,
   'PNTR': 34,
   'PPRS': 35,
   'PREL': 36,
   'PUNC': 37,
   'RPRE': 38,
   'VACT': 39,
   'VATT': 40,
   'VSTA': 41,

```

```
'XVAE': 42,
'XVAM': 43,
'XVBB': 44,
'XVBM': 45,
'XVMM': 46}
```

```
1 labels = [i for i in id2label.values()]
2 labels
```

```
['ADVI',
'ADVN',
'ADVP',
'ADVS',
'CFQC',
'CLTV',
'CMTR',
'CMTR@PUNC',
'CNIT',
'CVBL',
'DCNM',
'DDAC',
'DDAN',
'DDAQ',
'DDBQ',
'DIAC',
'DIAQ',
'DIBQ',
'DONM',
'EAFF',
'EITT',
'FIXN',
'FIXV',
'JCMP',
'JCRG',
'JSBR',
'NCMN',
'NCNM',
'NEG',
'NLBL',
'NONM',
'NPRP',
'NTTL',
'PDMN',
'PNTR',
'PPRS',
'PREL',
'PUNC',
'RPRE',
'VACT',
'VATT',
'VSTA',
'XVAE',
'XVAM',
'XVBB',
'XVBM',
'XVMM']
```

✓ Load pretrained model

Select a pretrained model for fine-tuning to develop a POS Tagger model using the Orchid corpus dataset.

- model: wangchanberta-base-att-spm-uncased
- Don't forget to update the num_labels.

You're ready to start training your model now! Load pretrained model with AutoModelForTokenClassification along with the number of expected labels, and the label mappings:

In the first part, we require you to select the wangchanberta-base-att-spm-uncased.

✓ Choose Pretrained Model

```
1 model_names = [
2     'wangchanberta-base-att-spm-uncased',
3     'wangchanberta-base-wiki-newmm',
4     'wangchanberta-base-wiki-ssg',
5     'wangchanberta-base-wiki-sefr',
6     'wangchanberta-base-wiki-spm',
7 ]
8
9 #@title Choose Pretrained Model
10 model_name = "wangchanberta-base-att-spm-uncased"
```

```

11
12 #create model
13 model = AutoModelForTokenClassification.from_pretrained(
14     f"airesearch/{model_name}",
15     revision='main',
16     num_labels=47, id2label=id2label, label2id=label2id
17 )
18

```

```

➦ /usr/local/lib/python3.11/dist-packages/huggingface_hub/file_download.py:795: FutureWarning: `resume_download` is deprecated and will be removed in version 0.11. The automatic download of files will always resume. You are encouraged to explicitly call `resume_download=True` if you wish to explicitly request to resume the download if the file exists, to avoid silently overriding possible metadata.
  warnings.warn(

model.safetensors: 100% 423M/423M [00:02<00:00, 231MB/s]

Some weights of the model checkpoint at airesearch/wangchanberta-base-att-spm-uncased were not used when initializing CamembertForTokenClassification:
- This IS expected if you are initializing CamembertForTokenClassification from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForTokenClassification model).
- This IS NOT expected if you are initializing CamembertForTokenClassification from the checkpoint of a model that you expect to be initialized from without any initialization (e.g. initializing a Pytorch model from a Pytorch checkpoint).
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```

✓ #TODO 2

- Configure your training hyperparameters using `**TrainingArguments**`. The only required parameter is `output_dir`, which determines the directory where your model will be saved. To upload the model to the Hugging Face Hub, set `push_to_hub=True` (note: you must be logged into Hugging Face for this). During training, the Trainer will compute seqeval metrics at the end of each epoch and store the training checkpoint.
- Provide the `**Trainer**` with the training arguments, as well as the model, dataset, tokenizer, data collator, and `compute_metrics` function.
- Use `**train()**` to fine-tune the model.

