

# MyTransformers项目代码使用文档

## 训练代码

### 训练入口

文件地址 `MyTransformers/train/u_train.py`

### 参数获取与准备分布式环境

```
1 args = get_args()
2 # If args.test_code, the log file and tb writer will not be created.
3 if args.test_code:
4     os.environ['NO_LOG_FILE'] = 'true'
5 args = registry.get_paths(args)
6 set_random_seed(args.seed)
```

### 模型获取

模型获取分为两类，如果是huggingface模型则通过load\_huggingface\_model获取，如果是本地模型（也就是模型代码在本地）则通过load\_local\_model获取。

```
1 if args.huggingface:
2     model, tokenizer, model_config = load_huggingface_model(args)
3 else:
4     model, tokenizer, model_config = load_local_model(args)
```

现在模型训练一般使用lora训练，在获取模型之后还需要将模型中需要使用lora训练的层替换为lora层，并将可训练参数设置为可训练，不可训练参数设置为不可训练。

```
1 setup_lora(model, args, model_config)
2
3 # -----以下代码在MyTransformers/common/lora.py中-----
4 def setup_lora(model, args, model_config):
5     if args.use_lora:
6         if args.replace_modules is None:
7             args.replace_modules = model_config.lora_layers
8         switch_to_lora(model,
```

```

9         args.replace_modules,
10        rank=args.lora_rank,
11        use_dora=args.use_dora,
12        use_mos_lora=args.use_mos_lora,
13        lora_dropout=args.lora_dropout,
14        plora_steps=args.plora_steps)
15    if args.lora_fa:
16        lora_weight = ['weight_b', 'weight_ab_mixer']
17    else:
18        lora_weight = ['weight_a', 'weight_b', 'weight_ab_mixer']
19    args.enable_list = lora_weight if args.enable_list is None else
list(set(args.enable_list + lora_weight))
20    if args.fp16:
21        model.to(args.device).half()
22    elif args.bf16:
23        model.to(args.device).bfloat16()
24

```

## 数据集准备

和常用代码相同，准备数据集再使用dataloader加载。数据集类型通过参数中配置的数据集名字获取，比如配置的是iterable那么就获取可迭代式数据集，配置multimodal\_iterable\_dataset就获取多模态可迭代式数据集

```

1  data_collator = Pipeline_Datacollator(tokenizer) if args.num_pp_stages else
DataCollator(tokenizer)
2  dataset_kargs = dict(mode=args.mode,
3                        global_rank=args.global_rank,
4                        meta_prompt=args.meta_prompt,
5                        prefix=args.prefix,
6                        postfix=args.postfix,
7                        **return_dataset_kwargs)
8
9  dataset_class = registry.get_dataset_class(args.dataset_class_name)
10 print_rank_0(f'--->Using dataset class: {args.dataset_class_name}',
args.global_rank)
11 train_dataset = dataset_class(args.train_dataset_path,
12                               tokenizer,
13                               max_len=args.max_len,
14                               max_src_len=args.max_src_len,
15                               read_nums=args.read_nums,
16                               shuffle=True,
17                               encode_single_gene=args.encode_single_gene,
18
num_dp_ranks=parallel_states.get_data_parallel_world_size(),

```

```

19 dp_rank=parallel_states.get_data_parallel_rank(),
20      **dataset_kargs)
21 g = torch.Generator()
22 train_dataloader = DataLoader(train_dataset,
23      collate_fn=data_collator,
24      shuffle=False,
25      drop_last=True,
26      batch_size=args.batch_size_per_gpu,
27      generator=g)

```

## 模型训练

本仓库仅允许使用deepspeed进行训练，所以首先需要将模型包装为deepspeed引擎。（会比Huggingface的accelerate库更快）

```

1 optimizer, lr_scheduler = get_optimizer(ds_config, args, model=model)
2 engine, optimizer, _, _ = deepspeed.initialize(model=model,
3      optimizer=optimizer,
4      lr_scheduler=lr_scheduler,
5      config=ds_config,
6      model_parameters=[p for p in
7      model.parameters() if p.requires_grad],
8      mpu=None if args.num_pp_stages
9      else parallel_states)

```

包装模型之后，需要准备forward\_step/backward\_step/eval\_step，以及特定任务需要的task\_print(训练默认只打印loss的信息，如果需要打印mcc等信息则需要准备)。

```

1 if args.num_pp_stages:
2     forward_step = forward_step_pipeline
3     eval_step = eval_step_pipeline
4     backward_step = None
5     task_print = task_print_pipeline
6 else:
7     forward_step = forward_step_deepspeed
8     eval_step = eval_step_deepspeed
9     backward_step = backward_step_deepspeed
10    task_print = task_print_deepspeed
11
12 trainer = Trainer(args, writer)
13 trainer.register_task_print(task_print)
14 args.gradient_accumulation_steps = 1 # For correctly print info.

