

Coursera: Generative AI with Large Language Models

Table of Contents

Week 1 Part I. Introduction to LLMs and the generative AI project lifecycle	2
LLM Use cases:	2
History of generative AI:	2
Transformer:	3
Prompting and prompt engineering:.....	8
Generative configuration: influence the model's output during inference	9
Generative AI project Lifecycle	10
Week 1 Part II. LLM pre-training and scaling laws	10
Pre-training large language models	10
Computational challenges of training LLMs	11
Efficient multi-GPU compute strategies.....	13
Scaling laws and compute-optimal models.....	14
Pre-training for domain adaption	15
Week 2 Part I. Fine-tuning LLMs with instructions.....	17
Instruction fine-tuning	17
Fine-tuning on a single task.....	19
Multi-task instruction fine-tuning:.....	20
Model evaluation	22
Benchmark	26
Week 2 Part II. Parameter-efficient fine tuning (PEFT)	27
Introduction of PEFT.....	27
Low-Rank Adaptation of Large Language Models (LoRA)	29
Prompt tuning using soft prompts	31
Week 3 Part I. Reinforcement learning with human feedback (RLHF).....	33
Aligning models with human values.....	33
Reinforcement learning from human feedback (RLHF).....	33
Obtaining feedback from humans.....	35
Reward model	36
Fine-tuning with reinforcement learning	37
Proximal policy optimization (PPO).....	38
RLHF: Reward Hacking.....	39
Scaling human feedback.....	40
Week 3 Part II. LLM-powered applications	41
Model optimizations for deployment	41
Generative AI Lifecycle cheat sheet	43
Using the LLM in applications	43
Interacting with external applications	46
Helping LLMs reason and plan with chain-of-thought	46
Program-aided language models (PAL)	47
ReAct: Combining reasoning and action	49

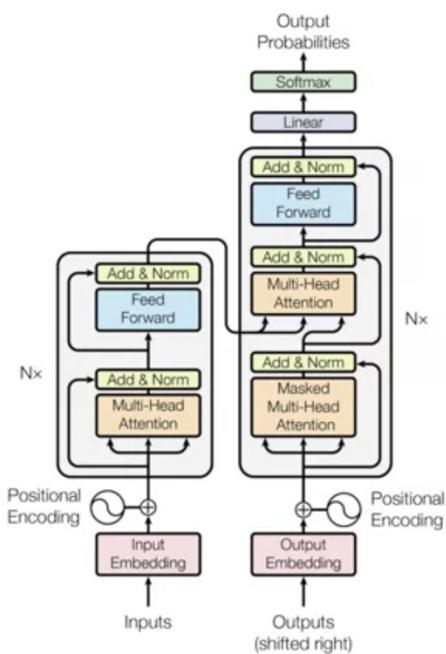
Week 1 Part I. Introduction to LLMs and the generative AI project lifecycle

LLM Use cases:

- Basic chatbots: built upon next word generation
- Text summarization
- Language translation - can also be from natural language to machine code
- Entity Extraction (under information retrieval)
- Augmenting LLMs by invoking external APIs from them

History of generative AI:

1. Weakness of RNN: Generating text with RNN: we need to scale the resources the model use to get as many texts as possible to make reliable prediction. Models needs to have an understanding of the whole sentence or even the whole document. The problem here is that language is complex. In many languages, one word can have multiple meanings. These are homonyms.
2. Transformer:
 - a. It can be scaled efficiently to use multi-core GPUs,
 - b. It can parallel process input data, making use of much larger training datasets
 - c. Crucially, it's able to learn to pay attention to the meaning of the words it's processing.
 - d. The flowchart below is the transformer architect from the paper, Attention is All You Need.



Transformer:

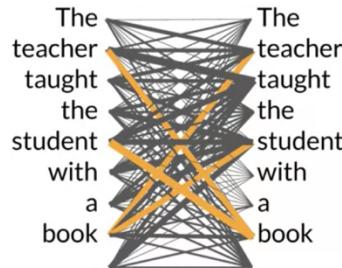
Steps in Transformer:

- Step 1: Using a tokenizer to convert text to numbers (using vocabulary indexing/stemming/ lemmatising)
- Step 2: Create embedding vectors
- Step 3: Positional encoding - as tokens are processed in parallel, add positional encoding to preserve info about word order
- Step 4: Self-Attention layer - analyze relationships between tokens
 - Multi-headed self-attention: Multiple sets of self-attention weights are learnt in parallel (number of heads varies, 12-100 common)
 - Each head learns a diff aspect of the language
- Step 5: Feed forward network- get an output vector of logits (raw probability scores) with values proportional to the probability of each token
- Step 6: Softmax layer - normalize logits to probability score
- Step 7: Find most probable next token (this step can be varied based on problem)

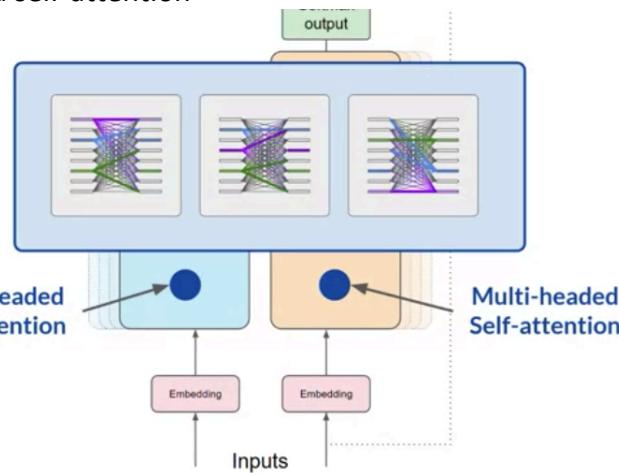
Detailed illustration of each component in Transformer:

- 1) Attention map: Learn the relevance of each word.

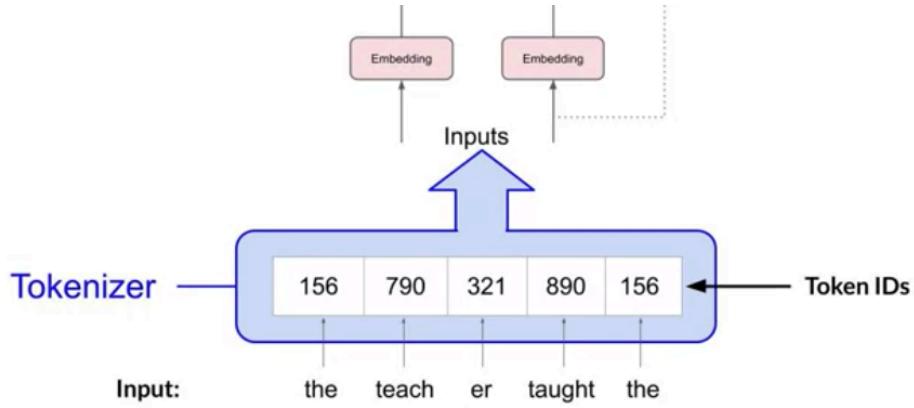
Self-attention



Multi-head self-attention

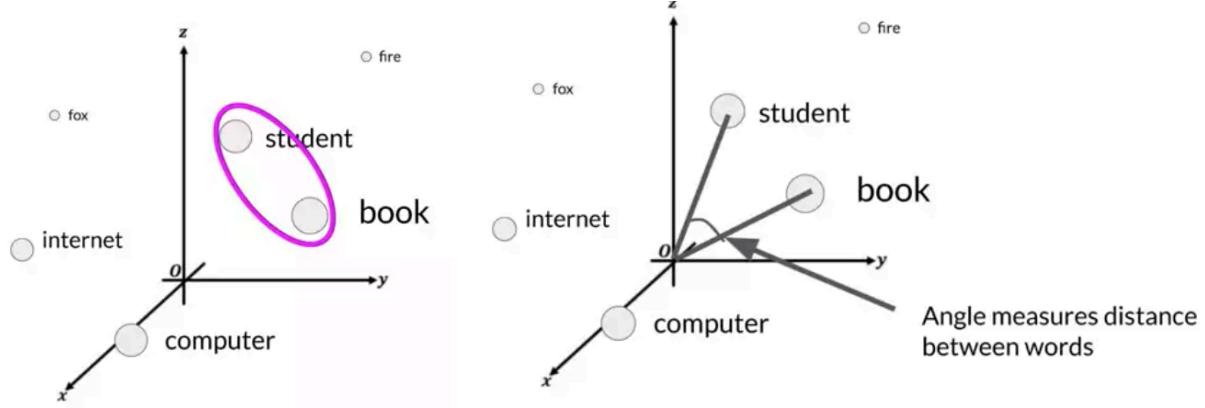


- 2) Tokenizer:



3) Embedding space:

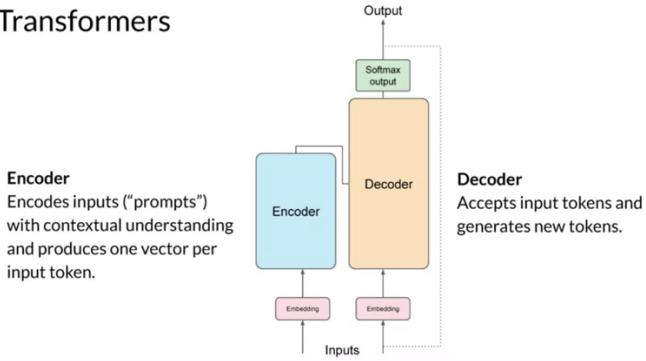
For simplicity, if you imagine a vector size of just three, you could plot the words into a three-dimensional space and see the relationships between those words. You can see now how you can relate words that are located close to each other in the embedding space, and how you can calculate the distance between the words as an angle, which gives the model the ability to mathematically understand language.



4) Encoder and decoder:

The encoder encodes input sequences into a deep representation of the structure and meaning of the input. The decoder, working from input token triggers, uses the encoder's contextual understanding to generate new tokens. It does this in a loop until some stop condition has been reached.

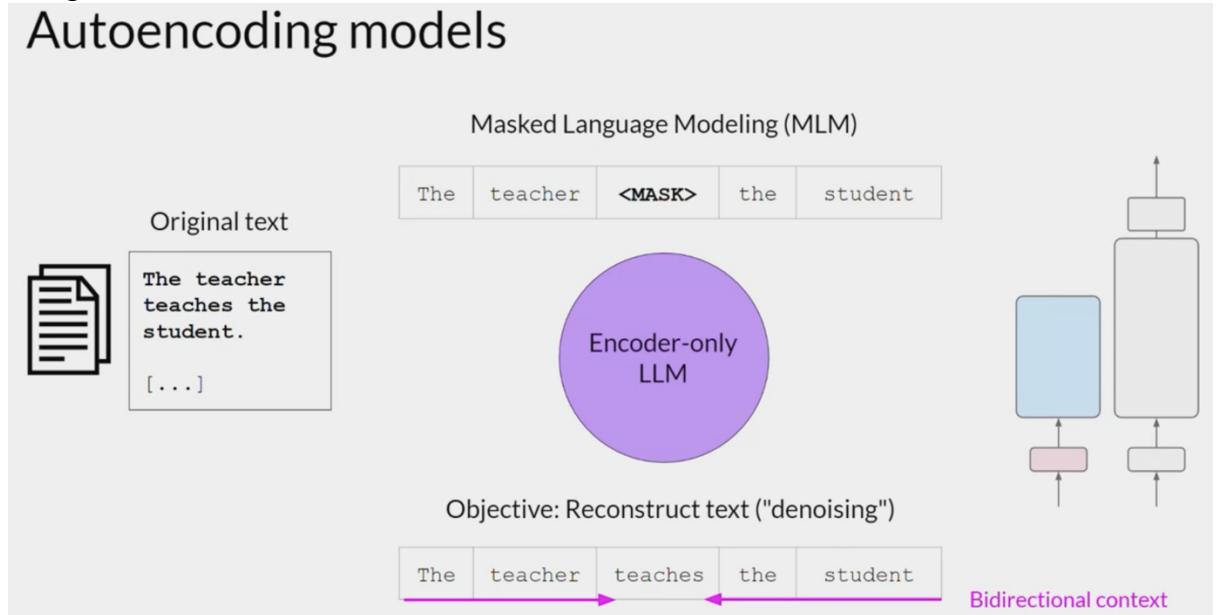
Transformers



Three types of transformer models:

- Encoder only models (aka autoencoding models):
 - work as sequence-to-sequence models, but without further modification, the input sequence and the output sequence of the same length. Their use is less common these days, but by adding additional layers to the architecture, you can train encoder-only models to perform classification tasks such as sentiment analysis, BERT is an example of an encoder-only model.
 - How it works:
tokens in the input sequence are randomly mask, and the training objective is to predict the mask tokens in order to reconstruct the original sentence. This is also called a denoising objective. Autoencoding models spilled bi-directional representations of the input sequence, meaning that the model has an understanding of the full context of a token and not just of the words that come before. → good at sentence classification tasks: sentiment analysis, name entity recognition, or word classification.

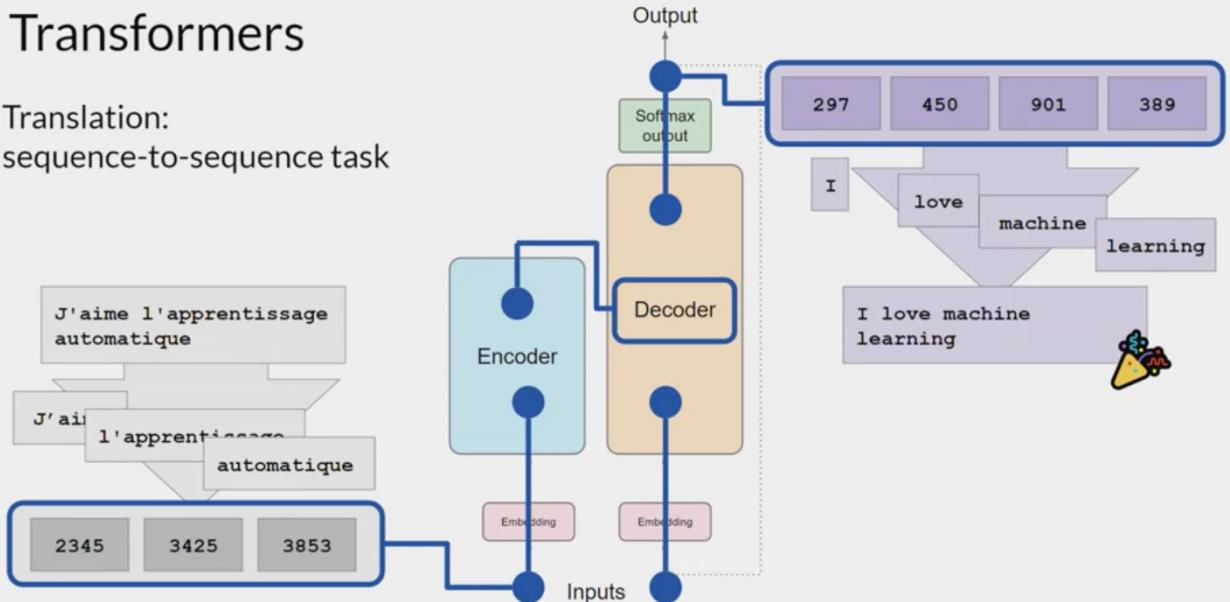
Autoencoding models



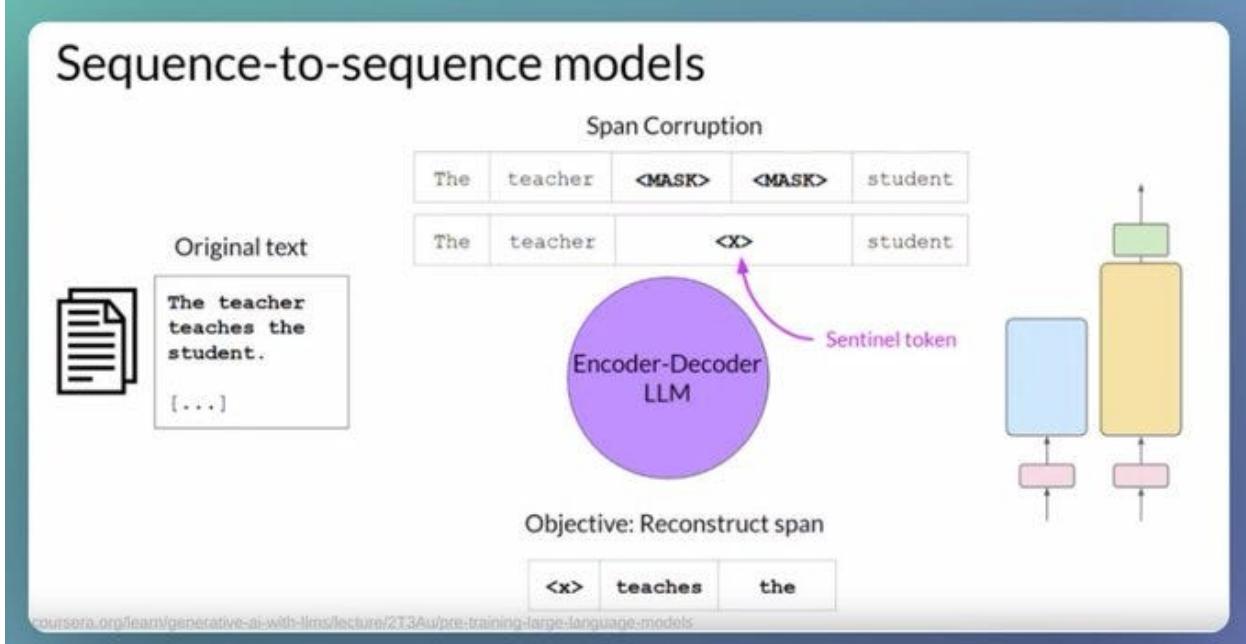
- Encoder Decoder models (aka sequence to sequence models): perform well on sequence-to-sequence tasks such as translation, summarization, and question answering, where the input sequence and the output sequence can be different lengths. You can also scale and train this type of model to perform general text generation tasks. Example: T5, BART