Read [huggingface's tutorial](#) for more details.

```

1 training_args = TrainingArguments(
2     #####
3     output_dir="pos-spm-uncased",
4     learning_rate=2e-5,
5     per_device_train_batch_size=32,
6     per_device_eval_batch_size=32,
7     num_train_epochs=2,
8     weight_decay=0.01,
9     push_to_hub=True
10    #####
11 )
12
13 trainer = Trainer(
14     #####
15     model=model,
16     args=training_args,
17     train_dataset=tokenized_orchid["train"],
18     eval_dataset=tokenized_orchid["test"],
19     data_collator=data_collator,
20     compute_metrics=compute_metrics,
21     #####
22 )
23
24 trainer.train()

```

```

➔ /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_deprecation.py:131: FutureWarning: 'Repository' (from 'hu
For more details, please read https://huggingface.co/docs/huggingface\_hub/concepts/git\_vs\_http.
warnings.warn(warning_message, FutureWarning)
Cloning https://huggingface.co/Jirayuwat12/pos-spm-uncased into local empty directory.
WARNING:huggingface_hub.repository:Cloning https://huggingface.co/Jirayuwat12/pos-spm-uncased into local empty directory
/usr/local/lib/python3.11/dist-packages/transformers/optimization.py:411: FutureWarning: This implementation of AdamW is
warnings.warn(
wandb: Using wandb-core as the SDK backend. Please refer to https://wandb.me/wandb-core for more information.
Tracking run with wandb version 0.19.5
Run data is saved locally in /content/wandb/run-20250201_083555-7n089k1e
Syncing run breezy-vortex-16 to Weights & Biases (docs)
View project at https://wandb.ai/myfisteam/huggingface
View run at https://wandb.ai/myfisteam/huggingface/runs/7n089k1e
[1158/1158 07:25, Epoch 2/2]

```

Step Training Loss

500	0.984700
1000	0.361100

```

TrainOutput(global_step=1158, training_loss=0.6244609681229929, metrics={'train_runtime': 450.5004,
'train_samples_per_second': 82.131, 'train_steps_per_second': 2.57, 'total_flos': 1207602506181672.0, 'train_loss':
0.6244609681229929, 'epoch': 2.0})

```

▼ Inference

With your model fine-tuned, you can now perform inference.

```
1 text = "การประชุมทางวิชาการ ครั้งที่ 1"
```

In the first part, we require you to select the wangchanberta-base-att-spm-uncased.

▼ Choose Pretrained Model

```

1 from transformers import AutoTokenizer
2
3 # Load pretrained tokenizer from Hugging Face
4 #@title Choose Pretrained Model
5 model_name = "airesearch/wangchanberta-base-att-spm-uncased"
6
7 tokenizer = Tokenizer(model_name).from_pretrained(model_name)
8 inputs = tokenizer(text, return_tensors="pt")

```

```

➔ /usr/local/lib/python3.11/dist-packages/huggingface_hub/file_download.py:795: FutureWarning: `resume_download` is deprec
warnings.warn(
The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may
The tokenizer class you load from this checkpoint is 'CamembertTokenizer'.
The class this function is called from is 'WangchanbertaTokenizer'.
The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may
The tokenizer class you load from this checkpoint is 'CamembertTokenizer'.
The class this function is called from is 'WangchanbertaTokenizer'.

```

```
1 inputs
```

```

➔ {'input_ids': tensor([[ 5, 10, 882, 8222, 8, 10, 1014, 8, 10, 59, 6]]), 'attention_mask':
tensor([[1, 1, 1, 1, 1, 1, 1, 1, 1, 1]])}

```

```

1 from transformers import AutoModelForTokenClassification
2
3 ## Load your fine-tuned model from Hugging Face
4 model = AutoModelForTokenClassification.from_pretrained("jirayuwat12/pos-spm-uncased") ## your model path from Hugging Fa
5 with torch.no_grad():
6     logits = model(**inputs).logits

```

```

➔ config.json: 100%                2.41k/2.41k [00:00<00:00, 71.1kB/s]

pytorch_model.bin: 100%            419M/419M [00:09<00:00, 42.7MB/s]
/usr/local/lib/python3.11/dist-packages/transformers/modeling_utils.py:463: FutureWarning: You are using `torch.load` wi
return torch.load(checkpoint_file, map_location="cpu")

```

```

1 predictions = torch.argmax(logits, dim=2)
2 predicted_token_class = [model.config.id2label[t.item()]] for t in predictions[0]]
3 predicted_token_class

```

```

→ ['PUNC',
   'PUNC',
   'VACT',
   'NCMN',
   'PUNC',
   'PUNC',
   'CFQC',
   'DONM',
   'DONM',
   'DONM',
   'PUNC']

