

JIT in MegEngine

王彪

wangbiao@megvii.com

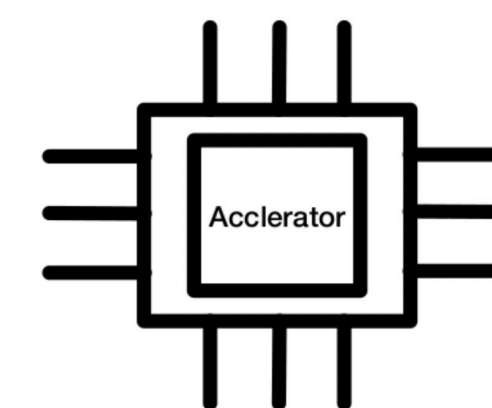
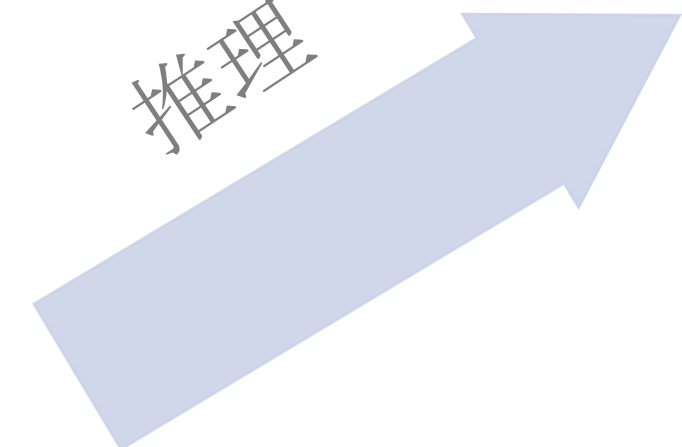




