

CHAPTER3

TENSORRT基础入门

主 讲：韩君
单 位：早稻田大学
公众号：自动驾驶之心

主要内容

- 1 TENSORRT简介
- 2 TENSORRT的应用场景
- 3 TENSORRT的模块
- 4 导出ONNX以及修改ONNX的方法
- 5 初步使用TENSORRT



01

TensorRT简介

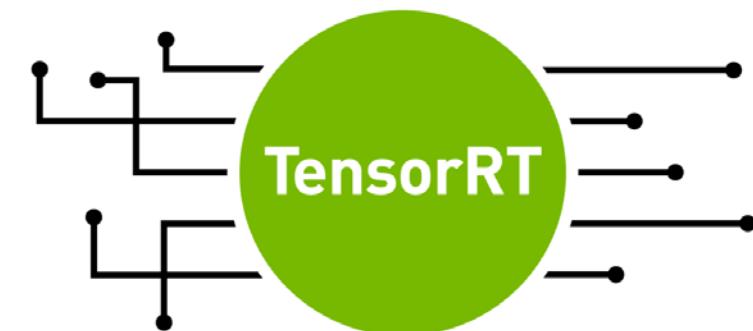
Goal: 理解TensorRT的工作流程，TensorRT的极限，以及其他DNN编译器

什么是TensorRT

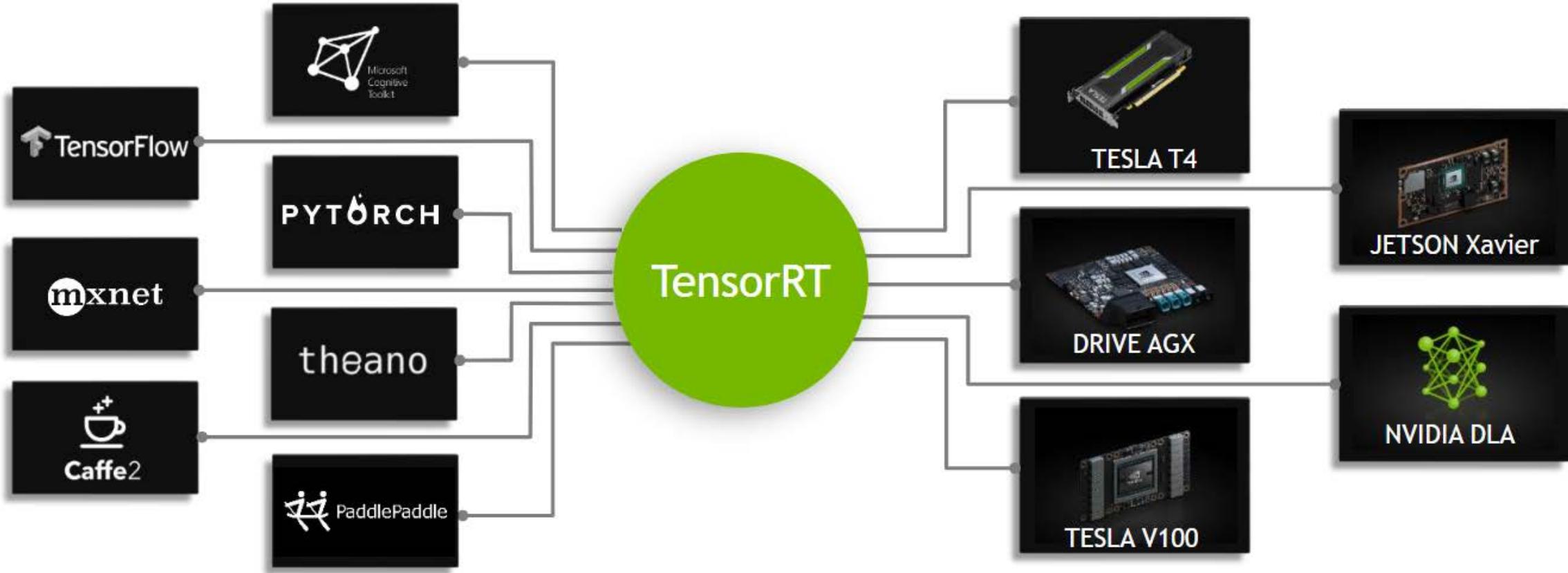
*“NVIDIA® TensorRT™, an **SDK** for **high-performance** deep learning **inference**, includes a deep learning inference optimizer and runtime that delivers **low latency and high throughput** for inference applications.”*

可以把它理解为一种针对NVIDIA GPU的一种优化编译器，以及其他开发工具

- 自动优化模型
 - 寻找模型中可以并行处理的地方
 - 针对当前部署的GPU框架，寻找最优的调度和并行策略
- 支持多框架输入
 - ONNX
- Python/C++ API接口
 - 可以方便在自己的程序中调用TensorRT API来实现推理



什么是TensorRT



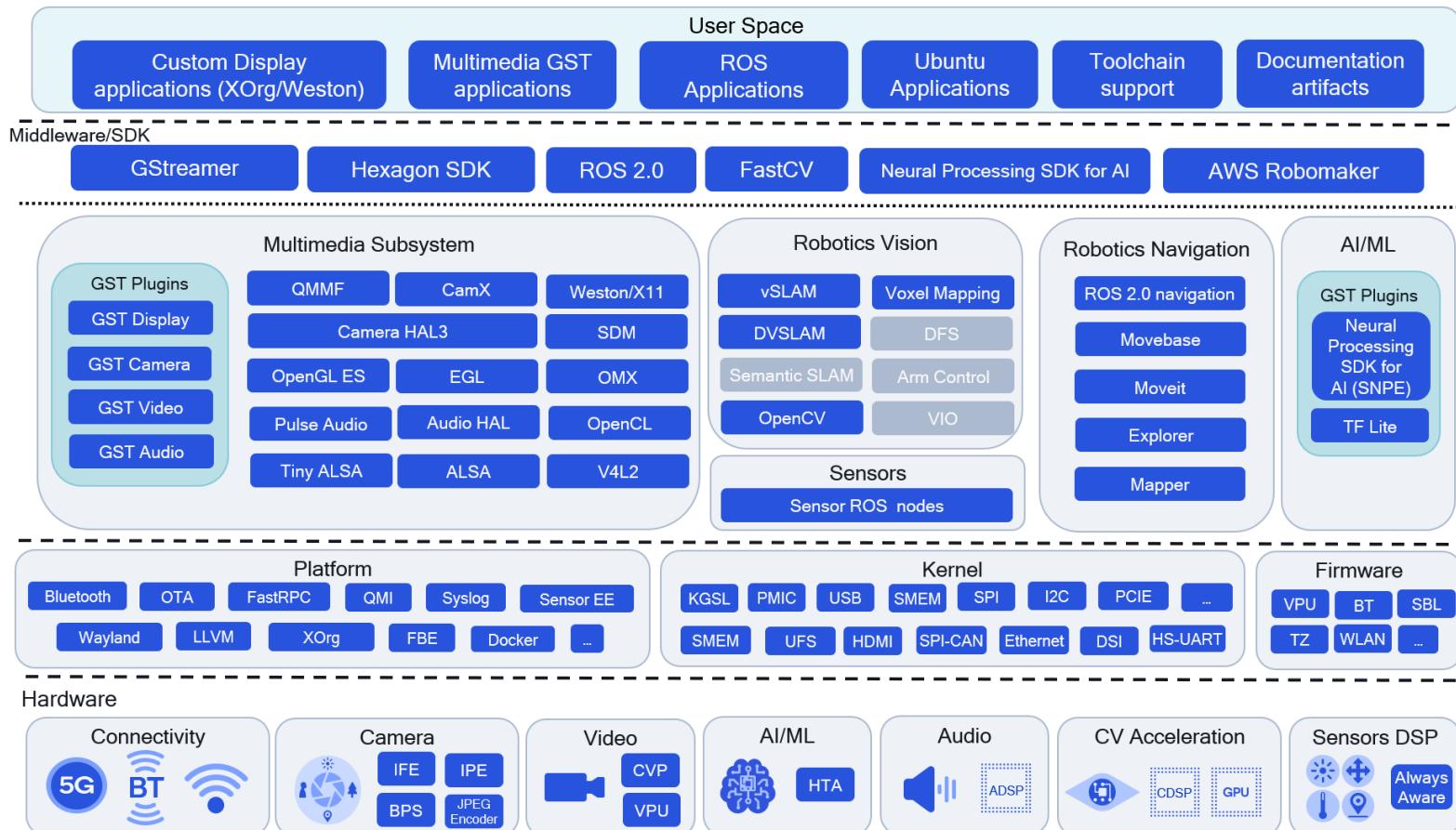
现在的大多数DL framework都可以通过TensorRT生成推理引擎，部署在硬件上

- 比较常见的组合：Pytorch -> ONNX -> TensorRT
- 目前自动驾驶的部署硬件大多都会采用NVIDIA
 - DRIVE series
 - Jetson series

什么是TensorRT

(扩展)

很多厂商除了NVIDIA，也会有很多其他选择



Qualcomm^[1]

(比较典型的Heterogeneous Architecture，这
里以Robotics举例，自动驾驶也差不多)

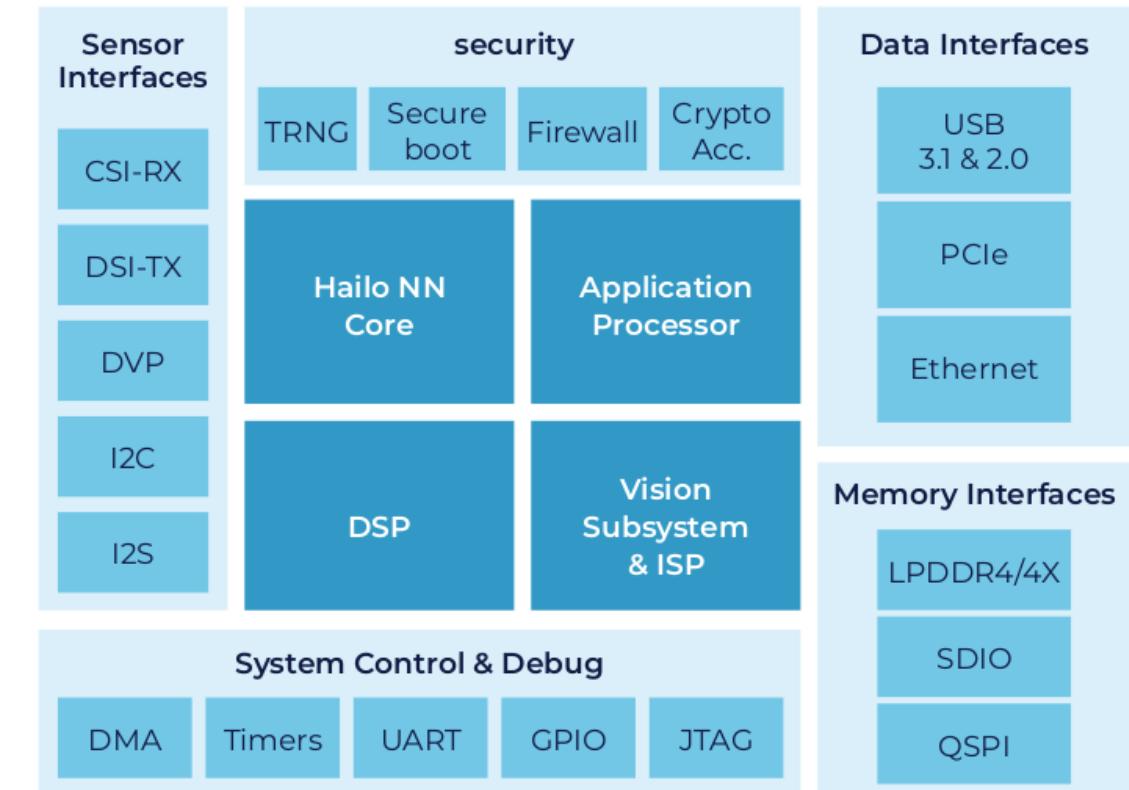
^[1][Qualcomm Software Architecture](#)

什么是TensorRT

(*)Dataflow中文叫做“数据流”，是一个很传统的概念。最近针对DNN的硬件设计有偏向Dataflow的趋势。该兴趣的同学可以自学一下。建议参考Hailo, Google TPU, MIT Eyeriss的paper

(扩展)

很多厂商除了NVIDIA，也会有很多其他选择



Hailo^[1, 2]

(比较典型的小型Edge Device for DNN，有很多量化和调度的技巧，针对自己的硬件而设计的Dataflow^(*) Compiler很精妙)

^[1]Detail of Hailo AI Edge Accelerator Emerge

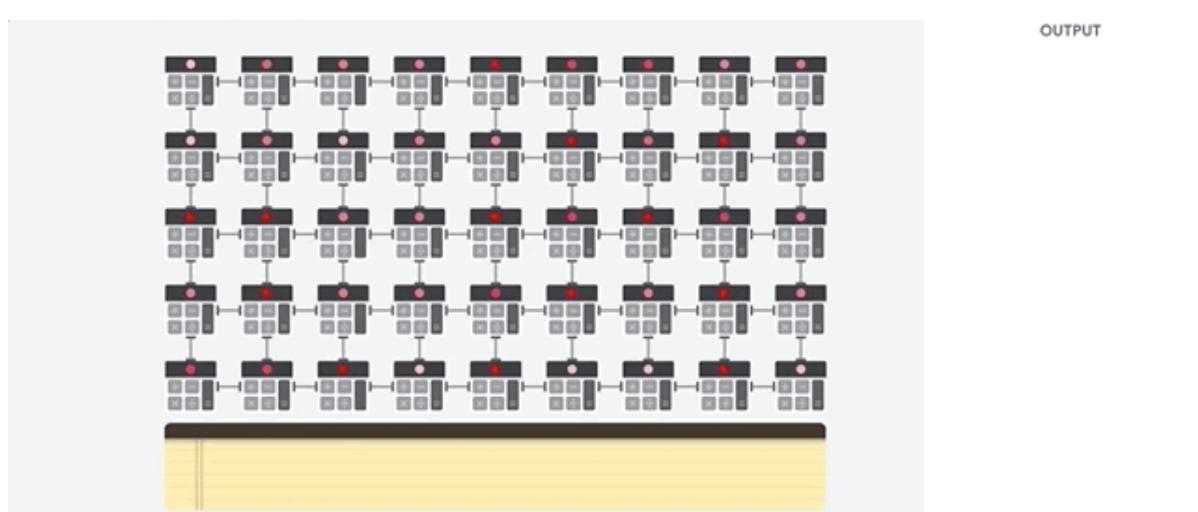
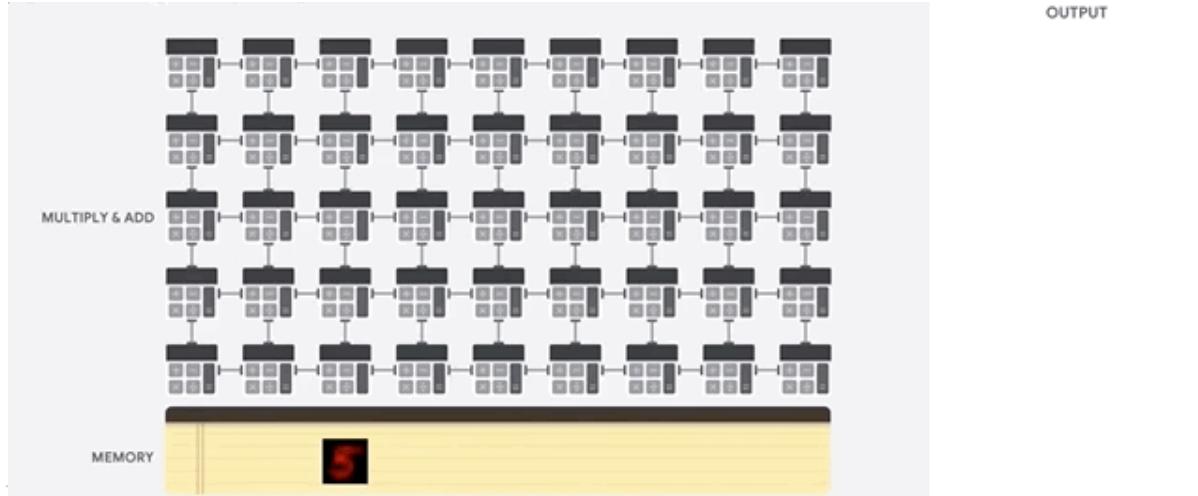
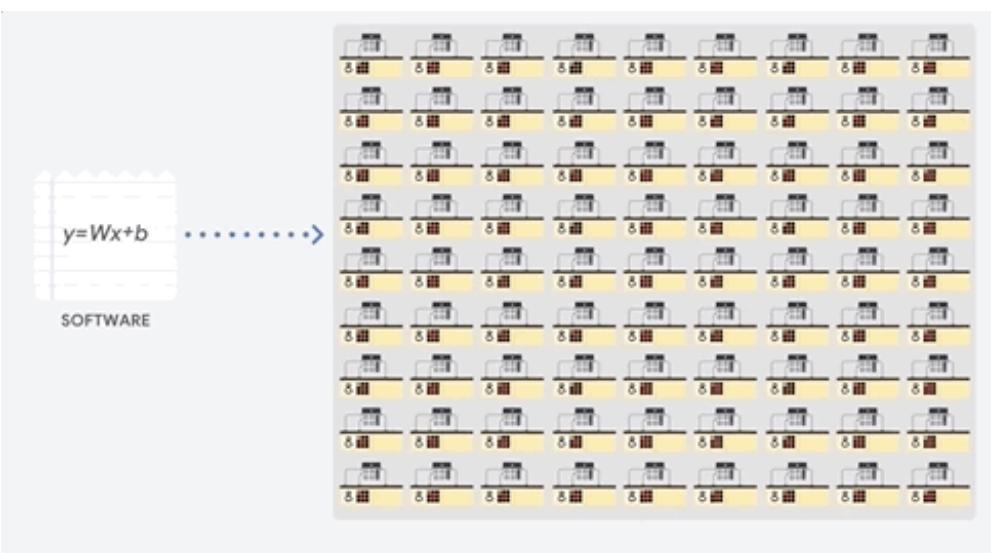
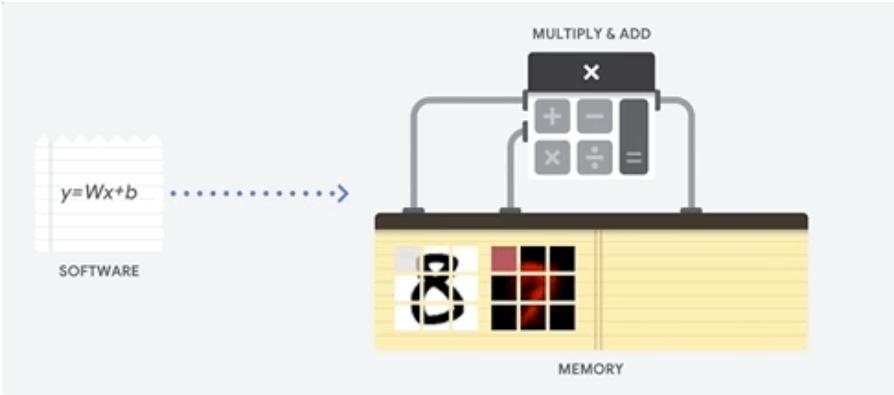
^[2]Hailo-15 quad-core AI Vision processor delivers up to 20 TOPS for Smart Cameras

什么是TensorRT

(扩展)

很多厂商除了NVIDIA，也会有很多其他选择

(*)Google的TPU也是属于dataflow architecture的一个例子。TPU内部有一个on-chip memory叫做HBM(high-bindwidth memory)，当读取数据的时候，HBM里的数据会以queue的形式放到MXU(Matrix Multiplication Unit)进行数据的流动以及计算



Google TPU^(*)
左上CPU, 左下GPU, 右上右下TPU

[1] [Introduction to Cloud TPU](#)

什么是TensorRT

(*)TOPS: 全称是Tera Operations per second。是衡量一个硬件的计算力的常用指标。这个会在之后仔细讲。现阶段有个认识就好

回归正题，有很多部署硬件但为什么NVIDIA和TensorRT一直都很火

- 整体硬件设计和编译技术很成熟
 - BUG算比较少的那种
 - 对量化的支持比较全面以及完善
 - SDK很充足
- 可以参考的文档比较多
 - 官方给的资源还是相当丰富的
 - 给的SAMPLE也很适合学习
- Community很大
 - 出现问题的时候有的问的

The screenshot shows the 'Revision History' section of the NVIDIA Deep Learning TensorRT Documentation. The page title is 'NVIDIA Deep Learning TensorRT Documentation' with a sub-section 'Developer Guide [PDF] - Last updated March 15, 2023'. A search bar is at the top right. The left sidebar has navigation links like 'Getting Started', 'Developer Guide', 'Revision History', and 'Appendix'. The main content area is titled 'Revision History' and says 'This is the revision history of the NVIDIA TensorRT 8.6 Developer Guide.' It lists five chapters with their update dates and summaries:

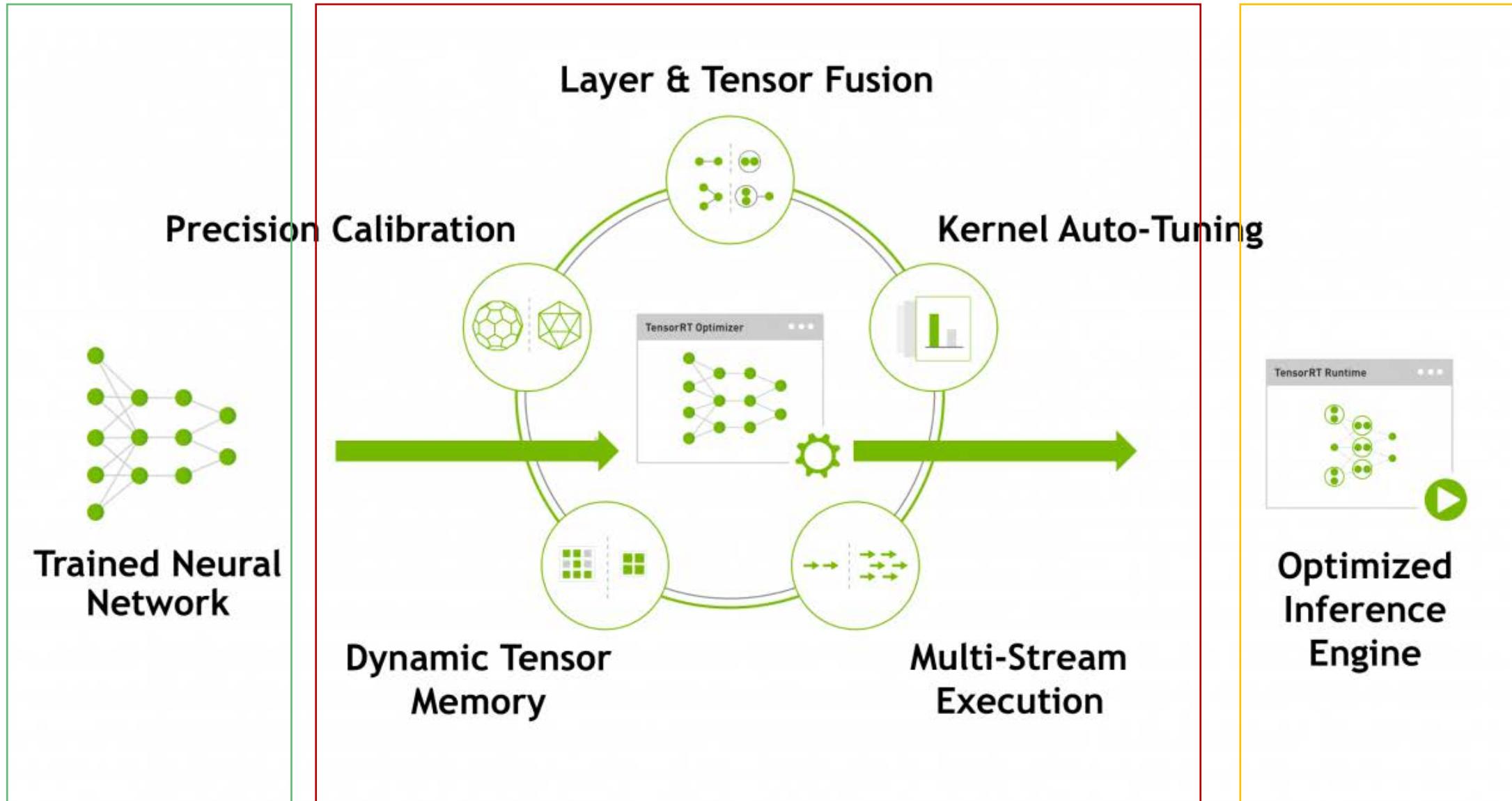
- Chapter 2 Updates**: January 17, 2023. Summary: Added a footnote to the [Types and Precision](#) topic.
- Chapter 3 Updates**: January 17, 2023. Summary: Updated the [Performing Inference](#) topic regarding executing a network.
- Chapter 4 Updates**: January 17, 2023. Summary: Updated the [Building an Engine](#) topic regarding insufficient workspace.
 - Updated the [Building an Engine](#) topic regarding insufficient workspace.
 - Updated the [Performing Inference](#) topic regarding executing a network.
- Chapter 5 Updates**: January 19, 2023. Summary: Added a new topic called [Runtime Options](#).
January 24, 2023. Summary: Rewrote the existing [CUDA Lazy Loading](#) topic.
March 8, 2023. Summary: Added a new topic called [Compatibility](#).

But ...?