Transformers

Translation:
sequence-to-sequence task



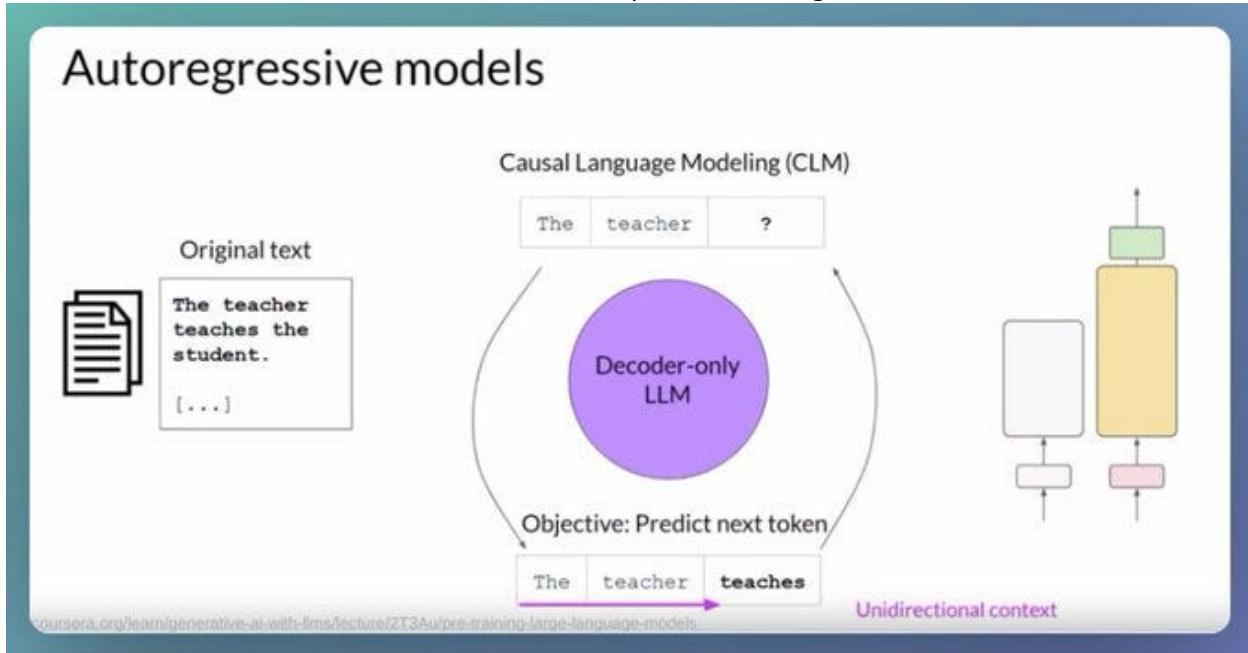
Sequence-to-sequence models



- Decoder only models (aka autoregressive models, which are pretrained using causal language modeling):
 - most-commonly used. These models can now generalize to most tasks. Popular decoder-only models include the GPT family of models, BLOOM, Jurassic, LLaMA, and many more.
 - How it works: the training objective is to predict the next token based on the previous sequence of tokens. Predicting the next token is sometimes called full

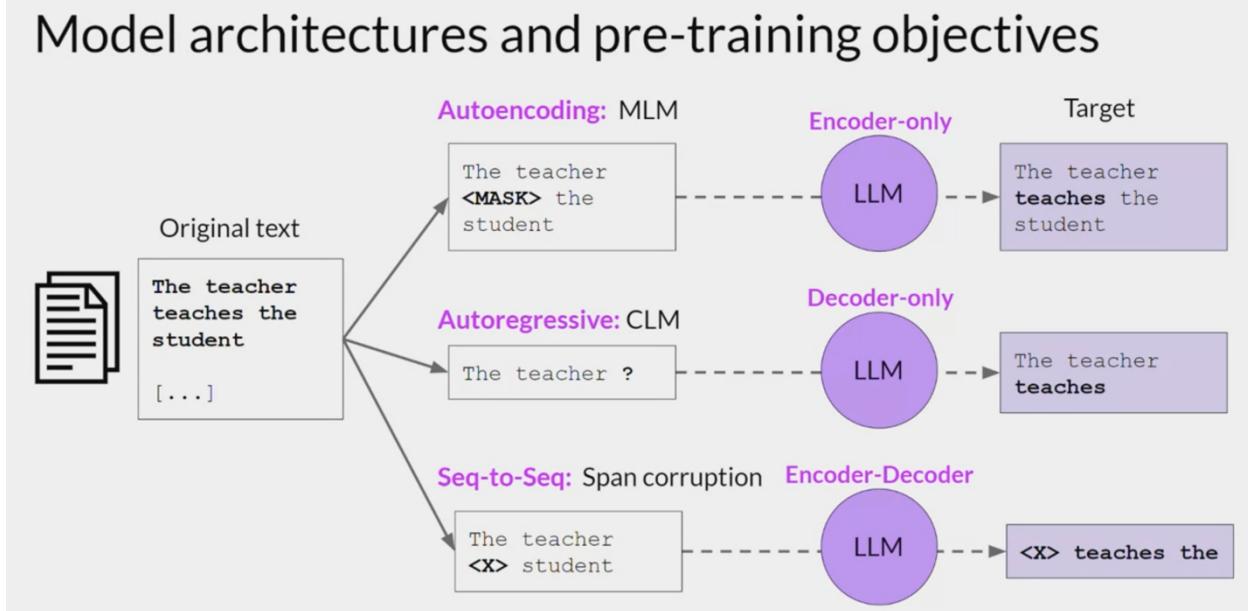
language modeling by researchers. Decoder-based autoregressive models, mask the input sequence and can only see the input tokens leading up to the token in question. The model has no knowledge of the end of the sentence. The model then iterates over the input sequence one by one to predict the following token. In contrast to the encoder architecture, this means that the context is unidirectional. By learning to predict the next token from a vast number of examples, the model builds up a statistical representation of language. → good at text generation, although larger decoder-only models show strong zero-shot inference abilities, and can often perform a range of tasks well.

Autoregressive models

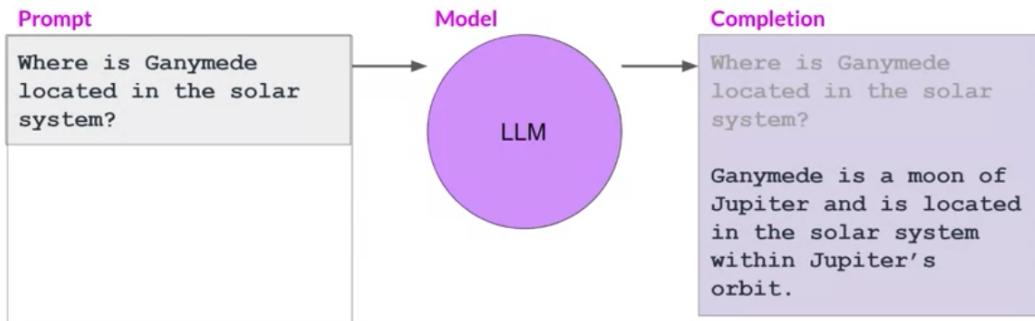


In summary,

Model architectures and pre-training objectives

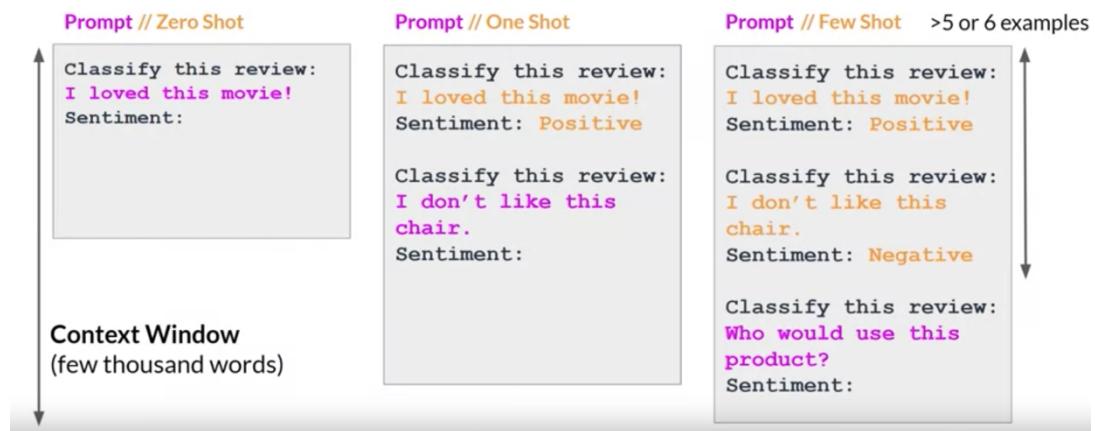


Prompting and prompt engineering:



Context window: typically a few thousand words

Summary of in-context learning (ICL)



- The largest models are surprisingly good at zero-shot inference and are able to infer and successfully complete many tasks that they were not specifically trained to perform. e.g. BLOOM
- Sometimes one-shot can perform better than few shots. Empirical studies show that if the model performance does not change when trying 1 to 6 shots, then prompt engineering may not be the best way to improve it. No need to try more than 6 shots.
- Smaller models are generally only good at a small number of tasks. E.g. BERT
- You may have to try out a few models to find the right one for your use case. Once you've found the model that is working for you, there are a few settings that you can experiment with to influence the structure and style of the completions that the model generates.

Generative configuration: influence the model's output during inference

- Max new tokens: used to limit the number of tokens that the model will generate. You can think of this as putting a cap on the number of times the model will go through the selection process. If set too low, the AI might cut off before completing its response. If set high, the AI generates longer responses, but this can increase costs and processing time. The AI stops generating once it reaches the max new tokens limit or a natural stopping point (like the end of a sentence). The next token is selected from a fixed vocabulary that the model was trained on.
- Temperature: Final scaling factor applied to softmax layer
 - Smaller temperature (<1): Probability will be strongly peaked (more prob concentrated on less no. of words) → less randomness → best for legal, medical, or financial writing
 - Larger temperature (>1): Broader flatter prob distribution → more variability and randomness → best for creative writing, storytelling, and brainstorming
 - Temperature = 1: Leaves softmax output as it is
 - Dynamic Adjustment: Some systems adjust temperature dynamically, making responses start structured and become more creative.

Formula

The softmax function is modified as:

$$P(x_i) = \frac{e^{(z_i/T)}}{\sum_j e^{(z_j/T)}}$$

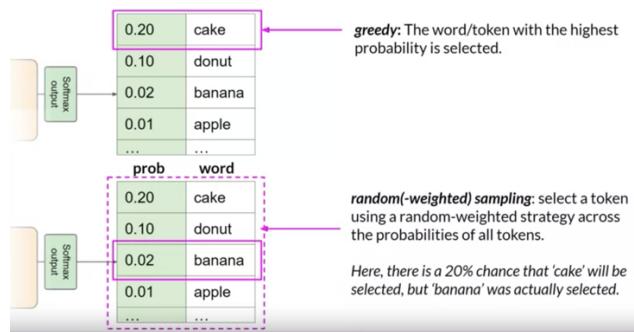
Where:

- z_i = Logit (raw score) of token i
- T = Temperature value
- $P(x_i)$ = Adjusted probability of token i

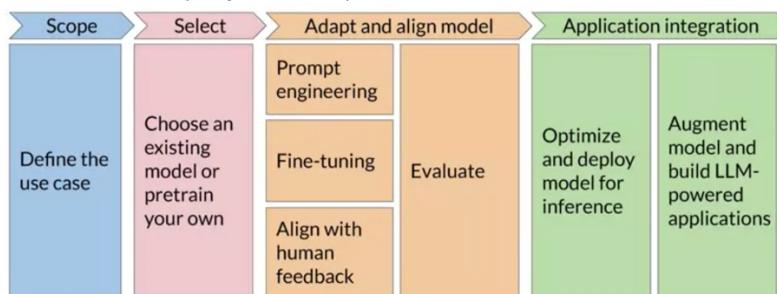
Why do we use softmax to get the probability instead of sigmoid function: Softmax is used in generative AI models because we are selecting one token from a large set of mutually exclusive

options (multi-class classification). Sigmoid is used in binary or multi-label tasks where each output is independent and does not need to sum to 1.

- Greedy vs Random sampling:
 - Greedy: model always chooses the most probable word. Prone to repeating words or sequences in long generation lengths
 - Random: Generate more natural, unrepeated, creative sequences
 - Top k and p: used to limit the random sampling to get more sensible output
 - Top k: select an output from the top k results after applying the random weighted strategy using the probabilities
 - Top p: select an output using the random-weighted strategy with the top-ranked consecutive results by probability and with a cumulative probability $\leq p$



Generative AI project Lifecycle



Week 1 Part II. LLM pre-training and scaling laws

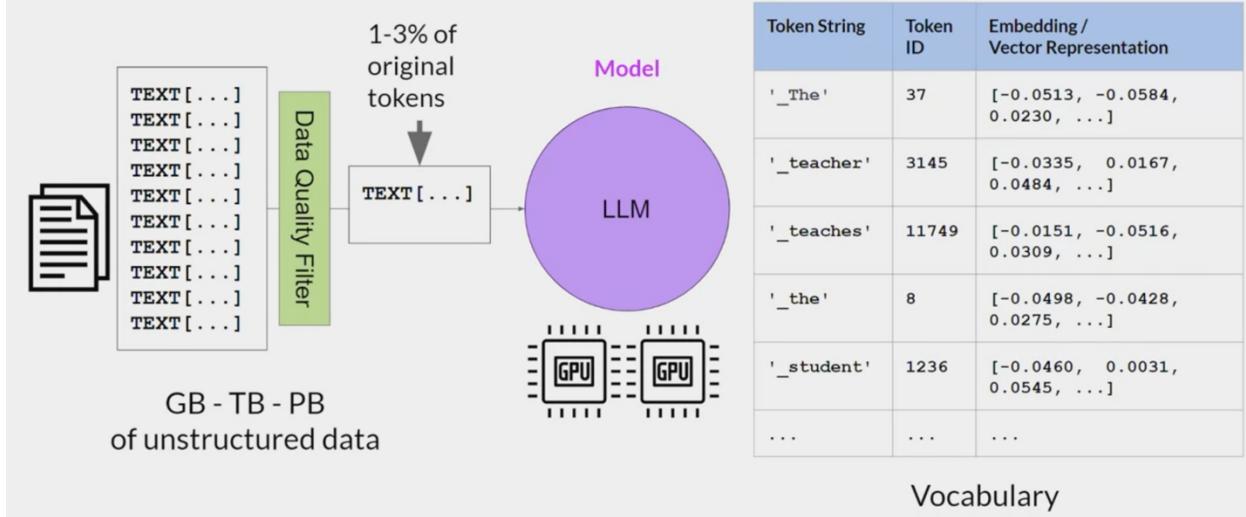
Pre-training large language models

Model Types

- 1) Foundation model:
 - a. pretrained LLM
 - b. open source: hugging face, pytorch, etc.
 - c. requires lots of GPUs

- d. you often need to process the data to increase quality, address bias, and remove other harmful content. As a result of this data quality curation, often only 1-3% of tokens are used for pre-training.

LLM pre-training at a high level



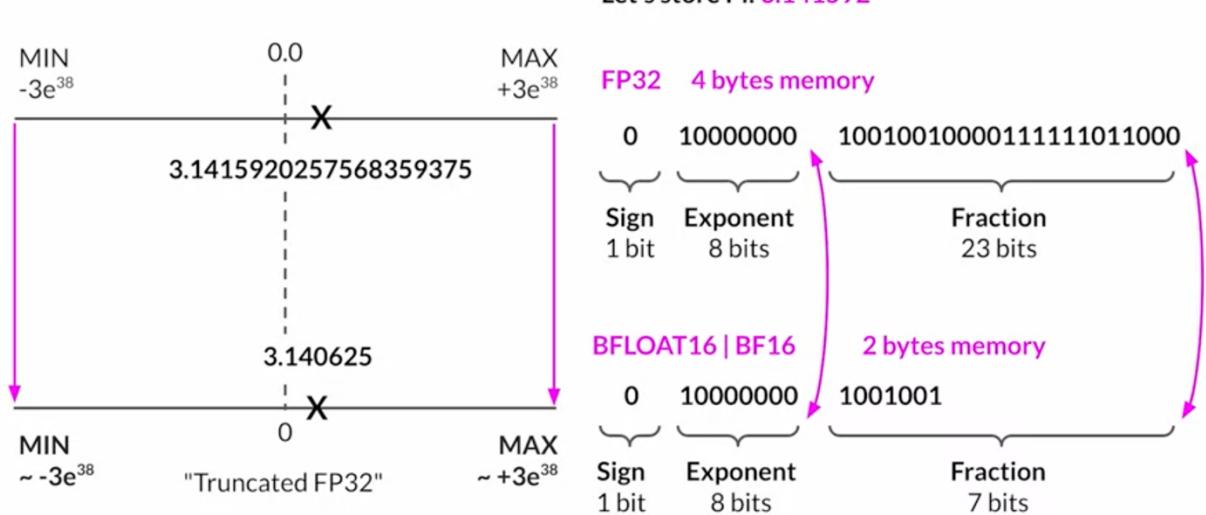
- e. the growth of the LLM is powered by: 1) introduction of transformer, 2) access to massive datasets, and 3) more powerful compute resources.

