edgerunner AI

Dr. Dylan
Acting Head of Research

Transformers
and LoRA

# What is a Transformer?

Transformers were first described in the seminal research paper by Google:

*Attention Is All You Need*

Initially described for language translation, they have become the basis for the majority of LLMs today (although other architectures exist!)

The key breakthrough was the implementation of the *attention* mechanism

Before we jump into the details, we need to understand *tokenizers*

# Tokenizers

# What are Tokens?

Despite often passing the Turing test, LLMs don't actually know what words are. Instead, they have a complicated mathematical representation of sub-word units called *tokens*.

Tokens can be anything:

- I n d i v i d u a l [space] l e t t e r s [period]
- Entire words.
- Usually words, but sometimes fractions of long words like photos-ynthesis.
- xxbos xxmaj whacky ways of xxcap tokenizing grammar.xxeos

# Tokenizers

Different LLMs have different tokenizers that *must* be used. Switching them leads to nonsensical results.

- GPT4: Edge Runner AI just raised $ 17 . 5 M to get 64 H 100 GPUs ! S woo get y

  [16577, 26032, 20837, 1327, 15478, 548, 1422, 13, 20, 44, 316, 717, 220, 2220, 487, 1353, 193432, 0, 336, 1338, 479, 6494]

- Llama: <|begin_of_text|> Edge Runner AI just raised $ 17 . 5 M to get 64 H 100 GPUs ! S woo get y <|end_of_text|>

  [128000, 11918, 20051, 15592, 1120, 9408, 400, 1114, 13, 20, 44, 311, 636, 220, 1227, 473, 1041, 71503, 0, 328, 49874, 456, 88, 128001]

- DeepSeek: Edge Runner AI just raised $ 17 . 5 M to get 64 H 100 GP Us ! S wo og ety

  [0, 35296, 67737, 7703, 1438, 9927, 957, 1002, 16, 23, 47, 304, 1178, 223, 2892, 437, 1457, 21845, 8095, 3, 327, 1015, 520, 1925]

Although tokens are *usually* similar, the mapping from tokens to numbers differs.

# Embeddings

Depending on the tokenizer, the model can understand between 32,000 and 200k+ different tokens (multilingual models, generally).

This includes "special tokens" like
`<|begin_of_text|><|start_header_id|>system<|end_header_id|>`
and wrappers around thinking, images, and tool calls, depending on the model.

Since computers don't know what a word is, each of these tokens is represented by a (very) high dimensional vector, called an *embedding*.

# Embeddings

All the tokens in a model's vocabulary exist in an embedding matrix.

# Model Sizes - GPT1 and GPT2

| Parameters | Layers | $d_{model}$ | $n_{heads}$ | $n_{neurons}$ |
|---|---|---|---|---|
| GPT1 - 117M | 12 | 768 | 12 | 3072 |
| GPT2 - 117M | 12 | 768 | 12 | 3072 |
| GPT2 - 345M | 24 | 1024 | 16 | 4096 |
| GPT2 - 762M | 36 | 1280 | 20 | 5120 |
| GPT2 - 1542M | 48 | 1600 | 25 | 6400 |

# Model Sizes - GPT3

| Parameters | Layers | $d_{model}$ | $n_{heads}$ | $d_{head}$ |
|---|---|---|---|---|
| 125M | 12 | 768 | 12 | 64 |
| 350M | 24 | 1024 | 16 | 64 |
| 760M | 24 | 1536 | 16 | 96 |
| 1.3B | 24 | 2048 | 24 | 128 |
| 2.7B | 32 | 2560 | 32 | 80 |
| 6.7B | 32 | 4096 | 32 | 128 |
| 13B | 40 | 5140 | 40 | 128 |
| 175B | 96 | 12288 | 96 | 128 |

# Model Sizes - Llama 3

| Parameters | Layers | $d_{model}$ | $n_{heads}$ | $n_{neurons}$* |
|---|---|---|---|---|
| 8B | 32 | 4096 | 32 | 14336 |
| 70B | 80 | 8192 | 64 | 28672 |
| 405B | 126 | 16384 | 128 | 53248 |

* in the GPT models, $n_{neurons}$ = 4 * $d_{model}$ . This is somewhat close to that.