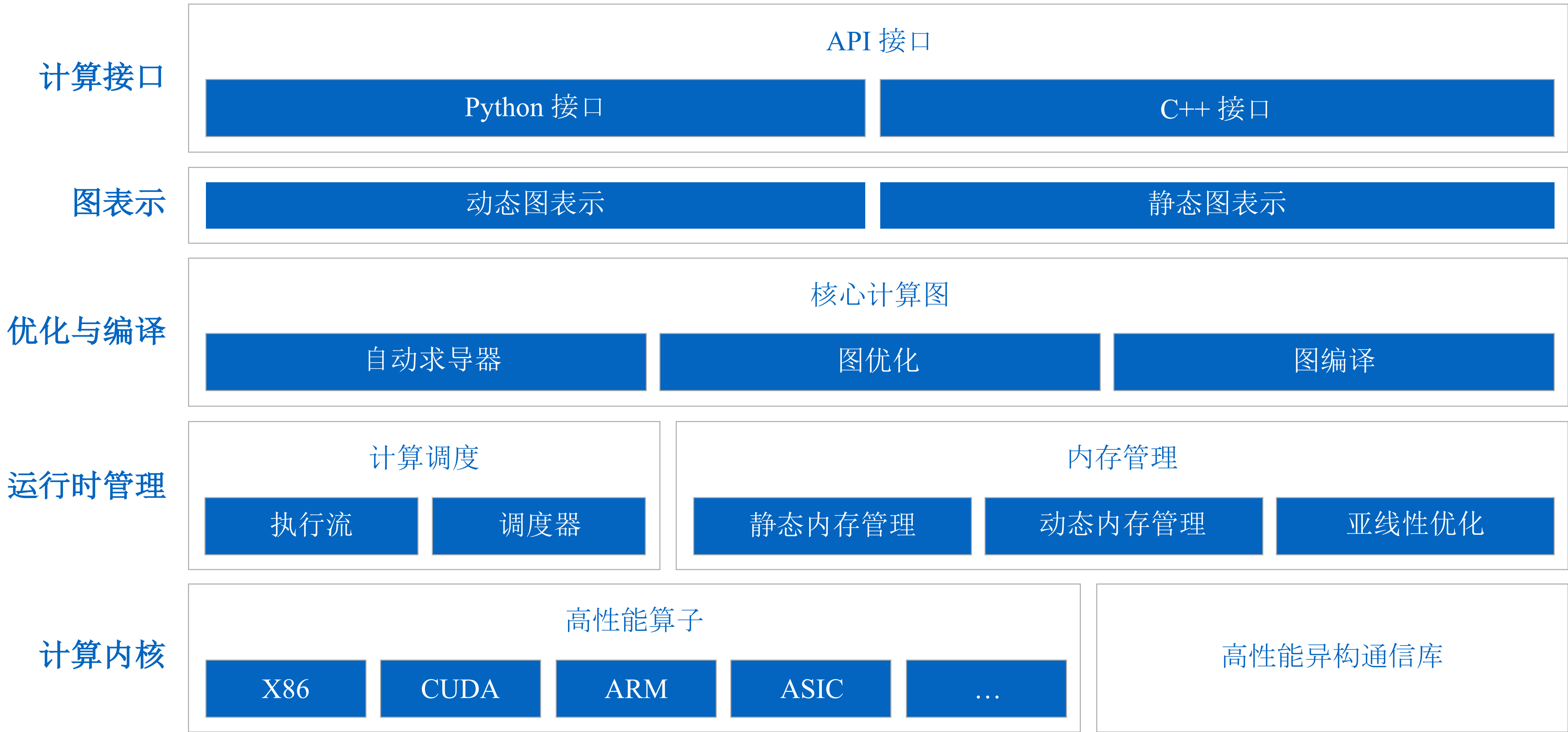
训练



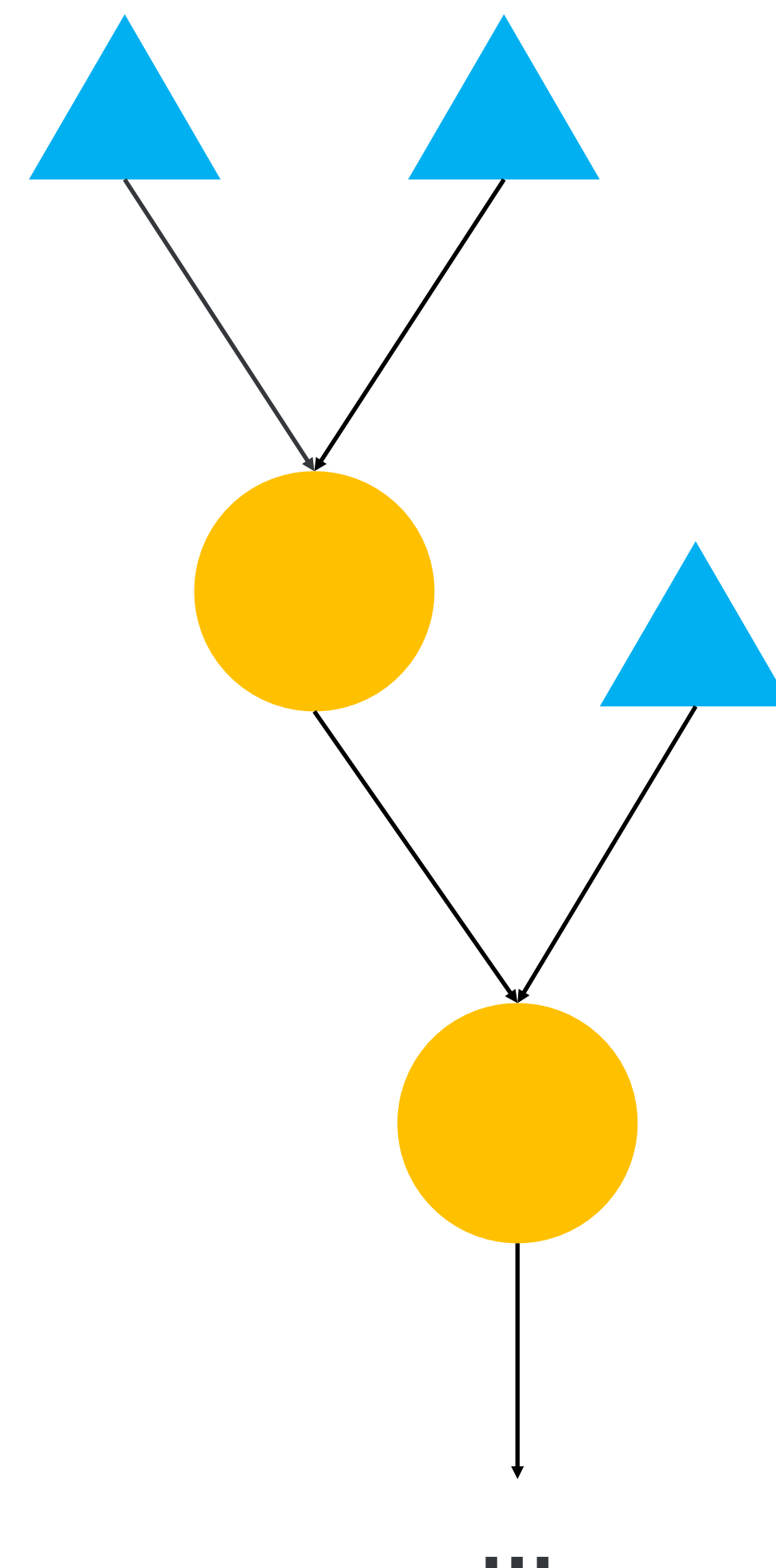
推理

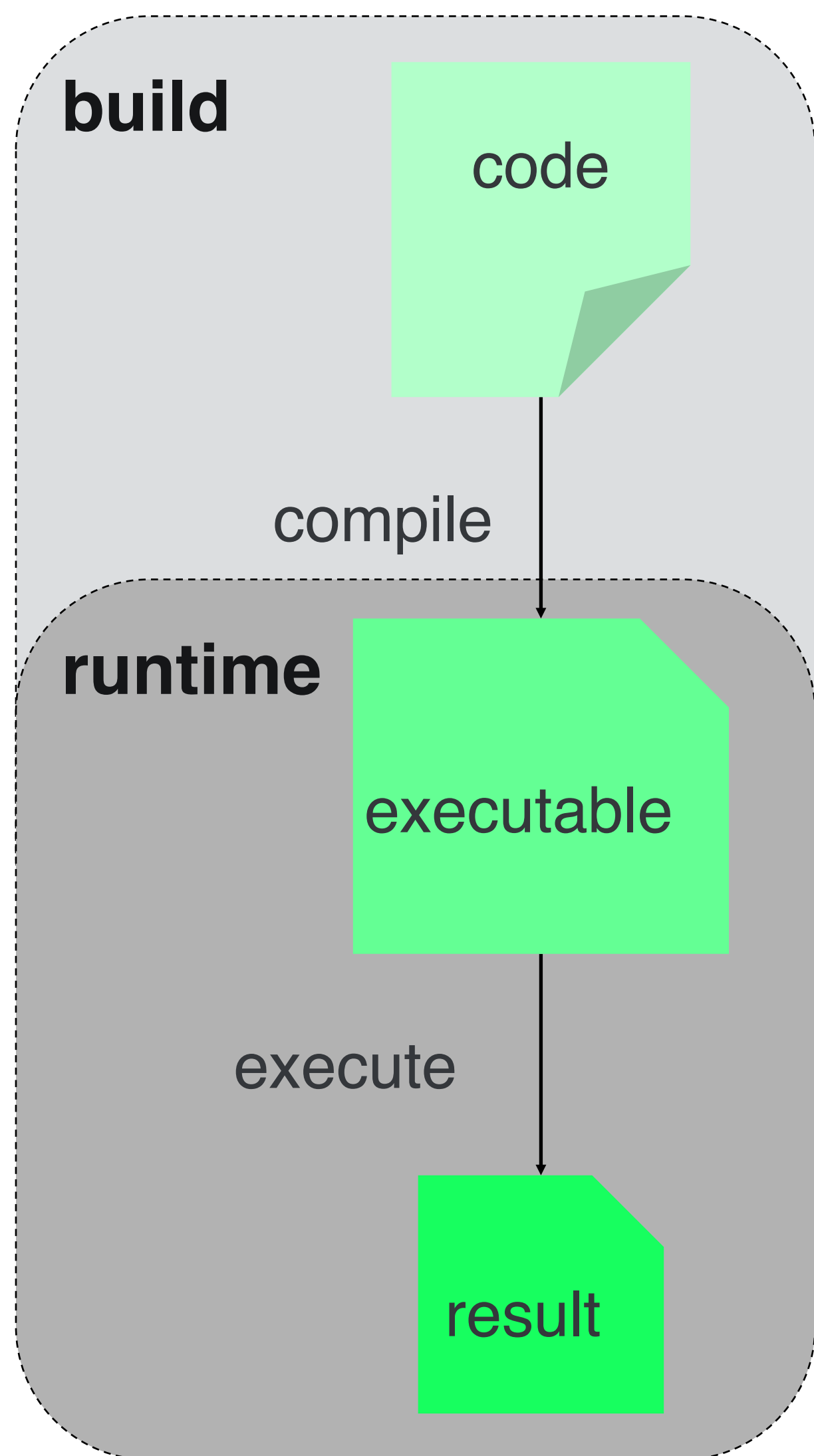


- ✓ 训推一体
- ✓ 动静结合
- ✓ 多平台高性能支持

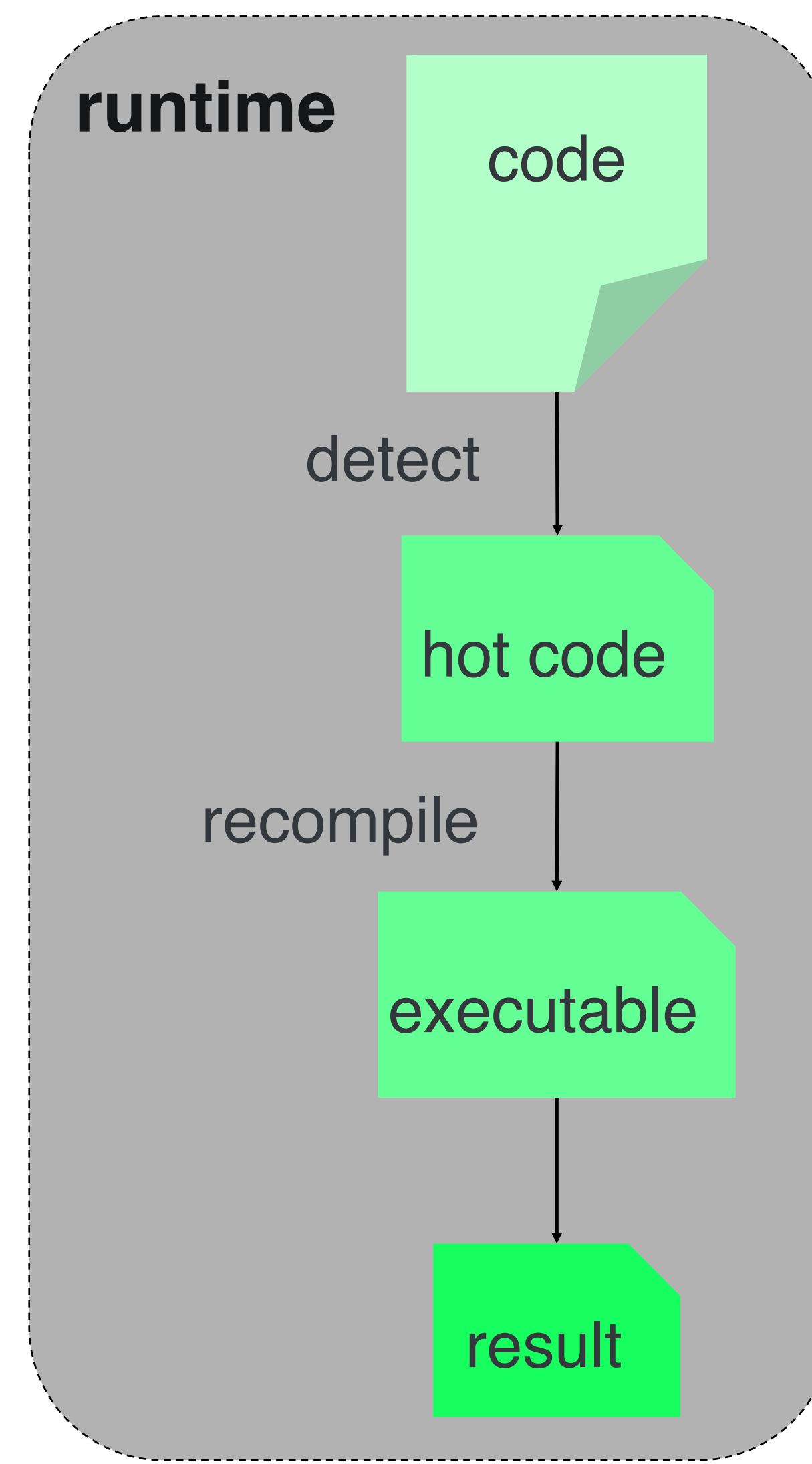


```
if __name__ == '__main__':  
    gm = ad.GradManager().attach(model.parameters())  
    opt = optim.SGD(model.parameters(), lr=0.0125, momentum=0.9,  
weight_decay=1e-4,  
    # 通过 trace 转换为静态图  
    @trace(symbolic=True)  
    def train():  
        with gm:  
            logits = model(image)  
            loss = F.loss.cross_entropy(logits, label)  
            gm.backward(loss)  
            opt.step()  
            opt.clear_grad()  
            return loss  
    loss = train()  
    loss.numpy()
```