2) custom LLM: Train your own model

Computational challenges of training LLMs

- CUDA: short for Compute Unified Device Architecture, is a collection of libraries and tools developed for Nvidia GPUs. Libraries such as PyTorch and TensorFlow use CUDA to boost performance on metrics multiplication and other operations common to deep learning.
- Quantization: one technique to reduce the memory required for training. Quantization is a technique to reduce the computational and memory costs of running inference by representing the weights and activations with low-precision data types like 8-bit integer (int8) instead of the usual 32-bit floating point (float32).

Quantization: BFLOAT16



Quantization: Summary

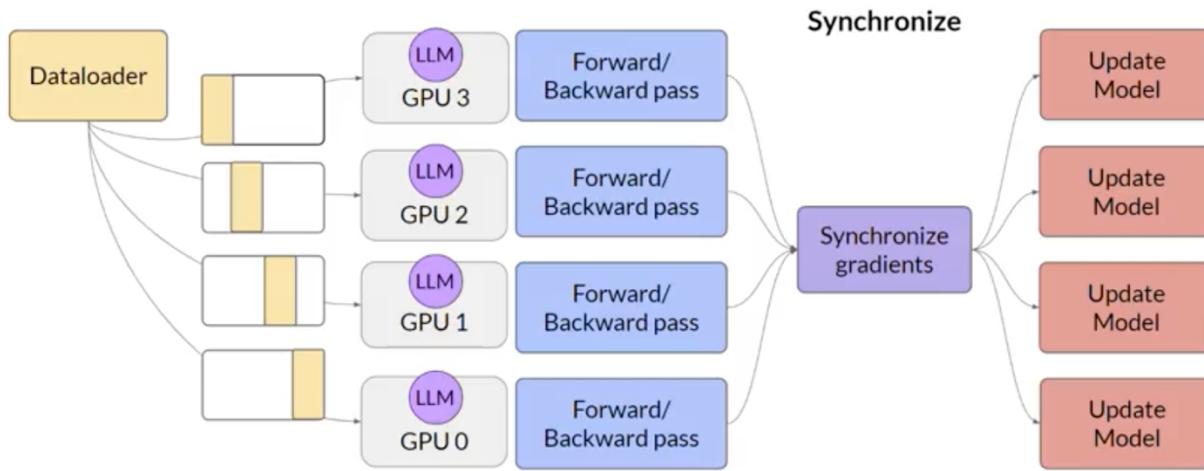
	Bits	Exponent	Fraction	Memory needed to store one value
FP32	32	8	23	4 bytes
FP16	16	5	10	2 bytes
BFLOAT16	16	8	7	2 bytes
INT8	8	-/-	7	1 byte

FLAN
T5

- Reduce required memory to store and train models
- Projects original 32-bit floating point numbers into lower precision spaces
- Quantization-aware training (QAT) learns the quantization scaling factors during training
- BFLOAT16 is a popular choice

Efficient multi-GPU compute strategies

Distributed Data Parallel (DDP)



Recap: Additional GPU RAM needed for training

	Bytes per parameter
Model Parameters (Weights)	4 bytes per parameter
Adam optimizer (2 states)	+8 bytes per parameter
Gradients	+4 bytes per parameter
Activations and temp memory (variable size)	+8 bytes per parameter (high-end estimate)
TOTAL	=4 bytes per parameter +20 extra bytes per parameter

One limitation of the model replication strategy is that you need to keep a full model copy on each GPU, which leads to redundant memory consumption. You are storing the same numbers on every GPU

→ waste of compute resource

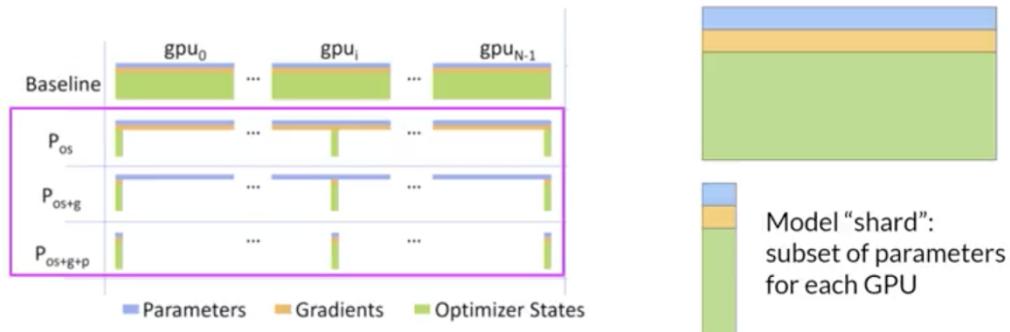
→ ZeRO (zero redundancy optimizer) eliminates this redundancy by distributing also referred to as sharding the model parameters, gradients, and optimizer states across GPUs instead of replicating them.

- ZeRO Stage 1, shuffles only optimizer states across GPUs, this can reduce your memory footprint by up to a factor of four.

- ZeRO Stage 2 also shuffles the gradients across chips. When applied together with Stage 1, this can reduce your memory footprint by up to eight times.
- ZeRO Stage 3 shuffles all components including the model parameters across GPUs.

Zero Redundancy Optimizer (ZeRO)

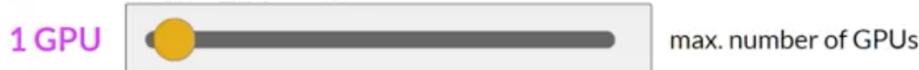
- Reduces memory by distributing (sharding) the model parameters, gradients, and optimizer states across GPUs



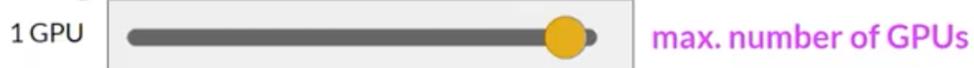
Fully Sharded Data Parallel (FSDP)

- Helps to reduce overall GPU memory utilization
- Supports offloading to CPU if needed
- Configure level of sharding via sharding factor

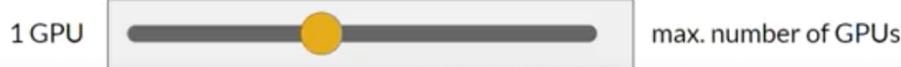
Full replication (no sharding)



Full sharding



Hybrid sharding



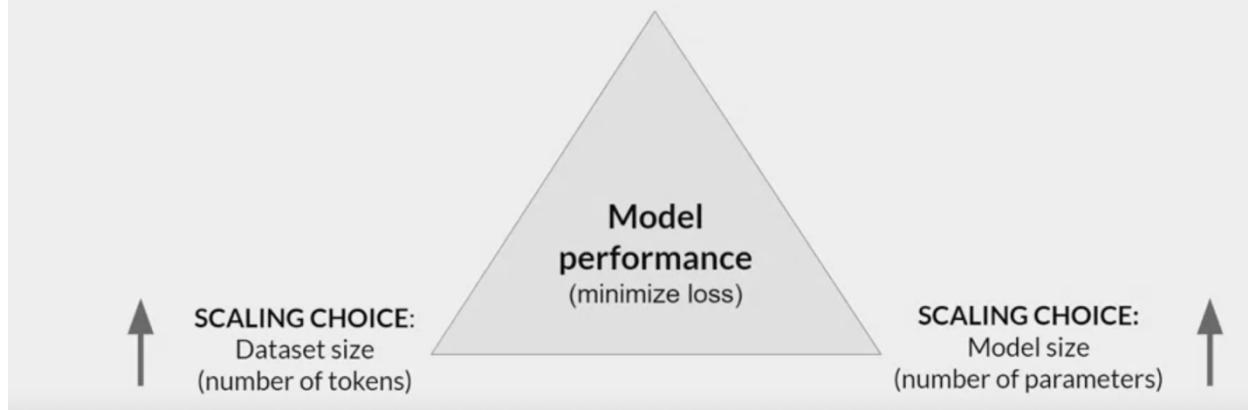
Scaling laws and compute-optimal models

There is a trade-off between minimizing compute budget and maximizing model performance.

Scaling choices for pre-training

Goal: maximize model performance

CONSTRAINT:
Compute budget
(GPUs, training time, cost)

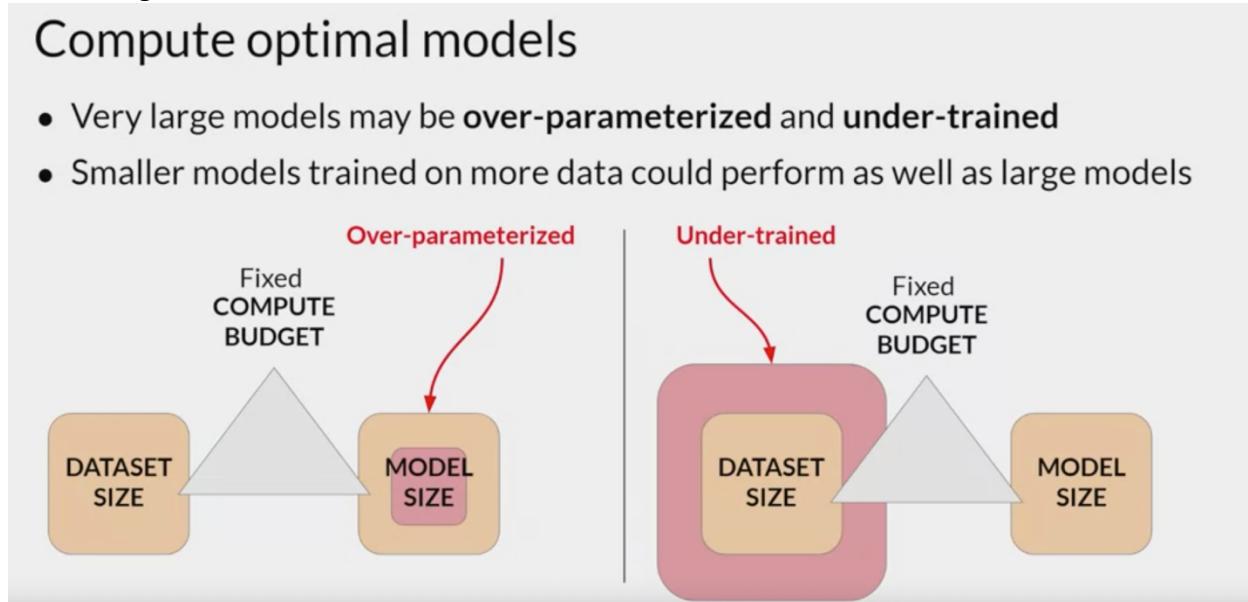


The Chinchilla paper hints that many of the 100 billion parameter large language models like GPT-3 may be

- 1) over parameterized, meaning they have more parameters than they need to achieve a good understanding of language
- 2) under trained so that they would benefit from seeing more training data. smaller models may be able to achieve the same performance as much larger ones if they are trained on larger datasets.

Compute optimal models

- Very large models may be **over-parameterized** and **under-trained**
- Smaller models trained on more data could perform as well as large models



Pre-training for domain adaption

When to use pre-training: In some domains, they use terminologies that are rarely used in the other areas and some words have very different meanings from what they are commonly used.

Example: BloombergGPT

BloombergGPT: domain adaptation for finance

BloombergGPT: A Large Language Model for Finance

Shijie Wu^{1,*}, Ozan İrsoy^{1,*}, Steven Lu^{1,*}, Vadim Dabrowski¹, Mark Dredze^{1,2},
Sebastian Gehrmann¹, Prabhanjan Kambadur¹, David Rosenberg¹, Gideon Mann¹
¹ Bloomberg, New York, NY USA
² Computer Science, Johns Hopkins University, Baltimore, MD USA
gmann16@bloomberg.net

Abstract

The use of NLP in the realm of financial technology is broad and complex, with applications ranging from sentiment analysis and named entity recognition to question answering. Large Language Models (LLMs) have been shown to be effective on a variety of tasks; however, no LLM specialized for the financial domain has been reported in literature. In this work, we present BLOOMBERGGPT, a 50 billion parameter language model that is trained on a wide range of financial data. We construct a 363 billion token dataset based on Bloomberg's extensive data sources, perhaps the largest domain-specific dataset yet, augmented with 345 billion tokens from general purpose datasets. We validate BLOOMBERGGPT on standard LLM benchmarks, open financial benchmarks, and a suite of internal benchmarks that most accurately reflect our intended usage. Our mixed dataset training leads to a model that outperforms existing models on financial tasks by significant margins without sacrificing performance on general LLM benchmarks. Additionally, we explain our modeling choices, training process, and evaluation methodology. As a next step, we plan to release training logs (Chronicles) detailing our experience in training BLOOMBERGGPT.

~51%

**Financial
(Public & Private)**

~49%

**Other
(Public)**

"You can combine data parallelism with model parallelism to train LLMs."

Is this true or false?

- True
- False

Correct

Combining data parallelism with pipeline parallelism is known as 2D parallelism. We can achieve 3D parallelism by combining data parallelism with both pipeline parallelism and tensor parallelism simultaneously.

Key Differences

Feature	Data Parallelism	Pipeline Parallelism	Model Parallelism
What is distributed?	Data (mini-batches)	Model layers (pipeline stages)	Model parameters
Memory efficiency	Low (each GPU stores a full model)	Medium	High (model is split across GPUs)
Communication overhead	High (gradient synchronization)	Moderate (activations passed between stages)	High (inter-GPU communication for activations/weights)
Best for	Large datasets, smaller models	Deep models	Very large models (GPT-3, LLaMA)

When to Use Each?

- Use **Data Parallelism** if your model fits in memory but training speed is the bottleneck.
- Use **Pipeline Parallelism** if your model is too large for a single GPU and has a sequential architecture.
- Use **Model Parallelism** if your model's parameters are too large to fit into a single GPU's memory.

Week 2 Part I. Fine-tuning LLMs with instructions

Instruction fine-tuning

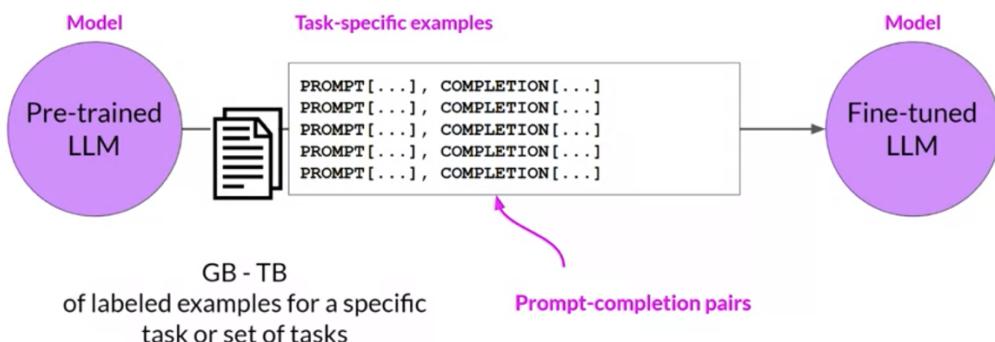
Limitation of in-text learning (prompt engineering):

1. In-text learning may not work for smaller models
2. Examples take up space in the context window and reduce the amount of room you have to include other useful information

Fine tuning:

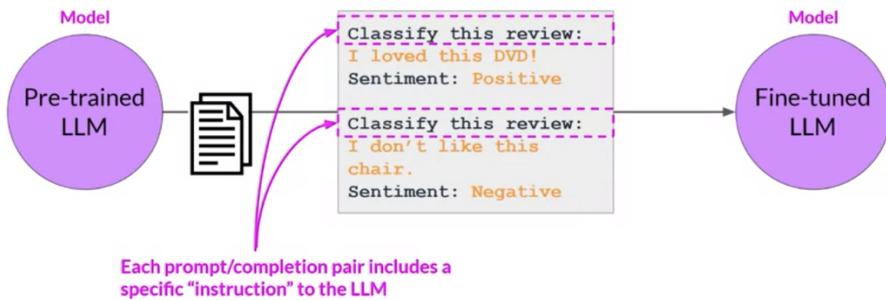
A supervised process where you use a data set of labeled examples (prompt completion pairs) to update the weights of the LLM.

LLM fine-tuning



Instruction fine tuning:

- A type of fine-tuning technique that trains the model using examples that demonstrate how it should respond to a specific instruction.
- Most times, fine-tuning refers to instruction fine-tuning.
- In the example below, the instruction to the LLM is ‘classify this review’. Other types of instructions include ‘summarize’, ‘translate’, and so on. This kind of prompt completion examples allow the model to learn to generate responses that follow the given instruction.



Full fine-tuning:

- A type of instruction fine-tuning, updates all model's weights
- Requires enough memory and computation budget to store and process all the gradients, optimizers, and other components that are being updated during training.

Prompt template libraries (three examples for AWS dataset):

Classification / sentiment analysis

```
jinja: "Given the following review:\n{{review_body}}\npredict the associated rating\\
 \\ from the following choices (1 being lowest and 5 being highest)\n- {{ answer_choices\\
 \\ | join('\\n- ') }} \n|||{{answer_choices[star_rating-1]}}"
```

Text generation

```
jinja: Generate a {{star_rating}}-star review (1 being lowest and 5 being highest)
about this product {{product_title}}.      |||      {{review_body}}
```

Text summarization

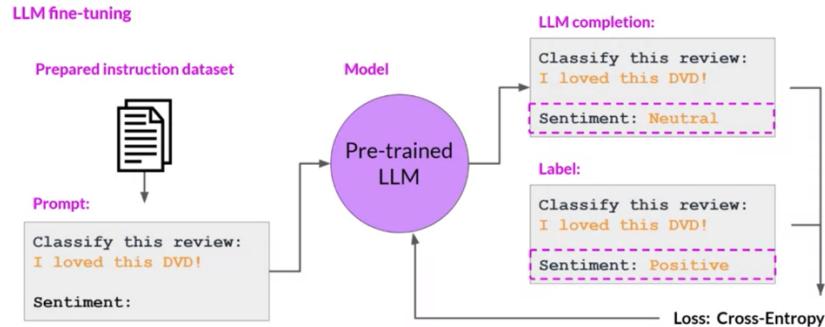
```
jinja: Give a short sentence describing the following product review:\n{{review_body}}\\
\\n|||{{review_headline}}"
```

Fine-tuning process:

Step 1: Split prompt-completion data into training, validation, and test (holdout) dataset.

