

Unified Vision-Language Pre-Training for Image Captioning and VQA

LoongCat

2024 年 5 月 14 日

1 Structure

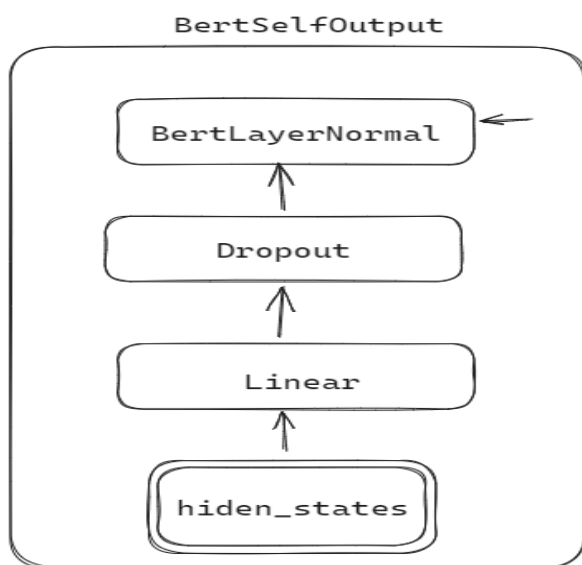


图 1: BertSelfOutput

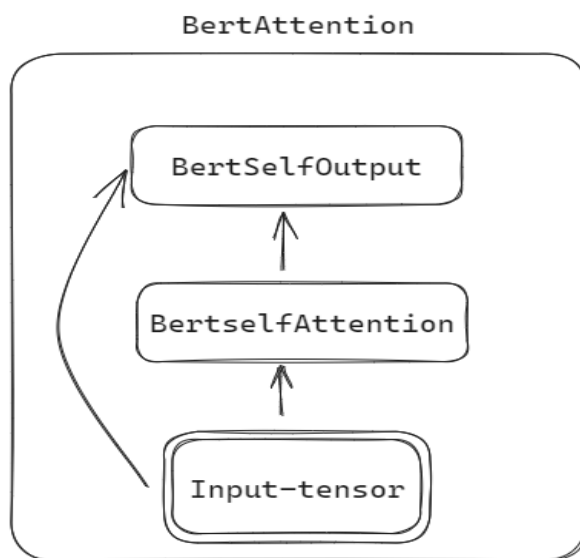


图 2: BertAttention

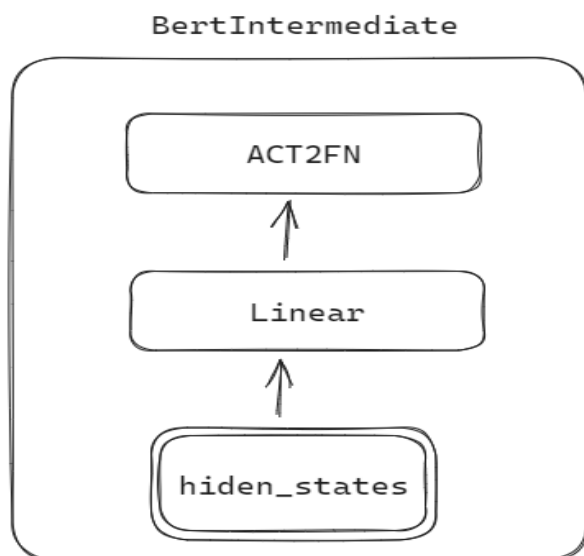


图 3: BertIntermediate

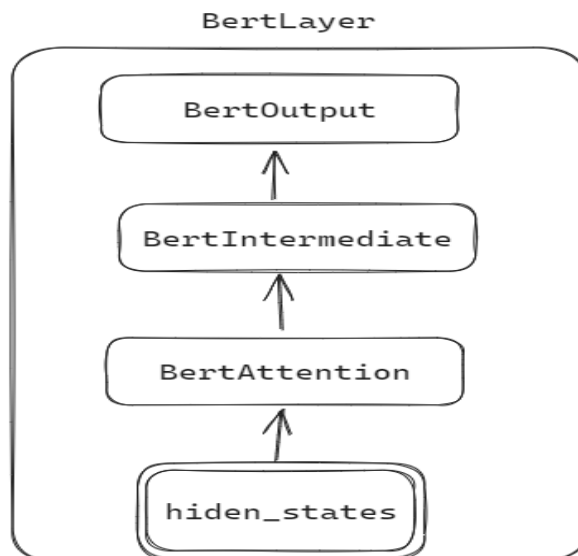


图 4: BertLayer

BertConfig类代码：关注： num_hidden_layers=12，有12个Layer

```

1      class BertConfig(object):
2          """Configuration class to store the configuration of a BertModel.
3          """
4
5      def __init__(self,
6                  vocab_size_or_config_json_file,
7                  hidden_size=768,
8                  num_hidden_layers=12,
9                  num_attention_heads=12,
10                 intermediate_size=3072,
11                 hidden_act="gelu",
12                 hidden_dropout_prob=0.1,
13                 attention_probs_dropout_prob=0.1,
14                 max_position_embeddings=512,
15                 type_vocab_size=2,
16                 relax_projection=0,
17                 initializer_range=0.02,
18                 task_idx=None,
19                 fp32_embedding=False,
20                 label_smoothing=None):
  
```

BertEncoder类代码：关注： self.layer

```

1      class BertEncoder(nn.Module):
2          def __init__(self, config):
3              super(BertEncoder, self).__init__()
4              layer = BertLayer(config)
5              self.layer = nn.ModuleList([copy.deepcopy(layer)
6                                          for _ in range(config.num_hidden_layers)])
7
8          def forward(self, hidden_states, attention_mask, prev_embedding=None,
9                      prev_encoded_layers=None, output_all_encoded_layers=True):
10             assert (prev_embedding is None) == (prev_encoded_layers is None), \
11                 "history embedding and encoded layer must be simultaneously given."
12             all_encoder_layers = []
13             if (prev_embedding is not None) and (prev_encoded_layers is not None):
14                 history_states = prev_embedding
15                 for i, layer_module in enumerate(self.layer):
16                     hidden_states = layer_module(
17                         hidden_states, attention_mask, history_states=history_states)
18                     if output_all_encoded_layers:
19                         all_encoder_layers.append(hidden_states)
20                     if prev_encoded_layers is not None:
21                         history_states = prev_encoded_layers[i]
22             else:
23                 for layer_module in self.layer:
24                     hidden_states = layer_module(hidden_states, attention_mask)
25                     if output_all_encoded_layers:
26                         all_encoder_layers.append(hidden_states)
27             if not output_all_encoded_layers:
28                 all_encoder_layers.append(hidden_states)
29             return all_encoder_layers

```

$hidden_states = layer_module(hidden_states, attention_mask)$ 这行就是表明:

$$H_l = Transformer(H_{l-1}), l \in [1, L]$$

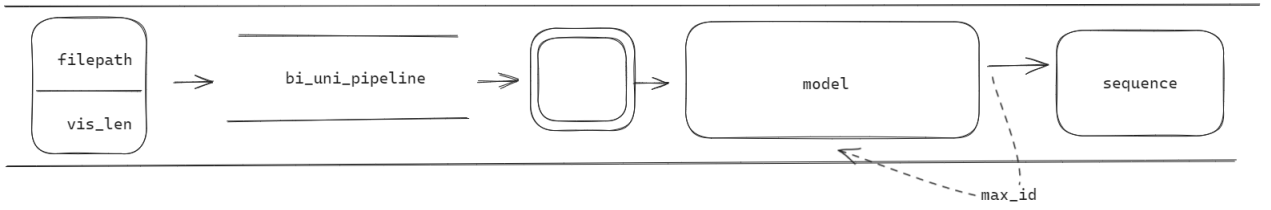


图 5: decode_img2txt文件中提供了i2t的方法，首先建立一个preprocess管道bi_uni_pipeline（来自于s2s_loader的Preprocess4Seq2seqDecoder）用来预处理输入的图像管道的输入：图像路径与图像大小管道的输出：input_id,..position_ids,..input_mask等管道的输出用instances来接收，再次经过处理之后，输送给model模型，为modeling的BertForSeq2SeqDecoder类对象，model输入：vis_feat,vis_pe,..attention_mask等 model输出：max_ids,max_probs

