

A practical introduction to Large Language Models

Nathanaël Fijałkow
CNRS, LaBRI, Bordeaux



LaBRI

université
de BORDEAUX

GOAL OF THIS COURSE

- Understand how LLMs work, what they can be used for
- Being able to deploy LLMs using Hugging Face
- Practice RAG and Agents

PART 1:

WHAT IS A LARGE LANGUAGE MODEL?

- Language models
- Tokenization
- Embeddings
- The attention mechanism
- Encoder / decoder
- Pre-training
- Fine-tuning

LANGUAGE MODELS

WHAT IS A LANGUAGE MODEL (LM)?

Input: a sentence (as a sequence of tokens)

Output: predict the next token

Basic examples:

- *Markov chain* is a LM, it gives a probabilistic distribution over the next token given the last token
- Naturally extended to *n-grams*: use the $(n-1)$ last tokens

N-GRAMS ARE LIMITED

Number of parameters: *vocab_size context_length*

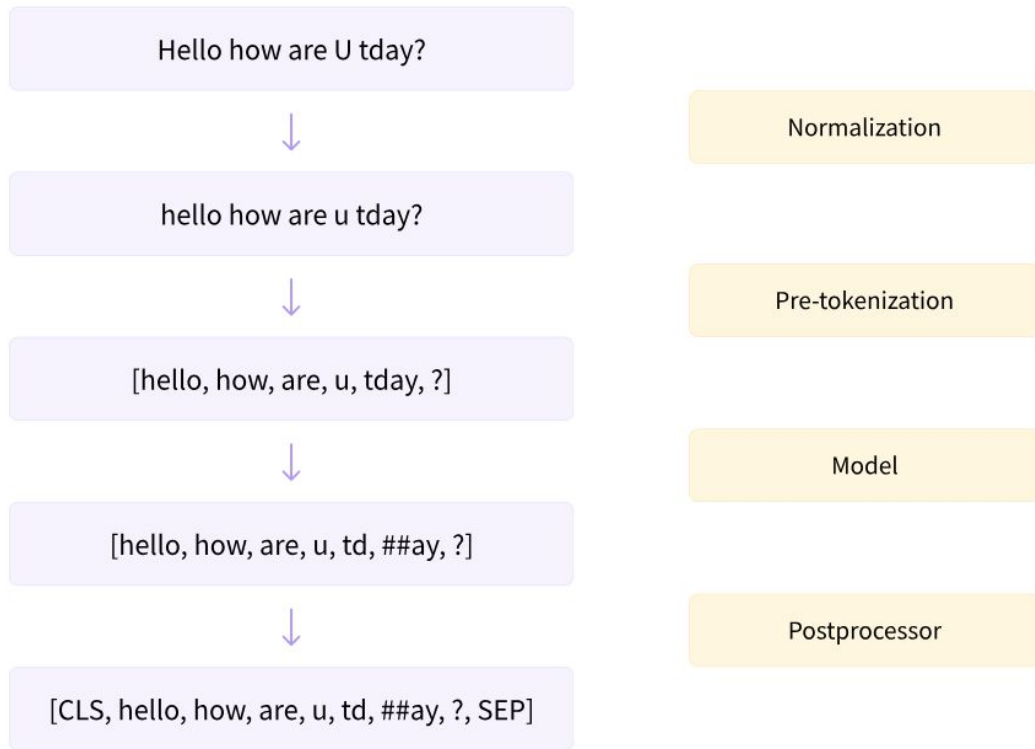
vocab_size: total number of tokens

context_length: number of tokens considered for prediction

Think of it as a very large matrix...

TOKENIZATION

TOKENIZATION: DECOMPOSING A SENTENCE INTO A SEQUENCE OF TOKENS



Every single explanation you will ever see about Language Models use **words**, **BUT** in reality the unit object is **tokens**

WORDS != TOKENS

WHAT IT ACTUALLY LOOKS LIKE:

```
test = "hello world"
test_encoded = tokenizer.encode(test)
test_encoded, [tokenizer.decode([x]) for x in test_encoded], tokenizer.decode(test_encoded)

([258, 285, 111, 492], ['he', 'll', 'o', ' world'], 'hello world')
```

Bottom line: at this point, we have converted a text into a sequence of integers (which represent tokens).

GPT-2 has 50,257 tokens

IN PRACTICE

Tiktoken implements BPE (one of the two most popular tokenization algorithms):