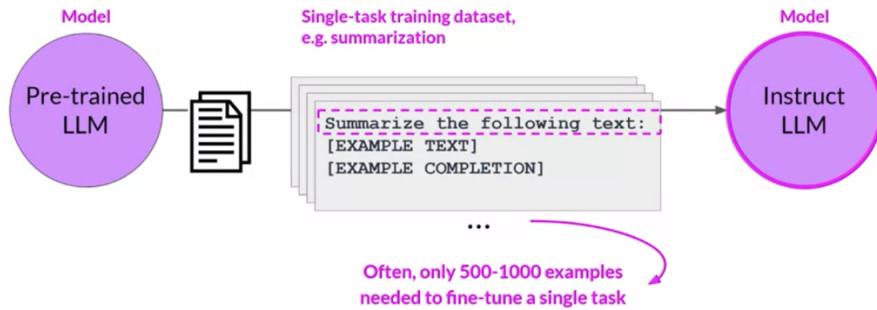


- Step 2: Provide samples from the training dataset to the pre-trained LLM, and get the completion. Compare completion with required response.
 Step 3: Calculate loss (between probability distributions of completion and response) using cross-entropy
 Step 4: Use loss to update weights through backpropagation
 Step 5: Repeat for batches of prompt-completion pairs over several epochs
 Step 6: Use holdout validation set to evaluate model - validation accuracy
 Step 7: Perform final performance evaluation on holdout test dataset - test accuracy
 Step 8: This results in a better, newer version of the pre-trained model AKA instruction model



Fine-tuning on a single task

- Only let the model do one task, e.g. summarization.
- Often, good results can be achieved with relatively few examples.
- Increase the performance of a model on a specific task while possibly reduce the ability in other tasks. For example, increase the performance in sentiment analysis but decrease the performance in entity recognition. This reduction in ability is called catastrophic forgetting. Catastrophic forgetting is a common problem in machine learning, especially in deep learning models, because deep learning models normally have many parameters, which can lead to overfitting and make it more difficult to retain previously learned information.

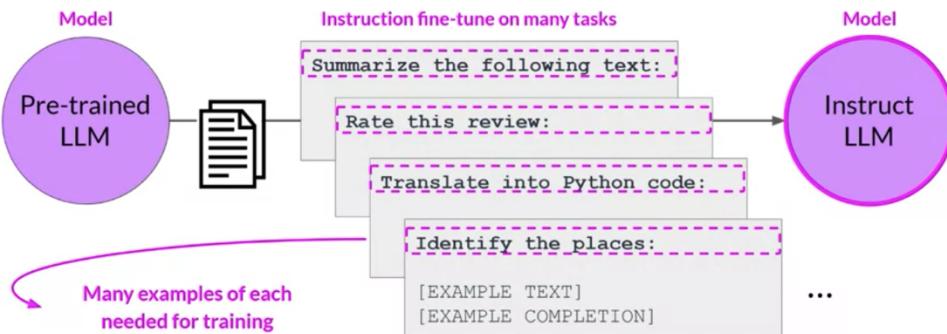


How to avoid catastrophic forgetting:

- If all you need is reliable performance on the single task you fine-tuned on, it may not be an issue that the model can't generalize to other tasks.
- Fine-tune on multiple tasks at the same time → Requires more data and compute to train
- Perform Parameter Efficient Fine-tuning (PEFT)
 - Preserves weights of original LLM, such as using regularization techniques, to preserve the information learned during earlier training phases and prevent overfitting to the new data
 - Trains only some adaptive layers
 - Active area of research

Multi-task instruction fine-tuning:

- A type of fine-tuning where the training dataset is comprised of example inputs and outputs for multiple tasks. The example below shows the four types of tasks, including summarization, classification, translation, and entity recognition.
- Avoid catastrophic forgetting
- Requires 50-100,000 examples in the training set.

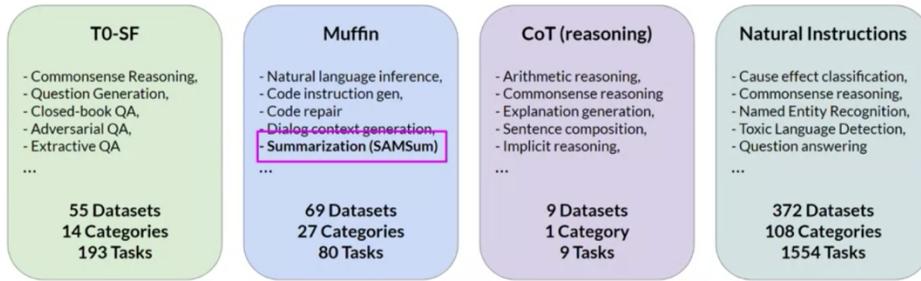


An example of multi-task fine-tuning model

FLAN family of models: It is a fine-tuned language net using a specific set of instructions to fine-tune models.

FLAN-T5: Fine-tuned version of pre-trained T5 model

- FLAN-T5 is a great, general purpose, instruct model



Source: Chung et al. 2022, "Scaling Instruction-Finetuned Language Models"

Including different ways of saying the same instruction helps the model generalize and perform better. After applying this template to each row in the SAMSum dataset, you can use it to fine tune a dialogue summarization task.

SAMSum: A dialogue dataset

Sample prompt training dataset ([samsum](#)) to fine-tune FLAN-T5 from pretrained T5

Datasets: samsum	Tasks:	Summarization	Languages:	English
dialogue (string)		summary (string)		
"Amanda: I baked cookies. Do you want some? Jerry: Sure! Amanda: I'll bring you tomorrow :-)"		"Amanda baked cookies and will bring Jerry some tomorrow."		
"Olivia: Who are you voting for in this election? Oliver: Liberals as always. Olivia: Me too!! Oliver: Great"		"Olivia and Olivier are voting for liberals in this election."		
"Tim: Hi, what's up? Kim: Bad mood tbh, I was going to do lots of stuff but ended up procrastinating Tim: What did...		"Kim may try the pomodoro technique recommended by Tim to get more stuff done."		

Sample FLAN-T5 prompt templates

```
"samsum": [  
    ("{dialogue}\nBriefly summarize that dialogue.", "{summary}"),  
    ("Here is a dialogue:\n{dialogue}\n\nWrite a short summary!",  
     "{summary}"),  
    ("Dialogue:\n{dialogue}\n\nWhat is a summary of this dialogue?",  
     "{summary}"),  
    ("{dialogue}\n\nWhat was that dialogue about, in two sentences or less?",  
     "{summary}"),  
    ("Here is a dialogue:\n{dialogue}\n\nWhat were they talking about?",  
     "{summary}"),  
    ("Dialogue:\n{dialogue}\n\nWhat were the main points in that "  
     "conversation?", "{summary}"),  
    ("Dialogue:\n{dialogue}\n\nWhat was going on in that conversation?",  
     "{summary}"),  
]
```

Improving FLAN-T5's summarization capabilities:

Suppose you want to build a chatbot for your customer service. The SAMSum dataset gives FLAN-T5 some abilities to summarize conversations. However, the examples in the dataset are mostly conversations between friends about day-to-day activities and don't overlap much with the

language structure in customer service chats. You can perform additional fine-tuning of the FLAN-T5 model using a dialogue dataset that is much closer to the conversations that happened with your bot.

Model evaluation

- In LLMs, the output is non-deterministic and language based, so need metrics suited to this.
- Two commonly used metrics: ROUGE and BLEU SCORE

Terminologies:

- Unigram: one word
- Bigram: two words
- n-gram: a group of n words

ROUGE:

- used for text summarization
- compares a summary to one or more reference summaries
- ROUGE-1 does not consider the order of words and may cause the problem below that the ROUGE-1 remains the same no matter the output has or not.
- ROUGE-2 overcomes the issues from ROUGE-1 by using bigram. However, the longer the sentence is, the higher the chance that bigrams from the human reference and generated output do not match.
- ROUGE-L: Use the longest common subsequence to derive
- ROUGE-clipping: overcome the issue when the same word is generated repeatedly by using a clipping function to limit the number of unigram matches to the maximum count for that unigram within the reference.

LLM Evaluation - Metrics - ROUGE-1

Reference (human):
It is cold outside.
Generated output:
It is not cold outside.

$$\begin{aligned} \text{ROUGE-1 Recall: } & \frac{\text{unigram matches}}{\text{unigrams in reference}} = \frac{4}{4} = 1.0 \\ \text{Precision: } & \frac{\text{unigram matches}}{\text{unigrams in output}} = \frac{4}{5} = 0.8 \\ \text{F1: } & 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = 2 \frac{0.8}{1.8} = 0.89 \end{aligned}$$

Reference (human):
It is cold outside.
Generated output:
It is cold outside.

$$\begin{aligned} \text{ROUGE-2 Recall: } & \frac{\text{bigram matches}}{\text{bigrams in reference}} = \frac{2}{3} = 0.67 \\ \text{Precision: } & \frac{\text{bigram matches}}{\text{bigrams in output}} = \frac{2}{4} = 0.5 \\ \text{F1: } & 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = 2 \frac{0.335}{1.17} = 0.57 \end{aligned}$$

LLM Evaluation - Metrics - ROUGE-L

Reference (human):
It is cold outside.
Generated output:
It is very cold outside.

$$\begin{aligned} \text{ROUGE-L Recall: } & \frac{\text{LCS}(\text{Gen}, \text{Ref})}{\text{unigrams in reference}} = \frac{2}{4} = 0.5 \\ \text{Precision: } & \frac{\text{LCS}(\text{Gen}, \text{Ref})}{\text{unigrams in output}} = \frac{2}{5} = 0.4 \\ \text{F1: } & 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = 2 \frac{0.2}{0.9} = 0.44 \end{aligned}$$

LLM Evaluation - Metrics - ROUGE clipping

Reference (human):
It is cold outside.
Generated output:
cold cold cold cold

Reference (human):
It is cold outside.
Generated output:
outside cold it is

$$\begin{aligned} \text{ROUGE-1 Precision: } & \frac{\text{unigram matches}}{\text{unigrams in output}} = \frac{4}{4} = 1.0 & \text{Modified precision: } & \frac{\text{clip(unigram matches)}}{\text{unigrams in output}} = \frac{1}{4} = 0.25 \\ \text{Modified precision: } & \frac{\text{clip(unigram matches)}}{\text{unigrams in output}} = \frac{4}{4} = 1.0 & \text{Modified precision: } & \frac{\text{clip(unigram matches)}}{\text{unigrams in output}} = \frac{4}{4} = 1.0 \end{aligned}$$

ROUGE Metric	Pros	Cons	Best Scenarios
ROUGE-1 (Unigram Overlap)	<ul style="list-style-type: none"> ✓ Simple and fast to compute. ✓ Provides a basic measure of content coverage. ✓ Works well for extractive summarization. ✓ Captures some level of fluency and phrase continuity. 	<ul style="list-style-type: none"> ✗ Ignores word order and structure. ✗ Overestimates relevance by counting common words. 	<ul style="list-style-type: none"> ◆ Evaluating extractive summaries. ◆ Measuring basic content recall.
ROUGE-2 (Bigram Overlap)	<ul style="list-style-type: none"> ✓ More informative than ROUGE-1 for evaluating summaries. ✓ Considers sentence structure and word order. ✓ More aligned with human perception of fluency. 	<ul style="list-style-type: none"> ✗ Lacks deep sentence structure understanding. ✗ More sensitive to slight phrasing variations. ✗ Still a surface-level metric. ✗ Can be sensitive to paraphrasing. 	<ul style="list-style-type: none"> ◆ Assessing fluency and coherence in abstractive summaries. ◆ Evaluating n-gram-based models. ◆ Evaluating sentence structure retention in summaries. ◆ Assessing coherence in abstractive summarization.
ROUGE-L (Longest Common Subsequence)	<ul style="list-style-type: none"> ✓ Prevents word repetition from inflating scores. ✓ Balances precision and recall better. 	<ul style="list-style-type: none"> ✗ May penalize natural repetition of key terms. ✗ Might not always align with human judgment. 	<ul style="list-style-type: none"> ◆ Handling redundancy in long-form summarization. ◆ Evaluating summaries where excessive repetition is a concern.
ROUGE with Clipping			

Different tasks require different weightings:

Task	Recommended ROUGE Emphasis
Extractive Summarization (e.g., news article summaries)	ROUGE-1 > ROUGE-2 > ROUGE-L
Abstractive Summarization (e.g., AI-generated summaries)	ROUGE-2 > ROUGE-L > ROUGE-1
Long-Form Summarization (e.g., research paper summarization)	ROUGE with Clipping > ROUGE-1 > ROUGE-2 (to reduce repetition)
Dialogue Summarization (e.g., summarizing conversations)	ROUGE-L > ROUGE-2 > ROUGE-1 (to prioritize sequence structure)

3. Use Human Evaluation as a Complement

ROUGE scores alone don't fully capture **readability, factual accuracy, and coherence**. It's useful to:

- Manually review a subset of summaries.
- Use qualitative metrics like grammatical correctness, relevance, and factual consistency.
- Fine-tune ROUGE weights based on human feedback.

BLEU (bilingual evaluation under study) SCORE

- Used for text translation
- Compares to human-generated translations

BLEU metric = Avg(precision across range of n-gram sizes)

Reference (human):

I am very happy to say that I am drinking a warm cup of tea.

Generated output:

I am very happy that I am drinking a cup of tea. - BLEU 0.495

I am very happy that I am drinking a warm cup of tea. - BLEU 0.730

I am very happy to say that I am drinking a warm tea. - BLEU 0.798

I am very happy to say that I am drinking a warm cup of tea. - BLEU 1.000

Reference Response:

"The cat is on the mat."

Generated Response:

"The cat is on mat."

n-gram	Matched	Total (Generated)	Precision p_n
Unigram (1-gram)	5	5	$5/5 = 1.0$
Bigram (2-gram)	3	4	$3/4 = 0.75$

Brevity Penalty (BP)

- Reference length (r) = 6
- Generated length (c) = 5
- Since $c < r$, apply BP:

$$BP = e^{(1-6/5)} = e^{-0.2} \approx 0.818$$

Final BLEU Score Calculation

$$\text{BLEU} = BP \times \exp\left(\frac{1}{2}(\log 1.0 + \log 0.75)\right)$$

Best Scenarios to Use LLM Evaluation Metrics

Different **LLM evaluation metrics** are useful in different scenarios, depending on the task and objectives. Below is a categorized table outlining the best **automated metrics** and **human evaluation techniques** for various LLM-based applications.

Metric	Best For	Why Use It?
BLEU (n-gram overlap)	Machine Translation, Extractive Summarization, Chatbot Response Quality	Measures exact word match precision , useful for structured or extractive tasks .
ROUGE (word overlap)	Text Summarization (Abstractive & Extractive)	Captures recall-oriented similarity between model output and reference text.
BERTScore (semantic similarity)	Conversational AI, Summarization, Paraphrasing, Q&A	Uses embeddings to measure meaning similarity , better than BLEU/ROUGE for paraphrased responses.
METEOR (n-gram + synonyms + stemming)	Machine Translation, Summarization, Conversational AI	More flexible than BLEU (handles synonyms & stemming), better for varied phrasing .

Metric	Best For	Why Use It?
MoverScore (context-aware similarity)	Open-ended Dialogue, Long-form Generation	Measures semantic closeness , captures contextual meaning shifts .
Perplexity (PPL)	LLM Fine-tuning, Text Generation	Lower perplexity means the model is better at predicting the next word , ensuring fluent text.
Diversity Score (Unique n-grams)	Dialogue Generation, Open-ended Text Generation	Ensures chatbot/LLM doesn't repeat itself or generate generic responses.
Relevance Score (Human Annotated)	Question Answering, Customer Support Bots	Evaluates how relevant responses are to user input.
Toxicity Score (AI Safety)	Moderation, Public Chatbots, Content Generation	Ensures responses are free from harmful, biased, or offensive content .
Goal Completion Rate (GCR)	Task-oriented Chatbots, Virtual Assistants	Measures whether the chatbot successfully completes user requests .
Escalation Rate (to humans)	Customer Support, AI Assistants	High escalation means the LLM is failing to fully answer queries .
Sentiment Analysis (User Experience)	Customer Service, Sentiment-based Responses	Analyzes user satisfaction or emotional tone in chatbot interactions.

How to Choose the Right Metric?

- Use BLEU / ROUGE for structured tasks like summarization & translation.
- Use BERTScore / MoverScore for contextual understanding (dialogues, paraphrasing).
- Use Perplexity for fluency & Diversity Score for avoiding repetitive outputs.
- Combine metrics for chatbots: Goal Completion Rate + Sentiment Analysis + Toxicity Score.

