

0、前置准备

需要根据《HW_6作业思路提示》，完成一些初步的工作。

1、进行 fp16 加速并测试速度

及格标准：设置build_config，对模型进行fp16优化；

思路提示：认真阅读老师提供的代码，修改builder.sh中的输入参数。

优秀标准：编写fp16 版本的layer_norm算子，使模型最后运行fp16版本的layer_norm算子。

思路提示：layer_norm 核函数的实现时，使用函数模板来处理，通过将类型作为参数传递给模板，可使编译器生成该类型的函数。

2、进行 int8 加速并测试速度

完善calibrator.py内的todo函数，使用calibrator_data.txt 校准集，对模型进行int8量化加速。

思路提示：将FP32降为INT8的过程相当于信息再编码（re-encoding information），就是原来使用32bit来表示一个tensor，现在使用8bit来表示一个tensor，还要求精度不能下降太多。

这部分整体还是基本的操作，主要完成TODO部分工作即可。

首先是读入calibrator_data.txt文件，转化为模型的输入数据。

```
1 def text2inputs(tokenizer, text):
2     encoded_input = tokenizer.encode_plus(text, return_tensors = "pt")
3
4     input_ids = encoded_input['input_ids'].int().detach().numpy()
5     token_type_ids = encoded_input['token_type_ids'].int().detach().numpy()
6     # position_ids = torch.arange(0, encoded_input['input_ids'].shape[1]).int().view(1,
7     -1).numpy()
8     seq_len = encoded_input['input_ids'].shape[1]
9     position_ids = np.arange(seq_len, dtype = np.int32).reshape(1, -1)
10    input_list = [input_ids, token_type_ids, position_ids]
11
12    return input_list
13
14 # TODO: your code, read inputs
15 with open(data_txt, "r") as f:
16     lines = f.readlines()
```

```

16         for i in range(0, num_inputs):
17             inputs = text2inputs(tokenizer, lines[i])
18             self.input_ids_list.append(inputs[0])
19             self.token_type_ids_list.append(inputs[1])
20             self.position_ids_list.append(inputs[2])
21             if i % 10 == 0:
22                 print("text2inputs:" + lines[i])

```

然后是get_batch()中的代码补全，这里贴出参考实现代码，如下：

```

1  # TODO your code, copy input from cpu to gpu
2  input_ids = self.input_ids_list[self.current_index]
3  token_type_ids = self.token_type_ids_list[self.current_index]
4  position_ids = self.position_ids_list[self.current_index]
5
6  seq_len = input_ids.shape[1]
7  if seq_len > self.max_seq_length:
8      print(seq_len)
9      print(input_ids.shape)
10     input_ids = input_ids[:, :self.max_seq_length]
11     token_type_ids = token_type_ids[:, :self.max_seq_length]
12     position_ids = position_ids[:, :self.max_seq_length]
13     print(input_ids.shape)
14
15     cuda.memcpy_htod(self.device_inputs[0], input_ids.ravel())
16     cuda.memcpy_htod(self.device_inputs[1], token_type_ids.ravel())
17     cuda.memcpy_htod(self.device_inputs[2], position_ids.ravel())
18
19     self.current_index += self.batch_size

```

参考资料：

- <https://github.com/NVIDIA/TensorRT> (TensorRT GitHub, 特别是demo 和samples部分)
- <https://docs.nvidia.com/deeplearning/tensorrt/index.html> (TensorRT 官方文档)