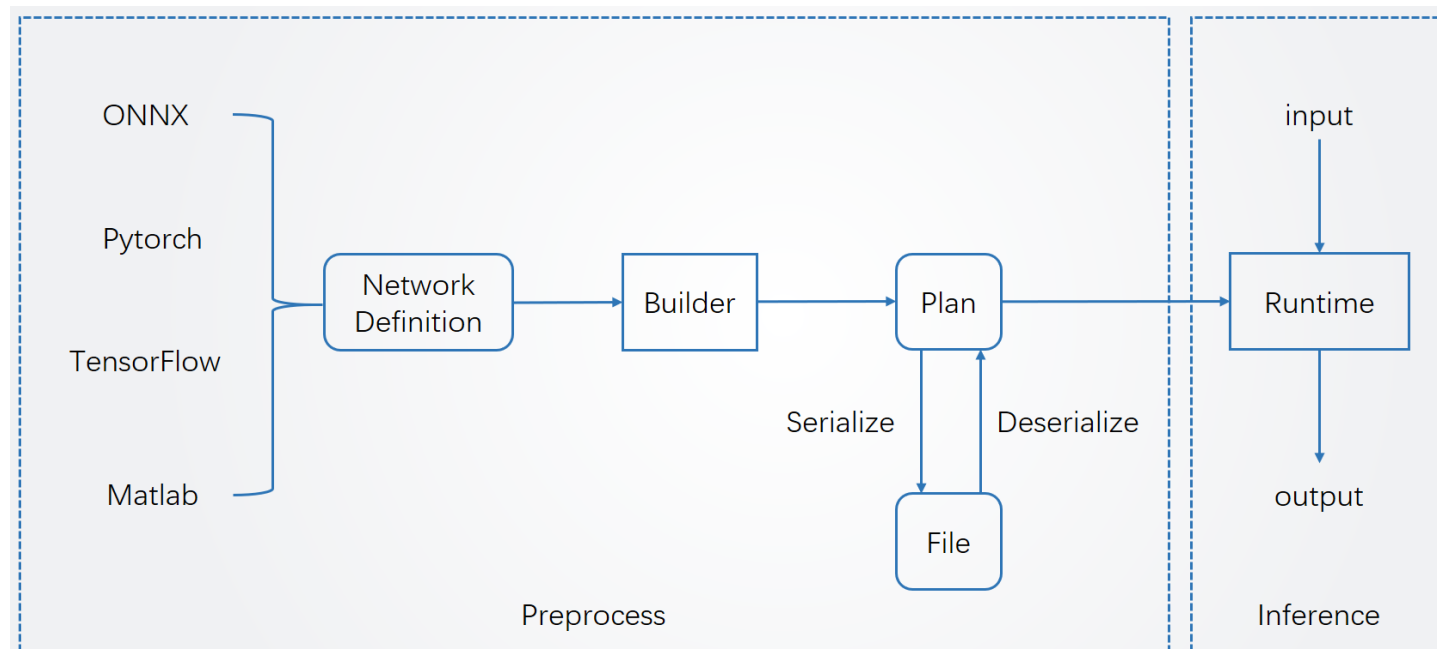


## 零、基础概念

TensorRT使用流程如下图所示，分为两个阶段：预处理阶段和推理阶段。其部署大致流程如下：  
1.导出网络定义以及相关权重；2.解析网络定义以及相关权重；3.根据显卡算子构造出最优执行计划；  
4.将执行计划序列化存储；5.反序列化执行计划；6.进行推理。

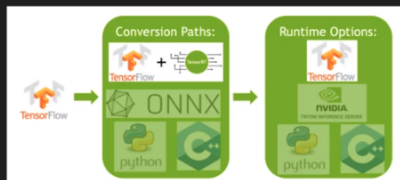


可以从步骤3可以得知，tensorrt实际上是和你的硬件绑定的，所以在部署过程中，如果你的硬件（显卡）和软件（驱动、cudatoolkit、cudnn）发生了改变，那么这一步开始就要重新走一遍了。换句话说，不同系统环境下生成Engine 包含硬件有关优化，不能跨硬件平台使用，注意环境统一(硬件环境+ CUDA/cuDNN/TensorRT 环境)，并且不同版本TensorRT 生成的 engine不能相互兼容，同平台同环境多次生成的engine可能不同。

