# pytorch导出onnx的原则-以SwinTransformer和 DETR在trt8.0.3.4部署为例

## 一、前言

pytorch导出onnx是通过追踪forward函数构建的计算图来实现的，然而这种构建的方式往往会由于模型的设计者没有考虑导出部署时的便捷性从而生成的onnx中算子比较冗余，比如 `x[x> 0.5] = 0` 这种inplace操作就会生成onnx的 `scatter` 算子，这个算子是不被tensorrt支持的。尽管可以采用plugin方式来实现，但其并不是一个较为有好的算法。`x[x > 0.5] = 0` 可以转换为 `y = x * (x > 0.5).to(torch.float32)` 这种形式来避免inplace操作。此外取shape操作一般用于动态shape，这种操作在部署框架中也是冗余的应该去掉，主要原因是大部分模型除了batch size可变以外，其余维度输入都是固定的，从而对于一般的取shape操作都是可以直接由常量计算得到，从而计算得到的shape实际上可以固化为一个常量。下面根据trtpy作者提出的建议，给出导出onnx时应遵守的原则。

## 二、原则

1. 对于任何用到shape、size返回值的参数时，例如：`tensor.view(tensor.size(0), -1)`，`B,C,H,W = x.shape` 这类操作，避免直接使用tensor.size的返回值，而是加上int转换，`tensor.view(int(tensor.size(0)), -1)`,`B,C,H,W = map(int, x.shape)`，断开跟踪。

2. 对于nn.Upsample或nn.functional.interpolate函数，一般使用scale_factor指定倍率，而不是使用size参数指定大小。如果源码中就是插值为固定大小，则该条忽略。

3. 对于reshape、view操作时，-1的指定请放到batch维度。其他维度计算出来即可。batch维度禁止指定为大于-1的明确数字。如果是一维，那么直接指定为-1就好。

4. torch.onnx.export指定dynamic_axes参数，并且只指定batch维度，禁止其他动态

5. 使用opset_version=11，不要低于11

6. 避免使用inplace操作，例如 `y[…, 0:2] = y[…, 0:2] * 2 - 0.5`，可以采用如下代码代替 `tmp = y[…, 0:2] * 2 - 0.5; y = torch.cat((y[..., 2:], tmp), dim = -1)`

7. 尽量少的出现5个维度，例如ShuffleNet Module，可以考虑合并wh避免出现5维

8. 尽量把让后处理部分在onnx模型中实现，降低后处理复杂度。比如在目标检测网络中最终输出设置为xywh或者xyxy，而不是一个中间结果。

## 三、常见问题方法

1. onnx的修改

使用python中 `onnx` 库或者TensorRT携带的更为方便的 `onnx_graphsurgeon`

```
1   # 安装onnx
2   pip install onnx
3   # 安装onnx_graphsurgeon
4   https://github.com/NVIDIA/TensorRT/tree/master/tools/onnx-graphsurgeon
    python3 -m pip install onnx_graphsurgeon --index-url
    https://pypi.ngc.nvidia.com
```

具体修改需要去官网查看相关案例，这里不再赘述。此外，onnx的本质时protobuf序列化后的二进制文件，所以其结构遵循 `.proto` 定义的文件，如果是使用 `onnx` 库修改时需要牢记这一点。

1. onnx的简化

`onnxsim` 可以简化onnx graph中一些常量的折叠运算，具体见官方github
https://github.com/daquexian/onnx-simplifier 。 即使在使用了以上原则以后，还是无法避免生成一些冗余的常量运算。比如 `d = a + b - c`

这里a,b,c都是常量且只有在计算d的时候用到，那么onnx计算图中只需要表示d即可，但是torch的api还做不到这一点。

1. 动态shape的问题

TensorRT是支持各个维度的动态shape的，不局限于动态的batchsize，其允许用户设定多组维度供TensorRT结构优化器来选择，但是在选择最优结构的时候需要预先指定一个固定维度。此外动态shape对于一般网络而言，仅仅是batch size可变，其余模型输入的长和宽一般是固定的，因此一般情况下只考虑动态batch size的问题，也就是**原则3**。

1. 获取onnx中间输入

onnxruntime是可以获得模型中间层结果的，具体方式还请读者搜索一下。本文给出一种更为优雅的方式，即采用TensorRT附带的 `polygraphy` ，具体网址
https://github.com/NVIDIA/TensorRT/tree/master/tools/Polygraphy 下文在DETR案例中将给出如何采用 `polygraphy` 定位哪个中间层出现不一致的情况，也就是 `onnxruntime` 结果与 `tensorrt` 结果不一致。

# 四、案例1-SwinTransformer在TensorRT上的部署

## 4.1 加载预训练权重

本文采用基于pytorch的 `timm` 库来加载分类模型和预训练权重，采用 `pip install timm` 安装，作者在做本次实验时 `timm` 版本为 `0.5.4` 。 采用如下代码记载模型

```
1   import timm
2   import torch
3   import os
4
```

```
swin_transformer_model = timm.create_model('swin_base_patch4_window7_224',
pretrained=True)
```

