

简单大模型推理系统实现

由于考研复试在即，本人时间过于仓促，故完成了基本功能，达到通过要求：

- 文本生成
- AI对话

该项目只需cargo run就可直接运行，下面讲述这两个功能的具体实现方法

文本生成

完成模型的推理功能就可以实现文本生成，推理过程被封装在 `forward` 前向传播函数中，该函数在初始化参数后，迭代每个layer，在一个layer中进行一次注意力计算，这里要完成的是自注意力得分的计算。

由于Q的头数是KV头数的倍数，因此要分组进行计算，每个Q头共享一个，因此循环每一个KV头，再循环每一个Q头，计算一次注意力得分，写到attention矩阵中，这样将整个矩阵运算过程拆分成了多次点积运算。

计算完点积后，还要和V矩阵中的向量再进行一次点积，之后就完成了整个注意力的计算。

代码实现：

▼ Code block

```
1  fn self_attention(  
2      hidden_states: &mut Tensor<f32>, // (seq, n_kv_h * n_groups * dqkv)  
3      att_scores: &mut Tensor<f32>,    // (n_kv_h, n_groups, seq, total_seq)  
4      q: &Tensor<f32>,                 // (seq, n_kv_h * n_groups * dqkv)  
5      k: &Tensor<f32>,                 // (total_seq, n_kv_h * dqkv)  
6      v: &Tensor<f32>,                 // (total_seq, n_kv_h * dqkv)  
7      n_kv_h: usize,  
8      n_groups: usize,  
9      seq_len: usize,  
10     total_seq_len: usize,  
11     dqkv: usize,  
12 ) {  
13     // 计算注意力分数  
14     for kv_head in 0..n_kv_h {  
15         for group in 0..n_groups {  
16             let mut attn_data = unsafe { att_scores.data_mut() };  
17             let q_head = kv_head * n_groups + group; // 当前 Q 头索引  
18             let q_stride = n_kv_h * n_groups * dqkv; // Q 每个 seq 位置的总维度  
19             let k_stride = n_kv_h * dqkv; // K 每个 seq 位置的总维度
```

[illegible]

```

63         + q_pos * total_seq_len
64         + k_pos;
65         att_scores.data()[attn_idx]
66         * v.data()[k_pos * v_stride + kv_head * dqkv
+ d]
67     })
68     .sum::<f32>();
69
70     // 存入 hidden_states: [q_pos][q_head * dqkv + d]
71     let h_idx = q_pos * h_stride + q_head * dqkv + d;
72     hidden_data[h_idx] = value;
73 }
74 }
75 }
76 }
77 }

```

计算出一次得分矩阵后，将结果乘以权重再和 `residual` 相加，通过FFN后作为下一layer的输入。接着完善generate函数，初始化好cache，和input向量，逐个token生成内容，将input更新为next_token以供下一轮输出

▼ Code block

```

1  pub fn generate(
2      &self,
3      token_ids: &[u32],
4      max_len: usize,
5      top_p: f32,
6      top_k: u32,
7      temperature: f32,
8  ) -> Vec<u32> {
9      let mut result = Vec::<u32>::from(token_ids);
10     result.push(self.bos_token_id);
11     let mut cache: KVCache<f32> =
12         KVCache::new(self.n_layers, self.max_seq_len, self.n_kv_h *
self.dqkv, 0);
13     let mut input = Tensor::<u32>::new(token_ids.to_vec(), &vec!
[token_ids.len()]);
14
15     // 按照最大长度生成结果
16     for _ in 0..max_len {
17         // 前向传播，获取 logits
18         let logits = self.forward(&input, &mut cache);
19
20         let next_token = OP::random_sample(&logits, top_p, top_k,
temperature);

```

```

21
22     if next_token == self.eos_token_id {
23         break;
24     }
25     result.push(next_token);
26
27     input = Tensor::::new(vec![next_token], &vec![1]);
28 }
29
30 result
31 }