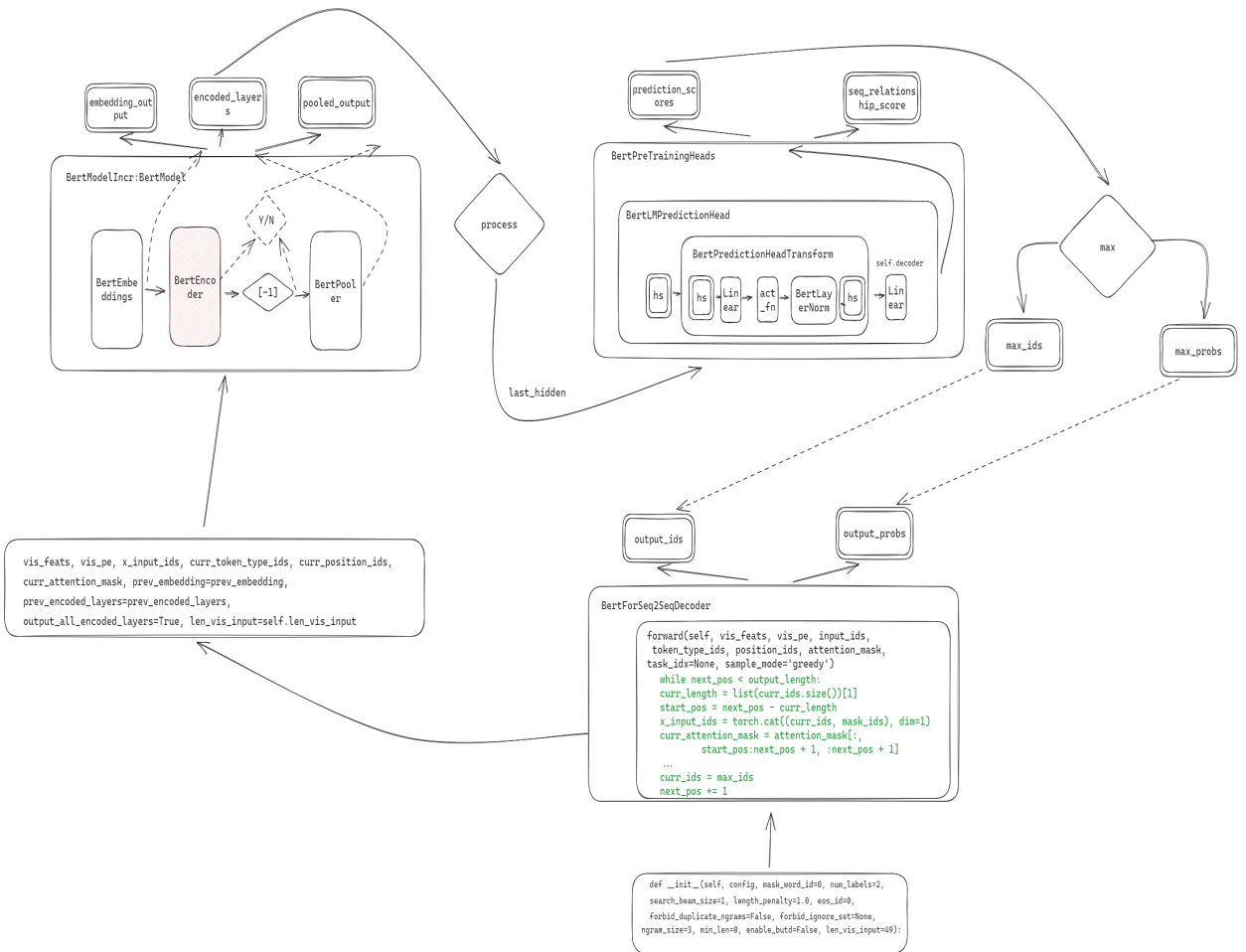


图 6: Bert4S2SDecoder, 生成过程由 next_pos < output_length: 该语句控制，没体现出来 STOP 标记的作用

decoder_img2txt.py:

```
1      #一个预处理图像的管道
2      #输入图像路径和图像大小
3      #返回input_ids, segment_ids(token_type_ids), position_ids, input_mask, task_idx, img,
      vis_pe
4      bi_uni_pipeline = []
5      bi_uni_pipeline.append(seq2seq_loader.Preprocess4Seq2seqDecoder(list(
6          tokenizer.vocab.keys()), tokenizer.convert_tokens_to_ids, args.max_seq_length,
7          max_tgt_length=args.max_tgt_length, new_segment_ids=args.new_segment_ids,
8          mode='s2s', len_vis_input=args.len_vis_input, enable_butd=args.enable_butd,
9          region_bbox_file=args.region_bbox_file,
      region_det_file_prefix=args.region_det_file_prefix))
10     ...
11
12     #用来解码的模型
13     #输入时
14     model = BertForSeq2SeqDecoder.from_pretrained(args.bert_model,
15         max_position_embeddings=args.max_position_embeddings,
16         config_path=args.config_path,
17         state_dict=model_recover, num_labels=cls_num_labels,
18         type_vocab_size=type_vocab_size, task_idx=3, mask_word_id=mask_word_id,
19         search_beam_size=args.beam_size, length_penalty=args.length_penalty,
20         eos_id=eos_word_ids, forbid_duplicate_ngrams=args.forbid_duplicate_ngrams,
21         forbid_ignore_set=forbid_ignore_set, ngram_size=args.ngram_size,
22         min_len=args.min_len,
23         enable_butd=args.enable_butd, len_vis_input=args.len_vis_input)
