Large
Language
Models

Questions

Flavio Giobergia





For the following decoder-only transformer model, compute:

- 1) the number of parameters for the embedding layer
- 2) The number of parameters of the transformer (excluding the embedding layer, or head layer(s))

Note: assume that none of the layers uses a bias term

- Vocabulary size V=30,000
- Embedding/model dimension d=768
- Positional embeddings: learned, for 512 positions
- Number of layers L=12
- Number of attention heads H=12
- FF-NN projection dimension: 3072
- Head dimension $d_k = d_v = 64$

- •Vocabulary size V=30,000
- •Embedding/model dimension d=768

— [Large Language Models] ————

- •Positional embeddings: learned, for 512 positions
- •Number of layers L=12
- •Number of attention heads H=12
- •FF-NN projection dimension: 3072
- •Head dimension $d_k = d_v = 64$

parameters (Embedding)

- (Token embedding size + positional embedding size)
- 30000 * 768 + 512 * 768 = 23,433,216

parameters (transformer)

12*(attention + FF)

Attention:

Wq, Wv, Wk, Wo = (768*64*12)*3

Wo = 64*12*768

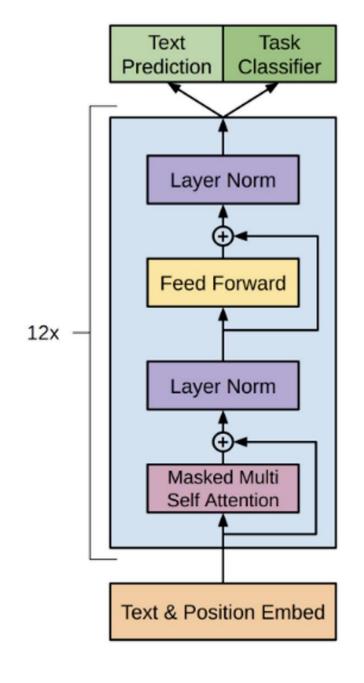
Layer norm(s): (768 * 2) * 2

FF-NN: 1st 768 => 3072, 2nd 3072 => 768

768 * 3072 + 3072 * 768 = (3072 * 768) * 2

= 7,080,960 for one layer => 7,080,960*12

84,971,520



Which of the following is NOT a Parameter-Efficient Fine-Tuning (PEFT) technique?

- Dynamic quantization
- b) Adapter layers
- c) Prompt tuning
- LoRA

Which of the following is NOT a Parameter-Efficient Fine-Tuning (PEFT) technique?

- a) **Dynamic quantization**
- b) Adapter layers
- c) Prompt tuning
- d) LoRA



- The following vectors for the input tokens are given
- Inputs: [2,0], [0,2]
- The following matrices W_q , W_k , W_v are given:

$$W_q = [[0,1,2],[0,0,2]]$$

 $W_k = [[1,2,2],[1,0,1]]$
 $W_v = [[1,2,0],[0,0,1]]$

• Compute Attention(Q,K,V) = $softmax\left(\frac{QK^T}{\sqrt{d_R}}\right)V$, for the first token

- T1 [2, 0]
- T2 [0, 2]

Inputs: [2,0], [0,2]

$$W_q = [[0,1,2], [0,0,2]]$$
 $W_k = [[1,2,2], [1,0,1]]$
 $W_v = [[1,2,0], [0,0,1]]$

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

- T1 [2, 0]
- T2 [0, 2]

$$QK^T$$

2 2

40

42

0 2 4 [24, 8] =>
$$\frac{[24, 8]}{\sqrt{3}} = \left[\frac{24}{\sqrt{3}}, \frac{8}{\sqrt{3}}\right]$$

Inputs: [2,0], [0,2]

$$W_q = [[0,1,2], [0,0,2]]$$
 $W_k = [[1,2,2], [1,0,1]]$
 $W_v = [[1,2,0], [0,0,1]]$

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



$$\frac{\mathbf{QK}^{\mathbf{T}}}{\sqrt{\mathbf{d_{\mathbf{K}}}}} = \left[\frac{24}{\sqrt{3}}, \frac{8}{\sqrt{3}}\right] = [\mathbf{a}, \mathbf{b}]$$

$$\frac{24}{\sqrt{3}} \Rightarrow \frac{e^{x}}{\sum e^{j}} = \frac{e^{a}}{e^{a} + e^{b}} \approx 0.999$$

$$\frac{8}{\sqrt{3}} = \frac{e^{b}}{e^{a} + e^{b}} \approx 0.001$$

$$W_q = [[0,1,2],[0,0,2]]$$

 $W_k = [[1,2,2],[1,0,1]]$
 $W_v = [[1,2,0],[0,0,1]]$

$$softmax \left(\frac{\mathbf{Q} \mathbf{K}^{\mathbf{T}}}{\sqrt{\mathbf{d}_{\mathbf{R}}}} \right) V$$

=
$$\left[\frac{e^{a}}{e^{a}+e^{b}}, \frac{e^{b}}{e^{a}+e^{b}}\right] * V \implies \left[2*\frac{e^{a}}{e^{a}+e^{b}}, 4*\frac{e^{a}}{e^{a}+e^{b}}, 0\right] + \left[0, 0, 2*\frac{e^{b}}{e^{a}+e^{b}}\right]$$

= $\left[2*\frac{e^{a}}{e^{a}+e^{b}}, 4*\frac{e^{a}}{e^{a}+e^{b}}, 2*\frac{e^{b}}{e^{a}+e^{b}}\right]$
= $\left[2, 4, 0\right]$

Which one of these models is a decoder-only?

