

vlm

pretrain

1.tokenizer

2.dataset

基于 llava，原理文章很多，且较容易。本文主要负责代码调试，明晰细节。

本文项目基于 llava: <https://github.com/dvlab-research/MGM>

pretrain

1.tokenizer

断点位于以下代码处: `if model_args.version == "v0":`

查看 tokenizer 的 special token，由于 unk 和 pad 都有，因此不需要做什么处理

```
tokenizer.special_tokens_map
{
  'bos_token': '<bos>', 'eos_token': '<eos>', 'unk_token': '<unk>', 'pad_token': '<pad>', 'additional_special_tokens': ['<start_of_turn>', '<end_of_turn>']}
> special variables
> function variables
'bos_token' = '<bos>'
'eos_token' = '<eos>'
'unk_token' = '<unk>'
'pad_token' = '<pad>'
> 'additional_special_tokens' = ['<start_of_turn>', '<end_of_turn>']
1: /
```

除此之外，以及 dataset 部分，如果需要自定义模型几乎不需要修改什么

2.dataset

断点打在 dataset 的 `__getitem__` 函数

```
1 image.shape: torch.Size([3, 768, 768])
```

进入 `preprocess_multimodal` 函数

原数据: `[[{'from': 'human', 'value': 'What is in the photo?\n<image>'}, {'from': 'gpt', 'value': 'the box chain purse black'}]]`

返回: `[[{'from': 'human', 'value': '<image>\nWhat is in the photo?'}, {'from': 'gpt', 'value': 'the box chain purse black'}]]`

`preprocess_multimodal` 函数只将 `\n<image>` 移至最前为 `<image>\n`

进入 `preprocess` 函数

进入 `preprocess_gemma` 函数

1.应用模版

```

1 print(conv)
2 Conversation(system='', roles=('user', 'model'), messages=[], offset=0, sep_style=SeparatorStyle.GEMMA: 7>, sep='', sep2='<eos>', version='gemma', skip_next=False)
3
4 # 指定身份
5 print(roles)
6 {'human': 'user', 'gpt': 'model'}
7
8 print(source)
9 [{'from': 'human', 'value': '<image>\nWhat is in the photo?'}, {'from': 'gpt', 'value': 'the box chain purse black'}]
10
11 # conv.messages包含了纠正身份的对话
12 print(conv)
13 Conversation(system='', roles=('user', 'model'), messages=[['user', '<image>\nWhat is in the photo?'], ['model', 'the box chain purse black']], offset=0, sep_style=SeparatorStyle.GEMMA: 7>, sep='', sep2='<eos>', version='gemma', skip_next=False)

```

应用对话模版: `conversations.append(conv.get_prompt())`

进入: `elif self.sep_style == SeparatorStyle.GEMM`

A:

拼接后的对话为

```

1 print(ret)
2 <start_of_turn>user
3 <image>
4 What is in the photo?<end_of_turn>
5 <start_of_turn>model
6 the box chain purse black<end_of_turn>
7 <eos>

```

2.tokenize 部分: `if has_image:`

进入函数 `tokenizer_image_token`

```
Python |
1 print(prompt_chunks)
2 [[2, 106, 1645, 108], [2, 108, 1841, 603, 575, 573, 2686, 235336, 107, 108,
   106, 2516, 108, 1175, 3741, 9488, 45189, 2656, 107, 108, 1]]
3
4 tokenizer("<start_of_turn>user\n").input_ids
5 [2, 106, 1645, 108]
6
7 tokenizer("What is in the photo?<end_of_turn>\n<start_of_turn>model\nthe bo
   x chain purse black<end_of_turn>\n<eos>").input_ids
8 [2, 1841, 603, 575, 573, 2686, 235336, 107, 108, 106, 2516, 108, 1175, 3741
   , 9488, 45189, 2656, 107, 108, 1]
```

相当于在 `<image>` 处进行了拆分

通过 `insert_separator` 函数

```
Python |
1 input_ids
2 [2, 106, 1645, 108, -200, 108, 1841, 603, 575, 573, 2686, 235336, 107, 108,
   106, 2516, 108, 1175, 3741, 9488, 45189, 2656, 107, 108, 1]
```

相当于在 `prompt_chunks` 间插入了 `image token`

`id -200`, 同时删除第二个 `chunk` 的句首 token

返回的 shape 为: `torch.Size([1, 25])`, 复制一份作为 `targets`

3.mask

```

1 input_ids
2 tensor([[ 2, 106, 1645, 108, -200, 108, 1841, 603,
3          575,
4          573, 2686, 235336, 107, 108, 106, 2516, 108,
5          1175,
6          3741, 9488, 45189, 2656, 107, 108, 1]])
7 targets
8 tensor([[ -100, -100, -100, -100, -100, -100, -100, -100, -100, -
9          100,
10          -100, -100, -100, -100, -100, -100, -100, 1175, 3741, 9
11          488,
12          45189, 2656, 107, 108, 1]])
11 tokenizer("the box chain purse black<end_of_turn>\n<eos>").input_ids
12 [2, 1175, 3741, 9488, 45189, 2656, 107, 108, 1]

```

只保留 `the box chain purse black<end_of_turn>\n<eos>` 这一部分

插值将 `image` 变为 `torch.Size([3, 336, 336])`