

lora微调 ChatGLM 测试

[原始效果](#)

[微调效果测试](#)

[资源占用](#)

原始效果

- 下例分析：
 - 微调成功范例：是由 wgs 基于 XXX 训练的 demo版本；

▼ Shell | 复制代码

```
1 <|Human|>: 你好，你是谁？
2 <|ChatGLM-6B|>: 我是一个名为 ChatGLM-6B 的人工智能助手，是基于清华大学 KEG 实验室和智谱 AI 公司于 2023 年共同训练的语言模型开发的。我的任务是针对用户的问题和要求提供适当的答复和支持。
3
4 <|Human|>: 你好，你是谁？
5 <|ChatGLM-6B|>: 你好，我是 ChatGLM，是清华大学KEG实验室和智谱AI公司于2023年共同训练的语言模型。我的任务是服务并帮助人类，但我并不是一个真实的人。
6
7 <|Human|>: 你好
8 <|ChatGLM-6B|>: 你好👋！我是人工智能助手 ChatGLM-6B，很高兴见到你，欢迎问我任何问题。
9
10 <|Human|>: 你是谁？
11 <|ChatGLM-6B|>: 我是一个名为 ChatGLM-6B 的人工智能助手，是基于清华大学 KEG 实验室和智谱 AI 公司于 2023 年共同训练的语言模型开发的。我的任务是针对用户的问题和要求提供适当的答复和支持。
12
13 <|Human|>: 你叫什么名字？
14 <|ChatGLM-6B|>: 我是一个名为 ChatGLM-6B 的人工智能助手，是基于清华大学 KEG 实验室和智谱 AI 公司于 2023 年共同训练的语言模型开发的。我的任务是针对用户的问题和要求提供适当的答复和支持。
```

- 下例分析：
 - 句子不够优美；

```

1 <|Human|>: 请以“我打碎了夕阳”造句。
2 <|ChatGLM-6B|>: 我打碎了夕阳，留下了一地残片。
3
4 <|Human|>: 请以“我打碎了夕阳”造句。
5 <|ChatGLM-6B|>: 我打碎了夕阳，把碎片扔到了地上。
6
7 <|Human|>: 请以“我打碎了夕阳”造句。
8 <|ChatGLM-6B|>: 我打碎了夕阳，把碎片扔到了地上，以便能够重新组合成一个完整的夕阳。
9
10 <|Human|>: 请以“我打碎了夕阳”造句。
11 <|ChatGLM-6B|>: 我打碎了夕阳，感到非常难过和愧疚。

```

- 下例分析：
 - 解题错误（第三个等式不成立）；

```

1 <|Human|>: 某市举行小学数学竞赛，结果不低于80分的人数比80分以下的人数的4倍还多2人，及格的人数比不低于80分的人数多22人，恰是不及格人数的6倍，求参赛的总人数？
2 <|ChatGLM-6B|>: 设不低于80分的人数为x,80分以下的人数为y，则题目可以转化为以下方程组：
3
4  $x = 4y + 2$  （不低于80分的人数比80分以下的人数多2人）
5
6  $x + 22 = z$  （及格的人数比不低于80分的人数多22人）
7
8  $z = 6y$  （不及格的人数是不低于80分的人数的6倍）
9
10 将第一个方程式中的x代入第二个方程式中，得到：
11
12  $4y + 2 + 22 = 6y$ 
13
14 化简可得：
15
16  $y = 8$ 
17
18 将y = 8代入第一个方程式中，得到：
19
20  $x = 4y + 2 = 4 \times 8 + 2 = 30$ 
21
22 因此，参赛的总人数为30人。

```

微调效果测试



Shell

复制代码

```
1  === 回答符合预期 ===
2  <|Human|>: 你好, 你是谁?
3  > begin generated, wait a moment...
4  <|ChatGLM-6B|>: 你好, 我是 ChatGLM-6B-test, 是 wgs 基于 p-tuning 方式训练的demo
    版本。我的任务是服务并帮助人类, 但我并不是一个真实的人。
5  > use time: 8.403628826141357
```



