

# Advanced Natural Language Processing

## 1 LLMs

### BERT

Bidirectional Encoder Representations from Transformers.

Basically a transformer encoder.

Training: BERT was trained on Books and Wikipedia.

Training objectives: Masked Language Modelling. Predict randomly masked tokens (by [MASK] token). Next Sentence Prediction (with [CLS] token).

Versions:

- BERT-Base: 12 transformer blocks, embedding dim: 768, 12 attention heads.
- BERT-Large: 24 transformer blocks, embedding dim: 1024, 16 attention heads.

Extensions:

- RoBERTa: Larger batches, removed NSP.
- SentenceBERT: BERT + siamese architecture (Sentence A and B in the same model, outputs are compared by some metric) + triplet loss (tune network such that distance between anchor sentence and positive sentence is smaller than anchor sentence and negative sentence)
- DistillBERT: Larger teacher network creates soft probability distribution labels. Smaller student model tries to replicate these distributions.

### GPT

Generative Pretrained Transformer.

Only a transformer decoder, since its task is not seq2seq, where you might want to refer to another part of the input sentence, but instead just generate next tokens (still the decoder receives already generated tokens as inputs).

Goal is to use few-shot learning and prompting instead of fine-tuning and task specific architectures.

Autoregressive: Feeds back generated tokens in decoder input.

### Other LLMs

#### T5

Text-to-Text-Transfer-Transformer.

Encoder and decoder block, similar to original Attention is all you need paper but with slight modifications.

## **RETRO**

Retrieval-Enhanced Transformer.

Enhanced by large text database.

In total 4% of GPT-3 size.

## **Flash Attention**

Self-attention is the bottleneck when training/inferencing transformers (quadratic time). During Attention, tensors are moved from HBM (High Bandwidth (GPU) Memory) to SRAM (Shared random-access (GPU) memory) multiple times: Before each operation, the tensors are moved from HBM to SRAM and then written back to HBM.

Flash Attention fuses multiple operations into one kernel, thus saving load and write operations. It is becoming the standard for transformers.

## **Fine-Tuning**

Traditional fine-tuning: Re-train all parameters with a small labelled dataset for example on domain knowledge.

Traditional fine-tuning with freezing: Freeze a part of the model.

## **Prompting**

Removes the need for additional layers for a downstream task. Uses the same format as the pretraining objective of LLMs: We ask the LLM to generate the answer based on a task-specific input.

Requires no new parameter, can be practiced on a closed-source LLM.

## **Zero-Shot**

Simple prompting is equivalent to zero-shot learning. The model predicts the answer based on the question.

Example: "Translate English to French: cheese = ".

## In-Context Learning and Few-Shot Learning

Learning that contains examples of solutions to the task.

Example: "Translate English to French: sea otter = loutre de mer. cheese = ".

## Prompting Terminology

Pattern:

A function that maps input to text (i.e. template for  $x$ ).

Example:  $f(x) = \text{"Review : } x\text{"}$

Verbalizer:

A function that maps a label to text (i.e. template for  $y$ ).

Example:  $v(y) = \text{"Sentiment : } y\text{"}$

Zero-shot prompting is using one pattern. Few-shot prompting is using many pairs of patterns and verbalizers and one final pattern.

Choosing a prompt is important and non-trivial, since experiments show different patterns and verbalizers exhibit large variance on the result accuracy. Additionally, it is uncertain how and why in-context learning works.

Reasons:

- Different input and output space distributions decrease/increase performance.
- LLM sees demonstrations not as ordered pairs.
- Highly dependent on choice, order and term frequency.

Prompting still works really well in practice and can perform better than smaller task-specific fine tuned models such as BERT.pretrained

Prompt-based fine tuning: [CLS] A \_ master class. It was [MASK].

## Parameter Efficient Fine-Tuning

### Prompt Search

Learns tokens in the prompt.

Instead of fine-tuning and predicting a [CLS] token for an input sentence, use the masked prompt as input sentence.

Classical fine-tuning: [CLS] A three-hour cinema master class.

Prompt search: [CLS] A \_ master class. It was [MASK].

### AutoPrompt

Prompt search method.

Update tokens in the pattern using gradient-guided search.

Prompt Template: {Original sentence}[T][T][T][T][T][MASK], where T are trigger Tokens that are determined using gradient-guided search.

Example: a real joy. atmosphere alot dialoge Clone totally [MASK].

### **Prompt Tuning**

Attaches a learnable embedding to the input pattern.

Example: [EMBED] [CLS] A \_ master class. [MASK].

Goal: Learn the [EMBED] token.

### **BitFit**

Only tunes bias in attention and linear layers.

Hypothesis: LLM does not need to learn new linguistic features, just domain knowledge.

### **Adapters**

In the Transformer block, after each feed-forward block, add an adapter block and a skip connection.

The adapter layer is built up like an autoencoder (down- and the up-projects the input). Then, only fine tune these adapter blocks.

### **LoRA**

Low-Rank Adaption of LLMs.

At each layer, substitute the weight update on  $W$  with an update on the low-rank decomposition of  $W$  ( $AB$ ).

$$W_{finetuned} = W_{pretrained} + \Delta W = W_{pretrained} + AB$$

Low rank means the matrix has much less params than  $W$ .

### **(IA)3**

Infused Adapter by Inhibiting and Amplifying Inner Activations.

Element-wise rescaling of model activations with a learned vector.

Separate learned vectors for each task.

## Human Preference Tuning

Fine tuning is not enough to make a model fully non-harmful/friendly and performant/helpful at the same time.

Training on human preferences (rankings, scores, etc.) is not differentiable, since no mathematically differentiable function was used to calculate the scores. Therefore, no supervised learning is possible.

Instead use Reinforcement Learning.

## Reinforcement Learning with Human Feedback

Create a smaller LM, the reward model, which is trained on labelled data (4 answers to a prompt ranked from best to worst).

The reward model then calculates a reward for each output of the LLM. This reward is used with the PPO algorithm to update the LLM weights.

To not overfit and let the LLM just "please" the reward model, we calculate a second loss from the LLM output and the old, frozen LLM output. This loss is added to the reward model output before being used in the PPO algorithm.

Observations:

- Smaller models get slightly worse, bigger models get even better.
- Compatible with specialized models (e.g. model for coding)
- Requires much more human annotation than fine tuning. Possible solution: Use LLMs themselves to identify bad answers (Constitutional AI/RLAIF)

## RLAIF

Reinforcement Learning with AI Feedback.

Uses a constitution (Constitutional AI) instead of humans.

1. Make the model revise harmful answers according to a random principle in the constitution. Then use fine tuning with the revised answers.
2. Use an off-the-shelf LLM to rate the answers of the fine-tuned model.
3. Train a reward model on these pairs.
4. Train the final model with reinforcement learning.

## DPO

Use Direct Policy Optimization instead of reinforcement learning.

Reformulates the problem into a single cross entropy loss.

So far better than PPO, but the debate is open.

## **2 Embeddings and XAI**

### **Embeddings**

#### **Static Embeddings**

Word2Vec, GloVe, FastText.

One embedding per word, the same for every context.

#### **Contextual Embeddings**

ELMo, BERT, GPT.

One embedding per word, different for every context.

Produced by a transformer encoder.

#### **Subword Embeddings**

##### **Byte-Pair Encoding (BPE)**

Merge characters into tokens until desired vocabulary size is reached.

### **2.0.1 Analyzing Contextual Embeddings**

Goal: Find information encoded in the embedding such as about language, the world or harmfulness.