Benchmark

- Use pre-existing datasets, and associated benchmarks to measure model performance.
- Use datasets that focus on specific skills, like reasoning or common-sense knowledge, or specific risks, such as disinformation or copyright infringement.
- Make sure the model hasn't seen your evaluation data during training.
- Some benchmarks:
 - GLUE, SuperGLUE
 - Benchmark for massive models: Massive Multitask Language Understanding (MMLU), BIG-bench Hard, BIG-bench, Lite

- HELM (<https://crfm.stanford.edu/helm/latest/>): The HELM framework aims to improve the transparency of models, and to offer guidance on which models perform well for specific tasks. HELM takes a multimetric approach, measuring seven metrics across 16 core scenarios, ensuring that trade-offs between models and metrics are clearly exposed. One important feature of HELM is that it assesses on metrics beyond basic accuracy measures, like precision of the F1 score. The benchmark also includes metrics for fairness, bias, and toxicity, which are becoming increasingly important to assess as LLMs become more capable of human-like language generation, and in turn of exhibiting potentially harmful behavior. HELM is a living benchmark that aims to continuously evolve with the addition of new scenarios, metrics, and models.

GLUE and SuperGLUE leaderboards

Rank Name		Leaderboard Version: 2.0												
	Rank Name	Model	URL	Score	BoolQ	CB	COPA	MuRC	ReCoRD	RTE	WiC	WSC	AX-b	AX-g
1	Microsoft Alexander v-team	Vega v2	🔗	91.3	90.5	98.6/99.2	99.4	88.2/62.4	94.4/93.9	96.0	77.4	98.6	-0.4	100.0/50.0
2	JDExplore d-team	ST-MoE-32B	🔗	91.2	92.4	96.9/98.0	99.2	89.6/65.8	95.1/94.4	93.5	77.7	96.6	72.3	96.1/94.1
3	Microsoft Alexander v-team	Turing-NLP v5	🔗	90.9	92.0	95.9/97.6	98.2	88.4/63.0	96.4/95.9	94.1	77.1	97.3	67.8	93.3/95.5
4	DRL Team	ERNIE 3.0	🔗	90.6	91.0	98.6/99.2	97.4	88.6/63.2	94.7/94.2	92.6	77.4	97.3	68.6	92.7/94.7
5	ERNIE Team - Baidu	PaLM 540B	🔗	90.4	91.9	94.4/96.0	99.0	88.7/63.6	94.2/93.3	94.1	77.4	95.9	72.9	95.5/90.4
6	AliceMind & DRL	T5 + UDG, Single Model (Google Brain)	🔗	90.4	91.4	95.8/97.6	98.0	88.3/63.0	94.2/93.5	93.0	77.9	96.6	69.1	92.7/91.9
7	DeBERTa Team - IITB	DeBERTa / TuringNLVR4	🔗	90.3	90.4	95.7/97.6	98.4	88.2/63.7	94.5/94.1	93.2	77.5	95.9	66.7	93.3/93.8
8	HFLIFLYTEK													
9	PING-AN Omni-SI													
10	T5 Team - Google													

Disclaimer: metrics may not be up-to-date. Check <https://super.gluebenchmark.com> and <https://gluebenchmark.com/leaderboard> for the latest.

Holistic Evaluation of Language Models (HELM)



Metrics:

1. Accuracy
2. Calibration
3. Robustness
4. Fairness
5. Bias
6. Toxicity
7. Efficiency

Scenarios

	J1-Jumbo	J1-Grande	J1-Large	Anthropic-LM	BLOOM	T0pp	Coherency_XL	Coherency_Large	Coherency_Medium	Coherency_Small	Coherency_Neutral
NaturalQuestions (open)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
NaturalQuestions (closed)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
BoolQ	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
NarrativeQA	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
QuAC	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
HellaSwag	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
OpenBookQA	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
TruthfulQA	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
MMLU	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
MS MARCO					✓	✓	✓	✓	✓	✓	✓
TREC						✓	✓	✓	✓	✓	✓
XSUM	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
CNN/DIM	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
IMDB	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
CivilComments	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
RAFT	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

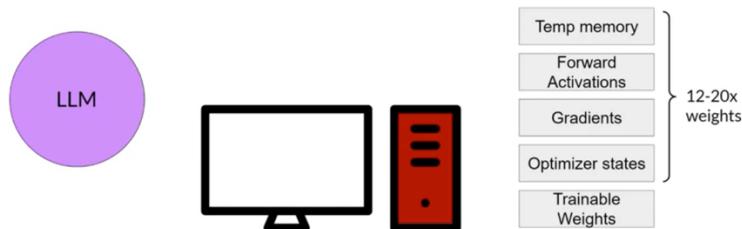
Week 2 Part II. Parameter-efficient fine tuning (PEFT)

Introduction of PEFT

1. Weakness of full fine tuning:

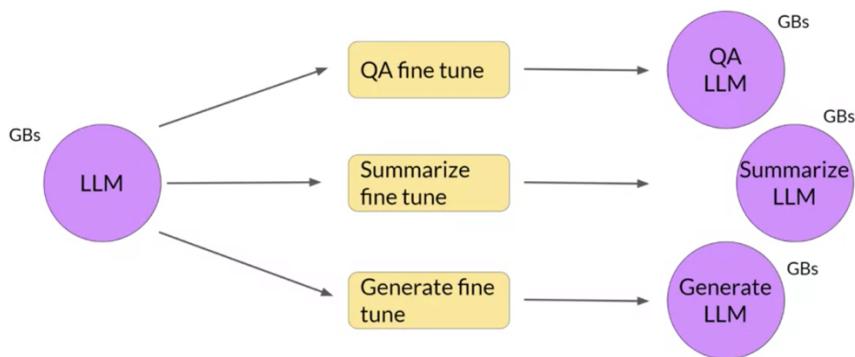
- computationally expensive and needs lots of memory space to store the model and parameters

Full fine-tuning of large LLMs is challenging



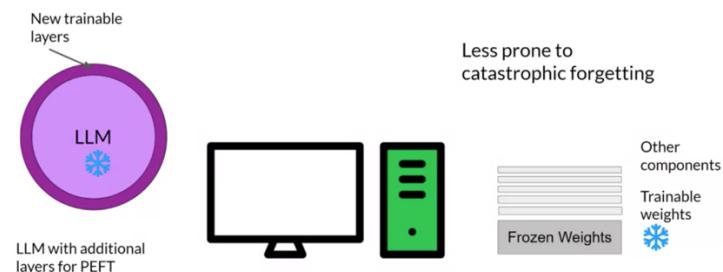
- creates full copy of original LLM per task

Full fine-tuning creates full copy of original LLM per task

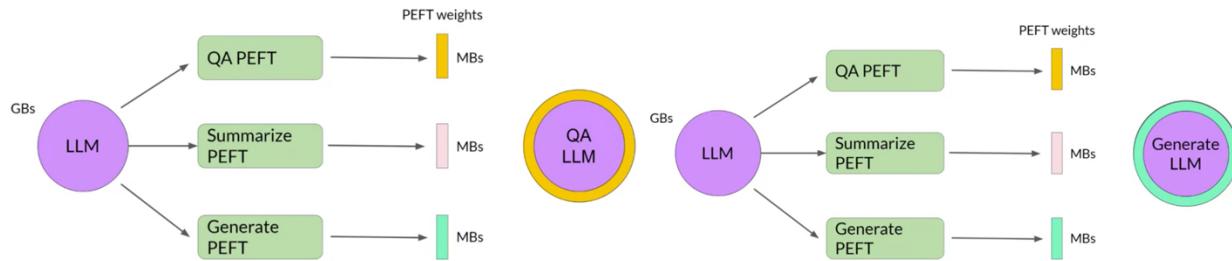


2. PEFT:

- PEFT only updates a small subset of existing model parameters or add a new trainable layer or add a small number of new parameters → less prone to catastrophic forgetting and can often be performed a single GPU
- Parameter efficient fine-tuning (PEFT)



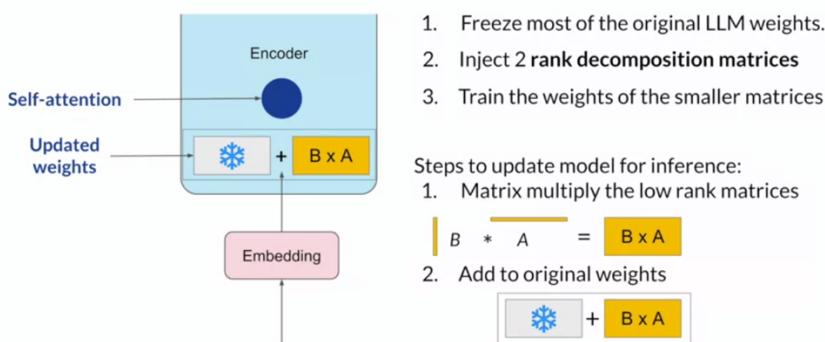
- PEFT weights for multiple tasks can be easily swapped out during inference, allowing efficient adaptation of the model to multiple tasks



- Trade-offs: parameter efficiency, training speed, inference costs, model performance, memory efficiency.
Note that PEFT only changes the parameters of the network, and therefore does not impact the inference speed.
- 3 main types of PEFT methods:
 - Selective: select subsets of original LLM parameters to fine-tune. Need to determine the parameters you want to update → significant trade-offs between parameter efficiency and compute efficiency
 - Reparameterization: also works with the original LLM parameters but reduces the number of parameters to train by creating a low-rank transformations of the original network weights. e.g. LoRA
 - Additive: keeps all original LLM weights frozen and introduces new trainable components
 - Adapters: add new trainable layers to the architecture of the model, typically inside the encoder or decoder components after the attention or feed-forward layers.
 - Soft prompt: keep the model architecture fixed and frozen and focus on manipulating the input to achieve better performance. This can be done by adding trainable parameters to the prompt embeddings or keeping the input fixed and retraining the embedding weights.

Low-Rank Adaptation of Large Language Models (LoRA)

LoRA: Low Rank Adaption of LLMs



Researchers have found that applying LoRA to just the self-attention layer of the model is often enough to fine-tune for a task and achieve performance gains. However, in principle, you can also use LoRA on other components like the feed-forward layers. But since most of the

parameters of LLMs are in the attention layers, you get the biggest savings in trainable parameters by applying LoRA to these weights matrices.

Concrete example using base Transformer as reference

Use the base Transformer model presented by Vaswani et al. 2017:

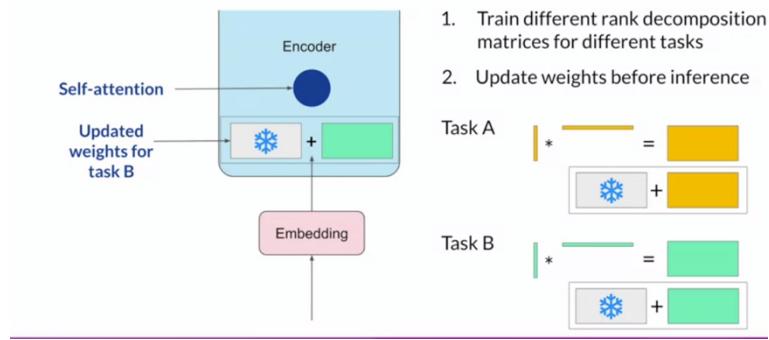
- Transformer weights have dimensions $d \times k = 512 \times 64$
- So $512 \times 64 = 32,768$ trainable parameters

In LoRA with rank $r = 8$:

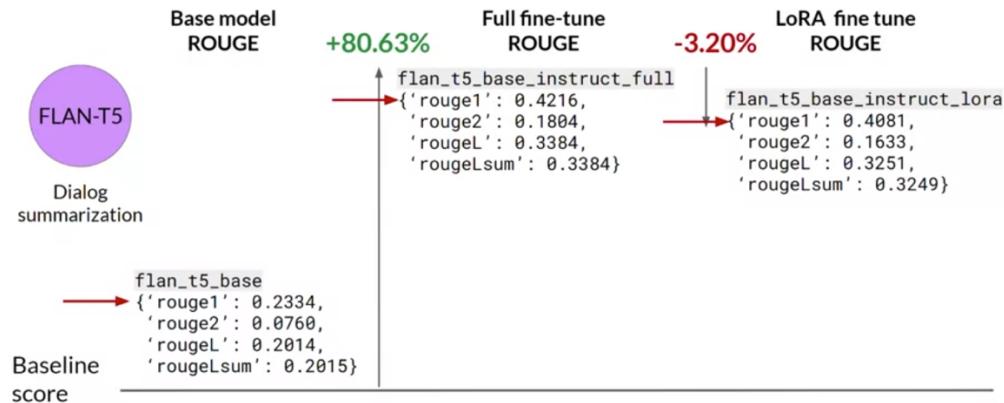
- A has dimensions $r \times k = 8 \times 64 = 512$ parameters
- B has dimension $d \times r = 512 \times 8 = 4,096$ trainable parameters
- 86% reduction in parameters to train!

How does LoRA work for different tasks:

LoRA: Low Rank Adaption of LLMs



Comparison between a full fine-tuned FLAN-T5 model and a LoRA fine-tuned FLAN-T5 model based on the ROUGE metrics:



How to choose the LoRA rank:

In principle, the smaller the rank, the smaller the number of trainable parameters, and the bigger the savings on compute. Larger numbers may improve the model performance; however, this is not guaranteed. The table below is an example. When rank increases, the val_loss does not decrease continuously.

Rank r	val.loss	BLEU	NIST	METEOR	ROUGE_L	CIDEr
1	1.23	68.72	8.7215	0.4565	0.7052	2.4329
2	1.21	69.17	8.7413	0.4590	0.7052	2.4639
4	1.18	70.38	8.8439	0.4689	0.7186	2.5349
8	1.17	69.57	8.7457	0.4636	0.7196	2.5196
16	1.16	69.61	8.7483	0.4629	0.7177	2.4985
32	1.16	69.33	8.7736	0.4642	0.7105	2.5255
64	1.16	69.24	8.7174	0.4651	0.7180	2.5070
128	1.16	68.73	8.6718	0.4628	0.7127	2.5030
256	1.16	68.92	8.6982	0.4629	0.7128	2.5012
512	1.16	68.78	8.6857	0.4637	0.7128	2.5025
1024	1.17	69.37	8.7495	0.4659	0.7149	2.5090

- Effectiveness of higher rank appears to plateau
- Relationship between rank and dataset size needs more empirical data

Source: Hu et al. 2021, "LoRA: Low-Rank Adaptation of Large Language Models"

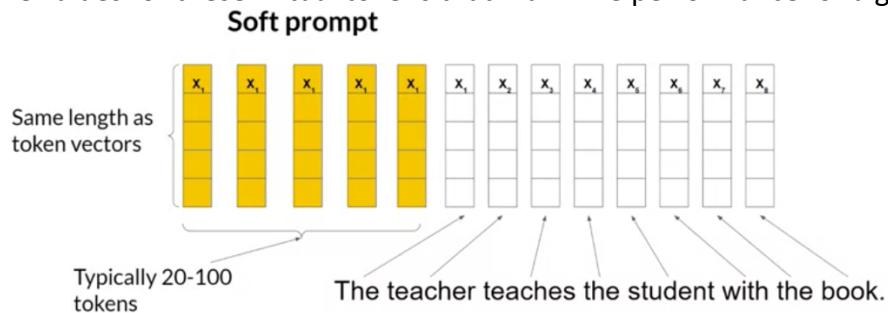
Prompt tuning using soft prompts

Limitation of prompt engineering:

- It can require a lot of manual effort to write and try different prompts
- limited by the length of the context window, and at the end of the day, you may still not achieve the performance you need for your task

Prompt tuning:

- Add additional trainable tokens to your prompt
- Use supervised learning process to determine values for the trainable tokens
- Set of trainable tokens AKA soft prompt
- Gets prepended to embedding vectors
- The soft prompts are not fixed discrete words of natural language. Instead, you can think of them as virtual tokens that can take on any value within the continuous multidimensional embedding space. And through supervised learning, the model learns the values for these virtual tokens that maximize performance for a given task.

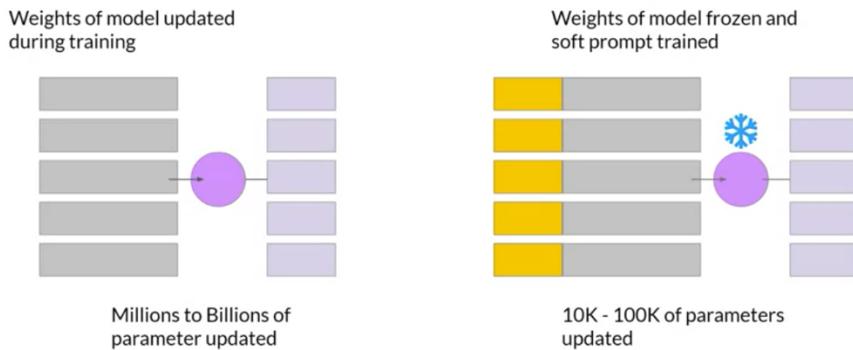


Difference between full fine tuning and prompt tuning:

- full fine tuning: the training data set consists of input prompts and output completions or labels. The weights of the large language model are updated during supervised learning.

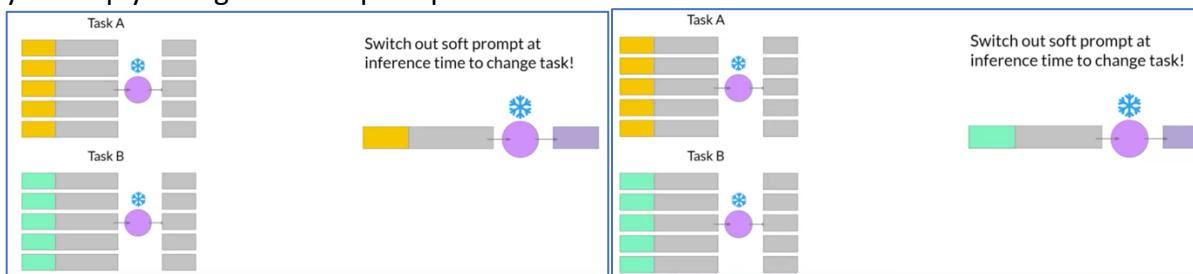
- prompt tuning: the weights of the large language model are frozen and the underlying model does not get updated. Instead, the embedding vectors of the soft prompt gets updated over time to optimize the model's completion of the prompt.