```

```
1 id2label
```

```

→ {0: 'ADVI',
   1: 'ADVN',
   2: 'ADVP',
   3: 'ADVS',
   4: 'CFQC',
   5: 'CLTV',
   6: 'CMTR',
   7: 'CMTR@PUNC',
   8: 'CNIT',
   9: 'CVBL',
   10: 'DCNM',
   11: 'DDAC',
   12: 'DDAN',
   13: 'DDAQ',
   14: 'DDBQ',
   15: 'DIAC',
   16: 'DIAQ',
   17: 'DIBQ',
   18: 'DONM',
   19: 'EAFF',
   20: 'EITT',
   21: 'FIXN',
   22: 'FIXV',
   23: 'JCMP',
   24: 'JCRG',
   25: 'JSBR',
   26: 'NCMN',
   27: 'NCNM',
   28: 'NEG',
   29: 'NLBL',
   30: 'NONM',
   31: 'NPRP',
   32: 'NTTL',
   33: 'PDMN',
   34: 'PNTR',
   35: 'PPRS',
   36: 'PREL',
   37: 'PUNC',
   38: 'RPRE',
   39: 'VACT',
   40: 'VATT',
   41: 'VSTA',
   42: 'XVAE',
   43: 'XVAM',
   44: 'XVBB',
   45: 'XVBM',
   46: 'XVMM'}

```

```

1 # Inference
2 # ignore special tokens
3 text = 'จะว่าไปแล้วเชิงเทียนของผมนี่สวยดีเหมือนกัน'
4 inputs = tokenizer(text, return_tensors="pt")
5 tokenized_input = tokenizer([text], is_split_into_words=True)
6 tokens = tokenizer.convert_ids_to_tokens(tokenized_input["input_ids"])
7 print('tokens : ', tokens)
8 with torch.no_grad():
9     logits = model(**inputs).logits
10 predictions = torch.argmax(logits, dim=2)
11 predicted_token_class = [model.config.id2label[t.item()] for t in predictions[0]]
12 print('predict pos : ', predicted_token_class)

```

```

→ tokens : ['<s>', ' ', 'จะว่าไป', 'แล้ว', 'เชิง', 'เทียน', 'ของ', 'ผม', 'นี่', 'สวยดี', 'เหมือนกัน', '</s>']
predict pos : ['PUNC', 'PUNC', 'JSBR', 'JSBR', 'NCMN', 'NCMN', 'RPRE', 'JSBR', 'ADVN', 'ADVN', 'PUNC']

```

Evaluate model :

The output from the model is a softmax over classes. We choose the maximum class as the answer for evaluation. Again, we will ignore the -100 labels.

```

1 import pandas as pd
2 from IPython.display import display
3
4 def evaluation_report(y_true, y_pred, get_only_acc=False):
5     # retrieve all tags in y_true
6     tag_set = set()
7     for sent in y_true:
8         for tag in sent:
9             tag_set.add(tag)
10    for sent in y_pred:
11        for tag in sent:
12            tag_set.add(tag)
13    tag_list = sorted(list(tag_set))
14
15    # count correct points
16    tag_info = dict()
17    for tag in tag_list:
18        tag_info[tag] = {'correct_tagged': 0, 'y_true': 0, 'y_pred': 0}
19
20    all_correct = 0
21    all_count = sum([len(sent) for sent in y_true])
22    speacial_tag = 0
23    for sent_true, sent_pred in zip(y_true, y_pred):
24        for tag_true, tag_pred in zip(sent_true, sent_pred):
25            # pass special token
26            if tag_true == -100 :
27                speacial_tag += 1
28                pass
29            if tag_true == tag_pred:
30                tag_info[tag_true]['correct_tagged'] += 1
31                all_correct += 1
32                tag_info[tag_true]['y_true'] += 1
33                tag_info[tag_pred]['y_pred'] += 1
34    print('speacial_tag :', speacial_tag) # delete number of special token from all_count
35    accuracy = (all_correct / (all_count - speacial_tag))
36
37    # get only accuracy for testing
38    if get_only_acc:
39        return accuracy
40
41    accuracy *= 100
42
43
44    # summarize and make evaluation result
45    eval_list = list()
46    for tag in tag_list:
47        eval_result = dict()
48        eval_result['tag'] = tag
49        eval_result['correct_count'] = tag_info[tag]['correct_tagged']
50        precision = (tag_info[tag]['correct_tagged'] / tag_info[tag]['y_pred']) * 100 if tag_info[tag]['y_pred'] else '-'
51        recall = (tag_info[tag]['correct_tagged'] / tag_info[tag]['y_true']) * 100 if (tag_info[tag]['y_true'] > 0) else 0
52        eval_result['precision'] = precision
53        eval_result['recall'] = recall
54        eval_result['f1_score'] = (2 * precision * recall) / (precision + recall) if (type(precision) is float and recall > 0) else 0
55
56        eval_list.append(eval_result)
57
58    eval_list.append({'tag': 'accuracy=%0.2f' % accuracy, 'correct_count': '', 'precision': '', 'recall': '', 'f1_score': ''})
59
60    df = pd.DataFrame.from_dict(eval_list)
61    df = df[['tag', 'precision', 'recall', 'f1_score', 'correct_count']]
62
63    display(df)
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956

