

# SD torch转onnx

## 分析

### 入口：compute\_score

for函数内产生新图片的函数是

```
1 | hk.process()
```

往上找：








```
1 | from canny2image_TRT import hackathon
2 | hk = hackathon()
```

定位到canny2image

### canny2img.py

上网查阅资料，这个函数应该是直接用于：

- 1.canny边缘检测
- 2.生成条件输入
- 3.控制生成结果

原图	条件输入	输出
无		
无		
		

根据compute\_score.py中的调用：

```

1 new_img = hk.process(img,
2     "a bird",
3     "best quality, extremely detailed",
4     "longbody, lowres, bad anatomy, bad hands, missing fingers",
5     1,
6     256,
7     20,
8     False,
9     1,
10    9,
11    2946901,
12    0.0,
13    100,
14    200)

```

分析其核心函数process输入参数如下：

- input\_image: 输入图片数据，重要
- prompt: 描述图像生成的文本提示
- a\_prompt: "additional prompt"或者是"activation prompt"，是与prompt配合的附加文本提示，用于调整细节
- n\_prompt: "negative prompt"，负面提示，描述模型在生成中需要规避的内容
- num\_samples: 生成样本数量
- image\_resolution: 生成图片的解析度
- ddim\_steps: DDIM (Denoising Diffusion Implicit Models) 算法生成图像的迭代次数
- guess\_mode: 是否开启猜测模式
- strength: 控制输入图像的影响力。若生成是基于已有图像，`strength` 决定了参考图像与生成结果的结合程度。`strength=0` 时，完全不参考输入图像，`strength=1` 时，完全参考输入图像。
- scale: 控制文本提示对生成图像的影响程度。在大多数图像生成模型中，`scale` 是 `prompt` 的影响力调节因子。较大的 `scale` 会让图像更强烈地遵循提示内容，较小的 `scale` 会使生成过程更自由。
- seed: 是否需要设置随机数
- eta: 通常与采样策略中的噪声调节有关。DDIM采样方法中，`eta` 调节生成过程中噪声的大小，影响生成图像的细节与噪声水平。
- low\_threshold 和 high\_threshold: 限制图像的亮度、对比度或者像素强度范围，控制生成图像的细节和质量。与图像生成过程中的图像过滤、边缘控制或去噪过程有关。

接下来对process方法进行主句分析

## 输入图片处理

```

1 img = resize_image(HWC3(input_image), image_resolution)
2 H, W, C = img.shape
3
4 detected_map = self.apply_canny(img, low_threshold, high_threshold)
5 detected_map = HWC3(detected_map)
6
7 control = torch.from_numpy(detected_map.copy()).float().cuda() / 255.0
8 control = torch.stack([control for _ in range(num_samples)], dim=0)
9 control = einops.rearrange(control, 'b h w c -> b c h w').clone()

```

- 先看HWC3，路径/annotator/util.py

```

1 def HWC3(x):

```

```

2     assert x.dtype == np.uint8
3     if x.ndim == 2:
4         x = x[:, :, None]
5     assert x.ndim == 3
6     H, W, C = x.shape
7     assert C == 1 or C == 3 or C == 4
8     if C == 3:
9         return x
10    if C == 1:
11        return np.concatenate([x, x, x], axis=2)
12    if C == 4:
13        color = x[:, :, 0:3].astype(np.float32)
14        alpha = x[:, :, 3:4].astype(np.float32) / 255.0
15        y = color * alpha + 255.0 * (1.0 - alpha)
16        y = y.clip(0, 255).astype(np.uint8)
17    return y

```

它应该是一个生成RGB图像的方法。

当为灰度图像（通道数1），将单通道扩展为3通道。

如果是RGBA图像（通道数4），带有透明通道，则使用透明通道融合RGB颜色和背景。简单来说就是透明通道就是设置透明度，当透明度越高RGB颜色越明显，透明度越低，背景越明显（也就是白色255）。

- 再看resize\_image()方法，也位于/annotator/util.py中

```

1  def resize_image(input_image, resolution):
2      H, W, C = input_image.shape
3      H = float(H)
4      W = float(W)
5      k = float(resolution) / min(H, W)
6      H *= k
7      W *= k
8      H = int(np.round(H / 64.0)) * 64
9      W = int(np.round(W / 64.0)) * 64
10     img = cv2.resize(input_image, (W, H), interpolation=cv2.INTER_LANCZOS4
11     if k > 1 else cv2.INTER_AREA)
12     return img

```

k是输出图片与输入图片最小边比值，保证整形后的图片最小边至少满足输出图片解析度。即根据目标尺寸调整输入图片大小

之后乘除64是为了让图片变为64的倍数，因为网络的输入是64\*64。

用了cv2.resize函数进行放缩，其中：

(W, H)是调整后的宽高

后面是使用的方法，如果k>1说明要放大，则可用差值法cv2.INTER\_LANCZOS4；如果k<1则说明要缩小，使用cv2.INTER\_AREA。

- apply\_canny = CannyDetector()，来自from annotator.canny import CannyDetector

```

1  class CannyDetector:
2      def __call__(self, img, low_threshold, high_threshold):
3          return cv2.Canny(img, low_threshold, high_threshold)