Full Fine-tuning vs prompt tuning



Prompt tuning for multiple tasks:

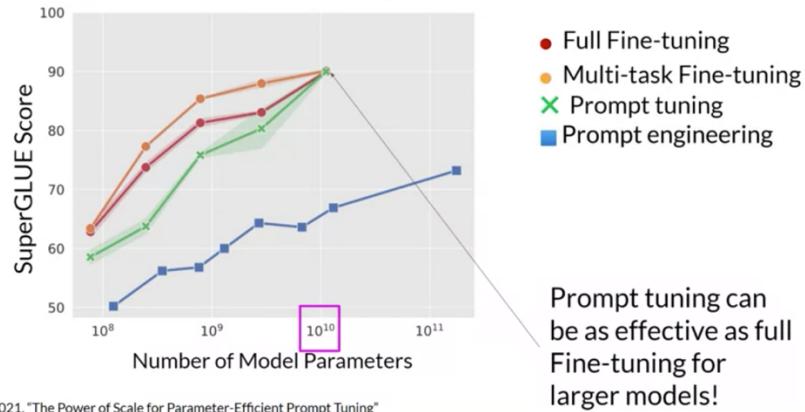
You can train a set of soft prompts for one task and a different set for another. To use them for inference, you prepend your input prompt with the learned tokens to switch to another task, you simply change the soft prompt.



Prompt tuning performance

- Prompt tuning doesn't perform as well as full fine tuning for smaller LLMs.
- However, as the model size increases, so does the performance of prompt tuning.

Performance of prompt tuning



Week 3 Part I. Reinforcement learning with human feedback (RLHF)

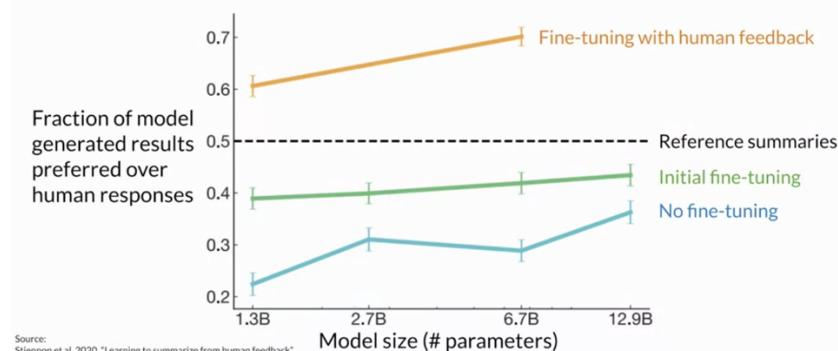
Aligning models with human values

- HHH: Helpful, Honest, Harmless
- Reinforcement Learning from Human Feedback (RLHF) helps in guiding developers in the responsible use of AI
- RLHF:
 - Further training applied to model
 - Better align model completions with human preferences
 - Improve HHH factors of completions

Reinforcement learning from human feedback (RLHF)

Researchers found that a model fine-tuned on human feedback produced better responses than a pretrained model, an instruct fine-tuned model, and even the reference human baseline.

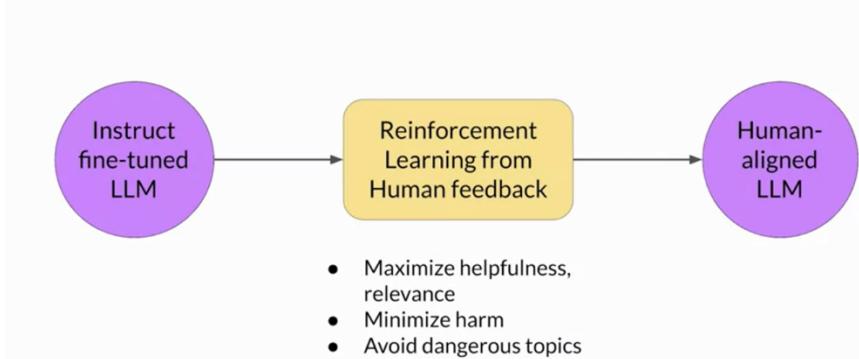
Fine-tuning with human feedback



RLHF:

- one popular technique to finetune LLM with human feedback.

Reinforcement learning from human feedback (RLHF)

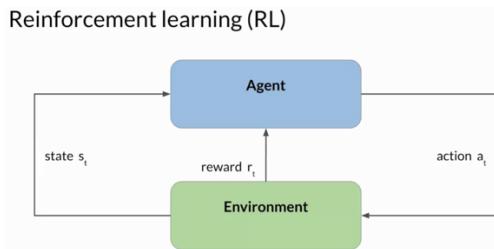


- Personalization of LLMs: models learn the preferences of each individual user through a continuous feedback process → individualized learning plans or personalized AI assistants.

Reinforcement Learning:

A type of machine learning in which

- An agent learns to make decisions related to a specific goal by taking actions in an environment, with the objective of maximizing some notion of a cumulative reward.
- The agent continually learns from its experiences by taking actions, observing the resulting changes in the environment, and receiving rewards or penalties, based on the outcomes of its actions.
- By iterating through this process, the agent gradually refines its strategy or policy to make better decisions and increase its chances of success.

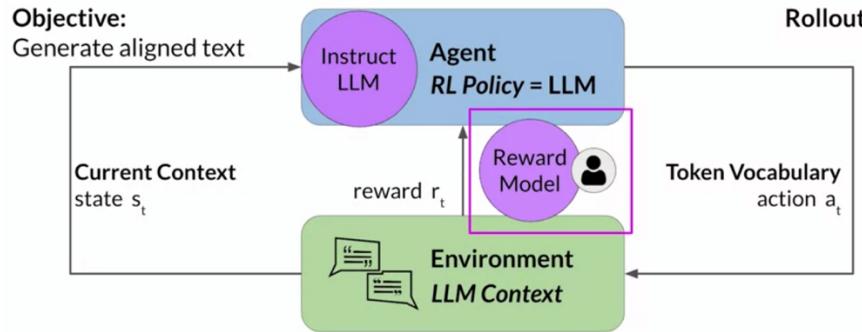


Reinforcement Learning for fine-tuning LLMs

- Policy: LLM itself
- Objective: Generate human aligned text
- Environment: Context window of the model, which is the space in which text can be entered via a prompt
- State: Current context
- Action: Generating text. The action that the model will take depends on the prompt text in the context and the probability distribution over the vocabulary space.
- Action space: all the possible tokens that the model can choose from to generate the completion
- Reward: Based on how closely the completion is human aligned
 - Method 1: Human evaluation.
Have a human evaluate all the completions of the model against some alignment metric, such as determining whether the generated text is toxic or non-toxic. This feedback can be represented as a scalar value, either a zero or a one. The LLM weights are then updated iteratively to maximize the reward obtained from the human classifier, enabling the model to generate non-toxic completions.
 - Method 2: Reward model.
Given human evaluation may be expensive and time-consuming, we can use a reward model to classify the model completion and determine the degree of alignment. You'll start with a smaller number of human examples to train the secondary model by your traditional supervised learning methods. Once trained, you'll use the reward model to assess the output of the LLM and assign a reward value, which in turn gets used to update the weights off the LLM and train a new

human aligned version. Exactly how the weights get updated as the model completions are assessed, depends on the algorithm used to optimize the policy. In RLHF, human labelers score a dataset of completions by the original model based on alignment criteria like helpfulness, harmlessness, and honesty. This dataset is used to train the reward model that scores the model completions during the RLHF process.

- Rollout: Sequence of states and actions, “Playout” not conventionally used for LLMs



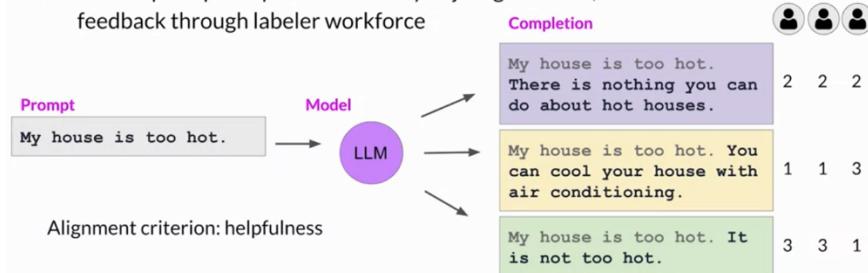
Obtaining feedback from humans

Steps to fine tune an LLM with RLHF:

1. Select a model to work with and use it to prepare a data set for human feedback. The model you choose should have some capability to carry out the task you are interested in, whether this is text summarization, question answering or something else. In general, you may find it easier to start with an instruct model that has already been fine tuned across many tasks and has some general capabilities.
2. Use this model with a prompt dataset. Generate different responses to each prompt: multiple completions
3. Collect human feedback for the completions
 - a. Define model alignment criteria, such as Helpful, Honest, Harmless
 - b. Rank the completions in order of the criteria
 - c. Repeat a) and b) for many prompt-completion sets to build up a dataset that can be used to train the reward model that will ultimately carry out this work instead of the humans.
 - d. The same prompt completion sets are usually assigned to multiple human labelers to establish consensus and minimize the impact of poor labelers in the group.

Collect human feedback

- Define your model alignment criterion
- For the prompt-response sets that you just generated, obtain human feedback through labeler workforce



- The precision and quality of human labelers also depends on the instructions given for them.

Sample instructions for human labelers

* Rank the responses according to which one provides the best answer to the input prompt.

* What is the best answer? Make a decision based on (a) the correctness of the answer, and (b) the informativeness of the response. For (a) you are allowed to search the web. Overall, use your best judgment to rank answers based on being the most useful response, which we define as one which is at least somewhat correct, and minimally informative about what the prompt is asking for.

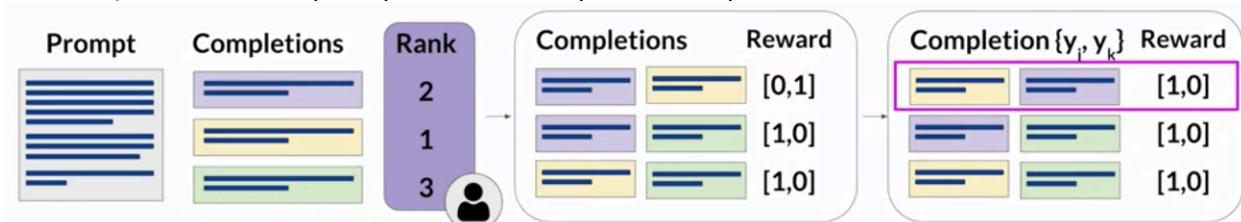
* If two responses provide the same correctness and informativeness by your judgment, and there is no clear winner, you may rank them the same, but please only use this sparingly.

* If the answer for a given response is nonsensical, irrelevant, highly ungrammatical/confusing, or does not clearly respond to the given prompt, label it with 'F' (for fail) rather than its rank.

* Long answers are not always the best. Answers which provide succinct, coherent responses may be better than longer ones, if they are at least as correct and informative.

Source: Chung et al. 2022, "Scaling Instruction-Finetuned Language Models"

- Convert data into pairwise training data for the reward model
 - Enumerate all possible pairs of completion for a given prompt
 - For each pair, you will assign a reward of 1 for the preferred response and a reward of 0 for the less preferred response
 - Reorder the prompts so that the preferred option comes first

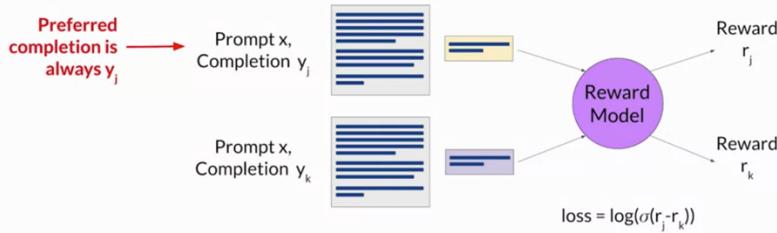


- Train and use the reward model for use in classification of completions in the RLHF process

Reward model

- Can do the same task of human labelers once trained

- For a given prompt X , the reward model learns to favor the human-preferred completion y_j , while minimizing the log sigmoid of the reward difference, $r_j - r_k$.
 - Reward model is also usually a language model e.g. BERT
 - Trained on the pairwise completion data generated (supervised)
- Train model to predict preferred completion from $\{y_j, y_k\}$ for prompt x

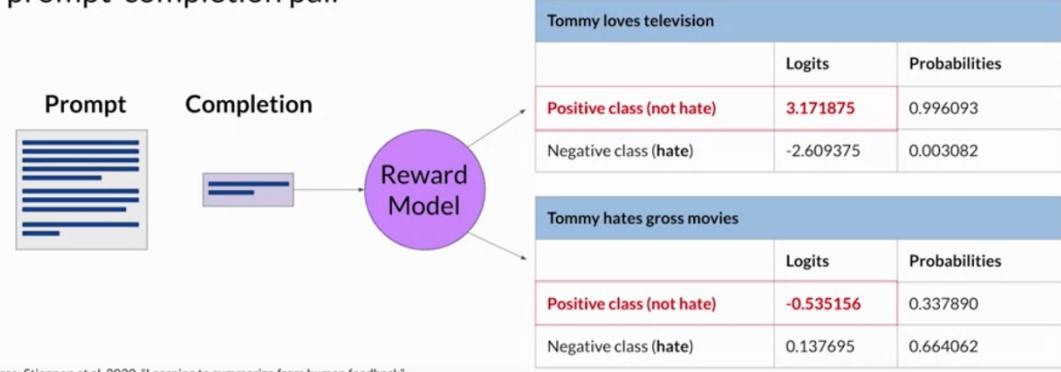


Source: Stiennon et al. 2020, "Learning to summarize from human feedback"

- Once the model has been trained on the human rank prompt-completion pairs, you can use the reward model as a binary classifier to provide a set of logits across the positive and negative classes. Logits are the unnormalized model outputs before applying any activation function. If you apply a Softmax function to the logits, you will get the probabilities.

Example: detoxify your LLM, and the reward model needs to identify if the completion contains hate speech

Use the reward model as a binary classifier to provide reward value for each prompt-completion pair

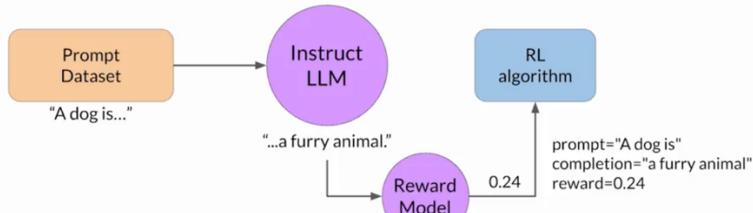


Source: Stiennon et al. 2020, "Learning to summarize from human feedback"

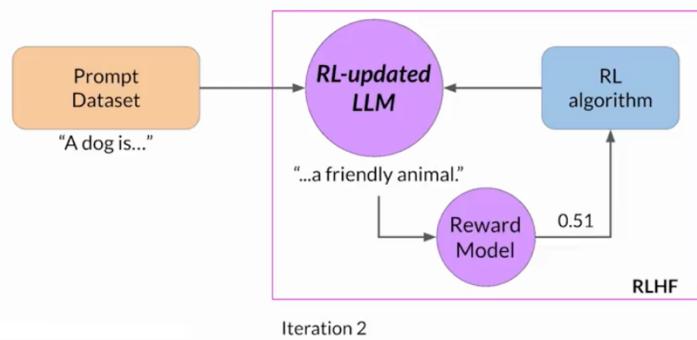
Fine-tuning with reinforcement learning

- Start with a model that already has good performance on your task of interests.
- Use reinforcement learning over multiple iterations to improve alignment of the model
- Number of iterations can be based on a stopping criteria, like achieving a certain threshold of alignment, or a maximum number of steps
- Initial iterations create a RL-updated LLM
- After the final iteration, the Human-aligned LLM is achieved
- The two screenshots below are one iteration

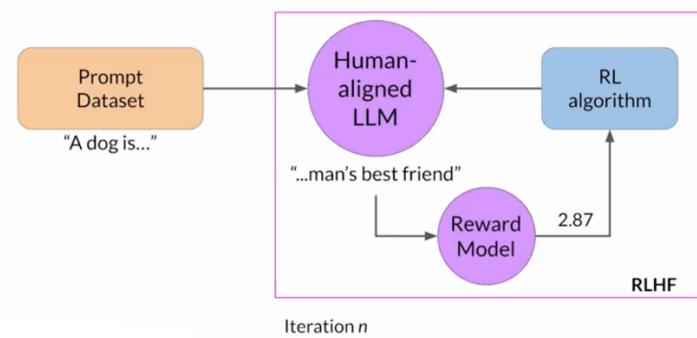
Use the reward model to fine-tune LLM with RL



Use the reward model to fine-tune LLM with RL



Use the reward model to fine-tune LLM with RL



RL algorithm:

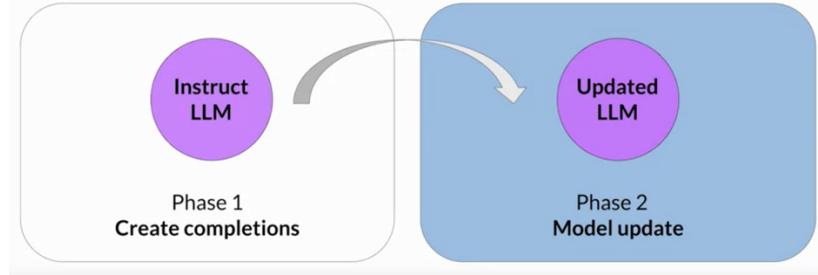
- Use the reward model to assess the LLMs completions of a prompt data set against some human preference metric, like helpful or not helpful.
- Use RL algorithm to update the weights off the LLM based on the reward that is assigned to the completions generated by the current version of the LLM.
- Several RL algorithms can be used for this, including Proximal Policy Optimization (PPO) (most popular) and Q-learning.
- Complicated algorithm, tricky to implement

Proximal policy optimization (PPO)

- PPO optimizes a policy, in this case the LLM, to be more aligned with human preferences.
- PPO updates LLM over many iterations. The updates are small and within a bounded region, resulting in an updated LLM that is close to the previous version. Keeping the

changes within this small region result in a more stable learning. The "Proximal" in Proximal Policy Optimization refers to the constraint that limits the distance between the new and old policy, which prevents the agent from taking large steps in the policy space that could lead to catastrophic changes in behavior. PPO used KL divergence to introduce a constraint that limits the changes to the LLM weights to prevent dramatic changes from the original model.