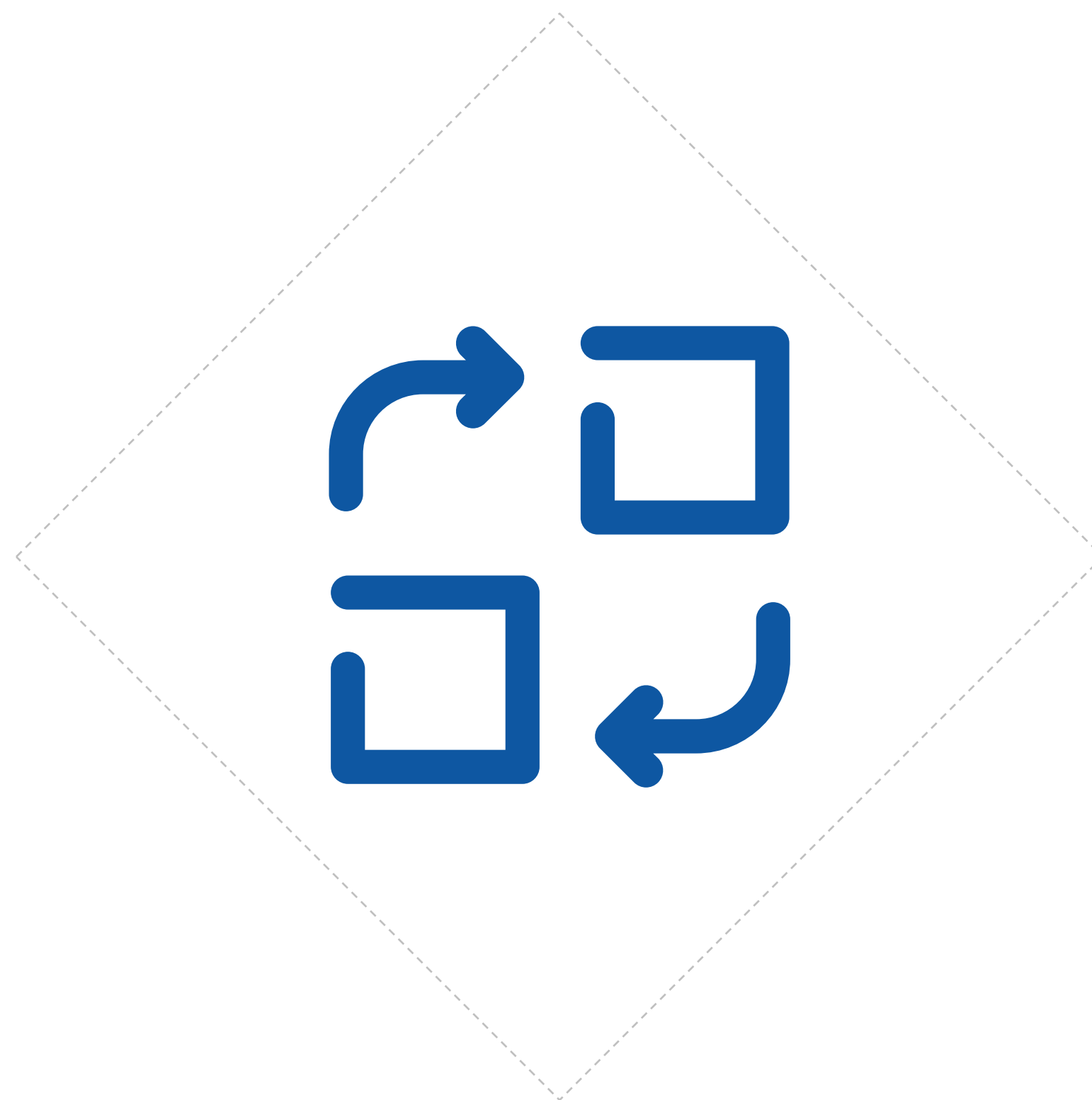
Ahead-of-time (AOT)



Just-in-time (JIT)

① Motivation

④ Evaluation



② Detect Fusion

③ JIT Compilation

Unary	
RELU	NEGATE
ABS	SIGMOID
ACOS	SIN
ASIN	TANH
CEIL	FAST_TANH
COS	ROUND
EXP	ERF
EXPM1	ERFINV
FLOOR	ERFC
LOG	ERFCINV
LOG1P	HSWISH

Binary	
ABS_GRAD	SWISH_GT0
ADD	TANH_GRAD
FLOOR_DIV	TRUE_DIV
MAX	LOG_SUM_EXP
MIN	LT
MOD	LEQ
EXP	EQ
MUL	SHL
POW	SHR
SIGMOID_GRAD	FAST_TANH_GRAD
SUB	RMULH
ATAN2	H_SWISH_GRAD

Fused Binary	
FUSE_ADD_RELU	max(x+y, 0)
FUSE_ADD_SIGMOID	1/(1+exp(-(x+y)))
FUSE_ADD_TANH	tanh(x+y)
FUSE_ADD_H_SWISH	hawish(x+y)

More	
COND_LEQ_MOV	x <= y ? z : 0
FUSE_MUL_ADD3	a * b + c
FUSE_MUL_ADD4	a * b + c * d

model	batchsize	elemwise computation / total computation(%)	elemwise time/total time(%)
resnet50	1	0.102054232	4.6
	8	0.102085366	10.8
	16	0.102080177	11.9
mobilenetV2	1	0.703333333	4.1
	8	0.704592902	8.9
	16	0.704627697	11.9
vgg16	1	0.02927408	5.8
	8	0.029277172	7.1
	16	0.029342568	9.4

training 图纯 device kernel 执行时间

设子图中的 operator 集合为 O_i ($0 \leq i \leq N$)