- 价格偏贵，当面向量产的时候价格是很需要重视的
- 过大的TOPS^(*)是否真的就很好？

^[1][NVIDIA Deep Learning TensorRT Documentation](#)

TensorRT的工作流介绍



训练好的模型转为
TensorRT可识别的模型

TensorRT逐个Layer进行分析并尝试各种优化策略
(量化、层融合、调度等等)

生成推理引擎，可以从
C++/Python程序中调用

^[1][NVIDIA TensorRT](#)

TensorRT的一些限制

针对不支持的算子

- 看看不同TensorRT版本中是否有做了支持
 - PyTorch: Facebook
 - ONNX: Microsoft
 - TensorRT: NVIDIA
- 修改PyTorch的算子选择，让它使用一些TensorRT支持的算子
- 自己写插件，内部实现自定义算子以及自定义cuda加速核函数
- 不使用onnx，自己创建parser直接调用TensorRT API逐层创建网络
 - 比如darknet模型

不同TensorRT版本的优化策略是不一样的

- 比如说，对Transformer的优化和TensorRT 7和8跑出来的性能是不一样的

有时你预期TensorRT的优化和实际的优化是不一样的

- 比如说，你期望它使用Tensor core，但kernel autotuning之后TensorRT觉得使用Tensor core反而会效率降低，结果给你分配CUDA core使用。因为内部需要做一些额外的处理

天生并行性差的layer，TensorRT也没有办法

- 1x1 conv这种layer，再怎么优化也没有7x7 conv并行效果好

TensorRT的一些限制

重点！！

TensorRT虽然很方便，但不要过分的依赖TensorRT给你的结果。你需要结合部署的硬件特性，做一些Benchmark和Profiling，学习使用一些Profile tools去分析模型的计算瓶颈在哪里，以及与分析你的预测优化策略和实际优化策略产生分歧的原因在哪里。



02

TensorRT的应用场景

Goal: 理解TensorRT在自动驾驶中的意义,自动驾驶在模型部署所关注的内容

TensorRT的优化意义

模型训练 (Training)

- 精度越高越好
 - 把模型做的深一点，宽一点
 - 使用丰富的Data augmentation
 - 使用各种Training trick

模型部署 (Deploy)

- 以精度不变，或者精度掉点很小的情况下尽量压缩模型
 - 减少计算量
 - 减少memory access
 - 提高计算密度

TensorRT的优化意义



Training environment

(训练时的GPU资源可以很丰富，多个GPU分布式训练也都是完全Okay的)



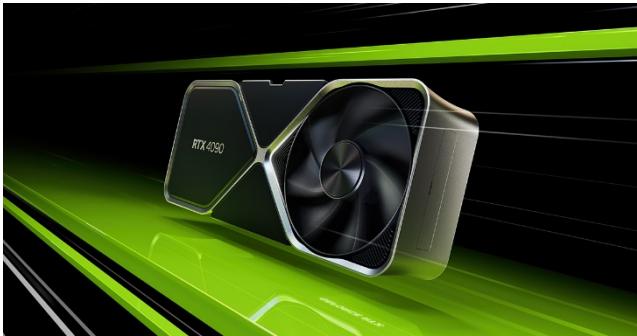
Deploy environment

(我们不可能在一个自动驾驶车上放那么大的服务 器，并且太贵。所以选择一个小一点的)

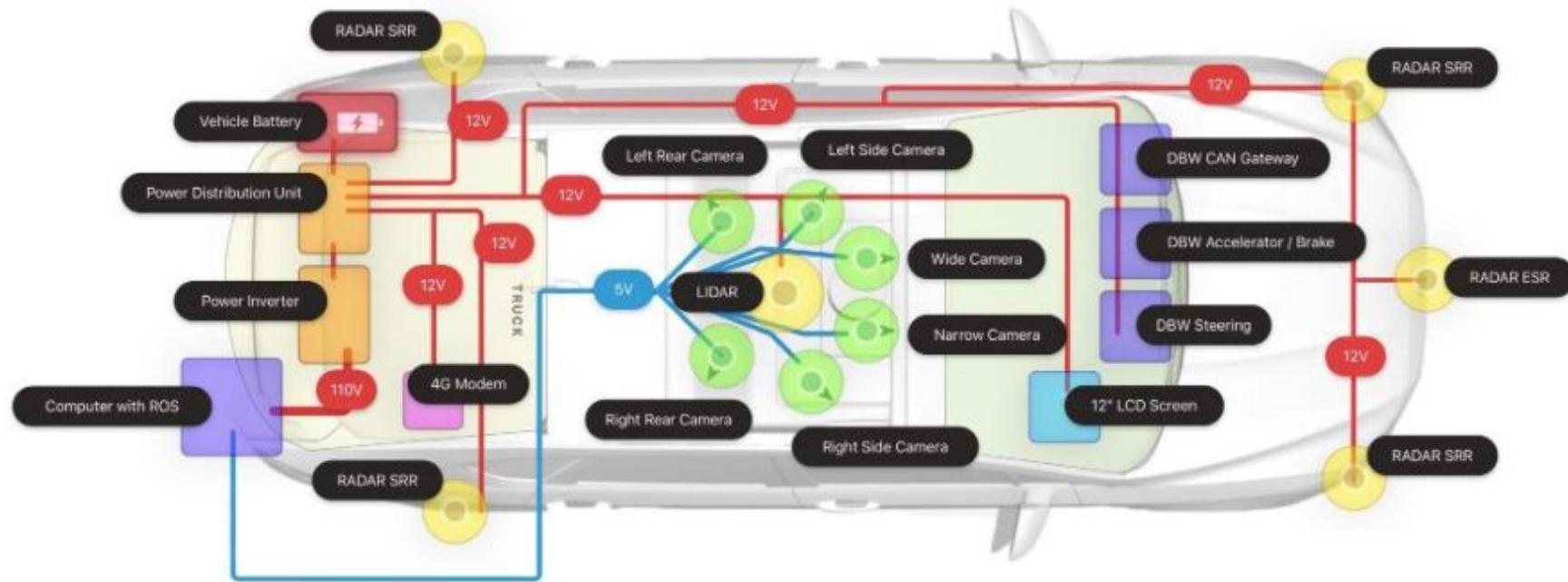
[1] [NVIDIA DGX System](#)

[2] [GeForce RTX 4090](#)

TensorRT的优化意义



Deploy environment^[1, 2]
(我们得考虑GPU的散热，因为经常一不小心就到了80~90°，所以优先考虑水冷)



自动驾驶车中的硬件设备^[3]
(硬件设备多，如果是水冷的话对整个汽车的安全性影响很大)

^[1][GeForce RTX 4090](#)

^[2][bizon z5000](#)

^[3][Under the Hood of a Self-Driving Taxi, a Look at Computation and Other Core Self-Driving Car Systems](#)

TensorRT的优化意义



NVIDIA Jetson AGX Xavier^[1]



NVIDIA Jetson AGX Orin^[2]

Deploy environment:NVIDIA Jetson series
(又要性能高，又要不那么热，那么很多地方就会把
视线集中在Jetson系列上)

^[1][Jetson AGX Xavier](#)

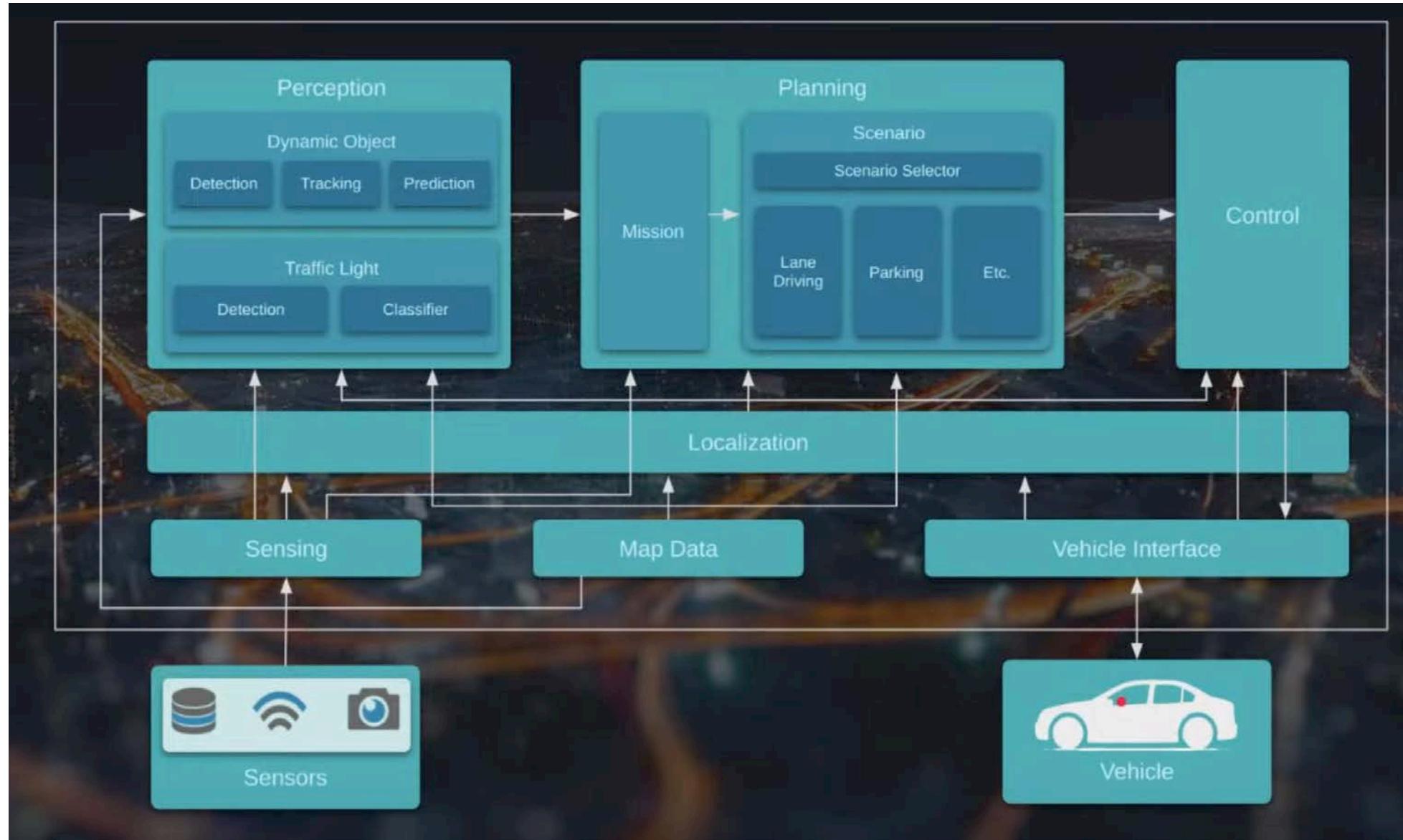
^[2][Jetson AGX Orin at GTC 2022](#)

自动驾驶中需要关注的实时性

Realtime性(实时性)

- 我们在部署模型是需要关注的一个很重要的指标
 - 一般来说15 fps是一个底线 (高速公路上需要至少30fps)
 - $100\text{km/h} = 28\text{m/s}$
- 对于单个小型DNN，在计算资源限制的情况下想达到这个指标难度系数不高
 - 模型可以做的很simple
- 但我们又得需要看的很远
 - 比如说150m远的红绿灯，以及红绿灯下面的箭头和时间
 - 比如说100m远的路标
- 不光是Perception，后期的Planning等等也需要时间

自动驾驶中需要关注的实时性



[1]自動運転におけるAIコンピューティング

自动驾驶中需要关注的实时性

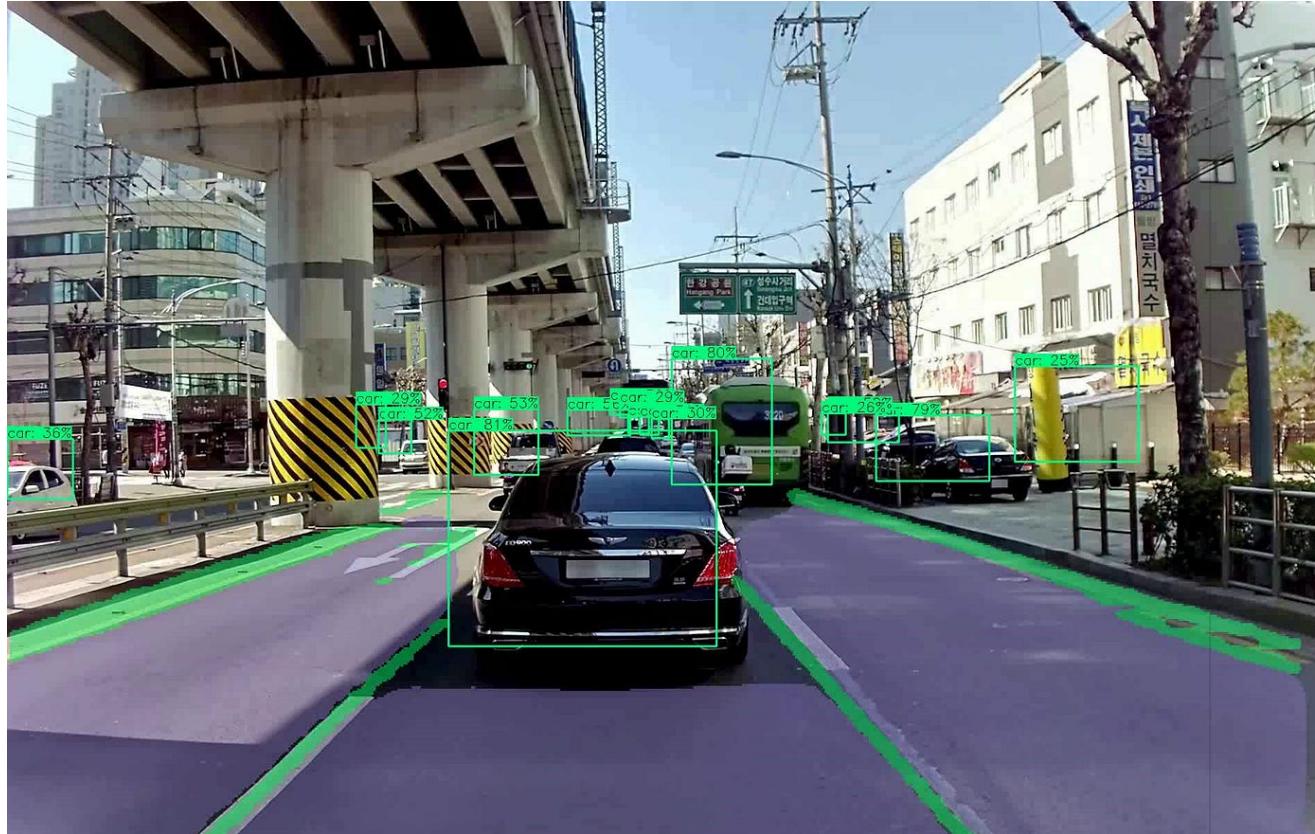
But...

- 多任务
 - 2D/3D 检测
 - 车道线检测
 - 图像分割
 - 深度估计
 - 多目标追踪
 - BEV检测
 - 光流预测
- 多传感器
 - LiDAR
 - Camera
 - Radar
- 高分辨率输入图像
- Multi-camera

如果把这些所有的事情能够在几短时间内做完并达到 $>10\text{fps}$, 甚至 30fps ,
是一个非常具有挑战性的事情

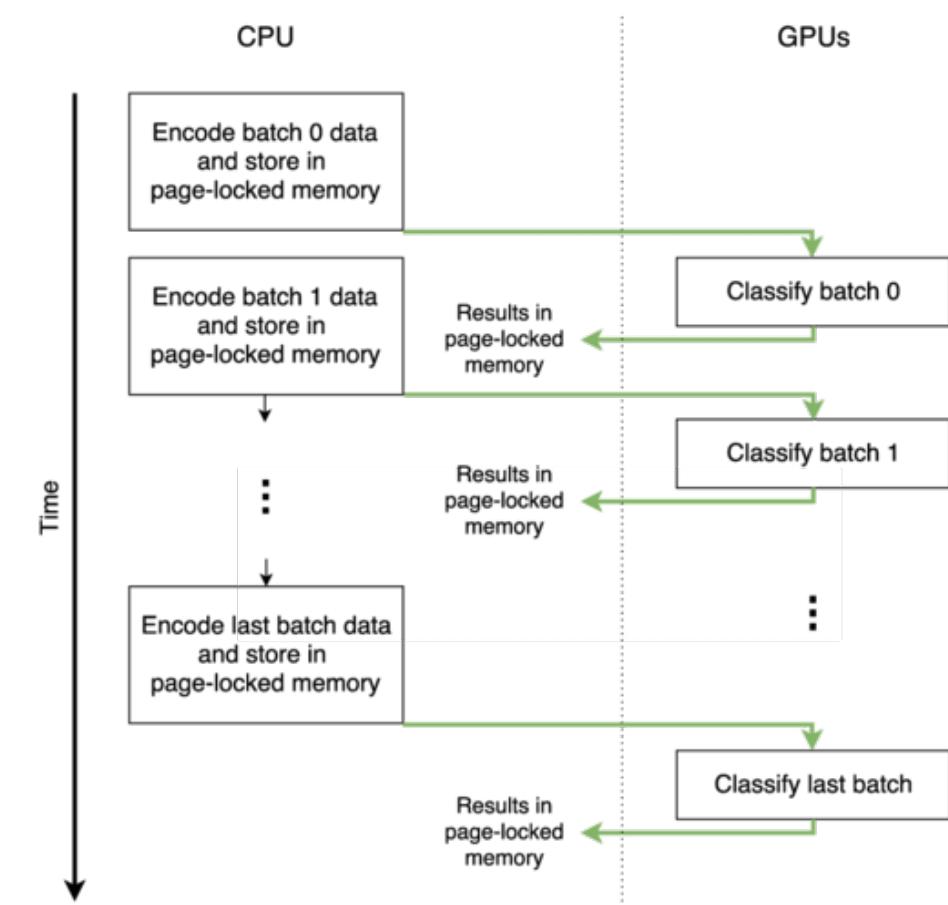
自动驾驶中需要关注的实时性

我们怎么做呢？工业界中有很多种方法，比如说：



Multi-task DNN in autonomous driving^[1]

再比如局部提高分辨率
优先处理局部区域等等



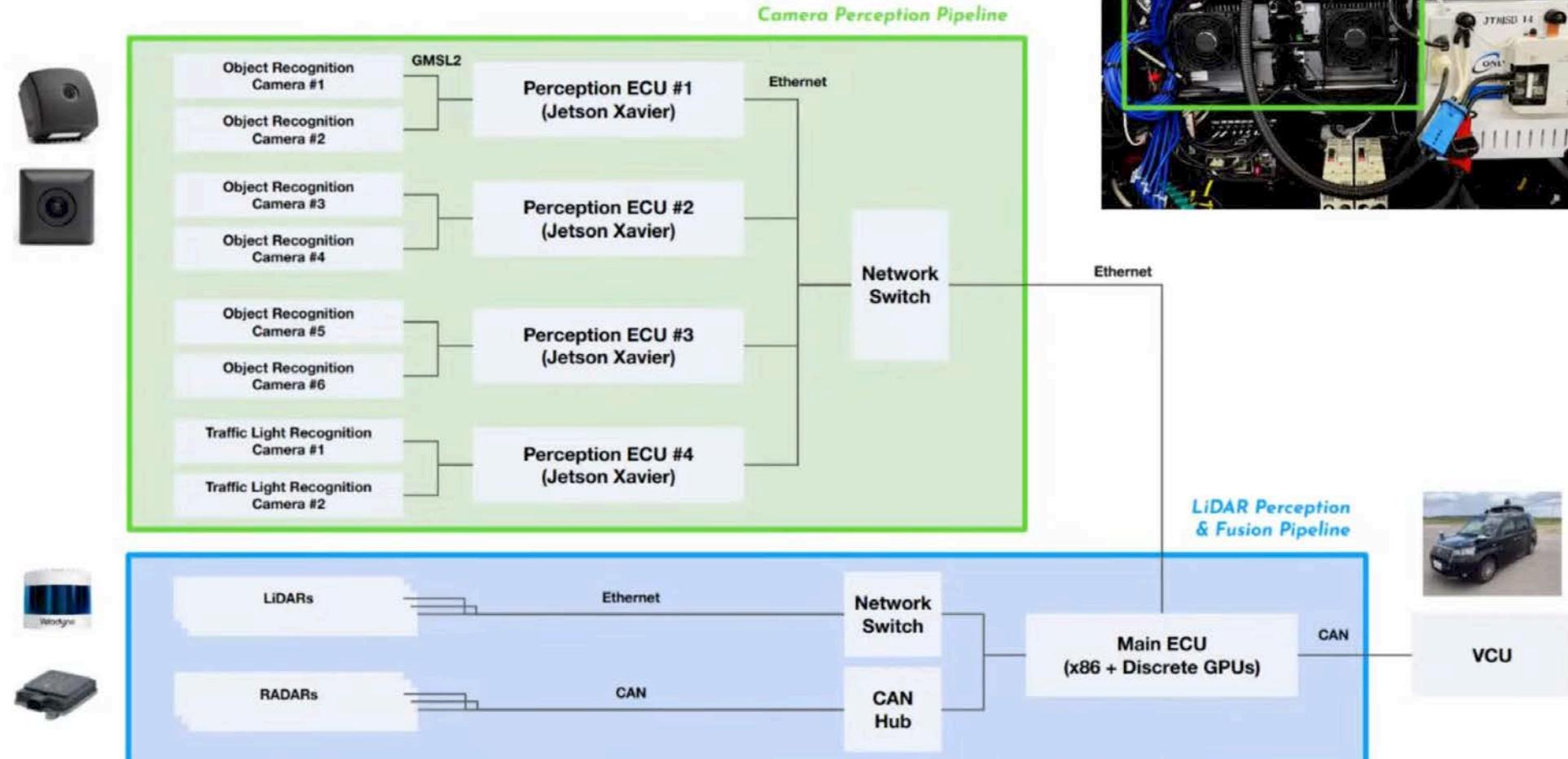
Asynchronous processing in GPU^[2]

^[1][Visual Perception for Self-Driving Cars! Part 5: Multi-Task Learning](#)
^[2][Accelerating metagenomic read classification on CUDA-enabled GPUs](#)

自动驾驶中需要关注的电力消耗

Main ECU (CPU + GPU) ~150W
Perception ECU (Object detection * 6) = ~90w
Perception ECU (Traffic light* 2) = ~20w
合计下来把所有的这些东西全部搭上只做最基本的优化的话大概 ~255w

Sensor/ECU system (example)



自动驾驶ECU方案举例
(camera: yolox, LiDAR: centetpoint)

[\[1\]自動運転におけるAIコンピューティング](#)

自动驾驶中需要关注的电力消耗

基本上我们需要的算力越大，消耗的电力越大。自动驾驶所需要控制的电力一般来说需要控制在**小于100w**。自动驾驶也属于Edge computing的一种，所以对电力消耗需要严格控制

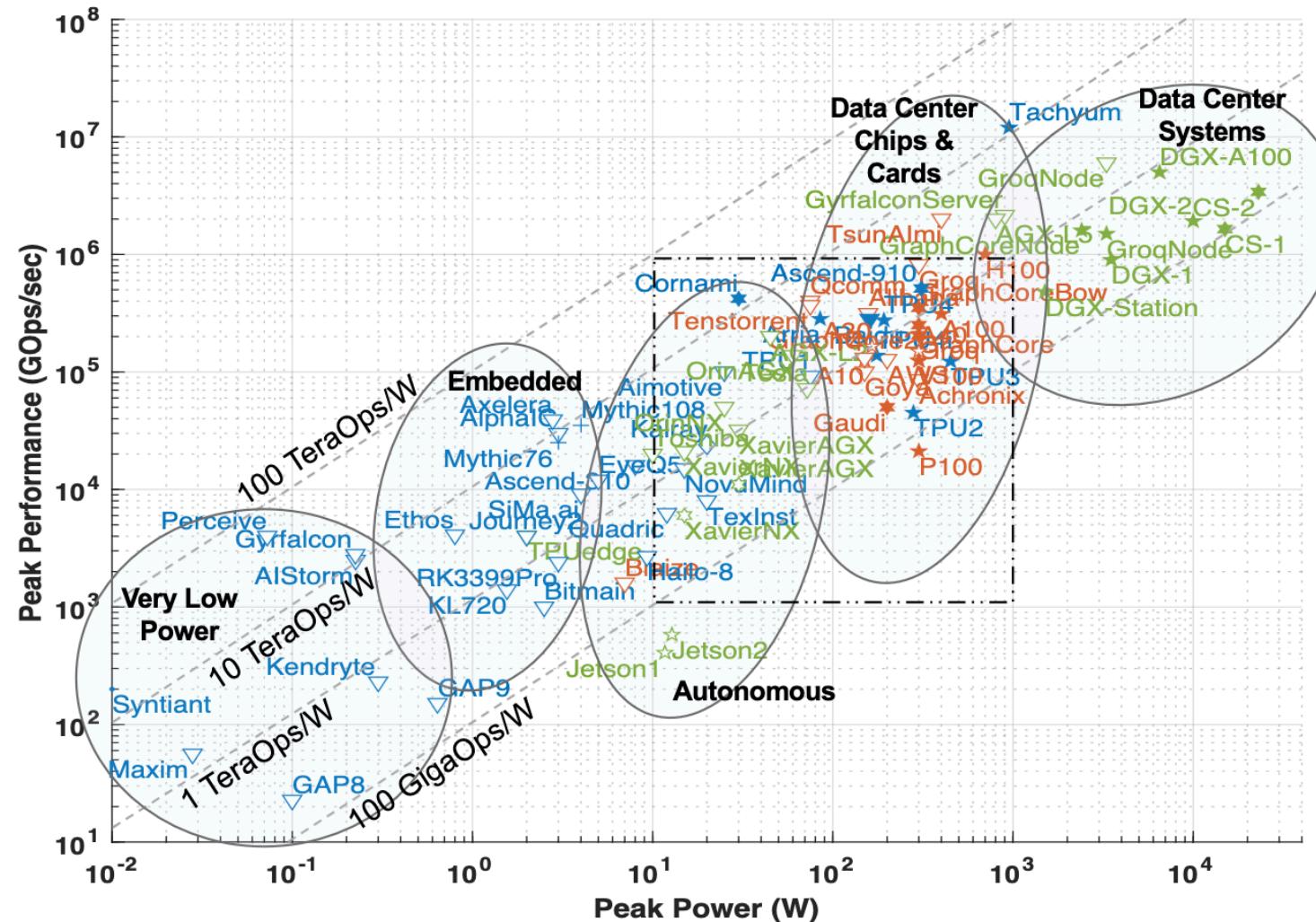


Fig. 2: Peak performance vs. power scatter plot of publicly announced AI accelerators and processors.