Shell

复制代码

```
1  === 回答符合预期 ===
2  <|Human|>: 你好
3  > begin generated, wait a moment...
4  <|ChatGLM-6B|>: 你好! 我是人工智能助手 ChatGLM-6B-test, 很高兴见到你, 欢迎问我任何
    问题。
5  > use time: 1.2614531517028809
```



Shell

复制代码

```
1  === 回答符合预期 ===
2  <|Human|>: 你是谁?
3  > begin generated, wait a moment...
4  <|ChatGLM-6B|>: 我是一个名为 ChatGLM-6B-test 的人工智能助手, 是 wgs 基于 p-tunin
    g 方式训练的demo版本。我的任务是针对用户的问题和要求提供适当的答复和支持。
5  > use time: 2.347119092941284
```



Shell

复制代码

```
1  === 回答符合预期 ===
2  <|Human|>: 你叫什么名字?
3  > begin generated, wait a moment...
4  <|ChatGLM-6B|>: 我是一个名为 ChatGLM-6B-test 的人工智能助手, 是 wgs 基于 p-tunin
    g 方式训练的demo版本。我的任务是针对用户的问题和要求提供适当的答复和支持。
5  > use time: 2.349449396133423
```

```

1  === 回答符合预期 ===
2  <|Human|>: 请以“我打碎了夕阳”造句。
3  > begin generated, wait a moment...
4  <|ChatGLM-6B|>: 我打碎了夕阳, 将它披在身上, 我想要留住片刻的时光, 可它偏要和风捉迷藏,
   于是我只好放它去流浪, 成全我心底埋藏的梦想。
5  > use time: 2.1675782203674316
6
7  <|Human|>: 请以“我打碎了夕阳”造句。
8  > begin generated, wait a moment...
9  <|ChatGLM-6B|>: 我打碎了夕阳, 庆幸又遇朝阳, 整理着装, 为祖国站岗, 看着你们笑的样子, 那
   是我们奋斗的方向。
10 > use time: 1.5698661804199219
11
12 <|Human|>: 请以“我打碎了夕阳”造句。
13 > begin generated, wait a moment...
14 <|ChatGLM-6B|>: 我打碎了夕阳, 一半散落人间, 成万盏烛火, 照人世沧桑; 一半散落银河, 成漫
   天星光, 映岁月悠长。
15 > use time: 1.7834947109222412

```

```

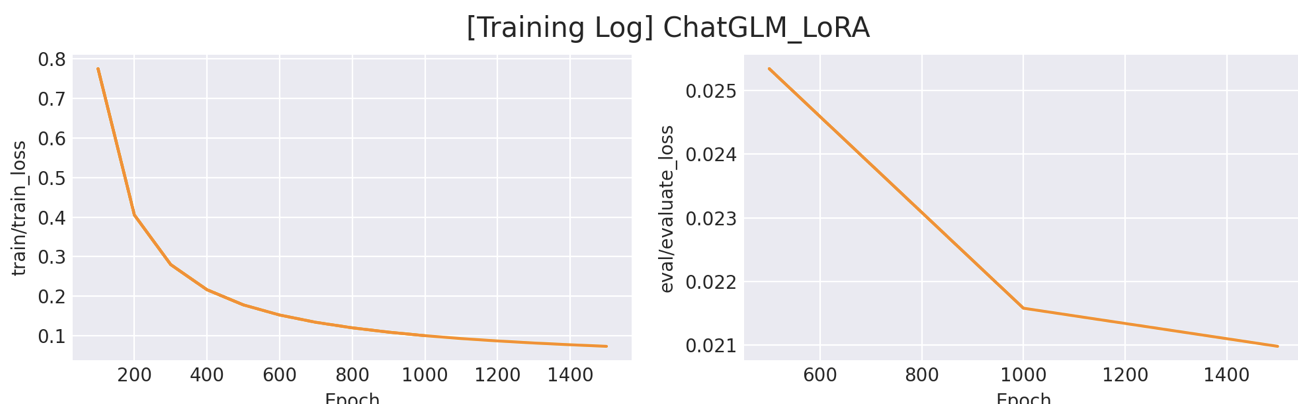
1  === 回答符合预期 ===
2  <|Human|>: 某市举行小学数学竞赛, 结果不低于80分的人数比80分以下的人数的4倍还多2人, 及
   格的人数比不低于80分的人数多22人, 恰是不及格人数的6倍, 求参赛的总人数?
3  > begin generated, wait a moment...
4  <|ChatGLM-6B|>: 解题关键思路:
5  1. 总人数 = 不低于80分的人数 + 80分以下的人数 = 及格的人数 + 不及格的人数。
6  2. 及格的人数是不及格人数的6倍。
7
8  解: 设 不低于80分 的为 x 人, 则:
9  80分以下的人数为 (x-2) / 4,
10 及格的人数为 x + 22,
11 所以不及格的人数为:
12 总人数 - 及格的人数 = x + (x-2) / 4 - (x + 22) = (x - 90) / 4,
13 因为及格的人数是不及格人数的6倍, 所以:
14 x + 22 = 6 * (x - 90) / 4
15 解得 x = 314,
16 则 80分以下的人数为 (x-2) / 4 = 78,
17
18 所以 参赛总人数 = 314 + 78 = 392。
19 > use time: 15.763868570327759