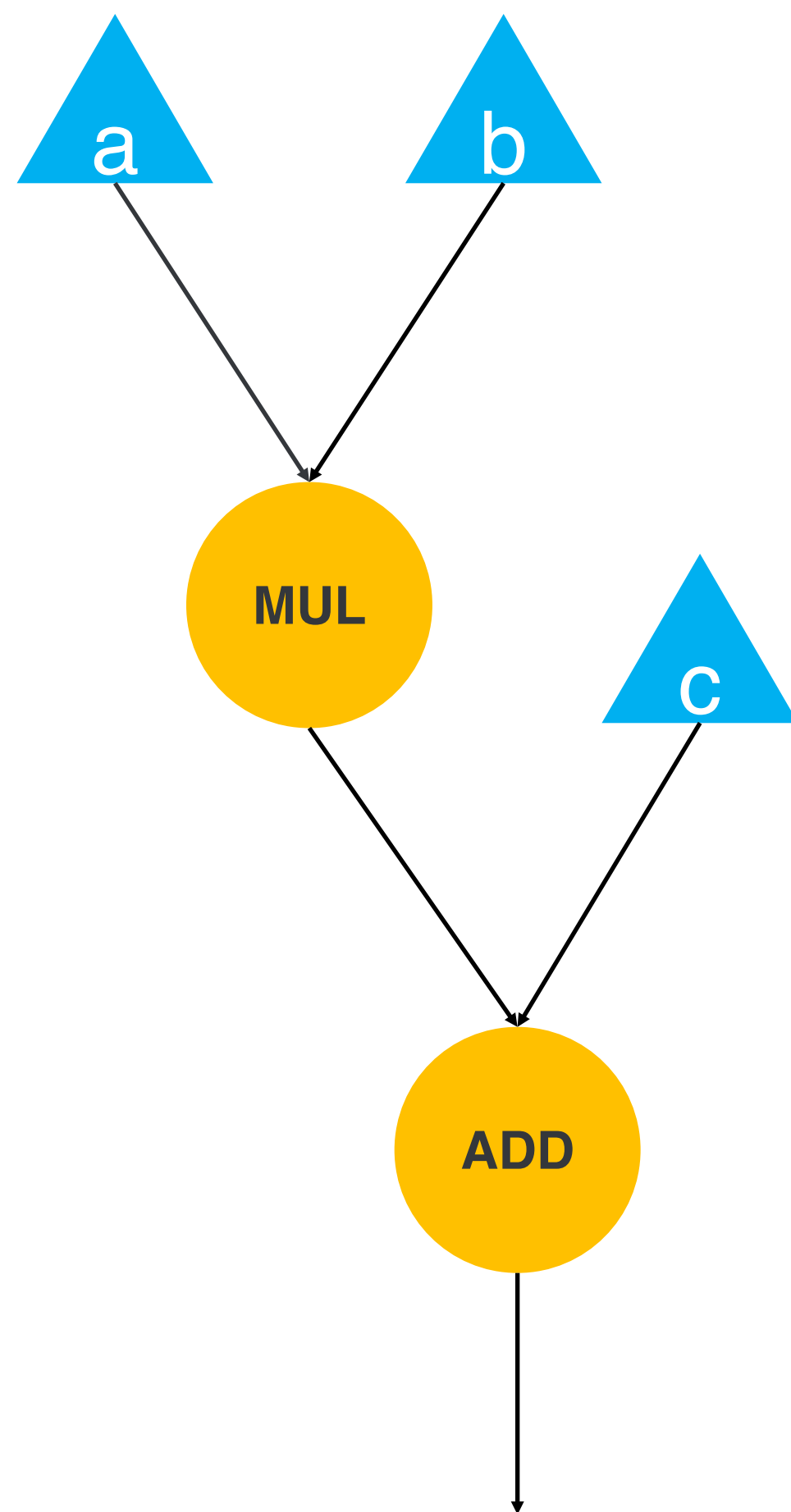
INP_i 表示第 i 个 operator 的 input 个数

OUP_i 表示第 i 个 operator 的 output 个数

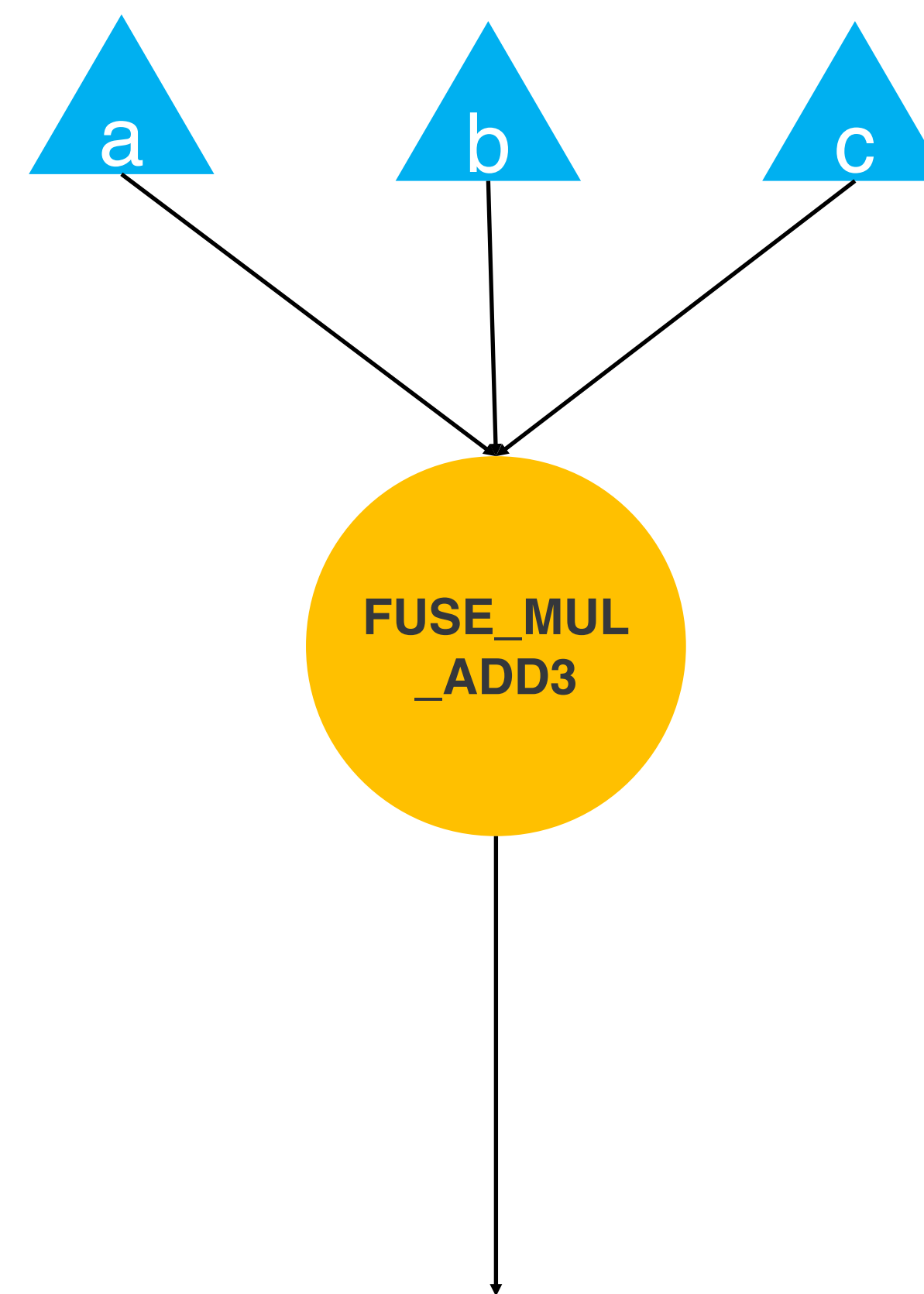
则整个子图的读次数为 $\sum_{i=0}^N INP_i$, 写次数为 $\sum_{i=0}^N OUP_i$.

Elemwise Fusion 可以减少子图中的 operator 的数量, 自然可以减少 operator 之间的读写次数.