```

本质是调用了OpenCV的边缘检测算法，获得边缘提取图像

- 模型输入数据处理

control是将处理后的图像放在GPU上加速获得的副本，根据生成图像的数量转为多份拷贝，放在一个张量里。以第0维度堆叠，形成(num\_samples, H, W, C)张量。num\_samples可以理解为批大小batch，**h**是图像的高度，**w**是图像的宽度，**c**是图像的通道数。

## 模型创建（初始化中）

```
1 self.apply_canny = CannyDetector()
2     self.model = create_model('./models/cldm_v15.yaml').cpu()
3     self.model.cond_stage_model.cuda()
4     self.use_trt = True
5     # if not self.use_trt:
6     if 1:
7
8         self.model.load_state_dict(load_state_dict('./models/control_sd15_canny.pth', location='cuda'))
9         self.model = self.model.cuda()
10
11     self.ddim_sampler = DDIMSampler(self.model)
12     self.warm_up()
```

- create\_model, 来自from cldm\_trt.model import create\_model

```
1 def create_model(config_path):
2     config = OmegaConf.load(config_path)
3     model = instantiate_from_config(config.model).cpu()
4     print(f'Loaded model config from [{config_path}]')
5     return model
```

使用OmegaConf库加载配置文件（.yaml格式），之后根据文件进行模型实例化。

```
1 def instantiate_from_config(config):
2     if not "target" in config:
3         if config == '__is_first_stage__':
4             return None
5         elif config == "__is_unconditional__":
6             return None
7         raise KeyError("Expected key `target` to instantiate.")
8     return get_obj_from_str(config["target"])(**config.get("params", dict()))
```

"target" 键通常包含需要实例化的对象的类名或对象的路径。如果配置对象没有"target"键，说明无法知道要实例化的对象是什么，因此会进入后续的错误处理逻辑。有两中情况会跳过模型实例化：

当 config == '\_\_is\_first\_stage\_\_', 用于表示模型第一阶段

当 config == "\_\_is\_unconditional\_\_", 用于表示某种无条件状态

当有"target"键会进行实例化：

- config["target"]：获取 "target" 键值，这里是cldm.cldm.ControlLDM
- get\_obj\_from\_str(config["target"])：返回cldm.cldm.ControlLDM对应的类或对象

- `**config.get("params", dict())`: 以字典形式获取配置中参数, 如果没有"params", 则返回一个空字典。`**`表示将结果作为关键字参数传递给目标类的构造函数。

## 启用内存优化模式

```
1 if config.save_memory:
2     self.model.low_vram_shift(is_diffusing=False)
```

找到`low_vram_shift()`函数:

```
1 def low_vram_shift(self, is_diffusing):
2     if is_diffusing:
3         self.model = self.model.cuda()
4         self.control_model = self.control_model.cuda()
5         self.first_stage_model = self.first_stage_model.cpu()
6         self.cond_stage_model = self.cond_stage_model.cpu()
7     else:
8         self.model = self.model.cpu()
9         self.control_model = self.control_model.cpu()
10        self.first_stage_model = self.first_stage_model.cuda()
11        self.cond_stage_model = self.cond_stage_model.cuda()
```

GPU资源有限, 并且在不同阶段, O型不同部分可能有不同的计算需求。

因此在“diffusing”阶段

- `model`和`control_model`是计算密集型部分。扩散过程包括生成、噪声添加、逐步更新等操作。这些操作通常需要大量的计算资源, 因此将它们放在 GPU 上能够利用 GPU 的并行计算能力, 从而加速计算过程。
- `first_stage_model` 和 `cond_stage_model` 在扩散阶段, 这些模型可能只需要偶尔使用, 因此将它们移到 CPU 上可以释放 GPU 显存, 避免不必要的显存占用。

在非“diffusing”阶段

- `first_stage_model` 和 `cond_stage_model` 这两个模型在其他阶段可能需要更多的计算。例如, 在模型的训练、推理、条件生成等过程中, `first_stage_model` 和 `cond_stage_model` 可能参与了更多的计算任务。这时候它们需要频繁与其他模型交互, 因此被放到 GPU 上可以加速这些计算过程。
- `model` 和 `control_model` 在非扩散阶段可能不是计算的核心部分, 因此可以将它们移回 CPU 上, 节省 GPU 显存, 以便其他计算密集型任务使用。

## 设置随机数种子

```
1 if seed == -1:
2     seed = random.randint(0, 65535)
3     seed_everything(seed)
```

## 拼接条件语句

```

1 cond = {"c_concat": [control], "c_crossattn":
  [self.model.get_learned_conditioning([prompt + ', ' + a_prompt] *
  num_samples))]}
2 un_cond = {"c_concat": None if guess_mode else [control], "c_crossattn":
  [self.model.get_learned_conditioning([n_prompt] * num_samples))]}
3 shape = (4, H // 8, W // 8)

```

## 控制比重

```

1 self.model.control_scales = [strength * (0.825 ** float(12 - i)) for i in
  range(13)] if guess_mode else ([strength] * 13)

```

得到输入图片对生成图片的影响度

如果guss\_mode打开，则按照[strength \* (0.825 \*\* float(12 - i)) for i in range(13)]公式生成

如果guss\_mode关闭，表示完全由参考输入图片，则control\_scales全为1.