基于TensorRT进行开发的三种工作流程如下图所示：（第5章作业是使用Parser，第6，7章作业是使用TensorRT API开发）

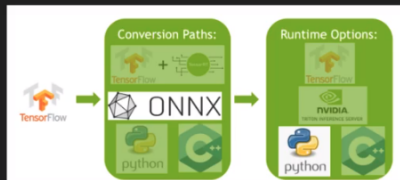
## ➤ 使用框架自带 TRT 接口 (TF-TRT, Torch-TensorRT)

- 简单灵活，部署仍在原框架中，无需书写 Plugin



## ➤ 使用 Parser (TF/Torch/... → ONNX<sup>[1]</sup> → TensorRT)

- 流程成熟，ONNX 通用性好，方便网络调整，兼顾效率性能



## ➤ 使用 TensorRT 原生 API 搭建网络

- 性能最优，精细网络控制，兼容性最好

方法	易用性	性能	兼容性	开发效率	遇到不支持的OP
框架自带接口	★★★★	★	★★	★★★★	返回原框架计算
使用Parser	★★	★★☆	★★☆	★★	改网/改Parser/写Plugin
API搭建	★	★★★★	★★★★	★	写Plugin

➤ [1] TensorRT 也支持 UFF 和 prototxt 格式的网络，但在未来版本中将被废弃

## 推荐学习资料：

- [https://www.bilibili.com/video/BV15Y4y1W73E?spm\\_id\\_from=333.337.search-card.all.click](https://www.bilibili.com/video/BV15Y4y1W73E?spm_id_from=333.337.search-card.all.click)
- <https://github.com/NVIDIA/trt-samples-for-hackathon-cn/tree/master/cookbook>
- <https://docs.nvidia.com/deeplearning/tensorrt/developer-guide/index.html> (TensorRT 文档)
- [https://docs.nvidia.com/deeplearning/tensorrt/api/c\\_api/](https://docs.nvidia.com/deeplearning/tensorrt/api/c_api/) (TensorRT C++ API)
- [https://docs.nvidia.com/deeplearning/tensorrt/api/python\\_api/](https://docs.nvidia.com/deeplearning/tensorrt/api/python_api/) (TensorRT Python API)
- <https://developer.nvidia.com/nvidia-tensorrt-download> (TensorRT 下载)
- <https://docs.nvidia.com/deeplearning/tensorrt/support-matrix/index.html> (TensorRT 版本支持列表)

# 一、背景介绍

深蓝学院《CUDA入门与深度神经网络加速》课程，TensorRT 部分的作业。

内容是使用TensorRT Python API 手动搭建图像识别Vision Transformer(ViT)模型。学员不需要从零开始搭建，这里已经搭好了一个模型结构框架，进行相应的填充即可。

ViT 模型知识给大家推荐一些博客文章，请大家自行学习：

- Vision Transformer 超详细解读系列文章 <https://zhuanlan.zhihu.com/p/340149804>
- 详解 Vision Transformer (ViT) [https://blog.csdn.net/qq\\_39478403/article/details/118704747](https://blog.csdn.net/qq_39478403/article/details/118704747)

**注意：** 整个作业只支持batch=1

## 二、模型数据介绍和下载

### 2.1 模型下载

github上有多个ViT项目，这里选择了一个最适合学习的项目（已经下载了，见压缩包ViT-pytorchmain.zip）。

```
1 https://github.com/jeonsworld/ViT-pytorch
2 model_type: ViT-B_16
3 dataset: cifar10-100_500
4 input: [batch, 3, 224, 224]
5 output: [batch, 10]
```

但这个项目有个比较大的缺点，只提供了pre-train和 fine-tuned 的numpy格式的模型，识别率为0，无法直接使用。有两个改进方案：

方案1，按照该项目readme中介绍的使用方法，训练出来一个模型，并转成onnx模型。

方案2，课程提供了一个只训练了100 step的模型，一个checkpoint模型和一个转好的onnx模型。

文件下载地址：<https://pan.baidu.com/s/1StW6Z4yy8uTklR52X9K8Cg?pwd=m3w7>，提取码：m3w7。

### 2.2 数据下载

使用 CIFAR10-100\_500 数据集。

1. github上的项目，训练时为了提高读取文件的效率，采用的是将多张图片合并成一个文件的形式。

<https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz>

2. 咱们作业使用的更直观的 jpg 格式的数据，共十类，每类1000张图片。<https://github.com/YoongiKim/CIFAR-10-images>