# Embeddings



Male-Female      Verb tense

Directions in the high-dimensional embedding space can carry semantic meaning
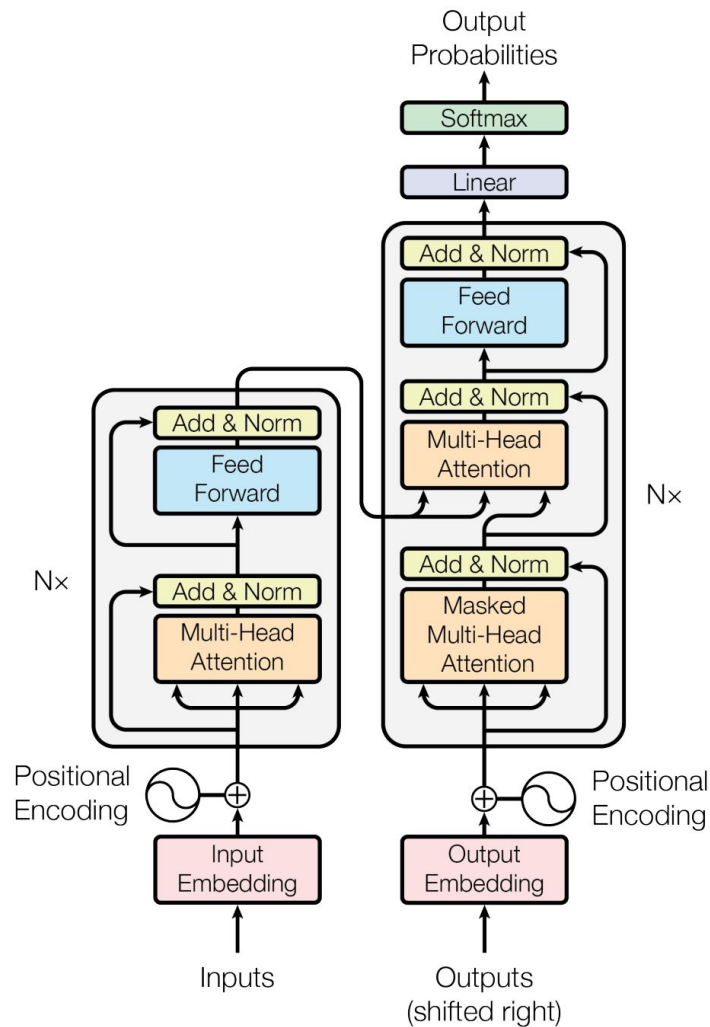
There may be a unique direction representing gender, size, or "capital-ness"