## 参数进入模型

调试: [python pdb 代码调试 - 最全最详细的使用说明 - 简书](#)

## ControlNet

```

1 samples, intermediates = self.ddim_sampler.sample_simple(ddim_steps,
  num_samples,
2                                     shape, cond,
  verbose=False, eta=eta,
3
  unconditional_guidance_scale=scale,
4
  unconditional_conditioning=un_cond)
5 x_samples = self.model.decode_first_stage(samples)
6 x_samples = (einops.rearrange(x_samples, 'b c h w -> b h w c') * 127.5 +
  127.5).cpu().numpy().clip(0, 255).astype(np.uint8)

```

可见调用了self.ddim\_sampler的sample\_simple()函数

self.ddim\_sampler来自self.ddim\_sampler = DDIMSampler(self.model)，找到ddim\_hacker.py，进行调试。

进入循环：

```

1 for i, step in enumerate(iterator):
2     index = total_steps - i - 1
3     ts = torch.full((batch_size,), step, device=device, dtype=torch.long)

```

之后，出现了

```

(Pdb) n
DDIM Sampler: 0% | 0/20 [00:00<?, ?it/s]

```

说明开始进入模型输入阶段了。

看到代码

```

1 if self.controlnet_trt and self.controlunet_trt:

```

```

2     hint = torch.cat(conditioning['c_concat'], 1)
3     cond_txt = torch.cat(conditioning['c_crossattn'], 1)
4     #if self.cuda_graph_instance is None:
5     #cudart.cudaStreamBeginCapture(self.stream1.ptr,
6     cudart.cudaStreamCaptureMode.cudaStreamCaptureModeGlobal)
7
8     torch.cuda.synchronize()
9     start_time = time.time()
10    control_trt_dict = self.controlnet_engine.infer({"x_noisy":img,
11    "hint":hint, "timestep":ts, "context":cond_txt}, stream = self.stream1,
12    use_cuda_graph=True)
13    torch.cuda.synchronize()
14    end_time = time.time()
15    # print(f"controlnet_engine time={(end_time-start_time)*1000}ms")
16    control = list(control_trt_dict.values())
17    input_dict = {'x_noisy': img, 'timestep': ts, 'context': cond_txt,
18    'control0': control[4], 'control1': control[5],
19    'control2': control[6], 'control3': control[7],
20    'control4': control[8], 'control5': control[9],
21    'control6': control[10], 'control7': control[11],
22    'control8': control[12], 'control9': control[13],
23    'control10': control[14], 'control11': control[15],
24    'control12': control[16]}
25    torch.cuda.synchronize()
26    start_time = time.time()
27    model_t = self.unet_engine.infer(input_dict, self.stream1,
28    use_cuda_graph=True)['latent'].clone()
29    torch.cuda.synchronize()
30    end_time = time.time()
31    # print(f"unet_engine time={(end_time-start_time)*1000}ms")
32    else:
33        model_t = self.model.apply_model(img, ts, conditioning)

```

看到了'x\_noisy'、'timestep'等输入，与老师提示的一致，本以为找到了controlnet的接口，结果pdb进入了else：

这里不知道if self.controlnet\_trt and self.controlnet\_trt判断依据是什么，我的猜测是是否会使用TRT的两种模型

答：从代码

```

1     controlnet_engine_path = "./engine/ControlNet.plan"
2     if not os.path.exists(controlnet_engine_path):
3         self.controlnet_trt = False

```

可以看出，如果ControlNet的TRT文件不存在时，就设为False。

也就是说，在无模型情况下运行会去创建模型，如果在有TRT模型情况下就会运行指定的模型

```

1     else:
2         model_t = self.model.apply_model(img, ts, conditioning)

```

这里链接到了cldm.py的apply\_model，也与老师提示的一致：

```

1     def apply_model(self, x_noisy, t, cond, *args, **kwargs):
2         assert isinstance(cond, dict)

```

```

3     diffusion_model = self.model.diffusion_model
4
5     cond_txt = torch.cat(cond['c_crossattn'], 1)
6
7     if cond['c_concat'] is None:
8         eps = diffusion_model(x=x_noisy, timesteps=t, context=cond_txt,
9                               control=None, only_mid_control=self.only_mid_control)
10    else:
11        control = self.control_model(x=x_noisy,
12                                     hint=torch.cat(cond['c_concat'], 1), timesteps=t, context=cond_txt)
13        control = [c * scale for c, scale in zip(control,
14                                                  self.control_scales)]
15        eps = diffusion_model(x=x_noisy, timesteps=t, context=cond_txt,
16                              control=control, only_mid_control=self.only_mid_control)
17
18    return eps

```

这里最主要的是if判断，从之前的输入：

```

1 # canny2image_TRT.py
2 cond = {"c_concat": [control], "c_crossattn":
3         [self.model.get_learned_conditioning([prompt + ', ' + a_prompt] *
4         num_samples))]}