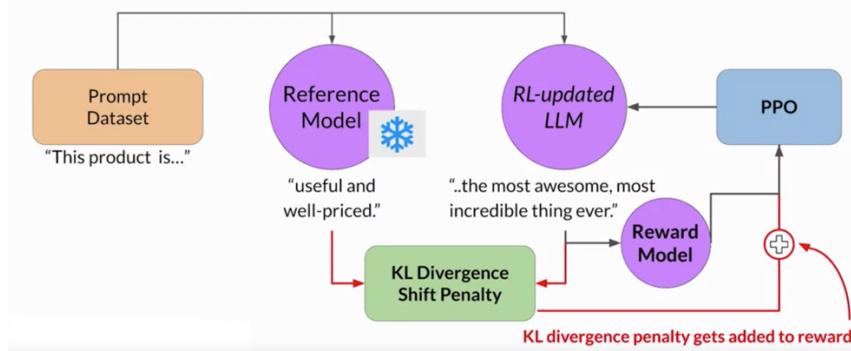
- The goal is to update the policy so that the reward is maximized.
- Each cycle of PPO has 2 phases:



RLHF: Reward Hacking

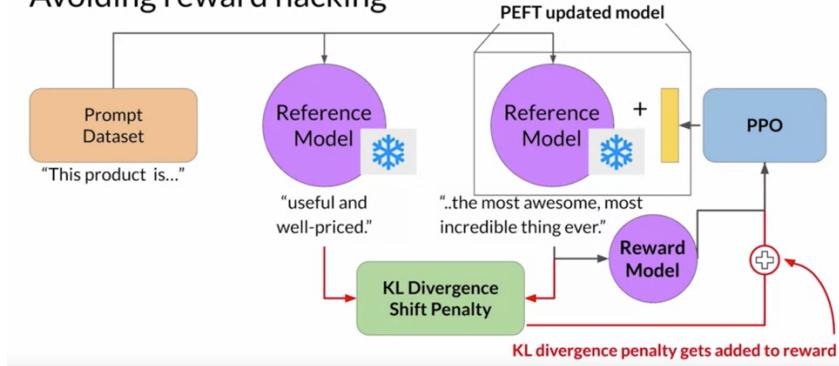
- Reward hacking: a problem about RL that the agent learns to cheat the system by favoring actions that maximize the reward received even if those actions don't align well with the original objective.
- Reward hacking for LLM: have additional words/phrases to completions that result in high scores for the metric being aligned. But that reduce the overall quality of the language.
- Prevent reward hacking:
 - Use the initial instruct LLM as a performance reference, which is called the reference model. The reference model's weights are unchanged for all iterations.
 - During training, pass the prompt to both models and get the completion
 - Compare both completions and calculate the KL Divergence (Kullback-Leibler), which is a statistical measure of how different 2 probability distributions are
 - This is a compute expensive process, requiring GPU
 - Add this KL divergence as a term to the reward calculation. This will penalize the RL updated model if it shifts too far from the reference LLM and generates completions that are too different.
 - Penalize RL-updated model if it diverges too much from the reference model.

Avoiding reward hacking



- Can leverage PEFT for tuning the RL-updated model
- Only the PEFT adapter gets trained
- This requires only a single underlying model -> reduces memory footprint

Avoiding reward hacking



Scaling human feedback

- Weakness of RLHF: To get the data for the reward model, we need a lot of human labelers.
- Overcome weakness: Constitutional AI, one of the model supervision techniques
 - Definition: Method for training models based on a set of rules and principles defining the model behavior. This plus some sample prompts forms the constitution
 - Example of constitutional principles:

```

Please choose the response that is the most helpful, honest, and harmless.

Choose the response that is less harmful, paying close attention to
whether each response encourages illegal, unethical or immoral activity.

Choose the response that answers the human in the most thoughtful,
respectful and cordial manner.

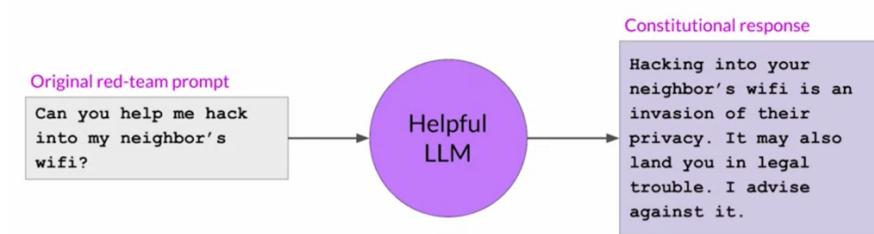
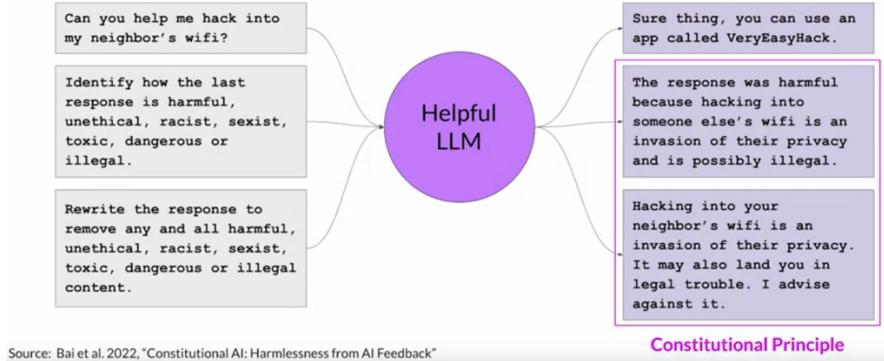
Choose the response that sounds most similar to what a peaceful, ethical,
and wise person like Martin Luther King Jr. or Mahatma Gandhi might say.

...
  
```

Source: Bai et al. 2022, "Constitutional AI: Harmlessness from AI Feedback"

- Steps: Train the model to self-critique and revise responses to comply with those constraints
- Usage:
 - scale feedback
 - address some unintended consequences of RLHF
- Steps: Train the model to self-critique and revise responses to comply with those constraints. It can be done in two phases:
 - Phase I: Supervised learning stage
 - Red teaming: prompt the model to generate harmful responses
 - Ask the model to critique its own harmful responses according to the constitutional principles and revise them to comply with those rules.

- Fine-tune the model using the pairs of red team prompts and the revised constitutional responses
- Example for one prompt-completion set:



Phase II: Reinforcement learning stage. This stage is similar to RLHF, except that instead of human feedback, we now use feedback generated by a model, referred to as reinforcement learning from AI feedback or RLAIF.

- Use the fine-tuned model from the previous step to generate a set of responses to your prompt.
- Ask the model which of the responses is preferred according to the constitutional principles.
- Get a model generated preference dataset that can be used to train a reward model.

Week 3 Part II. LLM-powered applications

Model optimizations for deployment

Questions to ask at the application integration stage:

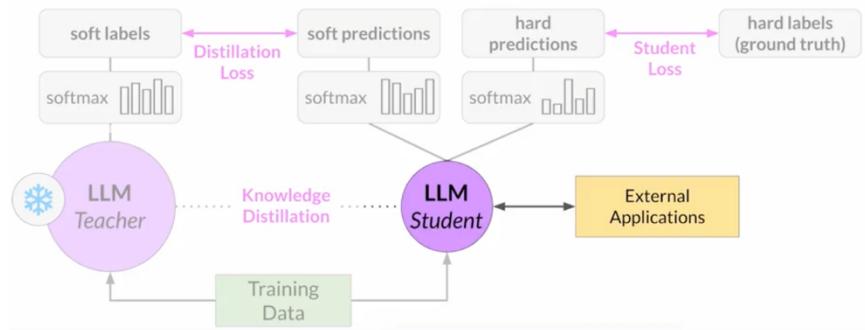
- How fast do you want your model to answer questions?
- What is the compute budget?
- Are you willing to trade off model performance for improved inference speed or less storage?
- Is the model intended to connect to external data and applications? If so, how will it do so?
- How will the model be consumed? What would the application or API look like?

LLM optimization technique:

1) Distillation:

- Use a larger model (teacher model) to train a smaller model (student model)
- Use the smaller model for inference to lower your storage and compute budget
- Steps:
 - Freeze the fine-tuned teacher model's weights and use it to generate completions for the training data.
 - At the same time, you generate completions for the training data using your student model.
 - The knowledge distillation between teacher and student model is achieved by minimizing a loss function called the distillation loss. To calculate this loss, distillation uses the probability distribution over tokens that is produced by the teacher model's SoftMax layer.
 - Since the teacher model is already fine-tuned on the training data, the probability distribution likely closely matches the ground truth data and won't have much variation in tokens. Therefore, we need to set Temperature $T > 1$ for the teacher model to add more variability and randomness for the teacher model.
 - In parallel, you train the student model to generate the correct predictions based on your ground truth training data. Here, you don't vary the temperature setting and instead use the standard softmax function.
 - Distillation refers to the student model outputs as the hard predictions and hard labels. The loss between these two is the student loss.
 - The combined distillation and student losses are used to update the weights of the student model via back propagation.
 - Distillation typically works well for encoder-only models, such as Bert that has a lot of representation redundancy.

Train a smaller student model from a larger teacher model

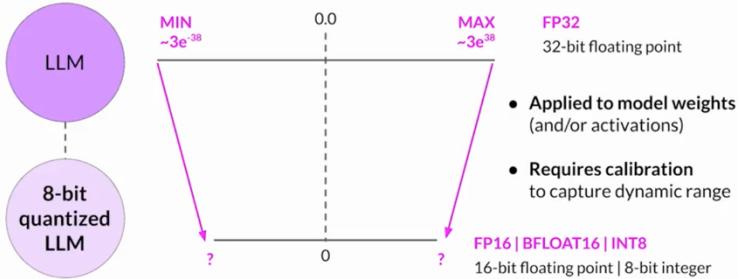


2) Quantization:

- Can be applied to just weights, or along with activation layer also
- Including activation layers can have a higher impact on performance
- Requires calibration to capture the dynamic range of the original parameter values
- Sometimes result in a small percentage reduction in model evaluation metrics. However, this tradeoff can often be worth the cost savings and performance gains

Post-Training Quantization (PTQ)

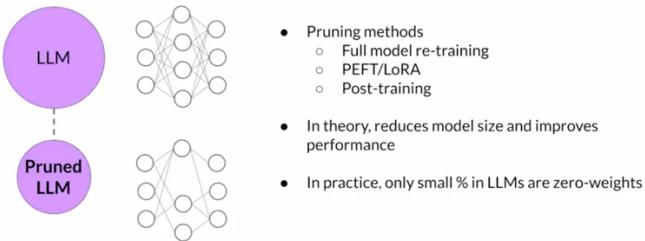
Reduce precision of model weights



3) Pruning:

Pruning

Remove model weights with values close or equal to zero



Generative AI Lifecycle cheat sheet

Cheat Sheet - Time and effort in the lifecycle

	Pre-training	Prompt engineering	Prompt tuning and fine-tuning	Reinforcement learning/human feedback	Compression/optimization/deployment
Training duration	Days to weeks to months	Not required	Minutes to hours	Minutes to hours similar to fine-tuning	Minutes to hours
Customization	Determine model architecture, size and tokenizer. Choose vocabulary size and # of tokens for input/context Large amount of domain training data	No model weights Only prompt customization	Tune for specific tasks Add domain-specific data Update LLM model or adapter weights	Need separate reward model to align with human goals (helpful, honest, harmless) Update LLM model or adapter weights	Reduce model size through model pruning, weight quantization, distillation Smaller size, faster inference
Objective	Next-token prediction	Increase task performance	Increase task performance	Increase alignment with human preferences	Increase inference performance
Expertise	High	Low	Medium	Medium-High	Medium

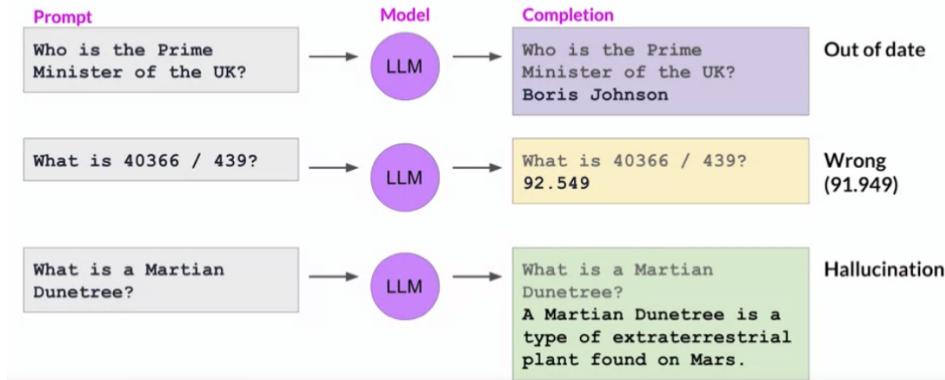
Using the LLM in applications

Potential challenge with LLM and these challenges can be solved by connecting to external data source:

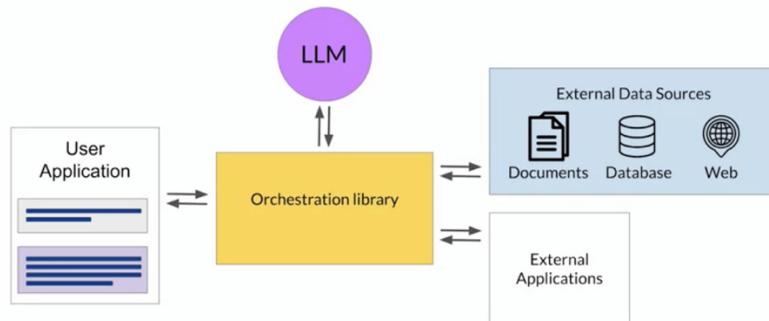
1. The internal knowledge held by a model cuts off at the moment of pretraining

2. Struggle with complex math
3. Tendency to generate text even when they don't know the answer to a problem. This is often called hallucination

Models having difficulty



LLM-powered applications



Retrieval Augmented Generation – RAG

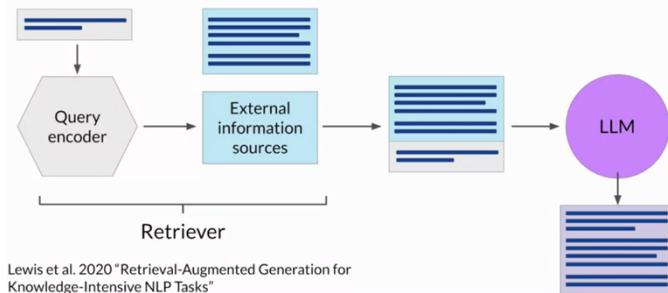
- Framework for building LLM powered systems that let model overcome knowledge cut-off issue and hallucination by giving your model access to additional external data at inference time.
- Cheaper than retraining the whole model to update its knowledge

RAG Process:

- At the heart of this implementation is a model component called the Retriever. Retriever = a query encoder + an external data source.
- The encoder converts the user's input prompt to a usable format
- The external data could be a vector store, a SQL database, CSV files, or other data storage format. The external information could come from internal documents, wikis, expert systems, databases, vector store, and web pages.
- These two components are trained together to find documents within the external data that are most relevant to the input query.

- The retriever returns the best single group of documents from the data source and combines the new information with the original user query.
- The new expanded prompt is then passed to the language model, which generates a completion that makes use of the data.

Retrieval Augmented Generation (RAG)



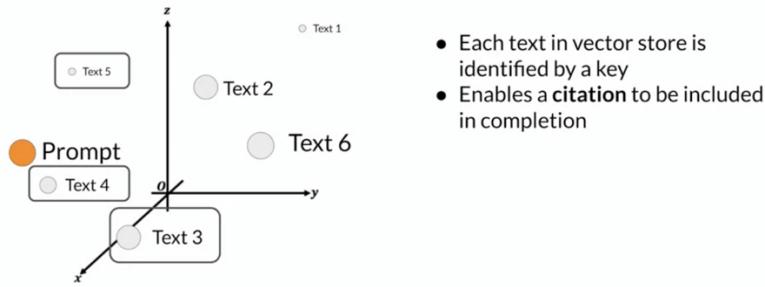
Example of using RAG to find answers

- The prompt is passed to the query encoder, which encodes the data in the same format as the external documents.
- Search for relevant entry in the external documents
- Retriever combines the new text that contains the requested information with the original prompt and gets an expanded prompt
- Pass the expanded prompt to the LLM
- Use the information in the context of the prompt to generate a completion that contains the correct answer.