```

```

15
16     try:
17         trainer.train(model=engine,
18                       train_data_loader=train_data_loader_iter,
19                       eval_data_loader=eval_data_loader_iter,
20                       optimizer=None,
21                       forward_step=forward_step,
22                       backward_step=backward_step,
23                       eval_step=eval_step,
24                       log_loss=True)
25     except Exception as e:
26         # When any error occurs during the training process, log the error.
27         traceback_info = traceback.format_exc()
28         print_rank_0(traceback_info, args.global_rank, logging.ERROR)

```

## Trainer代码

huggingface模型训练一般使用huggingface的trainer库进行训练，同样的，本仓库也有自己的trainer代码，便于需要时进行修改。

文件地址 [MyTransformers/train/trainer.py](#)

## trainer的初始化

初始化时需要传入args，args即是在训练代码中获取的参数。以及writer，该变量用于绘制tensorboard信息。

```

1     def __init__(self, args, writer=None):
2         self.args = args
3         self.end = False
4         self.writer = writer
5         self.all_loss = 0.0
6         self.global_step = 0
7         self.all_metric = []
8         self.eval_loss = None
9         self.eval_metric = []
10        self.best_eval_index = 0.0
11        self.wait = 0
12        self.save_floder = os.path.join(args.output_path, args.experiment_name)
13        ensure_directory_exists(self.save_floder, self.args.global_rank)
14

```

## 训练过程

训练时需要传入的必要参数有

model: 被训练的模型

forward\_step: 模型前向传播的过程

backward\_step: 模型反向传播的过程 (deepspeed将反向传播和前向传播封装在了一起, 所以这个参数可以不提供)

optimizer: 优化器 (使用deepspeed训练时不需要提供)

train\_data\_loader和eval\_data\_loader: 数据集

eval\_step: 模型进行评估时的过程

```
1 def train(self,
2           model:torch.nn.Module,
3           train_data_loader:DataLoader,
4           forward_step: Callable,
5           optimizer: Callable = None,
6           eval_data_loader:DataLoader = None,
7           backward_step: Callable = None,
8           eval_step: Callable = None,
9           log_loss: bool = False):
```

具体的训练代码可以简单拆分为, 前向传播->反向传播->在特定步数下跑验证集->打印训练信息->在特定步数下保存模型

```
1 with Timer(iterations=total_print) as timer:
2     for step in range(self.args.num_update_steps):
3         timer.average_time(entry='start')
4         # Execute the forward step of the model and calculate the loss and
metric.
5         loss, metric = forward_step(model, train_data_loader, self.args,
6                                     step)
7         # Accumulate loss and metric for gradient accumulation.
8         self.all_loss += loss.item()
9         self.all_metric.append(metric)
10        # Execute the backward step if provided (optional) to update model
parameters.
11        if backward_step:
12            backward_step(model, optimizer, loss)
13        # Manage and print training information.
14        # Evaluate the model at intervals specified by eval_interval.
15        if (step + 1) % self.args.eval_interval == 0 and eval_step is not
None:
16            assert eval_data_loader is not None, 'evaluation dataset can
not be None'
```

```

16         self.eval_loss, eval_metric = eval_step(model,
17           eval_data_loader, self.args, step)
18         self.eval_metric.append(eval_metric)
19         self.info_manager(step, timer, log_loss)
20         self.save_model(model, step)
21         if self.end:
22             print_rank_0("Early stopping triggered.",
self.args.global_rank)
23             break

```

## 模型权重保存

由于流水线并行时模型是分层保存的，而数据并行时模型是完整保存的。所以使用流水线并行训练模型时，模型的保存过程需要特殊处理。deepspeed为流水线并行模型提供了保存方法，即model.save\_checkpoint。至于流水线模型则直接使用torch\_save方法即可。

```

1     def save_model(self, model, step):
2         config_path = os.path.join(self.save_floder, 'config.json')
3         if not os.path.exists(config_path) and isinstance(self.args,
Namespace):
4             with open(config_path, 'w', encoding='utf-8') as f:
5                 print_rank_0(f'--->Saving training config at {step+1}th step
in {config_path}.', self.args.global_rank)
6                 json.dump(self.args.__dict__, f)
7
8         if self.args.num_pp_stages is not None:
9             if not self.end and (step+1) % self.args.save_interval == 0:
10                 print_rank_0(f'--->Start saving model at {step+1}th step in
{self.save_floder}.', self.args.global_rank)
11                 model.save_checkpoint(self.save_floder, tag=f'step_{step+1}')
12                 print_rank_0('--->Saved the model.', self.args.global_rank)
13             elif self.end and (step+1) % self.args.save_interval != 0:
14                 print_rank_0(f'--->Start saving model at final step in
{self.save_floder}.', self.args.global_rank)
15                 model.save_checkpoint(self.save_floder, tag='final')
16                 print_rank_0('--->Saved the model.', self.args.global_rank)
17             return
18
19         if self.args.global_rank <= 0:
20             if not self.end and (step+1) % self.args.save_interval == 0:
21                 save_path = os.path.join(self.save_floder,
f'step_{step+1}.ckpt')
22                 print_rank_0(f'--->Start saving model at {step+1}th step in
{save_path}.', self.args.global_rank)
23                 self.torch_save(model, save_path)

```

```

24     print_rank_0('--->Saved the model.', self.args.global_rank)
25     elif self.end:
26         save_path = os.path.join(self.save_folder, 'final.ckpt')
27         print_rank_0(f'--->Start saving model at final step in
{save_path}.', self.args.global_rank)
28         self.torch_save(model, save_path)
29         print_rank_0('--->Saved the model.', self.args.global_rank)

```

torch\_save方法将torch.save再次包装了一层，能够很方便的只保存可训练参数。（通过参数中的save\_trainable参数指定）

```

1     def torch_save(self, model:torch.nn.Module, save_path):
2         if self.args.save_trainable:
3             trainable_params = {}
4             for name, param in model.module.named_parameters():
5                 if param.requires_grad:
6                     trainable_params[name] = param.data
7             torch.save(trainable_params, save_path)
8         else:
9             torch.save(model.module.state_dict(), save_path)

```

## 模型代码

本处使用llama的代码为例

代码地址 [MyTransformers/model/llama](#)

