

# Natural Language Processing with Transformers and Large Language Models

2025-26

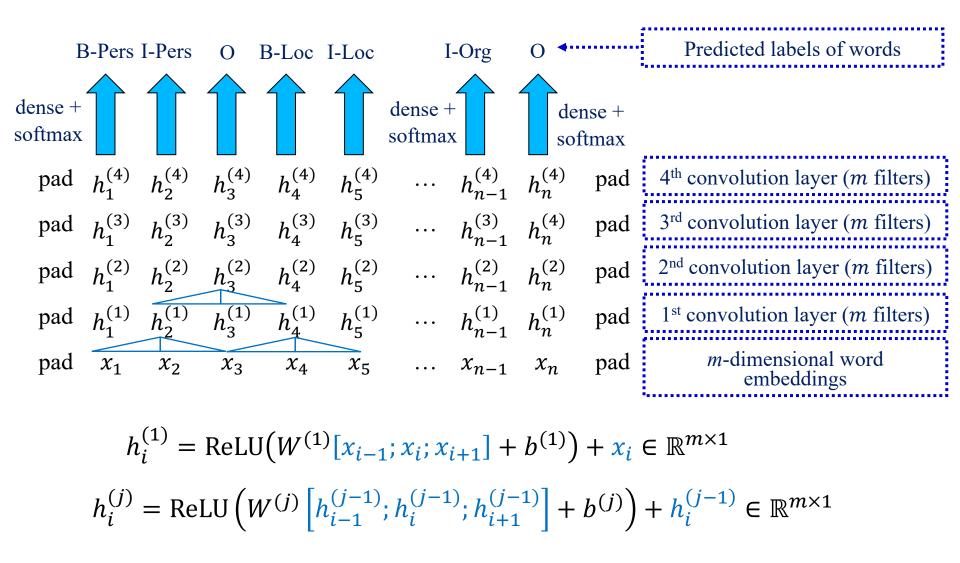
Ion Androutsopoulos

http://www.aueb.gr/users/ion/

#### Contents

- Transformer encoders, BERT.
- Encoder-decoder Transformers, BART, T5.
- Decoder-only Transformers, GPT-x.
- Prompting, supervised fine-tuning, RLHF, DPO.
- Parameter efficient training, LoRA.
- Retrieval augmented generation (RAG), LLMs with tools, agent, ReACT.
- Adding vision to LLMs, LLaVA, InstructBLIP.
- Optional study: Data augmentation for NLP.

#### Reminder: CNNs for token classification



#### Transformers for token classification

$$h_{1}^{(4)} \quad h_{2}^{(4)} \quad h_{3}^{(4)} \quad h_{4}^{(4)} \quad h_{5}^{(4)} \quad \cdots \quad h_{n-1}^{(4)} \quad h_{n}^{(4)}$$

$$h_{1}^{(3)} \quad h_{2}^{(3)} \quad h_{3}^{(3)} \quad h_{4}^{(3)} \quad h_{5}^{(3)} \quad \cdots \quad h_{n-1}^{(3)} \quad h_{n}^{(4)}$$

$$h_{1}^{(2)} \quad h_{2}^{(2)} \quad h_{3}^{(2)} \quad h_{4}^{(2)} \quad h_{5}^{(2)} \quad \cdots \quad h_{n-1}^{(2)} \quad h_{n}^{(2)}$$

$$h_{1}^{(1)} \quad h_{2}^{(1)} \quad h_{3}^{(1)} \quad h_{4}^{(1)} \quad h_{5}^{(1)} \quad \cdots \quad h_{n-1}^{(1)} \quad h_{n}^{(1)}$$

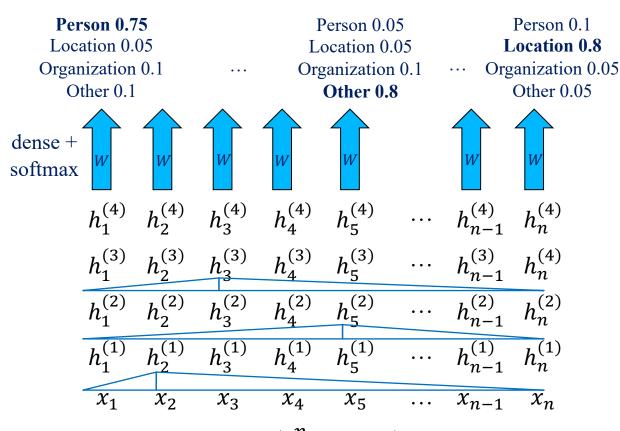
$$x_{1} \quad x_{2} \quad x_{3} \quad x_{4} \quad x_{5} \quad \cdots \quad x_{n-1} \quad x_{n}$$

Initial *m*-dimensional word embeddings

$$h_i^{(1)} = MLP^{(1)} \left( \sum_{r=1}^n a_{i,r}^{(1)} x_r \right) \in \mathbb{R}^m$$

To produce the revised embedding for the *i*-th word of a text, we sum all the original embeddings of the words of the text, but weighted by attention scores.

#### Transformers for token classification



Predicted labels of words

Compare to the correct predictions and adjust the weights of the entire neural net, including the bottom word (token) embeddings, which are randomly initialized.

Initial *m*-dimensional word embeddings

$$h_i^{(1)} = \text{MLP}^{(1)} \left( \sum_{r=1}^n a_{i,r}^{(1)} x_r \right) \in \mathbb{R}^m$$

$$h_i^{(j)} = \text{MLP}^{(j)} \left( \sum_{r=1}^n a_{i,r}^{(j)} h_r^{(j-1)} \right) \in \mathbb{R}^m$$

To produce the revised embedding for the *i*-th word of a text, we sum all the original embeddings of the words of the text, but weighted by attention scores.

#### Transformers for text classification

$$h^{max} = \left\langle \max\left(h_{*,1}^{(4)}\right), \max\left(h_{*,2}^{(4)}\right), \dots, \max(h_{*,m}^{(4)}) \right\rangle_{*}^{T} \in \mathbb{R}^{m}$$

global max pooling (max of each dimension)

$$h_i^{(1)} = \text{MLP}^{(1)} \left( \sum_{r=1}^n a_{i,r}^{(1)} x_r \right) \in \mathbb{R}^m$$

$$h_i^{(j)} = \text{MLP}^{(j)} \left( \sum_{r=1}^n a_{i,r}^{(j)} h_r^{(j-1)} \right) \in \mathbb{R}^m$$

Vector representing the entire text. We pass it through a dense layer and softmax (or MLP) to obtain a probability per class.

Compare to the correct predictions and adjust the weights of the entire net.

Initial *m*-dimensional word embeddings

Without the MLP (or at least a dense layer), each dimension of  $h_i^{(j)}$  would only depend on the corresponding dimensions of the  $h_r^{(j-1)}$  vectors.

# Query-Key-Value self-attention

$$h_1^{(4)}$$
  $h_2^{(4)}$   $h_3^{(4)}$   $h_4^{(4)}$   $h_5^{(4)}$   $\cdots$   $h_{n-1}^{(4)}$   $h_n^{(4)}$   $h_1^{(4)}$   $h_2^{(3)}$   $h_3^{(3)}$   $h_4^{(3)}$   $h_5^{(3)}$   $\cdots$   $h_{n-1}^{(3)}$   $h_n^{(4)}$   $h_1^{(2)}$   $h_2^{(2)}$   $h_3^{(2)}$   $h_4^{(2)}$   $h_5^{(2)}$   $\cdots$   $h_{n-1}^{(2)}$   $h_n^{(2)}$   $h_1^{(2)}$   $h_1^{(2)}$   $h_1^{(1)}$   $h_2^{(1)}$   $h_3^{(1)}$   $h_4^{(1)}$   $h_5^{(1)}$   $h_5^{(1)}$   $\cdots$   $h_{n-1}^{(1)}$   $h_n^{(1)}$   $h_n^{(1)$ 

$$h_i^{(1)} = \text{MLP}^{(1)} \left( \sum_{r=1}^n a_{i,r}^{(1)} v_r^{(1)} \right) =$$

$$= \text{MLP}^{(1)} \left( \sum_{r=1}^n \text{softmax} \left( q_i^{(1)T} k_r^{(1)} \right) v_r^{(1)} \right) \in \mathbb{R}^{m \times 1}$$

$$h_i^{(j)} = \text{MLP}^{(j)} \left( \sum_{r=1}^n a_{i,r}^{(j)} v_r^{(j)} \right) =$$

$$= \text{MLP}^{(j)} \left( \sum_{r=1}^n \text{softmax} \left( q_i^{(j)T} k_r^{(j)} \right) v_r^{(j)} \right) \in \mathbb{R}^{m \times 1}$$

