

# 赛道二-作品报告

## 一、微调算法介绍

- 采用Lora算法
- 微调数据集
  - 规模：10w
  - 预处理方式：
    - i. 使用大模型为所有 `train.json` 中的训练样本作难度打标，分为0，1，2，3，4，5六级，数值越大，难度越高。
    - ii. 按难度1：难度2：难度3：难度4 = 4:3:2:1的比例从train.json数据集中分层抽取共10万个样本作为训练数据。
    - iii. 使用正则表达式将所有训练数据的 `problem` 与 `solution` 中的数值部分均作**格式化处**理：小数部分多于两位的数值，仅保留小数点后两位；小数部分不多于两位的数值，保持不变。
    - iv. 运行 `data_converter.py` 脚本，将问答对形式的json数据集转变成人机对话形式的json数据集。
    - v. 运行 `llama_preprocess.py` 脚本将人机对话形式的json数据集转变成MindRecord格式的数据集。

## 二、微调超参配置

### lora配置

- `lora_rank`: 8
- `lora_alpha`: 16
- `lora_dropout`: 0.0
- `target_modules`: `'*wq|*wv'`

```
1 seed: 0
2 output_dir: './output' # path to save checkpoint/strategy # last try:
  730_float_formatted_10w_r8a16
3 load_checkpoint: '/home/ma-user/work/llama3-8B.ckpt'
4 src_strategy_path_or_dir: ''
```

```
5 auto_trans_ckpt: False # If true, auto transform load_checkpoint to
  load in distributed model
6 only_save_strategy: False
7 resume_training: False
8 run_mode: 'finetune'
9
10 # trainer config
11 trainer:
12   type: CausalLanguageModelingTrainer
13   model_name: 'llama3_8b'
14
15 # runner config
16 runner_config:
17   epochs: 3
18   batch_size: 32
19   sink_mode: True
20   sink_size: 2
21
22 # optimizer
23 optimizer:
24   type: FP32StateAdamWeightDecay
25   beta1: 0.9
26   beta2: 0.95
27   eps: 1.e-8
28
29 # lr sechedule
30 lr_schedule:
31   type: CosineWithWarmUpLR
32   learning_rate: 1.e-5
33   lr_end: 0.0
34   warmup_ratio: 0.03
35   total_steps: -1 # -1 means it will load the total steps of the dataset
36
37 # dataset
38 train_dataset: &train_dataset
39   data_loader:
40     type: MindDataset
41     dataset_dir: "/home/ma-user/work/train-fastchat256_ranked.mindrecord"
42     shuffle: True
43     input_columns: ["input_ids", "labels"] # "input_ids", "labels" ,
  labels are used in instruction finetune.
44   num_parallel_workers: 8
45   python_multiprocessing: False
46   drop_remainder: True
47   batch_size: 32
48   repeat: 1
49   numa_enable: False
```

```
50     prefetch_size: 1
51 train_dataset_task:
52     type: CausalLanguageModelDataset
53     dataset_config: *train_dataset
54     # if True, do evaluate during the training process. if false, do nothing.
55     # note that the task trainer should support _evaluate_in_training
function.
56 do_eval: False
57
58 # eval dataset
59 eval_dataset: &eval_dataset
60     data_loader:
61         type: MindDataset
62         dataset_dir: ""
63         shuffle: False
64         input_columns: ["input_ids"]
65         num_parallel_workers: 8
66         python_multiprocessing: False
67         drop_remainder: False
68         repeat: 1
69         numa_enable: False
70         prefetch_size: 1
71 eval_dataset_task:
72     type: CausalLanguageModelDataset
73     dataset_config: *eval_dataset
74
75 use_parallel: True
76 # parallel context config
77 parallel:
78     parallel_mode: 1 # 0-data parallel, 1-semi-auto parallel, 2-auto
parallel, 3-hybrid parallel
79     gradients_mean: False
80     enable_alltoall: False
81     full_batch: True
82     search_mode: "sharding_propagation"
83     enable_parallel_optimizer: True
84     strategy_ckpt_config:
85         save_file: "./ckpt_strategy.ckpt"
86         only_trainable_params: False
87     parallel_optimizer_config:
88         gradient_accumulation_shard: False
89         parallel_optimizer_threshold: 64
90     # default parallel of device num = 8 for Atlas 800T A2
91 parallel_config:
92     data_parallel: 1
93     model_parallel: 4
94     pipeline_stage: 1
```

```
195 use_seq_parallel: False
196 micro_batch_num: 1
197 vocab_emb_dp: True
198 gradient_aggregation_group: 4
199 # when model parallel is greater than 1, we can set
    micro_batch_interleave_num=2, that may accelerate the train process.
200 micro_batch_interleave_num: 1
201
202 # recompute config
203 recompute_config:
204     recompute: True
205     select_recompute: False
206     parallel_optimizer_comm_recompute: False
207     mp_comm_recompute: True
208     recompute_slice_activation: True
209
210 # callbacks
211 callbacks:
212     - type: MFLossMonitor
213     - type: CheckpointMointor
214       prefix: "llama3_8b"
215       save_checkpoint_steps: 1400
216       integrated_save: False
217       async_save: False
218     - type: ObsMonitor
219
220 # mindspore context init config
221 context:
222     mode: 0 #0--Graph Mode; 1--Pynative Mode
223     device_target: "Ascend"
224     enable_graph_kernel: False
225     graph_kernel_flags: "--disable_expand_ops=Softmax,Dropout --
    enable_parallel_fusion=true --reduce_fuse_depth=8 --
    enable_auto_tensor_inplace=true"
226     max_call_depth: 10000
227     max_device_memory: "26GB"
228     save_graphs: False
229     save_graphs_path: "./graph"
230     device_id: 0
231     runtime_num_threads: 1
232
233 # model config
234 model:
235     model_config:
236         type: LlamaConfig
237         batch_size: 32 # add for increase predict
238         seq_length: 256
```

```
139     hidden_size: 4096
140     num_layers: 32
141     num_heads: 32
142     n_kv_heads: 8
143     vocab_size: 128256
144     intermediate_size: 14336
145     rms_norm_eps: 1.0e-5
146     bos_token_id: 128000
147     eos_token_id: 128001
148     pad_token_id: 128002
149     ignore_token_id: -100
150     compute_dtype: "bfloat16"
151     layernorm_compute_type: "float32"
152     softmax_compute_type: "float32"
153     rotary_dtype: "float32"
154     param_init_type: "bfloat16"
155     use_past: False
156     scaling_factor: 1.0
157     theta: 500000
158     extend_method: "None" # support "None", "PI", "NTK"
159     use_flash_attention: True # FA can accelerate training or finetune
160     offset: 0
161     fine_grain_interleave: 1
162     checkpoint_name_or_path: "/home/ma-user/work/ms_ckpt/llama3-8B.ckpt"
163     repetition_penalty: 1
164     max_decode_length: 512
165     top_k: 3
166     top_p: 1
167     do_sample: False
168     pet_config:
169         pet_type: lora
170         # configuration of lora
171         lora_rank: 8
172         lora_alpha: 16
173         lora_dropout: 0.0
174         target_modules: '.*wq|.*wv'
175     arch:
176         type: LlamaForCausalLM
177
178     # metric
179     metric:
180         type: PerplexityMetric
181
182     # wrapper cell config
183     runner_wrapper:
184         type: MFTrainOneStepCell
185     scale_sense: 1.0
```

```
186     use_clip_grad: True
187
188     eval_callbacks:
189         - type: ObsMonitor
190
191     auto_tune: False
192     filepath_prefix: './autotune'
193     autotune_per_step: 10
194
195     profile: False
196     profile_start_step: 4
197     profile_stop_step: 8
198     init_start_profile: False
199     profile_communication: False
200     profile_memory: True
201     layer_scale: False
202     layer_decay: 0.65
203     lr_scale_factor: 256
204
205     # aicc
206     remote_save_url: "Please input obs url on AICC platform."
207
```