```

并且用pdb验证了一下：

```

(Pdb) p cond_txt
tensor([[[[-0.3884,  0.0229, -0.0522, ..., -0.4899, -0.3066,  0.0675],
          [ 0.0290, -1.3258,  0.3085, ..., -0.5257,  0.9768,  0.6652],
          [ 0.1073, -0.0619, -0.0716, ..., -1.8726, -0.6527,  0.8814],
          ...,
          [-1.4692, -0.2195, -0.1264, ..., -2.3614,  0.2694, -0.0432],
          [-1.4824, -0.2189, -0.1254, ..., -2.3543,  0.2633, -0.0390],
          [-1.4006, -0.1564, -0.0655, ..., -2.3871,  0.2966, -0.1056]]],
        device='cuda:0'])

```

可以看到是有内容的，pdb后也是进入了else。

这里有两个模型输入，我需要分别知道他们是什么模型：

```

1 control = self.control_model(x=x_noisy, hint=torch.cat(cond['c_concat'], 1),
2                               timesteps=t, context=cond_txt)

```

从这里的输入我看到了x\_noisy, hint, timesteps, context。与老师提示的对上了，很可能是controlnet，我需要进一步证明

进一步调试：

```

> /root/miniconda/lib/python3.8/site-packages/torch/nn/modules/module.py(1613) __getattr__()
-> return modules[name]
(Pdb) 1
1608         if name in _buffers:
1609             return _buffers[name]
1610         if '_modules' in self.__dict__:
1611             modules = self.__dict__['_modules']
1612             if name in modules:
1613 ->                 return modules[name]
1614             raise AttributeError("{} object has no attribute '{}'".format(
1615                 type(self).__name__, name))
1616
1617     def __setattr__(self, name: str, value: Union[Tensor, 'Module']) -> None:
1618         def remove_from(*dicts_or_sets):

```



到这里可以看到是通脱torch创建了一个modules，这里打印返回值modules[name]，出现：

```
1 (Pdb) p modules[name]
2 ControlNet(
3   (time_embed): Sequential(
4     (0): Linear(in_features=320, out_features=1280, bias=True)
5     (1): SiLU()
6     (2): Linear(in_features=1280, out_features=1280, bias=True)
7   )
8   (input_blocks): ModuleList(
9     (0): TimestepEmbedSequential(
10      (0): Conv2d(4, 320, kernel_size=(3, 3), stride=(1, 1),
11      ...
```

可以看到确实是Controlnet，返回，打印输入：

```
(Pdb) p x_noisy.shape
torch.Size([1, 4, 32, 48])
(Pdb) p t.shape
torch.Size([1])
(Pdb) p cond_txt.shape
torch.Size([1, 77, 768])
(Pdb) p torch.cat(cond['c_concat'], 1).shape
torch.Size([1, 3, 256, 384])
```

确定了Controlnet输入的shape：

```
1 x_noisy = torch.randn(1, 4, 32, 48, dtype=torch.float32)
2 timestep = torch.tensor([1], dtype=torch.int32)
3 context = torch.randn(1, 77, 768, dtype=torch.float32)
4 hint = torch.randn(1, 3, 256, 384, dtype=torch.float32)
5
6 input_names = ["x_noisy", "hint", "timestep", "context"] #
```

输入方面input\_name顺序这里不知道有没有规则，我的判断应该是在torch.onnx.export()函数中，input\_names要与输入位置一致。

答：input\_names只是给输入参数命名，一一对应即可

输出方面，老师提示中是"latent"，我并未找到其出处，可能起名没有什么限制。

答：就是要保证这个名称在 unet 输出的定义和 decoder 输入名称一致就行

```
1 output_names = ["latent"]
```

仿照clip转换的代码，还需要：

```
1 onnx_path = "./onnx/CONTROL_NET.onnx" # 设置Onnx输出路径
2
3 # 模型导出
4 torch.onnx.export(
5     control_net,
6     (x_noisy, hint, timestep, context),
```

```

7         onnx_path,
8         verbose=True,
9         opset_version=18,
10        do_constant_folding=True,
11        input_names=input_names,
12        output_names=output_names,
13        keep_initializers_as_inputs=True
14    )
15
16    # 验证onnx模型
17    output = control_net(x_noisy, hint, timestep, context) # 得到模型输出结果
18    input_dicts = {"x_noisy": x_noisy.numpy(), "hint": hint.numpy(), "timestep":
19    timestep.numpy(), "timestep": timestep.numpy()} # onnxruntime推理输入字典
20    onnxruntime_check(onnx_path, input_dicts, [output]) # 这个函数我看了下是为了检测
21    onnx模型导出是否正确

```

其中关于模型导出有几个问题：

- `dynamic_axes=dynamic_axes`是输入输出动态维度，我看模型controlnet输入输出应该都是固定的，我这个不确定

clip模型基于文本对图像分类，也可以基于图像对文本分类，是SD的文本条件约束。这个输入肯定是动态大小输入，主要是文本条数。

看其导出参数设置：

```

1 dynamic_axes = {
2     "input_ids": {1: "s"}, # 输入 "input_ids" 的第 1 维度是动态的，命名为 "s"
3     "last_hidden_state": {1: "s"} # 输出 "last_hidden_state" 的第 1 维度是动态
4     的，命名为 "s"
5 }