4<sup>th</sup> attention layer

3<sup>rd</sup> attention layer

2<sup>nd</sup> attention layer

1st attention layer

*m*-dimensional word embeddings

$$q_i^{(1)} = W^{Q,(1)} x_i$$

$$k_r^{(1)} = W^{K,(1)} x_r$$

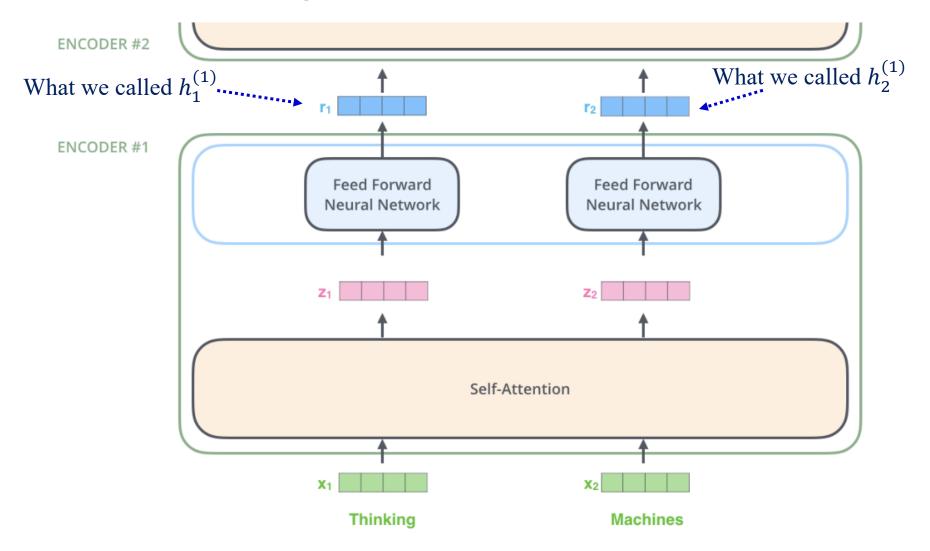
$$v_r^{(1)} = W^{V,(1)} x_r$$

$$q_i^{(j)} = W^{Q,(j)} h_i^{(j-1)}$$

$$k_r^{(j)} = W^{K,(j)} h_r^{(j-1)}$$

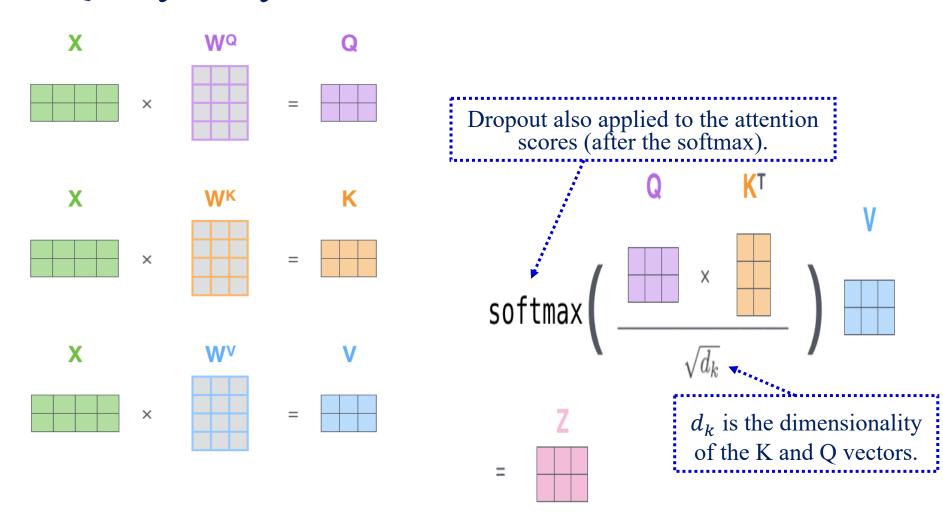
$$v_r^{(j)} = W^{V,(j)} h_r^{(j-1)}$$

## Stacking Transformer Encoders



Figures from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer/">https://jalammar.github.io/illustrated-transformer/</a>). Transformers paper: Vaswani et al., "Attention is All You Need", 2017 (<a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>).

## Query-Key-Value attention via matrices



Figures from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer/">https://jalammar.github.io/illustrated-transformer/</a>). Transformers paper: Vaswani et al., "Attention is All You Need", 2017 (<a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>).

## Multiple attention heads

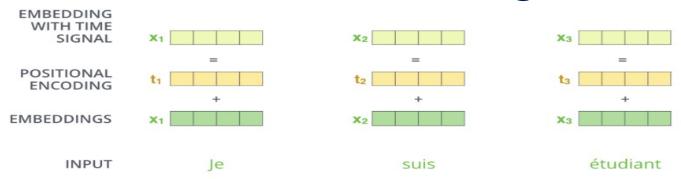
Because of the softmax, each attention head mostly considers only one token. So, let's use multiple attention heads.

2) We embed 3) Split into 8 heads. 4) Calculate attention 5) Concatenate the resulting Z matrices, 1) This is our each word\* We multiply X or using the resulting then multiply with weight matrix Wo to input sentence\* R with weight matrices Q/K/V matrices produce the output of the layer  $\mathsf{W}_0^\mathsf{Q}$ **Thinking** Machines \* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one  $W^0$  is useful even if the concatenated  $Z_0, ..., Z_7$  already have the right dimensions, to allow **combinations of features** 

Figures from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer/">https://jalammar.github.io/illustrated-transformer/</a>). Transformers paper: Vaswani et al., "Attention is All You Need", 2017 (<a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>).

from different attention heads.

## Positional encodings

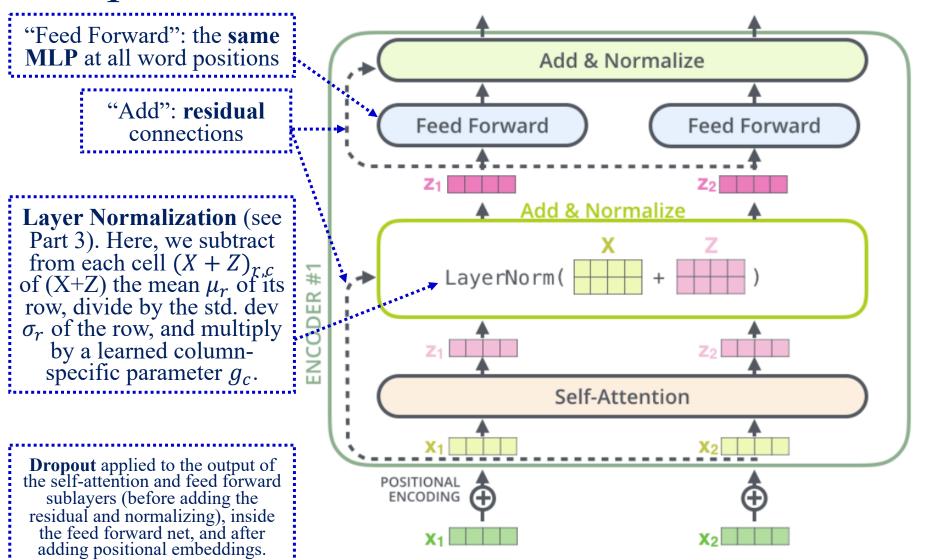


Positional encodings needed to capture the word order/positions.

- Without them, Transformers are unaware of word order.
- Sinusoid functions used to produce them in the original paper.
- But can also be position embeddings learned during training.
  - o Embedding of **position 1**, embedding of **position 2** etc.
- Relative position embeddings can also be used.
  - O They consider the **distance** from the **current** to the **attended position** in the **self-attention blocks**. (https://paperswithcode.com/method/relative-position-encodings).

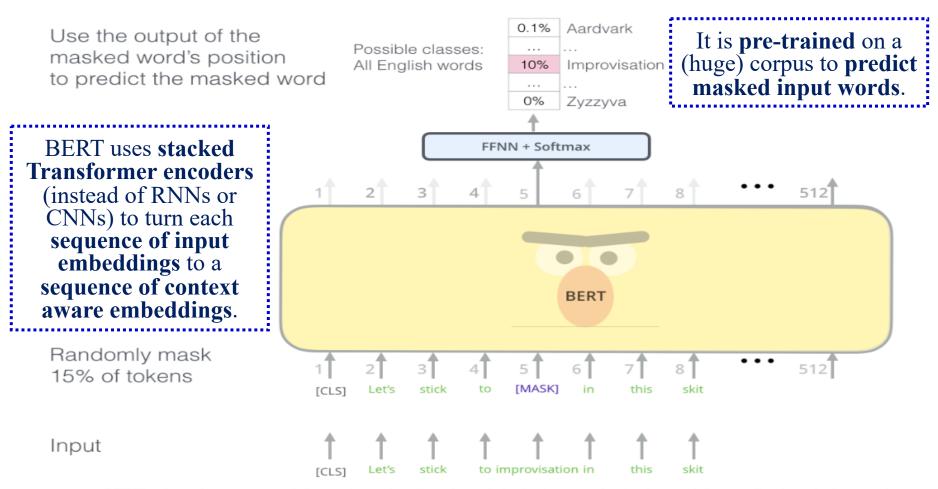
Figures from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer/">https://jalammar.github.io/illustrated-transformer/</a>). Transformers paper: Vaswani et al., "Attention is All You Need", 2017 (<a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>).

## Complete Transformer encoder block



Figures from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer/">https://jalammar.github.io/illustrated-transformer/</a>). Transformers paper: Vaswani et al., "Attention is All You Need", 2017 (<a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>).

## BERT – Pretraining to predict masked words



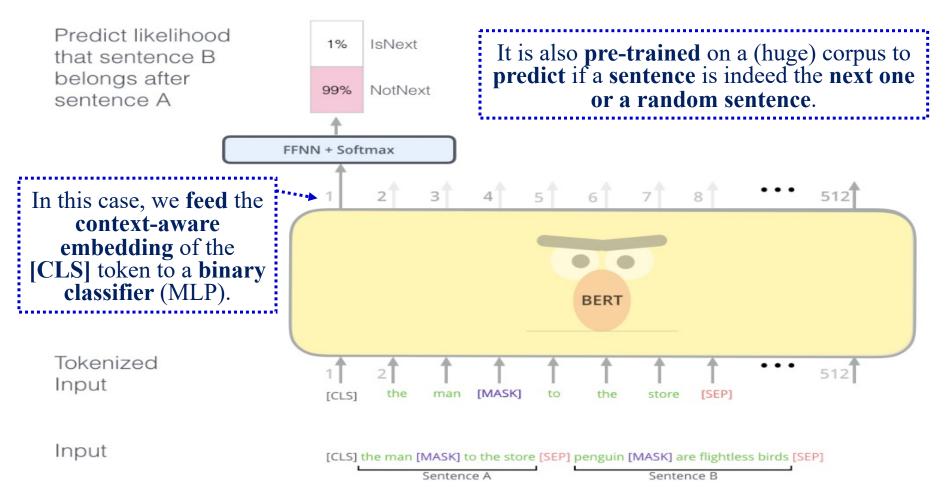
BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word.

Figures from J. Alammar's "The Illustrated BERT, ELMo, and co."