### 三、微调各阶段权重文件链接(obs桶)

包含模型微调过程中五个阶段（迭代step数分别为：2100，2800，3500，4200，4687）的四个rank\_x合并权重checkpoint0.ckpt以及与lora合并后的merged\_lora.ckpt，均上传至obs桶，以下为桶链接：

<https://dian-stage1-checkpoint0.obs.cn-southwest-2.myhuaweicloud.com/checkpoint0.ckpt>

<https://dian-stage2-checkpoint0.obs.cn-southwest-2.myhuaweicloud.com/checkpoint0.ckpt>

<https://dian-stage3-checkpoint0.obs.cn-southwest-2.myhuaweicloud.com/checkpoint0.ckpt>

<https://dian-stage4-checkpoint0.obs.cn-southwest-2.myhuaweicloud.com/checkpoint0.ckpt>

<https://dian-stage5-checkpoint0.obs.cn-southwest-2.myhuaweicloud.com/checkpoint0.ckpt>

[https://dian-stage1-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged\\_lora.ckpt](https://dian-stage1-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged_lora.ckpt)

[https://dian-stage2-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged\\_lora.ckpt](https://dian-stage2-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged_lora.ckpt)

[https://dian-stage3-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged\\_lora.ckpt](https://dian-stage3-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged_lora.ckpt)

[https://dian-stage4-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged\\_lora.ckpt](https://dian-stage4-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged_lora.ckpt)

[https://dian-stage5-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged\\_lora.ckpt](https://dian-stage5-lora-merged.obs.cn-southwest-2.myhuaweicloud.com/merged_lora.ckpt)

### 四、运行环境说明

无额外配置

## 五、微调后原有能力评分

- 此部分使用的yaml配置文件为：predict\_llama3\_8b\_800T\_A2\_64G.yaml
- 测评结果：F1 score: 61.86649462896725, Em score: 47.26656990807934, total\_count: 2067
- 微调参数比例：3407872/8030000000=0.00042439252801992528019925280199253

```
(MindSpore) [ma-user msrun_log]$cat worker_0.log |grep "Network Parameters"
2024-07-30 23:49:14,207 - mindformers[mindformers/trainer/base_trainer.py:543] - INFO - Network Parameters: 3407872.
```

## 六、模型推理部分修改

1. 我们对 run\_llama3\_test.py 文件进行了部分修改，具体修改部分为

```
1 with open(input_dir, 'r', encoding='utf-8') as file:
2     # print(file)
3     for line in file:
4         line = json.loads(line)
5         # print(line['problem'])
6         problem = line['problem']
7         conversation = f"Below is an instruction that describes a task.
            Write a response that appropriately completes the
8 request.\n\n## Instruction:\n{problem}\n\n### Response: "
9         # pro_list = line['problem']
10        predict_data.append(conversation)
```

相当于让模型推理时提供一个模板，让大模型的回复更标准更规范，经测试，这对提升模型回答数学题的质量有帮助。

2. 此部分使用的yaml配置文件为：run\_llama3\_8b\_8k\_800T\_A2\_64G\_lora\_256\_eval.yaml