```

```

8     predictions = torch.argmax(pred, dim=2)
9     # Append padded predictions to y_pred
10    y_pred.append(predictions.tolist()[0])

1 # check our prediction with label
2 # -100 is special tokens : [<s>, </s>, _]
3 print(y_pred[0])
4 print(y_test[0])

```

→ [37, 29, 39, 26, 26, 26, 37, 37, 26, 26, 26, 41, 37, 37, 26, 26, 39, 26, 37]
 [-100, 29, 39, 26, 26, 26, 37, -100, 26, 26, 26, 41, 37, -100, 26, 26, 39, 26, -100]

```

1 evaluation_report(y_test, y_pred)

```

→ speacial_tag : 21039

	tag	precision	recall	f1_score	correct_count
0	-100	-	0.0	-	0
1	0	-	0.0	-	0
2	1	69.254032	68.019802	68.631369	687
3	2	0.0	0.0	-	0
4	3	85.714286	10.344828	18.461538	6
5	4	85.714286	32.142857	46.753247	18
6	5	86.046512	21.387283	34.259259	37
7	6	55.775316	98.32636	71.176174	705
8	7	-	0.0	-	0
9	8	61.19403	72.84264	66.512167	287
10	10	76.292043	89.595376	82.410279	930
11	11	91.416309	93.421053	92.407809	426
12	12	67.1875	82.692308	74.137931	86
13	13	-	0.0	-	0
14	14	79.591837	75.728155	77.61194	78
15	15	87.987988	89.877301	88.92261	293
16	16	-	0.0	-	0
17	17	84.758364	91.935484	88.201161	228
18	18	69.768977	96.794872	81.089375	1057
19	19	-	0.0	-	0
20	20	100.0	58.823529	74.074074	10
21	21	91.994479	88.689288	90.311653	1333
22	22	70.754717	89.285714	78.947368	150
23	23	79.381443	81.052632	80.208333	77
24	24	95.357728	96.892342	96.11891	1746
25	25	83.167421	84.157509	83.659536	1838
26	26	86.697059	94.163685	90.276246	29477
27	27	76.826722	59.837398	67.276051	368
28	28	84.615385	75.862069	80.0	88
29	29	95.776772	98.755832	97.243492	635
30	30	74.888918	87.888784	78.58488	8888

✓ Other Pretrained model

```

-- -- -- -- --

```

In this section, we will experiment by fine-tuning other pretrained models, such as `airesearch/wangchanberta-base-wiki-newmm`, to see how about their performance.

Since each model uses a different word-tokenization method. for example, **`airesearch/wangchanberta-base-wiki-newmm` uses `newmm`**, while **`airesearch/wangchanberta-base-att-spm-uncased` uses `SentencePiece`**. please try fine-tuning and compare the performance of these models.

✓ #TODO 3

✓ Choose Pretrained Model

```
1 model_names = [  
2     'airesearch/wangchanberta-base-att-spm-uncased',  
3     'airesearch/wangchanberta-base-wiki-newmm',  
4     'airesearch/wangchanberta-base-wiki-ssg',  
5     'airesearch/wangchanberta-base-wiki-sefr',  
6     'airesearch/wangchanberta-base-wiki-spm',  
7 ]  
8  
9 #@title Choose Pretrained Model  
10 model_name = "airesearch/wangchanberta-base-wiki-newmm" #@  
11  
12 #create tokenizer  
13 tokenizer = Tokenizer(model_name).from_pretrained(  
14     f'{model_name}',  
15     revision='main',  
16     model_max_length=416,)  
17
```