(<a href="http://jalammar.github.io/illustrated-bert/">http://jalammar.github.io/illustrated-bert/</a>). BERT paper: Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 2018

(<a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>).

#### BERT – Pretraining to predict the next sentence



The second task BERT is pre-trained on is a two-sentence classification task. The tokenization is oversimplified in this graphic as BERT actually uses WordPieces as tokens rather than words --- so some words are broken down into smaller chunks.

14

Figures from J. Alammar's "The Illustrated BERT, ELMo, and co."

(<a href="http://jalammar.github.io/illustrated-bert/">http://jalammar.github.io/illustrated-bert/</a>). BERT paper: Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 2018

(<a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>).

#### BERT – Fine-tuning for sentence classification

We feed the context-aware embedding of the [CLS] token of each sentence to a task-specific classifier (e.g., MLP) that classifies the sentence (e.g., Positive, Neutral, Negative etc.)

Starting from the pre-trained BERT, we jointly train BERT (further) and the task-specific classifier on (possibly few) task-specific training examples (e.g., tweets + opinion labels).

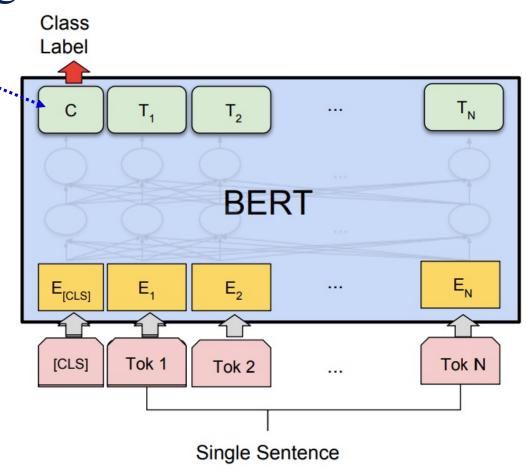


Figure from Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 2018 (https://arxiv.org/abs/1810.04805).

## BERT – Fine-tuning for token classification

We **feed** the **context-aware embeddings** of the sentence's words to a **classifier** (e.g., MLP) that classifies them as **B-Per**, **I-Per**, **B-Org**, **I-Org**, ..., **Other**.

Starting from the pre-trained BERT, we jointly train BERT (further) and the task-specific classifier on (possibly few) task-specific training examples (manually labeled sentences).

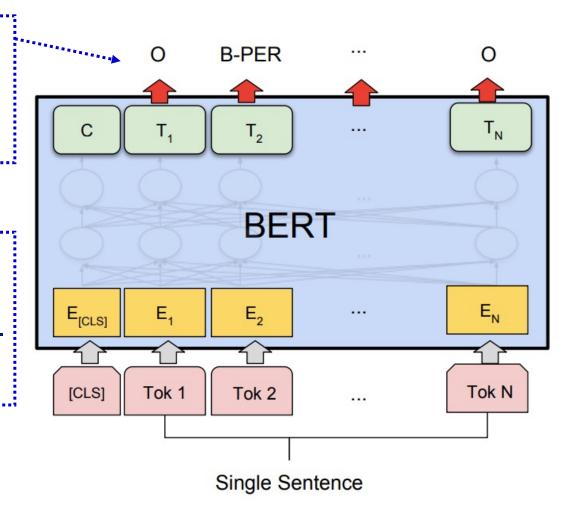


Figure from Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 2018 (<a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>).

## BERT – Fine-tuning for textual entailment

We feed the context-aware embedding of the [CLS] token of each sentence pair to a task-specific classifier (e.g., MLP) that classifies the pair as Entailment, Contradiction, Neutral. E.g., "Mary plays in the garden" entails "Mary is in the garden" but contradicts "Mary is asleep".

Starting from the pre-trained BERT, we jointly train BERT (further) and the task-specific classifier on (possibly few) task-specific training examples (annotated sentence pairs).

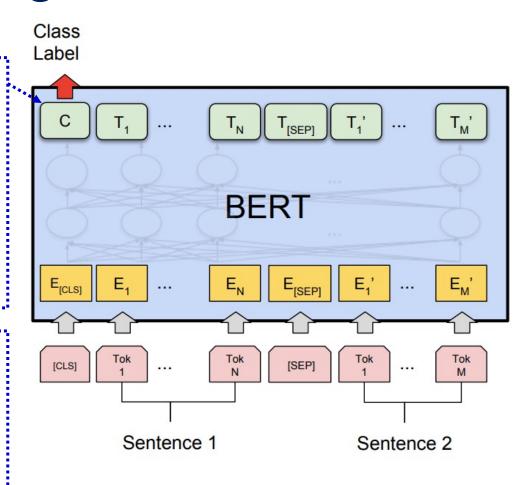


Figure from Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 2018 (<a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>).

## Machine Reading Comprehension (MRC)

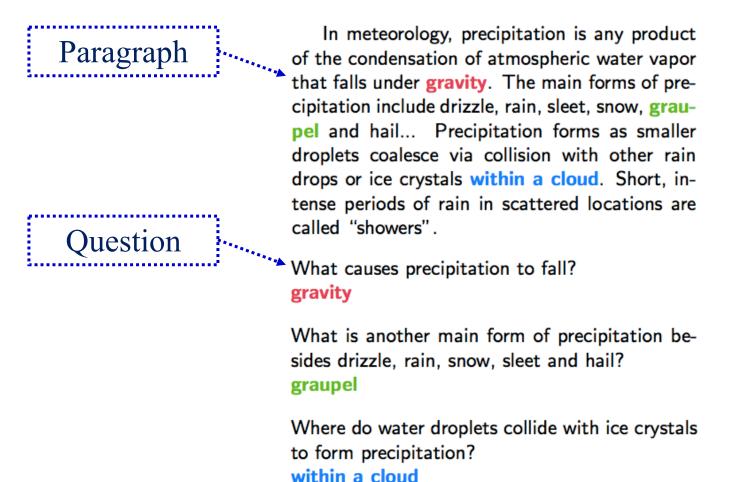


Figure from P. Rajpurkar et al., "SQuAD: 100,000+ Questions for Machine Comprehension of Text.", EMNLP 2016 (https://aclweb.org/anthology/D16-1264).

# BERT – Fine-tuning for MRC

We **feed** the **context-aware embeddings** of the paragraph's words to a **classifier** (e.g., MLP) that classifies them as **Start-Answer**, **End-Answer**, **Other**.

Starting from the pre-trained BERT, we jointly train BERT (further) and the task-specific classifier on (possibly few) task-specific training examples (paragraph-question pairs).

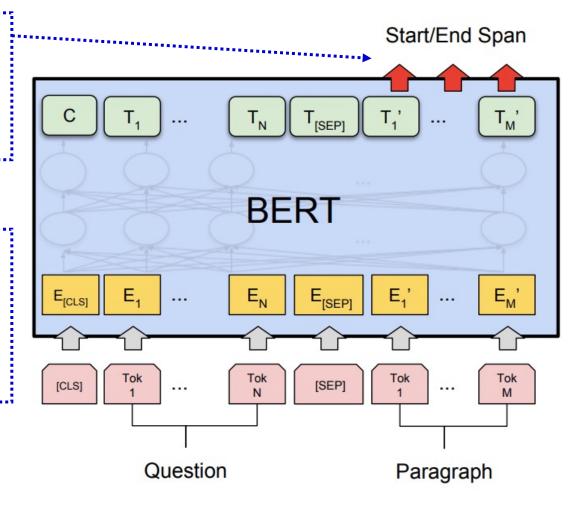


Figure from Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 2018 (https://arxiv.org/abs/1810.04805).

## SMITH (hierarchical BERT)

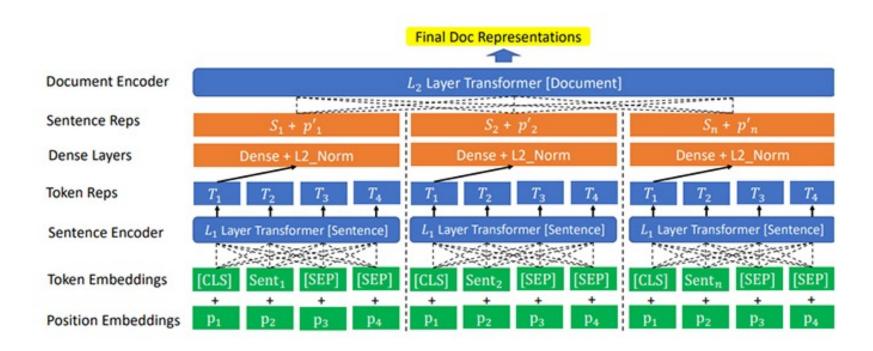
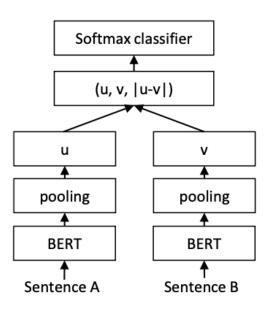


Figure from Yang et al., "Beyond 512 Tokens: Siamese Multi-depth Transformer-based Hierarchical Encoder for Long-Form Document Matching", CIKM 2020 (https://dl.acm.org/doi/10.1145/3340531.3411908).

BERT variants for long documents include, for example, also **Longformer** (<a href="https://arxiv.org/abs/2004.05150">https://arxiv.org/abs/2004.05150</a>) and **Big Bird** (<a href="https://arxiv.org/abs/2007.14062">https://arxiv.org/abs/2007.14062</a>), which are not hierarchical, but use **sparse attention** to avoid quadratic complexity (to the input length). See also, e.g., **FlashAttention** (<a href="https://arxiv.org/abs/2205.14135">https://arxiv.org/abs/2205.14135</a>).