### 2.3 验证base准确性

1. 使用 onnx 模型和 jpg 格式的数据 验证模型准确性。

2. 方案2提供的模型，识别准确率为89.72%（很低，因为只训练了100 step）

3. 验证命令

```
1 python valid.py -x ViT-B_16.onnx -d CIFAR-10-images-master/test/
```

## 2.4 测速

1. 测试 onnx cpu 速度 (选做)
2. 测试 onnx gpu 速度 (选做)
3. 测试 使用trt onnxparser 直接转出来的trt模型速度
4. 测试 手动转换出来的 不同精度 trt 模型的速度

trt 模型测速方式, 建议使用trtexec, 选择 Host Latency or GPU Compute 的 mean 时间即可:

```
1 ./trtexec --loadEngine=model.plan --shapes=input:1x1x224x224 --plugins=LayerNorm.so
```

## 三、目录文件介绍

```
1 ViT-pytorch # https://github.com/jeonsworld/ViT-pytorch
2 model2onnx.py # 将 npy 模型 转成 onnx 模型的代码, 注意重要的参数do_constant_folding=False
3
4 # layernorm 相关
5 LayerNormPlugin
6 test_nn_layer_norm.py
7 test_nn_layer_norm.sh
8
9 # 验证onnx和trt模型准确率的脚本
10 valid.sh
11 valid.py
12
13 # onnx 模型转成 trt 格式的相关
14 # builder.py 中的结构和 ViT-pytorch/models/modeling.py 差不多。
15 builder.sh
16 builder.py
17 trt_helper.py
18 calibrator.py
19 ViT-B_16.onnx
20
21 # 这个脚本很重要, 支持load checkpoint 格式模型, 并跑一张图片
22 # 当转出来的 trt 模型结果对不上时, 可以使用这个脚本进行调试。
23 # 使用方法自己琢磨吧。
24 test.py
```

## 四、作业内容

### 0. 使用trt onnx-parser 将 onnx 模型转成 trt 格式，并测速。

```
1 #使用 onnx-parser将onnx模型转换成trt plan格式
2 #sh builder.sh
3 #验证
4 python valid.py -p model.plan -d CIFAR-10-images-master/test/
```

