

Working with HF Transformers Library

CS XXX: Introduction to Large Language Models

Contents

- `pipeline` object
- Raw text to prediction
- Tokenizer
- Model
- Postprocessing
- Finetuning
- Loading custom model into `pipeline`

The pipeline Object

- The `pipeline` object in Hugging Face is a high-level API designed to simplify working with popular machine learning tasks. It abstracts away much of the complexity involved in loading models, tokenizing inputs, and post-processing outputs. With just a few lines of code, you can use state-of-the-art models for a wide variety of tasks.

```
0 from transformers import pipeline
```

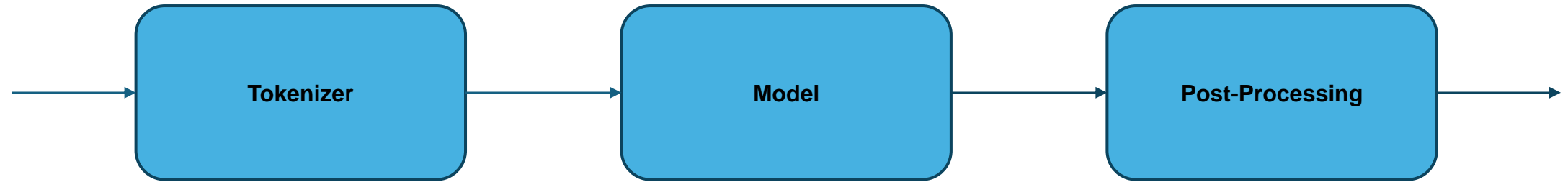
The pipeline Object

- With just a few lines of code, you can use state-of-the-art models for a wide variety of tasks.
- Sentiment Analysis

```
0 # Import the pipeline function from the transformers library
1 from transformers import pipeline
3
4 # Initialize a pipeline for sentiment analysis
5 classifier = pipeline("sentiment-analysis")
6
7 # Use the classifier to analyze the sentiment of a text
8 classifier("I love AI and want to know more about the current state of generative AI.")
```

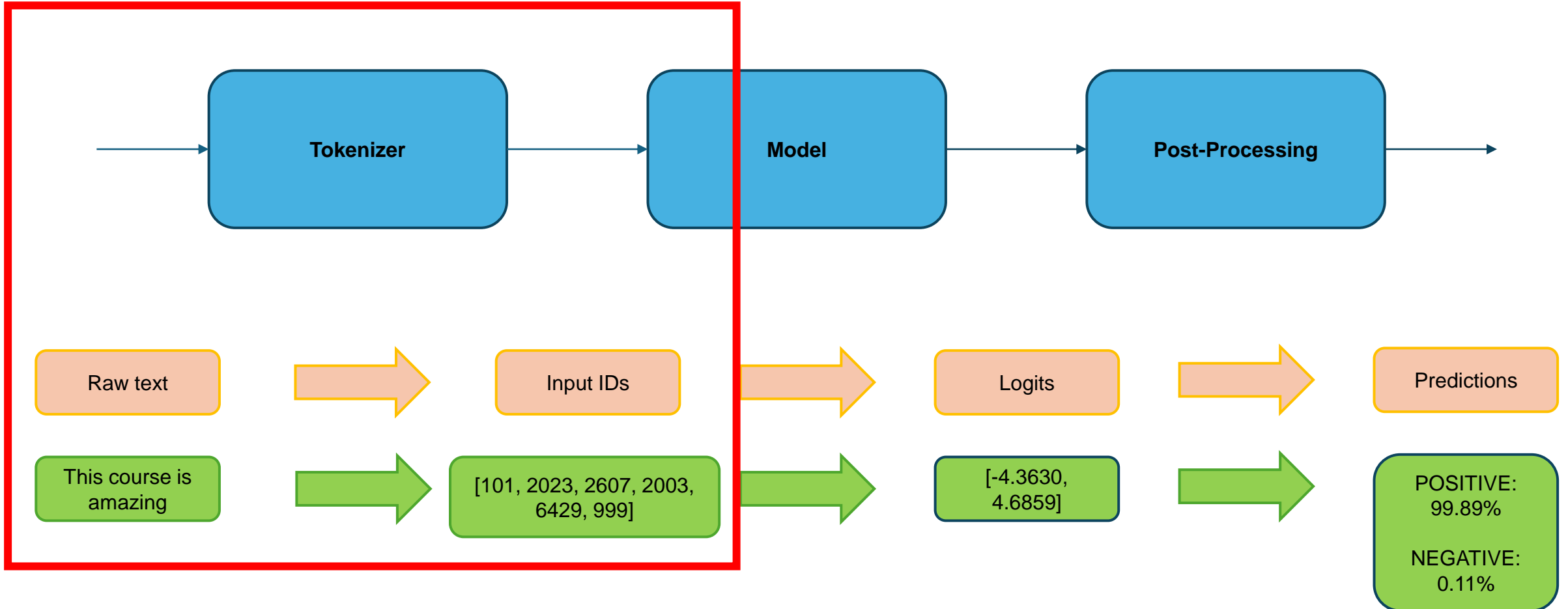
```
No model was supplied, defaulted to distilbert-base-uncased-finetuned-sst-2-English
[{'label': 'POSITIVE', 'score': 0.9966887831687927}]
```