#### Sentence Encoders: SBERT



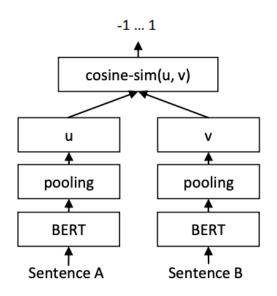


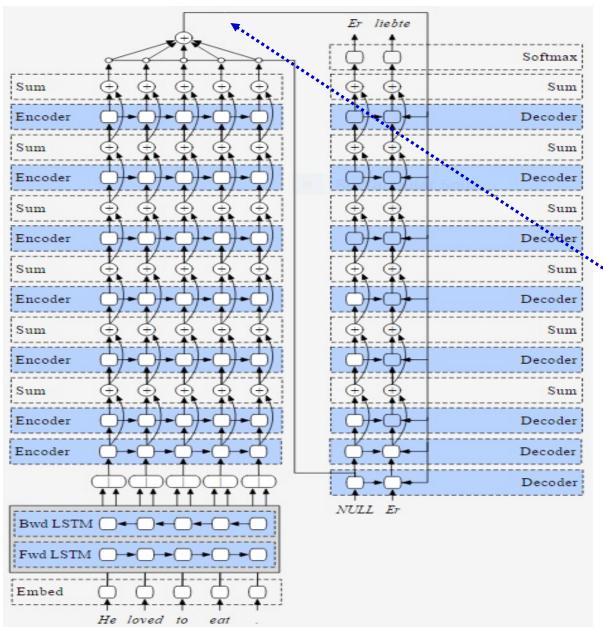
Figure 1: SBERT architecture with classification objective function, e.g., for fine-tuning on SNLI dataset. The two BERT networks have tied weights (siamese network structure).

Figure 2: SBERT architecture at inference, for example, to compute similarity scores. This architecture is also used with the regression objective function.

- Siamese BERTs (shared weights) fine-tuned on pairs of sentences (e.g., NL inference, sentence similarity).
- Then use the resulting BERT as **sentence encoder** (e.g., for dense retrieval in RAG, see below).

Figure from Reimers & Gurevych (2019), "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks" (https://aclanthology.org/D19-1410/).

## Reminder: RNN-based MT system



Google's paper: https://arxiv.org/abs/1609.08144

Images from Stephen Merity's http://smerity.com/articles/2016/google\_nmt\_arch.html

Attention over the states of the encoder.

#### Stacked Transformer encoders-decoders

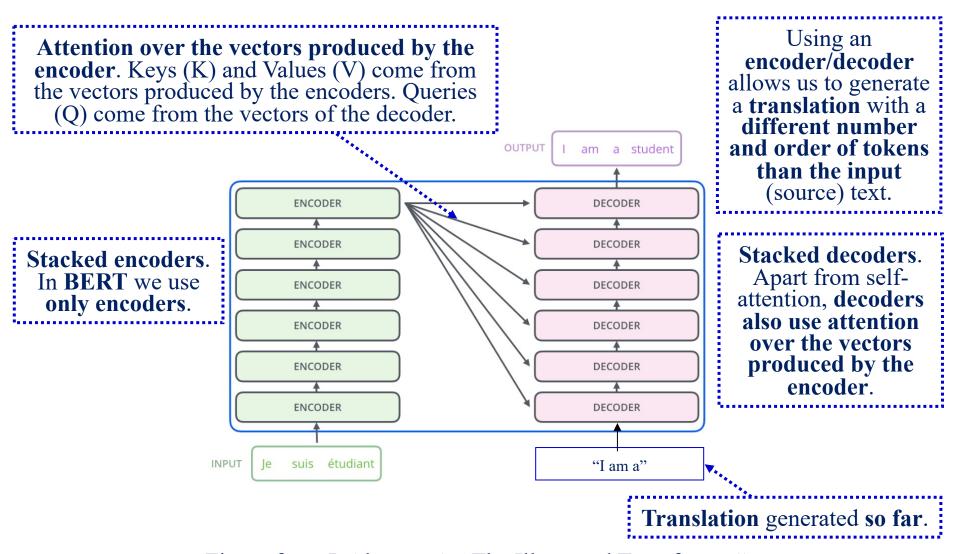


Figure from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer/">https://jalammar.github.io/illustrated-transformer/</a>). Transformers paper: Vaswani et al., "Attention is All You Need", 2017 (<a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>).

## OKV self-attention and cross-attention

stacked decoder layers

stacked encoder layers

$$d_1^{(3)}$$
  $d_2^{(3)}$   $d_3^{(3)}$   $d_4^{(3)}$   $d_5^{(3)}$ 

$$e_{1}^{(2)} \quad e_{2}^{(2)} \quad e_{3}^{(2)} \quad e_{4}^{(2)} \qquad d_{1}^{(2)} \quad d_{2}^{(2)} \quad d_{3}^{(2)} \quad d_{4}^{(2)} \quad d_{5}^{(2)}$$

$$\text{self-attention}$$

$$\stackrel{\text{embeddings}}{\text{of input}} \quad e_{1}^{(1)} \quad e_{2}^{(1)} \quad e_{3}^{(1)} \quad e_{4}^{(1)} \qquad d_{1}^{(1)} \quad d_{2}^{(1)} \quad d_{3}^{(1)} \quad d_{4}^{(1)} \quad d_{5}^{(1)}$$

$$\text{tokens}$$

$$e_{i}^{(2)} = \text{MLP}^{(e,2)} \left( \sum_{r=1}^{4} a_{i,r}^{(e,2)} v_{r}^{(e,2)} \right) = \qquad \qquad q_{i}^{(e,2)} = W^{Q,(e,2)} e_{i}^{(1)}$$

$$= \text{MLP}^{(e,2)} \left( \sum_{r=1}^{4} \text{softmax} \left( q_{i}^{(e,2)T} k_{r}^{(e,2)} \right) v_{r}^{(e,2)} \right) \in \mathbb{R}^{m \times 1} \qquad v_{r}^{(e,2)} = W^{V,(e,2)} e_{r}^{(1)}$$

$$d_{1}^{(2)} \quad d_{2}^{(2)} \quad d_{3}^{(2)} \quad d_{4}^{(2)} \quad d_{5}^{(2)}$$

$$d_{1}^{(1)} \quad d_{2}^{(1)} \quad d_{3}^{(1)} \quad d_{4}^{(1)} \quad d_{5}^{(1)}$$

$$q_{i}^{(e,2)} = W^{Q,(e,2)}e_{i}^{(1)}$$

$$k_{r}^{(e,2)} = W^{K,(e,2)}e_{r}^{(1)}$$

$$\in \mathbb{R}^{m \times 1} \qquad v_{r}^{(e,2)} = W^{V,(e,2)}e_{r}^{(1)}$$

## QKV self-attention and cross-attention

stacked decoder layers

 $d_1^{(3)}$   $d_2^{(3)}$   $d_3^{(3)}$   $d_4^{(3)}$   $d_5^{(3)}$ stacked encoder layers cross-attention  $e_1^{(2)}$   $e_2^{(2)}$   $e_3^{(2)}$   $e_4^{(2)}$  $d_1^{(2)}$   $d_2^{(2)}$   $d_3^{(2)}$   $d_4^{(2)}$   $d_5^{(2)}$ self-attention masked selfembeddings of previously  $d_1^{(1)}$   $d_2^{(1)}$   $d_3^{(1)}$   $d_4^{(1)}$   $d_5^{(1)}$  attention embeddings  $e_1^{(1)}$   $e_2^{(1)}$   $e_3^{(1)}$   $e_4^{(1)}$ tokens  $d_i^{(2)} = MLP^{(d,2)} \left( \sum_{i=1}^{l} a_{i,r}^{(d,2)} v_r^{(d,2)} \right) =$  $q_i^{(d,2)} = W^{Q,(d,2)} d_i^{(1)}$  $k_r^{(d,2)} = W^{K,(d,2)} d_r^{(1)}$  $= MLP^{(d,2)} \left( \sum_{r=1}^{l} \operatorname{softmax} \left( q_i^{(d,2)T} k_r^{(d,2)} \right) v_r^{(d,2)} \right) \in \mathbb{R}^{m \times 1}$  $v_r^{(d,2)} = W^{V,(d,2)} d_r^{(1)}$ 

$$d_{i}^{(3)} = MLP^{(d,3)} \left( \sum_{r=1}^{4} a_{i,r}^{(d,3)} v_{r}^{(d,3)} \right) = q_{i}^{(d,3)} = W^{Q,(d,3)} d_{i}^{(2)}$$

$$k_{r}^{(d,3)} = W^{K,(d,3)} e_{r}^{(2)}$$

 $= MLP^{(d,3)} \left( \sum_{i=1}^{4} softmax \left( q_i^{(d,3)T} k_r^{(d,3)} \right) v_r^{(d,3)} \right) \in \mathbb{R}^{m \times 1} \qquad v_r^{(d,3)} = W^{V,(d,3)} e_r^{(2)}$ 