如果下载速度较慢，请打开log中的下载网址采用浏览器下载好以后，再放到给出的 `cache` 文件夹中。 以上代码即完成了模型的预加载。

## 4.2 验证分类模型精度

```
1   valdir = 'YourPath/ImageNet/val'
2   normalize = transforms.Normalize(mean=[0.485, 0.456, 0.406],std=[0.229,
3   0.224, 0.225])
4   val_loader = torch.utils.data.DataLoader(
5       datasets.ImageFolder(valdir, transforms.Compose([
6           transforms.Resize(256, 3),
7           transforms.CenterCrop(224),
8           transforms.ToTensor(),
9           normalize,
10      ])),batch_size=128, shuffle=False, num_workers=12,
11  pin_memory=True)
12  import tqdm
13  def test(model, device, test_loader):
14      old_training_state = model.training
15      model.eval()
16      test_loss = 0
17      correct = 0
18      lossLayer = torch.nn.CrossEntropyLoss(reduction='sum')
19      for data, target in tqdm.tqdm(test_loader):
20          data, target = data.to(device), target.to(device)
21          with torch.no_grad():
22              output = model(data)
23          test_loss += lossLayer(output, target).item()
24          pred = output.argmax(dim=1, keepdim=True)
25          correct += pred.eq(target.view_as(pred)).sum().item()
26      test_loss /= len(test_loader.dataset)
27      model.train(old_training_state)
28      print('\nTest set: Average loss: {:.4f}, Accuracy: {:.3f}%\n'.format(
29          test_loss, 100. * correct / len(test_loader.dataset)
30
31      ))
32  device = torch.device('cuda')
33  swin_transformer_model.to(device)
    # Swin-B 85.2 22kto1k
    test(swin_transformer_model, device, val_loader)
```

经过测试，与原论文中给出精度仅存在细微差异，差了0.1个点，主要是输入图像插值方式带来的，这不是本文重点故不再修改。

## 4.3 初步导出

```
1   dummpy_input = torch.rand(1,3,224,224)
2   swin_transformer_model.to('cpu')
3   if not os.path.exists('swin-b-22kto1l.onnx'):
4       torch.onnx.export(swin_transformer_model, (dummpy_input, ), 'swin-b-
```

```
 5    22kto11.onnx',
 6                       input_names=['images',],
 7                       output_names=['output'],
 8                       dynamic_axes={
 9                           'images': {
10                               0:'batch',
11                           },
12                           'output':{
13                               0:'batch',
14                           }
15                       },
16                       opset_version=11,
                         verbose = True)
```

由于设置了 `verbose = True` ，以上代码会打印出很长的log。

## 4.4 修改模型代码

按照关键词 `onnx::Shape` 来搜索shape操作的代码位置，找到以后按照原则1，3修改代码。 在 `${PYENV}/lib/python3.9/site-packages/timm/models/layers/patch_embed.py` ，修改 `PatchEmbed` 的 `forward` 函数，主要是把 `flatten` 转换为 `view(-1, ...)` 。、 **${PYENV}表示 python环境所在目录** 原版为（下文不再展示原版）

```
 1   def forward(self, x):
 2           B, C, H, W = map(int, x.shape)
 3           _assert(H == self.img_size[0], f"Input image height ({H}) doesn't
 4   match model ({self.img_size[0]}).")
 5           _assert(W == self.img_size[1], f"Input image width ({W}) doesn't
 6   match model ({self.img_size[1]}).")
 7           x = self.proj(x)
 8           if self.flatten:
 9               x = x.flatten(2).transpose(1, 2)  # BCHW -> BNC
             x = self.norm(x)
             return x
```

修改为

```
 1   def forward(self, x):
 2           B, C, H, W = map(int,x.shape)
 3           _assert(H == self.img_size[0], f"Input image height ({H}) doesn't
 4   match model ({self.img_size[0]}).")
 5           _assert(W == self.img_size[1], f"Input image width ({W}) doesn't
 6   match model ({self.img_size[1]}).")
 7           x = self.proj(x)
 8           if self.flatten:
 9               _, newC, newH, newW = map(int, x.shape)
10               x = x.view(-1, newC, newH * newW).transpose(1, 2)  # BCHW ->
     BNC
             x = self.norm(x)
             return x
```

剩下shape操作基本都在 `${PYENV}/lib/python3.9/site-packages/timm/models/models/swin_transformer.py` ，先以关键词 `.shape` 在该文件中搜索，修改完毕后再看剩余 `window_partition` 和 `window_reverse` 修改如下

```
def window_partition(x, window_size: int):
    """
    Args:
        x: (B, H, W, C)
        window_size (int): window size
    Returns:
        windows: (num_windows*B, window_size, window_size, C)
    """
    B, H, W, C = map(int, x.shape)
    x = x.view(-1, H // window_size, window_size, W // window_size,
window_size, C)
    windows = x.permute(0, 1, 3, 2, 4, 5).contiguous().view(-1,
window_size, window_size, C)
    return windows
@register_notrace_function  # reason: int argument is a Proxy
def window_reverse(windows, window_size: int, H: int, W: int):
    """
    Args:
        windows: (num_windows*B, window_size, window_size, C)
        window_size (int): Window size
        H (int): Height of image
        W (int): Width of image
    Returns:
        x: (B, H, W, C)
    """
    B = int(windows.shape[0] / (H * W / window_size / window_size))
    C = int(windows.numel() // B // window_size // window_size // (H //
window_size * W // window_size))
    x = windows.view(-1, H // window_size, W // window_size, window_size,
window_size, C)
    x = x.permute(0, 1, 3, 2, 4, 5).contiguous().view(B, H, W, -1)
    return x
```