本仓库的模型需要提供四个文件：

## 模型配置

即config.py文件。该文件里需要提供模型的基本参数信息以dataclass类进行封装，形式如下：

```

1 @dataclass
2 class ModelArgs:
3     dim: int = 4096
4     head_dim: int = 128
5     hidden_size: int = 4096
6     n_layers: int = 32
7     num_hidden_layers: int = 32
8     n_heads: int = 32
9     num_attention_heads: int = 32
10    n_kv_heads: Optional[int] = None
11    vocab_size: int = -1 # defined later by tokenizer

```

```

12     multiple_of: int = 256 # make SwiGLU hidden layer size multiple of large
    power of 2
13     ffn_dim_multiplier: Optional[float] = None
14     norm_eps: float = 1e-5
15     rope_theta: float = 100000.0
16     max_batch_size: int = 32
17     max_seq_len: int = 2048
18     tokenizer: Optional[str] = ''
19     lora_layers: List[str] = field(default_factory=lambda: ['wk', 'wv', 'wq',
    'wo', 'w1', 'w2', 'w3'])
20     dtype: str = 'float16'
21     def get_dtype(self) -> Optional[torch.dtype]:
22         """Gets the torch dtype from the config dtype string."""
23         return STR_DTYPE_TO_TORCH_DTYPE.get(self.dtype, None)

```

还需要提供不同大小的模型需要的不同参数配置，并使用注册器进行注册，形式如下：

```

1 @registry.register_model_config("llama3_8b")
2 def get_config_for_llama3_8b() -> ModelArgs:
3     return ModelArgs(
4         n_kv_heads=8,
5         multiple_of=1024,
6         ffn_dim_multiplier=1.3,
7         vocab_size=128256,
8         rope_theta=500000.0
9     )

```

## 模型

即model.py文件。该文件里需要提供模型的定义代码和inference代码，一般来说从开源代码中copy并稍微进行改动即可，最后要使用注册器进行注册，形式如下：

```

1 @registry.register_model("llama3")
2 class Llama3(LlamaGenerate):
3     def __init__(self, model_args: ModelArgs):
4         try:
5             self.tokenizer = Llama3Tokenizer(model_args.tokenizer)
6             model_args.vocab_size = self.tokenizer.n_words
7         except:
8             pass
9         super().__init__(model_args=model_args)

```



## 模型数据并行类

即train\_model.py文件，由于开源模型代码一般仅提供定义代码和inference代码，所以要把模型用于训练的话还需要另外写一个训练类，将模型训练时的数据流动进行包装。

需要注意的是，不同模型的数据并行类都必须只接受model和args两个参数，如果需要修改，则要把每个模型的代码都进行修改。并在u\_train代码中做对应的修改。

该类需要继承BaseModel类，并使用注册器进行注册，形式如下：

```
1 @registry.register_train_model('llama3')
2 @registry.register_train_model('llama')
3 class LLaMaTrainModel(BaseModel):
4     """
5     Trainer class for llama, responsible for handling input and output during
6     training.
7     """
8     def __init__(self, model:LlamaGenerate, args):
9         """
10         Initializes basic attributes for the trainer class and precomputes
11         fixed values.
12
13         param model: llama model with pretrained weight.
14         param args: Arguments from argument parser.
15         """
16         super().__init__(args)
17         self.layers = model.model.layers
18         self.tok_embeddings = model.model.tok_embeddings
19         self.output = model.model.output
20         self.norm = model.model.norm
21         self.attention_mask = LLaMaTrainModel.get_masks(args.max_len)
22         self.freqs_cis = precompute_freqs_cis(args.head_dim,
23                                             args.max_len,
24                                             theta=args.rop_theta,
25                                             train_pi=args.train_pi,
26                                             train_pipeline=False)
```

前向过程直接继承BaseModel类即可：

```
1 def forward(self, **kwargs):
2     return super().forward(**kwargs)
```

但是该类需要额外提供三个方法：

embedding方法，用于语言模型的嵌入，该方法会在forward方法中被调用

```
1 def embedding(self, input_ids):
2     hidden_states = self.tok_embeddings(input_ids)
3     attention_mask =
4     self.attention_mask.to(hidden_states.device).to(hidden_states.dtype)
5     return hidden_states, attention_mask
```

model\_forward方法，用于模型的前向传播过程，该方法同样会在forward方法中被调用

```
1 def model_forward(self, logits, freqs_cis, attention_mask):
2     # Using activation checkpoint to reduce memory consumption or not.
3     for i in range(self.args.num_layers):
4         if self.args.activation_checkpoint:
5             logits = checkpoint(self.layers[i],
6                                 logits,
7                                 0,
8                                 freqs_cis,
9                                 attention_mask,
10                                self.args.attn_type,
11                                use_reentrant=False)
12         else:
13             logits = self.layers[i](x=logits,
14                                     start_pos=0,
15                                     freqs_cis=freqs_cis,
16                                     mask=attention_mask,
17                                     attn_type=self.args.attn_type)
18     logits = self.norm(logits)
19     logits = self.output(logits)
20     return logits
```

get\_masks静态方法，用于获取attention mask，该方法会在多个地方被使用（包括类的外部），所以需要写为静态方法

```
1 @staticmethod
2 def get_masks(seqlen, device='cpu', dtype=torch.float, start_pos=0):
3     if seqlen > 1:
4         mask = torch.full((seqlen, seqlen), float("-inf"), device=device)
5         mask = torch.triu(mask, diagonal=1)
6         mask = torch.hstack([torch.zeros((seqlen, start_pos),
7                                         device=device), mask]).to(dtype)
8     return mask
```