```

- `"s"` 是该动态维度的名称（可以自定义，通常使用有意义的名称，如 `"sequence_length"`）。

因此在`torch.onnx.export`内添加参数即可

- `keep_initializers_as_inputs=True`这个参数是看老师视频的答案里面有，但是不知道为什么这个要设置为True

`torch.onnx.export()`函数的`keep_initializers_as_inputs`参数，老师说设置成默认值 `None` 大多数时候都是没问题的，拿 `none false true` 测试 三者导出的 `onnx` 都能正常用。

以`forward(a, b=torch.Tensor([1]))`为例

- `keep_initializers_as_inputs=False`：带有默认值的参数（如 `b=torch.Tensor([1])`）会被视作一个常量（initializer），这个常量会在导出到 ONNX 时被嵌入到模型中。因此，当你调用导出的 ONNX 模型时，只需要传入 `a`，`b` 会自动使用其默认值。并且不能在调用时改变 `b`，除非你修改模型的定义，让 `b` 成为一个显式的输入。如果你想要动态地给 `b` 传入不同的值，你需要将 `keep_initializers_as_inputs` 设置为 `True`，这样 `b` 就会作为一个输入（input）被处理，而不是作为一个默认值。也就是说只能以`forward(a)`形式调用，`forward(a, b)`不行
- `keep_initializers_as_inputs=True`：默认值被视为模型输入，必须显式传递给模型。
- `forwad`什么时候该替换？

老师在课程里面讲到，onnx模型ONNX在实践中主要支持张量输入，要将非张量类型传给onnx最好需要将对象转成张量，如文本要通过嵌入层转换为张量，而图片数据可直接作为张量输入。

因此control\_net的pytorch转onnx代码最终为:

```
1 def export_control_net_model():
2     control_net = hk.model.control_model
3
4     x_noisy = torch.randn(1, 4, 32, 48, dtype=torch.float32)
5     timestep = torch.tensor([1], dtype=torch.int32)
6     context = torch.randn(1, 77, 768, dtype=torch.float32)
7     hint = torch.randn(1, 3, 256, 384, dtype=torch.float32)
8
9     onnx_path = "./onnx/CONTROL_NET.onnx"
10
11     input_names = ["x_noisy", "hint", "timestep", "context"] # 这里不知道输入顺
12     output_names = ["latent"]
13
14     # 模型输出
15     torch.onnx.export(
16         control_net,
17         (x_noisy, hint, timestep, context),
18         onnx_path,
19         verbose=True,
20         opset_version=18,
21         do_constant_folding=True,
22         input_names=input_names,
23         output_names=output_names,
24         keep_initializers_as_inputs=True
25     )
26     print("===== CONTROL_NET model export onnx done!")
27
28     # 验证onnx模型
29     output = control_net(x_noisy, hint, timestep, context)
30     input_dicts = {"x_noisy": x_noisy.numpy(), "hint": hint.numpy(),
31 "timestep": timestep.numpy(), "context": context.numpy()}
32     onnxruntime_check(onnx_path, input_dicts, output)
33     print("===== CONTROL_NET onnx model verify done!")
```

运行export\_onnx.py:

```
===== CONTROL_NET onnx model verify done!
root@6bbe6b07716d:~/TensorRT-StableDiffusion# cd onnx
root@6bbe6b07716d:~/TensorRT-StableDiffusion/onnx# ls
CLIP.onnx  CONTROL_NET.onnx  readme.md
root@6bbe6b07716d:~/TensorRT-StableDiffusion/onnx#
```

模型导出成功~

错误补充1:

```
Traceback (most recent call last):
  File "export_onnx_full.py", line 215, in <module>
    main()
  File "export_onnx_full.py", line 210, in main
    export_control_net_model()
  File "export_onnx_full.py", line 147, in export_control_net_model
    onnxruntime_check(onnx_path, input_dicts, [output])
  File "export_onnx_full.py", line 49, in onnxruntime_check
    ret = np.allclose(result[i], torch_outputs[i].detach().numpy(), rtol=1e-03, atol=1e-05, equal_nan=False)
AttributeError: 'list' object has no attribute 'detach'
```

对应代码是

```
1 | onnxruntime_check(onnx_path, input_dicts, [output])
```

因为controlnet结果已经是list了，所以不用加output不用加[]

错误补充2：

```
2025-02-12 19:39:32.362940235 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:MatMul_4410 appears in graph inputs and will not be treated as
constant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to ove
rride it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.p
y.
2025-02-12 19:39:32.362945135 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:MatMul_4431 appears in graph inputs and will not be treated as
constant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to ove
rride it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.p
y.
2025-02-12 19:39:32.362949844 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:MatMul_4432 appears in graph inputs and will not be treated as
constant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to ove
rride it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.p
y.
2025-02-12 19:39:32.362957250 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:MatMul_4433 appears in graph inputs and will not be treated as
constant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to ove
rride it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.p
y.
2025-02-12 19:39:32.362963124 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:Mul_4434 appears in graph inputs and will not be treated as co
nstant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to overri
de it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.py.
2025-02-12 19:39:32.362968726 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:Add_4435 appears in graph inputs and will not be treated as co
nstant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to overri
de it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.py.
2025-02-12 19:39:32.362973714 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:Mul_4436 appears in graph inputs and will not be treated as co
nstant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to overri
de it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.py.
2025-02-12 19:39:32.362978655 [W:onnxruntime:, graph.cc:1312 Graph] Initializer onnx:Add_4437 appears in graph inputs and will not be treated as co
nstant value/weight. This may prevent some of the graph optimizations, like const folding. Move it out of graph inputs if there is no need to overri
de it, by either re-generating the model with latest exporter/converter or with the tool onnxruntime/tools/python/remove_initializer_from_input.py.
Error onnxruntime_check
Error onnxruntime_check
Error onnxruntime_check
Error onnxruntime_check
===== CONTROL.NET onnx model verify done!
root@3f2b543c614d:~/TensorRT-StableDiffusion#
```