`WindowAttention` 类的 `forward` 函数修改为

```
def forward(self, x, mask: Optional[torch.Tensor] = None):
        """
        Args:
            x: input features with shape of (num_windows*B, N, C)
            mask: (0/-inf) mask with shape of (num_windows, Wh*Ww, Wh*Ww)
or None
        """
        B_, N, C = map(int, x.shape)
        qkv = self.qkv(x).reshape(B_, N, 3, self.num_heads, C //
self.num_heads).permute(2, 0, 3, 1, 4)
        q, k, v = qkv.unbind(0)  # make torchscript happy (cannot use
tensor as tuple)
        q = q * self.scale
        attn = (q @ k.transpose(-2, -1))
        relative_position_bias =
self.relative_position_bias_table[self.relative_position_index.view(-1)].vie
```

```
17              self.window_size[0] * self.window_size[1],
18    self.window_size[0] * self.window_size[1], -1)  # Wh*Ww,Wh*Ww,nH
19          relative_position_bias = relative_position_bias.permute(2, 0,
20    1).contiguous()  # nH, Wh*Ww, Wh*Ww
21          attn = attn + relative_position_bias.unsqueeze(0)
22          if mask is not None:
23              nW = int(mask.shape[0])
24              attn = attn.view(-1, nW, self.num_heads, N, N) +
25    mask.unsqueeze(1).unsqueeze(0)
26              attn = attn.view(-1, self.num_heads, N, N)
27              attn = self.softmax(attn)
28          else:
29              attn = self.softmax(attn)
          attn = self.attn_drop(attn)
          x = (attn @ v).transpose(1, 2)
          last_dim = int(x.numel() // B_ // N)
          x = x.reshape(-1, N, last_dim)
          x = self.proj(x)
          x = self.proj_drop(x)
          return x
```

`SwinTransformerBlock` 的 `forward` 函数修改如下:

```
1   def forward(self, x):
2          H, W = self.input_resolution
3          B, L, C = map(int, x.shape)
4          _assert(L == H * W, "input feature has wrong size")
5          shortcut = x
6          x = self.norm1(x)
7          x = x.view(-1, H, W, C)
8          # cyclic shift
9          if self.shift_size > 0:
10             shifted_x = torch.roll(x, shifts=(-self.shift_size, -
11   self.shift_size), dims=(1, 2))
12         else:
13             shifted_x = x
14         # partition windows
15         x_windows = window_partition(shifted_x, self.window_size)  #
16   nW*B, window_size, window_size, C
17         x_windows = x_windows.view(-1, self.window_size *
18   self.window_size, C)  # nW*B, window_size*window_size, C
19         # W-MSA/SW-MSA
20         attn_windows = self.attn(x_windows, mask=self.attn_mask)  # nW*B,
21   window_size*window_size, C
22         # merge windows
23         attn_windows = attn_windows.view(-1, self.window_size,
24   self.window_size, C)
25         shifted_x = window_reverse(attn_windows, self.window_size, H, W)
26   # B H' W' C
27         # reverse cyclic shift
28         if self.shift_size > 0:
29             x = torch.roll(shifted_x, shifts=(self.shift_size,
30   self.shift_size), dims=(1, 2))
         else:
             x = shifted_x
         x = x.view(-1, H * W, C)
```

```
                # FFN
                x = shortcut + self.drop_path(x)
                x = x + self.drop_path(self.mlp(self.norm2(x)))
                return x
```

`PatchMerging` 的 `forward` 函数修改如下：

```
 1  def forward(self, x):
 2          """
 3          x: B, H*W, C
 4          """
 5          H, W = self.input_resolution
 6          B, L, C = map(int, x.shape)
 7          _assert(L == H * W, "input feature has wrong size")
 8          _assert(H % 2 == 0 and W % 2 == 0, f"x size ({H}*{W}) are not
 9  even.")
10          x = x.view(B, H, W, C)
11          x0 = x[:, 0::2, 0::2, :]  # B H/2 W/2 C
12          x1 = x[:, 1::2, 0::2, :]  # B H/2 W/2 C
13          x2 = x[:, 0::2, 1::2, :]  # B H/2 W/2 C
14          x3 = x[:, 1::2, 1::2, :]  # B H/2 W/2 C
15          x = torch.cat([x0, x1, x2, x3], -1)  # B H/2 W/2 4*C
16          x = x.view(-1, H // 2 * W // 2, 4 * C)  # B H/2*W/2 4*C
17          x = self.norm(x)
18          x = self.reduction(x)
                return x
```