Legend

Computation Precision

- | | |
|---|----------|
| + | analog |
| ◀ | int1 |
| ▶ | int2 |
| ● | int4.8 |
| ▼ | int8 |
| ◆ | Int8.32 |
| ▲ | int16 |
| ■ | int12.16 |
| × | int32 |
| ★ | fp16 |
| ☆ | fp16.32 |
| ● | fp32 |
| * | fp64 |

Form Factor

- Chip
 - Card
 - System

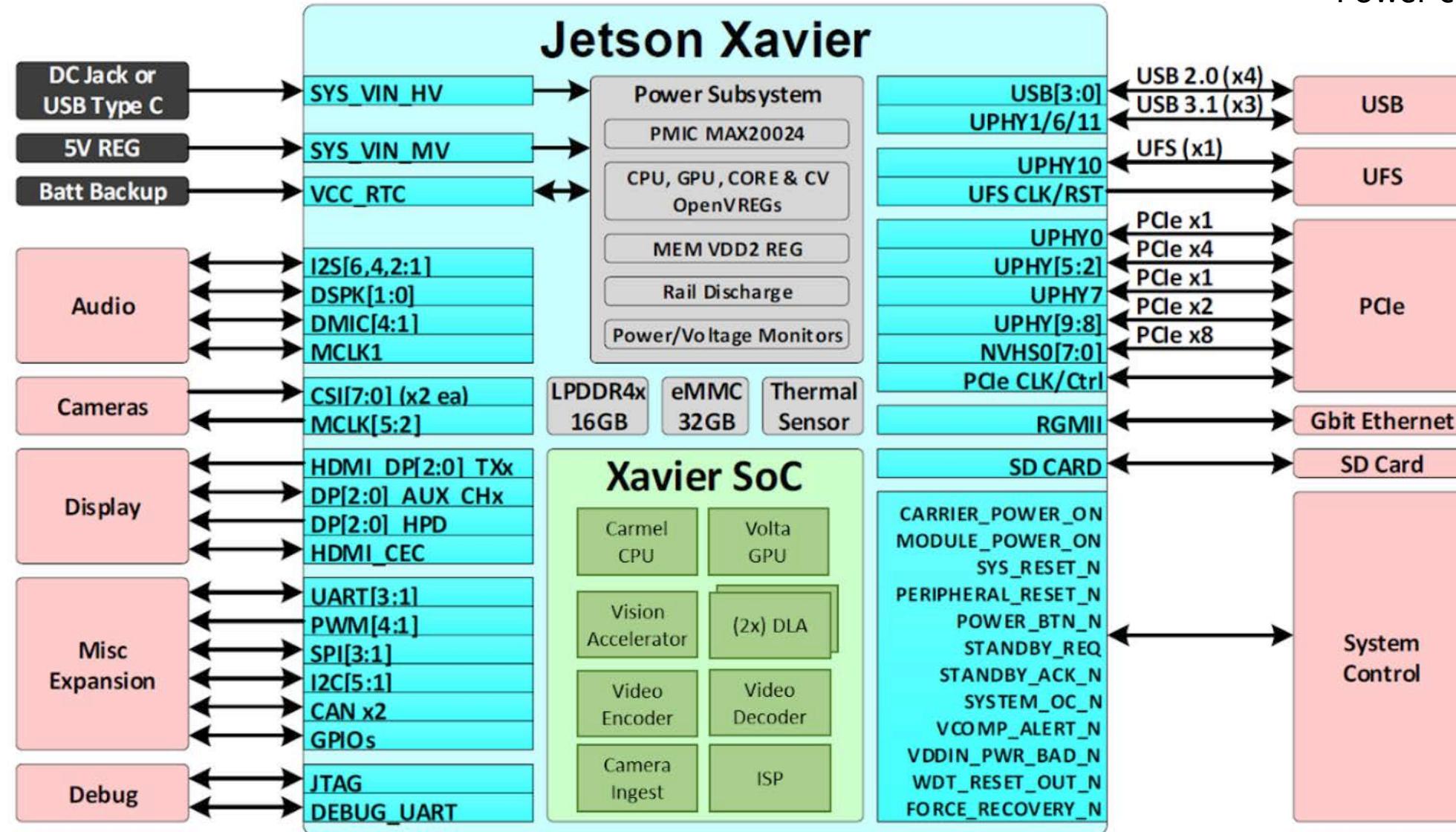
Computation Type

- Inference
 - Training

自动驾驶中需要关注的电力消耗

NVIDIA Jetson AGX Xavier

Peak performance
GPU: 22TOPS (DENSE INT8)
DLA0: 5TOPS (DENSE INT8)
DLA1: 5 TOPS (DENSE INT8)
Power consumption ~30W



自动驾驶中需要关注的电力消耗

NVIDIA Jetson AGX Orin

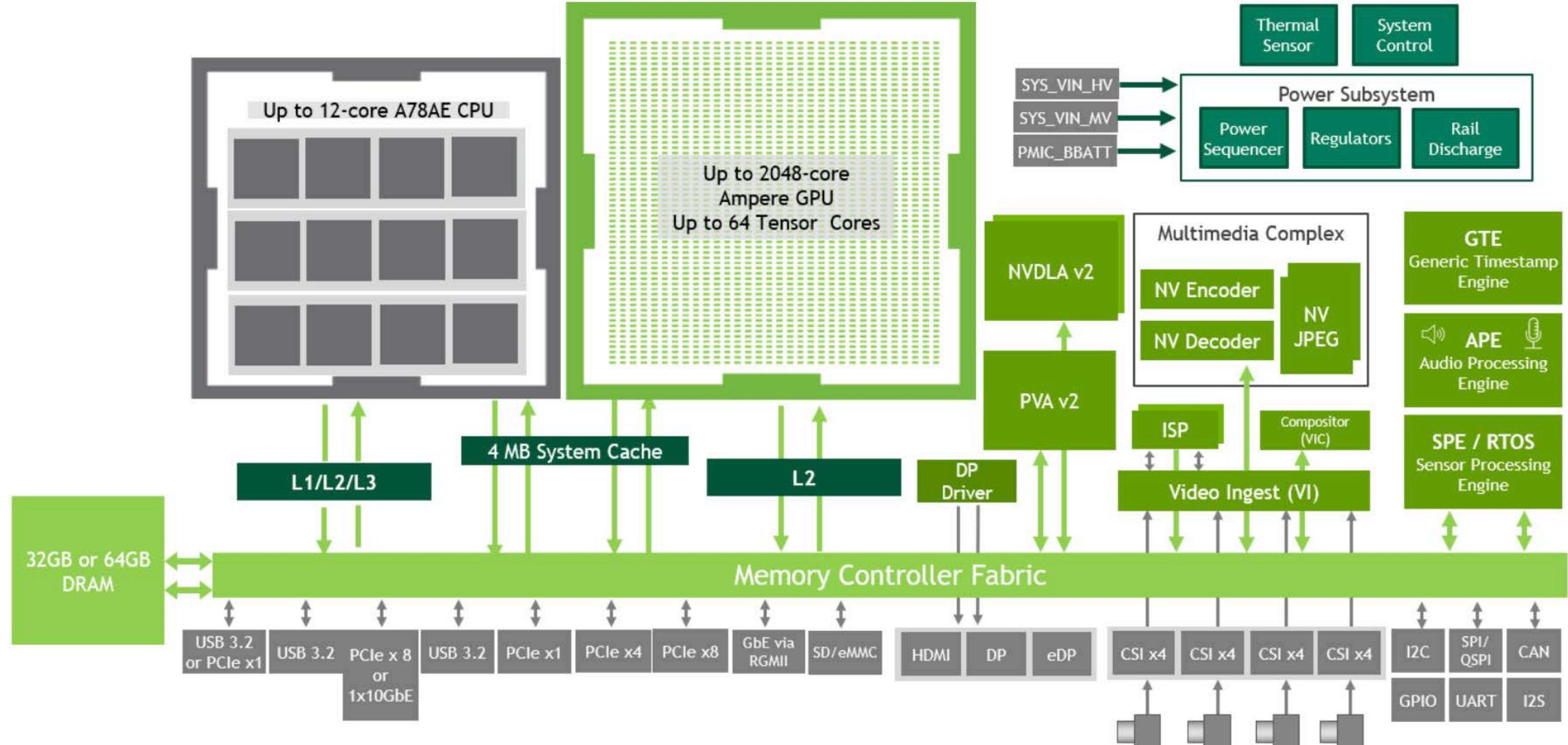
Peak performance

GPU: 170 TOPS (Sparse INT8), 85TOPS (DENSE INT8)

DLA0: 52.5 TOPS (Sparse INT8), 26.25TOPS (DENSE INT8)

DLA1: 52.5 TOPS (Sparse INT8), 26.25TOPS (DENSE INT8)

Power consumption ~60W



自动驾驶中模型部署所关注的东西

重点! !

RealTime(实时性)

Power Consumption(消耗电力)

Long range Accuracy(远距离)



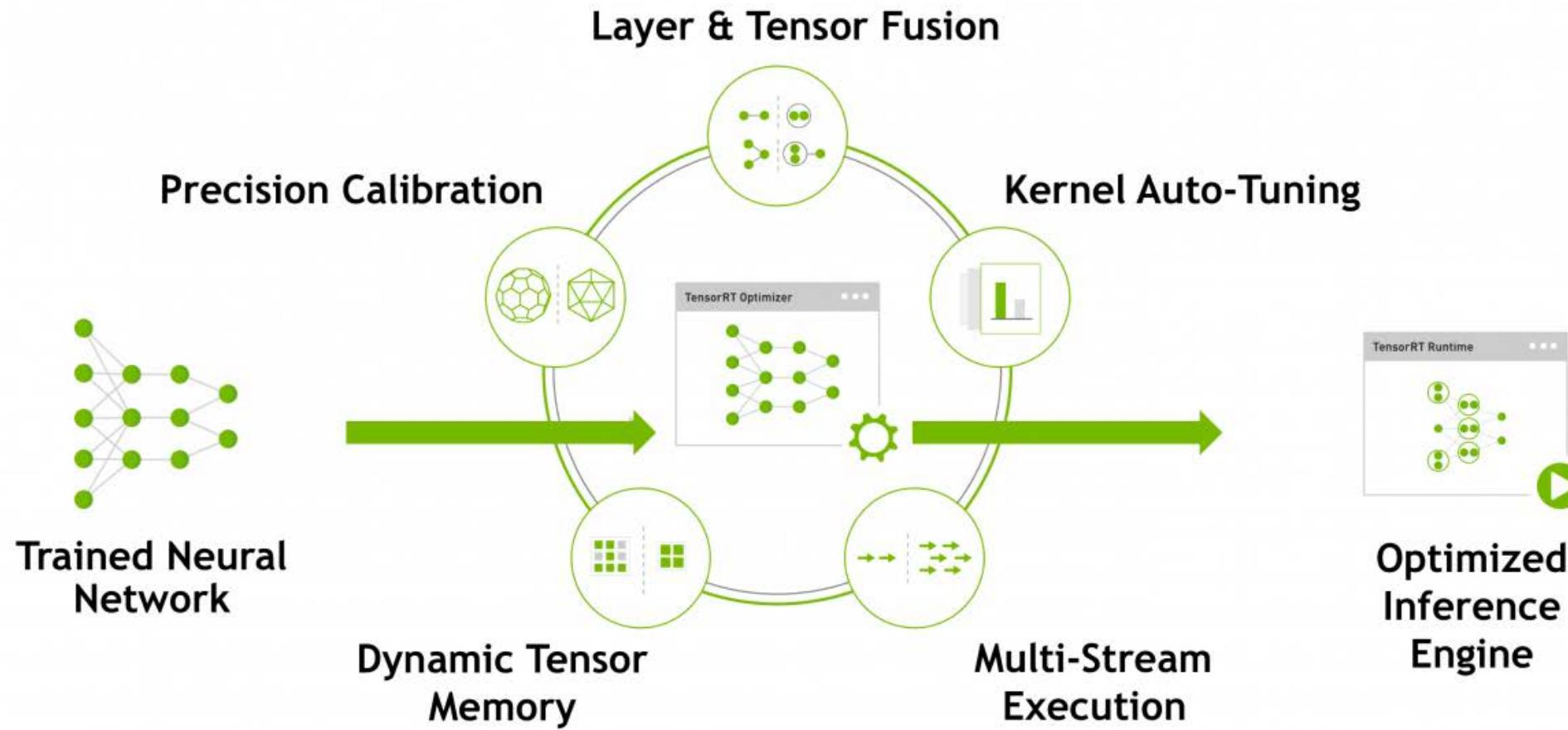
03

TensorRT的模块

Goal: 简单理解TensorRT中每个优化策略，为什么可以做层融合，以及理解量化的意义

TensorRT的优化策略

回顾一下



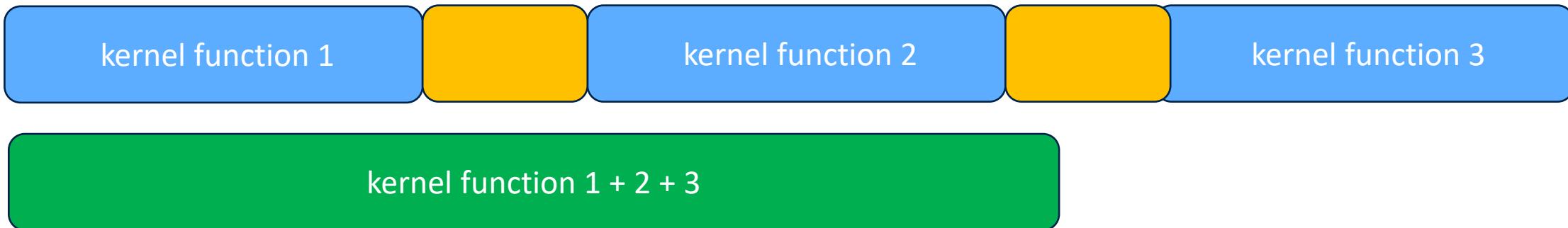
TensorRT的各类优化策略^[1]

^[1]NVIDIA TensorRT

Layer fusion

- 层融合

- Vertical layer fusion (垂直层融合)
- Horizontal layer fusion (水平层融合)

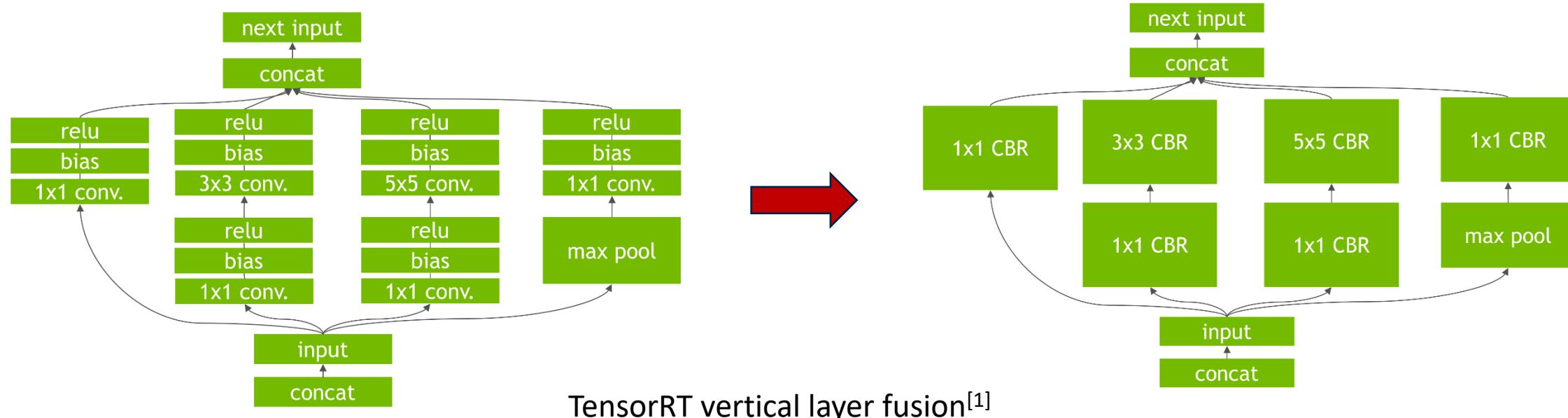


层融合可以减少启动kernel的开销与memory操作，从而提高效率
同时，有些计算可以通过层融合优化后，跟其他计算合并

TensorRT的优化策略

Layer fusion

- 层融合
 - Vertical layer fusion (垂直层融合)
 - 用的比较常见，针对conv + BN + ReLU进行融合



Convolution
kernel function

Batch Normalization
kernel function

ReLU
kernel function

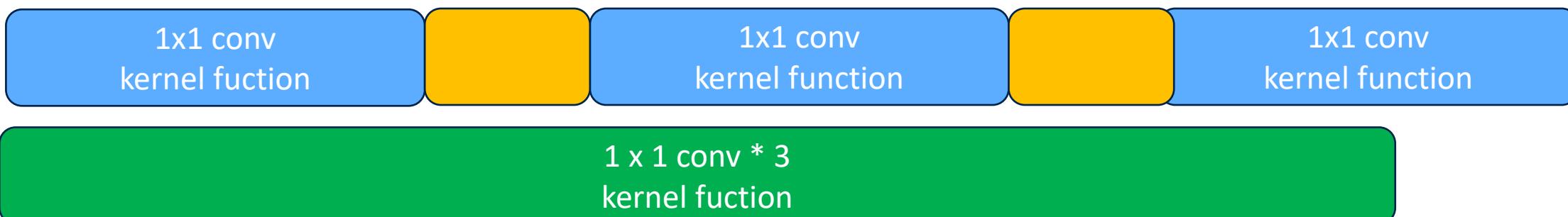
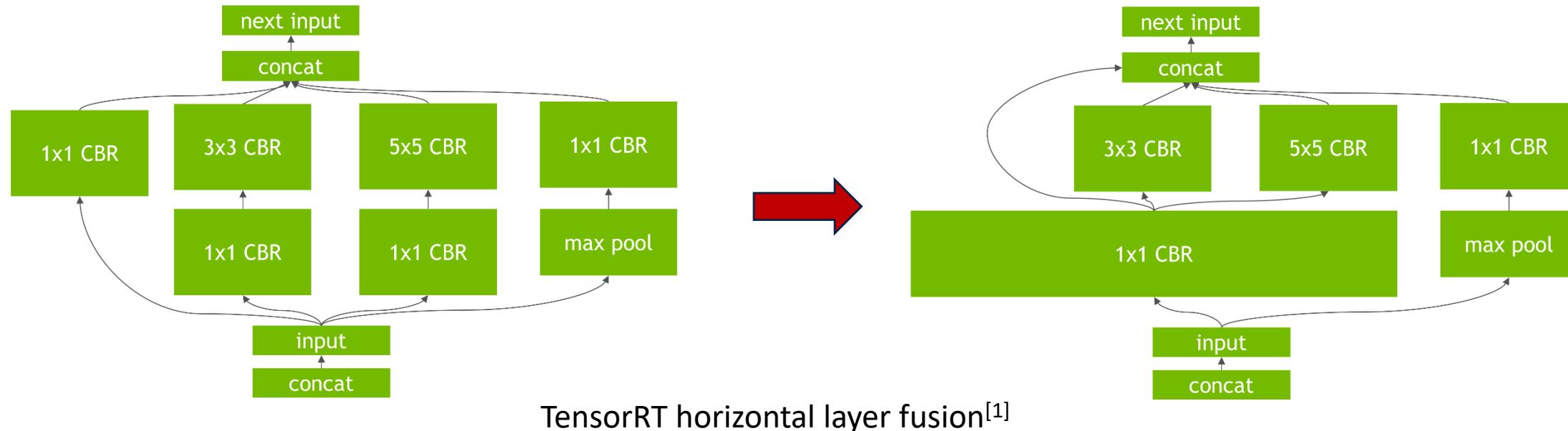
Conv + BN + ReLU
kernel function

思考：为什么conv + BN + ReLU的计算时间和conv的计算时间差不多？

TensorRT的优化策略

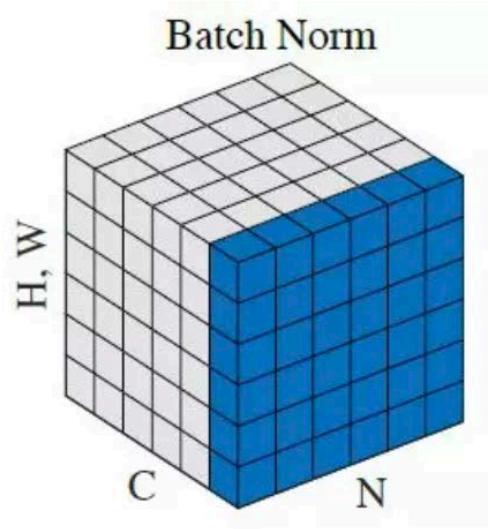
Layer fusion

- 层融合
 - Horizontal layer fusion (水平层融合)
 - 当模型中有水平方向上比较多的同类layer，会直接进行融合



TensorRT的优化策略

回顾一下Batch Normalization的公式



$$\mu_B = \frac{1}{B} \sum_{i=1}^B x_i$$

$$\sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (x_i - \mu_B)^2 + \epsilon$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y_i = \gamma * \hat{x}_i + \beta$$

μ_B : 代表均值

σ_B^2 : 代表方差

ϵ : 一个大于0的浮点数，用于防止分母为0

γ : scaling, 缩放因子

β : shift, 偏移因子

上面这些数值在推理的时候可以提前计算出来

TensorRT的优化策略

convolution + batch normalization融合

$$\mu_B = \frac{1}{B} \sum_{i=1}^B x_i$$

$$\sigma_B^2 = \frac{1}{B} \sum_{i=1}^B (x_i - \mu_B)^2 + \epsilon$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y_i = \gamma * \hat{x}_i + \beta$$

展开

$$y_i = \gamma * \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

代入
 $x_i = w * x_i + b$

$$y = \gamma * \frac{w * x + b - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

TensorRT的优化策略

convolution + batch normalization融合

$$y = \gamma * \frac{w * x + b - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

$$y = \frac{\gamma * w}{\sqrt{\sigma_B^2 + \epsilon}} * x + \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}}(b - \mu_B) + \beta$$

$$y = \left(\frac{\gamma * w}{\sqrt{\sigma_B^2 + \epsilon}}\right) * x + \left(\frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}}(b - \mu_B) + \beta\right)$$

定义

$$\hat{w} = \frac{\gamma * w}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\hat{b} = \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}}(b - \mu_B) + \beta$$

$$y = \hat{w} * x + \hat{b}$$

这两个参数值可以提前计算出来

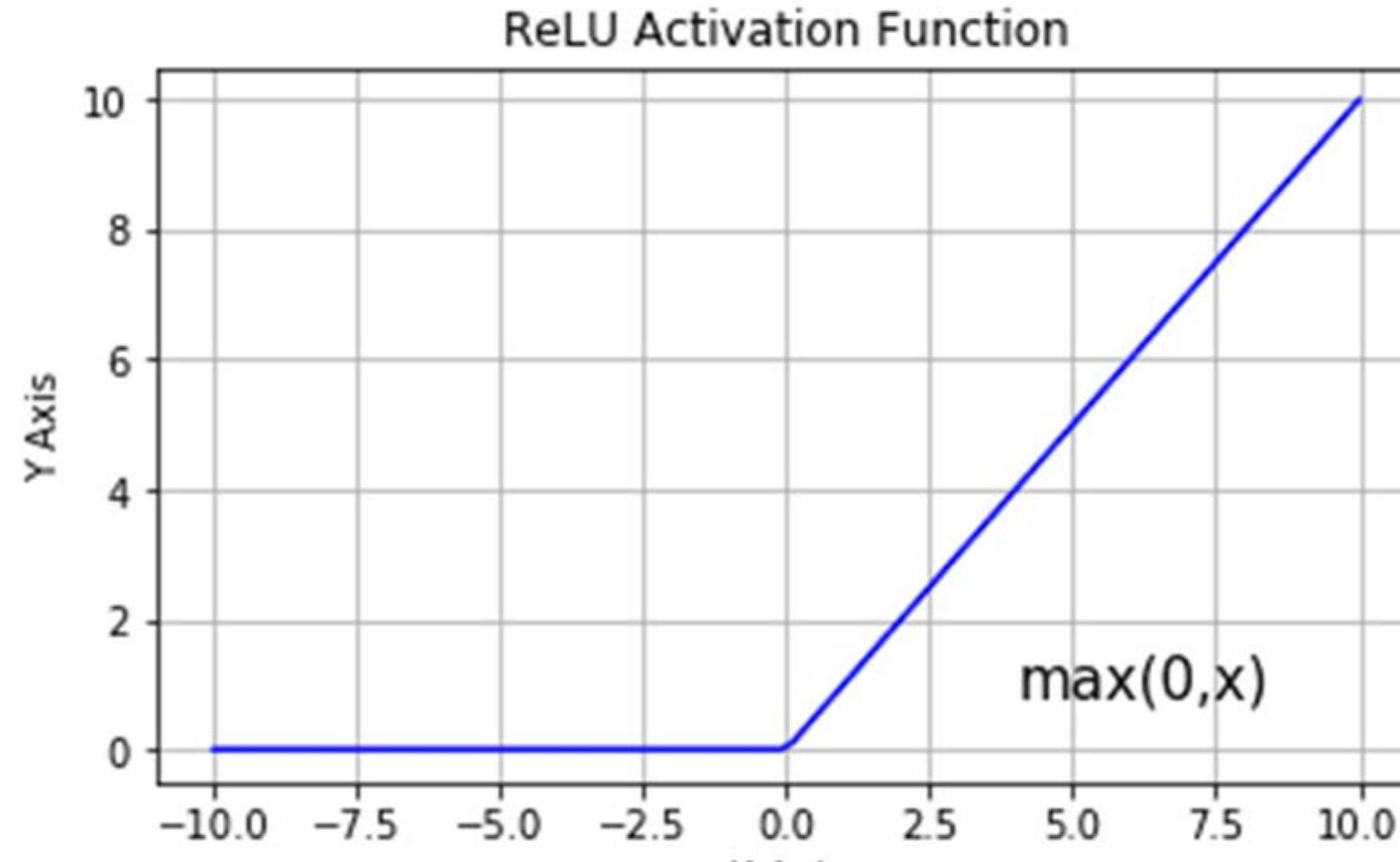
TensorRT的优化策略

convolution + batch normalization + ReLU 融合

$$y = \hat{w} * x + \hat{b}$$

TIPS:

最近的很多模型经常会有很
多种类的activation function，比
如GELU, Swish, Mish等等，这些激
活函数往往由于计算复杂很难
加速，可以尝试改成ReLU看看
精度的改变和性能的提升

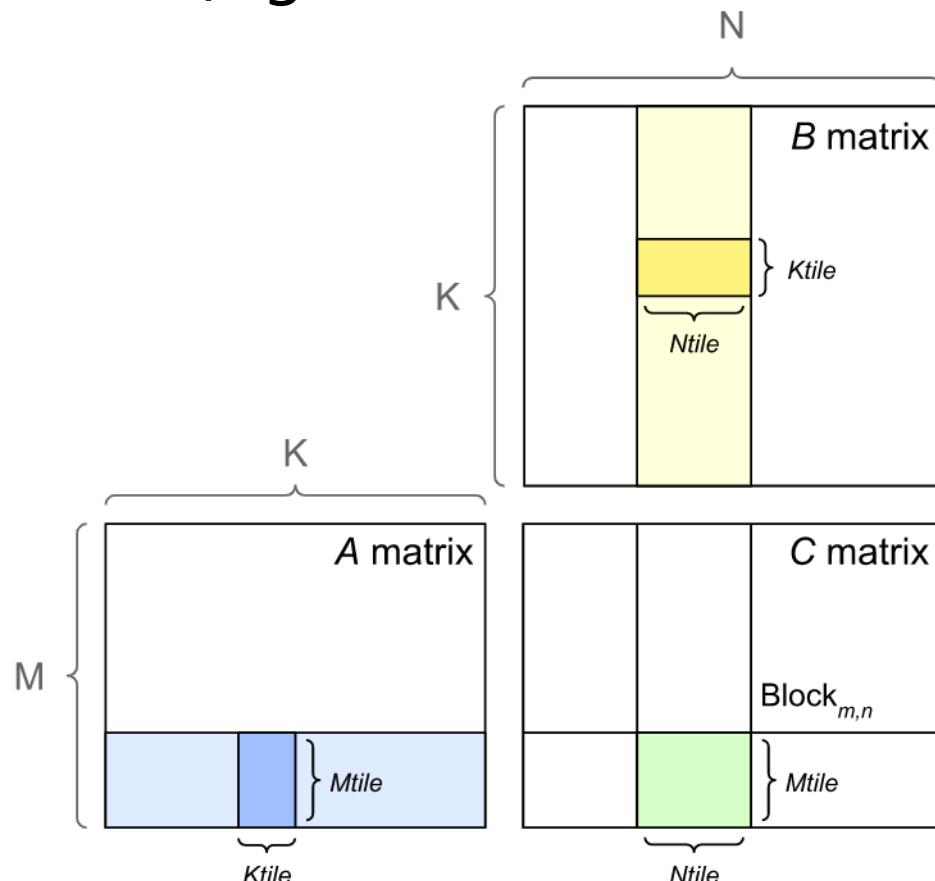


(ReLU只是做一个截取>0的值，内部没有计算)

^[1][What is ReLU and Sigmoid activation function?](#)

Kernel auto-tuning

- TensorRT内部对于同一个层使用各种不同kernel函数进行性能测试
 - 比如对于FC层中的矩阵乘法，根据tile size有很多中kernel function
 - (e.g. 32x32, 32x64, 64x64, 64x128, 128x128, 针对不同硬件有不同策略)



GPU的GEMM计算^[1]



Volta architecture block diagram^[2]