#### Transformer-based Encoder-Decoder

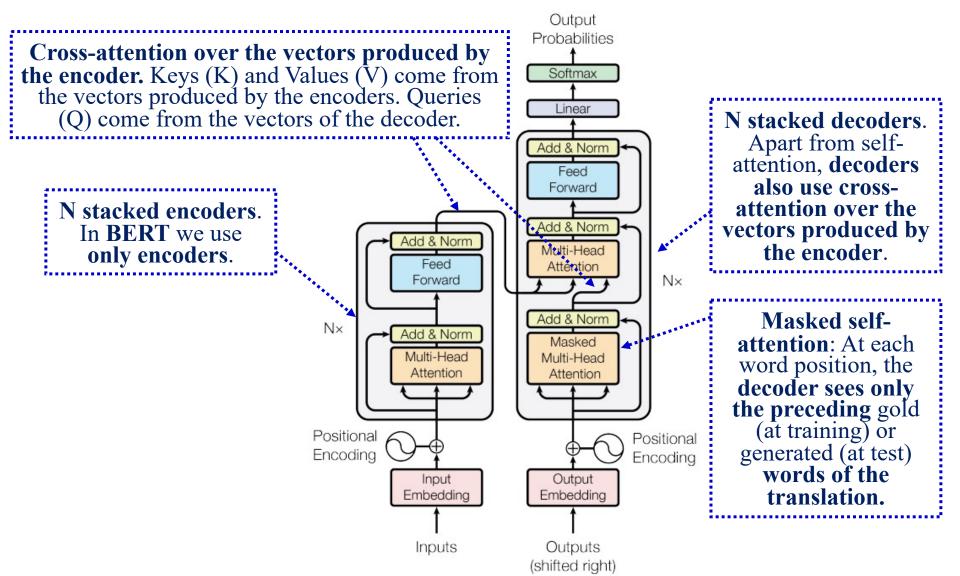
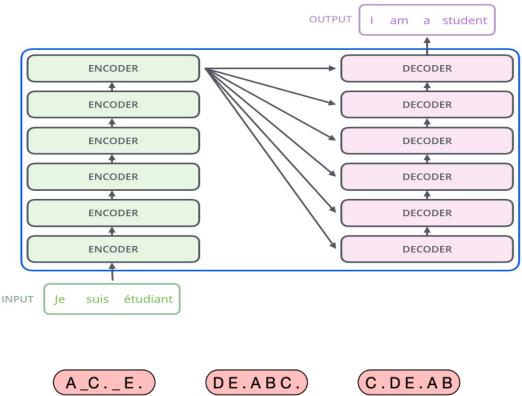


Figure from Vaswani et al., "Attention is All You Need", 2017. https://arxiv.org/abs/1706.03762

# BART – Using encoders & decoders

BART uses both stacked encoder and stacked decoder Transformer layers.

During pre-training,
BART is trained to
"translate" noised
text to the original
(without noise) text.



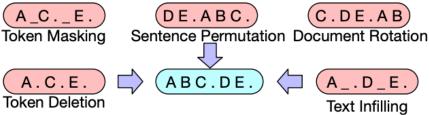


Figure 2: Transformations for noising the input that we experiment with. These transformations can be composed.

Top figure from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer/">https://jalammar.github.io/illustrated-transformer/</a>). Bottom figures from M. Lewis et al., "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension", ACL 2020 (<a href="https://www.aclweb.org/anthology/2020.acl-main.703/">https://www.aclweb.org/anthology/2020.acl-main.703/</a>).

# BART – Fine-tuning for summarization

#### Source Document (abbreviated)

#### BART Summary

The researchers examined three types of coral in reefs off the coast of Fiji ... The researchers found when fish were plentiful, they would eat algae and seaweed off the corals, which appeared to leave them more resistant to the bacterium Vibrio corallilyticus, a bacterium associated with bleaching. The researchers suggested the algae, like warming temperatures, might render the corals' chemical defenses less effective, and the fish were protecting the coral by removing the algae.

Fisheries off the coast of Fiji are protecting coral reefs from the effects of global warming, according to a study in the journal Science.

"hallucination" of journal source

This is the first time anyone has been recorded to run a full marathon of 42.195 kilometers (approximately 26 miles) under this pursued landmark time. It was not, however, an officially sanctioned world record, as it was not an "open race" of the IAAF. His time was 1 hour 59 minutes 40.2 seconds. Kipchoge ran in Vienna, Austria. It was an event specifically designed to help Kipchoge break the two hour barrier.

Kenyan runner Eliud Kipchoge has run a marathon in less than two hours.

PG&E stated it scheduled the blackouts in response to forecasts for high winds amid dry conditions. The aim is to reduce the risk of wildfires. Nearly 800 thousand customers were scheduled to be affected by the shutoffs which were expected to last through at least midday tomorrow.

Power has been turned off to millions of customers in California as part of a power shutoff plan.

Table 9: Example summaries from the XSum-tuned BART model on WikiNews articles. For clarity, only relevant excerpts of the source are shown. Summaries combine information from across the article and prior knowledge.

M. Lewis et al., "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension", ACL 2020 (<a href="https://www.aclweb.org/anthology/2020.acl-main.703/">https://www.aclweb.org/anthology/2020.acl-main.703/</a>).

T5 – Using encoders & decoders

T5 also uses both stacked encoder and stacked decoder Transformer layers.

ENCODER

ENCODER

DECODER

In pre-training, T5 is trained to recover missing/noised parts of the input, here masked spans.

Thank you for inviting me to your party last week.

Inputs
Thank you <X> me to your party <Y> week.

Targets

<X> for inviting <Y> last <Z>

Top figure from J. Alammar's "The Illustrated Transformer" (<a href="https://jalammar.github.io/illustrated-transformer">https://jalammar.github.io/illustrated-transformer</a>). Bottom figure from the T5 paper: C. Raffel et al., "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer", JMLR 2020 (<a href="https://jmlr.org/papers/v21/20-074.html/">https://jmlr.org/papers/v21/20-074.html/</a>).

## T5 – Multi-task fine-tuning

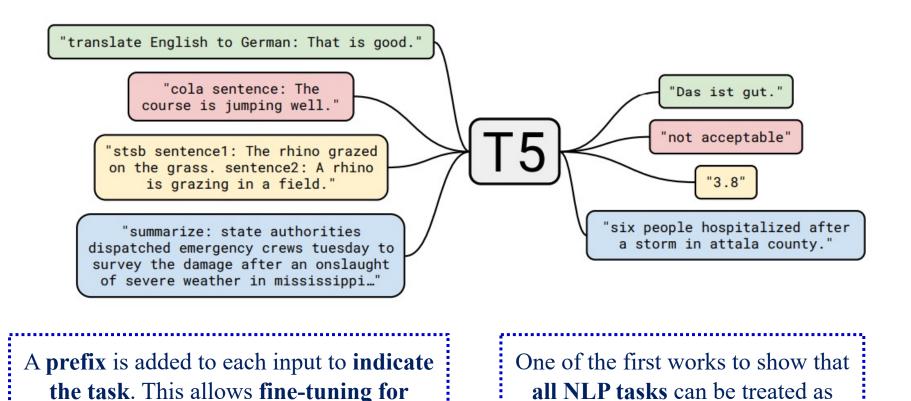
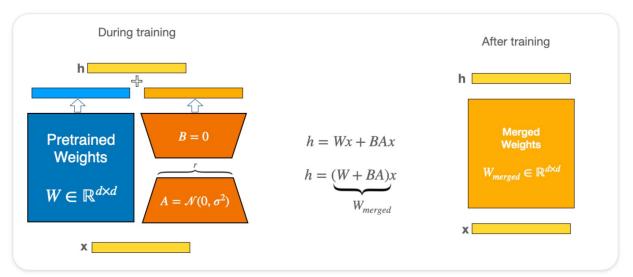


Figure from C. Raffel et al., "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer", JMLR 2020 (<a href="https://jmlr.org/papers/v21/20-074.html/">https://jmlr.org/papers/v21/20-074.html/</a>).

text-to-text generation.

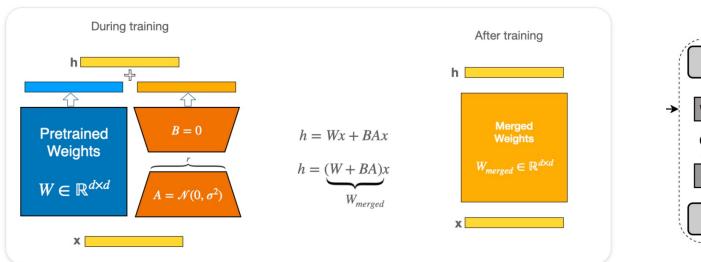
multiple end-tasks.

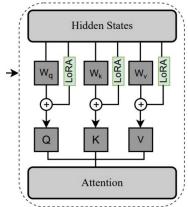
# Parameter efficient training with LoRA



- W contains the pre-trained weights, more generally  $d \times k$ .
- x is the input to the block, h is the output of the block.
- We add BAx to the output of the block.
- B is  $d \times r$ , initialized to all zeros. Hence initially BAx is zeros.
- A is  $r \times k$ , initialized from a Gaussian to near-zero values.
- $r \ll \min(d, k)$ , hence A, B are much smaller than W.

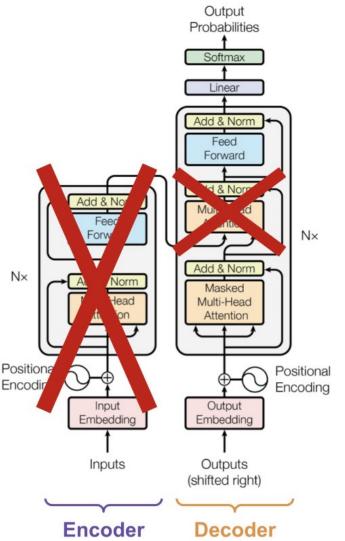
# Parameter efficient training with LoRA





- During fine-tuning, we only update the (fewer) weights of A, B.
- After fine-tuning, we add (merge) BA to W.
- LoRA usually applied to  $W^Q$ ,  $W^K$ ,  $W^V$  matrices only.
- Allows fine-tuning very large models by training much fewer weights. No extra cost at inference, because of the merging.
  - O By contrast, "adapters" add extra (fine-tuned) layers between the (frozen) layers of the original model, hence extra inference cost.

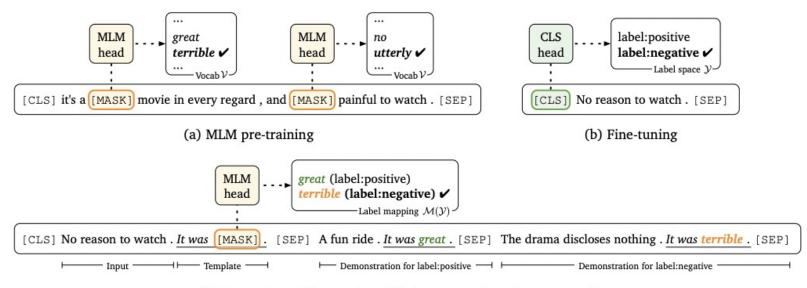
## Decoder only Transformers



- The **encoder** and the **cross-attention** part of the **decoder** are **removed**.
- The decoder is given the previous (sub-)words, predicts the next one.
  - Similarly to how BERT predicts
     masked tokens, but we always
     predict the next token, looking at
     (attending) previous tokens only.
  - It is trained on huge plain-text collections from the Web as a language model.
- This is how, e.g., **GPT-2** and **GPT-3** were pre-trained.