Elemwise Fusion

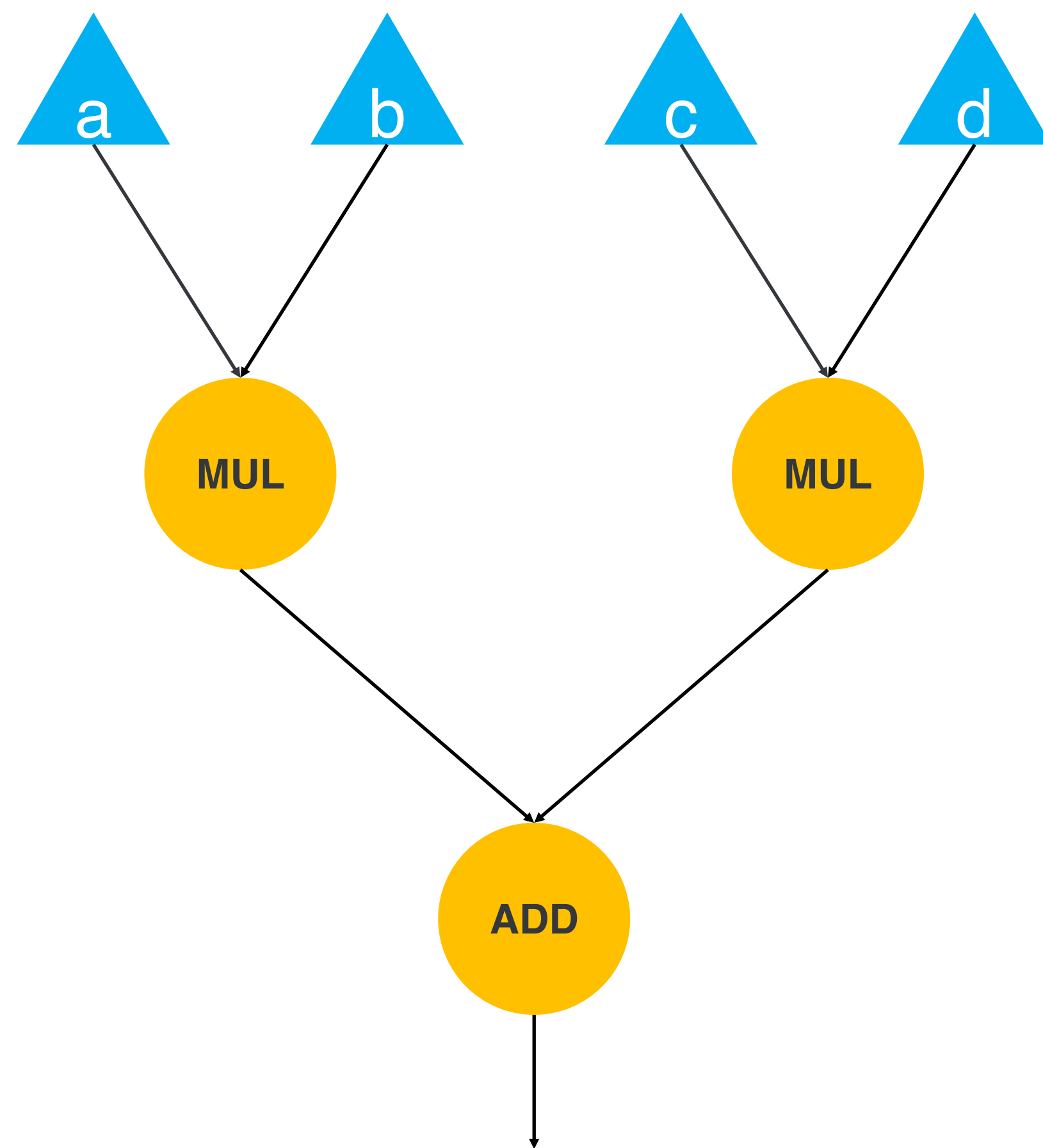


4 次读, 2 次写

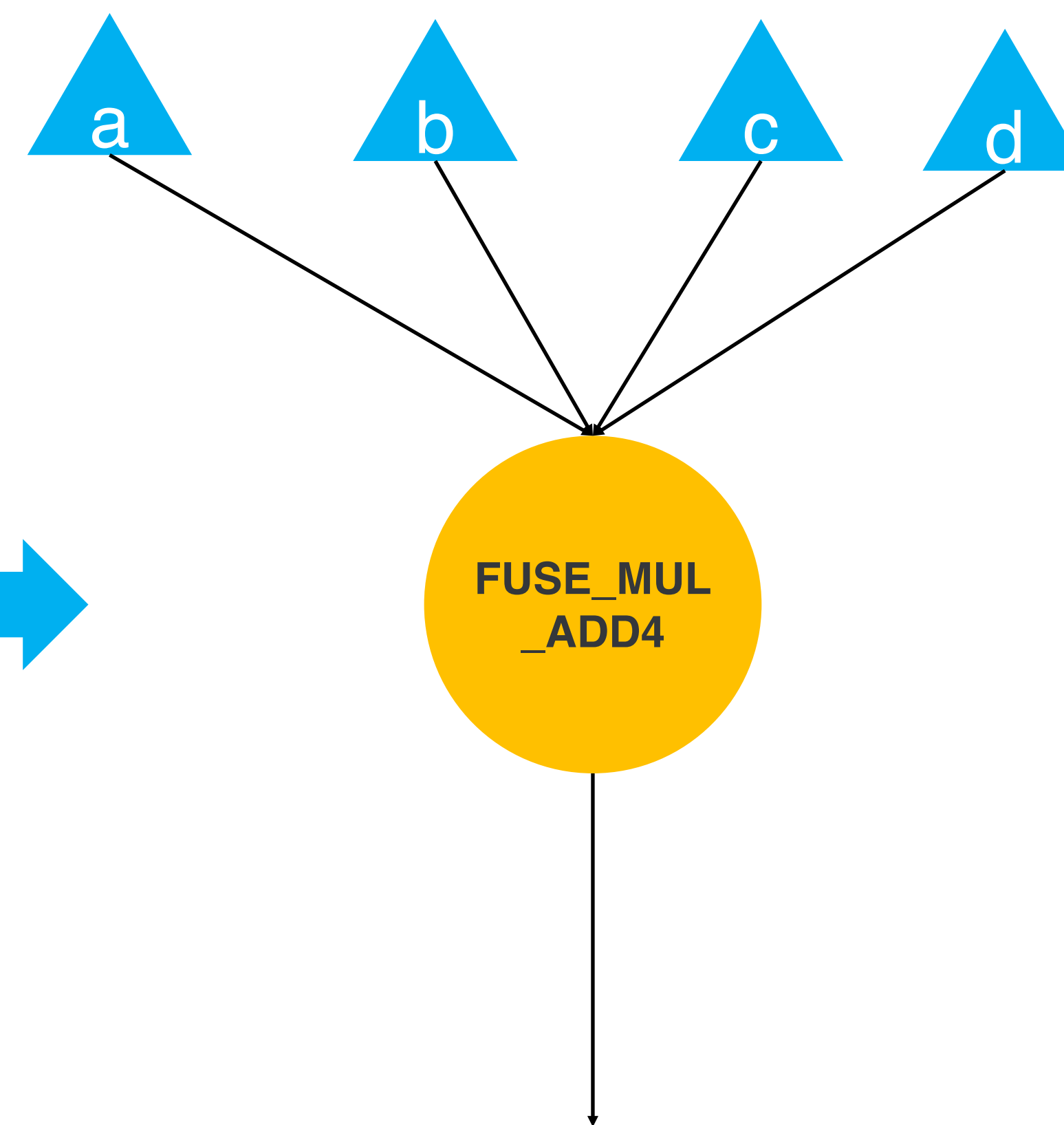


3 次读, 1 次写

Elemwise Fusion

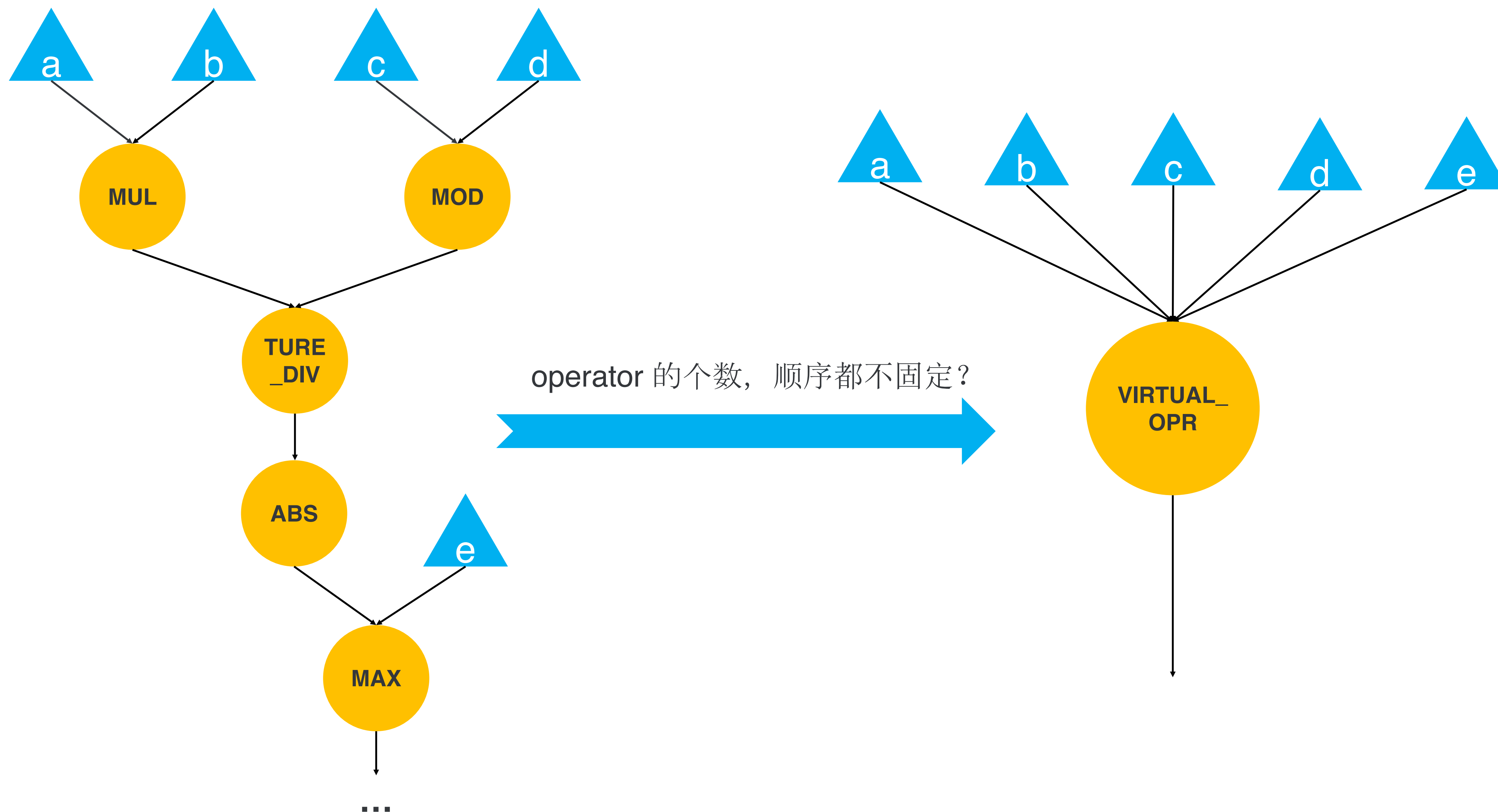


6 次读, 3 次写



4 次读, 1 次写

Elemwise Fusion



CNN 模型特征

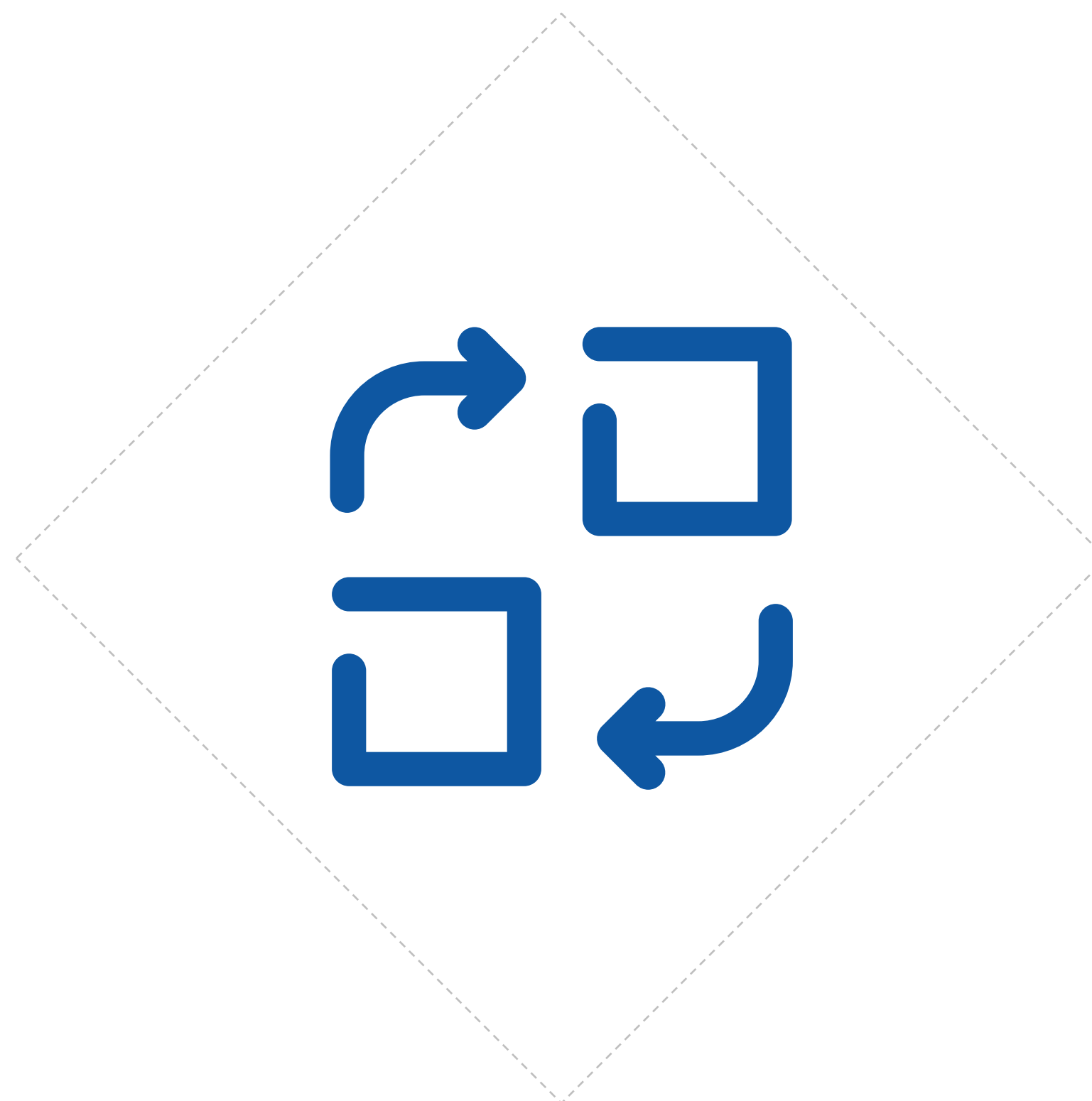
- 静态图模式下模型结构一般不变
- 模型训练过程历经很多个 `iter/min_batch`, 输入数据的 `shape` 不变