<https://github.com/openai/tiktoken>

EMBEDDINGS

THE 2003 (SILENT) BREAKTHROUGH

Journal of Machine Learning Research 3 (2003) 1137–1155

Submitted 4/02; Published 2/03

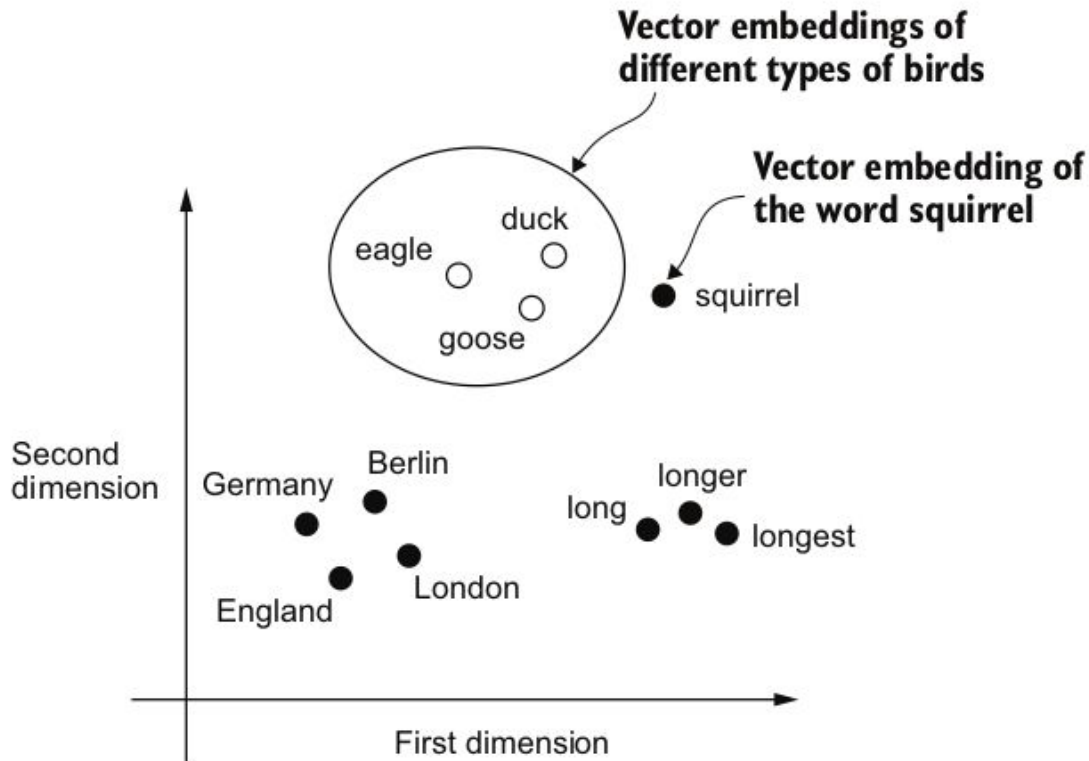
A Neural Probabilistic Language Model

Yoshua Bengio
Réjean Ducharme
Pascal Vincent
Christian Jauvin

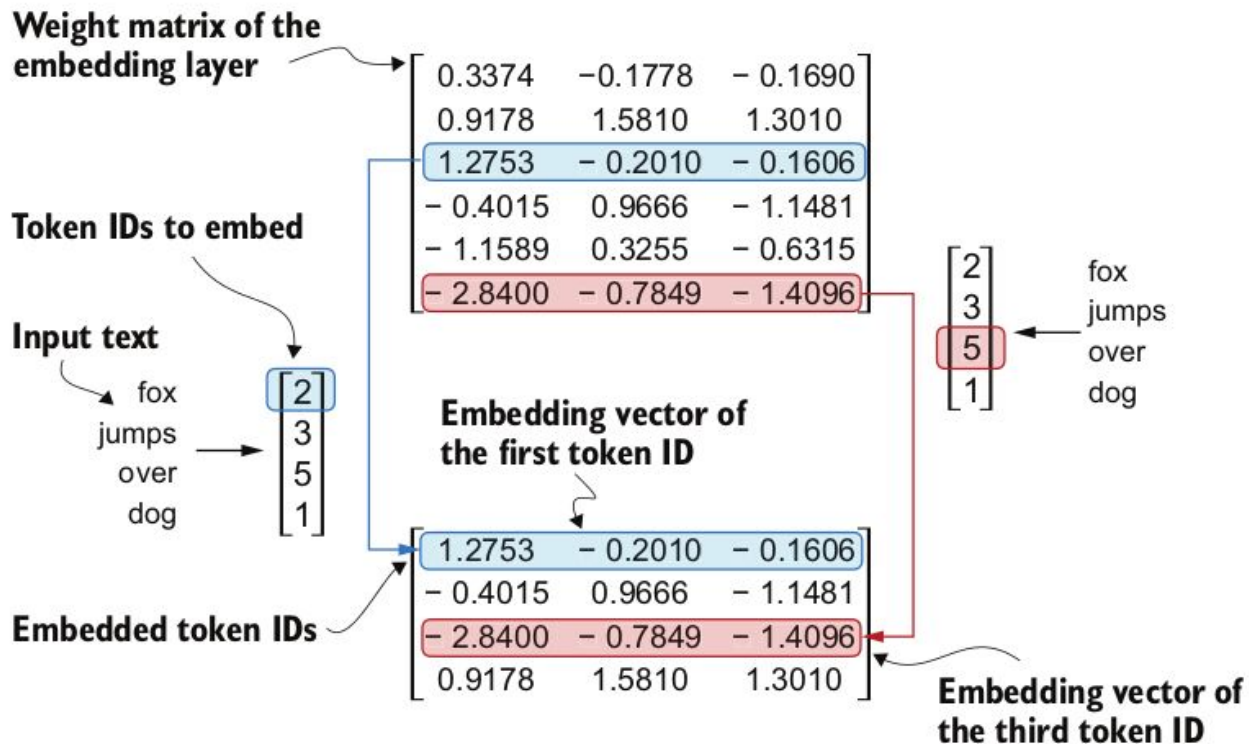
Département d'Informatique et Recherche Opérationnelle
Centre de Recherche Mathématiques
Université de Montréal, Montréal, Québec, Canada

BENGIOY@IRO.UMONTREAL.CA
DUCHARME@IRO.UMONTREAL.CA
VINCENTP@IRO.UMONTREAL.CA
JAUVINC@IRO.UMONTREAL.CA

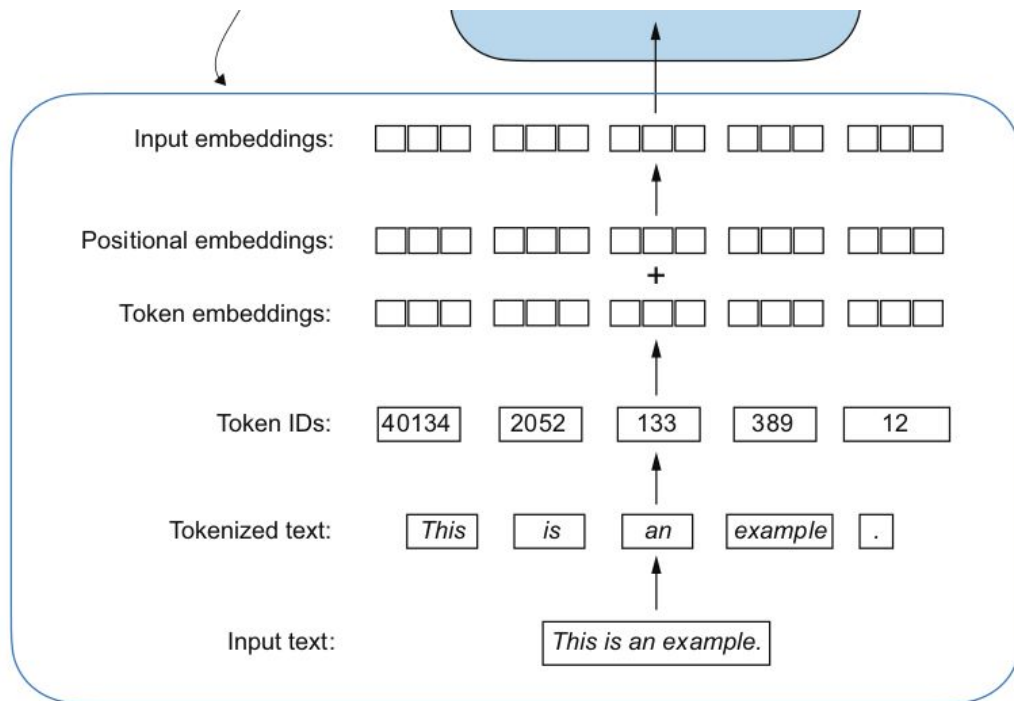
KEY IDEA: EMBEDDINGS



FROM TEXT TO VECTORS



Bottom line: at this point, we have converted a text into a sequence of (floating point) vectors. These are (almost) the inputs for our models.



STATISTICS

The smallest GPT-2 models (117M and 125M parameters) use an embedding size of 768 dimensions.