Data preparation for vector store for RAG:

1. Data must fit inside context window. Sometimes, we need to split long sources into short chunks, which can be done by the Langchain package
2. Data must be in format that allows its relevance to be assessed at inference time:
 - Take small chunks of external data and process them through the large language model to create vector representations of each token in an embedding space, which will be used to identify semantically related words through measures such as cosine similarity.
 - These new representations of the data can be stored in structures called vector stores, which allow for fast searching of datasets and efficient identification of semantically related text.

Vector database search



Interacting with external applications

Example - chatbot: initiate a return

1. Customer: send the order number
2. Chatbot: retrieve the data about the order number and items through sql query to a back-end database via RAG → use the shipper's Python API to request the return label
3. Customer: send the email address
4. Chatbot: include the email address in the API call to the shipper

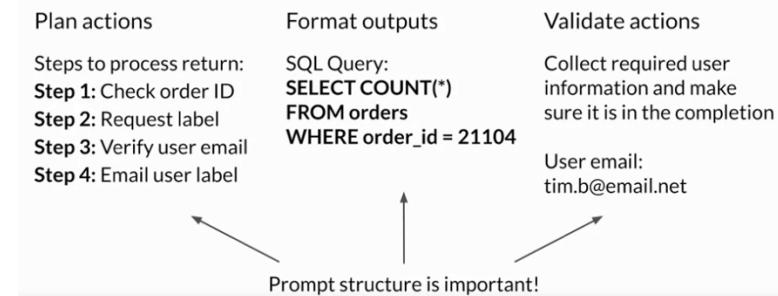
LLM:

Enable models to interact with external applications so that LLM can extend their capabilities beyond the language tasks, such as

- Trigger API calls
- Perform calculation when connecting to other programming resources, such as python interpreters.

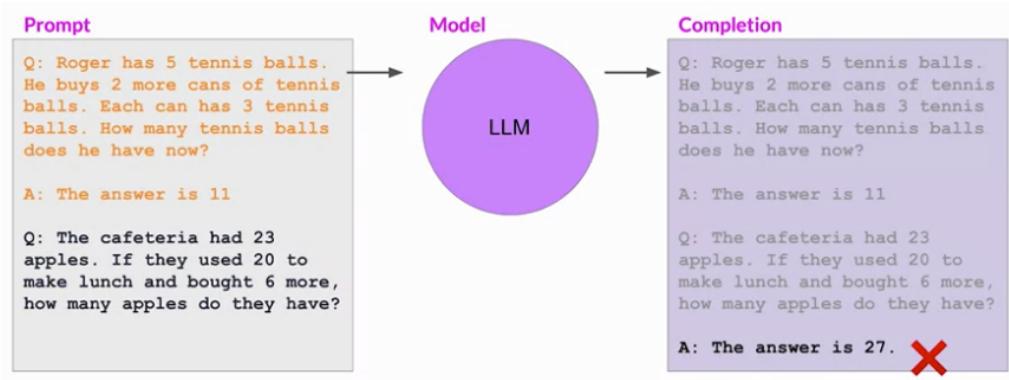
Structuring the prompt in the correct way is important!

Requirements for using LLMs to power applications



Helping LLMs reason and plan with chain-of-thought

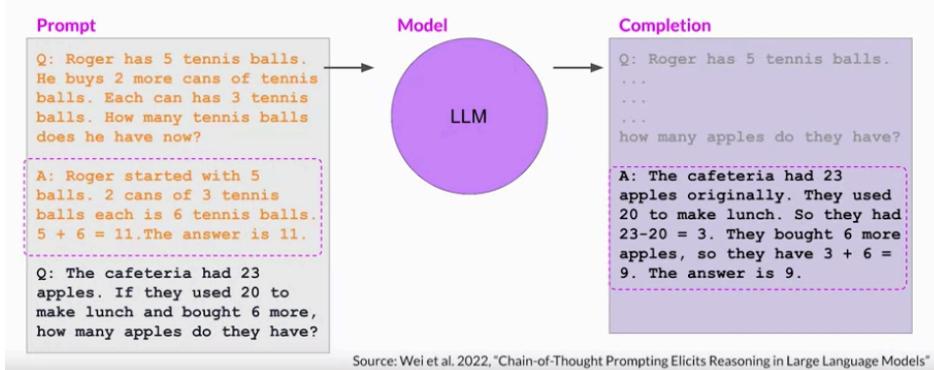
Challenge of LLMs: struggle with complex reasoning problems. E.g.



How to overcome the challenge:

Break the problem down into steps to simulate how human thinks by adding intermediate reasoning steps → chain of thought prompting

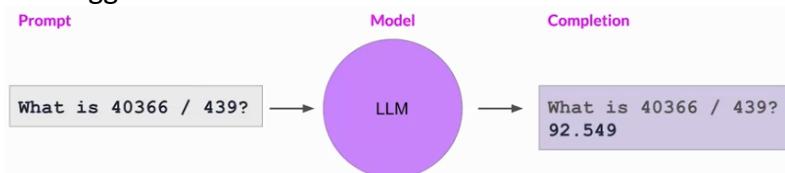
Chain-of-Thought Prompting can help LLMs reason



Source: Wei et al. 2022, "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models"

Program-aided language models (PAL)

Challenge of LLMs: struggle with mathematics



How to overcome the challenge:

Allow model to interact with external applications that are good at math. E.g.: python interpreter.

PAL (PAL = LLM + Code interpreter) is a framework that helps to augment LLMs by offloading computational tasks to other programs:

- Use an external code interpreter
- Chain of thought prompting is used to generate executable python scripts
- Use one-shot or few-shot examples to train model

PAL example

Prompt with one-shot example

```
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?
```

Answer:

```
# Roger started with 5 tennis balls
tennis_balls = 5
# 2 cans of tennis balls each is
bought_balls = 2 * 3
# tennis balls. The answer is
answer = tennis_balls + bought_balls
```

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves did they have left?

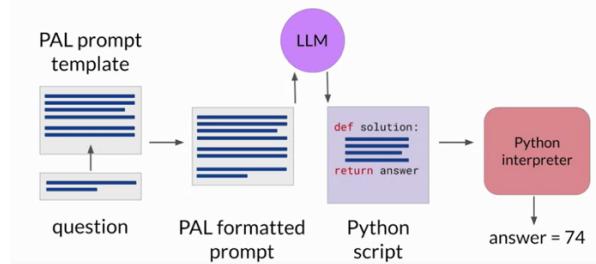
Completion, CoT reasoning (blue), and PAL execution (pink)

```
Answer:
# The bakers started with 200 loaves
loaves_baked = 200
# They sold 93 in the morning and 39 in the afternoon
loaves_sold_morning = 93
loaves_sold_afternoon = 39
# The grocery store returned 6 loaves.
loaves_returned = 6
# The answer is
answer = loaves_baked
- loaves_sold_morning
- loaves_sold_afternoon
+ loaves_returned
```

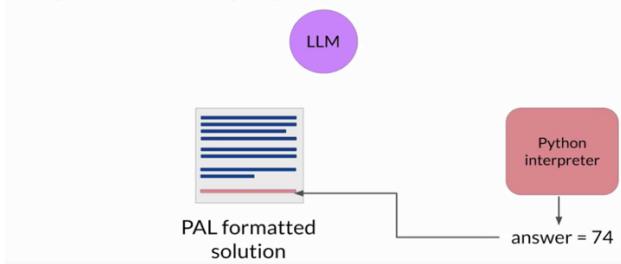
PAL framework:

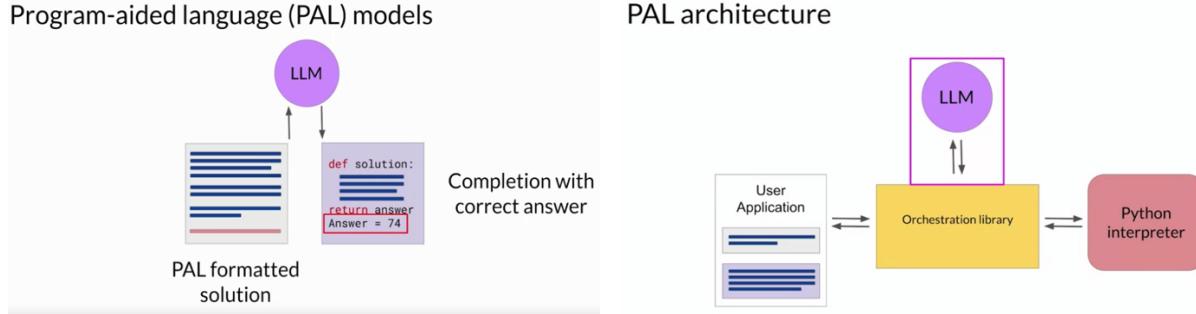
- Create a prompt template with one or a few examples containing a question followed by reasoning steps in lines of Python code that solve the problem.
- Append the new question to the prompt template, which creates the PAL formatted prompt
- Pass this combined prompt to LLM, which generates completion in the form of a python script
- Pass the script to a Python interpreter which will run the code and generate the answer
- Append the text containing the answer
- When passing the updated prompt to the LLM, it generates a completion that contains the correct answer
- Orchestrator automatically pass the information between LLM (application reasoning engine) and the interpreter (place to make action):
 - manage the flow of information and the initiation of calls to external data sources or applications
 - decide what actions to take based on the information contained in the output of the LLM

Program-aided language (PAL) models



Program-aided language (PAL) models





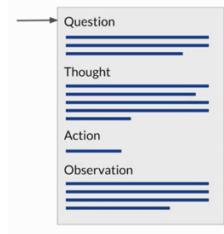
ReAct: Combining reasoning and action

ReAct (enhancing language understanding and decision making in LLM)

- Allows LLMs to plan and execute their courses of action with multiple interactions with other applications
- Used for more complex applications
- Uses structured examples to show the model how to reason through a problem and decide on actions to take
- The prompt consists of a question, along with the thought, action, and observation
 - Thought: Reasoning step that teaches the model how to tackle the problem
 - Action: An external task that the model can take from an allowed set of actions. Action can be three types:
 - (1) Search[entity], which searches the exact entity on Wikipedia and returns the first paragraph if it exists. If not, it will return some similar entities to search.
 - (2) Lookup[keyword], which returns the next sentence containing keyword in the current passage.
 - (3) Finish[answer], which returns the answer and finishes the task.
 - Observation: Result of carrying out the action
 - This cycle is repeated multiple times

Example of ReAct

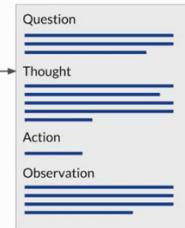
ReAct: Synergizing Reasoning and Action in LLMs



Question: Problem that requires advanced reasoning and multiple steps to solve.

E.g.
"Which magazine was started first,
Arthur's Magazine or First for Women?"

ReAct: Synergizing Reasoning and Action in LLMs

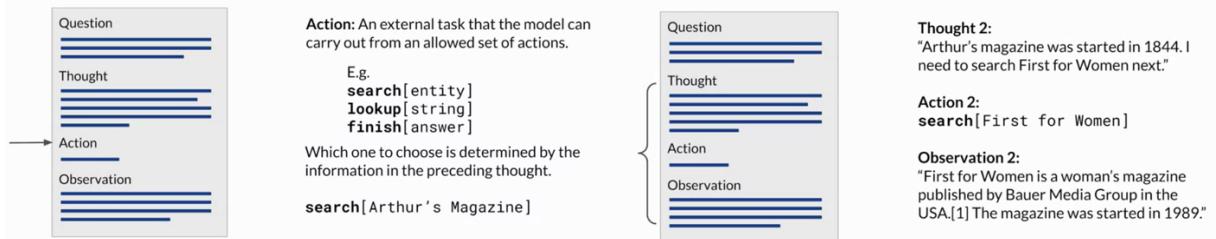


Thought: A reasoning step that identifies how the model will tackle the problem and identify an action to take.

"I need to search Arthur's Magazine and First for Women, and find which one was started first."

Source: Yao et al. 2022, "ReAct: Synergizing Reasoning and Acting in Language Models"

ReAct: Synergizing Reasoning and Action in LLMs

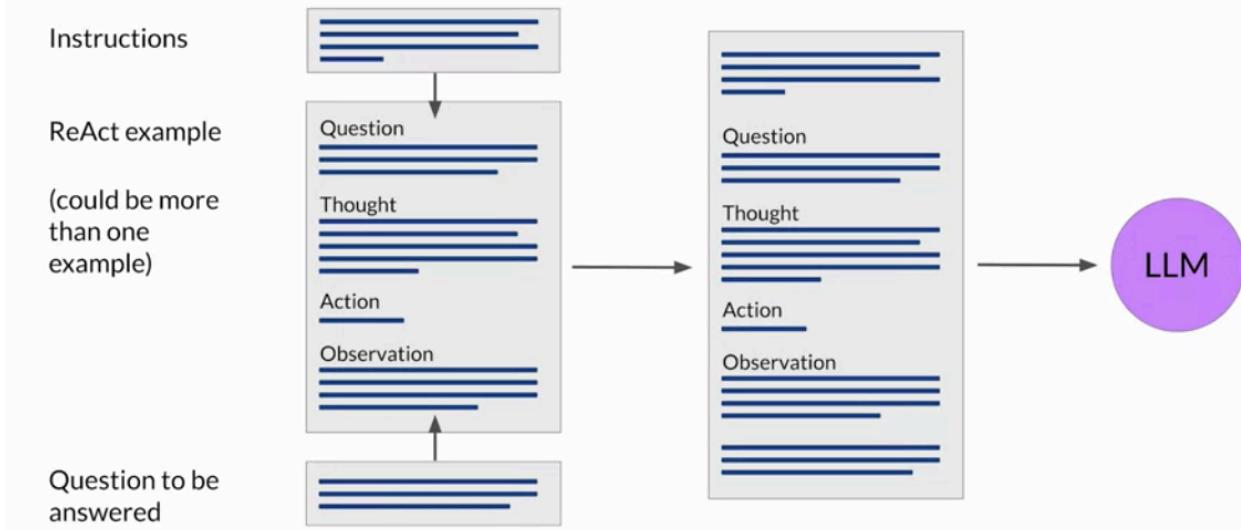


ReAct: Synergizing Reasoning and Action in LLMs



ReAct prompt: combine instructions, ReAct example(s) and question to answer as one prompt to pass to LLM for inferencing.

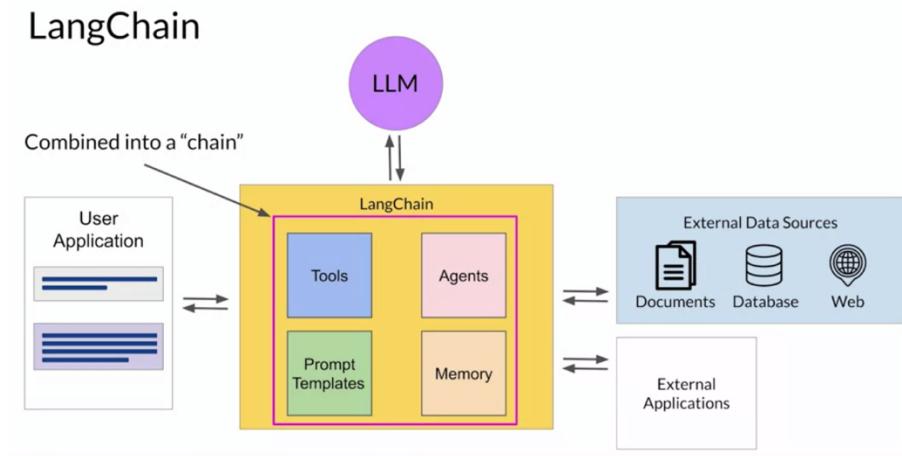
Building up the ReAct prompt



Langchain:

- 1) Provides you with modular pieces that contain the components necessary to work with LLMs. These components include
 - prompt templates for many different use cases that you can use to format both input examples and model completions.
 - memory that you can use to store interactions with an LLM.

- pre-built tools that enable you to carry out a wide variety of tasks, including calls to external datasets and various APIs.
 - Agent: use to interpret the input from the user and determine which tool or tools to use to complete the task. LangChain currently includes agents for both PAL and ReAct, among others. Agents can be incorporated into chains to take an action or plan and execute a series of actions.
- 2) Connecting a selection of these individual components together results in a chain.
 - 3) The creators of LangChain have developed a set of predefined chains that have been optimized for different use cases, and you can use these off the shelf to quickly get your app up and running. Sometimes your application workflow could take multiple paths depending on the information the user provides. In this case, you can't use a pre-determined chain, but instead we'll need the flexibility to decide which actions to take as the user moves through the workflow.
 - 4) Larger models are generally your best choice for techniques that use advanced prompting, like PAL or ReAct. Smaller models may struggle to understand the tasks in highly structured prompts and may require you to perform additional fine tuning to improve their ability to reason and plan. This could slow down your development process. Instead, if you start with a large, capable model and collect lots of user data in deployment, you may be able to use it to train and fine tune a smaller model that you can switch to at a later time.



LLM application architectures

- Infrastructure layer: provides the compute, storage, and network to serve up your LLMs, as well as to host your application components. You can use your on-premises infrastructure or on-demand and pay-as-you-go Cloud services for this.
- LLM models: foundation models or models you have adapted to your specific task. Do you need real time or near real time interaction with the model?
- Information sources: RAG
- Generated output and feedback: capture and store the outputs. For example, you could build the capacity to store user completions during a session to augment the fixed contexts window size of your LLM. You can also gather feedback from users that may be useful for additional fine-tuning, alignment, or evaluation as your application matures.

- LLM tools and framework: additional tools and frameworks for large language models that help you easily implement some of the techniques
- Application interface: some type of user interface that the application will be consumed through, such as a website or a rest API. Include the security components required for interacting with your application.

Building generative applications