## 模型流水线并行类

即pipeline\_model.py，该类专门用于流水线并行训练

需要编写模型中每一个逻辑层的代码，并封装在一起，最后使用注册器注册，形式如下：

```
1 @registry.register_pipeline_model("llama3")
2 @registry.register_pipeline_model("llama")
3 def get_pipeline_model(model, args):
4     layers = [LayerSpec(EmbeddingPipelineLayer, model=model, args=args),
5               *[LayerSpec(DecoderPipelineLayer, model=model, args=args,
6                           layer_idx=idx) for idx in
7                 range(args.num_layers)],
8               LayerSpec(FNormPipelineLayer, model=model),
9               LayerSpec(LossPipelineLayer, pad_id=args.pad_id)]
10    return PipelineModule(layers=layers, num_stages=args.num_pp_stages,
11                          partition_method='uniform')
```

逻辑层的代码形式如下，模型的每一个逻辑部分都需要写成这样的逻辑层，其输入是一个模型，并在这个模型中提取需要的部分，然后处理该部分的输入输出。其forward函数需要遵守deepspeed的协议，仅接受一个inputs输入元组，进入具体逻辑后再进行拆分：

```
1 class DecoderPipelineLayer(torch.nn.Module):
2     def __init__(self, model: LlamaGenerate, layer_idx, args):
3         super().__init__()
4         self.layer = model.model.layers[layer_idx]
5         self.args = args
6
7     def forward(self, inputs):
8         hidden_states, freqs_cis, attention_mask, labels = inputs
9         # [batch_size, input_len, hidden_dim]
10        if self.args.activation_checkpoint:
11            hidden_states = checkpoint(self.layer,
12                                       hidden_states,
13                                       0,
14                                       freqs_cis,
15                                       attention_mask,
16                                       self.args.attn_type,
17                                       use_reentrant=False)
18        else:
19            hidden_states = self.layer(hidden_states,
20                                       0,
21                                       freqs_cis,
```

```

22         attention_mask,
23         self.args.attn_type)
24     return hidden_states, freqs_cis, attention_mask, labels

```

## 通用代码

由于每个模型都有共通的部分，所以这些内容可以单独抽象出来

## 模型数据并行基类

即base\_model.py文件，该文件提供BaseModel类，该类提供了一个模型训练的基本流程，不同的模型只需要编写模型特定的子流程即可。

```

1     def forward(self, **kwargs):
2         input_ids, labels = kwargs["input_ids"], kwargs["labels"]
3         input_ids, labels, freqs_cis = self.cut_sequence(input_ids, labels)
4         hidden_states, attention_mask = self.embedding(input_ids)
5         logits = self.model_forward(hidden_states, freqs_cis, attention_mask)
6         loss = self.compute_loss(logits, labels)
7         return loss, self.compute_metric(logits, labels,
            kwargs["cal_metric_pos_tensor"])

```

具体的模型需要完成编写子流程，否则会导致NotImplementedError。

```

1     def embedding(self, input_ids):
2         raise NotImplementedError()
3
4     def model_forward(self, hidden_states, freqs_cis, attention_mask):
5         raise NotImplementedError()

```

该文件还提供了一些每个模型都要使用的子流程，如loss的计算。具体的模型如果需要修改loss的计算，那么直接覆盖这个方法即可。

```

1     def compute_loss(self, logits, labels):
2         shift_logits = logits[..., :-1, :].contiguous()
3         shift_labels = labels[..., 1:].contiguous()
4         loss = self.loss_fct(shift_logits.reshape(-1, shift_logits.size(-1)),
            shift_labels.reshape(-1))
5         return loss

```

## 注意力机制代码

即attention.py，在模型的训练过程中，可能会需要使用到flash\_attention等代码，而不同模型的注意力机制基本是相同的，所以把这部分抽象出来，方便对注意力机制进行不同的处理。

```
1 def attention_func(  
2     q: torch.Tensor,  
3     k: torch.Tensor,  
4     v: torch.Tensor,  
5     atten_mask: torch.Tensor,  
6     dropout_p: float,  
7     scaling: float,  
8     is_causal: bool,  
9     atten_type: str = ''  
10 ):  
11     """  
12     Attention function that supports different attention types.  
13  
14     Args:  
15         q (torch.Tensor): Query tensor of shape [batch_size, n_local_heads,  
input_len, head_dim].  
16         k (torch.Tensor): Key tensor of shape [batch_size, n_local_heads,  
input_len, head_dim].  
17         v (torch.Tensor): Value tensor of shape [batch_size, n_local_heads,  
input_len, head_dim].  
18         atten_mask (torch.Tensor): Attention mask tensor of shape [batch_size,  
1, input_len, input_len].  
19         dropout_p (float): Dropout probability.  
20         scaling (float): Scaling factor for the attention scores.  
21         is_causal (bool): Whether the attention is causal (only attend to past  
tokens).  
22         atten_type (str, optional): Type of attention to use. Can be  
'flash_attn', 'ulysses_flash_attn', 'ulysses_attn', or leave empty for the  
default naive_attention_func.  
23  
24     Returns:  
25         torch.Tensor: Output tensor of shape [batch_size, n_local_heads,  
input_len, head_dim].  
26     """  
27     if atten_type == 'flash_attn':  
28         require_version("torch>=2.0.0")  
29         with torch.backends.cuda.sdp_kernel(enable_flash=True):  
30             output = F.scaled_dot_product_attention(q, k, v,  
attn_mask=atten_mask, dropout_p=dropout_p, is_causal=is_causal)  
31     elif atten_type == 'ulysses_flash_attn':  
32         require_version("torch>=2.0.0")
```

```

33         with torch.backends.cuda.sdp_kernel(enable_flash=True):
34             flash_attn = F.scaled_dot_product_attention
35             dist_attn = DistributedAttention(flash_attn,
parallel_states.get_sequence_parallel_group(), scatter_idx=1, gather_idx=2)
36             output = dist_attn(q, k, v, attn_mask=attn_mask,
dropout_p=dropout_p, is_causal=is_causal)
37         elif atten_type == 'ulysses_attn':
38             dist_attn = DistributedAttention(naive_attention_func,
parallel_states.get_sequence_parallel_group(), scatter_idx=1, gather_idx=2)
39             output = dist_attn(q, k, v, attn_mask=attn_mask,
dropout_p=dropout_p, scaling=scaling, is_causal=is_causal)
40         else:
41             output = naive_attention_func(q, k, v, attn_mask=attn_mask,
dropout_p=dropout_p, scaling=scaling, is_causal=is_causal)
42         return output