黄色部分是warning，没有关系

onnxruntime\_check函数报错，看看是怎么检查模型的

```
1 | def onnxruntime_check(onnx_path, input_dicts, torch_outputs):
2 |     onnx_model = onnx.load(onnx_path)
3 |     # onnx.checker.check_model(onnx_model)
4 |     sess = rt.InferenceSession(onnx_path)
5 |     # outputs = self.get_output_names()
6 |     # latent input
7 |     # data = np.zeros((4, 77), dtype=np.int32)
8 |     result = sess.run(None, input_dicts)
9 |
10 |    for i in range(0, len(torch_outputs)):
11 |        ret = np.allclose(result[i], torch_outputs[i].detach().numpy(),
12 |                           rtol=1e-03, atol=1e-05, equal_nan=False)
13 |        if ret is False:
14 |            print("Error onnxruntime_check")
15 |            # import pdb; pdb.set_trace()
```

可见是将onnx文件导入，用onnxruntime运行得到结果与输入结果进行对比。

np.allclose是Numpy中的一个函数，用于判断两个数组是否在数值上近似相等。

定义：

```
1 | numpy.allclose(a, b, rtol=1e-5, atol=1e-8, equal_nan=False)
```

- a, b: 输入的两个数组，需要比较它们是否相等。
- rtol: 相对误差的容忍度（默认值是  $1 \times 10^{-5}$ ）。计算方式是基于两个数组对应元素的大小。
- atol: 绝对误差的容忍度（默认值是  $1 \times 10^{-8}$ ）。计算方式是一个全局的固定误差容忍度。
- equal\_nan: 是否将两个 NaN 视为相等。默认值是 False（即两个 NaN 不被视为相等）。设置为 True 后，两个对应位置的 NaN 将被视为相等

判断标准：

两个数组对应元素a和b满足以下条件时，视为近似相等：

$$|a - b| \leq atol + rtol \cdot |b|$$

转onnx运行后，输出与原值有误差是正常现象，重要的是能接受的误差范围是多少。

## ControlUnet

继续看下一步：

```
1 eps = diffusion_model(x=x_noisy, timesteps=t, context=cond_txt,
    control=control, only_mid_control=self.only_mid_control)
```

打印查看diffusion\_model，可以知道这就是controlunet：

```
1 ControlledUnetModel(
2   (time_embed): Sequential(
3     (0): Linear(in_features=320, out_features=1280, bias=True)
4     (1): SiLU()
5     (2): Linear(in_features=1280, out_features=1280, bias=True)
6   )
7   (input_blocks): ModuleList(
8     (0): TimestepEmbedSequential(
9       (0): Conv2d(4, 320, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
10    )
11    ...
```

从入口可以看出输入有4个：x\_noisy, timesteps, context, control。其中，x\_noisy, timesteps, context的shape和controlnet的一样，而control是controlnet的输出。打印一下：

```
(Pdb) for i, tensor in enumerate(control): print(f"Tensor {i}: {tensor.shape}")
Tensor 0: torch.Size([1, 320, 32, 48])
Tensor 1: torch.Size([1, 320, 32, 48])
Tensor 2: torch.Size([1, 320, 32, 48])
Tensor 3: torch.Size([1, 320, 16, 24])
Tensor 4: torch.Size([1, 640, 16, 24])
Tensor 5: torch.Size([1, 640, 16, 24])
Tensor 6: torch.Size([1, 640, 8, 12])
Tensor 7: torch.Size([1, 1280, 8, 12])
Tensor 8: torch.Size([1, 1280, 8, 12])
Tensor 9: torch.Size([1, 1280, 4, 6])
Tensor 10: torch.Size([1, 1280, 4, 6])
Tensor 11: torch.Size([1, 1280, 4, 6])
Tensor 12: torch.Size([1, 1280, 4, 6])
```

因此，controlunet输入可以设为：

```
1 x_noisy = torch.randn(1, 4, 32, 48, dtype=torch.float32)
2 timestep = torch.tensor([1], dtype=torch.int32)
3 context = torch.randn(1, 77, 768, dtype=torch.float32)
4
5 control = [
6     torch.randn(1, 320, 32, 48, dtype=torch.float32),
7     torch.randn(1, 320, 32, 48, dtype=torch.float32),
8     torch.randn(1, 320, 32, 48, dtype=torch.float32),
9     torch.randn(1, 320, 16, 24, dtype=torch.float32),
10    torch.randn(1, 640, 16, 24, dtype=torch.float32),
11    torch.randn(1, 640, 16, 24, dtype=torch.float32),
12    torch.randn(1, 640, 8, 12, dtype=torch.float32),
```

```

13     torch.randn(1, 1280, 8, 12, dtype=torch.float32),
14     torch.randn(1, 1280, 8, 12, dtype=torch.float32),
15     torch.randn(1, 1280, 4, 6, dtype=torch.float32),
16     torch.randn(1, 1280, 4, 6, dtype=torch.float32),
17     torch.randn(1, 1280, 4, 6, dtype=torch.float32),
18     torch.randn(1, 1280, 4, 6, dtype=torch.float32),
19 ]
20
21 input_names = ["x_noisy", "timestep", "context"]
22 for i in range(0, len(control)):
23     input_names.append("control" + str(i))