Raw text to Predictions



Tokenizer

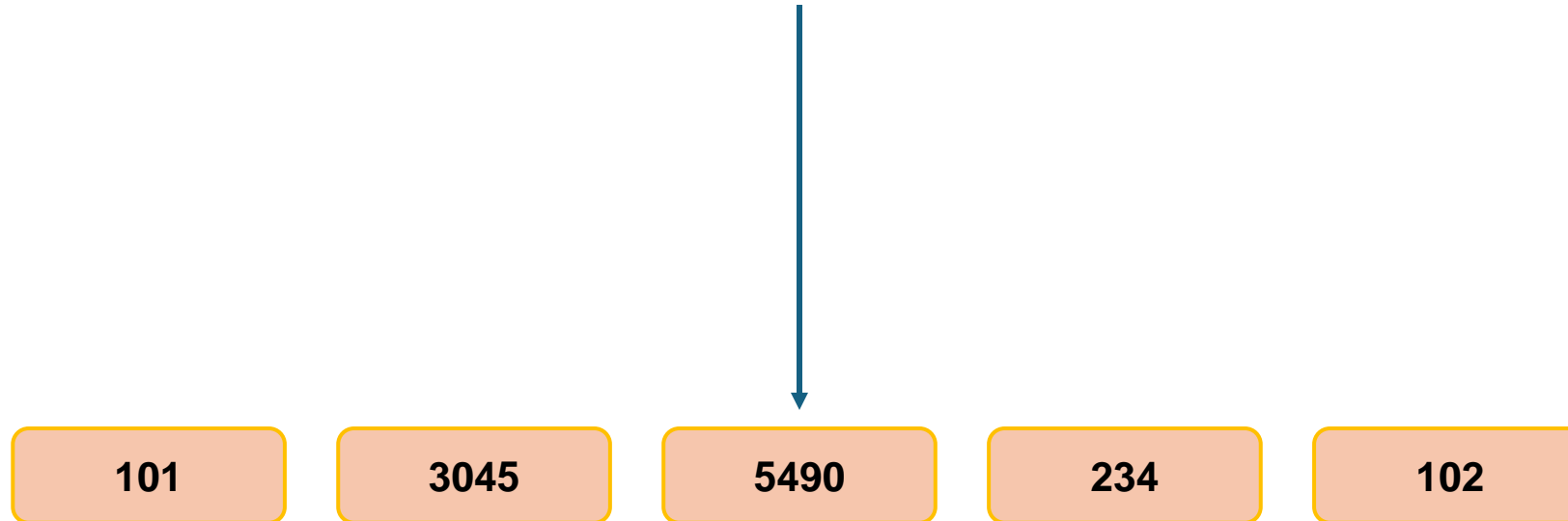
- The tokenizer is used to transform raw text to numbers



Tokenizer

- The tokenizer's objective is find a meaningful representation

“The HuggingFace headquarters are based in Brooklyn, New York”



Tokenizer

- Three different ways to tokenize
 - **Word Based:** Splitting a raw text into words
 - **Character Based:** Splitting a raw text into characters
 - **Sub-word Based:** Splitting a raw text into subwords

Tokenizer

- **Word Based:** Splitting a raw text into words

Split on spaces

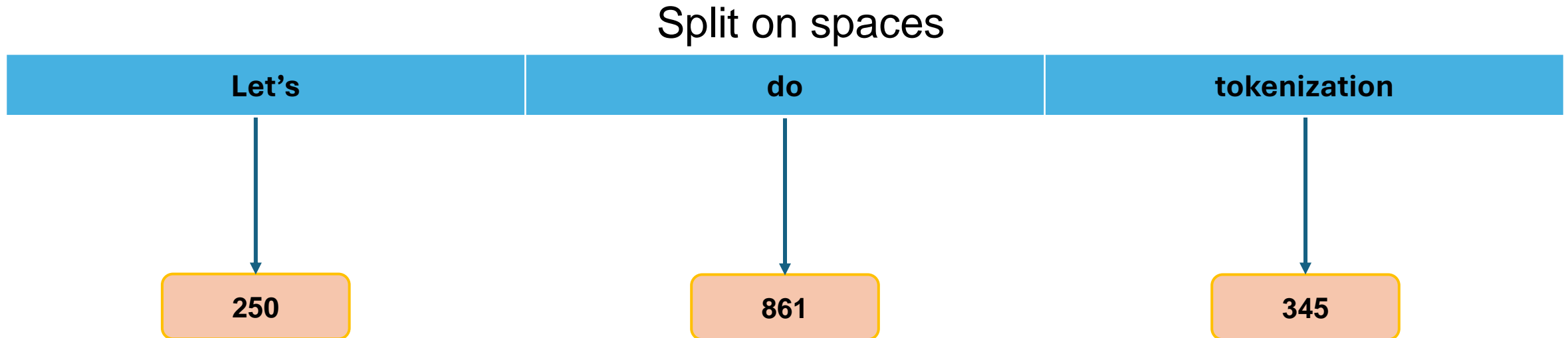
Let's

do

tokenization

Tokenizer

- **Word Based:** Splitting a raw text into words
 - Each word has a specific ID



Tokenizer

- **Word Based:** Splitting a raw text into words
 - Very similar words have entirely different meanings

the	→	1
of	→	2
and	→	3
to	→	4
in	→	5
was	→	6
the	→	7
is	→	8
for	→	9
as	→	10
on	→	11
with	→	12
that	→	13
dog	→	14
dogs	→	15

Tokenizer

- **Word Based:** Splitting a raw text into words

- The vocabulary can end up being very large

the	→	1
of	→	2
and	→	3
to	→	4
in	→	5
was	→	6
the	→	7
is	→	8
for	→	9
as	→	10
on	→	11
with	→	12
That	→	13
...	→	...

Malapropism	→	170,000
-------------	---	----------------

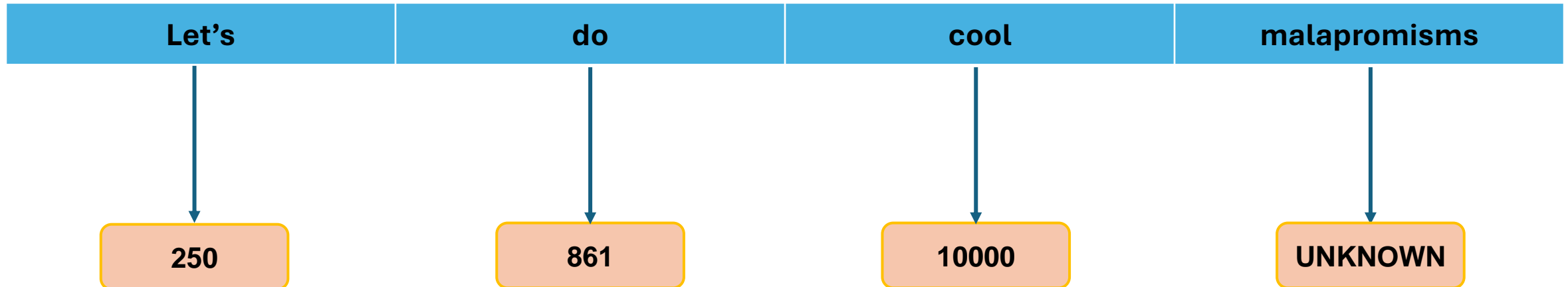
Tokenizer

- **Word Based:** Splitting a raw text into words
 - The vocabulary can end up being very large
 - Large vocabularies result in heavy models
 - We can limit the amount of words we add to the vocabulary