```

## Common库

### Lora

因为peft库用起来有点麻烦，而且也不得心应手，为了方便修改和使用，我自己写了一个lora的代码。

文件地址 [MyTransformers/common/lora.py](#)。

该库中提供了LinearWithLoRA类，可以将任何Linear层封装到这个类中，就能实现使用Lora训练。该类同时兼容了dora和plora的代码，如果有需要，可以将其他的lora变体的逻辑补充到这个类中。

```

1 class LinearWithLoRA(nn.Linear):
2     def __init__(
3         self,
4         in_features: int,
5         out_features: int,
6         lora_rank: int = 4,
7         lora_scaler: float = 32.0,
8         use_dora: bool = False,
9         quant: bool = False,
10        plora_steps: Union[int, None] = None
11    ):

```

具体模型训练时，直接调用该文件中提供的switch\_to\_lora函数即可，需要提供以下参数：

model: 需要使用lora训练的模型

replace\_names: 需要使用lora替换的层，如['wq','wv','wo']表示将模型中所有名字包含这三个字符串的线性层替换为lora层。

rank: lora的秩

lora\_scaler: lora的缩放因子

transposition: torch的线性层代码如linear\_1(x)的实际逻辑为 $x * \text{linear\_1.weight.T}$ ，但是有些模型的代码的线性层逻辑可能写为 $x * \text{linear\_1.weight}$ ，在这种时候，参数矩阵就需要转置，才能放入lora层。

该代码首先找到需要被替换的层，确保该层具有in\_features和out\_features两个参数，确保lora层的维度和原来的层的维数相同，再把相关的参数带入LinearWithLoRA类中构成一个lora层的实例。然后复制原来的参数到lora层中。

最后需要找到这个层的父层，再进行替换

```
1 def switch_to_lora(model:nn.Module,
2                     replace_names:Optional[Union[str, List[str]]] = None,
3                     rank:int = 4,
4                     lora_scaler:int = 32,
5                     transposition:bool = False,
6                     use_dora:bool = False,
7                     plora_steps: Optional[int] = None):
8     """
9     Switch function for lora, responsible for replacing Linear layer with
10    LinearWithLoRA layer
11
12    param model: Any pytorch model.
13    param replace_names: List of module names to be replaced by LoRA.
14    param rank: Rank for LoRA.
15    param lora_scaler: Scaler for LoRA.
16    param transposition: nn.Linear(x, w) compute  $xw^T$ , so the weight should in
17    shape [out_feature, in_feature]. Otherwise, transposition should be set to True
18    param use_dora: Weather to use dora
19    param plora_steps: The steps to merge and reset lora weight.
20    """
21    assert replace_names is not None, 'Replace names can not be None'
22    for name, module in model.named_modules():
23        replace_tag = False
24        for replace_name in replace_names:
25            if replace_name in name:
26                # Create LoRA layer instance.
27                replace_tag = True
28                if isinstance(module, LinearWithLoRA):
29                    module.merge_and_reset(new_rank=rank)
30                elif isinstance(module, nn.Module):
31                    assert all(hasattr(module, attr) for attr in
32                               ["in_features", "out_features", "weight"]), \
```

```

30         "Module is missing one or more of the required attributes:
    'in_features', 'out_features', 'weight'"
31         quant = getattr(module, "quant", False)
32         lora_layer = LinearWithLoRA(lora_rank=rank,
33                                     lora_scaler=lora_scaler,
34
35         in_features=module.in_features,
36
37         out_features=module.out_features,
38
39         use_dora=use_dora,
40         quant=quant,
41         plora_steps=plora_steps)
42
43         # Copy the original weight to the LoRA layer.
44         if transposition:
45             lora_layer.weight = nn.Parameter(module.weight.data.T)
46         else:
47             lora_layer.weight.data = module.weight.data
48         if quant:
49             lora_layer.weight_scaler = module.weight_scaler
50         # Replace the original layer with the LoRA layer.
51         parent = get_parent_model(model, module)
52         setattr(parent, list(parent._modules.items())[0], lora_layer)
53     [list(parent._modules.values()).index(module)][0], lora_layer)
54
55     if not replace_tag and isinstance(module, LinearWithLoRA):
56         # Merge weight to avoid unnecessary computing.
57         module.merge_and_del()