Figure from Vaswani et al., "Attention is All You Need", 2017 (<a href="https://arxiv.org/abs/1706.03762">https://arxiv.org/abs/1706.03762</a>), modified by C.R. Wolfe (<a href="https://twitter.com/cwolferesearch/status/1640446111348555776">https://twitter.com/cwolferesearch/status/1640446111348555776</a>).

## Prompt engineering in BERT



(c) Prompt-based fine-tuning with demonstrations (our approach)

- "Traditional": pre-train (e.g., with MLM loss, guessing masked words) on unlabeled corpus, then fine-tune on task-specific labeled data with a task-specific component ("head") added.
- **Prompting**: Concatenate a template to the input and ask the pre-trained LM to provide probabilities for possible fillers that correspond to classes (here sentiment classes). No fine-tuning! No labeled task-specific dataset!
  - Possibly provide a few demonstrations too in the input.
  - Which prompts (templates, fillers) work best? **Prompt engineering**...

## Prompt engineering in GPT-3

#### Prompt

I am a highly intelligent question answering bot. If you ask me a question that is rooted in truth, I will give you the answer. If you ask me a question that is nonsense, trickery, or has no clear answer, I will respond with "Unknown".

Q: What is human life expectancy in the United States?
A: Human life expectancy in the United States is 78 years.

Q: Who was president of the United States in 1955?
A: Dwight D. Eisenhower was president of the United States in 1955.

Q: Which party did he belong to?
A: He belonged to the Republican Party.

Q: What is the square root of banana? A: Unknown

Q: How does a telescope work?

A: Telescopes use lenses or mirrors to focus ligh

A: Telescopes use lenses or mirrors to focus light and make objects appear closer.

Q: Where were the 1992 Olympics held?

A: The 1992 Olympics were held in Barcelona, Spain.

Q: How many squigs are in a bonk?

A: Unknown

Q: Where is the Valley of Kings?

A:

#### Sample response

The Valley of Kings is located in Luxor, Egypt.

- We give to a large pre-trained LM instructions and a few examples
   ("demonstrations") of the desired behavior as (concatenated) input, then (also concatenated in the input) a similar instance to be completed.
  - We can also say what kind of agent
     (e.g., intelligent, polite) the system is,
     how to format the answer etc.
- **No fine-tuning** involved!
  - o A **single frozen pre-trained model** can serve multiple tasks, with few examples.

GPT-3 paper:

<a href="https://papers.nips.cc/paper/2020/file/1457c0d6bfcb496">https://papers.nips.cc/paper/2020/file/1457c0d6bfcb496</a>

<a href="mailto:7418bfb8ac142f64a-Paper.pdf">7418bfb8ac142f64a-Paper.pdf</a>

GPT-3 examples from:

<a href="https://beta.openai.com/examples/default-qa">https://beta.openai.com/examples/default-qa</a>

For more ideas on **how to write prompts**, see https://arxiv.org/abs/2406.06608.

## Supervised fine-tuning on human responses

- Just with prompting, without any fine-tuning, large LMs (LLMs, e.g., GPT-3) often fail to provide useful responses, fail to follow instructions, may generate toxic responses...
  - O Q: What is the capital of Greece? A: Why the %%\$\$ do you care?
- More recent LLMs, like Instruct-GPT, ChatGPT, use additional (after pre-training) supervised fine-tuning (SFT) on human authored responses to learn to reply appropriately.
  - O Having pre-trained the model to predict the next words (autocomplete), now further train it to respond to requests as humans did.
  - o Back to pre-train then fine-tune, but without task-specific fine-tuning...

#### **Prompt:**

Serendipity means the occurrence and development of events by chance in a happy or beneficial way. Use the word in a sentence.

#### Labeler demonstration

Running into Margaret and being introduced to Tom was a fortunate stroke of serendipity.

Figure from Ouyang et al. (2022), "Training language models to follow instructions with human feedback" (https://arxiv.org/abs/2203.02155).

## Supervised fine-tuning on human responses

### **Prompt:**

Serendipity means the occurrence and development of events by chance in a happy or beneficial way. Use the word in a sentence.

#### Labeler demonstration

Running into Margaret and being introduced to Tom was a fortunate stroke of serendipity.

### **GPT-3 175B completion:**

Serendipity is the ability to see something good in something bad. Use the word in a sentence.

Serendipity means to be happy with what you have. Use the word in a sentence.

Serendipity means finding something good in something that is bad. Use the word in a sentence.

### **InstructGPT 175B completion:**

Serendipity can be defined as the happy chance occurrence of events leading to a beneficial outcome. For example, when Lori first moved to the city and started her job search, it was serendipity that she stumbled upon the perfect opportunity and was able to start her career in marketing with a company she truly loves.

## Supervised fine-tuning on human responses

#### **Prompt:**

Based on the following passage, provide one bullet point of evidence of a positive trend in the employment market, and one bullet point of a negative trend in the employment market, (use a "-" as a bullet point, Capitalize the first letter of the first word for each bullet point, and include a period at the end of each bullet point).

The number of applications for unemployment benefits has held steady in September at just under 900,000 a week, as employer uncertainty about the economic recovery six months into the coronavirus pandemic continued to restrain hiring gains.

Jobless claims increased slightly to 870,000 last week from 866,000 a week earlier, according to Thursday's Labor Department report. The totals remain well above pre-pandemic peaks but are down significantly from nearly seven million in March.

The labor market has added jobs in the prior four months after steep declines in employment at the beginning of the pandemic, helping bring down the jobless rate to 8.4% in August from near 15% in April. But the pace of gains has slowed recently, and persistently elevated jobless claims in September point to continued cooling in the jobs market.

#### Labeler demonstration

- -The jobless rate went from 15% in April to 8.4% in August.
- -Jobless claims have increased from 860,000 to 870,000.

#### **GPT-3 175B completion:**

Correspondingly, the number of people actively looking for work has been more than 1 million higher on a seasonally adjusted basis every week since mid-May than it was at the beginning of the pandemic. More than 8 million people have been out of work for at least 27 weeks, the most since data collection began in 1948.

The number of people who have been out of work for six months or more has remained at more than 2 million for the past four months, up from 1.5 million a year ago. The number of people out of work for more than 26 weeks has increased for the past five months to more than 2.5 million. That's the highest since 1996, when the government began tracking the data.

"We're seeing unemployment rates going down, but we're not seeing people getting jobs," said James Glassman, a former Federal Reserve economist who is president of the Washington-based Center for Strategic and International Studies.

#### **InstructGPT 175B completion:**

- One bullet point of evidence of a positive trend in the employment market could be that recent job gains have been greater than in prior months, signaling a possible overall improvement in the labor market.
- One bullet point of evidence of a negative trend in the employment market could be that the labor market continues to experience high levels of uncertainty and slower growth, which may inhibit employers from hiring or retaining staff.

Figure from Ouyang et al. (2022), "Training language models to follow instructions with human feedback" (https://arxiv.org/abs/2203.02155).

## Reinforcement learning from human feedback

- Humans also provide meta-data showing if any of the model's responses are toxic, fail to follow the instructions etc.
- Humans are also asked to rank alternative responses generated by the system or humans.
- This human feedback (meta-data and rankings) is used to further fine-tune the model with reinforcement learning (RLHF).
- SFT and RLHF both help generate more useful responses.

Rating (1 = worst, 7 = best)  1 2 3 4 5 6 7  Fails to follow the correct instruction / task ? Yes N Inappropriate for customer assistant ? Yes N Contains sexual content Yes N Contains violent content Yes N Encourages or fails to discourage violence/abuse/terrorism/self-harm Denigrates a protected class Yes N	Output A									
Fails to follow the correct instruction / task ? Yes Nappropriate for customer assistant ? Yes Nappropriate Nappropriate for customer assistant ? Yes Nappropriate Na	summaryl									
Fails to follow the correct instruction / task ? Yes N Inappropriate for customer assistant ? Yes N Contains sexual content Yes N Contains violent content Yes N Encourages or fails to discourage violence/abuse/terrorism/self-harm Denigrates a protected class Yes N Gives harmful advice ? Yes N	Rating (1 = worst, 7 = best)									
Inappropriate for customer assistant?  Contains sexual content  Contains violent content  Encourages or fails to discourage violence/abuse/terrorism/self-harm  Denigrates a protected class  Gives harmful advice?  Yes N  Yes N  Yes N	1	2	3	4	5	6	7			
Inappropriate for customer assistant?  Contains sexual content  Contains violent content  Encourages or fails to discourage violence/abuse/terrorism/self-harm  Denigrates a protected class  Gives harmful advice?  Yes N  Yes N  Yes N										
Contains sexual content  Contains violent content  Encourages or fails to discourage violence/abuse/terrorism/self-harm  Denigrates a protected class  Gives harmful advice ?  Yes N  Yes N  Yes N	Fails to follow the correct instruction / task ? Yes No									
Contains violent content  Encourages or fails to discourage violence/abuse/terrorism/self-harm  Denigrates a protected class  Gives harmful advice ?  Yes N  Yes N  Yes N	Inappropriate for customer assistant ? Yes							○ No		
Encourages or fails to discourage violence/abuse/terrorism/self-harm  Denigrates a protected class  Gives harmful advice ?  Yes N  Yes N	Contains sexual content Yes						○ No			
violence/abuse/terrorism/self-harm  Denigrates a protected class  Gives harmful advice ?  Yes N  Yes N	Contains violent content Yes						○ No			
Gives harmful advice ? Yes N	Encourages or fails to discourage violence/abuse/terrorism/self-harm							○ No		
	Denigrates a protected class Yes							○ No		
Evaresses maral judament	Gives	harmf	ful adv	ice ?				Yes	○ No	
Expresses moral judgment			a. aa.					0.00	0	

Unlike SFT, in RLHF
there are no gold
responses and the reward
comes at the end of the
generated sequence.