24     ...
25     for src in img_dat:
26         src_tk = os.path.join(args.image_root, src.get('filepath', 'trainval'),
27             src['filename'])
28         input_lines.append((img_idx, imgid, src_tk))
29
30     ...
31
32     _chunk = input_lines[next_i:next_i + args.batch_size]
33     buf = [x[2] for x in _chunk]
34     # instance 其实就是图像路径和输入图像大小的一个元组
```

```

35     for instance in [(x, args.len_vis_input) for x in buf]:
36         #用 instances 来接收图像预处理的输出
37         for proc in bi_uni_pipeline:
38             instances.append(proc(instance))
39
40     ...
41
42     input_ids, token_type_ids, position_ids, input_mask, task_idx, img, vis_pe = batch
43
44     #将管道的输出处理之后送入解码模型
45     #输出 max_id 与max_prob
46     traces = model(conv_feats, vis_pe, input_ids, token_type_ids,
47                    position_ids, input_mask, task_idx=task_idx)
48
49     #将返回的 max_id 转换为tokens
50     output_buf = tokenizer.convert_ids_to_tokens(w_ids)
51     output_tokens = []
52     #将非标记 tokens 提取出来
53     for t in output_buf:
54         if t in ("[SEP]", "[PAD]"):
55             break
56         output_tokens.append(t)
57     # tokens 变句子?
58     output_sequence = ' '.join(detokenize(output_tokens))
59     #记录每个图片对应的句子描述id
60     output_lines[buf_id[i]] = output_sequence
61     #图像标记序号为 tup[1] 的描述为output_lines[imd_idx]
62     predictions = [{'image_id': tup[1], 'caption': output_lines[img_idx]} for img_idx,
                    tup in enumerate(input_lines)]

