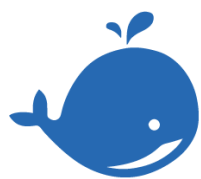


## 第三课-计算图的定义



# Datawhale

# KuiperInfer

## 计算图的格式

### 计算图的概念

KuiperInfer 使用的模型格式是 PNNX. PNNX 是 PyTorch Neural Network Exchange 的缩写，它的愿景是将 PyTorch 模型文件直接导出为高效、简洁的计算图。计算图的概念如第一章说的那样，一般包括了以下的几个部分：

1. **Operator**: 深度学习计算图中的计算节点。
2. **Graph**: 有多个 **Operator** 串联得到的有向无环图，规定了各个计算节点 (**Operator**) 执行的流程和顺序。
3. **Layer**: **计算节点中**运算的具体执行者，**Layer** 类先读取输入张量中的数据，然后对输入张量进行计算，得到的结果存放到计算节点的输出张量中，**当然，不同的算子中 Layer 的计算过程会不一致。**
4. **Tensor**: 用于存放**多维数据**的数据结构，方便数据在计算节点之间传递，同时该结构也封装矩阵乘、点积等与矩阵相关的基本操作。

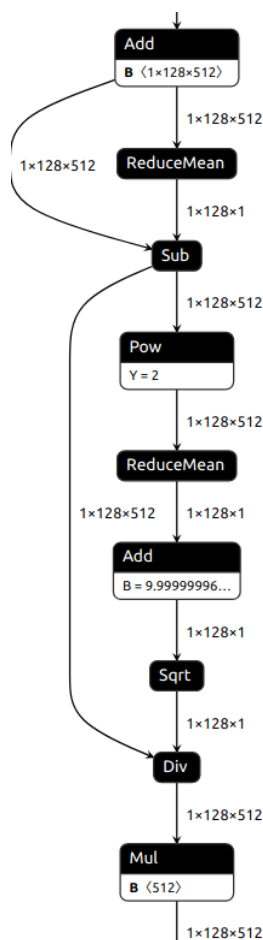
如果你对以上的内容已经没有印象了，可以自行回顾本课程的第一章。作为一种计算图形式，PNNX 自然也不例外。接下来，我们来探讨一下为什么在出现了 **ONNX** 计算图之后还需要 **PNNX**，以及它解决了什么问题？

# PNNX计算图的优势

参考资料：<https://zhuanlan.zhihu.com/p/427620428>

以往我们将训练好的模型导出为 ONNX 结构之后，模型中的一个复杂算子不仅经常会被拆分成多个细碎的算子，而且为了将这些细碎的算子拼接起来完成原有算子的功能，通常还需要一些称之为“胶水算子”的辅助算子，例如 Gather 和 Unsqueeze 等。

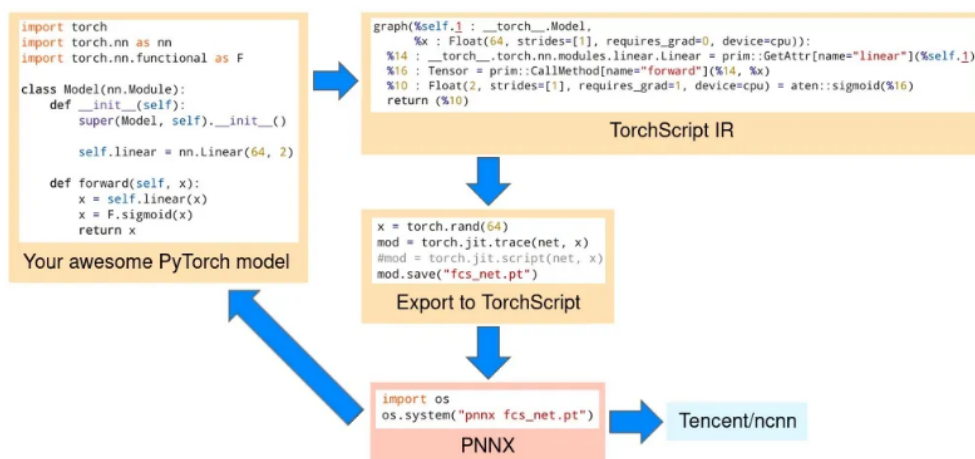
你还认得出下图的算子，其实是被拆分后 LayerNorm 算子吗？当然 ONNX 使用多个小算子去等价一个复杂算子的设计，也是为了用尽可能少的算子去兼容更多的训练框架。



但是过于细碎的计算图会不仅不利于推理的优化。另外，拆分的层次过于细致，也会导致算法工程师难以将导出的模型和原始模型进行结构上的相互对应。为了解决以上说到的问题，我们选用 NCNN 那推理框架的计算图格式之一 PNNX，那么 PNNX 给我们带来了什么呢？下图是它在模型部署中的位置。



## PNNX is yet another open standard



1. 使用模板匹配（pattern matching）的方法将匹配到的子图用对应等价的大算子替换掉，例如可以将上图子图中的多个小算子（已经在 TorchScript 中被拆分的）重新替换为 LayerNorm 算子。或者在对 PyTorch 模型导出时，也可以自定义某个 nn.Module 不被拆分；
2. 在 PyTorch 中编写的简单算术表达式在转换为 PNNX 后，会保留表达式的整体结构，而不会被拆分成许多小的加减乘除算子。例如表达式 `add(mul(@0, @1), add(@2, @3))` 不会被拆分为两个 add 算子和一个 mul 算子，而是会生成一个表达式算子 Expression；
3. PNNX 项目中有大量图优化的技术，包括了算子融合，常量折叠和消除，公共表达式消除等技术。
  - 算子融合优化是一种针对深度学习神经网络的优化策略，通过将多个相邻的计算算子合并为一个算子来减少计算量和内存占用。以卷积层和批归一化层为例，我们可以将两个算子合并为一个新的算子，也就是将卷积的公式带入到批归一化层的计算公式中，有如下公式：

$$Conv = w * x + b$$

$$BN = \gamma \frac{x - \hat{u}}{\sigma^2 + \epsilon} + \beta$$

其中  $w$  是卷积层的权重， $b$  是卷积层的偏移量， $\hat{u}$  和  $\sigma$  依次是样本的均值和方差， $\epsilon$  为一个极小值。带入后有：

$$Fused = \gamma \frac{(w * x + b) - \hat{u}}{\sigma^2 + \epsilon} + \beta$$

- 常量折叠是将在编译时期间将**表达式中的常量计算出来**，然后将**结果替换为一个等价的常量**，以减少模型在运行时的计算量。
- 常量移除就是将计算图中不需要的常数（**计算图推理的过程中未使用**）节点删除，从而减少计算图的文件和加载后的资源占用大小。
- 公共表达式消除优化是一种针对计算图中重复计算的优化策略，它**可以通过寻找并合并重复计算的计算节点，减少模型的计算量和内存占用**。

公共子表达式检测是指**查找计算图中相同的子表达式**，公共子表达式消除是指**将这些重复计算的计算节点合并为一个新的计算节点**，从而减少计算和内存开销。举个例子：

```
X = input(3, 224, 224)
A = Conv(X)
B = Conv(X)
C = A + B
```

在上方的代码中，`Conv(X)` 这个结果被计算了两次，公共子表达式消除可以将它优化为如下代码，这样一来就少了一次卷积的计算过程。

```
X = input(3, 224, 224)
T = Conv(X)
C = T + T
```

综上所述，如果在我们推理框架的底层用 `PNNX` 计算图，就可以吸收图优化和算子融合的结果，使得推理速度更快更高效。