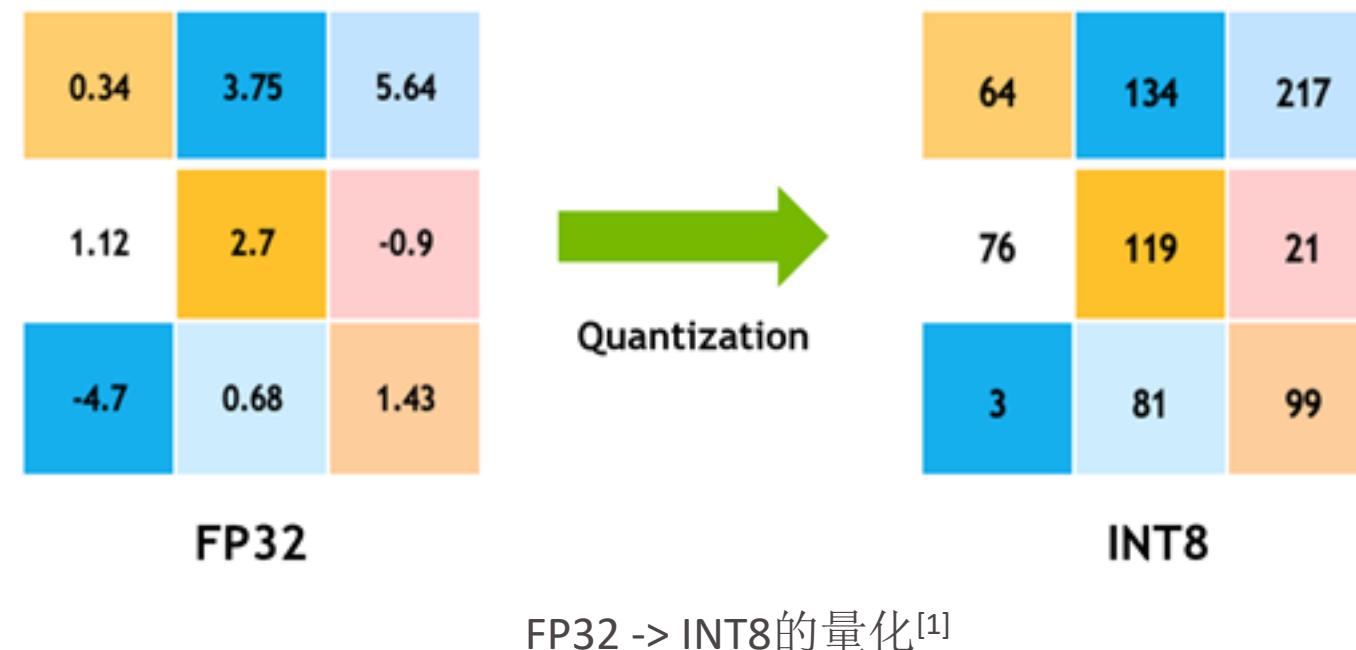
TIPS:

有时你的模型架构由于计算密度不够高的原因，在kernel auto tuning时TensorRT认为没必要选择tensor core而给你分配CUDA core使用

Quantization

- **量化**

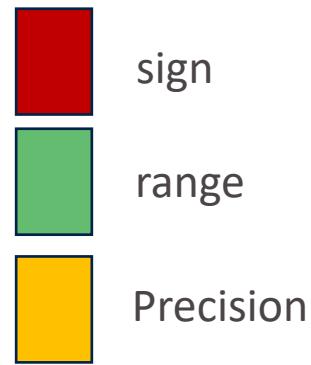
- 压缩模型的一个很重要的策略
- 将单精度类型(FP32)训练权重转变为半精度(FP16)或者整型(INT8, INT4)



^[1][Achieving FP32 Accuracy for INT8 Inference Using Quantization Aware Training with NVIDIA TensorRT](#)

TensorRT的优化策略

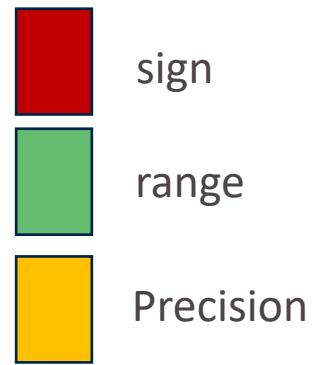
(*)IEEE 745 standard是IEEE于1985年制定的表示浮点数的一种规格



Quantization

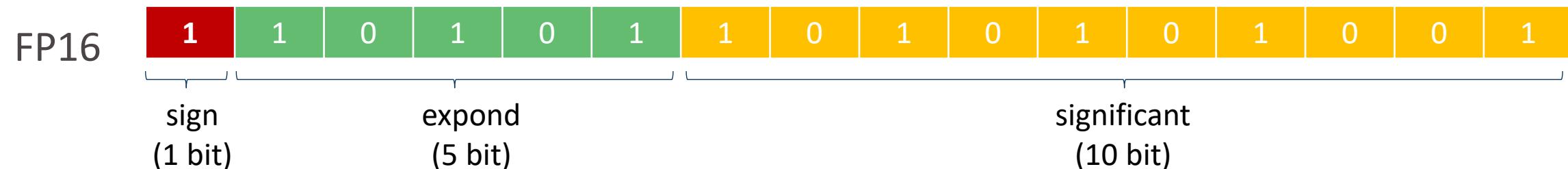
- 如何表示浮点数? (e.g. -6.650490625...)
 - 很显然, 我们用8bit是无法表达这种小数的
 - IEEE 745 standard^(*) 规定了一种表示浮点数(floating point)的方法
 - 一个浮点数可以分为三个部分:
 - sign
 - exponent (range)
 - significant (fraction, precision)





Quantization

- 如何表示浮点数? (e.g. -6.650490625...)



1

$$\text{sign} = -1^1 = -1$$

1	0	0	0	1
2^4	2^3	2^2	2^1	2^0

$$\text{exponent} = 2^4 + 2^0 = 17$$

$$\text{bias} = 2^{k-1} = 2^4 = 16$$

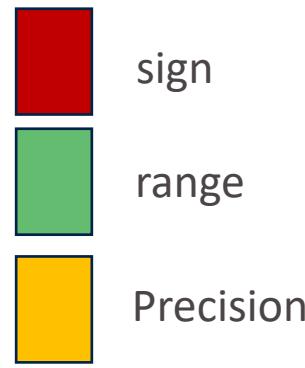
$$\text{biased exponent} = \text{exponent} - \text{bias} = 1$$

0	0	1	0	1	0	1	0	0	1
2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}	2^{-7}	2^{-8}	2^{-9}	2^{-10}

$$\text{significant} = 2^{-1} + 2^{-3} + 2^{-5} + 2^{-7} + 2^{-10} = 0.6650390625$$

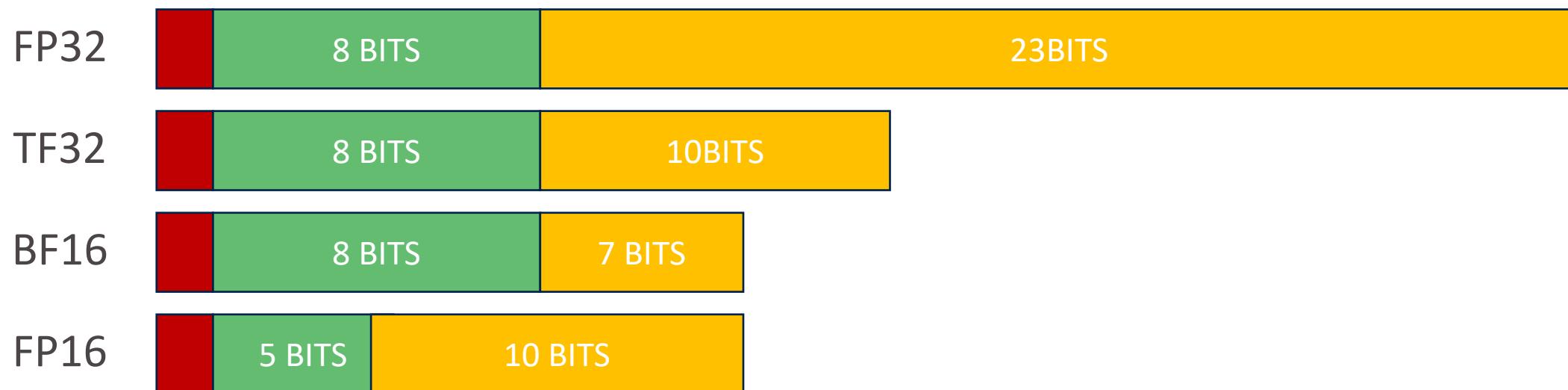
$$-1 \cdot 0.6650390625 \cdot 2^1 = -1.33078124$$

TensorRT的优化策略



Quantization

- 如何表示浮点数?
 - 表示浮点数的方式有很多。
 - FP32, TF32, BF16, FP16
 - 所占用的bits越多，表示的精度越高，但随之而来的就是占用的内存空间越大，计算时间越长



TensorRT的优化策略

(*)其实在Training的时候为了加快速度，也在使用量化的技巧。这个方法被称作Mixed precision training(混合精度学习)。感兴趣的同学可以阅读一下Baidu和NVIDIA发表在ICLR2018的“ Mixed precision training^[1] ” ,很有意思

Quantization

- 量化

- 训练的时候，因为需要优先考虑精度而不需要太重视速度^(*)，所以会使用FP32来表示权重和激活值
- 但部署的时候，我们需要想办法把FP32的数据尽量压缩，能够用**16bits, 8bits, 甚至4 bits**来表示它们
- 这个过程被称作量化(Quantization)



- 由于太过于重要，技巧很多，并且也是目前模型部署最优先考虑的策略，所以下一章单独拿出来详细介绍。目前只需要知道量化是干什么的就okay了

^[1][Mixed precision training](#)

04

导出ONNX，分析ONNX

Goal: 学习PyTorch导出ONNX的方法，以及onnx simplify的使用

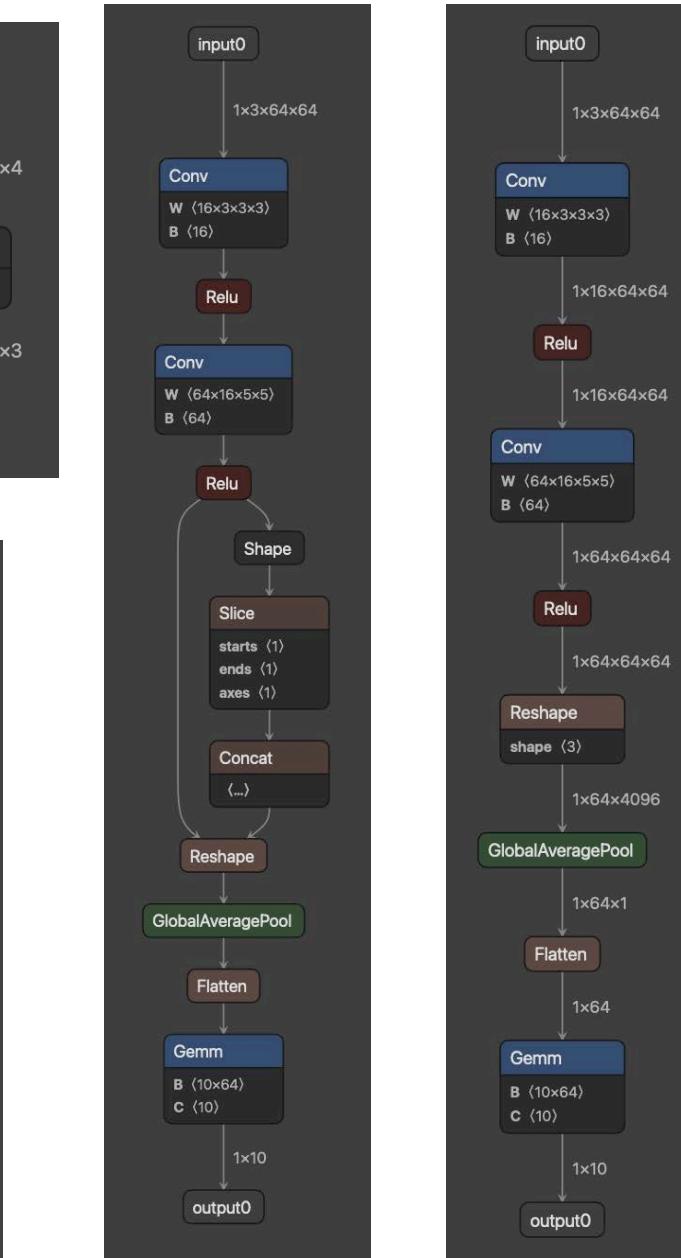
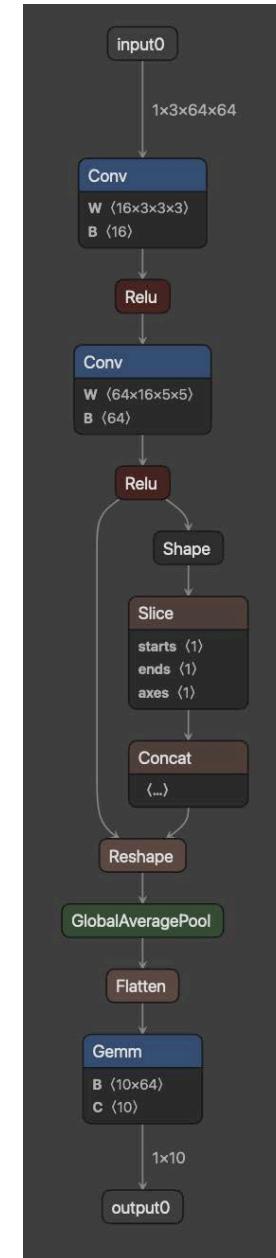
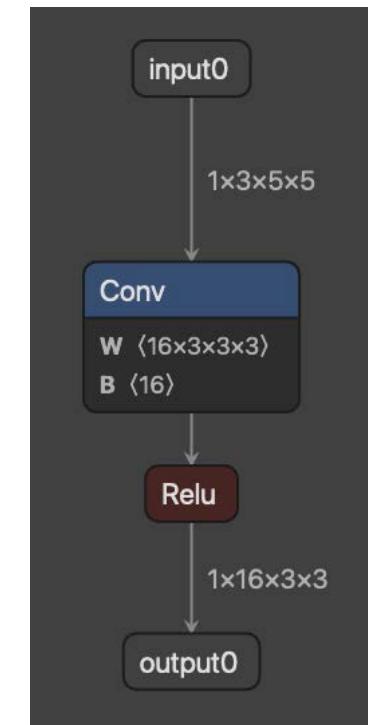
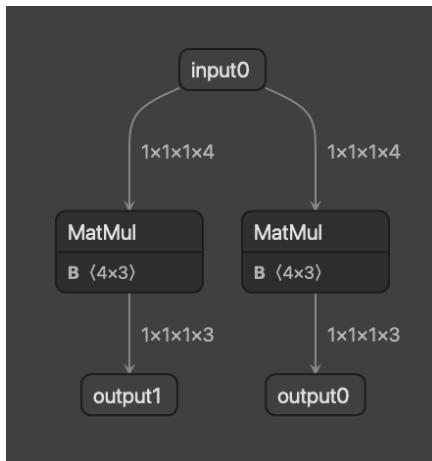
执行一下我们的python程序

```
3.1-generate-onnx/
model
├── example_dynamic_shape.onnx
├── example.onnx
└── example_two_head.onnx
src
├── example_dynamic_shape.py
└── example.py
example_two_head.py
```

```
3.2-export-onnx/
models
├── efficientnetb0.onnx
├── efficientnetV2.onnx
├── mobilenetV3.onnx
├── regnet1.6gf.onnx
├── resnet50.onnx
└── vgg11.onnx
src
├── load_torchvision.py
└── sample_cbr.py
sample_reshape.py
```

3.1-generate-onnx, 3.2-export-onnx

这两个我们今天一起看一下，熟悉导出onnx的步骤



执行结果的一部分展示

04

剖析onnx架构，理解ProtoBuf

Goal: 学习ONNX的Proto架构，使用onnx.helper创建onnx修改onnx

执行一下我们的python程序

```
3.3-read-and-parse-onnx/
└── models
    ├── sample-cbr.onnx
    ├── sample-convnet.onnx
    └── sample-linear.onnx
└── src
    ├── create_onnx_convnet.py
    ├── create_onnx_linear.py
    ├── parse_onnx_cbr.py
    ├── parse_onnx_convnet.py
    ├── parse_onnx_linear.py
    └── parser.py
```

```
*****parse input/output*****
Input info:
    name:      input0
    data Type: 1
    shape:     [1, 3, 64, 64]
Output info:
    name:      input0
    data Type: 1
    shape:     [1, 3, 64, 64]

*****parse node*****
node info:
    name:      conv2d_1
    op_type:   Conv
    inputs:    ['input0', 'conv2d_1.weight', 'conv2d_1.bias']
    outputs:   ['conv2d_1.output']
node info:
    name:      batchNorm1
    op_type:   BatchNormalization
    inputs:    ['conv2d_1.output', 'batchNorm1.scale', 'batchNorm1.bias', 'batchNorm1.mean', 'batchNorm1.var']
    outputs:   ['batchNorm1.output']
node info:
    name:      relu1
    op_type:   Relu
    inputs:    ['batchNorm1.output']
    outputs:   ['relu1.output']
node info:
    name:      avg_pool1
    op_type:   GlobalAveragePool
    inputs:    ['relu1.output']
    outputs:   ['avg_pool1.output']
node info:
    name:      conv2d_2
    op_type:   Conv
    inputs:    ['avg_pool1.output', 'conv2d_2.weight', 'conv2d_2.bias']
    outputs:   ['output0']
```

3.3-read-and-parse-onnx

一起学习onnx中ProtoBuf的框架，学习onnx的细节

执行结果的一部分展示

ONNX是什么？

(*)Protobuf: 全称叫做Protocol Buffer。是Google提出的一套表示和序列化数据的机制

ONNX是一种神经网络的格式，采用Protobuf^(*)二进制形式进行序列化模型。
Protobuf会根据用于定义的数据结构来进行序列化存储
同理，我们可以根据官方提供的数据结构信息，去修改或者创建onnx

1. Define proto

```
import "person_info.proto";

package persons;

message Person {
    PersonInfo info = 1; // characteristics of the person
    repeated Friend friends = 2; // friends of the person
}
```

2. Compile proto

```
_PERSON = _descriptor.Descriptor(
    name='Person',
    full_name='persons.Person',
    filename=None,
    file=DESCRIPTOR,
    containing_type=None,
    create_key=_descriptor._internal_create_key,
    fields=[
        ...
    ])
```

3. Serialize

```
info {
    age: 30
    height: 184
}
friends {
    friendship_duration: 365.1000061035156
    shared_hobbies: "books"
    shared_hobbies: "daydreaming"
    shared_hobbies: "unicorns"
}
person {
    info {
        age: 40
        height: 165
    }
}
```

```
graph {
node {
    input: "input0"
    input: "onnx::Conv_12"
    input: "onnx::Conv_13"
    output: "/conv1/Conv_output_0"
    name: "/conv1/Conv"
    op_type: "Conv"
    attribute {
        name: "dilations"
        type: INTS
        ints: 1
        ints: 1
    }
    attribute {
        name: "group"
        type: INTS
        i: 1
    }
    attribute {
        name: "kernel_shape"
        type: INTS
        ints: 3
        ints: 3
    }
    attribute {
        name: "pads"
        type: INTS
        ints: 0
        ints: 0
        ints: 0
        ints: 0
    }
    attribute {
        name: "strides"
        type: INTS
        ints: 1
        ints: 1
    }
}
```

Protobuf的编译以及序列化流程

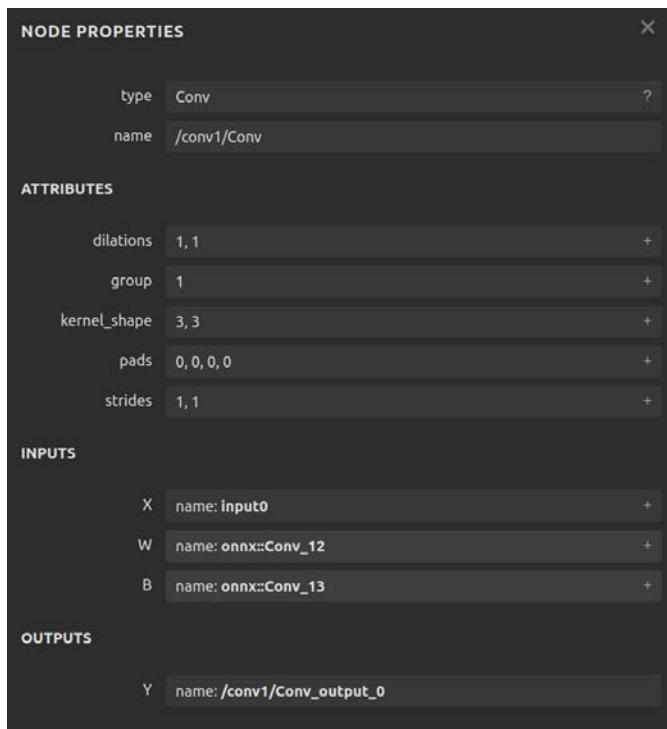
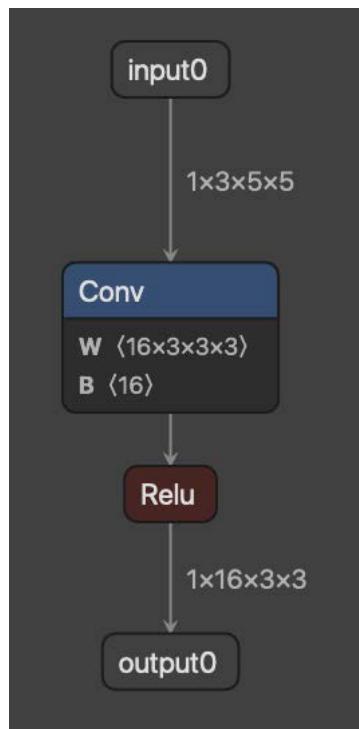
onnx中的各类Proto

onnx的各类proto的定义需要看官方文档 (<https://github.com/onnx/onnx/tree/main>)

这里面的onnx/onnx.in.proto定义了所有onnx的Proto

有关onnx的IR(intermediate representation)信息，看这里(<https://github.com/onnx/onnx/blob/main/docs/IR.md>)

```
4 # 理解 onnx中的组织结构
5 #   - ModelProto (描述的是整个模型的信息)
6 #     --- GraphProto (描述的是整个网络的信息)
7 #       ----- NodeProto (描述的是各个计算节点, 比如conv, linear)
8 #       ----- TensorProto (描述的是tensor的信息, 主要包括权重)
9 #       ----- ValueInfoProto (描述的是input/output信息)
10 #      -----
```



```
435 // Graphs
436 // A graph defines the computational logic of a model and is comprised of a parameter
437 // list of nodes that form a directed acyclic graph based on their inputs and outputs.
438 // This is the equivalent of the "network" or "graph" in many deep learning
439 // frameworks.
440 message GraphProto {
441   // The nodes in the graph, sorted topologically.
442   repeated NodeProto node = 1;
443
444   // The name of the graph.
445   optional string name = 2; // namespace Graph
446
447   // A list of named tensor values, used to specify constant inputs of the graph.
448   // Each initializer (both TensorProto as well SparseTensorProto) MUST have a name.
449   // The name MUST be unique across both initializer and sparse_initializer,
450   // but the name MAY also appear in the input list.
451   repeated TensorProto initializer = 5;
452
453   // Initializers (see above) stored in sparse format.
454   repeated SparseTensorProto sparse_initializer = 15;
455
456   // A human-readable documentation for this graph. Markdown is allowed.
457   optional string doc_string = 10;
458
459   // The inputs and outputs of the graph.
460   repeated ValueInfoProto input = 11;
461   repeated ValueInfoProto output = 12;
462
463   // Information for the values in the graph. The ValueInfoProto.name's
464   // must be distinct. It is optional for a value to appear in value_info list.
465   repeated ValueInfoProto value_info = 13;
466
467   // This field carries information to indicate the mapping among a tensor and its
468   // quantization parameter tensors. For example:
469   // For tensor 'a', it may have {'SCALE_TENSOR', 'a_scale'} and {'ZERO_POINT_TENSOR',
470   // which means, tensor 'a_scale' and tensor 'a_zero_point' are scale and zero point
471   // tensors respectively.
472   repeated TensorAnnotation quantization_annotation = 14;
473
474   reserved 3, 4, 6 to 9;
475   reserved "ir_version", "producer_version", "producer_tag", "domain";
476 }
```

理解onnx中的ValueInfoProto

```
173  
180 // Defines information on value, including the name, the type, and  
181 // the shape of the value.  
182 message ValueInfoProto {  
183     // This field MUST be present in this version of the IR.  
184     optional string name = 1;          // namespace Value  
185     // This field MUST be present in this version of the IR for  
186     // inputs and outputs of the top-level graph.  
187     optional TypeProto type = 2;  
188     // A human-readable documentation for this value. Markdown is allowed.  
189     optional string doc_string = 3;  
190 }
```

一般用来定义网络的input/output
(会根据input/output的type来附加属性)

```
173  
676 // Types  
677 //  
678 // The standard ONNX data types.  
679 message TypeProto {  
680  
681     message Tensor []  
682         // This field MUST NOT have the value of UNDEFINED  
683         // This field MUST have a valid TensorProto.DataType value  
684         // This field MUST be present for this version of the IR.  
685         optional int32 elem_type = 1;  
686         optional TensorShapeProto shape = 2;  
687     }  
688  
689     // repeated T  
690     message Sequence {  
691         // The type and optional shape of each element of the sequence.  
692         // This field MUST be present for this version of the IR.  
693         optional TypeProto elem_type = 1;  
694     };  
695  
696     // map<K,V>  
697     message Map {  
698         // This field MUST have a valid TensorProto.DataType value  
699         // This field MUST be present for this version of the IR.  
700         // This field MUST refer to an integral type ([U]INT{8|16|32|64}) or STRING  
701         optional int32 key_type = 1;  
702         // This field MUST be present for this version of the IR.  
703         optional TypeProto value_type = 2;  
704     };  
705  
706     // wrapper for Tensor, Sequence, or Map  
707     message Optional {  
708         // The type and optional shape of the element wrapped.  
709         // This field MUST be present for this version of the IR.  
710         // Possible values correspond to OptionalProto.DataType enum  
711         optional TypeProto elem_type = 1;  
712     };  
713 }
```