```

找父层的过程基本就是一个深度遍历的过程（比二叉树的遍历稍微简单一些），主要是需要找到该层的直接父层。

```

1 def get_parent_model(parent_model, module):
2     """
3     Find the parent module for the input module recursively.
4
5     param parent_model: Root model for the search.
6     param module: Submodule to find the parent module for.
7
8     Returns:
9     Parent module if found, None otherwise.
10    """
11    for _, sub_module in parent_model._modules.items():
12        # Direct sub modules of parent model.
13        if sub_module is module:
14            return parent_model
15    parent = get_parent_model(sub_module, module)

```



```
16         if parent:
17             return parent
18     return None
```

## Parser

parser代码用于准备模型训练过程中所需要的参数，如学习率等。parser中有三个子parser：

- base\_parser
- train\_parser
- ds\_parser

所以使用时需要这样使用：ds\_parser(train\_parser(base\_parser()))  
.get\_args()

## 注册器

注册器是本库中非常重要的特性，考虑到本库的目标是兼容多个不同的模型，所以为了避免适应不同的模型而在模型训练过程中产生过多冗余的代码，我采用了注册器机制。

如以下代码中，获取tokenizer，模型参数，模型定义，模型训练类，都需要通过注册器来完成。使用一个注册器来适应多个不同的模型，而不需要为每个模型准备一段代码。（但是每个模型的接口都需要一致才能使用）

```
1 def load_local_model(args):
2     # Train with local defined model.
3     tokenizer = registry.get_tokenizer_class(args.tokenizer_name)
4     (args.tokenizer_path)
5     print_rank_0(f'--->Using tokenizer: {args.tokenizer_name} with path:
6     {args.tokenizer_path}', args.global_rank)
7     config_type = '_'.join([args.model_name, args.variant])
8     model_config = registry.get_model_config_class(config_type)()
9     print_rank_0(f'--->Using model config: {config_type}', args.global_rank)
10    model_config.vocab_size = tokenizer.n_words
11    model = registry.get_model_class(args.model_name)(model_config)
12    print_rank_0(f'--->Using model: {args.model_name}, and loading its
13    training variant', args.global_rank)
14
15    # Load checkpoint if checkpoint path is provided.
16    if args.ckpt_path is not None and args.from_pretrained:
17        load_ckpt(model=model.model, ckpt_path=args.ckpt_path,
18        partial_ckpt_path=args.partial_ckpt_path, rank=args.global_rank)
19        print_rank_0(f'--->Using pretrained checkpoint at {args.ckpt_path}',
20        args.global_rank)
21    else:
22        print_rank_0('--->Not using pretrained checkpoint to start training.',
23        args.global_rank)
```

```

18
19     # Load config from model config to argument parser namespace.
20     args.head_dim = model_config.head_dim
21     args.head_num = model_config.num_attention_heads
22     args.hidden_size = model_config.hidden_size
23     args.num_layers = model_config.num_hidden_layers
24     args.rope_theta = model_config.rope_theta if args.rope_theta is None else
args.rope_theta
25     args.pad_id = tokenizer.pad_id
26
27     # Load model to training dtype.
28     if args.fp16:
29         model.half()
30     elif args.bf16:
31         model.bfloat16()
32
33     # Convert model to trainable model for given training type.
34     if args.num_pp_stages:
35         pipe_model_cls = registry.get_pipeline_model_class(args.model_name)
36         model = pipe_model_cls(model, args)
37     else:
38         train_model_cls = registry.get_train_model_class(args.model_name)
39         model = train_model_cls(model, args)
40     return model, tokenizer, model_config
41

```

注册器本身是一个包含了一个映射字典，和多个注册与取出的静态装饰器方法的类：

```

1 class Registry:
2     mapping = {
3         "model_mapping": {},
4         "pipeline_model_mapping": {},
5         "train_model_mapping": {},
6         "model_config_mapping": {},
7         "dataset_mapping": {},
8         "info_manager_mapping": {},
9         "tokenizer_mapping": {},
10        "paths_mapping": {}
11    }
12

```

注册装饰器如下：

```

1  @classmethod
2  def register_model(cls, name):
3      def wrap(model_cls):
4          if model_cls in cls.mapping['model_mapping']:
5              raise KeyError(
6                  "Name '{}' already registered for {}".format(
7                      name, cls.mapping["model_mapping"][name]
8                  )
9              )
10         cls.mapping['model_mapping'][name] = model_cls
11         return model_cls
12     return wrap
13

```

取定义的静态方法如下：

```

1  @classmethod
2  def get_model_class(cls, name):
3      result = cls.mapping["model_mapping"].get(name, None)
4      if result is None:
5          raise ValueError(f"Can not find name: {name} in model mapping, \
6 supported models are listed below:{cls.list_models()}")
7      else:
8          return result

```

## Utils

utils库中包含了许多有用的函数和类

Timer: 这是一个用于计时的环境，通常用在循环上面，可以用于计算单次循环的用时，和整个循环的用时。

```

1  class Timer(object):
2      def __init__(self, start=None, n_round=2, iterations: Optional[int] =
3          None):

```

logging配置：本库中的logging具有两个handler，一个是文件handler一个是控制台handler。文件handler负责将日志信息输出到文件中，而控制台handler负责将日志信息打印到控制台

handler的行为通过环境变量来控制，PRINTLEVEL负责控制控制台handler的级别，如设置为DEBUG，则级别低于DEBUG的日志不会被输出到控制台中