## 4.5 重新生成onnx并检查

1. 如果是 `.py` 就重新执行 `python xxx.py` ，如果是 `jupyter notebook` ，则需要重启内核。

2. `pip install netron` 安装 `netron`

3. 运行 `netron swin-b-22kto1l.onnx` 查看生成的onnx，查看模型整体简洁性（主观）。

经过查看onnx结构，发现修改的onnx结构较为简洁，且支持batch size的动态shape。

## 4.6 在TensorRT上验证

这里使用tensorRT的封装版本 `trtpy` ，主要解决版本隔离问题和安装繁琐问题。 安装和使用详情请查看https://zhuanlan.zhihu.com/p/462980738

```
 1  import trtpy
 2  import numpy as np
 3  def replace_suffix(path_name, new_extname = '.trtmodel'):
 4      return os.path.splitext(path_name)[0] + new_extname
 5  onnx_filepath = 'swin-b-22kto1l.onnx'
 6  engine_file = replace_suffix(onnx_filepath)
 7  if not os.path.exists(engine_file):
 8      trtpy.compile_onnx_to_file(128, onnx_filepath, engine_file,
 9  max_workspace_size=1 << 30) # 1 << 30 表示 1GiB
10  engine = trtpy.load(engine_file)
11  engine.print()
```

```python
12    def test_trt(engine, model, device, test_loader):
13        old_training_state = model.training
14        model.eval()
15        test_loss = 0
16        correct = 0
17        delta_max = 0
18        lossLayer = torch.nn.CrossEntropyLoss(reduction='sum')
19        for data, target in tqdm.tqdm(test_loader):
20            data, target = data.to(device), target.to(device)
21            with torch.no_grad():
22                output = model(data)
23            output_engine = engine(data)
24            delta_max = max(delta_max, torch.abs(output_engine -
25    output).max().item())
26
27            test_loss += lossLayer(output, target).item()
28            pred = output.argmax(dim=1, keepdim=True)
29            correct += pred.eq(target.view_as(pred)).sum().item()
30        print('delta_max:', delta_max)
31        test_loss /= len(test_loader.dataset)
32        model.train(old_training_state)
33        print('\nTest set: Average loss: {:.4f}, Accuracy: {:.3f}%\n'.format(
34            test_loss, 100. * correct / len(test_loader.dataset)
35
36        ))
      swin_transformer_model.to(device)
      test_trt(engine, swin_transformer_model, device, val_loader)
```

最后输出

```
1    delta_max: 0.0006237030029296875
2    Test set: Average loss: 0.6448, Accuracy: 85.126%
```

以上实验证明了新导出的onnx支持动态batch size，并且TensorRT与torch生成结果误差很小，abs diff 可以控制在0.0006以内。

## 4.7 小结

通过对SwinTransformer模型代码的修改，使得onnx导出更为简洁。因为该模型较为简单，在导出时并没有使用 `onnxsim` 简化模型，下一章对 `DETR` 的导出将会看到所有工具的综合运用进行Debug。 此外，Transformer在CV领域已经可以很方便的落地了，**并且是支持动态batchsize的**。个人的直观感受时较为部署时采用SwinTransformer占用显存比较大，不过 `85.2%` 的精度也算性价比较高的模型。

# 五、案例2-DETR在TensorRT上的部署

## 5.1 下载预训练模型并验证模型效果

这里采用集成的 `mmdetection` 库中的DETR。github地址为https://github.com/open-mmlab/mmdetection 为了方便修改代码， `mmdetection` 需要采用clone仓库， `python setup.py`

`develop` 方式安装，即

```
1   git clone https://github.com/open-mmlab/mmdetection.git
2   cd mmdetection
3   pip install -r requirements/build.txt
4   pip install -v -e .  # or "python setup.py develop"
```

1. 在此之前读者还请阅读 https://github.com/open-mmlab/mmdetection/blob/master/docs/en/get_started.md 了解全部安装步骤。

2. 从 https://github.com/open-mmlab/mmdetection/tree/master/configs/detr 下载预训练权重并记录path备用。

3. 安装onnxsim，polygraphy，onnxruntime库

修改 `mmdetection/tools/deployment/pytorch2onnx.py` ，让其支持再 `skip-postprocess` 条件下支持 `onnxsim` 在判断语句那里加上如下代码

```
1       if skip_postprocess:
2           warnings.warn('Not all models support export onnx without post '
3                         'process, especially two stage detectors!')
4           model.forward = model.forward_dummy
5           torch.onnx.export(
6               model,
7               one_img,
8               output_file,
9               input_names=['input'],
10              export_params=True,
11              keep_initializers_as_inputs=True,
12              do_constant_folding=True,
13              verbose=show,
14              opset_version=opset_version)
15          print(f'Successfully exported ONNX model without '
16                f'post process: {output_file}')
17          # 添加的代码 开始
18          # get the custom op path
19          ort_custom_op_path = ''
20          try:
21              from mmcv.ops import get_onnxruntime_op_path
22              ort_custom_op_path = get_onnxruntime_op_path()
23          except (ImportError, ModuleNotFoundError):
24              warnings.warn('If input model has custom op from mmcv, \\
25                  you may have to build mmcv with ONNXRuntime from
26  source.')
27          ort_custom_op_path = ort_custom_op_path if
28  len(ort_custom_op_path) > 0 else None
29          if do_simplify:
30              import onnxsim
31              from mmdet import digit_version
32              min_required_version = '0.3.0'
33              assert digit_version(onnxsim.__version__) >= digit_version(
34                  min_required_version
35              ), f'Requires to install onnx-simplify>=
36  {min_required_version}'
37              input_dic = {'input': img_list[0].detach().cpu().numpy()}
```

```
38              model_opt, check_ok = onnxsim.simplify(
39                  output_file,
40                  input_data=input_dic,
41                  custom_lib=ort_custom_op_path,
42                  dynamic_input_shape=dynamic_export)
43              if check_ok:
44                  onnx.save(model_opt, output_file)
45                  print(f'Successfully simplified ONNX model:
46  {output_file}')
47              else:
                    warnings.warn('Failed to simplify ONNX model.')
                print(f'Successfully exported ONNX model: {output_file}')
                # 添加的代码结束
            return
```