使用 JIT 的好处

- 仅需要在首次运行时编译一次
- 有较强的可移植性
- 解决 `elemwise` 模式组合爆炸的问题

① Motivation

④ Evaluation

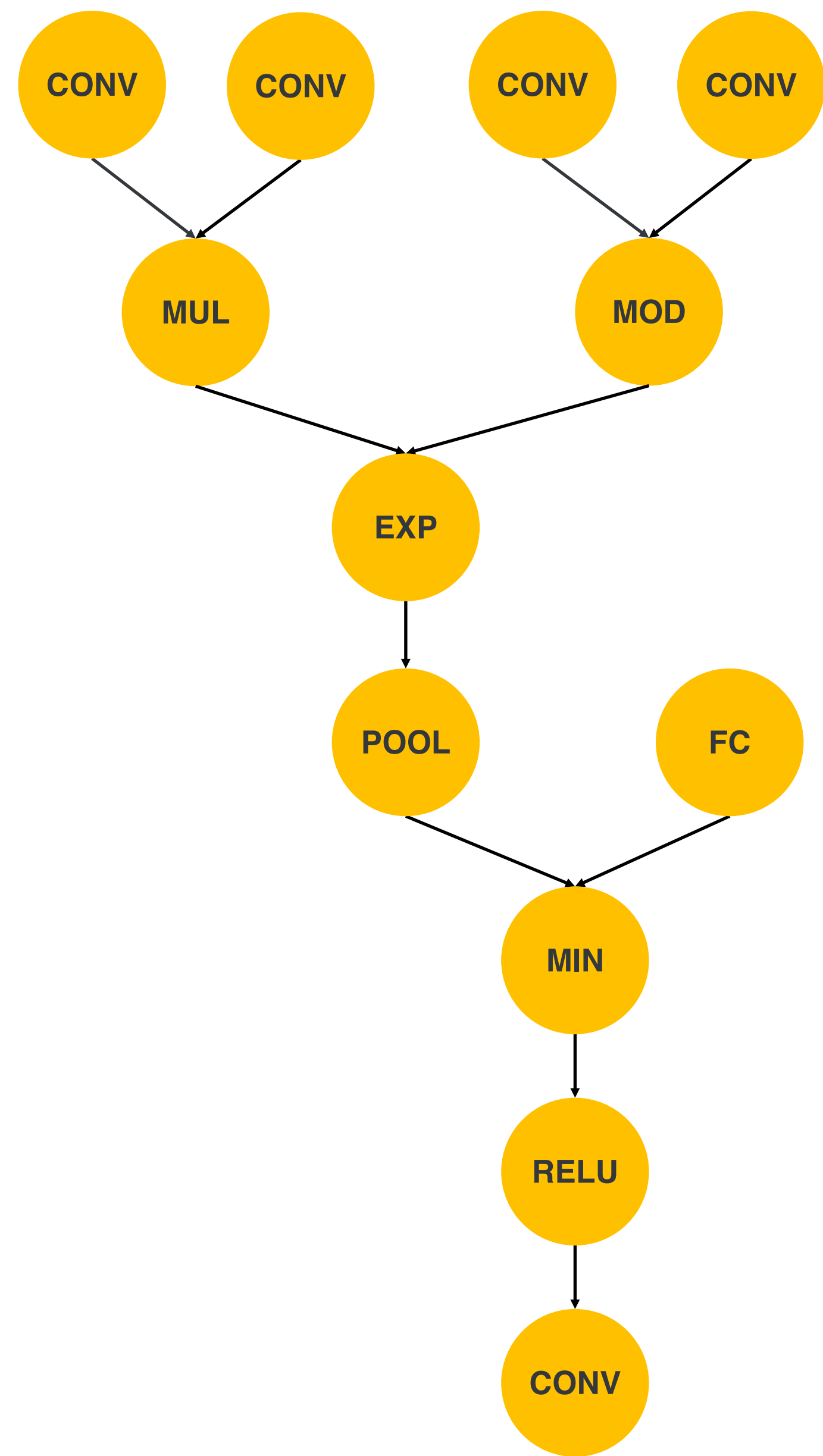


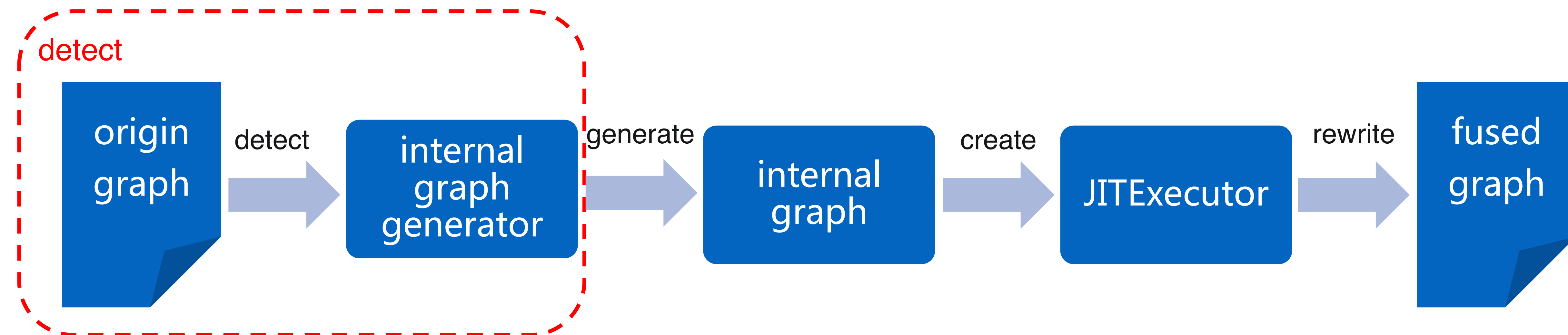
② Detect Fusion

③ JIT Compilation

| Detect Fusion

该计算图中可融合的子图是？





设

G : 计算图

opr : 图 **G** 中的算子

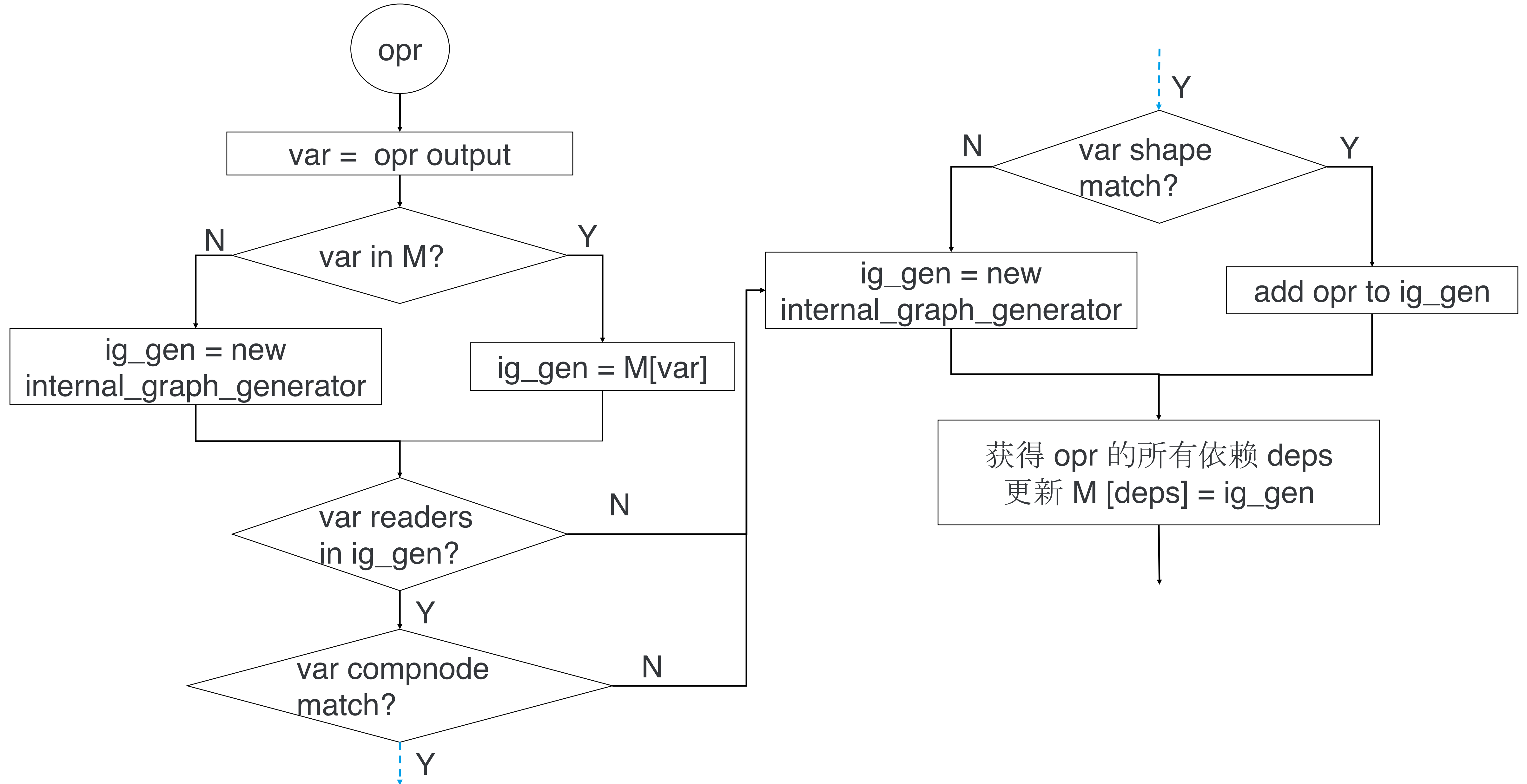
var : **opr** 的输入/输出

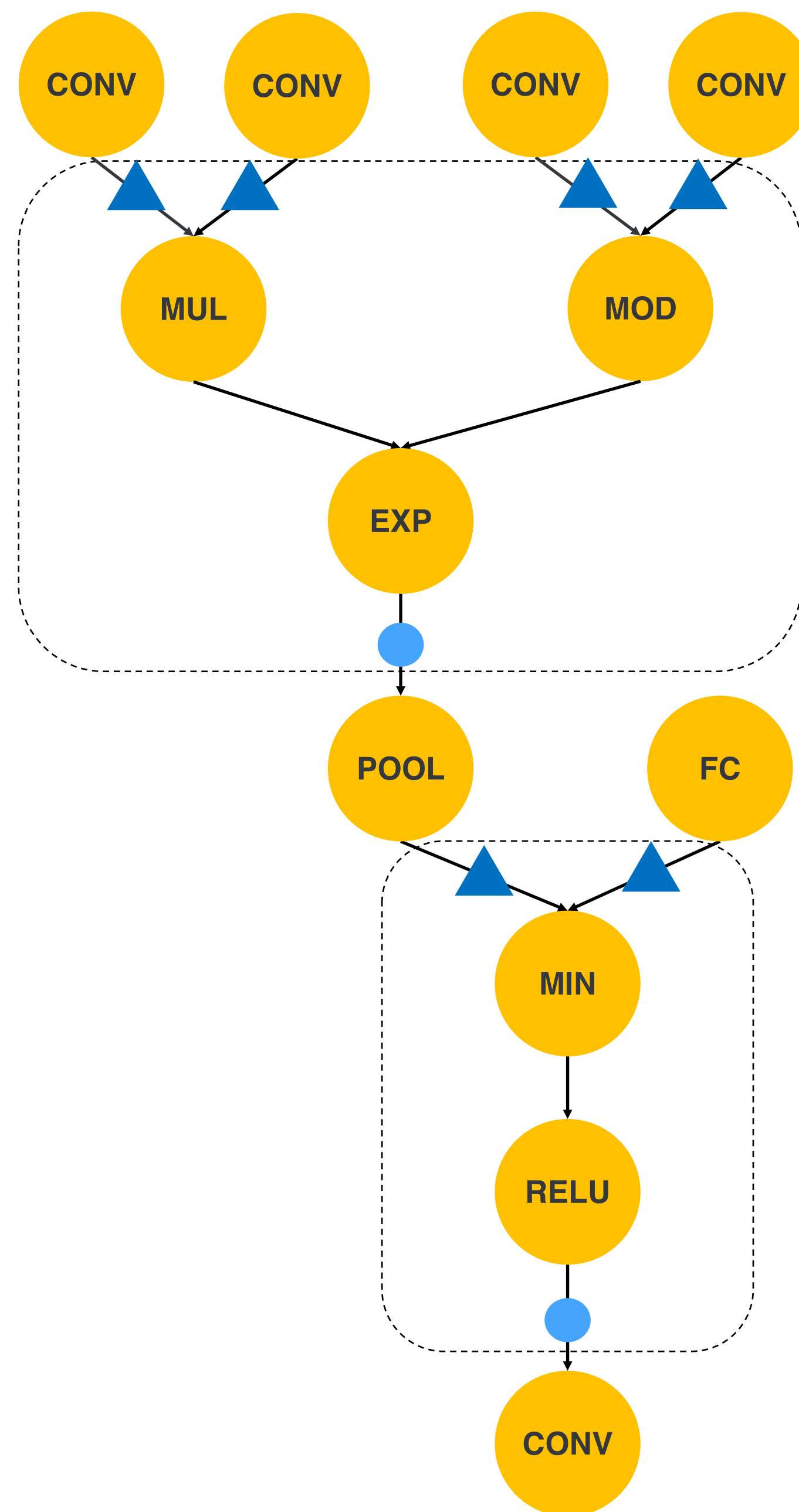
