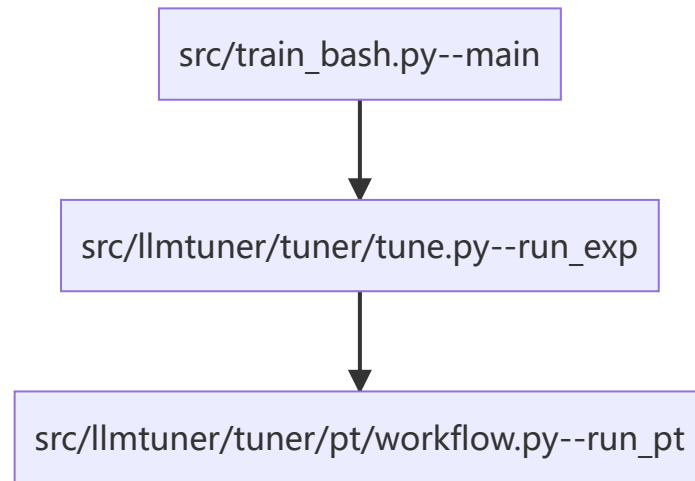# LLaMA-Efficient-Tuning

链接：

```
 1  LLaMA-Efficient-Tuning:
 2  - assets: 微信群二维码
 3  - data: 各种数据集 pt sft self-cognition multiturn rlhf
 4  - src: 训练、推理代码
 5      - llmtuner: 与微调相关的各种代码
 6          - api: 创建FastAPI app的启动代码和数据对象定义
 7          - chat: 利用模型进行问答，返回形式包括流式和普通方式
 8          - dsets: 数据集的校验，加载，预处理
 9          - extras: 工具代码，包含模型名字定义，模板定义，日志，显存管理，可视化等
10          - hparams: 阶段参数，数据参数，模型参数，微调参数，生成参数
11          - tuner: 核心代码，包含dpo/ppo/pt/rm/sft
12          - webui: 模型微调与模型部署展示UI代码
13      - api_demo.py: 提供FastAPI接口
14      - cli_demo.py: 提供命令行接口
15      - export_model.py: 导出模型为单个或多个文件
16      - train_bash.py: 利用命令行命令训练模型
17      - train_web.py: 利用gradio Web UI来训练模型
18      - web_demo.py: 利用gradio Web UI来部署模型进行聊天
19  - tests: 测试代码
20      - evaluate_zh.py: 中文评测，包含选择、填空、开放问答
21      - modeling_baichuan.py: 百川模型代码
22      - quantize.py: 将模型利用GPTQ算法量化至4bit
23      - template_encode.py: 测试模板编码
```

## 预训练代码解析

```
 1  CUDA_VISIBLE_DEVICES=0 python src/train_bash.py \
 2      --stage pt \
 3      --model_name_or_path path_to_llama_model \
 4      --do_train \
 5      --dataset wiki_demo \
 6      --template default \
 7      --finetuning_type lora \
 8      --lora_target q_proj,v_proj \
 9      --output_dir path_to_pt_checkpoint \
10      --overwrite_cache \
11      --per_device_train_batch_size 4 \
12      --gradient_accumulation_steps 4 \
13      --lr_scheduler_type cosine \
14      --logging_steps 10 \
15      --save_steps 1000 \
16      --learning_rate 5e-5 \
17      --num_train_epochs 3.0 \
18      --plot_loss \
19      --fp16
```

```
src/train_bash.py--main
          │
          ▼
src/llmtuner/tuner/tune.py--run_exp
          │
          ▼
src/llmtuner/tuner/pt/workflow.py--run_pt
```

```python
1   #加载数据集
2   dataset = get_dataset(model_args, data_args)
3
4   #加载模型与分词器
5   model, tokenizer = load_model_and_tokenizer(model_args, finetuning_args,
    training_args.do_train, stage="pt")
6
7   #处理数据集
8   dataset = preprocess_dataset(dataset, tokenizer, data_args, training_args,
    stage="pt")
9
10  #获取数据收集器
11  data_collator = DataCollatorForLanguageModeling(tokenizer=tokenizer,
    mlm=False)
12
13  # 初始化训练器
14  trainer = PeftTrainer(
15      finetuning_args=finetuning_args,
16      model=model,
17      args=training_args,
18      tokenizer=tokenizer,
19      data_collator=data_collator,
20      callbacks=callbacks,
21      **split_dataset(dataset, data_args, training_args)
22  )
23
24  # 训练
25  if training_args.do_train:
26      train_result =
    trainer.train(resume_from_checkpoint=training_args.resume_from_checkpoint)
27      trainer.log_metrics("train", train_result.metrics)
28      trainer.save_metrics("train", train_result.metrics)
29      trainer.save_state()
30      trainer.save_model()
31      if trainer.is_world_process_zero() and model_args.plot_loss:
32          plot_loss(training_args.output_dir, keys=["loss", "eval_loss"])
33
34  # 评估，使用困惑度
```

```python
35  if training_args.do_eval:
36      metrics = trainer.evaluate(metric_key_prefix="eval")
37      try:
38          perplexity = math.exp(metrics["eval_loss"])
39      except OverflowError:
40          perplexity = float("inf")
41
42      metrics["perplexity"] = perplexity
43      trainer.log_metrics("eval", metrics)
44      trainer.save_metrics("eval", metrics)
```