新建一个bash脚本 `export_detr_onnx.sh`

```
1  #!/bin/bash
2  python tools/deployment/pytorch2onnx.py
   configs/detr/detr_r50_8x2_150e_coco.py
   detr_r50_8x2_150e_coco_20201130_194835-2c4b8974.pth --skip-postprocess --
   show --simplify
```

执行如下指令：

```
1  bash export_detr_onnx.sh &> export_log.txt
```
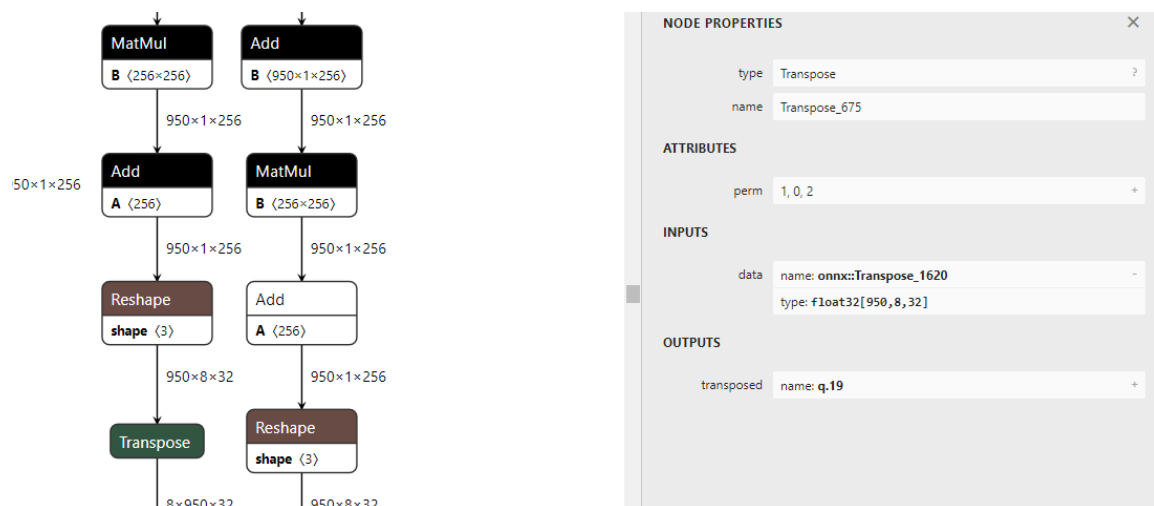
通过以上步骤即可得到一个精简版的onnx，此处并没有修改任何模型代码，因为修改模型代码的目的是为了处理inplace操作和实现动态batchsize，在mmdetection中DETR，并没有inplace操作，且动态batchsize操作实现起来很难，下面会给出解释。 使用netron可视化onnx后发现很多冗余算子，这里需要使用 `onnx_graphsurgeon` 清理一下。

```
1  import onnx_graphsurgeon as gs
2  import onnx
3  import numpy as np
4  def cleanOnnx():
5      onnx_save_path = "tmp.onnx"
6      graph = gs.import_onnx(onnx.load(onnx_save_path))
7      graph.cleanup()
8      onnx_save_path = "tmp_new.onnx"
9      onnx.save(gs.export_onnx(graph), onnx_save_path)
10 cleanOnnx()
```

## 5.2 动态batchsize实现很难

```
25                          [nb_dec, bs, num_query, 4].
26              """
              num_levels = len(feats)
              img_metas_list = [img_metas for _ in range(num_levels)]
              # return multi_apply(self.forward_single, feats, img_metas_list)
              all_cls_scores_list, all_bbox_preds_list =
    multi_apply(self.forward_single, feats, img_metas_list)
              image_metas_for_onnx = img_metas_list[0:1][0]
              image_metas_for_onnx[0]['img_shape_for_onnx'] =
    torch.as_tensor(img_metas[0]['img_shape'])[:2]
              return self.onnx_export(all_cls_scores_list, all_bbox_preds_list,
    image_metas_for_onnx)
```