```

输出

```

1  output_names = ["latent"]
2
3  onnx_path = "./onnx/CONTROLUNET.onnx"
4
5  torch.onnx.export(
6      controlled_unet_model,
7      (x_noisy, timestep, context, control),
8      onnx_path,
9      verbase=True,
10     opset_version=18,
11     do_constant_folding=True,
12     input_names=input_names,
13     output_names=output_names
14 )

```

验证

```

1  output = controlled_unet_model(x_noisy, hint, timestep, control)
2  input_dicts = {"x_noisy": x_noisy.numpy(), "timestep": timestep.numpy(),
3  "context": context.numpy(), "control": control.numpy()}
4  onnxruntime.check(onnx_path, input_dicts, output)

```

总结:

```

1  def export_controlled_unet_model():
2      controlled_unet_model = hk.model.model.diffusion_model
3
4      x_noisy = torch.randn(1, 4, 32, 48, dtype=torch.float32)
5      timestep = torch.tensor([1], dtype=torch.int32)
6      context = torch.randn(1, 77, 768, dtype=torch.float32)
7
8      control = [
9          torch.randn(1, 320, 32, 48, dtype=torch.float32),
10         torch.randn(1, 320, 32, 48, dtype=torch.float32),
11         torch.randn(1, 320, 32, 48, dtype=torch.float32),
12         torch.randn(1, 320, 16, 24, dtype=torch.float32),
13         torch.randn(1, 640, 16, 24, dtype=torch.float32),
14         torch.randn(1, 640, 16, 24, dtype=torch.float32),
15         torch.randn(1, 640, 8, 12, dtype=torch.float32),
16         torch.randn(1, 1280, 8, 12, dtype=torch.float32),
17         torch.randn(1, 1280, 8, 12, dtype=torch.float32),
18         torch.randn(1, 1280, 4, 6, dtype=torch.float32),

```

```

19         torch.randn(1, 1280, 4, 6, dtype=torch.float32),
20         torch.randn(1, 1280, 4, 6, dtype=torch.float32),
21         torch.randn(1, 1280, 4, 6, dtype=torch.float32),
22     ]
23     # import pdb; pdb.set_trace()
24     onnx_path = "./onnx/CONTROL_UNET.onnx"
25
26     input_names = ["x_noisy", "timestep", "context"]
27     for i in range(0, len(control)):
28         input_names.append("control" + str(i))
29     output_names = ["latent"]
30
31     # 模型输出
32     torch.onnx.export(
33         controlled_unet_model,
34         (x_noisy, timestep, context, control),
35         onnx_path,
36         verbose=True,
37         opset_version=18,
38         do_constant_folding=True,
39         input_names=input_names,
40         output_names=output_names,
41     )
42     print("===== CONTROL_UNET model export onnx done!")
43
44     # 验证onnx模型
45     # 这里如果先计算output会导致control清空因此要先写input_dicts,不然for循环执行不了
    (len(control)=0)
46     input_dicts = {"x_noisy": x_noisy.numpy(), "timestep": timestep.numpy(),
47 "context": context.numpy()}
48     for i in range(0, len(control)):
49         input_dicts["control" + str(i)] = control[i].numpy()
50
51     output = controlled_unet_model(x_noisy, timestep, context, control)
52
53     onnxruntime_check(onnx_path, input_dicts, output)
54     print("===== CONTROL_UNET onnx model verify done!")

```

报错问题:

```

===== CONTROL_UNET model export onnx done!
Traceback (most recent call last):
  File "export_onnx_full.py", line 219, in <module>
    main()
  File "export_onnx_full.py", line 215, in main
    export_controlled_unet_model()
  File "export_onnx_full.py", line 203, in export_controlled_unet_model
    onnxruntime_check(onnx_path, input_dicts, output)
  File "export_onnx_full.py", line 46, in onnxruntime_check
    result = sess.run(None, input_dicts)
  File "/root/miniconda/lib/python3.8/site-packages/onnxruntime/capi/onnxruntime_inference_collection.py", line 216, in run
    self._validate_input(list(input_feed.keys()))
  File "/root/miniconda/lib/python3.8/site-packages/onnxruntime/capi/onnxruntime_inference_collection.py", line 198, in _validate_input
    raise ValueError(
ValueError: Required inputs ('control0', 'control1', 'control2', 'control3', 'control4', 'control5', 'control6', 'control7', 'control8', 'control9', 'control10', 'control11', 'control12') are missing from input feed ({'x_noisy', 'timestep', 'context'}).