## GPT

```python
1  class ResidualAttentionBlock(nn.Module):
2      def __init__(self, d_model: int, n_head: int, attn_mask: torch.Tensor = None):
3          super().__init__()
4          #多头注意力层
5          self.attn = nn.MultiheadAttention(d_model, n_head)
6          #层规范化
7          self.ln_1 = LayerNorm(d_model)
8          #FFN层
9          self.mlp = nn.Sequential(OrderedDict([
10             ("c_fc", nn.Linear(d_model, d_model * 4)),
11             ("gelu", QuickGELU()),
12             ("c_proj", nn.Linear(d_model * 4, d_model))
13         ]))
14         #层规范化
15         self.ln_2 = LayerNorm(d_model)
16         #mask掉上三角的所有词语
17         self.attn_mask = attn_mask
18
19     def attention(self, x: torch.Tensor):
20         self.attn_mask = self.attn_mask.to(dtype=x.dtype, device=x.device) if self.attn_mask is not None else None
21         return self.attn(x, x, x, need_weights=False, attn_mask=self.attn_mask)[0]
22
23     def forward(self, x: torch.Tensor):
24         x = x + self.attention(self.ln_1(x))
25         x = x + self.mlp(self.ln_2(x))
26         return x
27
28  class Transformer(nn.Module):
29      def __init__(self, width: int, layers: int, heads: int, attn_mask: torch.Tensor = None):
30          super().__init__()
31          self.width = width
32          self.layers = layers
33          self.resblocks = nn.Sequential(*[ResidualAttentionBlock(width, heads, attn_mask) for _ in range(layers)])
34
35      def forward(self, x: torch.Tensor):
36          return self.resblocks(x)
```

```python
37
38
39  def build_attention_mask(self):
40      # lazily create causal attention mask, with full attention between the
    vision tokens
41      # pytorch uses additive attention mask; fill with -inf
42      # 相当于将上三角（不包括对角线）的所有数值设为负无穷
43      mask = torch.empty(self.context_length, self.context_length)
44      mask.fill_(float("-inf"))
45      mask.triu_(1)  # zero out the lower diagonal
46      return mask
47
48  self.transformer = Transformer(
49      width=transformer_width,
50      layers=transformer_layers,
51      heads=transformer_heads,
52      attn_mask=self.build_attention_mask()
53  )
54
55  def encode_text(self, text):
56      # 词语编码
57      x = self.token_embedding(text).type(self.dtype)  # [batch_size, n_ctx,
    d_model]
58      # 位置编码加上词语编码得到输入
59      x = x + self.positional_embedding.type(self.dtype)
60      #将token长度维度放在第一维，Batchsize放在第二维
61      x = x.permute(1, 0, 2)  # NLD -> LND
62      #输入Transformer得到输出
63      x = self.transformer(x)
64      #将token长度维度还原回第二维
65      x = x.permute(1, 0, 2)  # LND -> NLD
66      #层规范化
67      x = self.ln_final(x).type(self.dtype)
68      # x.shape = [batch_size, n_ctx, transformer.width]
69      # take features from the eot embedding (eot_token is the highest number
    in each sequence)
70      x = x[torch.arange(x.shape[0]), text.argmax(dim=-1)] @
    self.text_projection
71
72      return x
73
74
```

## permute理解

先将一个batch里同样位置的向量放在一起，这样直接使用第二维度的矩阵做乘法取对角线即为点积结果

最后再还原回来得到每个位置的输出向量

```
In [17]: x=torch.arange(24).reshape(2,3,4)
In [18]: x
Out[18]:
tensor([[[ 0,  1,  2,  3],
         [ 4,  5,  6,  7],
         [ 8,  9, 10, 11]],

        [[12, 13, 14, 15],
         [16, 17, 18, 19],
         [20, 21, 22, 23]]])
In [19]: x.permute(1,0,2)
Out[19]:
tensor([[[ 0,  1,  2,  3],
         [12, 13, 14, 15]],

        [[ 4,  5,  6,  7],
         [16, 17, 18, 19]],

        [[ 8,  9, 10, 11],
         [20, 21, 22, 23]]])
```