The largest GPT-3 model (175B parameters) uses an embedding size of 12,288 dimensions.

THE ATTENTION MECHANISM

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* †

University of Toronto

aidan@cs.toronto.edu

Łukasz Kaiser*

Google Brain

lukaszkaiser@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com

ATTENTION IS ALL YOU NEED

The paper came in 2017, in a wave of more and more complicated architectures around recurrent neural networks (RNNs), aiming at dealing with long contexts.

It does not do anything radically new: it says that “attention mechanism is enough to enable long contexts”.

A SIDE-NOTE

OpenAI scientist Noam Brown:

“The incredible progress in AI over the past five years
can be summarized in one word: scale.”

Recently, older architectures (LSTMs) reached similar performances as Transformers...

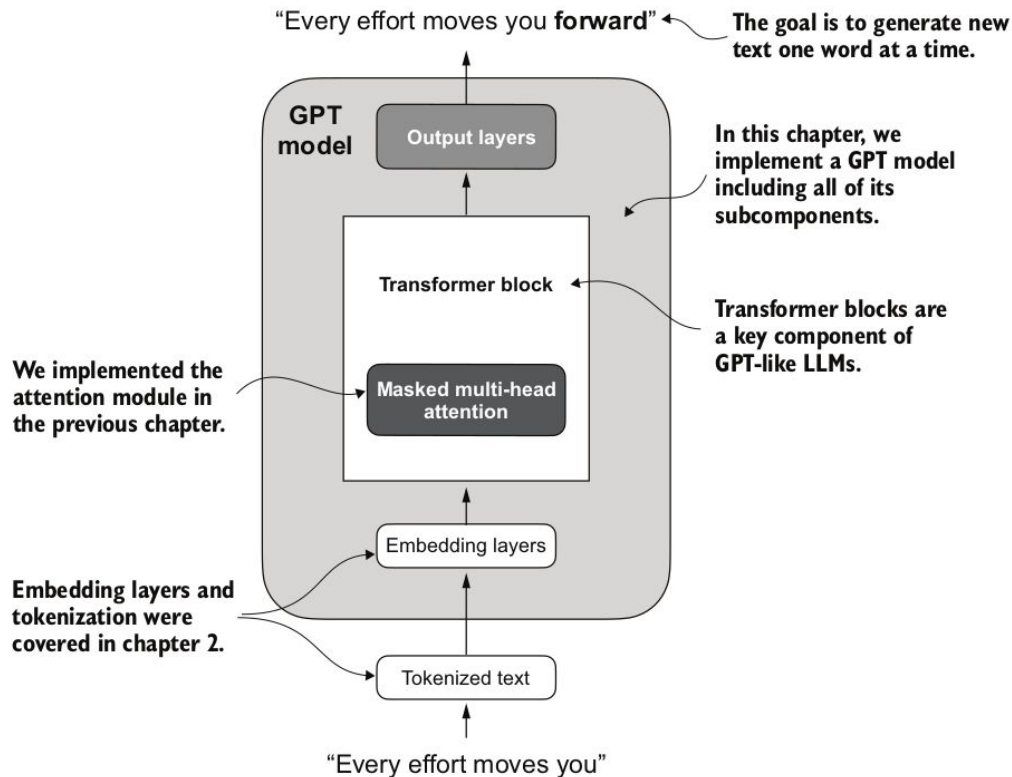
A SELF-ATTENTION HEAD

Input: an embedding vector $x(i)$ for each token i

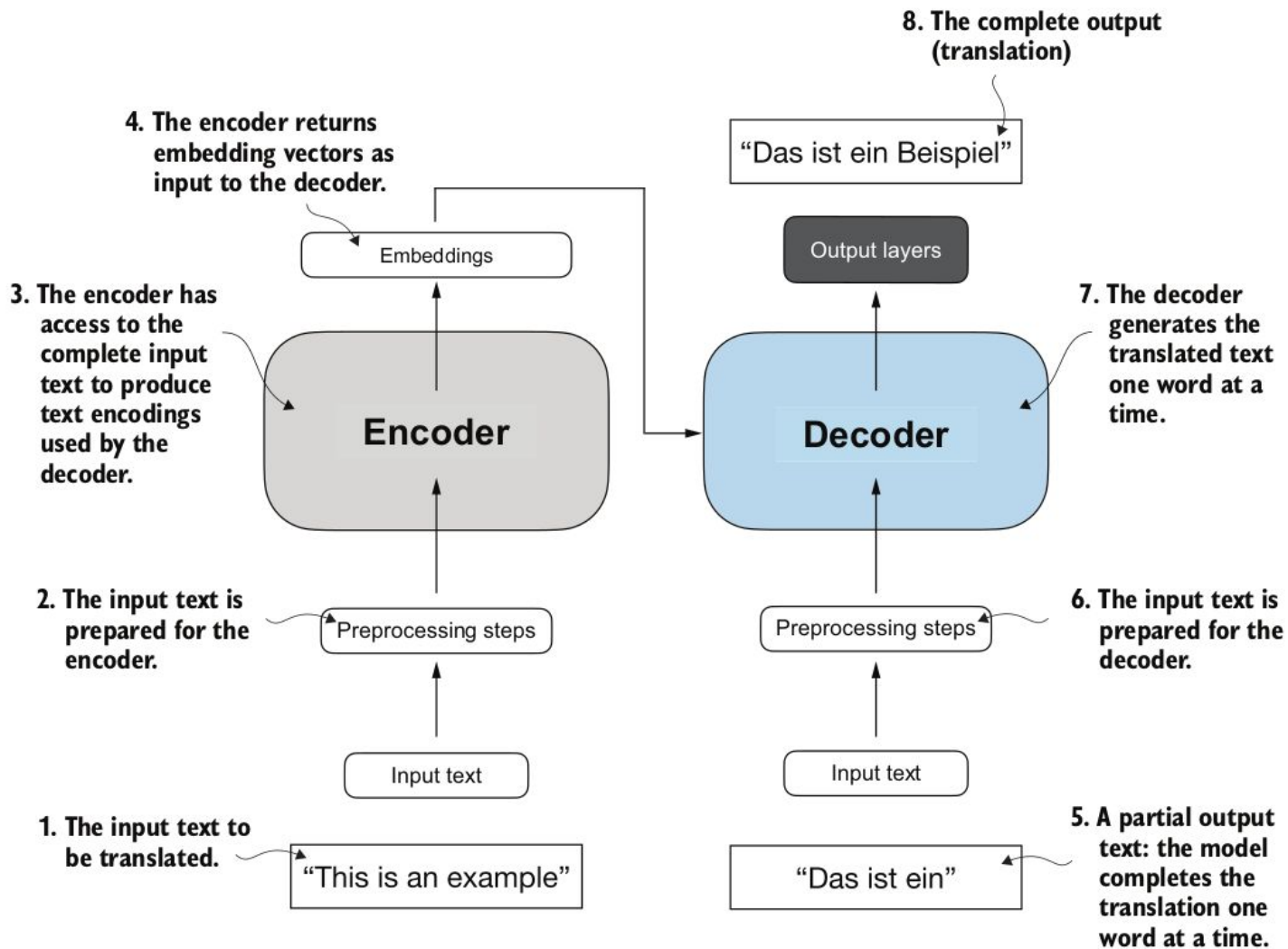
Output: a context vector $z(i)$ for each token i