#### 思路提示：

首先，获取onnx模型，有两种方式：

- 使用老师提供的onnx模型
- 或者执行model2onnx.py自己生成出onnx模型

然后，将 onnx 模型转成 trt 格式（可以参考第5章作业），并测速。

执行如下命令，实现转换。如需测试，则徐指定一张图片的路径。

```
1 python builder.py -x ./models/ViT-B_16.onnx -o ./plans/model.plan
2
3 python builder.py -x ./models/ViT-B_16.onnx -o ./plans/model.plan --img_path ViT-
  pytorch/00816.jpeg
```

```
[Constant]_output[Float(768)] -> (Unnamed Layer* 572) [PluginV2DynamicExt]_output_0[Float(1,197,768)]
Layer(Reformat): 575_nn.Linear_pre_reshape_copy_input, Tactic: 0x00000000000003e8, (Unnamed Layer* 572) [PluginV2DynamicExt]_output_0[Float(1,1,768)] -> 575_
,1,768))
Layer(NoOp): Reformating CopyNode for Input Tensor 0 to 575_nn.Linear_pre_reshape, Tactic: 0x0000000000000000, 575_nn.Linear_pre_reshape_copy_input[Float(1,
575_nn.Linear_pre_reshape[Float(1,1,768)]
Layer(NoOp): 575_nn.Linear_pre_reshape, Tactic: 0x0000000000000000, Reformatted Input Tensor 0 to 575_nn.Linear_pre_reshape[Float(1,1,768)] -> (Unnamed Layer
1)]
Layer(CublasConvolution): 576_nn.Linear, Tactic: 0x0000000000000001, (Unnamed Layer* 574) [Shuffle]_output[Float(1,1,768,1,1)] -> (Unnamed Layer* 575) [Fully
Layer(NoOp): 577_nn.Linear_after_reshape, Tactic: 0x0000000000000000, (Unnamed Layer* 575) [Fully Connected]_output[Float(1,1,10,1,1)] -> (Unnamed Layer* 576
[09/15/2022-18:39:41] [TRT] [I] [MemUsageChange] TensorRT-managed allocation in building engine: CPU +0, GPU +0, now: CPU 0, GPU 0 (MiB)
[09/15/2022-18:39:41] [TRT] [I] build engine in 12.403 Sec
[09/15/2022-18:39:41] [TRT] [V] Serializing Engine...
[09/15/2022-18:39:41] [TRT] [W] The getMaxBatchSize() function should not be used with an engine built from a network created with NetworkDefinitionCreationF
will always return 1.
[09/15/2022-18:39:41] [TRT] [W] The getMaxBatchSize() function should not be used with an engine built from a network created with NetworkDefinitionCreationF
will always return 1.
[09/15/2022-18:39:41] [TRT] [I] Saving Engine to ./plans/model.plan
[09/15/2022-18:39:41] [TRT] [I] Done.
root@fa90d08d8678:/workspace/ViT_TRT#
```

```

=====test_img=====
torch.Size([3, 224, 224])
tensor([[ 0.2157,  0.2157,  0.2157, ...,  0.2235,  0.1765,  0.0353],
        [ 0.2157,  0.2157,  0.2157, ...,  0.2000,  0.2000,  0.0980],
        [ 0.2157,  0.2157,  0.2157, ...,  0.2314,  0.2471,  0.1686],
        ...,
        [-0.0353, -0.0196,  0.0039, ..., -0.2706, -0.2706, -0.2549],
        [-0.0510, -0.0275, -0.0118, ..., -0.2784, -0.2627, -0.2314],
        [-0.0588, -0.0431, -0.0196, ..., -0.2863, -0.2627, -0.2157]],
        [[ 0.2078,  0.2078,  0.2078, ...,  0.1451,  0.0980, -0.0431],
         [ 0.2078,  0.2078,  0.2078, ...,  0.1373,  0.1373,  0.0353],
         [ 0.2078,  0.2078,  0.2078, ...,  0.1922,  0.2078,  0.1294],
         ...,
         [-0.1843, -0.1686, -0.1451, ..., -0.4745, -0.4745, -0.4588],
         [-0.2000, -0.1765, -0.1608, ..., -0.4824, -0.4667, -0.4353],
         [-0.2078, -0.1922, -0.1686, ..., -0.4902, -0.4667, -0.4196]],
        [[ 0.1765,  0.1765,  0.1765, ...,  0.1294,  0.0824, -0.0588],
         [ 0.1765,  0.1765,  0.1765, ...,  0.1137,  0.1137,  0.0118],
         [ 0.1765,  0.1765,  0.1765, ...,  0.1608,  0.1765,  0.0980],
         ...,
         [-0.3020, -0.2863, -0.2627, ..., -0.5765, -0.5765, -0.5608],
         [-0.3176, -0.2941, -0.2784, ..., -0.5843, -0.5686, -0.5373],
         [ 0.3255,  0.3090,  0.2863, ..., -0.5922, -0.5686, -0.5216]]])
time=49.89453700000013ms
outputs.shape:(1, 1, 10)
outputs.sum:0.0
[[[-0.1607827 -0.2510565 -0.08208042  0.9764596 -0.10142983
      0.3275517 -0.21142386 -0.16902082 -0.1875862 -0.140631 ]]]
root@fa90d08d8678:/workspace/ViT TRT#

```

## 1. 学习使用trt python api 搭建网络

填充trt\_helper.py 中的空白函数，包括Linear, LayerNorm, addSoftmax等。学习使用 api 搭建网络的过程。

**思路提示：**

这部分的主要工作是使用TensorRT Python API实现Bert网络中用到的基本单元操作。

那么大家首先要熟悉TensorRT Python API有哪些，这些函数的输入及输出是什么，都要有清晰的认知。详细的可以参看 [https://docs.nvidia.com/deeplearning/tensorrt/api/python\\_api/infer/Graph/Network.html](https://docs.nvidia.com/deeplearning/tensorrt/api/python_api/infer/Graph/Network.html)。