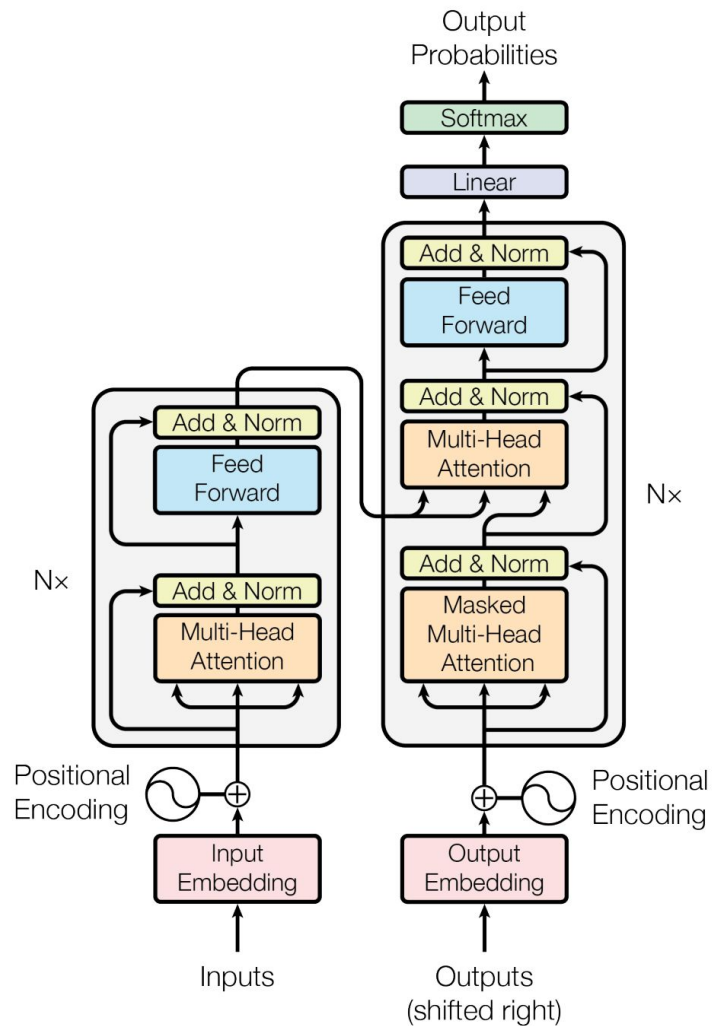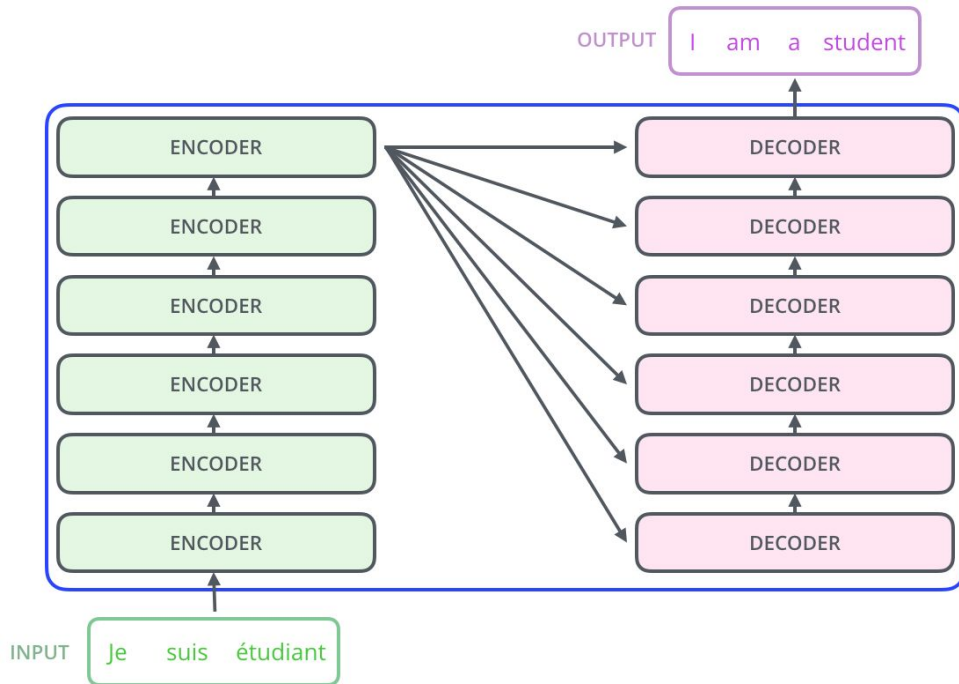
# Transformers

# So then: What is a Transformer?

A Transformer is a deep learning architecture making use of the multi-head *attention* mechanism.

