

# LANGUAGE TRANSLATION WITH NN.TRANSFORMER & TORCH TEXT

## Data Sourcing & Processing

- how to use torchtext's inbuilt datasets
- tokenize a raw text sentence
- build vocabulary
- numericalize tokens into tensor

```
from torchtext.data.utils import get_tokenizer
from torchtext.vocab import build_vocab_from_iterator
from torchtext.datasets import Multi30k
from typing import Iterable, List
```

```
SRC_LANGUAGE = 'de'
TGT_LANGUAGE = 'en'
```

# Place-holders

```
token_transform = {}
```

```
vocab_transform = {}
```

# Create source & target language tokenizer.

# Make sure to install the dependencies.

# pip install -U spacy

# python -m spacy download en\_core\_web\_sm

# python -m spacy download de\_core\_news\_sm

```
token_transform[SRC_LANGUAGE] = get_tokenizer('spacy', language='de_core_news_sm')
```

```
token_transform[TGT_LANGUAGE] = get_tokenizer('spacy', language='en_core_web_sm')
```

# helper function to yield list of tokens

```
def yield_tokens(data_iter: Iterable, language: str) -> List[str]:
```

```
    language_index = {SRC_LANGUAGE: 0, TGT_LANGUAGE: 1}
```

```
    for data_sample in data_iter:
```

```
        yield token_transform[language](data_sample[language_index[language]])
```

get\_tokenizer() 함수

SRC / TGT의  
data\_sample에  
tokenizing.

Language index

data\_sample[0] or data\_sample[1]

# Define special symbols and indices

UNK\_IDX, PAD\_IDX, BOS\_IDX, EOS\_IDX = 0, 1, 2, 3

special\_symbols = ['<unk>', '<pad>', '<bos>', '<eos>']

for ln in [SRC\_LANGUAGE, TGT\_LANGUAGE]:

# Training data Iterator

train\_iter = Multi30k (split = 'train', language\_pair = (SRC\_LANGUAGE, TGT\_LANGUAGE))

# Create torchtext's Vocab object

Vocab\_transform[ln] = build\_vocab\_from\_iterator (yield\_tokens (train\_iter, ln),  
 min\_freq = 1,  
 specials = special\_symbols,  
 special\_first = True)

# Set UNK\_IDX as the default index.

# This index is returned when the token is not found.

# If not set, it throws RuntimeError when the queried token is not found in the vocabulary.

for ln in [SRC\_LANGUAGE, TGT\_LANGUAGE]:

Vocab\_transform[ln].set\_default\_index (UNK\_IDX)