`Line24-Line28` 是新增的，用于解码中间信息，直接输出坐标信息和分类信息。

## 5.4 使用onnxruntime验证

1. 测试图片如下：

读者可以右键存为 `car.jpg`

1. 分析模型config文件，得到数据预处理步骤

2. 给出onnxruntime验证代码

```
1   import onnxruntime as rt
2   import cv2
3   def get_onnx_runner(onnx_filepath):
4       sess = rt.InferenceSession(onnx_filepath)
5       input_names = [item.name for item in sess.get_inputs()]
6       label_names = [item.name for item in sess.get_outputs()]
7       def runner(input_tensors):
8           nonlocal label_names
9           nonlocal input_names
10          pred_onnx = sess.run(label_names, dict(zip(input_names,
11  input_tensors)))
12          return dict(zip(label_names,pred_onnx))
13      return runner
14  def _normalize(img, mean, std):
15      img = img.astype(np.float32) / 255
16      mean = np.array(mean, dtype=np.float32).reshape(1, 1, 3) / 255
17      std = np.array(std, dtype=np.float32).reshape(1, 1, 3) / 255
18      img = (img - mean) / std
19
20      return img
21  def preprocess(img):
22      img = cv2.resize(img, dsize=(1216,800))
23      mean,std = [[103.53, 116.28, 123.675], [57.375, 57.12, 58.395]]
24      img = _normalize(img, mean, std)
25      img = np.transpose(img, axes=[2, 0, 1])
26      img = np.expand_dims(img, 0)
27      return img
28  runner = get_onnx_runner("./tmp_new.onnx")
29  ori_image = cv2.imread('../car.jpg')
30  img = cv2.cvtColor(ori_image, cv2.COLOR_BGR2RGB)
```

```python
31   img = preprocess(img)
32   outs = runner([img, ])
33   print(outs.keys())
34   # det_bboxes[x1,y1,x2,y2,scores], det_labels[]
35   det_bboxes, scores = outs['3288'][:, :, :4], outs['3288'][:, :, 4:]
36   det_labels = outs['3209']
37   det_bboxes = np.round(det_bboxes).astype(np.int32)
38   from torchvision.ops import nms as torch_nms
39   import torch
40   det_bboxes = det_bboxes[0]
41   det_labels = det_labels[0]
42   scores = scores[0]
43   scores = scores[:, 0]
44   def nms(ori_dets):
45       tmp_arr = np.array(ori_dets, dtype=np.float32)
46       dets, score = tmp_arr[:,:4], tmp_arr[:,4]
47       dets_tensor = torch.as_tensor(dets)
48       score_tensor = torch.as_tensor(score)
49       keep_index = torch_nms(dets_tensor, score_tensor, 0.4)
50       new_dets = [ori_dets[index] for index in keep_index]
51       return new_dets
52   def decode_infer(scores, det_bboxes, det_labels, score_threshold):
53       ret_dict = dict()
54       max_scores = scores
55       nonzero_index  = np.nonzero(max_scores > score_threshold)
56   [0].tolist()
57       if nonzero_index:
58           bbox_pred = det_bboxes[nonzero_index, :]
59           score_indexes = max_scores[nonzero_index]
60           score_indexes = score_indexes[..., None].astype(np.float32)
61           label_indexes = det_labels[nonzero_index]
62           print(label_indexes)
63           label_indexes = label_indexes[..., None].astype(np.float32)
64           max_label_indexes = label_indexes[:, 0]
65           bbox = np.concatenate([bbox_pred, score_indexes, label_indexes],
66   axis = -1)
67           for i,index in enumerate(max_label_indexes):
68               ret_list = ret_dict.setdefault(index, list())
69               ret_list.append(bbox[i])
70       for k,v in ret_dict.items():
71           ret_dict[k] = nms(v)
72       return ret_dict
73   ret_dict = decode_infer(scores, det_bboxes, det_labels, 0.5)
74   labels = ("person", "bicycle", "car", "motorcycle", "airplane", "bus",
75   "train", "truck", "boat", "traffic light",
76               "fire hydrant", "stop sign", "parking meter", "bench",
77   "bird", "cat", "dog", "horse", "sheep", "cow",
78               "elephant", "bear", "zebra", "giraffe", "backpack",
79   "umbrella", "handbag", "tie", "suitcase", "frisbee",
80               "skis", "snowboard", "sports ball", "kite", "baseball bat",
81   "baseball glove", "skateboard", "surfboard",
82               "tennis racket", "bottle", "wine glass", "cup", "fork",
83   "knife", "spoon", "bowl", "banana", "apple",
84               "sandwich", "orange", "broccoli", "carrot", "hot dog",
85   "pizza", "donut", "cake", "chair", "couch",
86               "potted plant", "bed", "dining table", "toilet", "tv",
87   "laptop", "mouse", "remote", "keyboard", "cell phone",
```

```
 88                "microwave", "oven", "toaster", "sink", "refrigerator",
 89    "book", "clock", "vase", "scissors", "teddy bear",
 90                "hair drier", "toothbrush")
 91    color_list = (
 92            (216 , 82 , 24),
 93            (236 ,176 , 31),
 94            (125 , 46 ,141),
 95            (118 ,171 , 47),
 96            ( 76 ,189 ,237),
 97            (238 , 19 , 46),
 98            ( 76 , 76 , 76),
 99            (153 ,153 ,153),
100            (255 ,  0 ,  0),
101            (255 ,127 ,  0),
102            (190 ,190 ,  0),
103            (  0 ,255 ,  0),
104            (  0 ,  0 ,255),
105            (170 ,  0 ,255),
106            ( 84 , 84 ,  0),
107            ( 84 ,170 ,  0),
108            ( 84 ,255 ,  0),
109            (170 , 84 ,  0),
110            (170 ,170 ,  0),
111            (170 ,255 ,  0),
112            (255 , 84 ,  0),
113            (255 ,170 ,  0),
114            (255 ,255 ,  0),
115            (  0 , 84 ,127),
116            (  0 ,170 ,127),
117            (  0 ,255 ,127),
118            ( 84 ,  0 ,127),
119            ( 84 , 84 ,127),
120            ( 84 ,170 ,127),
121            ( 84 ,255 ,127),
122            (170 ,  0 ,127),
123            (170 , 84 ,127),
124            (170 ,170 ,127),
125            (170 ,255 ,127),
126            (255 ,  0 ,127),
127            (255 , 84 ,127),
128            (255 ,170 ,127),
129            (255 ,255 ,127),
130            (  0 , 84 ,255),
131            (  0 ,170 ,255),
132            (  0 ,255 ,255),
133            ( 84 ,  0 ,255),
134            ( 84 , 84 ,255),
135            ( 84 ,170 ,255),
136            ( 84 ,255 ,255),
137            (170 ,  0 ,255),
138            (170 , 84 ,255),
139            (170 ,170 ,255),
140            (170 ,255 ,255),
141            (255 ,  0 ,255),
142            (255 , 84 ,255),
143            (255 ,170 ,255),
144            ( 42 ,  0 ,  0),
```

```
145              (  84 ,   0 ,   0),
146              (127 ,   0 ,   0),
147              (170 ,   0 ,   0),
148              (212 ,   0 ,   0),
149              (255 ,   0 ,   0),
150              (  0 ,  42 ,   0),
151              (  0 ,  84 ,   0),
152              (  0 , 127 ,   0),
153              (  0 , 170 ,   0),
154              (  0 , 212 ,   0),
155              (  0 , 255 ,   0),
156              (  0 ,   0 ,  42),
157              (  0 ,   0 ,  84),
158              (  0 ,   0 , 127),
159              (  0 ,   0 , 170),
160              (  0 ,   0 , 212),
161              (  0 ,   0 , 255),
162              (  0 ,   0 ,   0),
163              ( 36 ,  36 ,  36),
164              ( 72 ,  72 ,  72),
165              (109 , 109 , 109),
166              (145 , 145 , 145),
167              (182 , 182 , 182),
168              (218 , 218 , 218),
169              (  0 , 113 , 188),
170              ( 80 , 182 , 188),
171              (127 , 127 ,   0),
172          )
173    frame = np.copy(ori_image)
174    height = img.shape[2]
175    width = img.shape[3]
176    for v in ret_dict.values():
177        for left, top, right, bottom, score, label_index in v:
178            label_index = int(label_index)
179            new_left = int(left / width * frame.shape[1])
180            new_top = int(top / height * frame.shape[0])
181            new_right = int(right / width * frame.shape[1])
182            new_bottom = int(bottom / height * frame.shape[0])
183            cv2.rectangle(frame, (new_left, new_top), (new_right,
184    new_bottom), color = color_list[label_index])
185            # 写字
            ret, baseline = cv2.getTextSize(labels[label_index],
    cv2.FONT_HERSHEY_SIMPLEX, 0.8, 1)
            cv2.rectangle(frame, (new_left, new_top - ret[1] - baseline),

                          (new_left + ret[0], new_top),
    color_list[label_index], -1)
            cv2.putText(frame, labels[label_index], (new_left, new_top -
    baseline),

                        cv2.FONT_HERSHEY_SIMPLEX, 0.8, (255, 255, 255), 1)
    import matplotlib.pyplot as plt
    h,w = frame.shape[:2]
    plt.figure(figsize=(w / 50, h / 50))
    plt.imshow(frame[..., ::-1])
    plt.show()
```