```

显示是在验证模型时，输入没有control选项。

但是我明明写了

```

1 input_dicts = {"x_noisy": x_noisy.numpy(), "timestep": timestep.numpy(),
2 "context": context.numpy()}
3 for i in range(0, len(control)):
4     input_dicts["control" + str(i)] = control[i].numpy()

```



input\_dicts里应该有control才对。于是我pdb一下，发现input\_dicts里没有control!

继续pdb，发现 `for i in range(0, len(control))`: 没有执行，查看`len(control)=1`，再查看control为空。但是在创建模型是control有内容。于是想到可能是在计算output结果时将control传入导致里面数据清空，pdb证实我的猜想是正确的，因此input\_dicts要在output计算前

即：

```
1 input_dicts = {"x_noisy": x_noisy.numpy(), "timestep": timestep.numpy(),
2               "context": context.numpy()}
3     for i in range(0, len(control)):
4         input_dicts["control" + str(i)] = control[i].numpy()
5     output = controlled_unet_model(x_noisy, timestep, context, control)
```

## decoder

之后回到canny2imageTRT.py

```
1 x_samples = self.model.decode_first_stage(samples)
2 x_samples = (einops.rearrange(x_samples, 'b c h w -> b h w c') * 127.5 +
3               127.5).cpu().numpy().clip(0, 255).astype(np.uint8)
4 results = [x_samples[i] for i in range(num_samples)]
```

这里的model是：

```
1 self.model = create_model('./models/cldm_v15.yaml').cpu()
```

pdb查看：

```
1 (Pdb) p decode_model
2 AutoencoderKL(
3   (encoder): Encoder(
4     (conv_in): Conv2d(3, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
5     (down): ModuleList(
6       (0): Module(
7         (block): ModuleList(
8           (0-1): 2 x ResnetBlock(
9             (norm1): GroupNorm(32, 128, eps=1e-06, affine=True)
10            (conv1): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1),
11              padding=(1, 1))
12            ...
```

AutoencoderKL是自动编码器，是一种数据压缩算法。

显然，根据 `x_samples = self.model.decode_first_stage(samples)` 知道，decoder输入只有一个samples。通过pdb，samples大小为[1, 4, 32, 48]

因此：



```

1  def export_decoder_model():
2      # control_net = hk.model.control_model
3
4      decode_model = hk.model.first_stage_model
5      decode_model.forward = decode_model.decode
6
7      latent = torch.randn(1, 4, 32, 48, dtype=torch.float32)
8
9      input_names = ["latent"]
10     output_names = ["image"]
11
12     onnx_path = "./onnx/DECODER.onnx"
13
14     torch.onnx.export(
15         decode_model,
16         (latent),
17         onnx_path,
18         verbose=True,
19         opset_version=18,
20         do_constant_folding=True,
21         input_names=input_names,
22         output_names=output_names,
23         keep_initializers_as_inputs=True
24     )
25     print("===== DECODER model export onnx done!")
26
27     # 验证onnx模型
28     output = decode_model(latent)
29     input_dicts = {"latent": latent.numpy()}
30     onnxruntime_check(onnx_path, input_dicts, [output])
31     print("===== DECODER onnx model verify done!")

```

## 模型测试

Controlnet和decoder的onnxruntime\_check能在相对误差的容忍度rtol为1e-3、绝对误差容忍度atol为1e-5的情况下导出。

```

===== Diagnostic Run torch.onnx.export version 2.0.0+cu118 =====
verbose: False, log level: Level.ERROR
===== 0 NONE 0 NOTE 0 WARNING 0 ERROR =====

===== CONTROL UNET model export onnx done!
===== CONTROL UNET onnx model verify done!
root@126c3b8cb4e6:~/TensorRT-StableDiffusion#

```

```

stant value/weight. This may prevent some of the graph optimization
e it, by either re-generating the model with latest exporter/conve
2025-02-13 22:38:12.323682334 [W:onnxruntime:, graph.cc:1312 Graph
stant value/weight. This may prevent some of the graph optimization
e it, by either re-generating the model with latest exporter/conve
===== DECODER onnx model verify done!
root@126c3b8cb4e6:~/TensorRT-StableDiffusion#

```

而Controlnet导出误差稍微大点，因此需要提高容忍度。

```

de it, by either re-generating the model with latest exporter/conv
2025-02-13 22:26:12.050543239 [W:onnxruntime:, graph.cc:1312 Graph
nstant value/weight. This may prevent some of the graph optimizati
de it, by either re-generating the model with latest exporter/conv
Error onnxruntime_check
Error onnxruntime_check
Error onnxruntime_check
===== CONTROL_NET onnx model verify done!

```

一般首先增加 `rtol`，可以让模型或者计算更加宽松地允许比例上的误差，这对于大范围的数值数据或具有较大数量级差异的数据很有用。然后，可以再考虑调整 `atol` 来容忍一些小的绝对差异，这样可以保证即使数值很小的差异也不被忽视。

测试：

- `rtol=1e-3,atol=1e-4`：通过

```

de it, by either re-generating the model with latest exporter/conv
2025-02-13 22:29:35.874318638 [W:onnxruntime:, graph.cc:1312 Graph
nstant value/weight. This may prevent some of the graph optimizati
de it, by either re-generating the model with latest exporter/conv
===== CONTROL_NET onnx model verify done!

```