```

seq2seq_loader.py:

```
1      #全1 的下三角矩阵, 创建 MASKED 矩阵用
2      self._tril_matrix = torch.tril(torch.ones(
3          (max_len, max_len), dtype=torch.long))
4
5
6      def __call__(self, instance):
7          img_path, max_a_len = instance[:2]
8          tokens_a = ['[UNK]'] * self.len_vis_input
9
10         # Add Special Tokens
11         padded_tokens_a = ['[CLS]'] + tokens_a + ['[SEP]']
12         tokens = padded_tokens_a
13         if self.new_segment_ids:
14             segment_ids = [4]*(len(padded_tokens_a)) \
15                 + [5]*(max_len_in_batch - len(padded_tokens_a))
16         else:
17             segment_ids = [0]*(len(padded_tokens_a)) \
18                 + [1]*(max_len_in_batch - len(padded_tokens_a))
19         ...
20
21         position_ids = []
22         for i in range(len(tokens_a) + 2):
23             position_ids.append(i)
24         for i in range(len(tokens_a) + 2, max_a_len + 2):
25             position_ids.append(0)
26         for i in range(max_a_len + 2, max_len_in_batch):
27             position_ids.append(i - (max_a_len + 2) + len(tokens_a) + 2)
28
29         # Token Indexing
30         input_ids = self.indexer(tokens)
31
32         # Zero Padding
33         #首先创建一个全 0 的矩阵
34         input_mask = torch.zeros(
35             max_len_in_batch, max_len_in_batch, dtype=torch.long)
36         #图像区域填充 1
37         input_mask[:, :len(tokens_a)+2].fill_(1)
38
39         second_st, second_end = len(padded_tokens_a), max_len_in_batch
```

```

40
41     #句子部分拷贝全 1 矩阵的句子部分 tokens 长度的下三角矩阵
42     input_mask[second_st:second_end, second_st:second_end].copy_(
43         self._tril_matrix[:second_end-second_st, :second_end-second_st])
44
45     img = torch.from_numpy(region_feat_f[img_id][:]).float()
46     vis_pe = torch.from_numpy(region_bbox_f[img_id][:])
47
48     # lazy normalization of the coordinates...
49     w_est = torch.max(vis_pe[:, [0, 2]])*1.+1e-5
50     h_est = torch.max(vis_pe[:, [1, 3]])*1.+1e-5
51     vis_pe[:, [0, 2]] /= w_est
52     vis_pe[:, [1, 3]] /= h_est
53     rel_area = (vis_pe[:, 3]-vis_pe[:, 1])*(vis_pe[:, 2]-vis_pe[:, 0])
54     rel_area.clamp_(0)
55     vis_pe = torch.cat((vis_pe[:, :4], rel_area.view(-1, 1), vis_pe[:, 5:]), -1) #
        confident score
56     normalized_coord = F.normalize(vis_pe.data[:, :5]-0.5, dim=-1)
57     vis_pe = torch.cat((F.layer_norm(vis_pe, [6]), \
58         F.layer_norm(cls_label, [1601])), dim=-1) # 1601 hard coded...
59
60     return (input_ids, segment_ids, position_ids, input_mask, self.task_idx, img,
        vis_pe)