the	→	1
of	→	2
and	→	3
to	→	4
in	→	5
was	→	6
the	→	7
is	→	8
for	→	9
as	→	10
on	→	11
with	→	12
that	→	13
...	→	...
hug	→	10,000

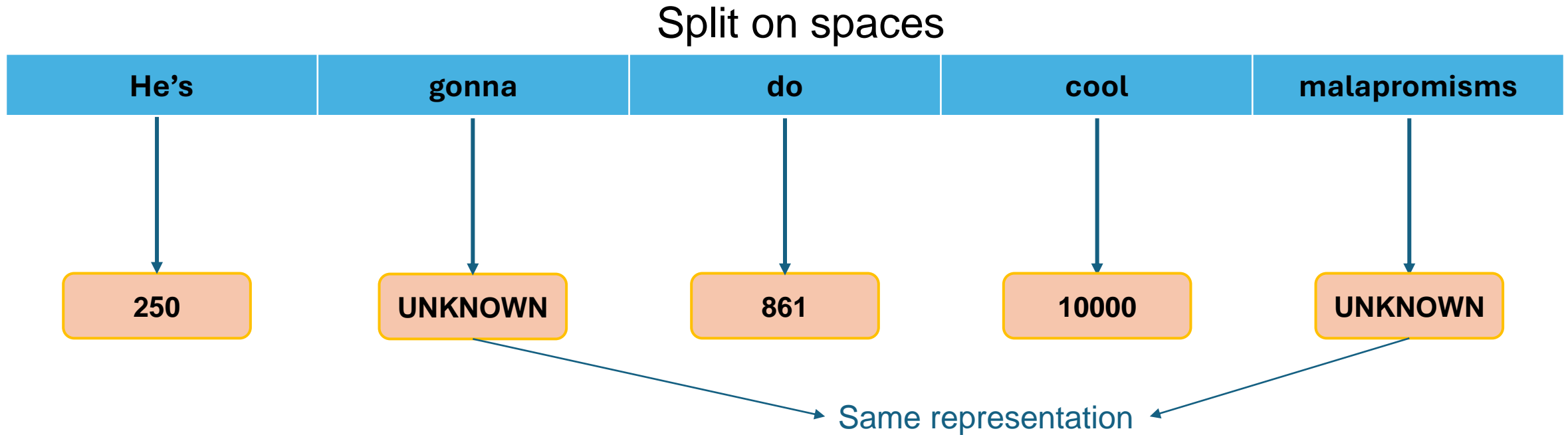
Tokenizer

- **Word Based:** Splitting a raw text into words
 - Out of vocabulary words result in a loss of information



Tokenizer

- **Word Based:** Splitting a raw text into words
 - Out of vocabulary words result in a loss of information



Tokenizer

- **Character Based:** Splitting a raw text into characters
 - Out of vocabulary words result in a loss of information

Split on characters

L	e	t	'	s	d	o	t	o	k	e	n	i	z	t	i	o	n	!
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Tokenizer