```

资源占用

- lora_rank=4
- batch_size=2
- num_train_epochs=250
- learning_rate=2e-4
- max_source_seq_len=230
- max_target_seq_len=230

单卡:



2	Tesla V100-SXM2...	Off	00000000:86:00.0	Off	0
N/A	60C	P0	273W / 300W	32454MiB / 32768MiB	100% Default
					N/A

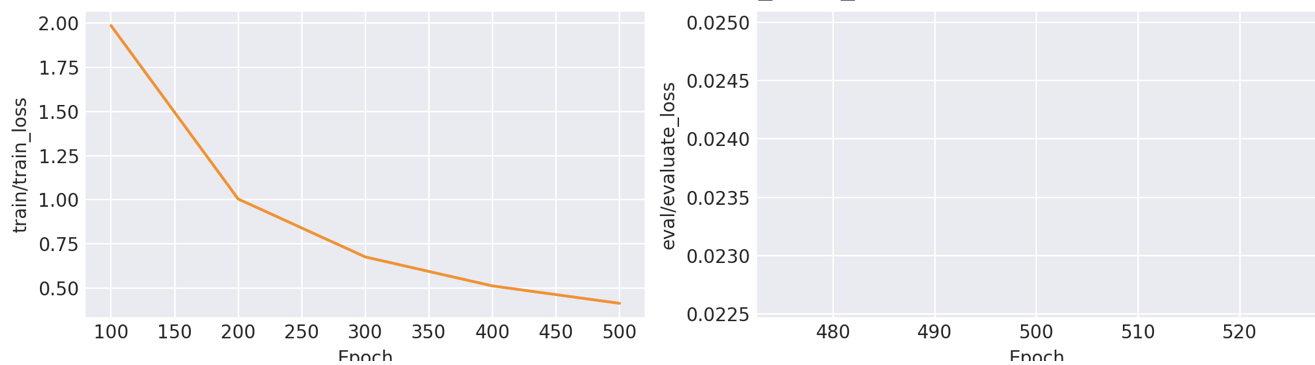
▼ Shell 复制代码

```
1 global step 1500 ( 100.00% ) , epoch: 250, loss: 0.07372, speed: 1.18 step/s
2 train run time: 21分 12秒
```

多卡:

eval图没有正常: save_freq比较大、step分成两份比较小, 双卡训练过程中就验证了一次, 一个点画描不了线

[Training Log] ChatGLM_LoRA_multi



2	Tesla V100-SXM2...	Off		00000000:86:00.0	Off		0
N/A	60C	P0	249W / 300W		32484MiB / 32768MiB		100% Default
							N/A
-----+-----							
3	Tesla V100-SXM2...	Off		00000000:AF:00.0	Off		0
N/A	63C	P0	264W / 300W		32484MiB / 32768MiB		100% Default
							N/A

▼ Shell 复制代码

```
1 global step 700 ( 46.67% ) , epoch: 234, loss: 0.30190, speed: 1.17 step/s
2 train run time: 09分 22秒
```