- **BERT** a)
- GPT-2
- T5
- All options are decoder-only models

Which one of these models is a decoder-only?

- **BERT** a)
- GPT-2 b)
- T5
- All options are decoder-only models

In Causal Language Modeling, the word "causal" means that:

- The output is generated considering only the current token
- The output is generated considering only the previous tokens
- The output is generated considering only the previous and current tokens
- Only the decoder is used during inference
- The attention mechanism is adopted
- The surrounding context must also be used for token generation

In Causal Language Modeling, the word "causal" means that:

- a) The output is generated considering only the current token
- b) The output is generated considering only the previous tokens
- c) The output is generated considering only the previous and current tokens
- d) Only the decoder is used during inference
- e) The attention mechanism is adopted
- f) The surrounding context must also be used for token generation



- A model generates the sentence:
 - "the fast brown fox jumped the bored dog"
- Whereas, the target sentence is:
 - "the quick brown fox jumps over the lazy dog"
- Compute the BLEU-2 score for this generation.
- Remember

$$BLEU-n = BP \cdot \exp(\frac{1}{n} \sum_{i} \log(precision_{i}))$$

$$Brevity\ Penalty = BP = \begin{cases} 1 & \text{if } g > r \\ e^{\left(1 - \frac{r}{g}\right)} & \text{if } g \le r \end{cases}$$

- Where $precision_i$ is the precision for the specific i-gram, g and r are the lengths of the generated and reference texts, respectively
- Consider each word to be a token in this case

- "the fast brown fox jumped the bored dog"
- "the quick brown fox jumps over the lazy dog"

1)

G: the, fast, brown, fox, jumped, the, bored, dog

R: the, quick, brown, fox, jumps, over, the, lazy, dog

Precision 1 = 5/8

2)

G: (the, fast), (fast, brown), (brown, fox), (fox, jumped), (jumped, the), (the bored), (bored, dog)

R: (the, quick), (quick brown), (brown, fox), (fox, jumps), (jumps, over), (over, the), (the, lazy), (lazy, dog)

Precision 2 = 1/7

BLEU-2 = BP * $\exp(\frac{1}{2})$ * $(\log(5/8) + \log(1/7))$

Generated (g) = 8 words Reference (r) = 9 words

 $BLEU-n = BP \cdot \exp(\frac{1}{n}\sum_{i}\log(precision_{i}))$ $BP = \begin{cases} 1 & \text{if } g > r \\ e^{\left(1 - \frac{r}{g}\right)} & \text{if } g \le r \end{cases}$ $RP = e^{\left(1 - \frac{r}{g}\right)} = 0.8825$

$$BLEU - 2 = 0.8825 * \exp(\frac{1}{2} * \log(5/8) + \frac{1}{2} * \log(1/7))) = 0.2637$$

What is the primary difference between static and dynamic quantization?

- Dynamic quantization computes scaling values for each activation during inference
- Dynamic quantization applies to weights only, while static applies to weights and activations
- Static quantization is slower but less accurate than dynamic quantization
- Dynamic Quantization requires a calibration step before the quantization
- None of the others

What is the primary difference between static and dynamic quantization?

- Dynamic quantization computes scaling values for each activation during inference
- Dynamic quantization applies to weights only, while static applies to weights and activations
- Static quantization is slower but less accurate than dynamic quantization
- Dynamic Quantization requires a calibration step before the quantization
- None of the others

- You are given the following corpus of text: abababab
- Apply BPE to extract 3 tokens (excluding a and b). Write the tokens and their extraction frequency (i.e., the frequency computed at the step when they were merged).
- Note: when doing substitution, work from left to right

Abababab

|a|b|a|b|a|b|a|b|

- 1) ab: 4 ba: 3
- (extract ab:4 → X)

Abababab

|a|b|a|b|a|b|a|b|

- 1) ab: 4 ba: 3
- (extract ab:4 → X)
- |ab|ab|ab|ab|
- abab: 3
- (extract abab:3 → Y)

|a|b|a|b|a|b|a|b|

- 1) ab: 4 ba: 3
- (extract ab:4 → X)
- |ab|ab|ab|ab|
- abab: 3
- (extract abab:3 → Y)

Abababab

```
|a|b|a|b|a|b|a|b|
```

- 1) ab: 4 ba: 3
- (extract ab:4 → X)
- |ab|ab|ab|ab|
- abab: 3
- (extract abab:3 → Y)
- |abab|abab|
- abababab: 1
- Extract abababab: 1 -> Z

In LoRA, what does the process assume about the weight update matrix?

- That it is low rank or is near low rank
- That it replaces the initial weight matrix
- That it is initialized with numbers taken from a normal distribution $\mathcal{N}(0,1)$
- That it is fixed (or frozen) throughout training

In LoRA, what does the process assume about the weight update matrix?

- That it is low rank or is near low rank
- That it replaces the initial weight matrix
- That it is initialized with numbers taken from a normal distribution $\mathcal{N}(0,1)$
- That it is fixed (or frozen) throughout training

Which one of these is a deterministic sampling approach? Select all correct answers

- Temperature sampling with T=0 a)
- Temperature sampling, with 0 < T < 1 b)
- Temperature sampling, with T > 1c)
- d) Beam Search
- **Greedy sampling**
- Top-k sampling
- Top-p sampling g)

Which one of these is a deterministic sampling approach? Select all correct answers

- Temperature sampling with T=0 a)
- Temperature sampling, with 0 < T < 1b)
- Temperature sampling, with T > 1c)
- **Beam Search** d)
- **Greedy sampling** e)
- Top-k sampling
- Top-p sampling g)