Intuition: $z(i)$ gathers *contextual* information

ATTENTION HEADS AS KEY COMPONENTS IN A TRANSFORMER



ENCODER / DECODER



DECODERS USE CAUSAL ATTENTION

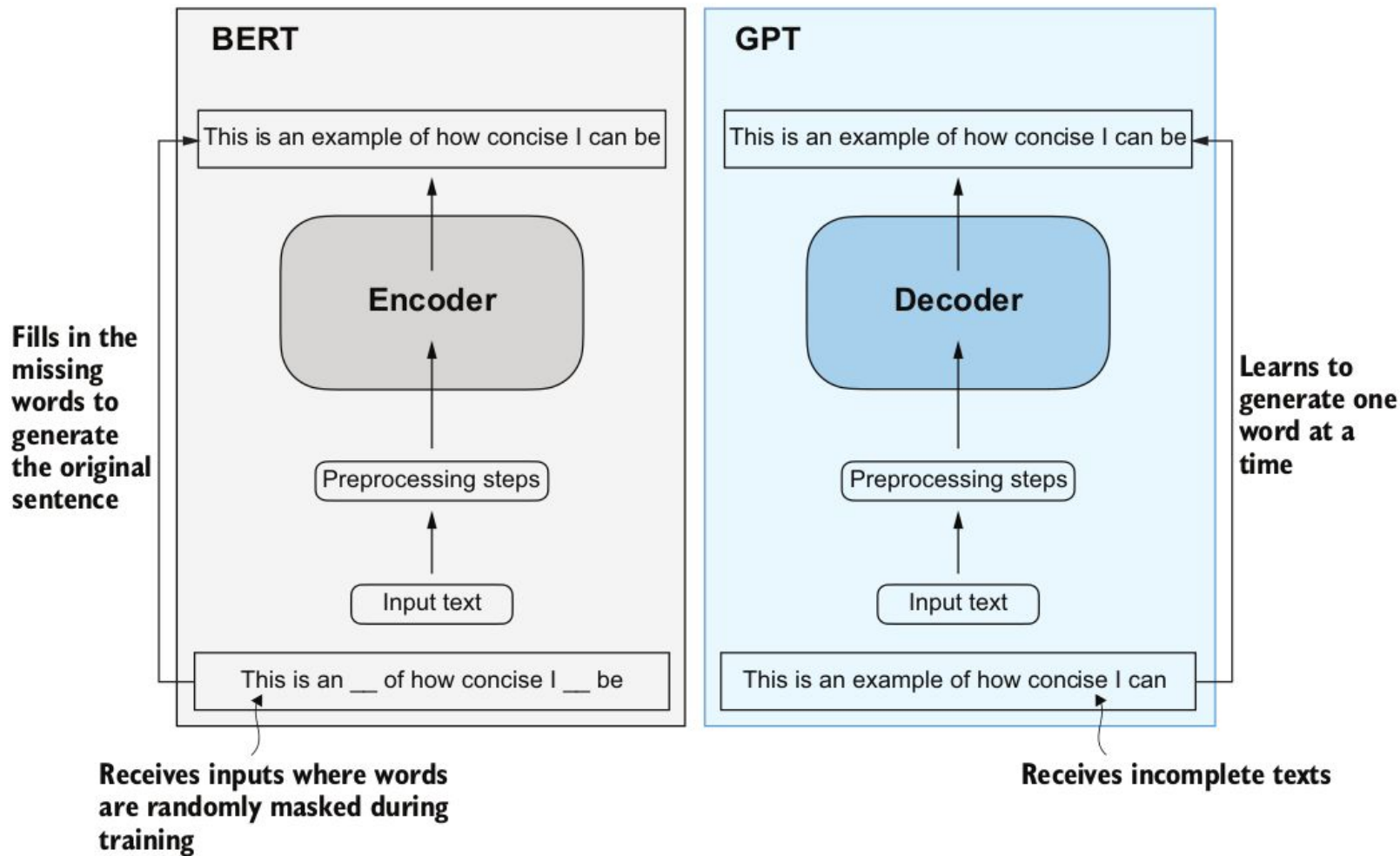
	Your	journey	starts	with	one	step
Your	0.19	0.16	0.16	0.15	0.17	0.15
journey	0.20	0.16	0.16	0.14	0.16	0.14
starts	0.20	0.16	0.16	0.14	0.16	0.14
with	0.18	0.16	0.16	0.15	0.16	0.15
one	0.18	0.16	0.16	0.15	0.16	0.15
step	0.19	0.16	0.16	0.15	0.16	0.15

Attention weight for input tokens
corresponding to “step” and “Your”



	Your	journey	starts	with	one	step
Your	1.0					
journey	0.55	0.44				
starts	0.38	0.30	0.31			
with	0.27	0.24	0.24	0.23		
one	0.21	0.19	0.19	0.18	0.19	
step	0.19	0.16	0.16	0.15	0.16	0.15

Masked out
future tokens
for the “Your”
token

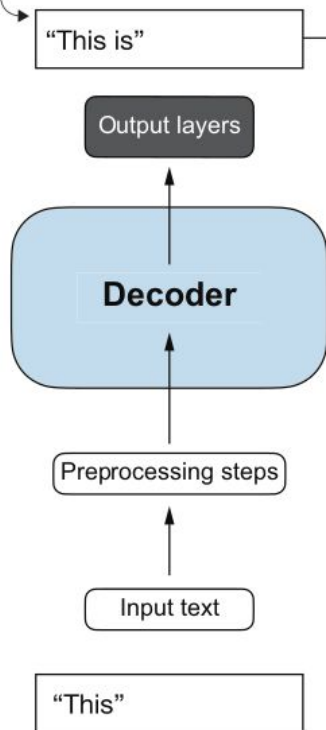


WHAT DOES "AUTO-REGRESSIVE" MEAN?