# Transformer Architecture

OUTPUT    I am a student

ENCODER → DECODER

INPUT    Je suis étudiant

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Decoder-Only Architecture (GPTs)



**Output Token Vectors**

Decoder Block
Decoder Block
Decoder Block

Position Embedding

**Input Token Vectors**

FFNN

Layer Norm

Masked Self-Attention

Layer Norm

# Attention

# What is Attention?

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

*Q* = Query

*K* = Key

*V* = Value

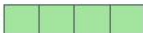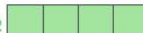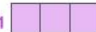$d_k$ = Length of each vector *Q*, *K*, *V*

# Attention: A Breakdown

$QK^T$: This is just the dot product of the Query and Key matrices, measuring how similar they are.

Query: What information do you have?
Key: I have this information!

This is large when the Query and Key values are similar.

The animal didn't cross the street because it was too tired

Layer: 5 ⬍ Attention: Input - Input ⬍

| | |
|---|---|
| The_ | The_ |
| animal_ | animal_ |
| didn_ | didn_ |
| '_ | '_ |
| t_ | t_ |
| cross_ | cross_ |
| the_ | the_ |
| street_ | street_ |
| because_ | because_ |
| it_ | it_ |
| was_ | was_ |
| too_ | too_ |
| tire | tire |
| d_ | d_ |

# Attention

$$\mathrm{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

In practice, this is done many tokens at a time by multiplying matrices instead of vectors.

| | Edge | Runner |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Attention

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

In practice, this is done several tokens at a time by multiplying matrices instead of vectors.

# Multi-Headed Attention

**X**

**Edge Runner**

Calculating attention separately in eight different attention heads

ATTENTION HEAD #0    ATTENTION HEAD #1    ...    ATTENTION HEAD #7

$Z_0$    $Z_1$    $Z_7$

2) Multiply with a weight matrix $W^O$ that was trained jointly with the model

X

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

$W^O$

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

This Z matrix is used to update the embedding in each position

**Z**

=

# Putting it all together

1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer

Edge Runner

**X**

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

**R**

$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

...

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

...

$Q_7$
$K_7$
$V_7$

$Z_0$

$Z_1$

...

$Z_7$

$W^O$

$Z$

# Linear Layers



Layer 1

After going through the Attention block, the outputs are added to the original embeddings, which are normalized and passed through a standard one-hidden-layer feedforward neural network.



These layers take up ⅔ of the parameters in LLMs, and is where the "facts" in the network are stored.

This process is repeated $n_{layers}$ times.

# What's the output?

After a final linear layer, the embeddings are passed through an "unembedding" matrix.

This produces a probability distribution of most likely next tokens.

Which word in our vocabulary is associated with this index?

am

Get the index of the cell with the highest value (argmax)

5

log_probs

0 1 2 3 4 5 … vocab_size

Softmax

logits

0 1 2 3 4 5 … vocab_size

Linear

Decoder stack output

We sample a token from the distribution, append it to the text, and repeat the whole process with the new (slightly longer) input.

# Hallucinations

They're a fundamental part the way the next token is selected:



**Temperature setting**

Softmax output

**Cooler temperature (e.g <1)**

| prob | word |
|------|------|
| 0.001 | apple |
| 0.002 | banana |
| 0.400 | cake |
| 0.012 | donut |
| … | … |

Strongly peaked probability distribution

**Higher temperature (>1)**

| prob | word |
|------|------|
| 0.040 | apple |
| 0.080 | banana |
| 0.150 | cake |
| 0.120 | donut |
| … | … |

Broader, flatter probability distribution

# Further Optimizations

Reducing the size of the attention scores to reduce memory/compute and increase sequence length: sparse attention, ring attention, blockwise attention

Linformers, Structured State Space Sequence models, Mamba/Jamba

**KV-caching** - storing previously computed key and value matrices

**Mixture of Experts (MoE)** - fewer active parameters at a time

**Multi-Token Prediction (MTP)** - what it says on the tin

**Multi-head Latent Attention (MLA)** - low-dimensional representation of attention

**Rotary Position Embeddings (RoPE)** - enables longer sequences

# Training a Transformer

Unsurprisingly, all of the weights in a LLM are stored as matrices.