LOGLEVEL负责控制文件handler的级别，设置为INFO，则级别低于INFO的信息不会被写入文件

NO\_LOG\_FILE用于负责文件的创建，当测试代码时，可以通过这个变量避免产生文件

日志文件会根据当前时间进行创建，文件名即为时间。

```
1 def configure_logging(log_path, rank: Optional[int]=0):
2     """
3     Configure logging functionality.
4
5     Args:
6     log_path (str): Path where the log files will be stored.
7     rank (optional[int]): Level used for creating directories, default is 0.
8
9     Returns:
10    loggegr (logging.Logger): Configured logger object.
11    """
12    # Get environment variables.
13    sh_level = os.environ.get("PRINTLEVEL", logging.DEBUG)
14    fh_level = os.environ.get("LOGLEVEL", logging.INFO)
15    fh_disable = os.environ.get("NO_LOG_FILE", "false") == 'true' # Convert
16    string variable to boolean
17    logger = logging.getLogger("MyTransformers")
18    logger.setLevel(logging.DEBUG)
19    formatter = logging.Formatter('%(asctime)s [%(levelname)s] %(message)s')
20    # Define the log format
21    # Create a console log handler
22    sh = logging.StreamHandler()
23    sh.setFormatter(formatter)
24    sh.setLevel(sh_level)
25    logger.addHandler(sh)
26
27    timezone = pytz.timezone('Asia/Shanghai')
28    # Get current date and time strings, formatted as 'yy-mm-dd' and 'HH-
29    MM.log'
30    date_string, hour_string = datetime.now().strftime('%y-%m-%d'),
31    datetime.now().strftime('%H-%M') + '.log'
32    log_path = os.path.join(log_path, date_string)
33    ensure_directory_exists(log_path, rank)
34    if not fh_disable:
35        fh = logging.FileHandler(os.path.join(log_path, hour_string))
36        fh.setFormatter(formatter)
37        fh.setLevel(fh_level)
38        logger.addHandler(fh)
```

```
38
39     return logger
```

打印函数：本库中的print\_rank\_0会自动使用日志的形式进行打印，当log的级别合适时，print的内容会被输出到控制台，并写入文件中。

```
1 def print_rank_0(msg, rank=0, level=logging.INFO, flush=True):
2     if "logger" not in globals() and rank<=0:
3         global logger
4         # Create logger when print_rank_0 being used.
5         logger = configure_logging(os.environ.get("LOG_FOLDER", "log"), rank)
6         if isinstance(level, str):
7             level = getattr(logging, level.upper(), logging.INFO)
8         if rank <= 0:
9             logger.log(msg=msg, level=level)
10        if flush:
11            logger.handlers[0].flush()
```

加载模型的函数，load\_ckpt可以结合预训练权重和后训练权重进行加载。后训练权重会被优先加载到模型中（覆盖掉预训练权重）。这个函数可以配合保存可训练参数的函数使用，这样训练时就能够只保存可训练函数，并正常加载模型。

并且该函数能够支持加载具有unexpected keys和missing keys的权重而不报错，但是会将相关的信息打印出来。

```
1 def get_merged_state_dict(ckpt_path:str=None,
2                             partial_ckpt_path:str=None):
3     if ckpt_path:
4         ckpt = torch.load(ckpt_path, map_location='cpu')
5     else:
6         ckpt = {}
7
8     if partial_ckpt_path:
9         partial_ckpt = torch.load(partial_ckpt_path, map_location='cpu')
10
11        for name, param in partial_ckpt.items():
12            ckpt[name] = param
13    return ckpt
14
15 def load_ckpt(model:Module,
16               ckpt_path:str=None,
17               partial_ckpt_path:str=None,
18               rank:int=0):
```

```

19     """
20     load model checkpoint safely.
21
22     param model: Any pytorch model.
23     param ckpt_path: Path of checkpoint of all params.
24     param partial_ckpt_path: Path of partial model checkpoint. Must be
provided if load trainable params and pretrained params.
25     """
26
27     rank_0 = rank == 0
28     model_state_dict = model.state_dict()
29     ckpt = get_merged_state_dict(ckpt_path, partial_ckpt_path)
30
31     try:
32         model.load_state_dict(ckpt, strict=False)
33     except Exception:
34         if rank_0:
35             print_exc()
36
37     missing_keys = [key for key in model_state_dict.keys() if key not in
ckpt.keys()]
38     unexpected_keys = [key for key in ckpt.keys() if key not in
model_state_dict.keys()]
39
40     print_rank_0(f'--->Missing keys:{missing_keys}. Unexpected keys:
{unexpected_keys}.', rank=rank)

```

## 初始化分布式环境

```

1 def init_dist(args):
2     if args.device == 'cuda':
3         if args.local_rank == -1:
4             device = torch.device("cuda")
5         else:
6             torch.cuda.set_device(args.local_rank)
7             device = torch.device("cuda", args.local_rank)
8             deepspeed.init_distributed(dist_backend="nccl")
9             args.world_size = dist.get_world_size()
10            args.global_rank = dist.get_rank()
11        else:
12            device = 'cpu'
13        if args.num_sp_stages:
14            assert args.attn_type == 'ulysses_attn', 'when using sequence
parallism, the attention type must be `ulysses_attn`'

```

```

15 parallel_states.initialize_model_parallel(sequence_model_parallel_size=args.num
    _sp_stages)
16     elif args.num_pp_stages:
17
18 parallel_states.initialize_model_parallel(pipeline_model_parallel_size=args.num
    _pp_stages)
19 else:
20     parallel_states.initialize_model_parallel()
21     return device, args

```

## 随机种子

```

1 def set_random_seed(seed):
2     if seed is not None:
3         random.seed(seed)
4         np.random.seed(seed)
5         torch.manual_seed(seed)
6         torch.cuda.manual_seed_all(seed)