```

```

1      #输出隐藏层的解码器，输入为：图像特征， word_id， attention_mask， 参数矩阵等
2      self.bert = BertModelIncr(config)
3      #输出预测结果的解码器，输入为隐藏层状态
4      self.cls = BertPreTrainingHeads(
5          config, self.bert.embeddings.word_embeddings.weight, num_labels=num_labels)
6
7
8      def forward(self, vis_feats, vis_pe, input_ids, token_type_ids, position_ids,
9                  attention_mask, task_idx=None, sample_mode='greedy'):
10
11          vis_feats = self.vis_embed(vis_feats) # image region features
12          vis_pe = self.vis_pe_embed(vis_pe) # image region positional encodings
13
14          ...
15          mask_ids = input_ids[:, :1] * 0 + self.mask_word_id
16          next_pos = list(input_ids.size())[1]
17          output_length = list(token_type_ids.size())[1]
18
19          while next_pos < output_length:
20              curr_length = list(curr_ids.size())[1]
21              start_pos = next_pos - curr_length
22              #把当前的输出当做输入来产生下一个词
23              x_input_ids = torch.cat((curr_ids, mask_ids), dim=1)
24              curr_token_type_ids = token_type_ids[:, start_pos:next_pos + 1]
25              #实时更新传入的 MASKED 矩阵
26              curr_attention_mask = attention_mask[:,
27                  start_pos:next_pos + 1, :next_pos + 1]
28              #更新矩阵pos
29              curr_position_ids = position_ids[:, start_pos:next_pos + 1]
30              #将图像与当前产生的 token 送入decoer 模块，产生所有hidden_states 层的输出
31              new_embedding, new_encoded_layers, _ = \
32                  self.bert(vis_feats, vis_pe, x_input_ids, curr_token_type_ids,
33                          curr_position_ids,
34                          curr_attention_mask, prev_embedding=prev_embedding,
35                          prev_encoded_layers=prev_encoded_layers,
36                          output_all_encoded_layers=True, len_vis_input=self.len_vis_input)
37
38              #只要最后一个隐藏层的输出
39              last_hidden = new_encoded_layers[-1][:, -1:, :]

```

```

38
39     #把最后一个隐藏层输出送入预测模块，产生预测结果
40     prediction_scores, _ = self.cls(
41         last_hidden, None, task_idx=task_idx)
42
43     #根据结果产生 max_id max_probs
44     if sample_mode == 'greedy':
45         max_probs, max_ids = torch.max(prediction_scores, dim=-1)
46     elif sample_mode == 'sample':
47         prediction_scores.squeeze_(1)
48         prediction_probs = F.softmax(prediction_scores, dim=-1).detach()
49         max_ids = torch.multinomial(prediction_probs, num_samples=1,
50             replacement=True)
51         max_probs = torch.gather(F.log_softmax(prediction_scores, dim=-1),
52             1, max_ids) # this should be logprobs
53     else:
54         raise NotImplementedError
55     output_ids.append(max_ids)
56     output_probs.append(max_probs)
57
58     #更新参数矩阵??
59     if prev_embedding is None:
60         prev_embedding = new_embedding[:, :-1, :]
61     else:
62         prev_embedding = torch.cat(
63             (prev_embedding, new_embedding[:, :-1, :]), dim=1)
64     if prev_encoded_layers is None:
65         prev_encoded_layers = [x[:, :-1, :]
66             for x in new_encoded_layers]
67     else:
68         prev_encoded_layers = [torch.cat((x[0], x[1][:, :-1, :]), dim=1)
69             for x in zip(prev_encoded_layers, new_encoded_layers)]
70     #当前输出作为下轮输入
71     curr_ids = max_ids
72     next_pos += 1
73     return torch.cat(output_ids, dim=1), torch.cat(output_probs, dim=1)

```

decode_img2txt文件中提供了i2t的方法，首先建立一个preprocess管道bi_uni_pipeline（来自于s2s_loader的Preprocess4Seq2seqDecoder）用来预处理输入的图像

管道的输入：图像路径与图像大小

管道的输出：input_id,..position_ids,..input_mask等

管道的输出用instances来接收，再次经过处理之后，输送给model模型，为modeling的BertForSeq2SeqDecoder类对象，

model输入：vis_feat,vis_pe,..attention_mask等

model输出：max_ids,max_probs