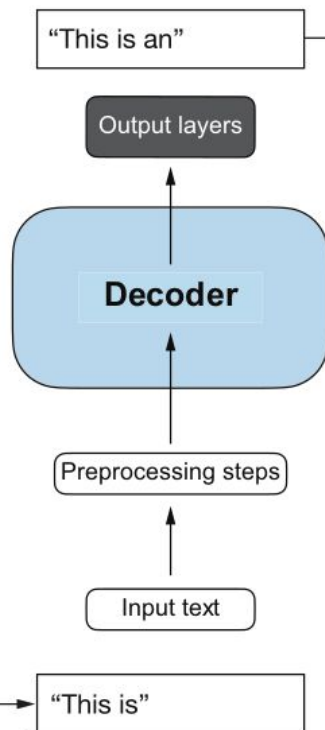
It means that for generating a single new token we feed the model with the input + all tokens generated so far.

Creates the next word based on the input text

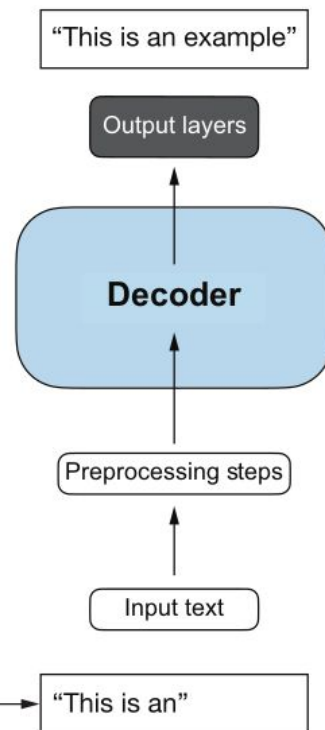
Iteration 1



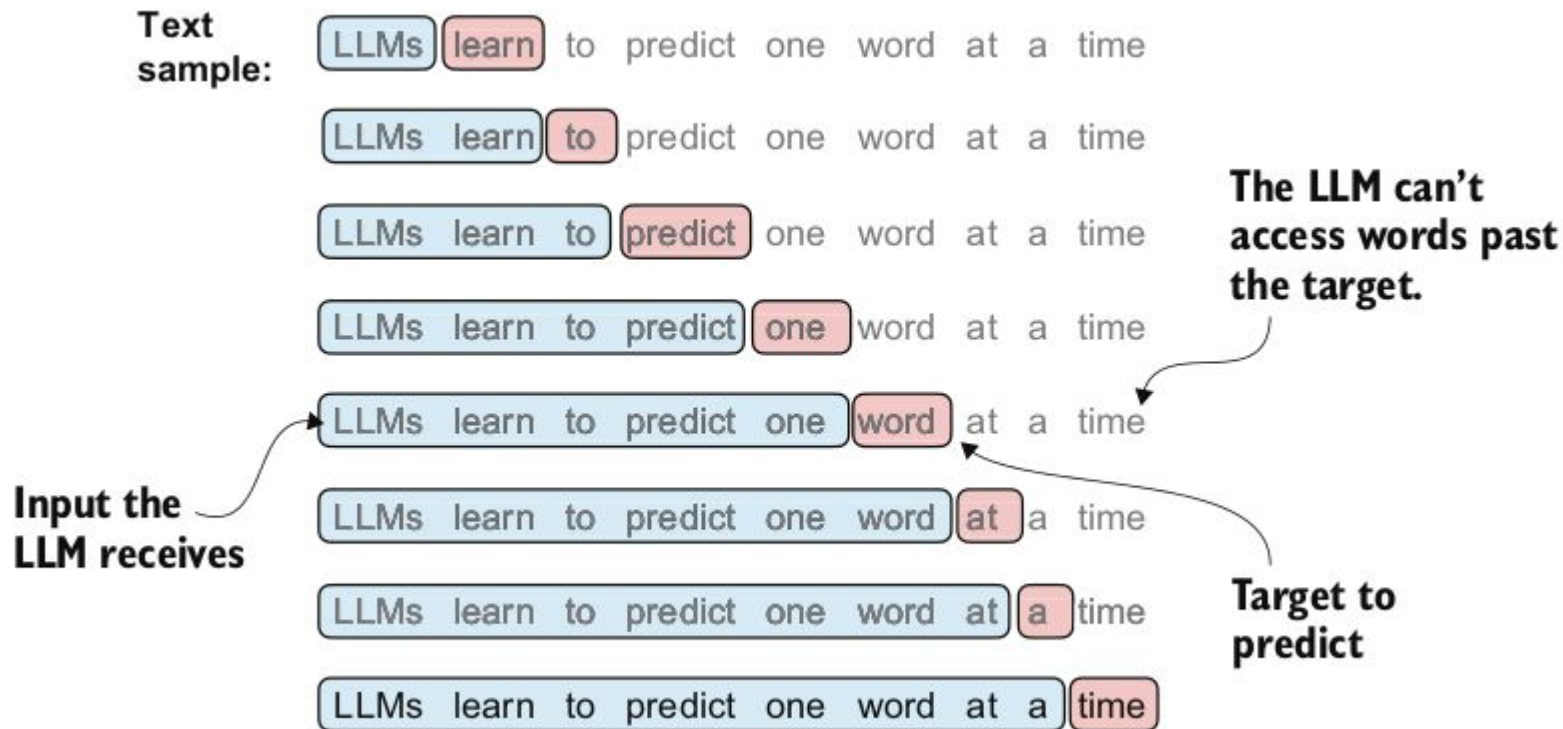
Iteration 2



Iteration 3



SLIDING WINDOWS



MODELS' SIGNATURES

Input: `x` of shape `(context_length)`, `y` of shape `(context_length)`

Output: `model(x,y) = (logits, loss)` where

- `logits` has shape `(context_length, vocab_size)`
- `loss` has shape `(context_length)`

For each window, make the prediction and compute the loss

MODELS' SIGNATURES WITH BATCHING

Input: X of shape (batch_size, context_length), Y of shape (batch_size, context_length)

Output: model(X,Y) = (logits, loss) where

- logits has shape (batch_size, context_length, vocab_size)
- loss has shape (batch_size, context_length)

PRE-TRAINING

TRAINING A LANGUAGE MODEL

Input: a sentence (as a sequence of tokens)

Output: predict the next token

Training loop:

- Feed the model a piece of text, using the next token as target
- Apply gradient descent to increase the probability to generate the target

WHAT IS CROSS ENTROPY LOSS?

Cross entropy measures the difference between probability distributions: it quantifies the dissimilarity between the predicted probability distribution and the true probability distribution.

In language modelling we do not have the true distribution of words, it is approximated from a training set:

$$H(T, q) = - \sum_{i=1}^N \frac{1}{N} \log_2 q(x_i)$$

Where N is the number of tokens in the training set and $q(x_i)$ is the probability that the model outputs x_i .