- **Character Based:** Splitting a raw text into characters

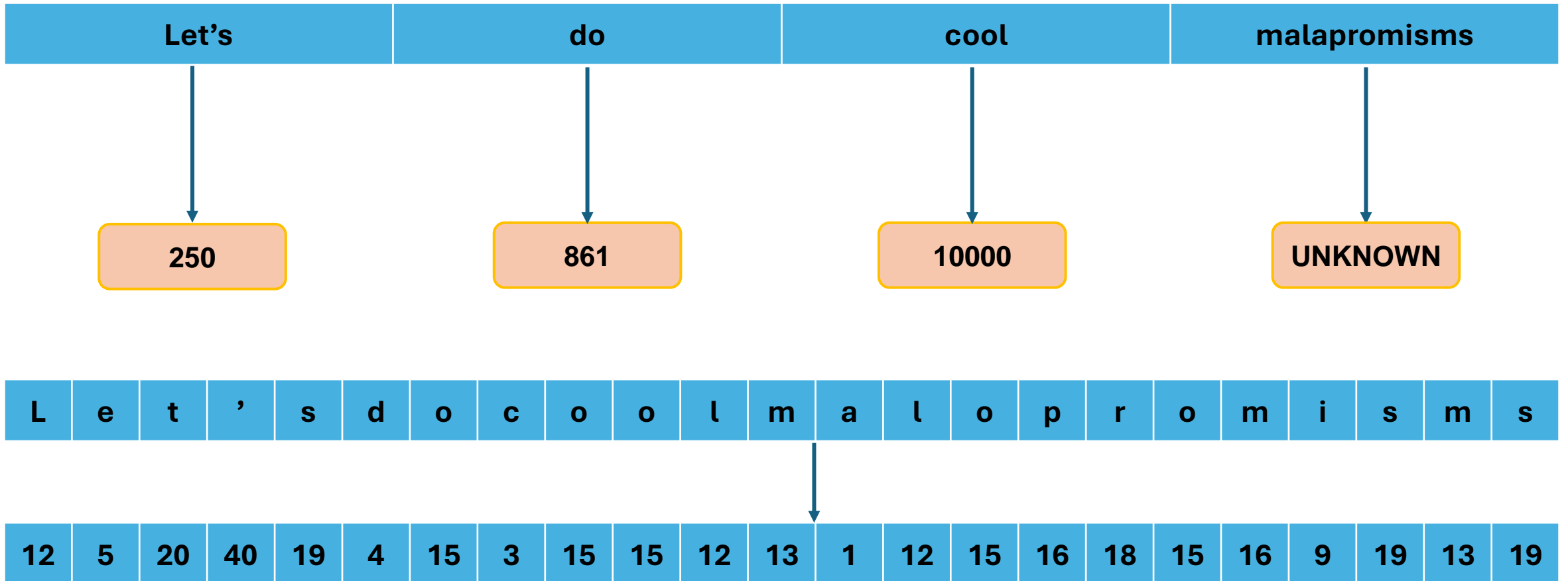
- Vocabularies are slimmer

Character-based vocabulary

a	→	1
b	→	2
c	→	3
d	→	4
e	→	5
f	→	6
g	→	7
...	→	...
1	→	27
2	→	28
3	→	29
...	→	...
!	→	37
...	→	...
à	→	256

Tokenizer

- **Character Based:** Splitting a raw text into characters
 - Fewer out-of-vocabulary words



Tokenizer

- **Sub-word Based:** Splitting a raw text into subwords
 - Middle ground between word and character based algorithms

Word-based vocabulary

Very large vocabularies

Large quantity of out of vocabulary tokens

Loss of meaning across very similar words

Char-based vocabulary

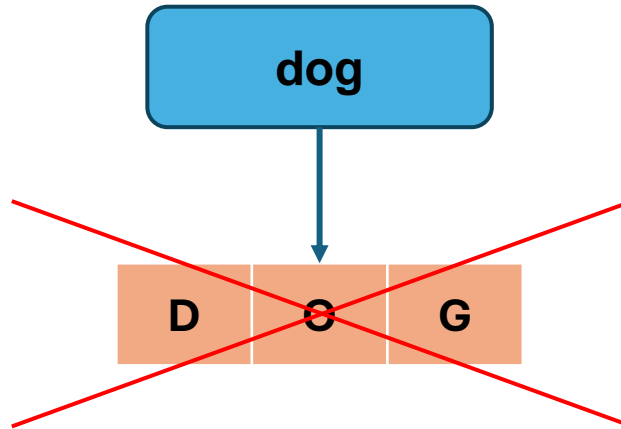
Very long sequences

Less meaningful individual tokens

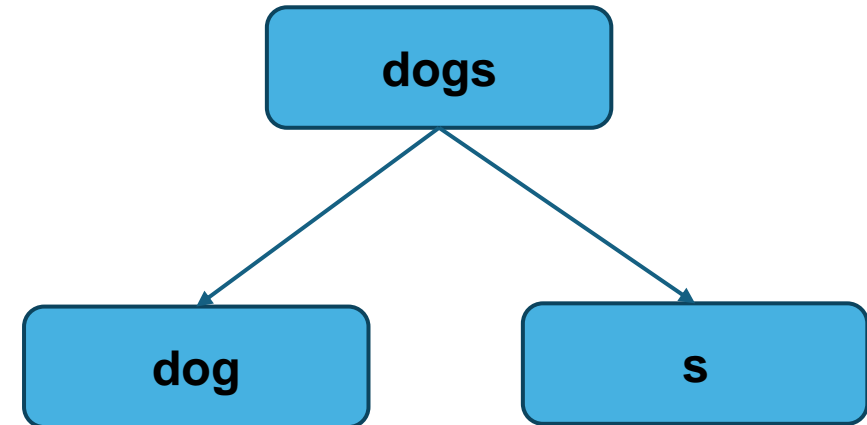
Tokenizer

- **Sub-word Based:** Splitting a raw text into subwords
 - Middle ground between word and character based algorithms

Frequently used words should not be split into smaller subwords



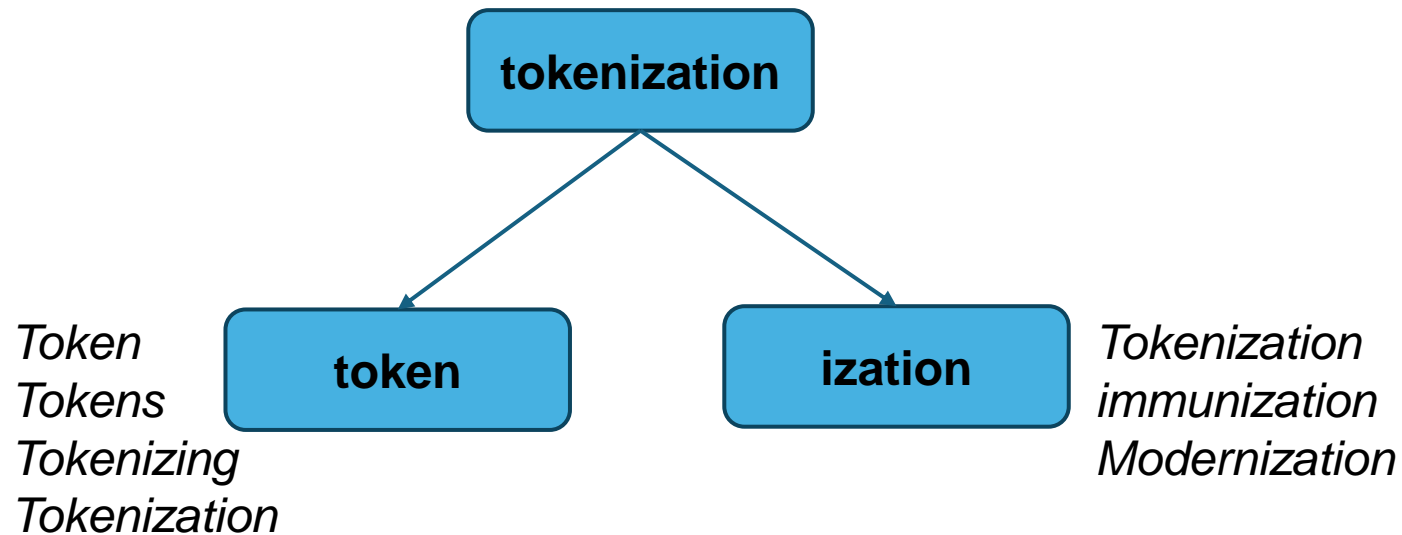
Rare words should be decomposed into meaningful subwords