model_name: airesearch/wangchanberta-base-wiki-newmm

```
➔ /usr/local/lib/python3.11/dist-packages/huggingface_hub/file_download.py:795: FutureWarning: `resume_download` is deprecated  
warnings.warn(  
    newmm.json: 100% 3.56M/3.56M [00:00<00:00, 17.7MB/s]  
    config.json: 100% 559/559 [00:00<00:00, 43.5kB/s]  
  
The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may  
The tokenizer class you load from this checkpoint is 'RobertaTokenizer'.  
The class this function is called from is 'ThaiWordsNewmmTokenizer'.  
The tokenizer class you load from this checkpoint is not the same type as the class this function is called from. It may  
The tokenizer class you load from this checkpoint is 'RobertaTokenizer'.  
The class this function is called from is 'ThaiWordsNewmmTokenizer'.
```

```
1 example = orchidl["train"][1899]  
2 print('sentence :', example["sentence"])  
3 tokenized_input = tokenizer([example["sentence"]], is_split_into_words=True)  
4 tokens = tokenizer.convert_ids_to_tokens(tokenized_input["input_ids"])  
5 print('tokens :', tokens)  
6 print('label tokens :', example["label_tokens"])
```

```
➔ sentence : โดยพิจารณาจากพจนานุกรมภาษาคู่ (bilingual transfer dictionary)  
tokens : ['<s>', 'โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', '<_>', '<unk>', '<_>', 'transfer', '<_>', 'dictionary']  
label tokens : ['โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', ' ', '(', 'bilingual transfer dictionary', ')']
```

It's the same problem as above.

****Warning: Can we use same function as above ?****

****Warning: Please beware of <unk>, an unknown word token.****

****Warning: Please be careful of "ํ", the 'am' vowel. WangchanBERTa's internal preprocessing replaces all "ํ" to "ิ" and "ำ" ****

```
1 def majority_vote_pos(examples):  
2  
3     #####  
4     # TO DO: Since the tokens from the output of the pretrained tokenizer  
5     # do not match the tokens in the label tokens of the dataset,  
6     # the task is to create a function to determine the POS tags of the tokens generated by the pretrained tokenizer.  
7     # This should be done by referencing the POS tags in the label tokens. If a token partially overlaps with others,  
8     # the POS tag from the segment with the greater number of characters should be assigned.  
9     #  
10    # Example :  
11    # "การประชุม" (9 chars) is formed from "การ" (3 chars) + "ประชุม" (6 chars).  
12    # "การ" has a POS tag of 21,  
13    # and "ประชุม" has a POS tag of 39.  
14    # Therefore, the POS tag for "การประชุม" is 39,  
15    # as "การประชุม" is derived more from the "ประชุม" part than from the "การ" part.  
16    #  
17    # 'ทางวิชาการ' (10 chars) is formed from 'ทาง' (3 chars) + 'วิชาการ' (7 chars)  
18    # "ทาง" has a POS tag of 26,  
19    # and "วิชาการ" has a POS tag of 2.  
20    # Therefore, the POS tag for "ทางวิชาการ" is 2,  
21    # as "ทางวิชาการ" is derived more from the "ทาง" part than from the "วิชาการ" part.  
22  
23    # FILL CODE HERE  
24    # tokenize word by pretrained tokenizer
```

```

25     tokenized_inputs = tokenizer([examples["sentence"]], is_split_into_words=True)
26     label_tokens = examples["label_tokens"]
27     pos_tags = examples["pos_tags"]
28     new_pos_result = []
29
30     new_tokens = tokenizer.convert_ids_to_tokens(tokenized_inputs["input_ids"])
31
32     label_idx = 0
33     i = 0
34
35     # for (token, tag) in zip(label_tokens, pos_tags):
36     #     print(f"'{token}': {tag}")
37     # print("-----")
38
39     try:
40         for t in new_tokens:
41             if t in ["<s>", "</s>", "_"]:
42                 new_pos_result.append(-100)
43                 continue
44             if t == "<unk>":
45                 new_pos_result.append(-100)
46                 # label_idx += 1
47                 continue
48
49             buffer = ""
50             weights = {}
51             t = t.replace(' ', '')
52             t = t.replace("<_>", " ")
53             if t[0] == "_":
54                 t = t[1:]
55             while label_tokens[label_idx][i] != t[0]:
56                 i += 1
57                 if i == len(label_tokens[label_idx]):
58                     label_idx += 1
59                     i = 0
60             j = 1
61             # print(f"BEGIN t: '{t}'")
62             for a in range(30):
63                 buffer += label_tokens[label_idx][i]
64                 # print(f"buffer: {buffer}, t: {t[:j]}")
65
66                 if buffer != t[:j]:
67                     buffer = ""
68                     j = 1
69                     continue
70
71                 j += 1
72                 if pos_tags[label_idx] not in weights:
73                     weights[pos_tags[label_idx]] = 0
74                 weights[pos_tags[label_idx]] += 1
75                 i += 1
76                 if i == len(label_tokens[label_idx]):
77                     label_idx += 1
78                     i = 0
79                 if buffer == t:
80                     break
81
82             # for (token, tag) in zip(new_tokens, new_pos_result):
83             #     print(f"'{token}': {tag}")
84             # print(weights)
85             max_key = max(weights, key=weights.get)
86             new_pos_result.append(max_key)
87     except IndexError:
88         print(examples)
89         for (token, tag) in zip(label_tokens, pos_tags):
90             print(f"'{token}': {tag}")
91         print("-----")
92         for (token, tag) in zip(new_tokens, new_pos_result):
93             print(f"'{token}': {tag}")
94         # print(new_tokens, "\n", label_tokens, "\n", pos_tags)
95         # print(new_pos_result)
96         raise ValueError("Age cannot be negative!")
97
98     tokenized_inputs['tokens'] = new_tokens
99     tokenized_inputs['labels'] = new_pos_result
100
101     return tokenized_inputs
102     #####