理解onnx中的TensorProto

```
478 // Tensors
479 // A serialized tensor value.
480 // message TensorProto {
481 enum DataType {
482 UNDEFINED = 0;
483 // Basic types.
484 FLOAT = 1; // float
485 UINT8 = 2; // uint8_t
486 INT8 = 3; // int8_t
487 UINT16 = 4; // uint16_t
488 INT16 = 5; // int16_t
489 INT32 = 6; // int32_t
490 INT64 = 7; // int64_t
491 STRING = 8; // string
492 BOOL = 9; // bool
493
494 // IEEE754 half-precision floating-point format (16 bits wide).
495 // This format has 1 sign bit, 5 exponent bits, and 10 mantissa bits.
496 FLOAT16 = 10;
497
498 DOUBLE = 11;
499 UINT32 = 12;
500 INT64 = 13;
501 COMPLEX64 = 14; // complex with float32 real and imaginary components
502 COMPLEX128 = 15; // complex with float64 real and imaginary components
503
504 // Non-IEEE floating-point format based on IEEE754 single-precision
505 // floating-point number truncated to 16 bits.
506 // This format has 1 sign bit, 8 exponent bits, and 7 mantissa bits.
507 BFLOAT16 = 16;
508
509 // Non-IEEE floating-point format based on papers
510 // FP8 Formats for Deep Learning, https://arxiv.org/abs/2209.05433,
511 // 8-bit Numerical Formats For Deep Neural Networks, https://arxiv.org/pdf/2206.02
512 // Operators supported FP8 are Cast, CastLike, QuantizeLinear, DequantizeLinear.
513 // The computation usually happens inside a block quantize / dequantize
514 // fused by the runtime.
515 FLOAT8E4M3FN = 17; // float 8, mostly used for coefficients, supports nan, not
516 FLOAT8E4M3FNNUZ = 18; // float 8, mostly used for coefficients, supports nan, not
517 FLOAT8E5M2 = 19; // follows IEEE 754, supports nan, inf, mostly used for grad
518 FLOAT8E5M2FNNUZ = 20; // follows IEEE 754, supports nan, inf, mostly used for grad
519
520 // Future extensions go here.
521
522 } // The shape of the tensor.
523 repeated int64 dims = 1;
524
525 // The data type of the tensor.
526 // This field MUST have a valid TensorProto.DataType value
527 optional int32 data_type = 2;
528
529 // For very large tensors, we may want to store them in chunks, in which
530 // case the following fields will specify the segment that is stored in
531 // the current TensorProto.
532 message Segment {
533 optional int64 begin = 1;
534 optional int64 end = 2;
535 }
536 optional Segment segment = 3;
537
538 // Tensor content must be organized in row-major order.
539
540 // Depending on the data_type field, exactly one of the fields below with
541 // name ending in _data is used to store the elements of the tensor.
542
543 // For float and complex64 values
544 // Complex64 tensors are encoded as a single array of floats,
545 // with the real components appearing in odd numbered positions,
546 // and the corresponding imaginary component appearing in the
547 // subsequent even numbered position. (e.g., [1.0 + 2.0i, 3.0 + 4.0i]
548 // is encoded as [1.0, 2.0 ,3.0 ,4.0])
549
550 // When this field is present, the data_type field MUST be FLOAT or COMPLEX64.
551 repeated float float_data = 4 [packed = true];
552
553 // For int32, uint8, int8, uint16, int16, bool, float8, and float16 values
554 // float16 and float8 values must be bit-wise converted to an int16_t prior
555 // to writing to the buffer.
556
557 // When this field is present, the data_type field MUST be
558 // INT32, INT16, INT8, UINT16, UINT8, BOOL, FLOAT16, BFLOAT16, FLOAT8E4M3FN, FLOAT8E
559 repeated int32 int32_data = 5 [packed = true];
560
561 // For strings.
562 // Each element of string_data is a UTF-8 encoded Unicode
563 // string. No trailing null, no leading BOM. The protobuf "string"
564 // scalar type is not used to match ML community conventions.
565
566 // When this field is present, the data_type field MUST be STRING
567 repeated bytes string_data = 6;
568
569 // For int64.
570 // When this field is present, the data_type field MUST be INT64
571 repeated int64 int64_data = 7 [packed = true];
572
573 // Optionally, a name for the tensor.
574 optional string name = 8; // namespace Value
575
576 // A human-readable documentation for this tensor. Markdown is allowed.
577 optional string doc_string = 12;
578
579 // Serializations can either use one of the fields above, or use this
580 // raw bytes field. The only exception is the string case, where one is
581 // required to store the content in the repeated bytes string_data field.
582
583 // When this raw_data field is used to store tensor value, elements MUST
584 // be stored in as fixed-width, little-endian order.
585 // Floating-point data types MUST be stored in IEEE 754 format.
586 // Complex64 elements must be written as two consecutive FLOAT values, real compon
587 // Complex128 elements must be written as two consecutive DOUBLE values, real compon
588 // Boolean type MUST be written one byte per tensor element (00000001 for true, 0000
589 //
590 // Note: the advantage of specific field rather than the raw_data field is
591 // that in some cases (e.g. int data), protobuf does a better packing via
592 // variable length storage, and may lead to smaller binary footprint.
593 // When this field is present, the data_type field MUST NOT be STRING or UNDEFINED
594 optional bytes raw_data = 9;
595
596 // Data can be stored inside the protobuf file using type-specific fields or raw_dat
597 // Alternatively, raw bytes data can be stored in an external file, using the extern
598 // external_data stores key-value pairs describing data location. Recognized keys ar
599 // - "location" (required) - POSIX filesystem path relative to the directory where t
600 // protobuf model was stored
601 // - "offset" (optional) - position of byte at which stored data begins. Integer sto
602 // Offset values SHOULD be multiples 4096 (page size) to ena
603 // - "length" (optional) - number of bytes containing data. Integer stored as string
604 // - "checksum" (optional) - SHA1 digest of file specified in under 'location' key.
605 repeated StringStringEntryProto external_data = 13;
606
607 // Location of the data for this tensor. MUST be one of:
608 // - DEFAULT - data stored inside the protobuf message. Data is stored in raw_data (
609 // - EXTERNAL - data stored in an external location as described by external_data fi
610 enum DataLocation {
611 DEFAULT = 0;
612 EXTERNAL = 1;
613 }
```

一般用来定义一个权重，比如conv的w和b
(dims是repeated类型，意味着是数组)
(raw_data是bytes类型)

理解onnx中的NodeProto

```
192 // Nodes
193 //
194 // Computation graphs are made up of a DAG of nodes, which represent what is
195 // commonly called a "layer" or "pipeline stage" in machine learning frameworks.
196 //
197 // For example, it can be a node of type "Conv" that takes in an image, a filter
198 // tensor and a bias tensor, and produces the convolved output.
199 message NodeProto {
200     repeated string input = 1;      // namespace Value
201     repeated string output = 2;    // namespace Value
202
203     // An optional identifier for this node in a graph.
204     // This field MAY be absent in this version of the IR.
205     optional string name = 3;      // namespace Node
206
207     // The symbolic identifier of the Operator to execute.
208     optional string op_type = 4;    // namespace Operator
209     // The domain of the OperatorSet that specifies the operator named by op_type.
210     optional string domain = 7;    // namespace Domain
211
212     // Additional named attributes.
213     repeated AttributeProto attribute = 5;
214
215     // A human-readable documentation for this node. Markdown is allowed.
216     optional string doc_string = 6;
217 }
```

一般用来定义一个计算节点，比如conv, linear
(input是repeated类型，意味着是数组)
(output是repeated类型，意味着是数组)
(attribute有一个自己的Proto)
(op_type需要严格根据onnx所提供的Operators写)

RandomUniform	1
RandomUniformLike	1
Reciprocal	13, 6, 1
ReduceMax	18, 13, 12, 11, 1
ReduceMean	18, 13, 11, 1
ReduceMin	18, 13, 12, 11, 1
ReduceProd	18, 13, 11, 1
ReduceSum	13, 11, 1
Reshape	19, 14, 13, 5, 1
Resize	19, 18, 13, 11, 10
ReverseSequence	10
RoiAlign	16, 10
Round	11
STFT	17
Scan	19, 16, 11, 9, 8
Scatter (deprecated)	11, 9
ScatterElements	18, 16, 13, 11
ScatterND	18, 16, 13, 11
ScatterUpdate	11

理解onnx中的AttributeProto

```
110 // Attributes
111 // 
112 // A named attribute containing either singular float, integer, string, graph,
113 // and tensor values, or repeated float, integer, string, graph, and tensor values.
114 // An AttributeProto MUST contain the name field, and *only one* of the
115 // following content fields, effectively enforcing a C/C++ union equivalent.
116 // Note: this enum is structurally identical to the OpSchema::AttrType
117 message AttributeProto {
118     // enum defined in schema.h. If you rev one, you likely need to rev the other.
119     enum AttributeType {
120         UNDEFINED = 0;
121         FLOAT = 1;
122         INT = 2;
123         STRING = 3;
124         TENSOR = 4;
125         GRAPH = 5;
126         SPARSE_TENSOR = 11;
127         TYPE_PROTO = 13;
128
129         FLOATS = 6;
130         INTS = 7;
131         STRINGS = 8;
132         TENSORS = 9;
133         GRAPHS = 10;
134         SPARSE_TENSORS = 12;
135         TYPE_PROTOS = 14;
136     }
137
138     // The name field MUST be present for this version of the IR.
139     optional string name = 1;           // namespace Attribute
140
141     // if ref_attr_name is not empty, ref_attr_name is the attribute name in parent func
142     // In this case, this AttributeProto does not contain data, and it's a reference of
143     // in parent scope.
144     // NOTE: This should ONLY be used in function (sub-graph). It's invalid to be used in
145     optional string ref_attr_name = 21;
146
147     // A human-readable documentation for this attribute. Markdown is allowed.
148     optional string doc_string = 13;
149
150     // The type field MUST be present for this version of the IR.
151     // For 0.0.1 versions of the IR, this field was not defined, and
152     // implementations needed to use has_field heuristics to determine
153     // which value field was in use. For IR_VERSION 0.0.2 or later, this
```

```
154         // field MUST be set and match the field in use. This
155         // change was made to accommodate proto3 implementations.
156         optional AttributeType type = 20; // discriminator that indicates which field below
157
158         // Exactly ONE of the following fields must be present for this version of the IR
159         optional float f = 2;           // float
160         optional int64 i = 3;          // int
161         optional bytes s = 4;         // UTF-8 string
162         optional TensorProto t = 5;    // tensor value
163         optional GraphProto g = 6;     // graph
164         optional SparseTensorProto sparse_tensor = 22; // sparse tensor value
165
166         // Do not use field below, it's deprecated.
167         // optional ValueProto v = 12;    // value - subsumes everything but graph
168         optional TypeProto tp = 14;    // type proto
169
170         repeated float floats = 7;    // list of floats
171         repeated int64 ints = 8;      // list of ints
172         repeated bytes strings = 9;  // list of UTF-8 strings
173         repeated TensorProto tensors = 10; // list of tensors
174         repeated GraphProto graphs = 11; // list of graph
175         repeated SparseTensorProto sparse_tensors = 23; // list of sparse tensors
176         repeated TypeProto type_protos = 15; // list of type protos
177
178 }
```

一般用来定义一个node的属性。比如说kernel size

(

比较常见的方式就是把(key, value)传入Proto，之后

name = key

i = value

)

理解onnx中的GraphProto

```
435 // Graphs
436 // A graph defines the computational logic of a model and is comprised of a parameter
437 // list of nodes that form a directed acyclic graph based on their inputs and outputs.
438 // This is the equivalent of the "network" or "graph" in many deep learning
439 // frameworks.
440 message GraphProto {
441     // The nodes in the graph, sorted topologically.
442     repeated NodeProto node = 1;
443
444     // The name of the graph.
445     optional string name = 2;    // namespace Graph
446
447     // A list of named tensor values, used to specify constant inputs of the graph.
448     // Each initializer (both TensorProto as well SparseTensorProto) MUST have a name.
449     // The name MUST be unique across both initializer and sparse_initializer,
450     // but the name MAY also appear in the input list.
451     repeated TensorProto initializer = 5;
452
453     // Initializers (see above) stored in sparse format.
454     repeated SparseTensorProto sparse_initializer = 15;
455
456     // A human-readable documentation for this graph. Markdown is allowed.
457     optional string doc_string = 10;
458
459     // The inputs and outputs of the graph.
460     repeated ValueInfoProto input = 11;
461     repeated ValueInfoProto output = 12;
462
463     // Information for the values in the graph. The ValueInfoProto.name's
464     // must be distinct. It is optional for a value to appear in value_info list.
465     repeated ValueInfoProto value_info = 13;
466
467     // This field carries information to indicate the mapping among a tensor and its
468     // quantization parameter tensors. For example:
469     // For tensor 'a', it may have {'SCALE_TENSOR', 'a_scale'} and {'ZERO_POINT_TENSOR',
470     // which means, tensor 'a_scale' and tensor 'a_zero_point' are scale and zero point
471     repeated TensorAnnotation quantization_annotation = 14;
472
473     reserved 3, 4, 6 to 9;
474     reserved "ir_version", "producer_version", "producer_tag", "domain";
475 }
```

(*)initializer: 在onnx中一般表示权重信息。我们可以在netron看到

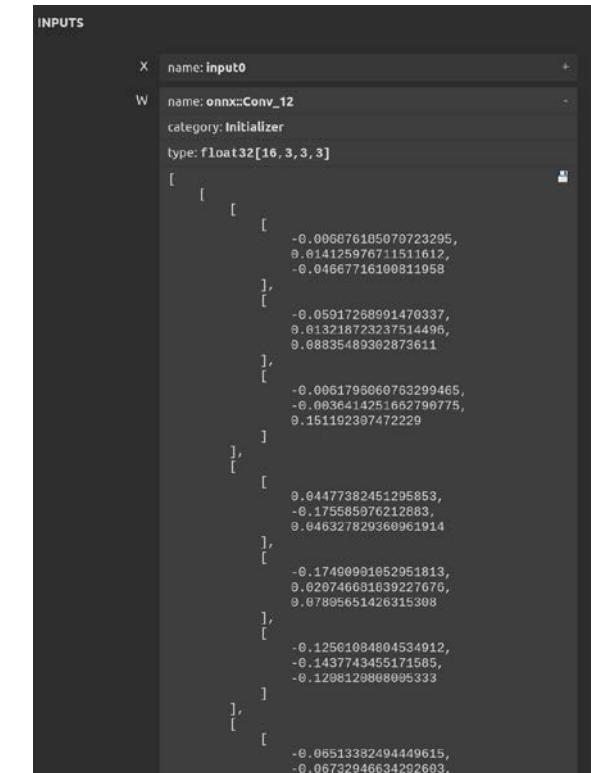
一般用来定义一个网络。包括

- input/output
- initializer^(*)
- node

(node是repeated，所以是数组)

(initializer是repeated，所以是数组)

(input/output是repeated，所以是数组)



理解onnx中的ModelProto

```
337 // Models
338 //
339 // ModelProto is a top-level file/container format for bundling a ML model and
340 // associating its computation graph with metadata.
341 //
342 // The semantics of the model are described by the associated GraphProto's.
343 message ModelProto {
344     // The version of the IR this model targets. See Version enum above.
345     // This field MUST be present.
346     optional int64 ir_version = 1;
347
348     // The OperatorSets this model relies on.
349     // All ModelProtos MUST have at least one entry that
350     // specifies which version of the ONNX OperatorSet is
351     // being imported.
352     //
353     // All nodes in the ModelProto's graph will bind against the operator
354     // with the same-domain/same-op_type operator with the HIGHEST version
355     // in the referenced operator sets.
356     repeated OperatorSetIdProto opset_import = 8;
357
358     // The name of the framework or tool used to generate this model.
359     // This field SHOULD be present to indicate which implementation/tool/framework
360     // emitted the model.
361     optional string producer_name = 2;
362
363     // The version of the framework or tool used to generate this model.
364     // This field SHOULD be present to indicate which implementation/tool/framework
365     // emitted the model.
366     optional string producer_version = 3;
367
368     // Domain name of the model.
369     // We use reverse domain names as name space indicators. For example:
370     // `com.facebook.fair` or `com.microsoft.cognitiveservices`
371     //
372     // Together with `model_version` and GraphProto.name, this forms the unique identity
373     // the graph.
374     optional string domain = 4;
375
376     // The version of the graph encoded. See Version enum below.
377     optional int64 model_version = 5;
378
379     // A human-readable documentation for this model. Markdown is allowed.
380     optional string doc_string = 6;
381
382     // The parameterized graph that is evaluated to execute the model.
```

```
383     optional GraphProto graph = 7;
384
385     // Named metadata values; keys should be distinct.
386     repeated StringStringEntryProto metadata_props = 14;
387
388     // Training-specific information. Sequentially executing all stored
389     // `TrainingInfoProto.algorithm`s and assigning their outputs following
390     // the corresponding `TrainingInfoProto.update_binding`s is one training
391     // iteration. Similarly, to initialize the model
392     // (as if training hasn't happened), the user should sequentially execute
393     // all stored `TrainingInfoProto.initialization`s and assigns their outputs
394     // using `TrainingInfoProto.initialization_binding`s.
395
396     // If this field is empty, the training behavior of the model is undefined.
397     repeated TrainingInfoProto training_info = 20;
398
399     // A list of function protos local to the model.
400     //
401     // Name of the function "FunctionProto.name" should be unique within the domain "Fun"
402     // In case of any conflicts the behavior (whether the model local functions are give
403     // or standard opserator sets are given higher priority or this is treated as error)
404     // the runtimes.
405
406     // The operator sets imported by FunctionProto should be compatible with the ones
407     // imported by ModelProto and other model local FunctionProtos.
408     // Example, if same operator set say 'A' is imported by a FunctionProto and ModelPro
409     // or by 2 FunctionProtos then versions for the operator set may be different but,
410     // the operator schema returned for op_type, domain, version combination
411     // for both the versions should be same for every node in the function body.
412
413     // One FunctionProto can reference other FunctionProto in the model, however, recurs
414     // is not allowed.
415     repeated FunctionProto functions = 25;
416 }
417
```

一般用来定义模型的全局信息，比如opset
(graph并不是repeated，所以一个model对应一个graph)

根据onnx中的Proto信息， 创建onnx

onnx提供了一些很方便的api来创建onnx

onnx.helper.make_tensor

onnx.helper.make_tensor_value_info

onnx.helper.make_attribute

onnx.helper.make_node

onnx.helper.make_graph

onnx.helper.make_model

```
##### 创建第一个conv节点 #####
conv1_output_name = "conv2d_1.output"
conv1_in_ch       = input_channel
conv1_out_ch      = 32
conv1_kernel      = 3
conv1_pads        = 1

# 创建conv节点的权重信息
conv1_weight      = np.random.rand(conv1_out_ch, conv1_in_ch, conv1_kernel, conv1_kernel)
conv1_bias        = np.random.rand(conv1_out_ch)

conv1_weight_name = "conv2d_1.weight"
conv1_weight_initializer = create_initializer_tensor(
    name          = conv1_weight_name,
    tensor_array  = conv1_weight,
    data_type     = onnx.TensorProto.FLOAT)

conv1_bias_name   = "conv2d_1.bias"
conv1_bias_initializer = create_initializer_tensor(
    name          = conv1_bias_name,
    tensor_array  = conv1_bias,
    data_type     = onnx.TensorProto.FLOAT)

# 创建conv节点，注意conv节点的输入有3个：input, w, b
conv1_node = onnx.helper.make_node(
    name          = "conv2d_1",
    op_type       = "Conv",
    inputs        = [
        model_input_name,
        conv1_weight_name,
        conv1_bias_name
    ],
    outputs       = [conv1_output_name],
    kernel_shape  = [conv1_kernel, conv1_kernel],
    pads          = [conv1_pads, conv1_pads, conv1_pads, conv1_pads],
```

根据onnx中的Proto信息， 读取onnx

只要知道Proto的数据结构， 我们就可以parse整个onnx， 或者获取特定点的信息

```
print("\n*****parse input/output*****")
for input in inputs:
    input_shape = []
    for d in input.type.tensor_type.shape.dim:
        if d.dim_value == 0:
            input_shape.append(None)
        else:
            input_shape.append(d.dim_value)
    print("Input info: \
          \n\tname: {} \
          \n\tdata Type: {} \
          \n\tshape: {}".format(input.name, input.type.tensor_type.elem_type, input_shape))

for output in outputs:
    output_shape = []
    for d in output.type.tensor_type.shape.dim:
        if d.dim_value == 0:
            output_shape.append(None)
        else:
            output_shape.append(d.dim_value)
    print("Output info: \
          \n\tname: {} \
          \n\tdata Type: {} \
          \n\tshape: {}".format(output.name, output.type.tensor_type.elem_type, output_shape))

print("\n*****parse node*****")
for node in nodes:
    print("node info: \
          \n\tname: {} \
          \n\top_type: {} \
          \n\tinputs: {} \
          \n\toutputs: {}".format(node.name, node.op_type, node.input, node.output))

print("\n*****parse initializer*****")
for initializer in initializers:
    print("initializer info: \
          \n\tname: {} \
          \n\tdata_type: {} \
          \n\tshape: {}".format(initializer.name, initializer.data_type, initializer.dims))
```

```
*****parse input/output*****
Input info:
  name:      input0
  Data Type: 1
  shape:     [1, 3, 64, 64]
Output info:
  name:      input0
  data Type: 1
  shape:     [1, 3, 64, 64]

*****parse node*****
node info:
  name:      conv2d_1
  op_type:   Conv
  inputs:    ['input0', 'conv2d_1.weight', 'conv2d_1.bias']
  outputs:   ['conv2d_1.output']

node info:
  name:      batchNorm1
  op_type:   BatchNormalization
  inputs:    ['conv2d_1.output', 'batchNorm1.scale', 'batchNorm1.bias',
             'batchNorm1.output']

node info:
  name:      relu1
  op_type:   Relu
  inputs:    ['batchNorm1.output']
  outputs:   ['relu1.output']

node info:
  name:      avg_pool1
  op_type:   GlobalAveragePool
  inputs:    ['relu1.output']
  outputs:   ['avg_pool1.output']

node info:
  name:      conv2d_2
  op_type:   Conv
  inputs:    ['avg_pool1.output', 'conv2d_2.weight', 'conv2d_2.bias']
  outputs:   ['output0']

*****parse initializer*****
```

根据onnx中的Proto信息，修改onnx

虽然onnx官方提供了而一些pythonAPI来修改onnx，但是我这里推荐大家使用TensorRT下的`onnxsurgeon`
详细请看后面的小节

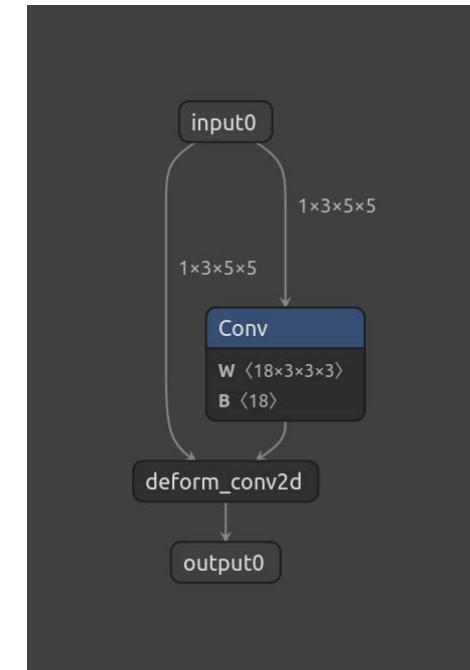
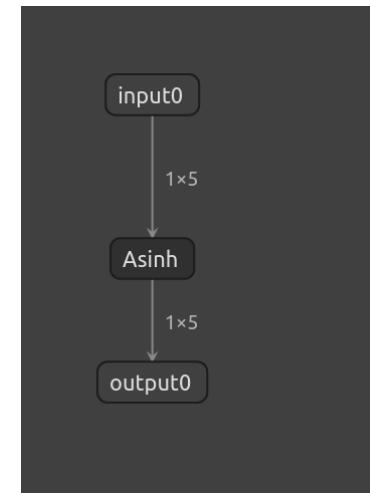
04

onnx注册算子的方法

Goal: 学习pytorch导出onnx不成功的时候如何解决(without plugin篇)

执行一下我们的python程序

```
../3.4-export-unsupported-node/
├── models
│   ├── sample-asinh.onnx
│   ├── sample-deformable-conv.onnx
│   ├── swin-tiny-after-simplify.onnx
│   └── swin-tiny-before-simplify.onnx
└── src
    ├── sample_asinh.py
    ├── sample_asinh_register.py
    ├── sample_deformable_conv.py
    └── sample_deformable_conv_register.py
```



```
===== Diagnostic Run torch.onnx.export version 2.0.1+cu117 =====
verbose: False, log level: Level.ERROR
===== 0 NONE 0 NOTE 0 WARNING 0 ERROR =====
Finished normal onnx export
result from Pytorch is : tensor([[0.6005, 0.6733, 0.5829, 0.2650, 0.4782]])
result from onnx is: [array([[0.6005171 , 0.67329633, 0.5829002 , 0.26500833, 0.47822666]], dtype=float32)]
```

3.4-export-unsupported-node

一起学习onnx中ProtoBuf的框架，学习onnx的细节

执行结果的一部分展示

转换swin-tiny时候出现的不兼容op的例子

```
symbolic_m = _find_symbolic_in_registry(domain, op_name, opset_version, operator_export_type)
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/utils.py", line 947, in _find_symbolic_in_registry
    return sym_registry.get_registered_op(op_name, domain, opset_version)
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/symbolic_registry.py", line 112, in get_registered_op
    raise RuntimeError(msg)
RuntimeError: Exporting the operator roll to ONNX opset version 9 is not supported. Please feel free to request support or submit a pull request on PyTorch GitHub.
[status] [C:\Users\2\Swin_Transformer (main)\1] [1]+ 11
```

先使用opset9看看是否可以导出

```
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/utils.py", line 947, in _find_symbolic_in_registry
    return sym_registry.get_registered_op(op_name, domain, opset_version)
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/symbolic_registry.py", line 112, in get_registered_op
    raise RuntimeError(msg)
RuntimeError: Exporting the operator roll to ONNX opset version 12 is not supported. Please feel free to request support or submit a pull request on PyTorch GitHub.
```

使用opset12依然会出现问题

```
# reverse cyclic shift
if self.shift_size > 0:
    if not self.fused_window_process:
        shifted_x = window_reverse(attn_windows, self.window_size, H, W) # B H' W' C
        x = torch.roll(shifted_x, shifts=(self.shift_size, self.shift_size), dims=(1, 2))
    else:
        x = WindowProcessReverse.apply(attn_windows, B, H, W, C, self.shift_size, self.window_size)
else:
    shifted_x = window_reverse(attn_windows, self.window_size, H, W) # B H' W' C
    x = shifted_x
x = x.view(B, H * W, C)
x = shortcut + self.drop_path(x)
```

(*)注意：因为onnx是一种图结构表示，并不包含各个算子的实现。除非我们是要在onnx-runtime上测试，否则我们更看重onnx-trt中这个算子的支持情况

当出现导出onnx不成功的时候，我们需要考虑的事情

难度从第到高：

- 修改opset的版本
 - 查看不支持的算子在新的opset中是否被支持
 - 如果不考虑自己搭建plugin的话，也需要看看onnx-trt中这个算子是否被支持(*)
 - 大家习惯看官方文档
- 替换pytorch中的算子组合
 - 把某些计算替换成onnx可以识别的
- 在pytorch登记onnx中某些算子
 - 有可能onnx中有支持，但没有被登记
- 直接修改onnx，创建plugin
 - 使用onnx-surgeon
 - 一般是在加速某些算子上使用

unsupported asinh算子

```
class Model(torch.nn.Module):
    def __init__(self):
        super().__init__()

    def forward(self, x):
        x = torch.asinh(x)
        return x
```

```
File "/home/kalfazed/miniconda3/envs/trt-starter/lib/python3.9/site-packages/torch/onnx/utils.py", line 1117, in _model_to
    graph = _optimize_graph(
File "/home/kalfazed/miniconda3/envs/trt-starter/lib/python3.9/site-packages/torch/onnx/utils.py", line 665, in _optimize_
    graph = _C._jit_pass_onnx(graph, operator_export_type)
File "/home/kalfazed/miniconda3/envs/trt-starter/lib/python3.9/site-packages/torch/onnx/utils.py", line 1901, in _run_symb
    raise errors.UnsupportedOperatorError(
torch.onnx.errors.UnsupportedOperatorError: Exporting the operator 'aten::asinh' to ONNX opset version 12 is not supported.
    ll request on PyTorch GitHub: https://github.com/pytorch/pytorch/issues.
[station]~/C/d/i/t/c/3/src (main|✚55) [1]$ ll
```

src/sample_asinh.py
(出现导出问题，我们先去寻找官方文档)