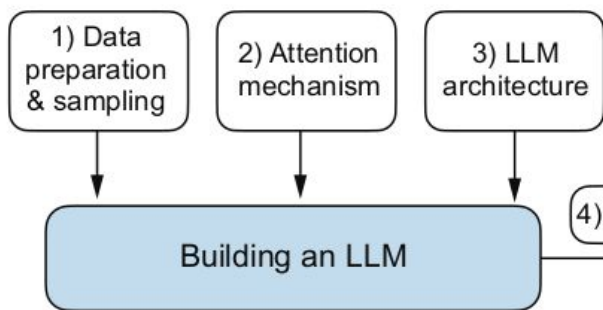
FINE-TUNING

FOUNDATION MODELS

Language Models are not very useful, they randomly generate texts... But this means that they somehow capture some information from natural language! They are also called *foundation models*.

Fine-tuning is about making Language Models solve concrete tasks, like classification, question answering, named-entity recognition...

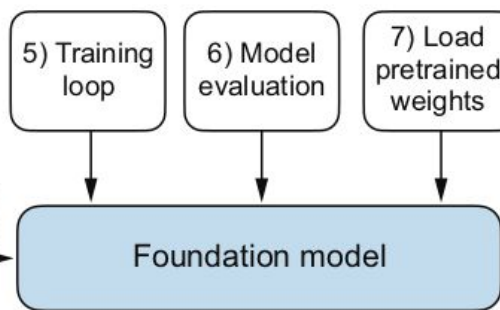
STAGE 1



Implements the data sampling and understand the basic mechanism

4) Pretraining

STAGE 2

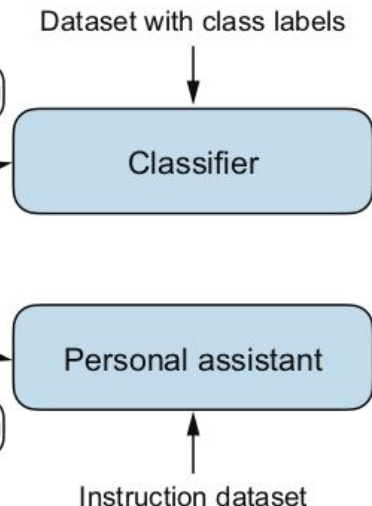


Pretrains the LLM on unlabeled data to obtain a foundation model for further fine-tuning

Fine-tunes the pretrained LLM to create a classification model

8) Fine-tuning

STAGE 3



9) Fine-tuning

Fine-tunes the pretrained LLM to create a personal assistant or chat model

TRAINING IS EXPENSIVE

We often cannot afford updating the *whole* model!

Most of us will not train foundation models... Rather
fine-tune existing ones.

LOW-RANK ADAPTATION (LORA)

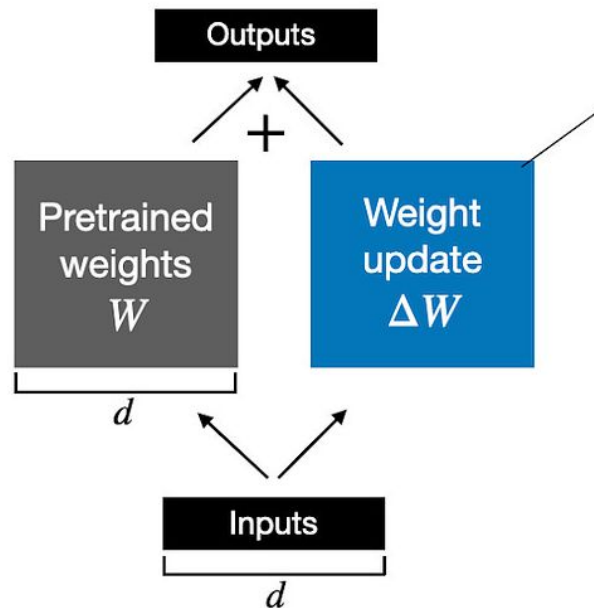
Two key ideas:

- (1) We only store the changes, not a new model
- (2) We only update a small number of parameters

IDEA: STORING WEIGHT UPDATES

Say we consider a linear layer with matrix W . We keep the matrix W fixed and store ΔW

Weight update in **regular finetuning**

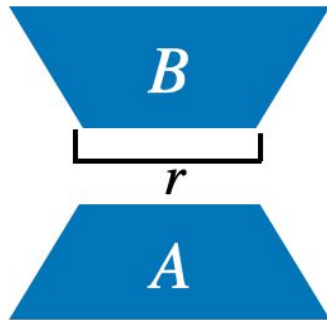


RANK APPROXIMATIONS

A matrix W of dimension $d \times d$ contains $d \times d$ parameters. It can be **rank- r approximated** by two matrices $A \times B$ with:

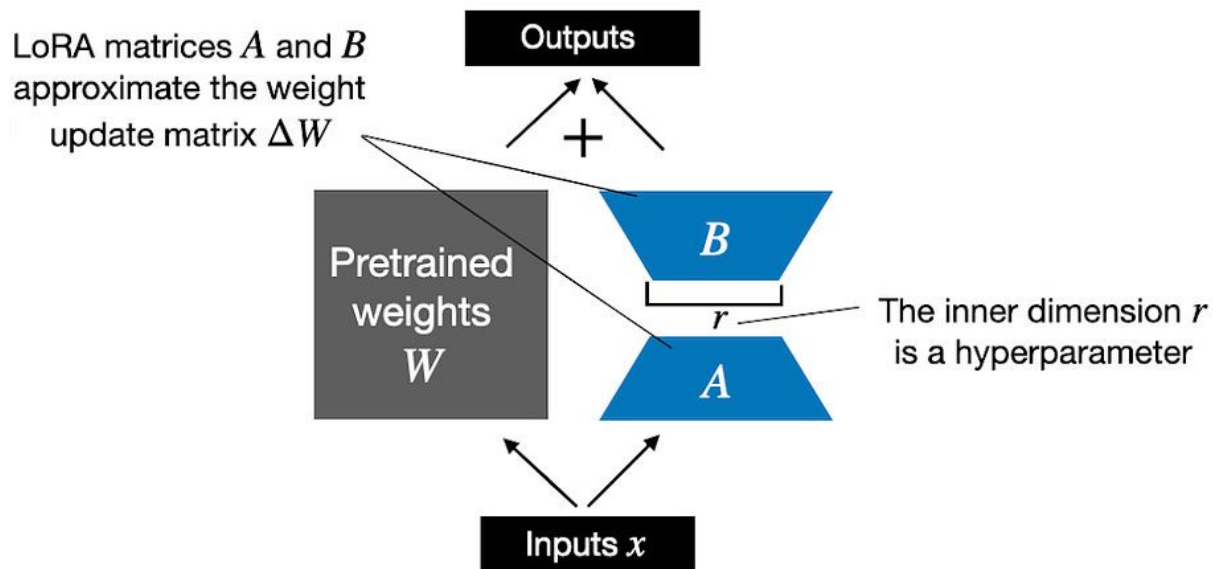
- A of dimension $d \times r$
- B of dimension $r \times d$

Instead of $d \times d$ parameters we now have $2 \times d \times r$ parameters.