```

```

1 tokenized_orchid = orchidl.map(majority_vote_pos)

```

Parameter 'function'=<function majority_vote_pos at 0x7bf5800cb600> of the transform datasets.arrow_dataset.Dataset._map
WARNING:datasets.fingerprint:Parameter 'function'=<function majority_vote_pos at 0x7bf5800cb600> of the transform dataset

Map: 100% 18500/18500 [00:17<00:00, 1065.59 examples/s]

Map: 100% 4625/4625 [00:04<00:00, 1227.81 examples/s]

```
1 # hard test case
2 example = tokenized_orchid["train"][1899]
3 for i in example :
4     print(i, ":", example[i])
```

id : 1899
label_tokens : ['โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', ' ', '(', 'bilingual transfer dictionary', ')']
pos_tags : [25, 39, 38, 26, 26, 5, 37, 37, 26, 37]
sentence : โดยพิจารณาจากพจนานุกรมภาษาคู่ (bilingual transfer dictionary)
input_ids : [0, 80, 3973, 45, 12252, 3496, 592, 5, 3, 5, 30055, 5, 63190, 178, 2]
token_type_ids : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
attention_mask : [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]
tokens : ['<s>', 'โดย', 'พิจารณา', 'จาก', 'พจนานุกรม', 'ภาษา', 'คู่', '<_>', '<unk>', '<_>', 'transfer', '<_>', 'dictionary']
labels : [-100, 25, 39, 38, 26, 26, 5, 37, -100, 26, 26, 26, 26, 37, -100]

Choose Pretrained Model

```
1 model_names = [
2     'wangchanberta-base-att-spm-uncased',
3     'wangchanberta-base-wiki-newmm',
4     'wangchanberta-base-wiki-ssg',
5     'wangchanberta-base-wiki-sefr',
6     'wangchanberta-base-wiki-spm',
7 ]
8
9 #@title Choose Pretrained Model
10 model_name = "wangchanberta-base-wiki-newmm" #@param ["wan
11
12 #create model
13 model = AutoModelForTokenClassification.from_pretrained(
14     f"airesearch/{model_name}",
15     revision='main',
16     num_labels=47, id2label=id2label, label2id=label2id
17 )
18
```

model_name: wangchanberta-base-wiki-newmm

pytorch_model.bin: 100% 646M/646M [00:06<00:00, 160MB/s]

/usr/local/lib/python3.11/dist-packages/transformers/modeling_utils.py:463: FutureWarning: You are using `torch.load` wi
return torch.load(checkpoint_file, map_location="cpu")
Some weights of the model checkpoint at airesearch/wangchanberta-base-wiki-newmm were not used when initializing Roberta
- This IS expected if you are initializing RobertaForTokenClassification from the checkpoint of a model trained on anoth
- This IS NOT expected if you are initializing RobertaForTokenClassification from the checkpoint of a model that you exp
Some weights of RobertaForTokenClassification were not initialized from the model checkpoint at airesearch/wangchanberta
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
1 data_collator = DataCollatorForTokenClassification(tokenizer=tokenizer)
```

#TODO 4


Fine-tuning other pretrained model with our orchid corpus.