最终效果如下图：

## 5.5 TensorRT上验证

这里本文继续采用trtpy，注意由于本文的代码时放在一个 `jupyter notebook` 中，所以复用了上一小节的变量。 `det_labels, det_bboxes_with_scores` 这两个顺序可能会反，读者需要自行调换。

```
import trtpy
import numpy as np
import os
def replace_suffix(path_name, new_extname = '.trtmodel'):
    return os.path.splitext(path_name)[0] + new_extname
onnx_filepath = 'tmp_new.onnx'
engine_file = replace_suffix(onnx_filepath)
trtpy.compile_onnx_to_file(1, onnx_filepath, engine_file,
max_workspace_size=1 << 30)
engine = trtpy.load(engine_file)
engine.print()
data = torch.as_tensor(img).to(torch.float32).cuda()
output_engine = engine(data)
print(output_engine)
det_labels, det_bboxes_with_scores = output_engine
det_bboxes, scores = det_bboxes_with_scores[:, :, :4],
det_bboxes_with_scores[:, :, 4:]
det_bboxes = det_bboxes[0].round().to(torch.int32)
det_labels = det_labels[0].round().to(torch.int32)
scores = scores[0]
scores = scores[:, 0]
ret_dict_trt = decode_infer(scores.cpu().numpy(),
det_bboxes.cpu().numpy(), det_labels.cpu().numpy(), 0.5)
frame = np.copy(ori_image)
height = img.shape[2]
width = img.shape[3]
for v in ret_dict_trt.values():
    for left, top, right, bottom, score, label_index in v:
        label_index = int(label_index)
        new_left = int(left / width * frame.shape[1])
        new_top = int(top / height * frame.shape[0])
        new_right = int(right / width * frame.shape[1])
        new_bottom = int(bottom / height * frame.shape[0])
        cv2.rectangle(frame, (new_left, new_top), (new_right,
new_bottom), color = color_list[label_index])
        ret, baseline = cv2.getTextSize(labels[label_index],
cv2.FONT_HERSHEY_SIMPLEX, 0.8, 1)
        cv2.rectangle(frame, (new_left, new_top - ret[1] - baseline),

                      (new_left + ret[0], new_top),
color_list[label_index], -1)
        cv2.putText(frame, labels[label_index], (new_left, new_top -
baseline),

                    cv2.FONT_HERSHEY_SIMPLEX, 0.8, (255, 255, 255), 1)
h,w = frame.shape[:2]
```

```
        plt.figure(figsize=(w / 50, h / 50))
        plt.imshow(frame[..., ::-1])
        plt.show()
```