Tokenizer

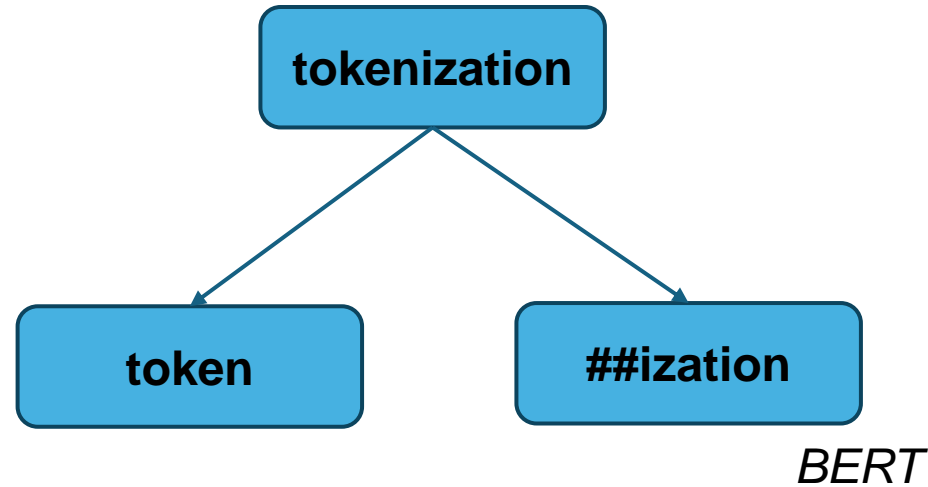
- **Sub-word Based:** Splitting a raw text into subwords
 - Subwords help identify similar syntactic or semantic situations in text

Rare words should be decomposed into meaningful subwords



Tokenizer

- **Sub-word Based:** Splitting a raw text into subwords
 - Subword tokenization algorithms can identify start of word tokens



- **##** is a marker to indicate that the subword is part of a word and not a standalone token.

Tokenizer

- **Sub-word Based:** Splitting a raw text into subwords
 - Most models obtaining state-of-the-art results in English today use some kind of subword-tokenization algorithm

WordPiece

BERT, DistilBERT

Unigram

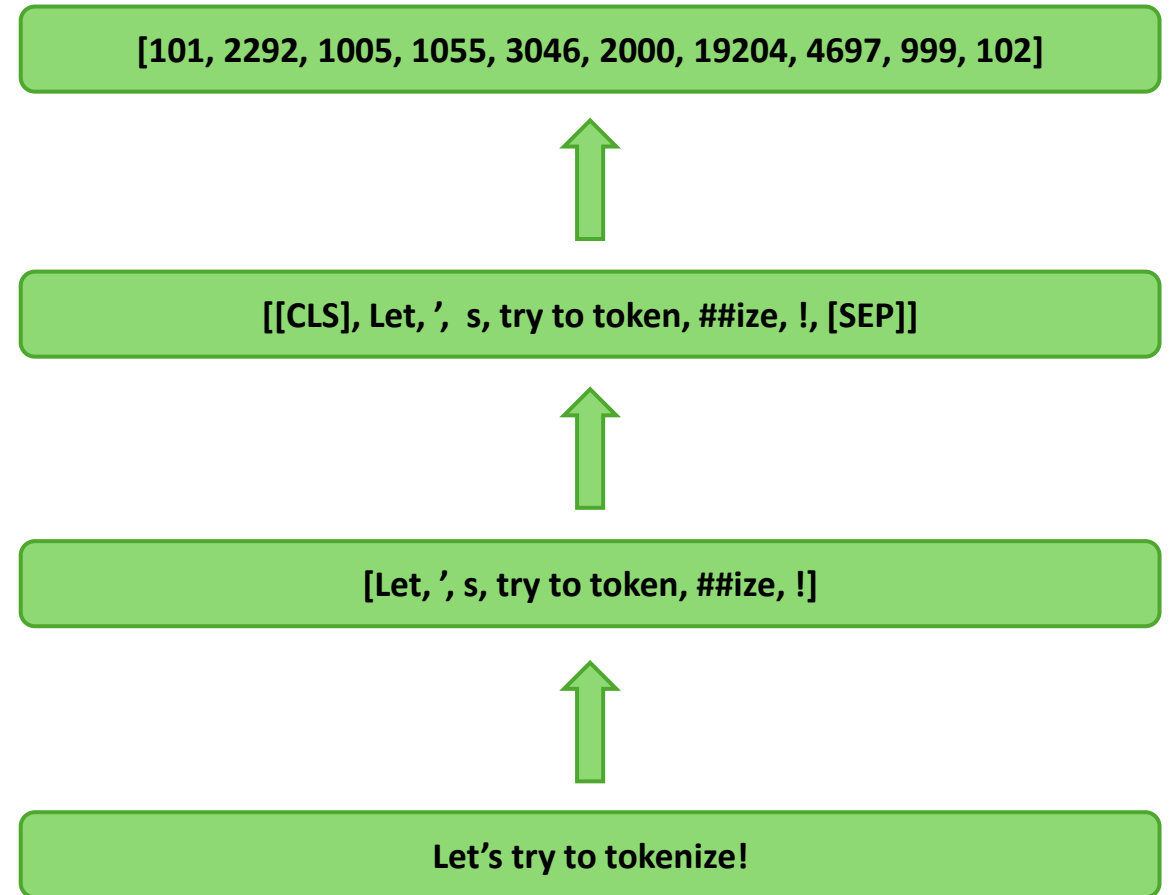
XLNet, ALBERT

**Byte-Pair
Encoding**

GPT2, RoBERTa

Tokenizer

- Tokenization Pipeline



Tokenizer

- The first step of the pipeline is to split the text into tokens

```
0 from transformers import AutoTokenizer
1 tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
2 tokens = tokenizer.tokenize("Let's try to tokenize!")
3 print(tokens)
```

```
['let', "'", 's', 'try', 'to', 'token', '##ize', '!']
```