During backpropagation, we update the individual weights so the outputs better match the actual next word in our corpus.

However, updating (hundreds of) billions of parameters can be expensive and time-consuming - not to mention requiring a shitload of data to do a good job.

GPT-3

Total weights:
175,181,291,520

| | | |
|---|---|---|
| Embedding | d_embed * n_vocab <br> 12,288    50,257 | = 617,558,016 |
| Key | d_query * d_embed * n_heads * n_layers <br> 128    12,288    96    96 | = 14,495,514,624 |
| Query | d_query * d_embed * n_heads * n_layers <br> 128    12,288    96    96 | = 14,495,514,624 |
| Value | d_value * d_embed * n_heads * n_layers <br> 128    12,288    96    96 | = 14,495,514,624 |
| Output | d_embed * d_value * n_heads * n_layers <br> 12,288    128    96    96 | = 14,495,514,624 |
| Up-projection | n_neurons * d_embed * n_layers <br> 49,152    12,288    96 | = 57,982,058,496 |
| Down-projection | d_embed * n_neurons * n_layers <br> 12,288    49,152    96 | = 57,982,058,496 |
| Unembedding | n_vocab * d_embed <br> 50,257    12,288 | = 617,558,016 |

# LoRA

# Low-Rank Adaptation (LoRA)

- Fine-tune models to be better at a specific domain or task

- Can train models with less VRAM, less time, and less data

- Popular for Language and Diffusion models

- Easier to transmit and download

- Can store and swap LoRAs

Because full fine-tunes are so expensive, I've only ever made (and merged) LoRAs

# A little linear algebra

The *rank* of a matrix is the number of columns (or rows) that are linearly independent

$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

What is the rank of this matrix?

If the rank of a matrix is less than min($m,n$) it is *low-rank*

We can decompose a large matrix into two matrices with lower rank.

$$n\begin{array}{c}\overset{m}{\boxed{\phantom{XXX}X\phantom{XXX}}}\end{array} = n\begin{array}{c}\overset{r}{\boxed{\phantom{x}A\phantom{x}}}\end{array}r\begin{array}{c}\overset{m}{\boxed{\phantom{XXX}B\phantom{XXX}}}\end{array}$$
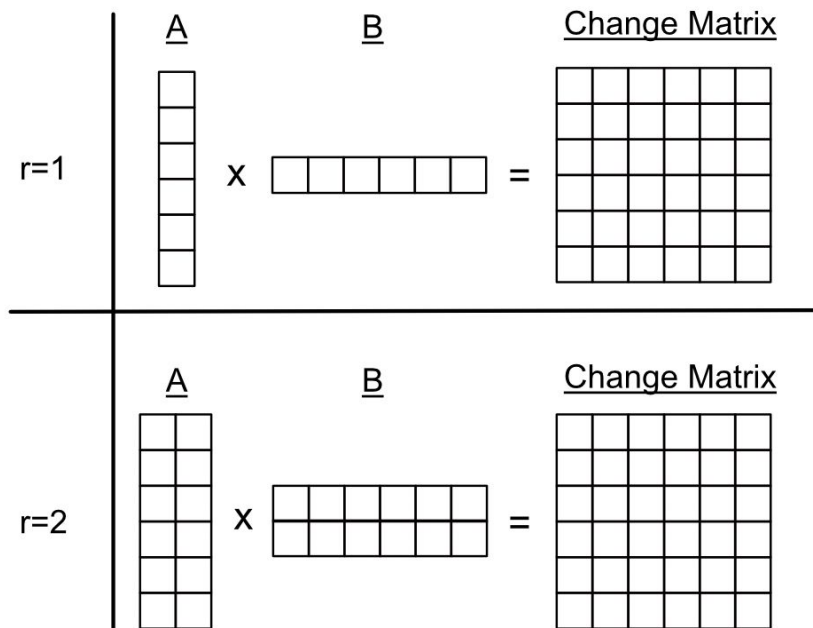
# The idea of Low-Rank Adaptation (LoRA)