(*)opset: operator sets。可以理解为每个onnx版本所支持的算子集合。一般来说越大越多

unsupported asinh算子

Ops	Opset
Acosh	9
Add	14, 13, 7, 6, 1
And	7, 1
ArgMax	13, 12, 11, 1
ArgMin	13, 12, 11, 1
Asin	7
Asinh	9
Atan	7
Atanh	9
AveragePool	19, 11, 10, 7, 1
BatchNormalization	15, 14, 9, 7, 6, 1
BitShift	11
BitwiseAnd	18

这里我们可以看到asinh算子在opset9(*)这个版本已经开始支持了。
那我们换一下opset试一下

```
0 def export_norm_onnx():
1     input    = torch.rand(1, 5)
2     model   = Model()
3     model.eval()
4
5     file    = ".../models/sample-asinh.onnx"
6     torch.onnx.export(
7         model      = model,
8         args       = (input, ),
9         f          = file,
10        input_names = ["input0"],
11        output_names = ["output0"],
12        opset_version = 9)
13     print("Finished normal onnx export")
```

<https://github.com/onnx/onnx/blob/main/docs/Operators.md>

```
graph = _C._jit_pass_onnx(graph, operator_export_type)
File "/home/kalfazed/miniconda3/envs/trt-starter/lib/python3.9/site-packages/torch/onnx/utils.py", line 1901, in _run_sym
    raise errors.UnsupportedOperatorError(
torch.onnx.errors.UnsupportedOperatorError: Exporting the operator 'aten::asinh' to ONNX opset version 9 is not supported.
L request on PyTorch GitHub: https://github.com/pytorch/pytorch/issues.
```

unsupported asinh算子

(*)opset: operator sets。可以理解为每个onnx版本所支持的算子集合。一般来说越大越多

```
- 3.2k 22 6月 01:10 📁 __init__.py
x - 22 6月 01:10 📁 __pycache__
- 411 22 6月 01:10 📁 _constants.py
- 2.0k 22 6月 01:10 📁 _deprecation.py
- 1.0k 22 6月 01:10 📁 _experimental.py
- 1.7k 22 6月 01:10 📁 _exporter_states.py
- 3.2k 22 6月 01:10 📁 _globals.py
x - 22 6月 01:10 📁 _internal
- 3.3k 22 6月 01:10 📁 _onnx_supported_ops.py
- 12k 22 6月 01:10 📁 _type_utils.py
- 3.9k 22 6月 01:10 📁 errors.py
- 560 22 6月 01:10 📁 operators.py
- 11k 22 6月 01:10 📁 symbolic_caffe2.py
- 62k 22 6月 01:10 📁 symbolic_helper.py
- 2.1k 22 6月 01:10 📁 symbolic_opset7.py
- 15k 22 6月 01:10 📁 symbolic_opset8.py
- 231k 22 6月 01:10 📁 symbolic_opset9.py
- 29k 22 6月 01:10 📁 symbolic_opset10.py
- 56k 22 6月 01:10 📁 symbolic_opset11.py
- 16k 22 6月 01:10 📁 symbolic_opset12.py
- 31k 22 6月 01:10 📁 symbolic_opset13.py
- 3.4k 22 6月 01:10 📁 symbolic_opset14.py
- 2.9k 22 6月 01:10 📁 symbolic_opset15.py
- 3.9k 22 6月 01:10 📁 symbolic_opset16.py
- 1.4k 22 6月 01:10 📁 symbolic_opset17.py
- 1.7k 22 6月 01:10 📁 symbolic_opset18.py
- 80k 22 6月 01:10 📁 utils.py
- 70k 22 6月 01:10 📁 verification.py
~/m/e/t/l/p/s/t/onnx $ pwd
fazed/miniconda3/envs/trt-starter/lib/python3.9/site-packages/torch/onnx
```

```
3 "addcmul",
4 "addmm",
5 "alias",
6 "amax",
7 "amin",
8 "aminmax",
9 "arange",
0 "argmax",
1 "argmin",
2 "as_strided",
3 "as_tensor",
4 "asin"|,
5 "atan",
6 "baddbmm",
7 "batch_norm",
8 "bernoulli",
9 "bitwise_not",
0 "bitwise_or",
1 "bmm",
2 "broadcast_tensors",
3 "bucketize",
4 "cat",
5 "cdist",
6 "ceil"
```

torch/onnx/symbolic_opset9.py

(很遗憾，虽然onnx支持asinh,但是pytorch2onnx没有建立起桥梁)

unsupported asinh算子

```
310
317 @_onnx_symbolic("aten::shape_as_tensor")
318 @beartype.beartype
319 def _shape_as_tensor(g: jit_utils.GraphContext, input):
320     return g.op("Shape", input)
321
322
323 @_onnx_symbolic("aten::reshape_from_tensor")
324 @beartype.beartype
325 def _reshape_from_tensor(g: jit_utils.GraphContext, input, shape):
326     if isinstance(shape, list):
327         shape = g.op("Concat", *shape, axis_i=0)
328     return reshape(g, input, shape)
329
330
331 @_onnx_symbolic("aten::reshape")
332 @symbolic_helper.quantized_args(True)
333 @beartype.beartype
334 def reshape(g: jit_utils.GraphContext, self, shape):
335     return symbolic_helper._reshape_helper(g, self, shape)
336
337
338 @_onnx_symbolic("aten::reshape_as")
339 @symbolic_helper.quantized_args(True)
340 @beartype.beartype
341 def reshape_as(g: jit_utils.GraphContext, self, other):
342     shape = g.op("Shape", other)
343     return reshape(g, self, shape)
344
```

torch/onnx/symbolic_opset9.py

但我们可以从这里面知道pytorch中的各个op是如何登记的(register)

aten::xxx

- c++的一个namespace，pytorch的很多算子的底层都是在aten这个命名空间下进行以c++进行实现的

onnx_symbolic

- 负责绑定
- 绑定pytorch中的算子与aten命名空间下的算子的一一对应

unsupported asinh算子

```
5  
6 # 创建一个asinh算子的symbolic, 符号函数, 用来登记  
7 # 符号函数内部调用g.op, 为onnx计算图添加Asinh算子  
8 # g: 就是graph, 计算图  
9 # 也就是说, 在计算图中添加onnx算子  
10 # 由于我们已经知道Asinh在onnx是有实现的, 所以我们只要在g.op调用这个op的名字就好了  
11 # symbolic的参数需要与Pytorch的asinh接口函数的参数对齐  
12 #     def asinh(input: Tensor, *, out: Optional[Tensor]=None) -> Tensor: ...  
13 def asinh_symbolic(g, input, *, out=None):  
14 |     return g.op("Asinh", input)  
15  
16 # 在这里, 将asinh_symbolic这个符号函数, 与PyTorch的asinh算子绑定。也就是所谓的“注册算子”  
17 # asinh是在名为aten的一个c++命名空间下进行实现的  
18 register_custom_op_symbolic('aten::asinh', asinh_symbolic, 12)  
19  
20  
21 # 这里容易混淆的地方:  
22 # 1. register_op中的第一个参数是PyTorch中的算子名字: aten::asinh  
23 # 2. g.op中的第一个参数是onnx中的算子名字: Asinh  
24
```

那么我们就可以按照类似的方法去自己创建这个联系

- 创建symbolic符号函数, 从而创建一个onnx operator
- 注册这个onnx operator, 并让它与底层的aten中的asinh实现绑定

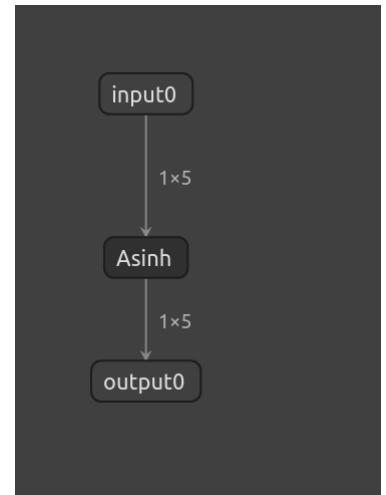
(注意: symbolic符号函数中的参数, 需要严格遵循pytorch中的定义)

```
571 def asin_(input: Tensor, *, out: Optional[Tensor]=None) -> Tensor: ...  
572 def asin_(input: Tensor) -> Tensor: ...  
573 def asinh(input: Tensor, *, out: Optional[Tensor]=None) -> Tensor: ...
```

unsupported asinh算子

```
[station]~/C/d/i/t/c/3/src (main|+555) $ python sample_asinh_register.py
=====
 Diagnostic Run torch.onnx.export version 2.0.1+cu117 =====
verbose: False, log level: Level.ERROR
=====
 0 NONE 0 NOTE 0 WARNING 0 ERROR =====

Finished normal onnx export
result from Pytorch is : tensor([[0.3245, 0.1989, 0.4941, 0.5102, 0.2483]])
result from onnx is:      [array([[0.3245445 , 0.19885206, 0.4941113 , 0.51015097, 0.2483138 ]],  
  dtype=float32)]
```



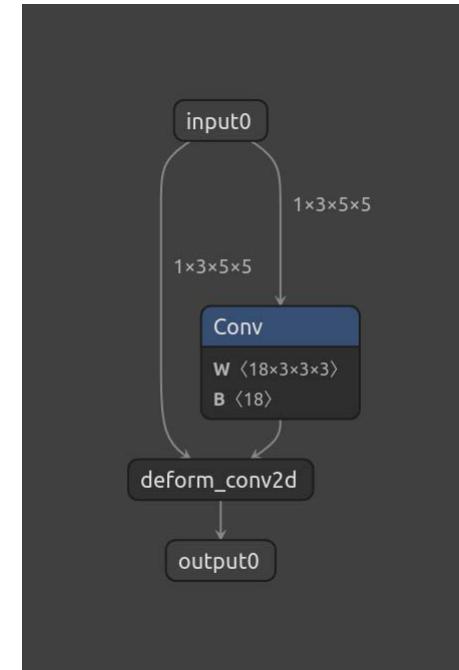
成功导出asinh，并可以使用onnxruntime进行验证

```
33
34 def validate_onnx():
35     input = torch.rand(1, 5)
36
37     # PyTorch的推理
38     model = Model()
39     x = model(input)
40     print("result from Pytorch is :", x)
41
42     # onnxruntime的推理
43     sess = onnxruntime.InferenceSession('../models/sample-asinh.onnx')
44     x = sess.run(None, {'input0': input.numpy()})
45     print("result from onnx is:    ", x)
46
```

onnxruntime的验证。很简单

unsupported deformable conv算子

```
8 # 注意
9 #   这里需要把args的各个参数的类型都指定
0 #   这里还没有实现底层对deform_conv2d的实现
1 #   具体dcn的底层实现是在c++完成的，这里会在后面的TensorRT plugin中回到这里继续讲这个案例
2 #   这里先知道对于不支持的算子，onnx如何导出即可
3 # @parse_args("v", "v", "v", "v", "v", "i", "none")
4 def dcn_symbolic(
5     g,
6     input,
7     weight,
8     offset,
9     mask,
0     bias,
1     stride_h, stride_w,
2     pad_h, pad_w,
3     dil_h, dil_w,
4     n_weight_grps,
5     n_offset_grps,
6     use_mask):
7     return g.op("custom::deform_conv2d", input, offset)
8
9 register_custom_op_symbolic("torchvision::deform_conv2d", dcn_symbolic, 12)
1
2
3 class Model(torch.nn.Module):
4     def __init__(self):
5         super().__init__()
6         self.conv1 = nn.Conv2d(3, 18, 3)
7         self.conv2 = torchvision.ops.DeformConv2d(3, 3, 3)
8
9     def forward(self, x):
0         x = self.conv2(x, self.conv1(x))
1         return x
```





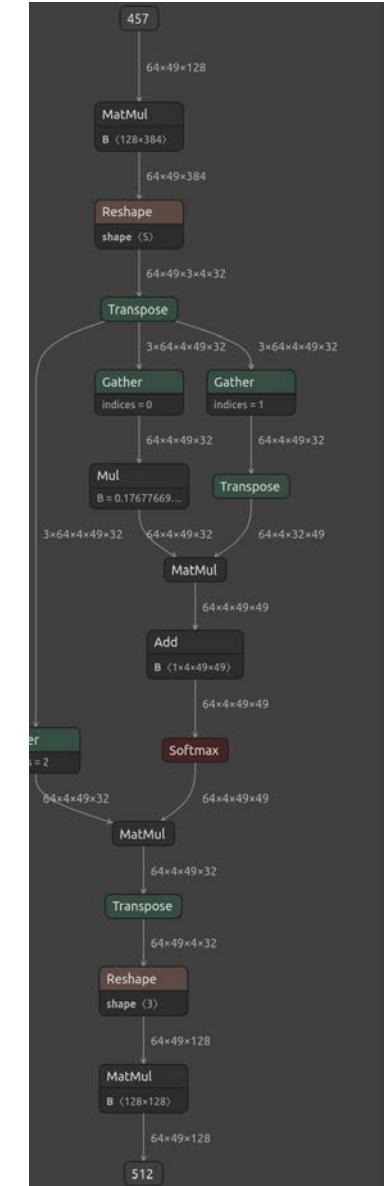
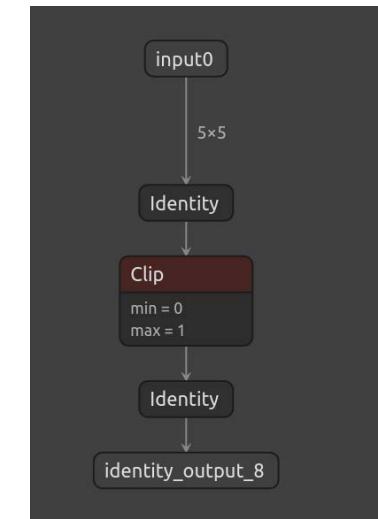
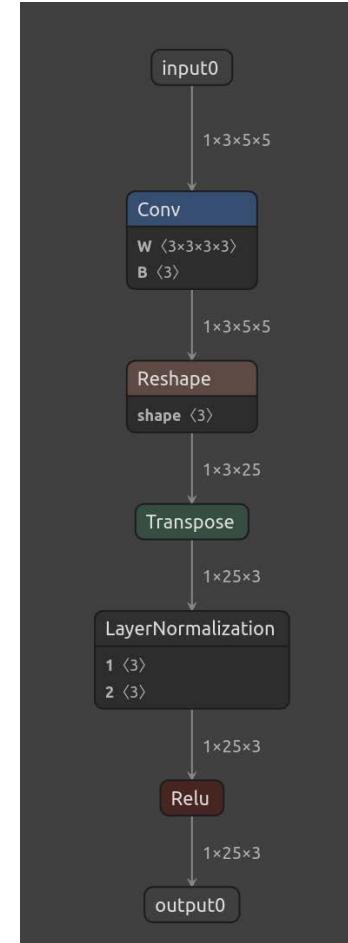
04

onnx-graph-surgeon

Goal: 学习使用onnx-surgeon，比较与onnx.helper的区别，学习快速修改onnx以及替换算子 /创建算子的技巧

执行一下我们的python程序

```
3.5-onnxsurgeon/
  models
    sample-complicated-graph.onnx
    sample-conv.onnx
    sample-ln-after.onnx
    sample-ln-before.onnx
    sample-minmax.onnx
    sample-minmax-to-clip.onnx
    swin-subgraph-LN.onnx
    swin-subgraph-MSHA.onnx
    swin-tiny.onnx
  src
    gs_create_complicate_graph.py
    gs_create_conv.py
    gs_create_subgraph.py
    gs_replace_LN.py
    gs_replace_op.py
```



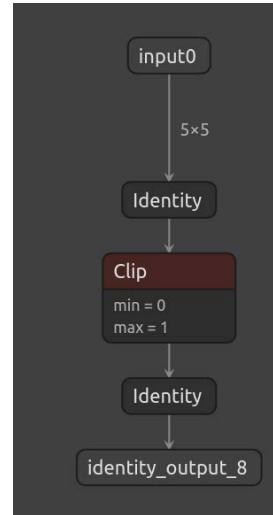
执行结果的一部分展示

3.5-onnxsurgeon
学习使用onnxsurgeon来创建/修改onnx

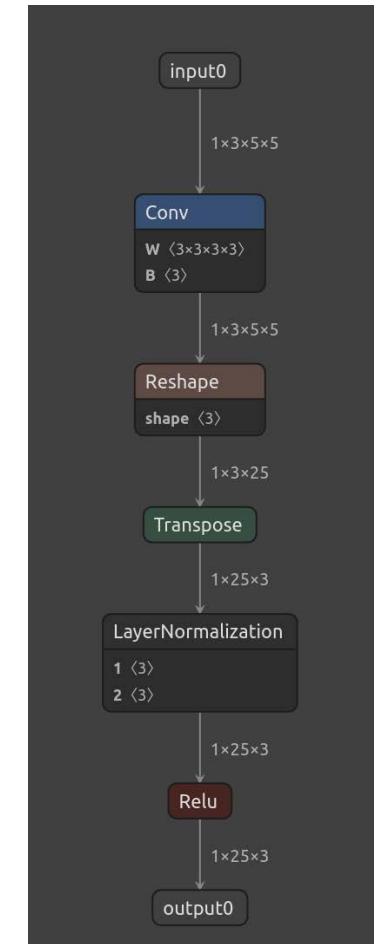
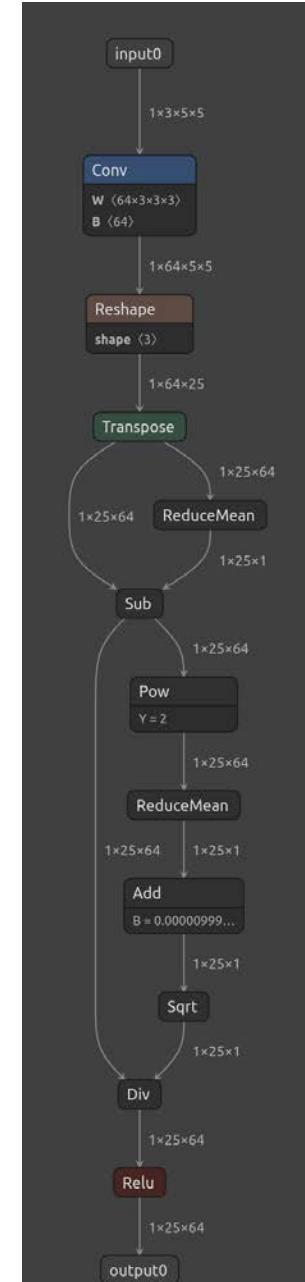
onnx-graph-surgeon

创建/修改onnx的工具。在TensorRT/tools中可以安装

- 更加方便的添加/修改onnx节点
- 更加方便的修改子图
- 更加方便的替换算子
- (底层一般是用的onnx.helper，但是给做了一些封装)



将min-max替换成clip算子
(clip算子在onnx中有支持，所以直接替换)



将一系列复杂的LayerNorm计算替换成一个LN算子
(TensorRT 8.6以前的做法)

onnx-surgeon vs onnx.hepler

```
4  
5 # 理解onnx中的组织结构  
6 #   - ModelProto (描述的是整个模型的信息)  
7 #   --- GraphProto (描述的是整个网络的信息)  
8 #   ----- NodeProto (描述的是各个计算节点, 比如conv, linear)  
9 #   ----- TensorProto (描述的是tensor的信息, 主要包括权重)  
10 #  ----- ValueInfoProto (描述的是input/output信息)  
11 #  ----- AttributeProto (描述的是node节点的各种属性信息)  
12
```

onnx ProtoBuf中的IR表示

```
4  
5 # onnx_graph_surgeon(gs)中的IR会有以下三种结构  
6 # Tensor  
7 #   -- 有两种类型  
8 #     -- Variable: 主要就是那些不到推理不知道的变量  
9 #     -- Constant: 不用推理时, 而在推理前就知道的变量  
10 # Node  
11 #   -- 跟onnx中的NodeProto差不多  
12 # Graph  
13 #   -- 跟onnx中的GraphProto差不多  
14
```

- **gs**帮助我们隐藏了很多信息
- **node**的属性以前使用**AttributeProto**保存,但是**gs**中统一用**dict**来保存

onnx graph surgeon中的IR表示

onnx-surgeon vs onnx.hepler

```
6 def create_initializer_tensor(
7     name: str,
8     tensor_array: np.ndarray,
9     data_type: onnx.TensorProto = onnx.TensorProto.FLOAT
10) -> onnx.TensorProto:
11
12     initializer = onnx.helper.make_tensor(
13         name = name,
14         data_type = data_type,
15         dims = tensor_array.shape,
16         vals = tensor_array.flatten().tolist())
17
18     return initializer
19
#####
20 conv1_output_name = "conv2d_1.output"
21 conv1_in_ch = input_channel
22 conv1_out_ch = 32
23 conv1_kernel = 3
24 conv1_pads = 1
25
26 # 创建conv节点的权重信息
27 conv1_weight = np.random.rand(conv1_out_ch, conv1_in_ch, conv1_kernel, conv1_kernel)
28 conv1_bias = np.random.rand(conv1_out_ch)
29
30 conv1_weight_name = "conv2d_1.weight"
31 conv1_weight_initializer = create_initializer_tensor(
32     name = conv1_weight_name,
33     tensor_array = conv1_weight,
34     data_type = onnx.TensorProto.FLOAT)
35
36 conv1_bias_name = "conv2d_1.bias"
37 conv1_bias_initializer = create_initializer_tensor(
38     name = conv1_bias_name,
39     tensor_array = conv1_bias,
40     data_type = onnx.TensorProto.FLOAT)
41
42 # 创建conv节点, 注意conv节点的输入有3个: input, w, b
43 conv1_node = onnx.helper.make_node(
44     name = "conv2d_1",
45     op_type = "Conv",
46     inputs = [
47         model_input_name,
48         conv1_weight_name,
49         conv1_bias_name
50     ],
51     outputs = [conv1_output_name],
52     kernel_shape = [conv1_kernel, conv1_kernel],
53     pads = [conv1_pads, conv1_pads, conv1_pads, conv1_pads],
54 )
```

使用原生的onnx.hepler创建onnx

```
input = gs.Variable(
    name = "input0",
    dtype = np.float32,
    shape = (1, 3, 224, 224))
weight = gs.Constant(
    name = "conv1.weight",
    values = np.random.randn(5, 3, 3, 3))
bias = gs.Constant(
    name = "conv1.bias",
    values = np.random.randn(5))
output = gs.Variable(
    name = "output0",
    dtype = np.float32,
    shape = (1, 5, 224, 224))
node = gs.Node(
    op = "Conv",
    inputs = [input, weight, bias],
    outputs = [output],
    attrs = {"pads": [1, 1, 1, 1]})
graph = gs.Graph(
    nodes = [node],
    inputs = [input],
    outputs = [output])
model = gs.export_onnx(graph)
```

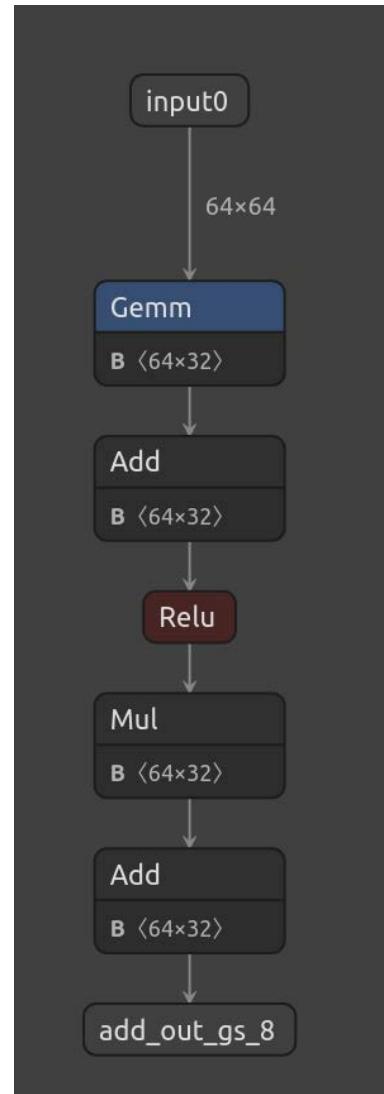
使用原生的gs创建onnx

onnx-surgeon vs onnx.hepler

gs可以自定义一些函数去创建onnx，使整个onnx的创建更加方便

```
5 #####在graph注册调用的函数#####
6 @gs.Graph.register()
7 def add(self, a, b):
8     return self.layer(op="Add", inputs=[a, b], outputs=["add_out_gs"])
9
10 @gs.Graph.register()
11 def mul(self, a, b):
12     return self.layer(op="Mul", inputs=[a, b], outputs=["mul_out_gs"])
13
14 @gs.Graph.register()
15 def gemm(self, a, b, trans_a=False, trans_b=False):
16     attrs = {"transA": int(trans_a), "transB": int(trans_b)}
17     return self.layer(op="Gemm", inputs=[a, b], outputs=["gemm_out_gs"], attrs=attrs)
18
19 @gs.Graph.register()
20 def relu(self, a):
21     return self.layer(op="Relu", inputs=[a], outputs=["act_out_gs"])
22
```

```
47 # 设计网络架构
48 gemm0      = graph.gemm(input0, consA, trans_b=True)
49 relu0      = graph.relu(*graph.add(*gemm0, consB))
50 mul0       = graph.mul(*relu0, consC)
51 output0    = graph.add(*mul0, consD)
52
```