另外，builder.py中有关于模型的实现，可以加深对于模型的理解。

```

1 def addLinear(self, x, weight, bias, layer_name=None, precision=None):
2     input_len = len(x.shape)
3     if input_len < 3:
4         raise RuntimeError("addLinear x.shape.size must >= 3")
5

```

```

6         if layer_name is None:
7             layer_name = "nn.Linear"
8
9         # calc pre_reshape_dims and after_reshape_dims
10        pre_reshape_dims = trt.Dims()
11        after_reshape_dims = trt.Dims()
12        if input_len == 3:
13            pre_reshape_dims = (0, 0, 0, 1, 1)
14            after_reshape_dims = (0, 0, 0)
15        elif input_len == 4:
16            pre_reshape_dims = (0, 0, 0, 0, 1, 1)
17            after_reshape_dims = (0, 0, 0, 0)
18        elif input_len == 5:
19            pre_reshape_dims = (0, 0, 0, 0, 0, 1, 1)
20            after_reshape_dims = (0, 0, 0, 0, 0)
21        else:
22            raise RuntimeError("addLinear x.shape.size >5 not support!")
23
24        # add pre_reshape layer
25        trt_layer = self.network.add_shuffle(x)
26        trt_layer.reshape_dims = pre_reshape_dims
27
28        self.layer_post_process(trt_layer, layer_name+"_pre_reshape", precision)
29
30        x = trt_layer.get_output(0)
31
32        # add Linear layer
33        out_features = weight.shape[0]
34        weight = trt.Weights(weight)
35        if bias is not None:
36            bias = trt.Weights(bias)
37
38        trt_layer = self.network.add_fully_connected(x, out_features, weight, bias)
39        self.layer_post_process(trt_layer, layer_name, precision)
40        x = trt_layer.get_output(0)
41
42        # add after_reshape layer
43        trt_layer = self.network.add_shuffle(x)
44        trt_layer.reshape_dims = after_reshape_dims

```

```

45         self.layer_post_process(trt_layer, layer_name+"_after_reshape", precision)
46         x = trt_layer.get_output(0)
47
48         return x

```

```

1  def addSoftmax(self, x: trt.ITensor, dim: int = -1, layer_name=None, precision=None) ->
    trt.ITensor:
2      trt_layer = self.network.add_softmax(x)
3
4      input_len = len(x.shape)
5      if dim is -1:
6          dim = input_len
7      trt_layer.axes = int(math.pow(2, input_len-1))
8
9      layer_name_prefix = "nn.Softmax[dim=" + str(dim) + "]"
10     if layer_name is None:
11         layer_name = layer_name_prefix
12     else:
13         layer_name = layer_name_prefix + "." + layer_name
14
15     self.layer_post_process(trt_layer, layer_name, precision)
16
17     x = trt_layer.get_output(0)
18     return x

```

```

1  def addAdd(self, a, b, layer_name=None, precision=None):
2      trt_layer = self.network.add_elementwise(a, b, trt.ElementWiseOperation.SUM)
3      if layer_name is None:
4          layer_name = "elementwise.sum"
5      else:
6          layer_name = "elementwise.sum." + layer_name
7
8      self.layer_post_process(trt_layer, layer_name, precision)
9
10     x = trt_layer.get_output(0)
11     return x

```



```

1  def addScale(
2      self,
3      x: trt.ITensor,
4      scale: float,
5      layer_name: str = None,
6      precision: trt.DataType = None
7  ) -> trt.ITensor:
8      """scale"""
9      input_len = len(x.shape)
10     if input_len < 3:
11         raise RuntimeError("input_len < 3 not support now! ")
12
13     if layer_name is None:
14         layer_name = "Scale"
15
16     # The input dimension must be greater than or equal to 4
17     if input_len is 3:
18         trt_layer = self.network.add_shuffle(x)
19         trt_layer.reshape_dims = (0, 0, 0, 1)
20         self.layer_post_process(trt_layer, layer_name+".3dto4d", precision)
21         x = trt_layer.get_output(0)
22
23     np_scale = trt.Weights(np.array([scale], dtype=np.float32))
24     trt_layer = self.network.add_scale(x, mode=trt.ScaleMode.UNIFORM,
25                                       shift=None, scale=np_scale, power=None)
26     self.layer_post_process(trt_layer, layer_name, precision)
27     x = trt_layer.get_output(0)
28
29     if input_len is 3:
30         trt_layer = self.network.add_shuffle(x)
31         trt_layer.reshape_dims = (0, 0, 0)
32         self.layer_post_process(trt_layer, layer_name+".4dto3d", precision)
33         x = trt_layer.get_output(0)
34
35     return x