Instead of updating all 49,152 × 12,288 ≈ 604 million parameters in one of the linear layers, we can instead update as few as 49,152 × 1 + 1 × 12,288 = 61,440

$m \times n \approx (m \times r) \times (r \times n)$

$r << (m,n)$

$W_0 + \Delta W = W_0 + AB$
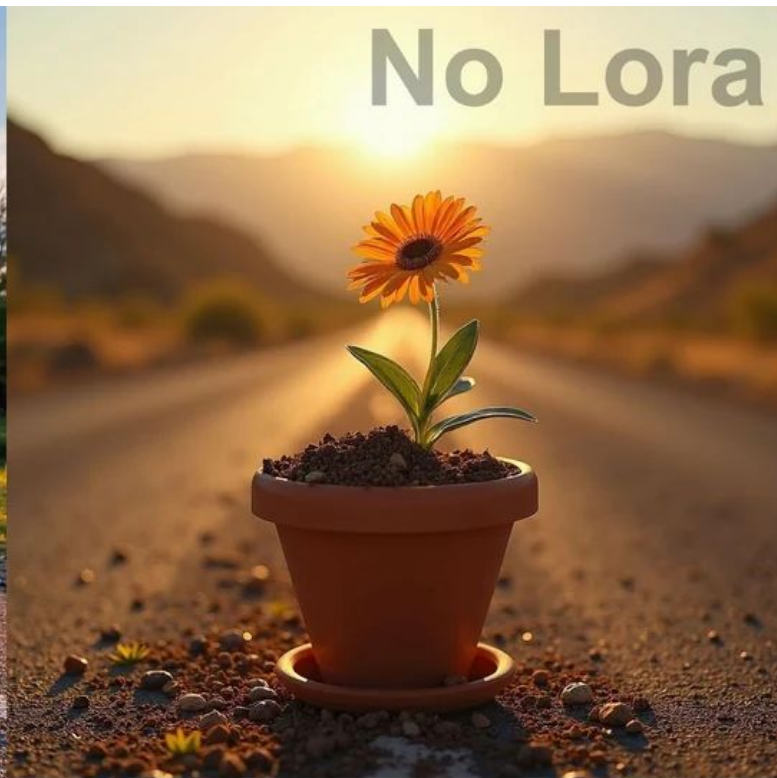
Usually, we set $r$ = 16 up to $r$ = 256.

# Examples of LoRA

# Examples of LoRA



clay pot full of dirt with a beautiful daisie planted in it, shining in the autumn sun on an abandoned, wallpaper, no blur

# Examples of LoRA



Without LoRA

With LoRA

# LLM Examples of LoRA

In the same vein, we can modify the outputs of language models by training them on a particular type of text.

For example, you could create a Shakespeare LoRA, or a Trump LoRA.

In practice, we've been building LoRAs for task-specific or role-specific purposes:

- Acquisitions agent
- Logistics officer
- Scriptwriting assistant

This post-training LoRA requires data with "instruction" "input" and "output"

# Practical LoRA Training - How to do it?

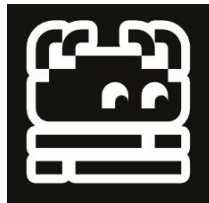There are several frameworks that allow you to train LoRAs:

- HuggingFace PEFT
- Unsloth
- Axolotl

Images only:

- Kohya-trainer
- Dreambooth

You can get away with just a .yaml file and data

# Finetuning with Axolotl

I wrote a [whole document](#) on training with Axolotl.