Tokenizer

- Each token has a unique ID

```
0 from transformers import AutoTokenizer
1 tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
2 tokens = tokenizer.tokenize("Let's try to tokenize!")
3 input_ids = tokenizer.convert_tokens_to_ids(tokens)
4 print(input_ids)
```

```
[2292, 1005, 1055, 3046, 2000, 19204, 4697, 999]
```

```
0 decoded_string = tokenizer.decode([2292, 1005, 1055, 3046, 2000, 19204, 4697, 999])
1 print(decoded_string)
```

```
['let', "'", 's', 'try', 'to', 'token', '##ize', '!']
```

Tokenizer

- The tokenizer adds special tokens the model expects

```
0 from transformers import AutoTokenizer
1 tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
2 tokens = tokenizer.tokenize("Let's try to tokenize!")
3 input_ids = tokenizer.convert_tokens_to_ids(tokens)
4 print(input_ids)
6 final_inputs = tokenizer.prepare_for_model(input_ids)
7 print(final_inputs["input_ids"])
```

```
[2292, 1005, 1055, 3046, 2000, 19204, 4697, 999]
[101, 2292, 1005, 1055, 3046, 2000, 19204, 4697, 999, 102]
```

```
0 print(tokenizer.decode([101, 2292, 1005, 1055, 3046, 2000, 19204, 4697, 999, 102]))
```

```
['[CLS]', 'let', "'", 's', 'try', 'to', 'token', '##ize', '!', '[SEP]']
```

Tokenizer

- A tokenizer takes texts as inputs and outputs numbers the associated model can make sense of

```
0 from transformers import AutoTokenizer
1 tokenizer = AutoTokenizer.from_pretrained("bert-base-cased")
2
3 print(tokenizer("Using a Transformer network is simple"))
```

```
{'input_ids': [101, 7993, 170, 11303, 1200, 2443, 1110, 3014, 102],
 'token_type_ids': [0, 0, 0, 0, 0, 0, 0, 0, 0],
 'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1]}
```

Tokenizer

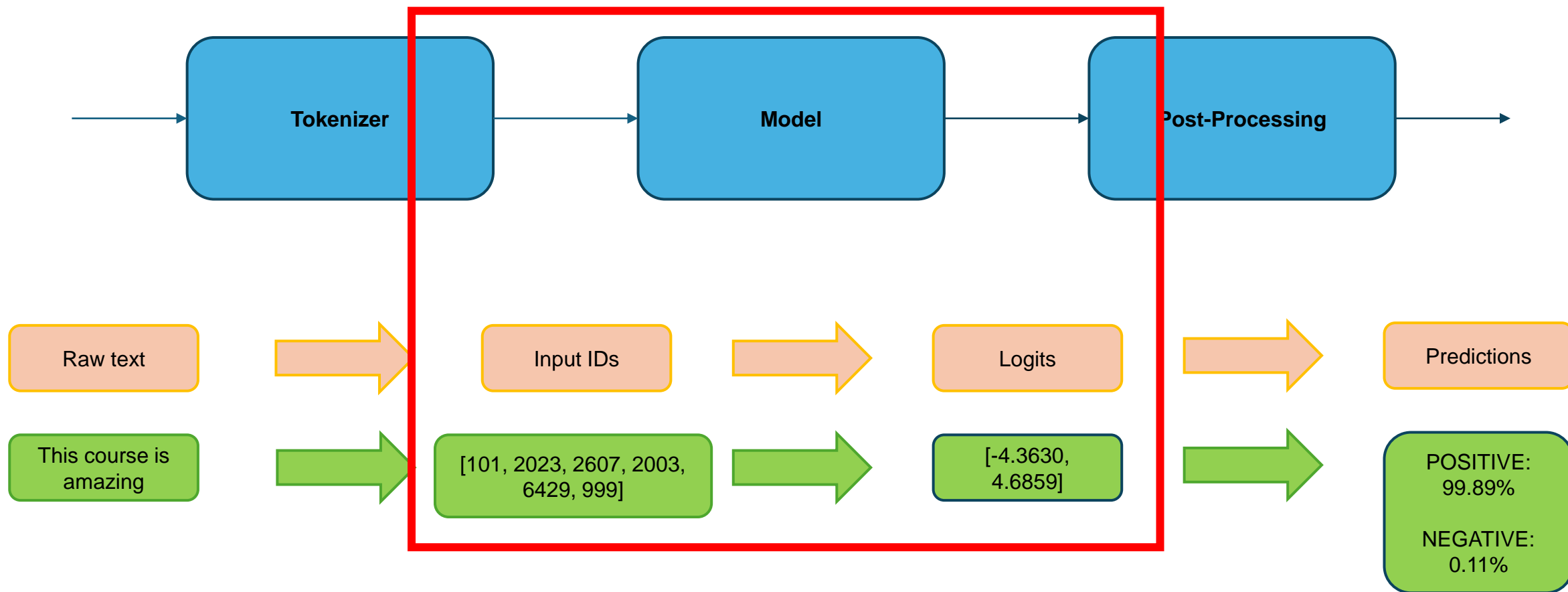
- A tokenizer takes texts as inputs and outputs numbers the associated model can make sense of

```
0 from transformers import AutoTokenizer
1 checkpoint = "distilbert-base-uncased-finetuned-sst-2-english"
2 tokenizer = AutoTokenizer.from_pretrained(checkpoint)
3 raw_inputs = [
4     "I've been waiting for a HuggingFace course my whole life.",
5     "I hate this so much!",
6 ]
7 inputs = tokenizer(raw_inputs, padding=True, truncation=True, return_tensors="pt")
8 print(inputs)
```

```
{
  'input_ids': tensor([
    [ 101, 1045, 1005, 2310, 2042, 3403, 2005, 1037, 17662, 12172, 2607, 2026, 2878, 2166, 1012, 102],
    [ 101, 1045, 5223, 2023, 2061, 2172, 999, 102, 0, 0, 0, 0, 0, 0, 0, 0]
  ]),
  'attention_mask': tensor([
    [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1],
    [1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0]
  ])
}
```

Model

- The model takes the input IDs and processes them through its layers to generate logits. Logits are unnormalized scores representing the likelihood of different outcomes.