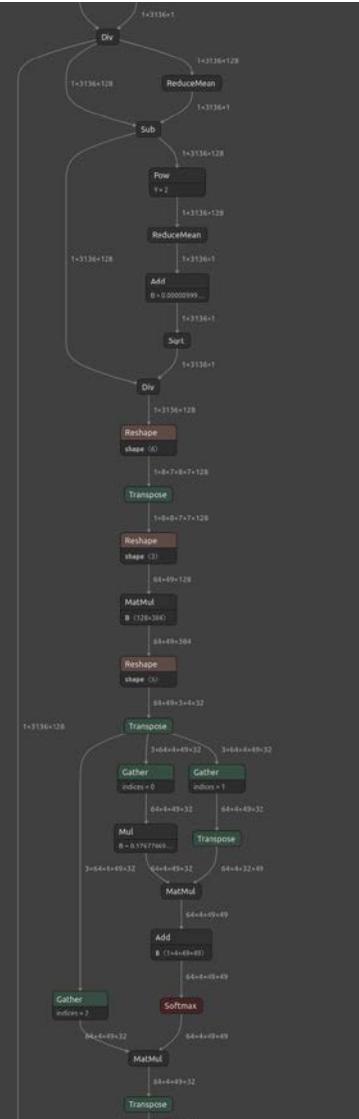


类似于onnx中symbolic符号函数来注册算子一样
我们完全可以自己创建一些算子在这里使用

(*)polygraphy是TensorRT用来分析模型的一个非常非常重要的一个软件，会在后面的章节详细展开讲

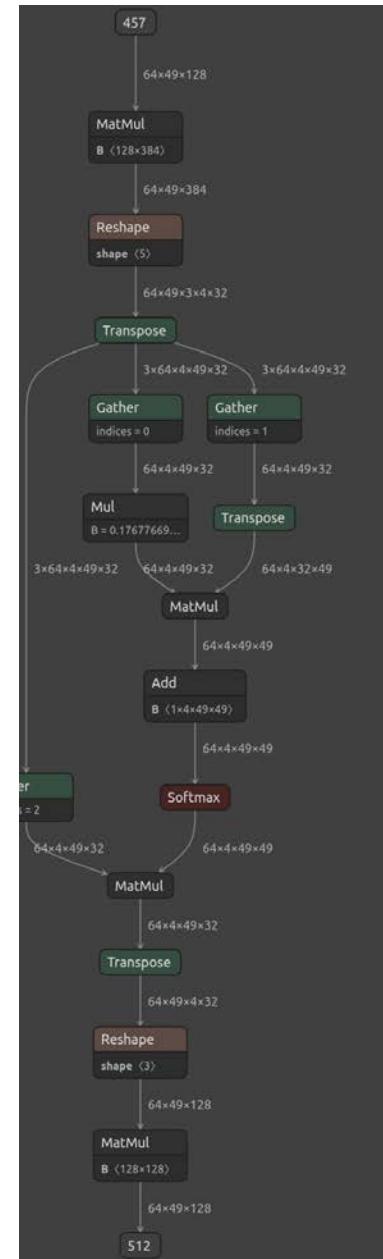
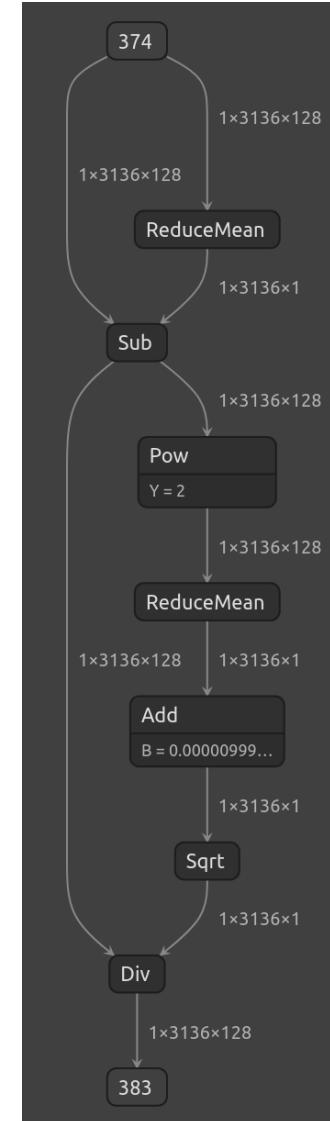
onnx-surgeon vs onnx.hepler

gs可以方便我们把整个网络中的一些子图给“挖”出来，以此来分析细节
(一般配合polygraphy^(*)使用，去寻找量化掉精度严重的子图)



```
# LayerNorm部分
print(tensors["374"]) # LN的input1: 1 x 3136 x 128
print(tensors["375"]) # LN的input2: 1 x 3136 x 1
print(tensors["383"]) # LN的输出: 1 x 3136 x 128
graph.inputs = [
    tensors["374"].to_variable(dtype=np.float32, shape=(1, 3136, 128))]
graph.outputs = [
    tensors["383"].to_variable(dtype=np.float32, shape=(1, 3136, 128))]
graph.cleanup()
onnx.save(gs.export_onnx(graph), "../models/swin-subgraph-LN.onnx")
```

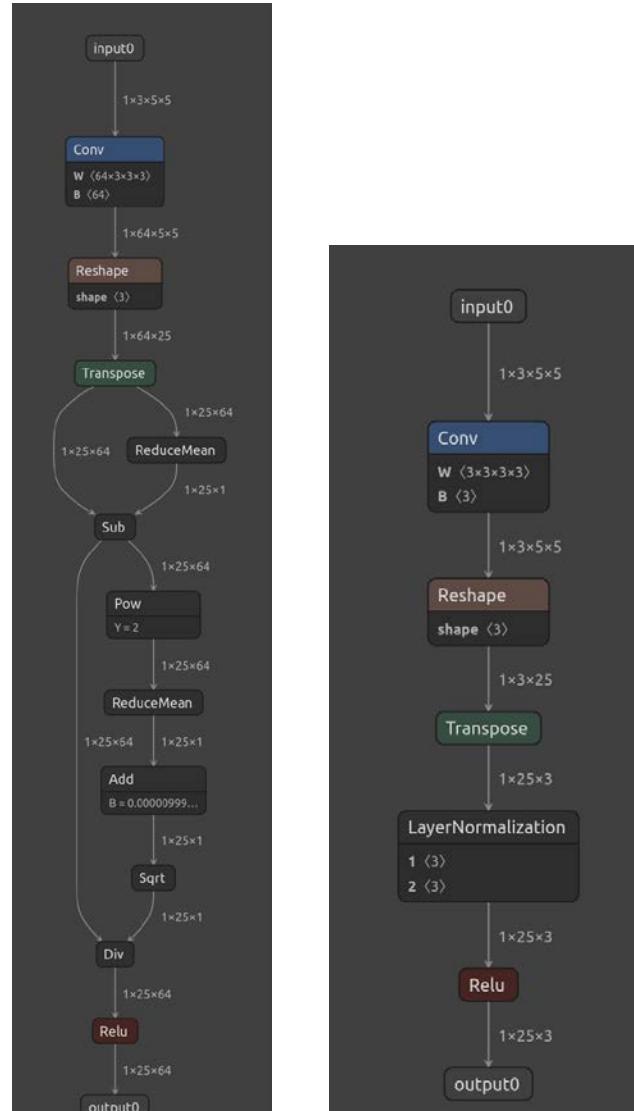
```
# MHSA部分
graph = gs.import_onnx(model)
tensors = graph.tensors()
print(tensors["457"]) # MHSA输入matmul: 64 x 49 x 128
print(tensors["5509"]) # MHSA输入matmul的权重: 128 x 384
print(tensors["5518"]) # MHSA输出matmul的权重: 128 x 128
print(tensors["512"]) # MHSA输出: 64 x 49 x 128
graph.inputs = [
    tensors["457"].to_variable(dtype=np.float32, shape=(64, 49, 128))]
graph.outputs = [
    tensors["512"].to_variable(dtype=np.float32, shape=(64, 49, 128))]
graph.cleanup()
onnx.save(gs.export_onnx(graph), "../models/swin-subgraph-MSHA.onnx")
```



onnx-surgeon vs onnx.hepler

gs中最重要的一个特点，在于我们可以使用gs来替换算子或者创建算子
(这个会直接跟后面的TensorRT plugin绑定，实现算子的加速或者不兼容算子的实现)

```
15 @gs.Graph.register()
16 def layerNorm(self, inputs, outputs, axis, epsilon):
17     attrs = {'axis': np.int64(axis), 'epsilon': np.float(epsilon)}
18     return self.layer(op="LayerNormalization", inputs=inputs, outputs=outputs, attrs=attrs)
19
20 @gs.Graph.register()
21 def layerNorm_default(self, inputs, outputs):
22     return self.layer(op="LayerNormalization", inputs=inputs, outputs=outputs)
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79 def change_onnx_graph():
80     graph = gs.import_onnx(onnx.load_model('../models/sample-ln-before.onnx'))
81     tensors = graph.tensors()
82
83     norm_scale = gs.Constant(name="norm.weight", values=np.ones(shape=[3], dtype=np.float32))
84     norm_bias = gs.Constant(name="norm.bias", values=np.zeros(shape=[3], dtype=np.float32))
85
86     inputs = [tensors["/Transpose_output_0"]]
87     outputs = [tensors["/norm/Div_output_0"]]
88
89     # 因为要替换子网，所以需要把子网和周围的所有节点都断开联系
90     for item in inputs:
91         item.outputs.clear()
92
93     for item in outputs:
94         item.inputs.clear()
95
96     # 为了迎合onnx中operator中的设计，这里把scale和bias给加上
97     inputs = [tensors["/Transpose_output_0"],
98               norm_scale,
99               norm_bias]
100
101     # 这个onnx中的epsilon，我们给加上。当然，我们也可以选择默认的值
102     epsilon = [tensors["/norm/Constant_1_output_0"]]
103     print(type(epsilon[0].values))
104
105     # 通过注册的LayerNorm，重新把断开的联系链接起来
106     graph.layerNorm(inputs, outputs, axis=-1, epsilon=epsilon[0].values)
107
108     # 删除所有额外的节点
109     graph.cleanup()
110
111     onnx.save(gs.export_onnx(graph), "../models/sample-ln-after.onnx")
```



onnx修改前(左), onnx修改后(右)



04

快速分析开源代码 并导出onnx

Goal: 以Swin Transformer为例学习快速导出onnx并分析onnx的方法

执行一下我们的python程序

```
C/d/i/t/chapter3-tensorrt-basics-and-onnx (main|+555) $ cd 3.6-export-onnx-from-oss/
- 23 7月 17:56 ↵models
- 16 7月 13:50 ↵Swin-Transformer
- 16 7月 12:54 ↵weights
C/d/i/t/c/3.6-export-onnx-from-oss (main|+555) $ tree models/
    .onnx-after-simplify.onnx
    .onnx-after-simplify-opset17.onnx
    .onnx-before-simplify.onnx
    .onnx
    .
    4 files
C/d/i/t/c/3.6-export-onnx-from-oss (main|+555) $ tree weights/
    .onnx_patch4_window7_224.pth
```



3.6-export-onnx-from-oss
(快速分析开源代码)

执行结果的一部分展示

转换swin-tiny时候出现的不兼容op的解决方案

```
symbolic_m = _find_symbolic_in_registry(domain, op_name, opset_version, operator_export_type)
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/utils.py", line 947, in _find_symbolic_in_registry
    return sym_registry.get_registered_op(op_name, domain, opset_version)
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/symbolic_registry.py", line 112, in get_registered_op
    raise RuntimeError(msg)
RuntimeError: Exporting the operator roll to ONNX opset version 9 is not supported. Please feel free to request support or submit a pull request on PyTorch GitHub.
[status] [C:\Users\2\Swin_Transformer (main)\M] [4]: 11
```

先使用opset9看看是否可以导出

```
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/utils.py", line 947, in _find_symbolic_in_registry
    return sym_registry.get_registered_op(op_name, domain, opset_version)
File "/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx/symbolic_registry.py", line 112, in get_registered_op
    raise RuntimeError(msg)
RuntimeError: Exporting the operator roll to ONNX opset version 12 is not supported. Please feel free to request support or submit a pull request on PyTorch GitHub.
```

使用opset12依然会出现问题

```
# reverse cyclic shift
if self.shift_size > 0:
    if not self.fused_window_process:
        shifted_x = window_reverse(attn_windows, self.window_size, H, W) # B H' W' C
        x = torch.roll(shifted_x, shifts=(self.shift_size, self.shift_size), dims=(1, 2))
    else:
        x = WindowProcessReverse.apply(attn_windows, B, H, W, C, self.shift_size, self.window_size)
else:
    shifted_x = window_reverse(attn_windows, self.window_size, H, W) # B H' W' C
    x = shifted_x
x = x.view(B, H * W, C)
x = shortcut + self.drop_path(x)
```

转换swin-tiny时候出现的不兼容op的解决方案

TORCH.ROLL

`torch.roll(input, shifts, dims=None) → Tensor`

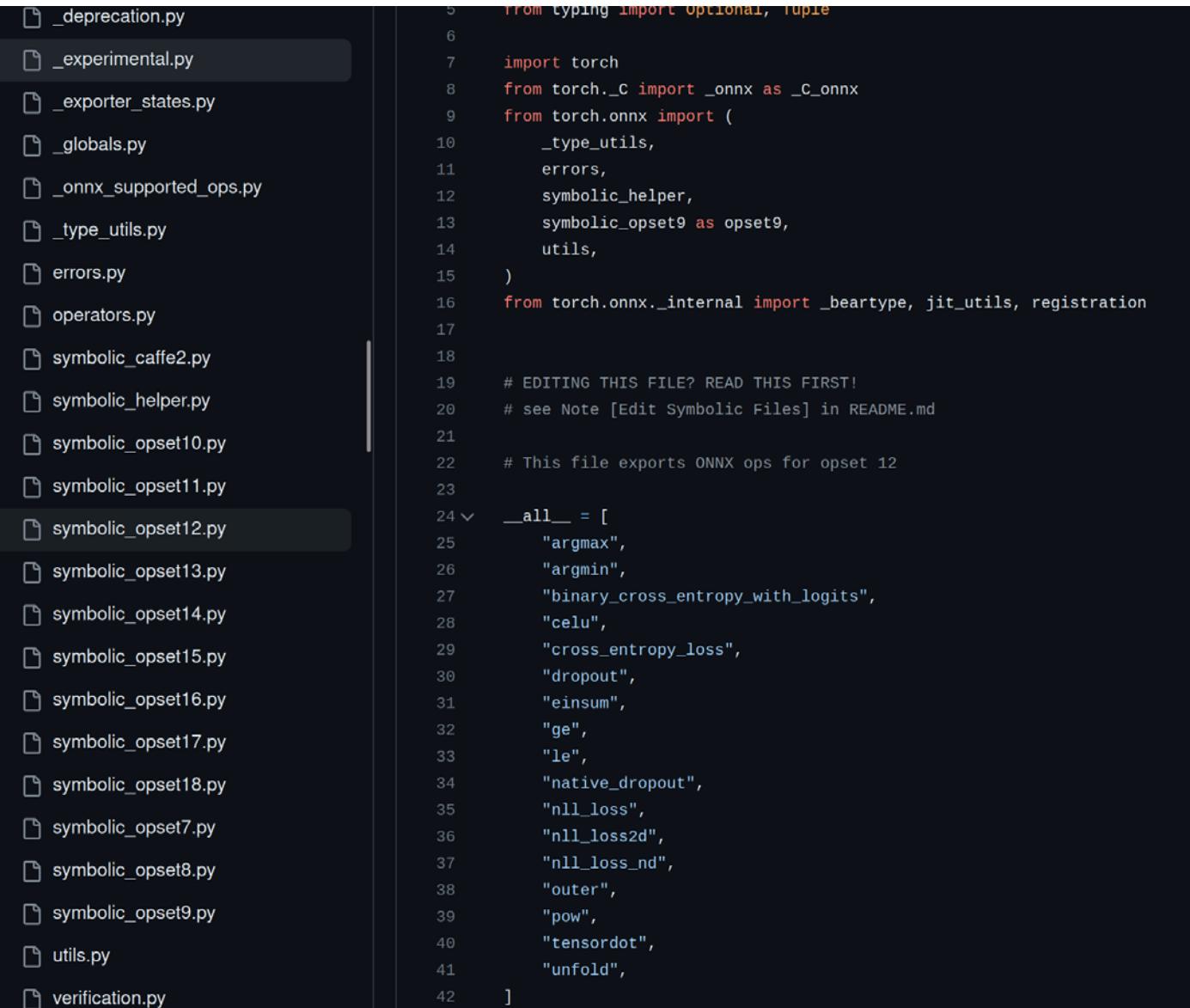
Roll the tensor `input` along the given dimension(s). Elements that are shifted beyond the last position are re-introduced at the first position. If `dims` is `None`, the tensor will be flattened before rolling and then restored to the original shape.

Parameters:

- `input` (`Tensor`) – the input tensor.
- `shifts` (`int` or `tuple of ints`) – The number of places by which the elements of the tensor are shifted. If `shifts` is a tuple, `dims` must be a tuple of the same size, and each dimension will be rolled by the corresponding value
- `dims` (`int` or `tuple of ints`) – Axis along which to roll

RandomUniform	1
RandomUniformLike	1
Reciprocal	13, 6, 1
ReduceMax	18, 13, 12, 11, 1
ReduceMean	18, 13, 11, 1
ReduceMin	18, 13, 12, 11, 1
ReduceProd	18, 13, 11, 1
ReduceSum	13, 11, 1
Reshape	19, 14, 13, 5, 1
Resize	19, 18, 13, 11, 10
ReverseSequence	10
RoiAlign	16, 10
Round	11
STFT	17
Scan	19, 16, 11, 9, 8
Scatter (deprecated)	11, 9
ScatterElements	18, 16, 13, 11
ScatterND	18, 16, 13, 11
ScatterUpdate	11

转换swin-tiny时候出现的不兼容op的解决方案



```
5     from typing import Optional, Tuple
6
7     import torch
8     from torch._C import _onnx as _C_onnx
9     from torch.onnx import (
10         _type_utils,
11         errors,
12         symbolic_helper,
13         symbolic_opset9 as opset9,
14         utils,
15     )
16     from torch.onnx._internal import _beartype, jit_utils, registration
17
18
19     # EDITING THIS FILE? READ THIS FIRST!
20     # see Note [Edit Symbolic Files] in README.md
21
22     # This file exports ONNX ops for opset 12
23
24     __all__ = [
25         "argmax",
26         "argmin",
27         "binary_cross_entropy_with_logits",
28         "celu",
29         "cross_entropy_loss",
30         "dropout",
31         "einsum",
32         "ge",
33         "le",
34         "native_dropout",
35         "nll_loss",
36         "nll_loss2d",
37         "nll_loss_nd",
38         "outer",
39         "pow",
40         "tensordot",
41         "unfold",
42     ]
```

https://github.com/pytorch/pytorch/blob/main/torch/onnx/symbolic_opset12.py

转换swin-tiny时候出现的不兼容op的解决方案

```
cd -> directory, it does not exist
[station]~/m/e/s/l/p/s/torch [1]$ cd onnx/
.rw-rw-r-- 16k 27 2月 2021 __init__.py
drwxrwxr-x - 16 7月 12:58 __pycache__
.rw-rw-r-- 579 27 2月 2021 operators.py
.rw-rw-r-- 9.6k 27 2月 2021 symbolic_caffe2.py
.rw-rw-r-- 35k 27 2月 2021 symbolic_helper.py
.rw-rw-r-- 1.8k 27 2月 2021 symbolic_opset7.py
.rw-rw-r-- 11k 27 2月 2021 symbolic_opset8.py
.rw-rw-r-- 120k 27 2月 2021 symbolic_opset9.py
.rw-rw-r-- 12k 27 2月 2021 symbolic_opset10.py
.rw-rw-r-- 40k 27 2月 2021 symbolic_opset11.py
.rw-rw-r-- 7.7k 27 2月 2021 symbolic_opset12.py
.rw-rw-r-- 8.0k 27 2月 2021 symbolic_opset13.py
.rw-rw-r-- 4.9k 27 2月 2021 symbolic_registry.py
.rw-rw-r-- 60k 27 2月 2021 utils.py
[station]~/m/e/s/l/p/s/t/onnx $ pwd
/home/kalfazed/miniconda3/envs/swin/lib/python3.7/site-packages/torch/onnx
[station]~/m/e/s/l/p/s/t/onnx $
```

因为是opset的问题，所以我们在这里面找

```
6 from sys import maxsize as maxsize
7
185 @parse_args('v', 'is', 'is')
186 def roll(g, self, shifts, dims):
187     assert len(shifts) == len(dims)
188
189     result = self
190     for i in range(len(shifts)):
191         shapes = []
192         shape = sym_help._slice_helper(g,
193                                         result,
194                                         axes=[dims[i]],
195                                         starts=[-shifts[i]],
196                                         ends=[maxsize])
197         shapes.append(shape)
198         shape = sym_help._slice_helper(g,
199                                         result,
200                                         axes=[dims[i]],
201                                         starts=[0],
202                                         ends=[-shifts[i]])
203         shapes.append(shape)
204         result = g.op("Concat", *shapes, axis_i=dims[i])
205
206     return result
```

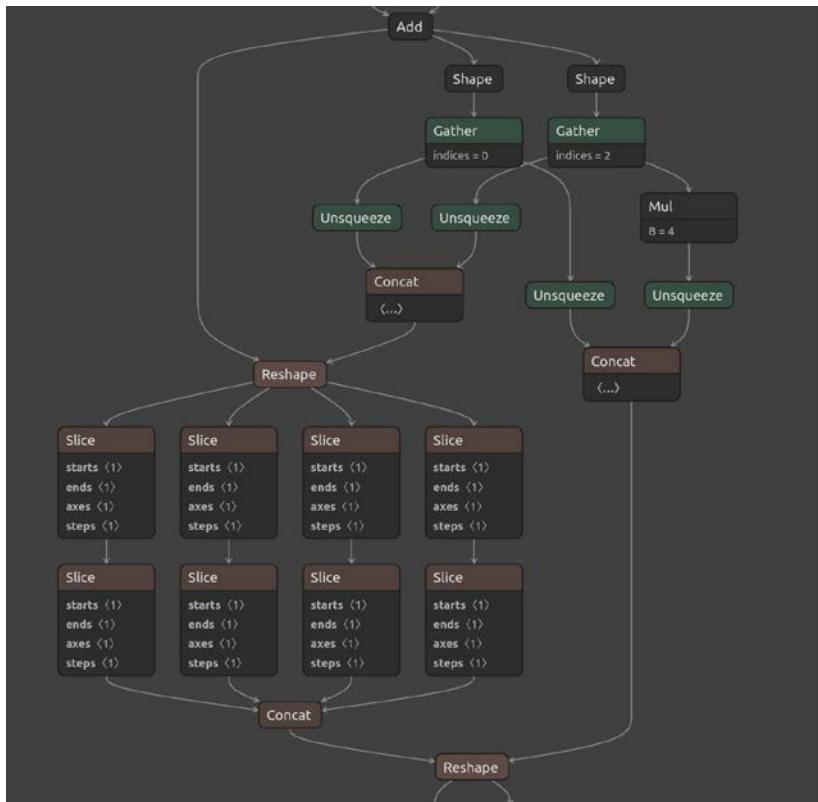
添加pytorch2onnx的roll算子实现

swin-tiny的onnx分析

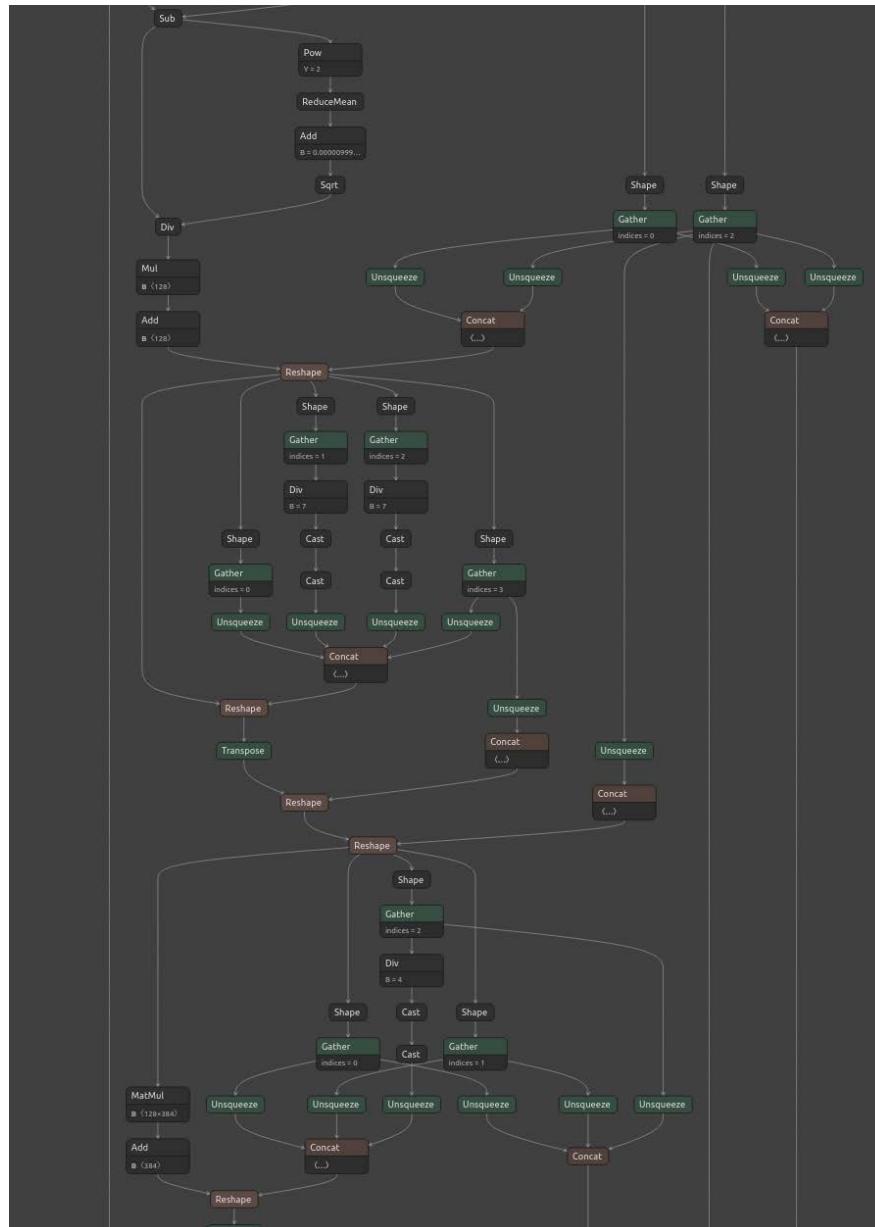
Finished normal onnx export
[station]~/C/d/i/t/c/3/Swin-Transformer (