```

```

1  def addMatMul(self, a: trt.ITensor, b: trt.ITensor, layer_name: Optional[str] = None) ->
    trt.ITensor:

```

```

2         trt_layer = self.network.add_matrix_multiply(a, trt.MatrixOperation.NONE,
3                                                       b, trt.MatrixOperation.NONE)
4
5         if layer_name is None:
6             layer_name = "torch.matmul"
7         else:
8             layer_name = "torch.matmul." + layer_name
9
10        self.layer_post_process(trt_layer, layer_name, None)
11
12        x = trt_layer.get_output(0)
13        return x

```

## 2. 编写layernorm plugin

trt不支持layer\_norm算子，编写layer\_norm plugin，并将算子添加到网络中，进行验证。

1. 及格：将“基础款LayerNormPlugin.zip”中实现的基础版 layer\_norm算子 插入到 trt\_helper.pyaddLayerNorm函数中。
2. 优秀：将整个layer\_norm算子实现到一个kernel中，并插入到 trt\_helper.py addLayerNorm函数中。可以使用testLayerNormPlugin.py对合并后的plugin进行单元测试验证。
3. 进阶：在2的基础上进一步优化，线索见 [https://www.bilibili.com/video/BV1i3411G7vN?spm\\_id\\_from=333.999.0.03](https://www.bilibili.com/video/BV1i3411G7vN?spm_id_from=333.999.0.03)。

### 思路提示：

首先，要了解LayerNorm是啥，什么原理，是怎么操作实现的。例如对于输入 x 形状为 (N, C, H, W)，normalized\_shape 为 (H, W) 的情况，可以理解为输入 x 为 (N\*C, H\*W)，在 N\*C 个行上，每行有 H\*W 个元素，对每行的元素求均值和方差，得到 N\*C 个 mean 和 inv\_variance，再对输入按如下 LayerNorm 的计算公式计算得到 y。若 elementwise\_affine=True，则有 H\*W 个 gamma 和 beta，对每行 H\*W 个的元素做变换。

然后，要编译基础版LayerNormPlugin编译，注意修改Makefile中 CUDA\_PATH，TRT\_PATH的值，对应于自己的环境进行配置。执行编译后，得到如下的情况

然后，再将算子插入到addLayerNorm函数中。

```

1 def addLayerNorm(self, x, gamma, beta, layer_name=None, precision=None):
2     plg_creator = self.plugin_registry.get_plugin_creator("LayerNorm", "1", "")
3     if not plg_creator:
4         raise RuntimeError("Could not find LayerNorm")

```

```

5
6     # pfc = trt.PluginFieldCollection([data_type, dim, eps, gamma_w, beta_w])
7     pfc = trt.PluginFieldCollection([])
8     plugin = plg_creator.create_plugin("LayerNorm", pfc)
9     if not plugin:
10         raise RuntimeError("Could not create_plugin LayerNormPluginDynamic")
11
12     gamma = self.network.add_constant(gamma.shape, gamma).get_output(0)
13     beta = self.network.add_constant(beta.shape, beta).get_output(0)
14
15     trt_layer = self.network.add_plugin_v2([x, gamma, beta], plugin)
16
17     if layer_name is None:
18         layer_name = "nn.LayerNorm"
19
20     self.layer_post_process(trt_layer, layer_name, precision)
21
22     return trt_layer.get_output(0)

```

### 3. 观察GELU算子的优化过程

GELU算子使用一堆基础算子堆叠实现的（详细见trt\_helper.py addGELU函数），直观上感觉很分散，计算量比较大。

但在实际build过程中，这些算子会被合并成一个算子。build 过程中需要设置log为trt.Logger.VERBOSE，观察build过程。

**思路提示：**理解addGELU函数的操作，然后对应于Build过程中的相关操作。

### 4. 学习 builder.py

重点是梳理build\_embeddings\_layer函数逻辑，体会pytorch api 和 trt api 的差异。

### 5. 进行 fp16 加速，观察模型大小，准确率和速度

需要注意plugin 是否支持 fp16? 是否设置了fp16?