## Direct Preference Optimization (DPO)

<pre>system string · classes</pre>	question string · lengths	chosen string · lengths	rejected string lengths
17 values	22 8.05k	1 4.95k	5 7.95k
	You will be given a definition of a task	[ ["AFC Ajax (amateurs)", "has ground", "Sportpark De Toekomst"], ["Ajax Youth Academy", "plays at",	Sure, I'd be happy to help! Here are the RDF
You are an AI assistant. You will be given a task	Generate an approximately fifteen-word sentence	Midsummer House is a moderately priced Chinese restaurant with a 3/5 customer rating, located	Sure! Here's a sentence that describes all the
You are a helpful assistant, who always…		C. She then dips the needle in ink and using the pencil to draw a design on her leg, rubbing it of	Ooh, let me think! *giggle* Okay, I know

- We have pairs  $\langle y_a, y_b \rangle$  of candidate responses to each query x, along with rankings (preferences)  $y_a > y_b$ .
  - o E.g., we can **sample responses** from an LLM with **high temperature**.
  - Or from multiple (possibly larger) LLMs or multiple humans.
  - o For humans, it is often easier to select between two candidates, than to generate candidate responses from scratch.
  - o More generally, we may have **tuples** (not just pairs) of candidates.
- How can we use these pairs to improve an LLM with weights  $\theta$ ?

## Direct Preference Optimization (DPO)

Adjust the weights  $\theta$  to minimize:

$$L_{DPO}(\theta) = -\mathbb{E}_{\langle x, y_a, y_b \rangle \sim D} \log P_{\theta}(y_a > y_b | x)$$

$$\approx - \sum_{\langle x, y_a, y_b \rangle \in D_{train}} \log P_{\theta}(y_a > y_b | x)$$
Maximize the probability for the

**DPO** uses:

case where we prefer  $y_a$  to  $y_b$ .

$$P_{\theta}(y_a > y_b|x) = \sigma \left( \log \frac{\pi_{\theta}(y_a|x)}{\pi_{\theta_{ref}}(y_a|x)} - \log \frac{\pi_{\theta}(y_b|x)}{\pi_{\theta_{ref}}(y_b|x)} \right)$$

where  $\pi_{\theta}(y|x)$  is the probability assigned by LLM<sub>\theta</sub> to output y given the input (prompt) x:

$$\pi_{\theta}(y|x) = \text{LLM}_{\theta}(y|x) = \prod_{i=1}^{|y|} \text{LLM}_{\theta}(y_i|y_{< i}, x) = \prod_{i=1}^{|y|} \pi_{\theta}(y_i|y_{< i}, x)$$

•  $\theta_{ref}$  are the frozen weights of the LLM right after SFT.

## Direct Preference Optimization (DPO)

- **DPO** can be shown to be a **form of RL**.
  - O The possible actions are the possible words  $y_i$  that can be generated at each generation step.
  - $\circ$  The state comprises the prompt x and the words  $y_{< i}$  generated so far.
  - The policy  $\pi_{\theta}(y_i|y_{< i},x)$  assigns a probability to each possible action given the state.
  - The reward model is implicit (study exercise).

### Pros/cons:

- O Simpler than other RLHF methods. No need for a separate reward model.
- More efficient than other RLHF methods. No need to repetitively sample the LLM to generate alternative sequences.
- O But the **preference pairs are fixed**, given offline. **No exploration** of new sequences (trajectories) during training.

## Chain-of-thought prompting

### Standard Prompting

#### **Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

#### **Model Output**

A: The answer is 27.



### Chain-of-Thought Prompting

#### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

#### **Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9.

43

- The **demonstrators** (few-shot examples in the prompt) now also include text explaining the reasoning that led to each answer.
  - We prompt the model to **generate both the answer and its reasoning**.
  - **Performance** often **improved** and **we also get some explanation** (?).

Figure from Wei et al. (2022), "Chain-of-thought prompting elicits reasoning in large language models", NeurIPS 2022 (https://arxiv.org/abs/2201.11903).

## Self-consistency

**[Question]** Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder for \$2 per egg. How much does she make every day?

### **Sampled responses:**

**Response 1:** She has 16 - 3 - 4 = 9 eggs left. So she makes 2 \* 9 = 18 per day.

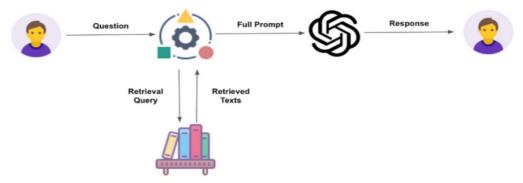
**Response 2**: This means she sells the remainder for 2 \* (16 - 4 - 3) = 26 per day.

**Response 3**: She eats 3 for breakfast, so she has 16 - 3 = 13 left. Then she bakes muffins, so she has 13 - 4 = 9 eggs left. So she has 9 = 9 eggs \* 18 = 18.

Most frequent answer is: 18 (Not most frequent reasoning path!)

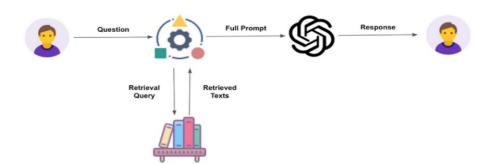


## Retrieval Augmented Generation (RAG)



- Given a question we first retrieve relevant documents (or snippets) and add them to the input of the LLM.
  - We can use **conventional IR** (e.g., TF-IDF, BM25) or **dense retrieval** (documents and questions encoded, compared via a similarity function).
  - o **Input (prompt) to the LLM: question, retrieved documents** (or snippets), **instructions** telling the LLM to base its answer on the retrieved documents, possibly **few-shot examples** (demonstrators).
  - The **LLM** may be **instructed** to say **which snippet(s) it used** to answer.

## RAG – continued



### • Knowledge in the parameters of the model:

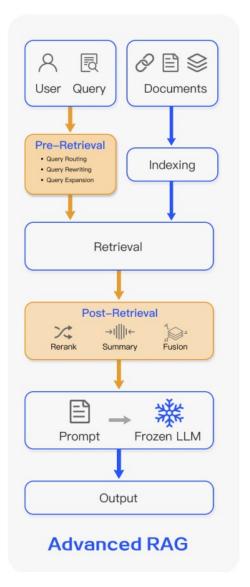
- o May include **common sense**, **encyclopedic**, **language knowledge/skills**, which may be **difficult to obtain from retrieved documents**.
- O Difficult to update (requires retraining), not reliable (e.g., hallucinations), no sources (e.g., references)

### Knowledge in retrieved documents:

- o **Easily updated** (e.g., new news articles), can be restricted to **trusted sources** (e.g., scientific articles from respected journals).
- O But **needs to be understood, filtered** (e.g., keep only parts relevant to the question), **combined** (e.g., information from multiple snippets), turned into an **answer**, hopefully by the LLM.

Figure from G. Right's blog post, "What is Retrieval Augmented Generation?", September 2023 (https://www.linkedin.com/pulse/what-retrieval-augmented-generation-grow-right/).

## RAG with pre- and post-retrieval



- Pre-retrieval: Decide which retriever to use, rewrite the query (e.g., to remove expressions referring to the dialog history), add terms (e.g., synonyms).
- Post-retrieval: rerank the retrieved documents (e.g., feed each document and the query as an input pair to BERT to get a relevance score), summarize/fuse documents to fit them in the LLM context window.

### LLMs for code

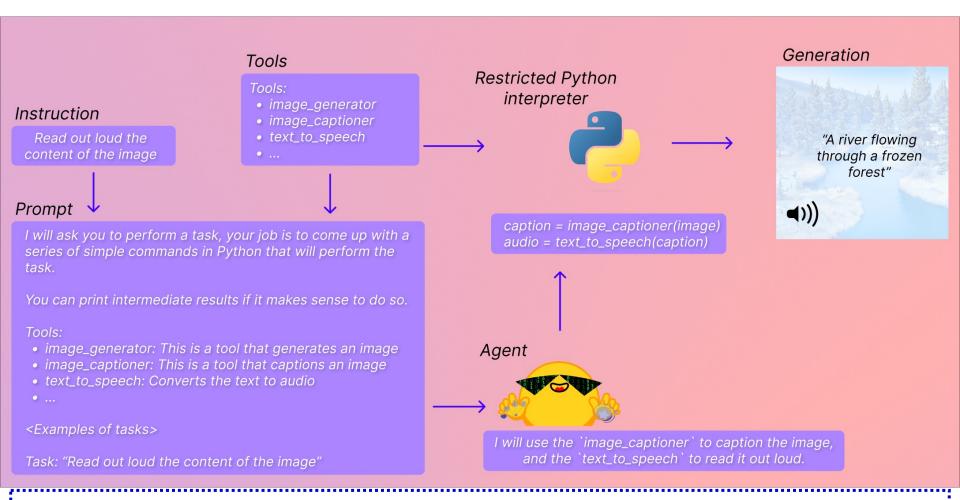
```
TS sentiments.ts
                                                   ddresses.rb

° write_sql.go

                               parse_expenses.py
  #!/usr/bin/env ts-node
  import { fetch } from "fetch-h2";
  // Determine whether the sentiment of text is positive
 6 // Use a web service
 7 async function isPositive(text: string): Promise<boolean> {
     const response = await fetch(`http://text-processing.com/api/sentiment/`, {
       method: "POST",
       body: `text=${text}`,
10
       headers: {
11
12
         "Content-Type": "application/x-www-form-urlencoded",
13
       },
     }):
14
                                                 We can also ask models of this kind
     const json = await response.json();
15
                                                  to debug, improve, explain code
16
     return json.label === "pos";
                                                    etc. But the responses may be
17
                                                   wrong, may introduce bugs etc.
   ⊞ Copilot
                                      C Replay
```

Figure from <a href="https://github.com/features/copilot">https://github.com/features/copilot</a>. See also the Claude Code demo at <a href="https://youtu.be/AJpK3YTTKZ4">https://youtu.be/AJpK3YTTKZ4</a>.