You need:

- A dataset (local or HuggingFace)
- A `config.yaml` file like this
- At least one GPU

Then type:

`axolotl train my_config.yaml`

```yaml
base_model: meta-llama/Llama-3.3-70B-Instruct
model_type: LlamaForCausalLM
tokenizer_type: AutoTokenizer

load_in_8bit: false
load_in_4bit: true # Note, we load the model in 4bits for training
strict: false

chat_template: llama3

datasets:
  - path: json
    type: alpaca_chat.load_qa
    ds_type: json
    data_files: Acquisitions_Data_Full.jsonl

dataset_prepared_path: last_run_prepared
val_set_size: 0.02
output_dir: ./outputs/70B_v1

sequence_len: 4096
sample_packing: true
eval_sample_packing: false # For multi-GPU setups, this has to be false
pad_to_sequence_len: true

# Same settings as last time
adapter: lora
lora_model_dir:
lora_r: 256
lora_alpha: 128
lora_dropout: 0.05
lora_target_linear: true
```

# Axolotl Optimizations

When we were GPU poor (or when training 70B+ models… ever), you will OOM

- `load_in_8bit`: easy, recommended
- `load_in_4bit`: easy, not recommended unless you have to
- Reduce `sequence_len`: 2048 should be good enough for most models.
- [LoRA Optimizations](#): Can use with 4bit, but not 8bit for some reason.
  - `lora_mlp_kernel: true`
  - `lora_qkv_kernel: true`
  - `lora_o_kernel: true`
- Single GPU only: [unsloth optimizations](#)
  - Requires pytorch < 2.6, which introduces other conflicts with the latest transformers, which is needed for training Mistral. Should be directly incorporated into axolotl soon.
- Decrease `lora_r` (and `lora_alpha` proportionally). I don't think this makes a huge difference.
- Reduce `micro_batch_size`. Increase `gradient_accumulation_steps` proportionally.
- Switch optimizer to `adamw_bnb_8bit`
- Enable `flash_attention` (you should be doing this anyway - not for Mistral)
- Increase DeepSpeed level (need to download new json file)

# Quantized LoRA

By setting `load_in_4bit: True` and `adapter: qlora`, you can quantize a LLM down to 4-bits and freeze the weights, reducing VRAM usage by ~75%.

QLoRA backpropagates gradients through a frozen, 4-bit quantized pretrained language model using LoRA, getting nearly the same results as LoRA or a full fine-tune, but on considerably smaller (cheaper) hardware.

If you have the horsepower, I'd avoid it. But it's very useful when you need it.

# Weight-Decomposed LoRA

DoRA first decomposes the weight matrix W into magnitude/direction components:

$$W = m\frac{V}{||V||_c}$$

It then applies LoRA fine-tuning to the direction only:

$$W' = m\frac{V + \Delta V}{||V + \Delta V||_c} = m\frac{W_0 + \underline{BA}}{||W_0 + \underline{BA}||_c}$$

Although this extra math results in slightly more compute (longer training),

empirical results show the output is similar (or better than) full fine-tunes.

You can use the flag `peft_use_dora: true` to activate DoRA training.

You can combine Quantized LoRA with Weight-Decomposed LoRA ⯈ QDoRA.

# Helpful References

Attention is All You Need, Vaswani et. al. (Original Transformers paper)

Language Models are Unsupervised Multitask Learners, Radford et. al. (GPT2)

Language Models are Few-Shot Learners, Brown et. al. (GPT3)

The Illustrated Transformer, Jay Alammar (Visuals from this presentation)

The Llama 3 Herd of Models, Llama Team

Neural Networks, 3Blue1Brown (specifically videos 5-8)

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

QLoRA: Efficient Finetuning of Quantized LLMs, Dettmers et. al.

DoRA: Weight-Decomposed Low-Rank Adaptation, Liu et. al.