Model

- The AutoModel class loads a model without a task-specific pretrained head

```
0 from transformers import AutoModel
1 checkpoint = "distilbert-base-uncased-finetuned-sst-2-english"
2 model = AutoModel.from_pretrained(checkpoint)
3 outputs = model(**inputs)
4 print(outputs.last_hidden_state.shape)
```

`torch.Size([2, 16, 768])`

batch size

sequence length

hidden size

Model

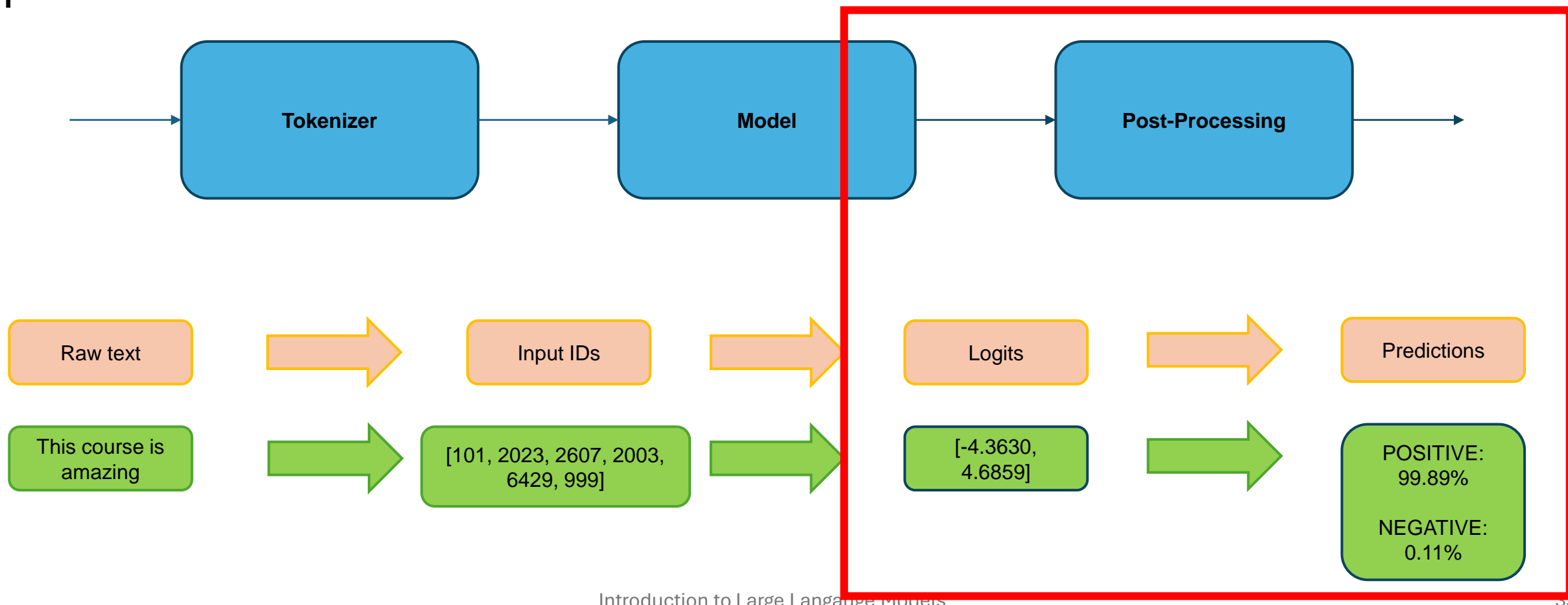
- Each AutoModelFor **X** class loads a model suitable for a specific task

```
0 from transformers import AutoModelForSequenceClassification
1 checkpoint = "distilbert-base-uncased-finetuned-sst-2-english"
2 model = AutoModelForSequenceClassification.from_pretrained(checkpoint)
3 outputs = model(**inputs)
4 print(outputs.logits)
```