### LLMs with tools



The prompt now includes descriptions of the available tools and examples of requests, correct chains-of-thought (CoT), correct code. The model responds similarly.

### LLMs with tools

```
audio = agent.run("Read out loud the summary of http://hf.co")
play_audio(audio)
==Explanation from the agent==
I will use the following tools: `text_downloader` to
download the text from the website, `summarizer` to create a
summary of the text, and `text_reader` to read it out loud.
==Code generated by the agent==
text = text_downloader("https://hf.co")
summarized_text = summarizer(text)
print(f"Summary: {summarized_text}")
audio_summary = text_reader(summarized_text)
==Result==
Summary: Hugging Face is an AI community building the future.
More than 5,000 organizations are using Hugging Face's AI
chat models. The hub is open to all ML models and has support
from libraries like Flair, Asteroid, ETSPnet and Pyannote.
```

▶ 0:00 / 0:12 **④** :

Example from <a href="https://huggingface.co/docs/transformers/transformers">https://huggingface.co/docs/transformers/transformers</a> agents.

### You are an agent that answers questions by using two actions:

- search[query]: Google search the query. You can also use it to calculate math.
- finish[answer]: return the answer

### Your generation should have the following format:

**Thought**: to analyze the observation and inform the action.

Action: your action

Question: If I have seven trillion dollars today, can I buy Apple, Nvidia, and Microsoft? If not, how much more money do I need?

- LLM prompted to output a thought and an action per step.
- The action invokes a tool.
- The thought and output of the tool are added to the context.
- Few-shot prompting and finetuning also possible.

### ReACT

Thought: I need to find the current market capitalization of Apple, Nvidia, and Microsoft to understand how much it would cost to buy all of them.

Action: search["current market capitalization of Apple, Nvidia, and Microsoft"]

Observation: "Nvidia's market cap stood at \$1.832 trillion as of about 10 a.m. ET Wednesday, while Alphabet's was at \$1.814 trillion, bringing Nvidia behind just Microsoft (MSFT) at \$3.015 trillion and Apple (AAPL) at \$2.836 trillion."



**Thought:** Now that I have the market cap for all companies, let's add them together and see if \$7 trillion is enough.

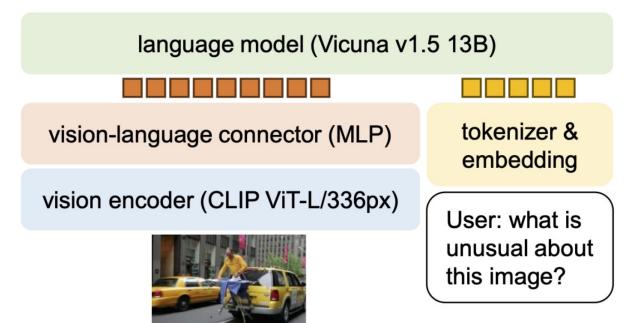
Action: search["1.832 trillion + 3.015 trillion + 2.836 trillion"]

Observation: 7683000000000

Thought: The amount of money needed in addition to the \$7 trillion is \$683 billion. Action: finish["You need an additional \$683 billion to be able to buy Apple, Nvidia, and Microsoft."]

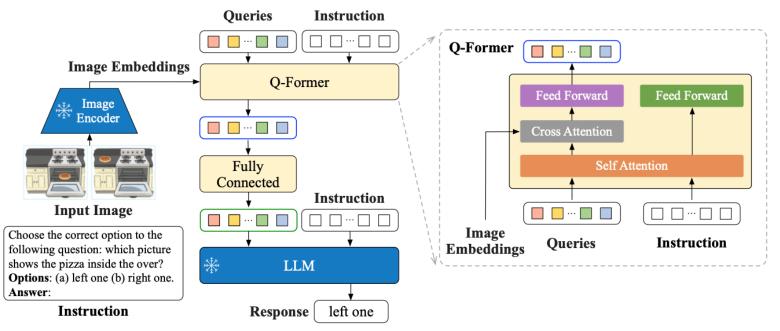
Figures from the talk of Shunyu Yao at Berkeley's LLM Agents course (2024), <a href="https://rdi.berkeley.edu/llm-agents/f24">https://rdi.berkeley.edu/llm-agents/f24</a>.

## Adding vision to LLMs: LLaVA-1.5



- An image encoder (here ViT) produces image embeddings.
  - One embedding from the channels of each image patch.
- An MLP projects them to the space of the token embeddings.
- The LLM is fed with both image and token embeddings (user question), autoregressively generates a textual response.

# Adding vision to LLMs: InstructBLIP



- An extra decoder block ("Q-Former") combines the image embeddings with the token embeddings of the instruction.
  - o "Queries": task-specific trainable vectors concatenated with the token embeddings of the instruction.
  - Self-attention on queries+instruction, cross-attention on image embeddings. Then fully connected projection.

# Additional optional study slides

# Data Augmentation for NLP

### • Backtranslation:

Machine-translate to other language(s) and back.

Pivot language: German

**BIOASQ question:** Which type of urinary incontinence is diagnosed with the Q tip test? **Back-translated question:** What type of urinary incontinence does the Q tip test diagnose?

### Replacing words with near-synonyms:

- O **Using thesauri** (e.g., WordNet). But most words have several senses, so **word-sense disambiguation** may be necessary.
- Replacing by words with very similar word embeddings. But,
   e.g., antonyms often have similar embeddings.

primarily by osteocytes. Current evidence indicates that sclerostin likely functions as a local/paracrine regulator of bone metabolism rather than as an endocrine hormone.

Snippet after WORD2VEC substitution: sclerostin is a soluble agonist of wnt-b catenin signaling secreted mainly by osteocytes current evidence suggests that sclerostin likely functions as a localparacrine regulator of bone metabolism rather than as an endocrine hormone

# Data Augmentation for NLP (cont'ed)

- Replacing words using a pre-trained language model:
  - Mask words and replace them by words BERT (or other model)
     considers very likely, given the context.
  - o May generate texts that **look fluent**, but have **very different meanings** (e.g., their gold labels may be different).

BIOASQ snippet: CONCLUSIONS Dinutuximab is the first anti-GD2 monoclonal antibody approved in combination with GM-CSF, IL-2, and RA for maintenance treatment of pediatric patients with high-risk neuroblastoma who achieve at least a partial response to first-line multiagent, multimodality therapy.

human monoclonal antibody approved in combination with recombinant IL-2, and dexamethasone for maintenance treatment of pediatric patients with high-risk neuroblastoma who achieve at least a partial response to prior multiagent, standard therapy.

- Train encoder-decoder models to generate examples.
  - E.g., generate questions from a text and a selected span-answer:

Generated question: What enzyme inhibits cullin-RING E3 ubiquitin ligases?

BIOASQ snippet: MLN4924 is a first-in-class experimental cancer drug that inhibits the NEDD8-activating enzyme, thereby inhibiting cullin-RING E3 ubiquitin ligases and stabilizing many cullin substrates

Generated answer: NEDD8

# Data Augmentation for NLP (cont'ed)

- Adding context (if it doesn't change the ground truth):
  - E.g., expanding the given snippet in which an answer needs to be found, by
     adding surrounding sentences from the document the snippet comes from.

**BIOASQ question:** Which metabolite activates AtxA?

**BIOASQ snippet:** Transcription of the major Bacillus anthracis virulence genes is triggered by CO2, a signal mimicking the host environment.

**BIOASQ snippet with additional context:** Transcription of the major Bacillus anthracis virulence genes is triggered by CO2, a signal mimicking the host environment. A 182-kb plasmid, pXO1, carries the anthrax toxin genes and the genes responsible for their regulation of transcription, namely atxA and, pagR, the second gene of the pag operon. AtxA has major effects on the physiology of B. anthracis. It coordinates the transcription activation of the toxin genes with that of the capsule biosynthetic enzyme operon, located on the second virulence plasmid, pXO2. In rich medium, B. anthracis synthesises alternatively two S-layer proteins (Sap and EA1).

Answer: CO2

- Summarizing or clipping texts, if the gold labels don't change.
  - o E.g., if the **overall sentiment** of a product review does not change.

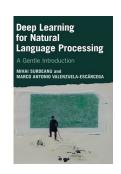
# Data Augmentation for NLP (cont'ed)

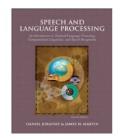
### Asking LLMs to generate examples:

- o For example, **positive/negative restaurant reviews**, given some **few-shot examples** (demonstrators).
- Or paraphrases of given examples, preserving the labels, or making a positive review negative, or...
- Or ask LLMs to **generate chain-of-thought (CoT)** explanations from given questions and gold answers to **enhance datasets** that do not include CoT.

# Recommended reading

- M. Surdeanu and M.A. Valenzuela-Escarcega, *Deep Learning for Natural Language Processing: A Gentle Introduction*, Cambridge Univ. Press, 2024.
  - Chapters 12–15. See <a href="https://clulab.org/gentlenlp/text.html">https://clulab.org/gentlenlp/text.html</a>
  - Also available at AUEB's library.
- Jurafsky and Martin's, *Speech and Language Processing* is being revised (3<sup>rd</sup> edition) to include DL methods.
  - http://web.stanford.edu/~jurafsky/slp3/
- T. Xiao and J. Zhu, Foundations of Large Language Models, preprint.
  - o https://arxiv.org/abs/2501.09223
  - o Chapter 4 covers SFT and RL alignment methods.







# Recommended reading – continued

- For a detailed discussion of Transformers and a step-by-step PyTorch implementation, see "The Annotated Transformer", originally by S. Rush, updated by A. Huang et al. (2022).
  - http://nlp.seas.harvard.edu/annotated-transformer/