1. 及格：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。

2. 优秀：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。

## 6. 进行 int8 加速，观察模型大小，准确率和速度

完善calibrator.py内的todo函数

1. int8 出来的模型，准确率可能会掉很多，为什么

2. int8 和 fp16 可以都 enable，观察 模型大小，准确率和速度

**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。

```
1 classes = ["airplane", "automobile", "bird", "cat", "deer", "dog", "frog", "horse",  
2 "ship", "truck"]  
3  
4 def read_all_test_imgs(args, num_inputs_per_class):  
5     imgs = []  
6     labels = []  
7     label_idx = 0  
8     class_count = 0  
9     for c in classes:  
10         path = args.calib_path+ "/" + c  
11         for i in os.listdir(path):  
12             # print(path + "/" + i)  
13             # assert 0  
14             img = image_preprocess(args, path + "/" + i).numpy()  
15             imgs.append(img)  
16             labels.append(label_idx)  
17             class_count = class_count + 1  
18             if class_count >= num_inputs_per_class:  
19                 break  
20             label_idx = label_idx + 1  
21     print(f"read {c} done ..., img num = {class_count}")  
22     class_count = 0
```

```
21
22     return imgs, labels
23
```

```
1 class ViTCalibrator(trt.IInt8LegacyCalibrator):
2     def __init__(self, args, cache_file, batch_size, num_inputs):
3         # Whenever you specify a custom constructor for a TensorRT class,
4         # you MUST call the constructor of the parent explicitly.
5         trt.IInt8LegacyCalibrator.__init__(self)
6
7         self.imgs, _ = read_all_test_imgs(args, num_inputs)
8
9         self.cache_file = cache_file
10
11        self.batch_size = batch_size
12        self.current_index = 0
13        if num_inputs > len(self.imgs):
14            self.num_inputs = len(self.imgs)
15        else:
16            self.num_inputs = num_inputs
17
18        # Allocate enough memory for a whole batch.
19        self.device_input = cuda.mem_alloc(self.batch_size * self.imgs[0].nbytes)
20
```

```
1 def get_batch(self, names):
2     if self.current_index + self.batch_size > self.num_inputs:
3         print("Calibrating index {:} batch size {:} exceed max input limit {:}
4 sentences".format(self.current_index, self.batch_size, self.num_inputs))
5
6         return None
7
8         current_batch = int(self.current_index / self.batch_size)
9         if current_batch % 10 == 0:
10            print("Calibrating batch {:}, containing {:} imgs".format(current_batch,
11 self.batch_size))
12
13
14        # import pdb
```

```

11         # pdb.set_trace()
12         cuda.memcpy_htod(self.device_input, self.imgs[self.current_index].ravel())
13
14         self.current_index += self.batch_size
15         return [self.device_input]

```

```

1 if __name__ == '__main__':
2     data_txt = "calibrator_data.txt"
3     bert_path = "bert-base-uncased"
4     cache_file = "bert_calibrator.cache"
5     batch_size = 1
6     max_seq_length = 200
7     num_inputs = 100
8     cal = BertCalibrator(data_txt, bert_path, cache_file, batch_size, max_seq_length,
9                           num_inputs)
9
10    cal.get_batch("input")
11    cal.get_batch("input")
12    cal.get_batch("input")

```

## 7. 支持batch\_size > 1(选做)

可以按照以下提示进行修改

- 1 1. 修改 builder.py 中的 profile 和valid.py
- 2 2. builder.py build\_embeddings\_layer 函数中, cls\_token 的batch 维度只有1
- 3 3. builder.py build\_vision\_transformer\_layer函数中, Slice batch维度
- 4 4. calibrator.py 中, 需要进行 batch 拼接。