WEIGHT UPDATE

Weight update in LoRA



```
def replace_linear_with_lora(model, rank, alpha):
    for name, module in model.named_children():
        if isinstance(module, torch.nn.Linear):
            # Replace the Linear layer with LinearWithLoRA
            setattr(model, name, LinearWithLoRA(module, rank, alpha))
        else:
            # Recursively apply the same function to child modules
            replace_linear_with_lora(module, rank, alpha)
```

- We then freeze the original model parameter and use the `replace_linear_with_lora` to replace the said `Linear` layers using the code below
- This will replace the `Linear` layers in the LLM with `LinearWithLoRA` layers

```
total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print(f"Total trainable parameters before: {total_params:,}")

for param in model.parameters():
    param.requires_grad = False

total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print(f"Total trainable parameters after: {total_params:,}")
```

Total trainable parameters before: 124,441,346
 Total trainable parameters after: 0

```
replace_linear_with_lora(model, rank=16, alpha=16)

total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print(f"Total trainable LoRA parameters: {total_params:,}")
```

Total trainable LoRA parameters: 2,666,528

PART 2:

HOW TO DEPLOY A LARGE LANGUAGE MODEL?

- Models
- Tokenizers
- Datasets
- Fine-tuning
- Inference

MODELS

THREE OPTIONS

- (1) **Inference-as-a-service:** through an API
- (2) **On the cloud:** as managed service or custom deployment
- (3) **Locally:** possible for small enough models (Ollama)

Hugging Face enables all three options!

[HTTPS://HUGGINGFACE.CO/](https://huggingface.co/)



It is a primarily a GitHub for models, but also develops a lot of useful packages and resources!

ARCHITECTURES

- **CPUs:** While generally slower for LLMs, they are more accessible and cost-effective for smaller models or less demanding tasks.
- **GPUs:** The most common choice for LLMs, offering significant performance improvements due to their parallel processing capabilities.
- **TPUs:** Google's specialized hardware designed for machine learning, providing even faster performance than GPUs for certain models.
- **Distributed Systems:** Multiple processors (CPUs, GPUs, or TPUs) working together to handle large models or high inference demands. Use `accelerate`: <https://huggingface.co/docs/accelerate/index>
- **Edge Devices:** Smaller, less powerful devices like smartphones and IoT devices can run optimized LLMs for specific tasks.

HOW TO GET GPU RESOURCES FOR FREE?

Anyone: Google colab, Kaggle, Codesphere, Sagemaker...

Locally: LaBRI has some (small-ish) GPU servers

Academics:

- grid5000 <https://www.grid5000.fr/w/Grid5000:Home>
- Jean-Zay <https://www.edari.fr/>

HOW MUCH MEMORY FOR A 3B MODEL?

The memory required to hold a 3B parameter LLM in memory depends heavily on the **data type** used to store the model weights:

Full Precision (FP32):

- Each parameter requires 32 bits (4 bytes)
- Total memory: 3 billion parameters * 4 bytes/parameter = 12 GB

Heavily Quantized (INT4):

- Each parameter requires 4 bits (0.5 bytes)
- Total memory: 3 billion parameters * 0.5 bytes/parameter = 1.5 GB

For comfortable performance and headroom, **8 GB or more** is recommended.

PIPELINE AND TASKS

TASKS

1. Text Classification

- Sentiment analysis: Determining the emotional tone of a text (positive, negative, neutral).
- Topic classification: Categorizing text into predefined topics.
- Spam detection: Identifying unsolicited or unwanted messages.
- Natural language inference: Determining the relationship between two sentences (entailment, contradiction, neutral).

2. Token Classification

- Named entity recognition (NER): Identifying and classifying named entities in text (people, organizations, locations, etc.).
- Part-of-speech (POS) tagging: Assigning grammatical tags to words (noun, verb, adjective, etc.).

3. Question Answering

- Extractive question answering: Finding the answer to a question within a given text.
- Multiple choice question answering: Selecting the best answer from a set of options.

MORE TASKS

4. Text Generation

- Text summarization: Generating a concise summary of a longer text.
- Translation: Translating text from one language to another.
- Dialogue generation: Creating conversational responses in a chatbot.
- Code generation: Generating code in various programming languages.

5. Text2Text Generation

- Paraphrasing: Rewriting a text while preserving its meaning.
- Summarization: Generating a concise summary of a longer text.
- Translation: Translating text from one language to another.

6. Fill-Mask

- Masked language modeling: Predicting missing words in a text.

7. Feature Extraction

- Generating embeddings: Creating numerical representations of text for use in other machine learning tasks.

PIPELINE: CONCISE BUT LITTLE CONTROL

```
from transformers import pipeline

# Create a sentiment analysis pipeline
classifier = pipeline(task="sentiment-analysis",
                      model="distilbert/distilbert-base-uncased-finetuned-sst-2-english")

# Run inference
result = classifier("This course is f***ing great!")

# Print the result
print(result)

[{'label': 'POSITIVE', 'score': 0.9998470544815063}]
```

YOU CAN CHOOSE THE "DEVICE" (CPU, GPU)

```
from transformers import pipeline

# Create a sentiment analysis pipeline
classifier = pipeline(task="sentiment-analysis",
                      model="distilbert/distilbert-base-uncased-finetuned-sst-2-english",
                      device="cuda")

# Run inference
result = classifier("This course is f***ing great!")

# Print the result
print(result)
```


The same model can be used for different tasks!

!! IMPORTANT !! Set `torch_dtype="auto"` to load the weights in the data type defined in a model's `config.json` file to automatically load the most memory-optimal data type.

[illegible][illegible]

WTF?

```
from transformers import AutoModelForTokenClassification

model = AutoModelForTokenClassification.from_pretrained("distilbert/distilbert-base-uncased",
                                                         num_labels = 5,
                                                         torch_dtype="auto")
```

config.json: 100%  483/483 [00:00<00:00, 9.73kB/s]

model.safetensors: 100%  268M/268M [00:05<00:00, 51.5MB/s]

Some weights of DistilBertForTokenClassification were not initialized from the model checkpoint at distilbert/distilbert-base-uncased and are newly initialized: ['classifier.bias', 'classifier.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Hugging Face: “Don’t worry, this is completely normal! The pretrained head of the BERT model is discarded, and replaced with a randomly initialized classification head. You will fine-tune this new model head on your sequence classification task, transferring the knowledge of the pretrained model to it.”