```

## 优化器

由于在一些情况下需要对优化器进行修改，使用deepspeed提供的优化器功能不太完整，需要自己写一些补充功能。

所以本库还提供了一个优化器的代码，地址 [MyTransformers/common/optimizer.py](#)

主要是一个获取优化器的函数，能够使用args获取不同的优化器。

```

1 def get_optimizer(ds_config, args, model):
2     # TODO: Add support to PLoRA and so no ...
3     if args.diy_optimizer:
4
5         # Acquire optimizer type.
6         if args.optim_type is not None:
7             optim_type = args.optim_type.lower()
8         elif 'optimizer' in ds_config:
9             optim_type = ds_config['optimizer'].get('type', 'adamw').lower()
10
11         # Acquire optimizer.
12         if args.use_galore:
13             message = 'galore cannot be used with the current DeepSpeed
                version, and running it will result in an error.'

```

```

14         warnings.warn(message, UserWarning)
15         isSuccess, optimizer = get_galore_optimizer(optim_type, args,
model)
16     else:
17         isSuccess, optimizer = get_regular_optimizer(optim_type, args,
model)
18
19     # If successful acquire optimizer, overwrite the default setting in
the ds config.
20     if isSuccess:
21         del ds_config['optimizer']
22         print_rank_0(F'--->deepspeed optimizer setting have been
overwritten', args.global_rank)
23     else:
24         print_rank_0(f'--->try to use diy optimizer failed, use the ds
setting', args.global_rank)
25
26     lr_scheduler = get_learning_rate_scheduler(optimizer, 0, args)
27     return optimizer, lr_scheduler
28
29     return None, None

```

## 学习率规划器

这里的学习率规划器代码主要是从SwissArmyTransformer库中copy的，SAT库又是从英伟达那copy的

现在一般是使用cosine的方式对学习率进行规划，如果有新的学习率规划的方式也可以在这个代码上进行补充修改。

```

1 class AnnealingLR(_LRScheduler):
2     """Anneals the learning rate from start to zero along a cosine curve."""
3
4     DECAY_STYLES = ['linear', 'cosine', 'exponential', 'constant', 'None']
5
6     def __init__(self, optimizer, start_lr, warmup_iter, num_iters,
decay_style=None, last_iter=-1, decay_ratio=0.5, auto_warmup_steps=50,
auto_warmup_rate=0.05):
7         """
8         params start_lr: base learning rate for the scheduler
9         params num_iters: total steps
10        params warmup_iter: how much iters to use warmup
11        params last_iters: how much iters already trained
12        params auto_warmup_steps: the fix warmup step
13        params auto_warmup_rate: the lr radio for fix warmup step
14        params decay_style: learning rate decay style after warmup

```



```

15         """
16         assert warmup_iter <= num_iters
17         self.optimizer = optimizer
18         self.lr_scale = deepcopy([x['lr'] if 'lr' in x else 1. for x in
optimizer.param_groups])
19         self.start_lr = start_lr
20         self.warmup_iter = warmup_iter
21         self.init_step = last_iter
22         self.num_iters = last_iter + 1
23         self.end_iter = num_iters
24         self.decay_style = decay_style.lower() if isinstance(decay_style, str)
else None
25         self.decay_ratio = 1 / decay_ratio
26         self.auto_warmup_steps = auto_warmup_steps
27         self.auto_warmup_rate = auto_warmup_rate
28         self.step(self.num_iters)
29         if not torch.distributed.is_initialized() or
torch.distributed.get_rank() == 0:
30             print_rank_0(f'--->learning rate decaying style
{self.decay_style}, ratio {self.decay_ratio}')
31
32     def get_lr(self):
33         # auto_warmup_steps并不取决于warmup的设置，而是固定的进行warmup
34         if self.num_iters <= self.init_step + self.auto_warmup_steps:
35             auto_lr = float(self.start_lr) * self.auto_warmup_rate
36             scheduled_lr = float(self.start_lr) * self.num_iters /
self.warmup_iter
37             return min(auto_lr, scheduled_lr)
38
39         # 根据warmup设置进行warmup
40         if self.warmup_iter > 0 and self.num_iters <= self.warmup_iter:
41             return float(self.start_lr) * self.num_iters / self.warmup_iter
42         else:
43             # 进行learning rate decay
44             if self.decay_style == self.DECAY_STYLES[0]:
45                 return self.start_lr*((self.end_iter-(self.num_iters-
self.warmup_iter))/self.end_iter)
46             elif self.decay_style == self.DECAY_STYLES[1]:
47                 decay_step_ratio = min(1.0, (self.num_iters -
self.warmup_iter) / self.end_iter)
48                 return self.start_lr / self.decay_ratio * (
49                     (math.cos(math.pi * decay_step_ratio) + 1) *
(self.decay_ratio - 1) / 2 + 1)
50             elif self.decay_style == self.DECAY_STYLES[2]:
51                 return self.start_lr
52             else:
53                 return self.start_lr

```

```
54
55     def step(self, step_num=None):
56         if step_num is None:
57             step_num = self.num_iters + 1
58         self.num_iters = step_num
59         new_lr = self.get_lr()
60         for group, scale in zip(self.optimizer.param_groups, self.lr_scale):
61             group['lr'] = new_lr * scale
62
63     def state_dict(self):
64         sd = {
65             'start_lr': self.start_lr,
66             'warmup_iter': self.warmup_iter,
67             'num_iters': self.num_iters,
68             'decay_style': self.decay_style,
69             'end_iter': self.end_iter,
70             'decay_ratio': self.decay_ratio
71         }
72         return sd
```