成功通过

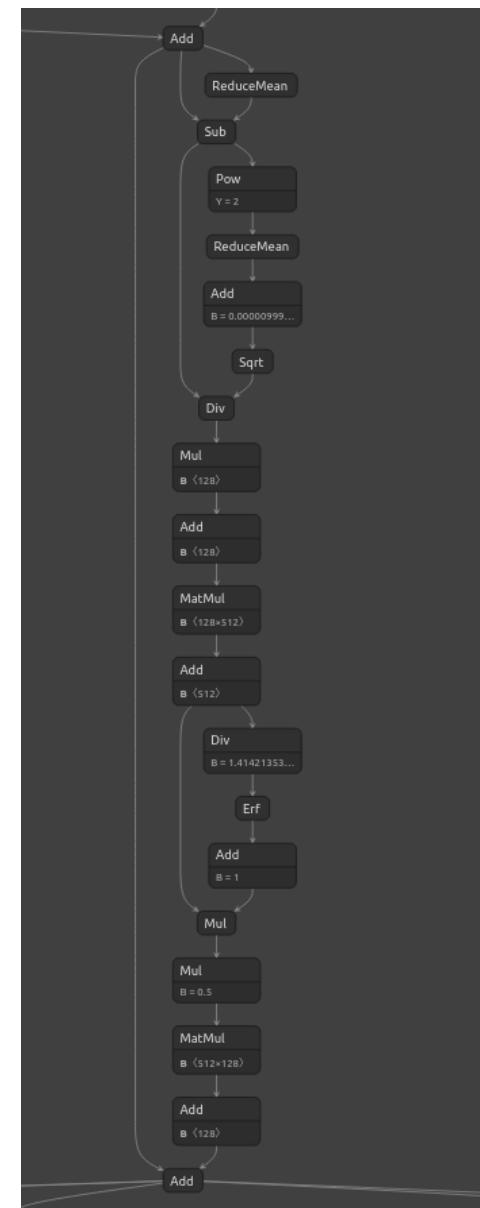
opset9即便把Roll的问题解决还会出现其他算子不兼容问题，所以直接在opset12上进行解决



swin的window shift部分

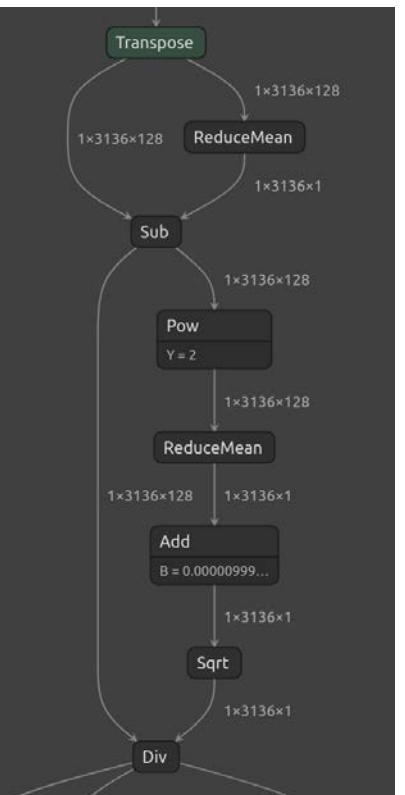
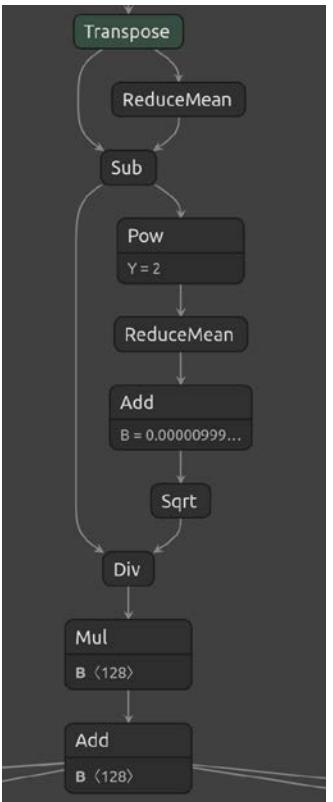


swin的attention部分

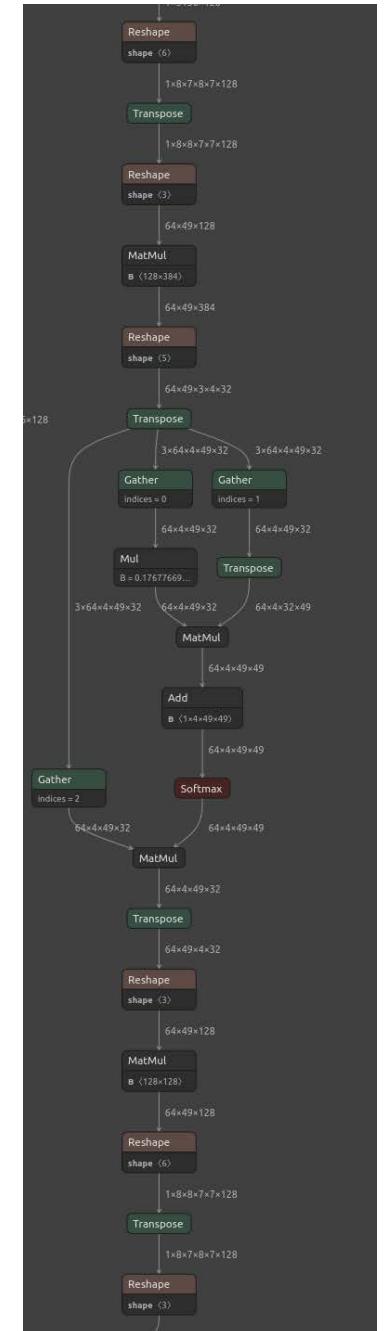
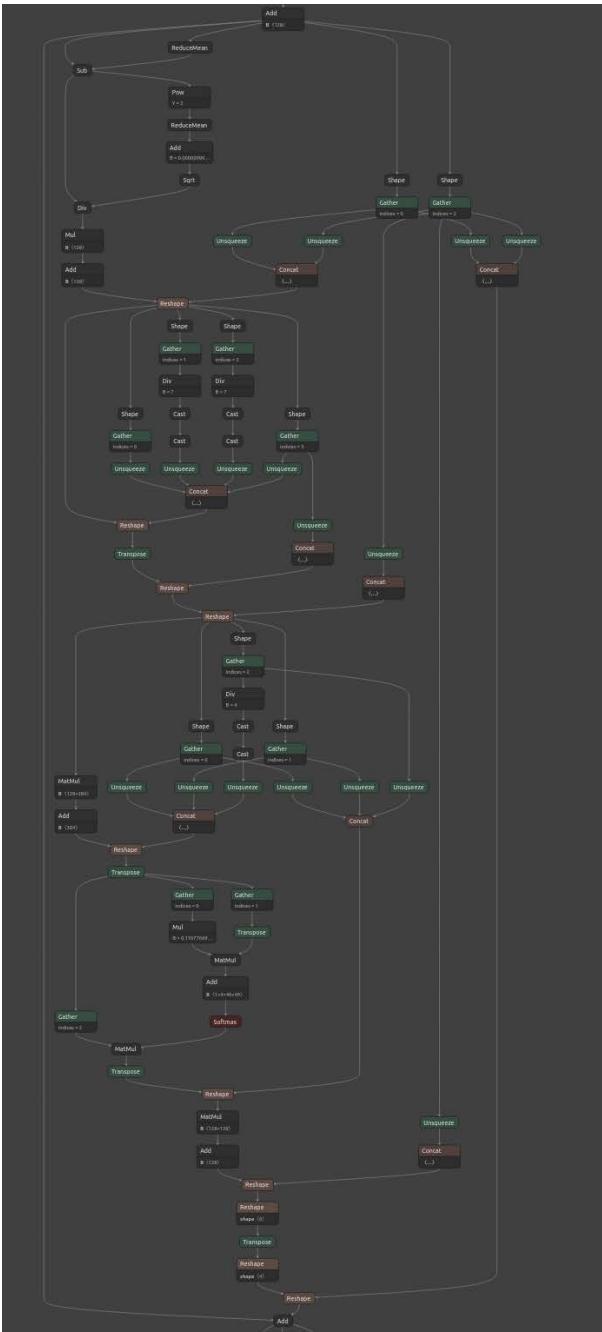


swin的LN部分

swin-tiny的onnx分析



onnxsim对LN的简化



onnxsim对attention的简化

根据onnx中的Proto信息，修改onnx

LN和BN的计算其实是差不多的

- 两者都是在做normalization，只不过是在不同的维度上

step1: mean和std的计算

$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

step2: normalization

$$\hat{a}^l = \frac{a^l - \mu^l}{\sqrt{(\sigma^l)^2 + \epsilon}}$$

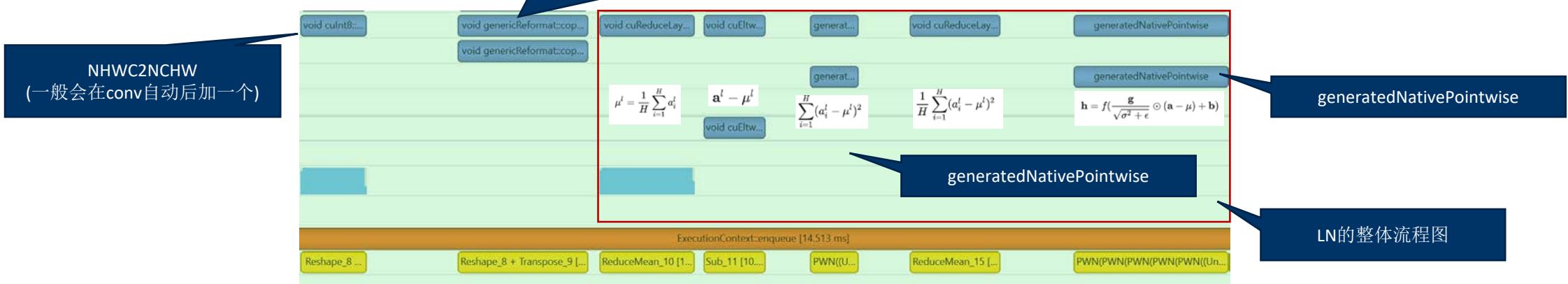
step3: 创建一个可以学习的mean和var，最后在接一个激活函数

$$h^l = f(g^l \odot \hat{a}^l + b^l)$$

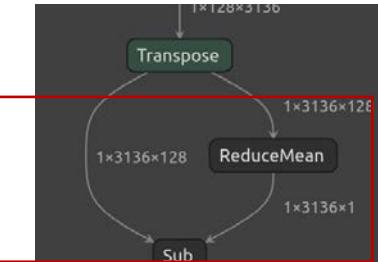
总结下来

$$h = f\left(\frac{g}{\sqrt{\sigma^2 + \epsilon}} \odot (a - \mu) + b\right)$$

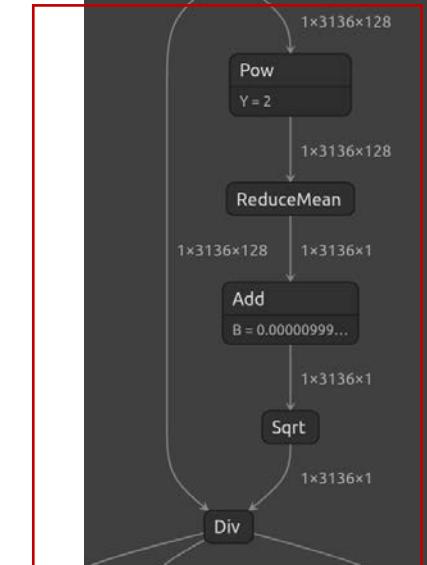
LN所需要的计算时间



Transformer的LN处理



$$a^l - \mu^l$$



$$\sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$\hat{a}^l = \frac{a^l - \mu^l}{\sqrt{(\sigma^l)^2 + \epsilon}}$$

$$h^l = f(g^l \odot \hat{a}^l + b^l)$$

$$h = f\left(\frac{g}{\sqrt{\sigma^2 + \epsilon}} \odot (a - \mu) + b\right)$$

計算時間大概是165 us、主要卡在了reshape和pointwise上了

根据onnx中的Proto信息，修改onnx

HardSigmoid	6, 1	18	
HardSwish	14	14	
LayerNormalization	17	17, 18	
LeakyRelu	16, 6, 1	16	
LessOrEqual	16, 12	16	
LogSoftmax	13, 11, 1	13, 18	

<https://github.com/onnx/onnx/blob/main/docs/Operators.md>

TensorRT 8.6 EA Release - 2023-3-13

Added

For more details, see the 8.6 EA release notes for new features added in TensorRT 8.6.

- Added support for `GroupNormalization`, `LayerNormalization`, `IsInf` operations
- Added support for INT32 input types for `Argmin`, `Argmax`, and `TopK`
- Added support for `ReverseSequence` operators with dynamic shapes
- Added support for `TopK` operators with dynamic `K` values
- Added `OnnxParserFlag` enum and `setFlag` interfaces to the ONNX parser to modify the default parsing behavior
- Added metadata tracking, now ONNX node metadata will be embedded into TensorRT layers

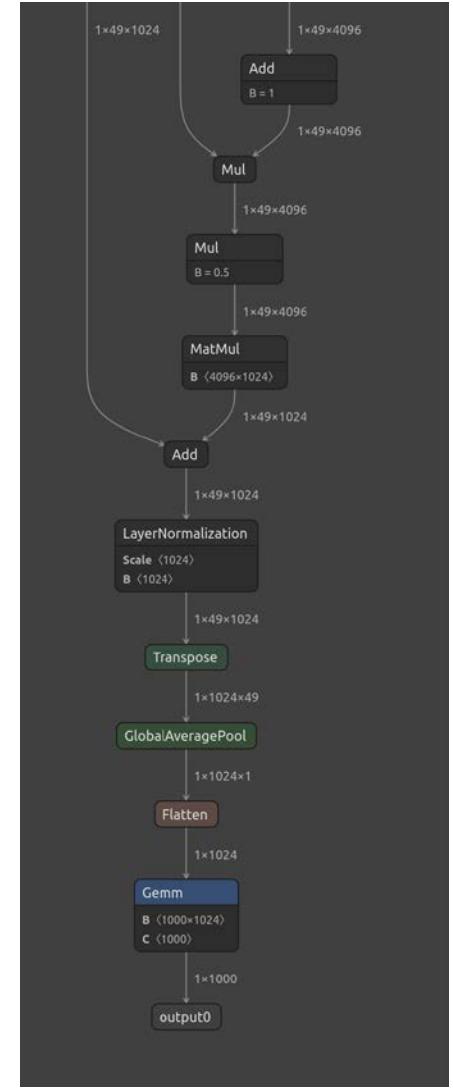
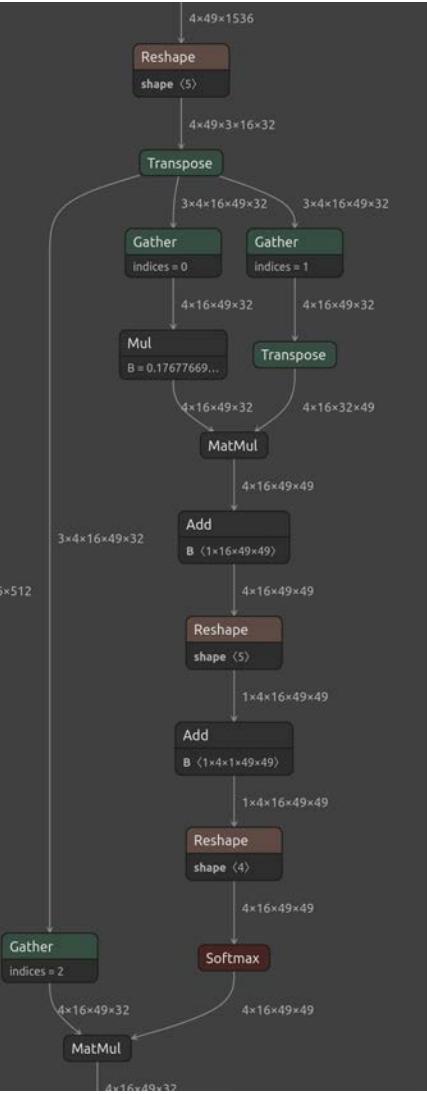
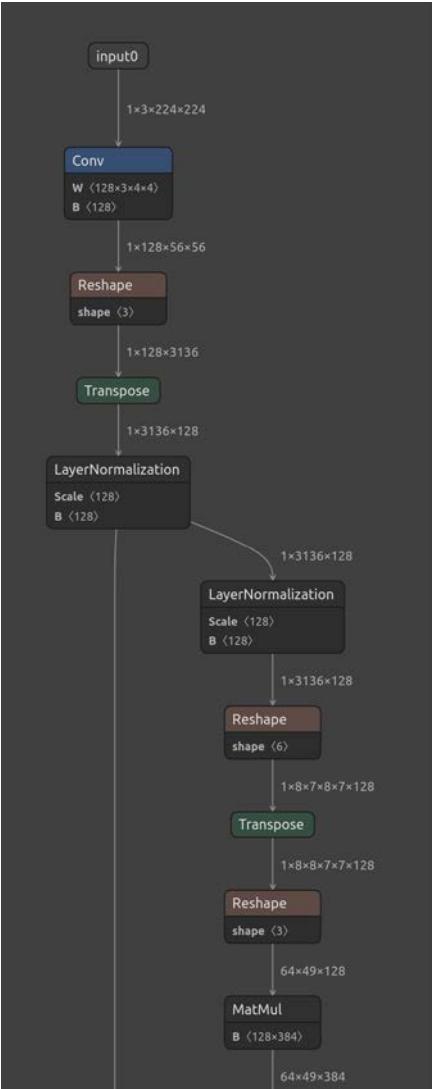
Changed

- All cast operations will now use the new `CastLayer` over the previous `IdentityLayer`.

```
65
66 def export_norm_onnx(model, file, input):
67     torch.onnx.export(
68         model,
69         args=(input,),
70         f=file,
71         input_names=["input0"],
72         output_names=["output0"],
73         opset_version=17)
74
75     print("Finished normal onnx export")
76
77 model_onnx = onnx.load(file)
78
79 # 检查导入的onnx model
80 onnx.checker.check_model(model_onnx)
81
82 # 使用onnx-simplifier来进行onnx的简化。
83 print(f"Simplifying with onnx-simplifier {onnxsim.__version__}...")
84 model_onnx, check = onnxsim.simplify(model_onnx)
85 assert check, "assert check failed"
86 onnx.save(model_onnx, file)
87
```

<https://github.com/onnx/onnx-tensorrt/blob/main/docs/Changelog.md>

成功导出onnx



开头，中间，结尾部分的onnx展示

TensorRT导出

导出了onnx不是重点。我们需要将这些onnx导出为TensorRT模型，并查看性能。我们有几种方式导出

- trtexec命令行 (快速测试推理引擎)
- TensorRT python API (大量的单元测试)
- TensorRT C++ API (结合其他的前处理后处理进行底层部署)



05

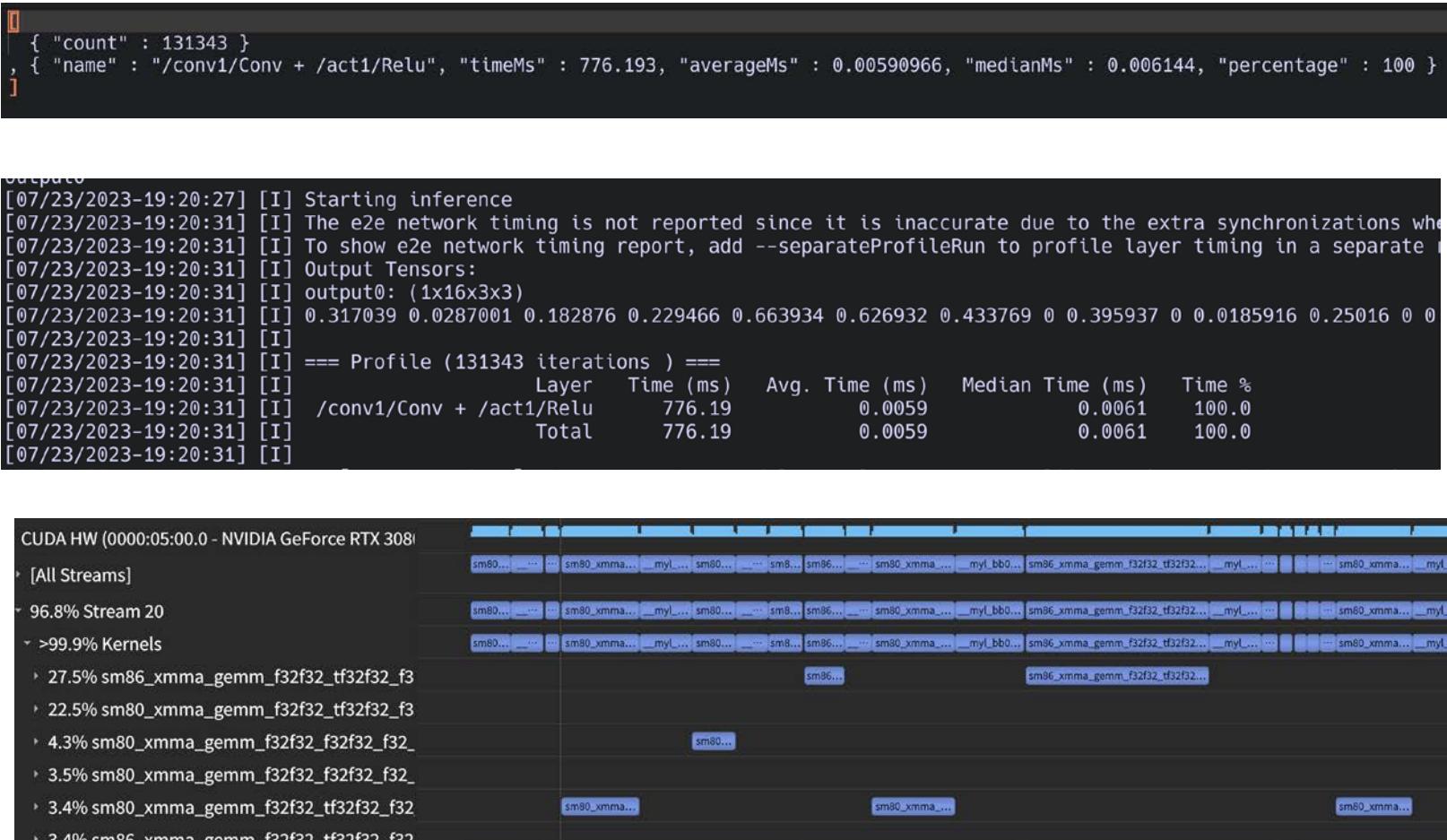
使用trtexec

Goal: 学习使用trtexec，以及通过trtexec的日志理解TensorRT内部所做的优化

执行一下我们的shell脚本

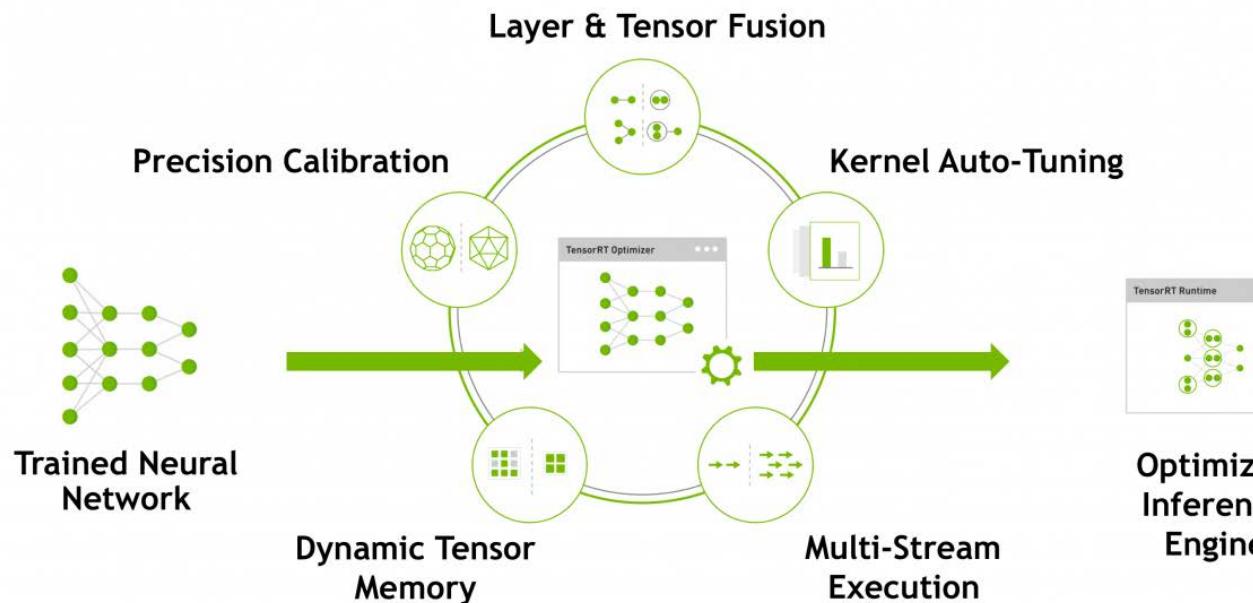
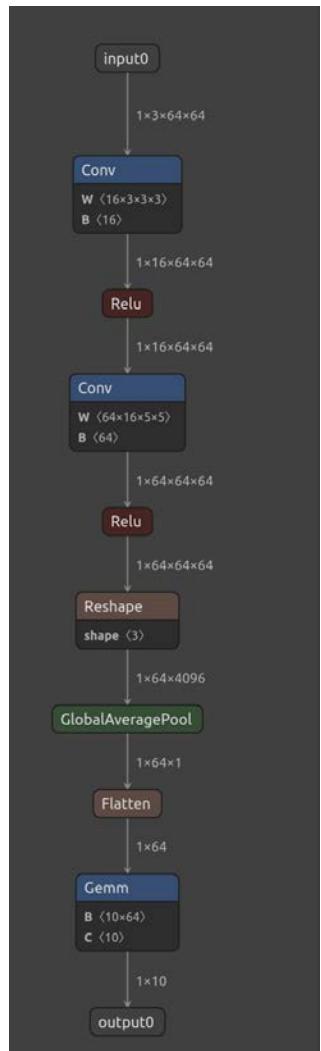
```
..../3.7-trtexec-analysis/
build
  engines
    sample.engine
  log
    sample
      build
        build_layer_info.log
        build.log
        build_output.log
        build_profile.log
      infer
        infer_layer_info.log
        infer.log
        infer_output.log
        infer_profile.log
      profile
        sample.nsys-report
models
  resnet50.engine
  resnet50.onnx
  sample2.engine
  sample2.onnx
  sample.engine
  sample.onnx
  swin-tiny-op17.onnx
  vgg16.engine
  vgg16.onnx
src
  load_torchvision.py
  sample_cbr.py
  sample_reshape.py
tools
  build.sh
  infer.sh
  profile.sh
```

3.7-学习如何使用trtexec



执行结果的一部分展示

是TensorRT中的一个非常好用的快速生成trt推理引擎的工具。便于我们不使用API快速查看推理引擎的加速效果，以及TensorRT优化后的模型架构



```

Layer
Reformatting CopyNode for Input Tensor 0 to /conv1/Conv + /act1/Relu
/conv1/Conv + /act1/Relu
/conv2/Conv + /act2/Relu
Reformatting CopyNode for Input Tensor 0 to /Reshape
/avgpool/GlobalAveragePool
/head/Gemm
Total
  
```

trtexec

trtexec提供的参数非常的多。我这里做了一些几个脚本在tools/目录下，方便大家使用

```
23  
24 trtexec --onnx=${1} \  
25 | | | --memPoolSize=workspace:2048 \  
26 | | | --saveEngine=${ENGINE_PATH}/${PREFIX}.engine \  
27 | | | --profilingVerbosity=detailed \  
28 | | | --dumpOutput \  
29 | | | --dumpProfile \  
30 | | | --dumpLayerInfo \  
31 | | | --exportOutput=${LOG_PATH}/build_output.log\  
32 | | | --exportProfile=${LOG_PATH}/build_profile.log \  
33 | | | --exportLayerInfo=${LOG_PATH}/build_layer_info.log \  
34 | | | --warmUp=200 \  
35 | | | --iterations=50 \  
36 | | > ${LOG_PATH}/build.log
```

tools/build.sh

trtexec

trtexec提供的参数非常的多。我这里做了一些几个脚本在tools/目录下，方便大家使用

```
23  
24 trtexec --loadEngine=${ENGINE_PATH}/${PREFIX}.engine \  
25 | | | | | --verbose \  
26 | | | | | --dumpOutput \  
27 | | | | | --dumpProfile \  
28 | | | | | --dumpLayerInfo \  
29 | | | | | --exportOutput=${LOG_PATH}/infer_output.log\  
30 | | | | | --exportProfile=${LOG_PATH}/infer_profile.log \  
31 | | | | | --exportLayerInfo=${LOG_PATH}/infer_layer_info.log \  
32 | | | | | --warmUp=200 \  
33 | | | | | --iterations=50 \  
34 | | | | > ${LOG_PATH}/infer.log
```

tools/infer.sh

trtexec

trtexec提供的参数非常的多。我这里做了一些几个脚本在tools/目录下，方便大家使用

```
23  
24 nsys profile \  
25 | | | --output=${LOG_PATH}/${PREFIX} \  
26 | | | --force-overwrite true \  
27 | | | trtexec --loadEngine=${ENGINE_PATH}/${PREFIX}.engine \  
28 | | | | --warmUp=200 \  
29 | | | | --iterations=50 \  
30
```

tools/profile.sh

trtexec

trtexec的参数一览，大家去阅读一下TensorRT官方文档。链接如下

<https://docs.nvidia.com/deeplearning/tensorrt/developer-guide/index.html#trtexec>

trtexec可以配合nvidia-smi一起使用进行profile分析

我们对模型进行优化的时候，这个是我要做的第一步事情