`print(output_engine)` 输出全0，会让人很头疼，怀疑是不是哪一步搞错了。这个时候就需要输出每一层的中间结果打印了。

## 5.6 调试模型中间结果

由于onnxruntime输出结果时正确的，那么我们有理由相信其中间结果也是正确的。进而可以转换为TensorRT与onnxruntime中间层不一致问题。 此时，我们可以借助 `polygraphy` 工具来完成对中间层输出的校对。 `polygraphy` 安装命令， `polygraphy` 作者安装的版本为 `0.35.2`

```
1   python -m pip install colored polygraphy --extra-index-url
    https://pypi.ngc.nvidia.com
```

将如下代码存为 `debug_detr.sh`

```
1   polygraphy run tmp_new.onnx --onnxrt --onnx-outputs mark all --save-
2   results=onnx_out.json
    polygraphy run tmp_new.onnx --trt --validate --trt-outputs mark all --
    save-results=trt_out.json
```

运行时你会发现报错，因为这个虚拟环境并没有安装原版的TensorRT，需要做一定修改。错误如下

```
1   [!] Module: 'tensorrt' is required but could not be imported.
2       You can try setting POLYGRAPHY_AUTOINSTALL_DEPS=1 in your environment
3   variables to allow Polygraphy to automatically install missing modules.
        Note that this may cause existing modules to be overwritten - hence,
    it may be desirable to use a Python virtual environment or container.
```

打开 `${PYENV}/lib/python3.9/site-packages/polygraphy/mod/importer.py` ，找到 `lazy_import` 函数，在第一句assert语句下面 加上这句话

```
1   if name == "tensorrt":
2           name = "trtpy.tensorrt"
```

此外还可能碰到说tensorrt没有 `__version__` ，这个时候需要修改trtpy中的 `tensorrt` 。 找到 `${PYENV}/lib/python3.9/site-packages/trtpy/tensorrt/__init__.py` 搜索 `__version__` 该行然后去掉注释。 经过以上过程，重新执行上面的生成脚本得到两个json文件，采用如下代码比较中间结果

```
1   import numpy as np
2   from polygraphy.json import load_json
3   info_onnx = load_json('onnx_out.json')
4   # print(info_onnx)
5   info_trt = load_json('trt_out.json')
```

```
 6    # print(info)
 7    runners_trt = list(info_trt.keys())
 8    runners_onnx = list(info_onnx.keys())
 9    print('onnx:', len(info_onnx.__getitem__(runners_onnx[0])[0]))
10    print('tensorrt:', len(info_trt.__getitem__(runners_trt[0])[0]))
11    outputs_trt = info_trt['lst'][0][1][0]['outputs']
12    outputs_onnx = info_onnx['lst'][0][1][0]['outputs']
13    for layer in outputs_onnx:
14        if layer in outputs_trt:     # 只分析相同名称的节点输出
15            value_onnx = outputs_onnx[layer]['values']
16            value_trt = outputs_trt[layer]['values']
17            try:
18                np.testing.assert_allclose(value_onnx, value_trt, 0.001,
19    0.001)
20            except Exception as e:
21                print('-' * 20)
22                print(layer, value_onnx.shape, value_trt.shape)
                  print(e)
```

经过以上步骤，会输出很多不一致的中间层，经过作者的搜索引擎大法和经验，首先定位到
Gather算子，可能是TensorRT的Gather算子不支持负值索引，因此修改cleanOnnx函数修改原来
的onnx。至于为什么是5，因为之前的形状是 `6x?x?x?x` ，6-1 = 5.

```
 1    import onnx_graphsurgeon as gs
 2    import onnx
 3    import numpy as np
 4    def cleanOnnx():
 5        onnx_save_path = "tmp.onnx"
 6        graph = gs.import_onnx(onnx.load(onnx_save_path))
 7        for node in graph.nodes:
 8            if node.name.startswith('Gather'):
 9                indices = node.inputs[1]
10                if isinstance(indices, gs.Constant):
11                    print(node.name)
12                    print(indices.values)
13                # 去掉负索引
14                if isinstance(indices, gs.Constant) and indices.values == -1:
15                    node.inputs[1].values = np.int64(5)
16                if isinstance(indices, gs.Constant):
17                    print(indices.values)
18        graph.cleanup()
19        onnx_save_path = "tmp_new.onnx"
20        onnx.save(gs.export_onnx(graph), onnx_save_path)
```

## 5.7 trtpy的小\bug

经过上面的修改，tensorRT矩形框输出正常了，但是分类结果全为0，原因在于trtpy将输出的
int32类型按照float32解析了，从而造成解码错误。一个讨巧的方法是在
`mmdetection/mmdet/models/dense_heads/detr_head.py` 中 `DETRHead` 类的 `onnx_export` 方法的
返回值加上类型转换，即将 `return det_bboxes, det_labels` 修改为 `return det_bboxes,`
`det_labels.to(torch.float32)` ，然后重新生成 `tmp.onnx` 和 `tmp_new.onnx` ，使用原生
tensorRT的用户可以自行忽略这一点。

## 六、 总结

本文系统梳理了onnx导出的原则和具体调试技巧，希望能给遇到相似错误的读者给以启发。作者会继续研究DETR的动态batch size问题，欢迎有兴趣的小伙伴一起讨论。

## 七、 参考链接

- https://zhuanlan.zhihu.com/p/436017991

- https://github.com/DataXujing/TensorRT-DETR

---

最后更新: March 21, 2024

创建日期: March 21, 2024