```
1 training_args = TrainingArguments(
2     #####
3     output_dir="pos-base-wiki-newmm",
4     learning_rate=2e-5,
5     per_device_train_batch_size=32,
6     per_device_eval_batch_size=32,
7     num_train_epochs=2,
8     weight_decay=0.01,
9     push_to_hub=True
10    #####
11 )
12
13 trainer = Trainer(
14     #####
15     model=model,
16     args=training_args,
17     train_dataset=tokenized_orchid["train"],
18     eval_dataset=tokenized_orchid["test"],
```

```

19 data_collator=data_collator,
20 compute_metrics=compute_metrics,
21 #####
22 )
23
24 trainer.train()

```

 /usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_deprecation.py:131: FutureWarning: 'Repository' (from 'hu
For more details, please read https://huggingface.co/docs/huggingface_hub/concepts/git_vs_http.
warnings.warn(warning_message, FutureWarning)
Cloning <https://huggingface.co/Jirayuwat12/pos-base-wiki-newmm> into local empty directory.
WARNING:huggingface_hub.repository:Cloning <https://huggingface.co/Jirayuwat12/pos-base-wiki-newmm> into local empty direc
/usr/local/lib/python3.11/dist-packages/transformers/optimization.py:411: FutureWarning: This implementation of AdamW is
warnings.warn(

[1158/1158 07:57, Epoch 2/2]

Step Training Loss

500	0.592600
1000	0.302100

TrainOutput(global_step=1158, training_loss=0.42407236774556584, metrics={'train_runtime': 477.395,
'train_samples_per_second': 77.504, 'train_steps_per_second': 2.426, 'total_flos': 912923649812664.0, 'train_loss':
0.42407236774556584, 'epoch': 2.0})

```

1 ##### EVALUATE YOUR MODEL #####
2 from transformers import AutoModelForTokenClassification
3
4 ## Load your fine-tuned model from Hugging Face
5 model = AutoModelForTokenClassification.from_pretrained("jirayuwat12/pos-base-wiki-newmm") ## your model path from Huggir
6
7 # prepare test set
8 test_data = tokenized_orchid["test"]
9
10 # labels for test set
11 y_test = []
12 for inputs in test_data:
13     y_test.append(inputs['labels'])
14 y_pred = []
15 device = 'cuda' if torch.cuda.is_available() else 'cpu'
16 for inputs in test_data:
17     text = inputs['sentence']
18     inputs = tokenizer(text, return_tensors="pt")
19     with torch.no_grad():
20         pred = model(**inputs).logits
21         predictions = torch.argmax(pred, dim=2)
22         # Append padded predictions to y_pred
23         y_pred.append(predictions.tolist()[0])
24
25 ##### EVALUATE YOUR MODEL #####
26 evaluation_report(y_test, y_pred)

```

↔ speacial_tag : 11485

	tag	precision	recall	f1_score	correct_count
0	-100	-	0.0	-	0
1	0	10.0	6.666667	8.0	1
2	1	76.66999	69.845595	73.098859	769
3	2	50.0	0.877193	1.724138	1
4	3	45.16129	25.0	32.183908	14
5	4	94.230769	79.032258	85.964912	49
6	5	78.030303	59.537572	67.540984	103
7	6	77.995643	91.094148	84.037559	358
8	7	-	0.0	-	0
9	8	70.241287	71.389646	70.810811	262
10	10	91.971665	88.422247	90.162037	779
11	11	93.461538	92.395437	92.92543	486

✓ #TODO 5

Compare the results between both models. Are they comparable? (Think about the ground truths of both models).

Propose a way to fairly evaluate the models.

Write your answer here :

It's not comparable because each model is trained using different dataset and different size

there are many approaches to make a fair comparison e.g.