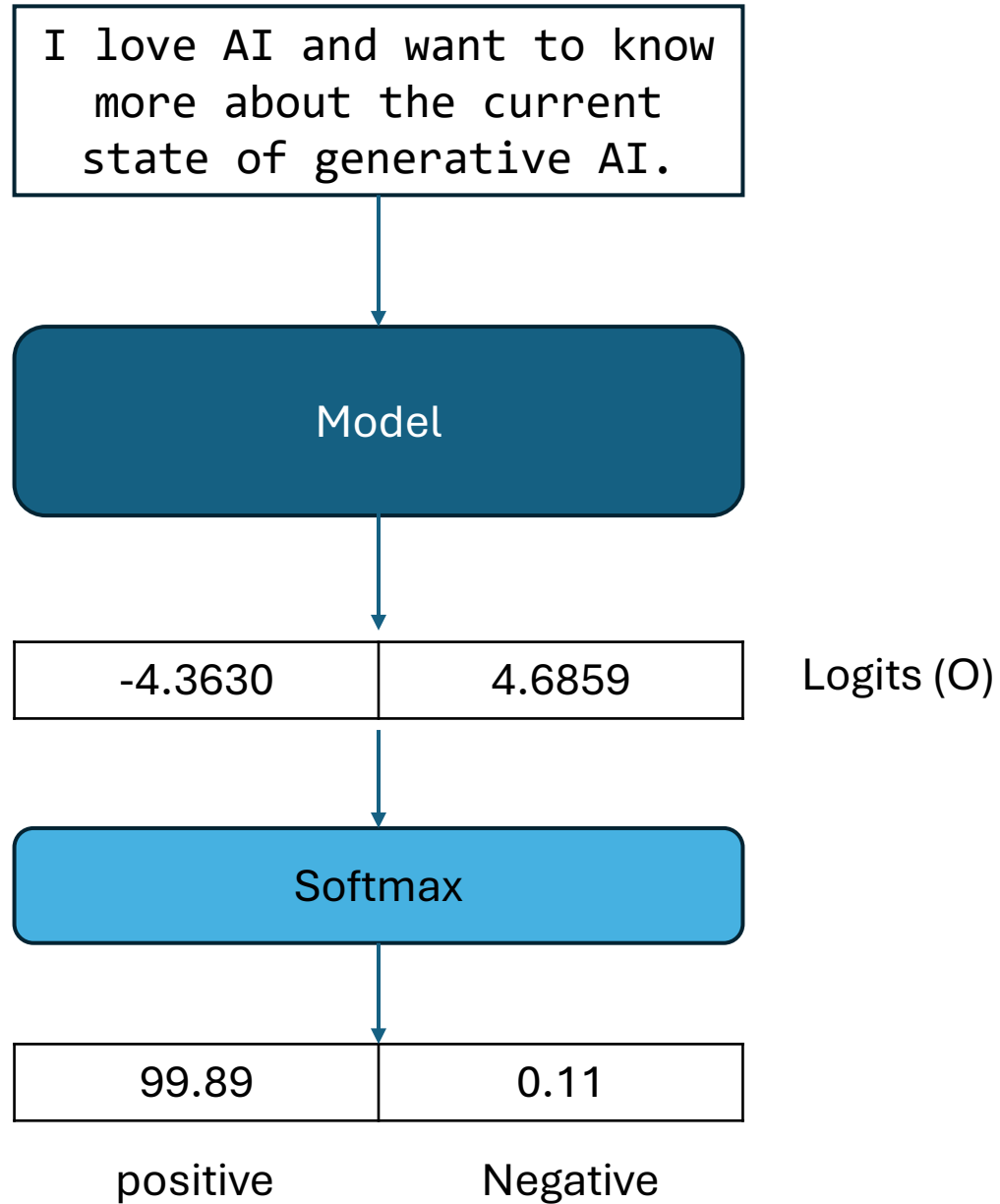
```
tensor([[ -1.5607,   1.6123],
        [  4.1692,  -3.3464]], grad_fn=<AddmmBackward>)
```


Post-Processing

- The logits are transformed into meaningful predictions. This involves normalizing the logits using methods like softmax to produce probabilities



Post-Processing



Finetuning

1. Load model for finetuning

```
0 from transformers import AutoModelForSequenceClassification, AutoTokenizer
1 # 1. Load the model and tokenizer
2 model_name = "bert-base-uncased"
3 model = AutoModelForSequenceClassification.from_pretrained(model_name, num_labels=2) #
4 Binary classification
5 tokenizer = AutoTokenizer.from_pretrained(model_name)
```

Finetuning

2. Prepare Data

```
0 from datasets import load_dataset
1 # 2. Prepare the dataset
2 # Example dataset: SST2 for sentiment classification
3 dataset = load_dataset("glue", "sst2")
4
5 # Tokenize the dataset
6 def tokenize_function(examples):
7     return tokenizer(examples['sentence'], padding="max_length", truncation=True)
8
9 tokenized_datasets = dataset.map(tokenize_function, batched=True)
10
11 # Split into train, validation, and test datasets
12 train_dataset = tokenized_datasets["train"]
13 val_dataset = tokenized_datasets["validation"]
14 test_dataset = tokenized_datasets["test"]
```

Finetuning

3. Define Training Arguments

```
0 from transformers import TrainingArguments
1 # 3. Define Training Arguments
2 training_args = TrainingArguments(
3     output_dir="./results",           # output dir for model predictions & checkpoints
4     evaluation_strategy="epoch",       # Evaluate every epoch
5     learning_rate=2e-5,                # Learning rate
6     per_device_train_batch_size=16,    # Batch size per GPU/CPU for training
7     per_device_eval_batch_size=64,     # Batch size for evaluation
8     num_train_epochs=3,                # Number of training epochs
9     weight_decay=0.01,                 # Strength of weight decay
10    logging_dir="./logs",               # Directory for storing logs
11    logging_steps=10,                   # Log every 10 steps
12    load_best_model_at_end=True,         # Load the best model when finished training
13    metric_for_best_model="accuracy",   # Metric to use for best model
14 )
```

Finetuning

4. Train the model and save

```
0 from transformers import Trainer
1 # Define Trainer
2 trainer = Trainer(
3     model=model,                # The model to train
4     args=training_args,         # Training arguments
5     train_dataset=train_dataset, # Training dataset
6     eval_dataset=val_dataset,   # Evaluation dataset
7     tokenizer=tokenizer,        # Tokenizer
8 )
9
10 # 4. Train the model
11 trainer.train()
12
13 # Save the fine-tuned model
14 trainer.save_model("./fine_tuned_bert")
15 tokenizer.save_pretrained("./fine_tuned_bert")
```

Finetuning

5. Test the model

```
0 # 5. Testing the model on the test set
1 test_results = trainer.evaluate(eval_dataset=test_dataset)
2 print("Test Results:", test_results)
```

Loading custom model or tokenizer into pipeline

- If you have a custom model or tokenizer, you can load it into pipeline for inference as follows

```
0 from transformers import pipeline, AutoModelForSequenceClassification, AutoTokenizer
1
2
3 # Load a custom model and tokenizer
4 model_path = "./fine_tuned_bert" # This is where your model and tokenizer are saved
5 model = AutoModelForSequenceClassification.from_pretrained(model_path)
6 tokenizer = AutoTokenizer.from_pretrained(model_path)
7
8 # Initialize the pipeline with the custom model and tokenizer
9 classifier = pipeline("text-classification", model=model, tokenizer=tokenizer)
10
11 # Use the pipeline
12 text = "This is a fantastic example!"
13 output = classifier(text)
14 print(output)
```


References

- Hugging Face. The Hugging Face Natural Language Processing Course.