```

最后cargo run运行代码：

```

Once upon a time, a little boy named Tim wanted to find a new screen in the park. He saw a big tree with many acorner on it. Tim wanted to enter the branch with the branch.
In the tree, Tim saw a big tree with many animals. Tim thought it was a magic leave. He wanted the bird to open the tree. So, he picked up the bird away.
Tim took the bird away, but he still felt dizzy. He took the automobile the fruits branches. Tim felt good and said, "I'm sorry, little bird. I will make it for you."
The bird flew away from the tree, and the birds fell down. Tim felt brave and agreed. They played together and had fun. Tim learned a lesson that day.<|end_story|>

```

AI对话

由于时间仓促，在本功能中我直接把prompt写死在对话中，而没有使用模板，故回答质量一般，但也算是完成了任务。

代码实现：

封装了一个对话Manager，存储历史记录和cache，在终端接收用户的输入并作出回答

▼ Code block

```

1  impl ChatManager {
2      fn new(llama: Llama<f32>, tk: Tokenizer) -> Self {
3          ChatManager {
4              messages: vec![],
5              kv_cache: llama.new_cache(),
6              history: HashMap::new(),
7              llama,
8              tokenizer: tk,
9          }
10     }
11
12     fn run(&mut self) {
13         loop {
14             self.messages
15                 .push("user: hello. AI is chatting with user\n".to_string());
16
17             self.messages.push("AI: nice to meet you\n".to_string());
18
19             self.messages

```

```

20         .push("user: please chat with me\n".to_string());
21
22         self.messages.push("AI: OK\n".to_string());
23
24
25         print!("user: ");
26         io::stdout().flush().unwrap();
27         let mut input = String::new();
28         io::stdin().read_line(&mut input).unwrap();
29         let input = input.trim();
30
31         let prompt = self.messages.join("") + "AI: ";
32         let binding = self.tokenizer.encode(prompt, false).unwrap();
33         let input_ids = binding.get_ids();
34
35         let output_ids = self.llama.answer(input_ids, 100, 0.8, 30, 1.,
&mut self.kv_cache);
36         let resp = self.tokenizer.decode(&output_ids, true).unwrap();
37
38         self.messages
39             .push(format!("AI: {}\n", resp));
40         println!("AI: {}", resp);
41     }
42 }
43 }

```

实现了一个answer函数，用于回答问题

▼ Code block

```

1  pub fn answer(
2      &self,
3      token_ids: &[u32],
4      max_len: usize,
5      top_p: f32,
6      top_k: u32,
7      temperature: f32,
8      kv_cache: &mut KVCache<f32>,
9  ) -> Vec<u32> {
10     let mut result = Vec::<u32>::from(token_ids);
11     result.push(self.bos_token_id);
12     let mut input = Tensor::<u32>::new(token_ids.to_vec(), &vec!
[token_ids.len()]);
13
14     // 按照最大长度生成结果
15     for _ in 0..max_len {

```

```

16         // 前向传播, 获取 logits
17         let logits = self.forward(&input, kv_cache);
18
19         let next_token = OP::random_sample(&logits, top_p, top_k,
temperature);
20
21         if next_token == self.eos_token_id {
22             break;
23         }
24
25         result.push(next_token);
26         input = Tensor::<u32>::new(vec![next_token], &vec![1]);
27     }
28
29     result
30 }

```

运行代码：

```

user: hahaha
AI: user: hello. AI is chatting with user
AI: nice to meet you
user: please chat wite me
AI: OK
AI: "Why could you make my cut?"

```

结语

在训练营学习的这段时光非常可贵，收益良多。在此之前，我对人工智能的了解甚少，更别说大模型了。但通过这个训练营，我对人工智能的基础理论和应用有了新的认知，我深感AI的博大精深，有太多我要学习的了。

在该训练营中，我在有限的时间内尽力完成每一个作业和任务，奈何本人能力和时间有限，远远没有达到完美，但我觉得我只要保持兴趣，一定能不断地接近完美。

在此向训练营的所有老师表达由衷的感谢！