SERVERLESS INFERENCE API: SLOW BUT FREE

```
from huggingface_hub import InferenceClient

client = InferenceClient(
    "cardiffnlp/twitter-roberta-base-sentiment-latest",
    token="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",
)

client.text_classification("Today is a great day")

[TextClassificationOutputElement(label='positive', score=0.9836677312850952),
 TextClassificationOutputElement(label='neutral', score=0.01135887298732996),
 TextClassificationOutputElement(label='negative', score=0.004973393864929676)]
```

PIPELINE IS GOOD TO GET STARTED

But soon you feel limited.

Let's see how to get a bit more control!

TOKENIZERS

ONLY DEALING WITH TEXT HERE...

The tokenizer is for dealing with texts. For other formats:

- Speech and audio, use a Feature extractor to extract sequential features from audio waveforms and convert them into tensors.
- Image inputs use a ImageProcessor to convert images into tensors.
- Multimodal inputs, use a Processor to combine a tokenizer and a feature extractor or image processor.

```
from transformers import AutoTokenizer
```

```
tokenizer = AutoTokenizer.from_pretrained("google-bert/bert-base-uncased")  
tokenizer
```

```
BertTokenizerFast(name_or_path='google-bert/bert-base-uncased', vocab_size=30522, model_max_length=512, is_fast=True, padding_side='right', truncation_side='right', special_tokens={'unk_token': '[UNK]', 'sep_token': '[SEP]', 'pad_token': '[PAD]', 'cls_token': '[CLS]', 'mask_token': '[MASK]'}, clean_up_tokenization_spaces=True), added_tokens_decoder={  
    0: AddedToken("[PAD]", rstrip=False, lstrip=False, single_word=False, normalized=False, special=True),  
    100: AddedToken("[UNK]", rstrip=False, lstrip=False, single_word=False, normalized=False, special=True),  
    101: AddedToken("[CLS]", rstrip=False, lstrip=False, single_word=False, normalized=False, special=True),  
    102: AddedToken("[SEP]", rstrip=False, lstrip=False, single_word=False, normalized=False, special=True),  
    103: AddedToken("[MASK]", rstrip=False, lstrip=False, single_word=False, normalized=False, special=True),  
}
```

```
sequence = "In a hole in the ground there lived a hobbit."  
print(tokenizer(sequence))
```

```
{'input_ids': [101, 1999, 1037, 4920, 1999, 1996, 2598, 2045, 2973, 1037, 7570, 10322, 4183, 1012, 102], 'token_type_ids': [0, 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, 1]}
```

THE NASTY BUSINESS OF DATA PREPROCESSING

You want to create batches (a lot more efficient!). This means that you may need to:

- Pad: add a special token `[PAD]` to make sure all inputs have the same size
- Truncate: if some inputs are larger than the context length of your model, you need to break them up into more inputs (but it's not that obvious: better introduce some overlapping!)

Good news: Tokenizer does that for you!


```
batch_sentences = [  
    "But what about second breakfast?",  
    "Don't think he knows about second breakfast, Pip.",  
    "What about elevensies?",  
]  
encoded_input = tokenizer(batch_sentences, padding=True, truncation=True)  
print(encoded_input)
```

```
{'input_ids': [[101, 2021, 2054, 2055, 2117, 6350, 1029, 102, 0, 0, 0, 0, 0, 0], [101, 2123, 1005, 1056, 2228, 200  
2, 4282, 2055, 2117, 6350, 1010, 28315, 1012, 102], [101, 2054, 2055, 5408, 14625, 1029, 102, 0, 0, 0, 0, 0, 0,  
0]], 'token_type_ids': [[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],  
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]], 'attention_mask': [[1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0], [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]]}
```

DATASETS

```
from datasets import load_dataset
```

```
dataset = load_dataset("yelp_review_full")  
dataset["train"][100]
```

```
{'label': 0,  
  'text': 'My expectations for McDonalds are t rarely high. But for one to still fail so spectacularly...that takes something special!\n\nThe cashier took my friends\'s order, then promptly ignored me. I had to force myself in front of a cashier who opened his register to wait on the person BEHIND me. I waited over five minutes for a gigantic order that included precisely one kid\'s meal. After watching two people who ordered after me be handed their food, I asked where mine was. The manager started yelling at the cashiers for \\'serving off their orders\' when they didn\'t have their food. But neither cashier was anywhere near those controls, and the manager was the one serving food to customers and clearing the boards.\n\nThe manager was rude when giving me my order. She didn\'t make sure that I had everything ON MY RECEIPT, and never even had the decency to apologize that I felt I was getting poor service.\n\nI\'ve eaten at various McDonalds restaurants for over 30 years. I\'ve worked at more than one location. I expect bad days, bad moods, and the occasional mistake. But I have yet to have a decent experience at this store. It will remain a place I avoid unless someone in my party needs to avoid illness from low blood sugar. Perhaps I should go back to the racially biased service of Steak n Shake instead!'}]
```

FINE-TUNING

TRAINING ARGUMENTS

```
from transformers import TrainingArguments

training_args = TrainingArguments(
    output_dir="checkpoints",
    learning_rate=2e-5,
    eval_strategy="epoch",
    save_strategy="epoch",
    logging_strategy="epoch",
    per_device_train_batch_size=32,
    per_device_eval_batch_size=16,
    num_train_epochs=10,
    weight_decay=0.01,
    report_to="none",
)
```

TRAINER

```
from transformers import Trainer

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=small_train_dataset,
    eval_dataset=small_eval_dataset,
    compute_metrics=compute_metrics,
)
```

```
trainer.train()
```

LORA

LOAD MODELS, TWO OPTIONS

Option 1: load a PEFT adapter

```
from transformers import AutoModelForCausalLM, AutoTokenizer  
  
peft_model_id = "ybelkada/opt-350m-lora"  
model = AutoModelForCausalLM.from_pretrained(peft_model_id)
```

Option 2: Load the model and its adapter

```
from transformers import AutoModelForCausalLM, AutoTokenizer  
  
model_id = "facebook/opt-350m"  
peft_model_id = "ybelkada/opt-350m-lora"  
  
model = AutoModelForCausalLM.from_pretrained(model_id)  
model.load_adapter(peft_model_id)
```



```
from transformers import AutoModelForSeq2SeqLM
```

```
model = AutoModelForSeq2SeqLM.from_pretrained("bigscience/mt0-small")
```

```
from peft import LoraConfig
```

```
peft_config = LoraConfig(  
    lora_alpha=16,  
    lora_dropout=0.1,  
    r=64,  
    bias="none",  
    task_type="CAUSAL_LM",  
)
```

```
from peft import get_peft_model
```

```
model = get_peft_model(model, peft_config)  
model.print_trainable_parameters()
```

```
trainable params: 2,752,512 || all params: 302,929,280 || trainable%: 0.9086
```

The rest is as before: TrainingArguments and Trainer

PART 3:

ADVANCED TOPICS

- Retrieval-augmented Generation (RAG)
- Agents

RETRIEVAL-AUGMENTED GENERATION (RAG)

AGENTS