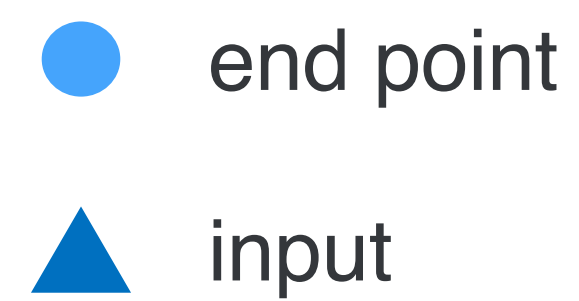
Algorithm : detect – 检测可以 **fuse** 的子图

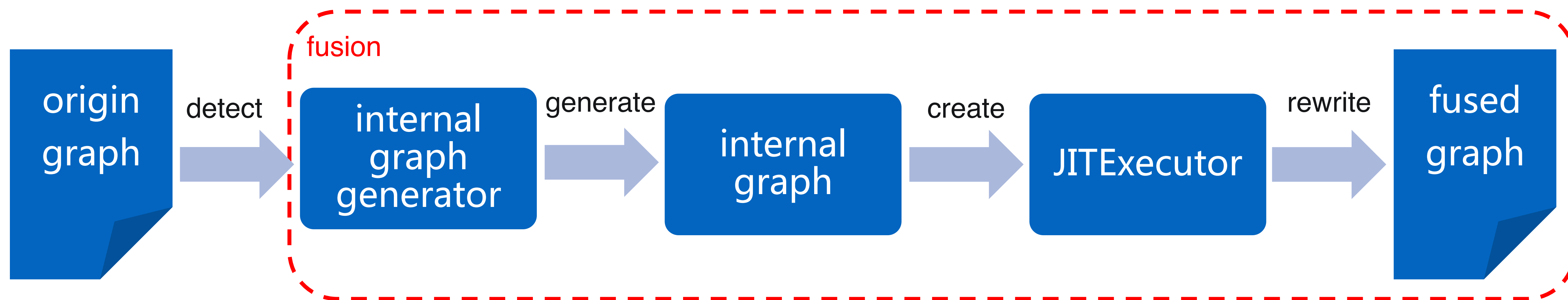
input : G – 原始计算图

output : M – 子图的 **output var** 与对应的 **internal graph generator**

1. 按照逆拓扑序列遍历图 G 中的算子 opr
2. 如果 opr 不是 Elemwise/PowC/TypeCvt/Reduce/Dimshuffle/JITExecutor, 返回步骤1
3. 如果 opr 的 input/output 数据类型不是 float32/float16, 返回步骤1
4. **process_opr(opr)**
5. 转到步骤 1







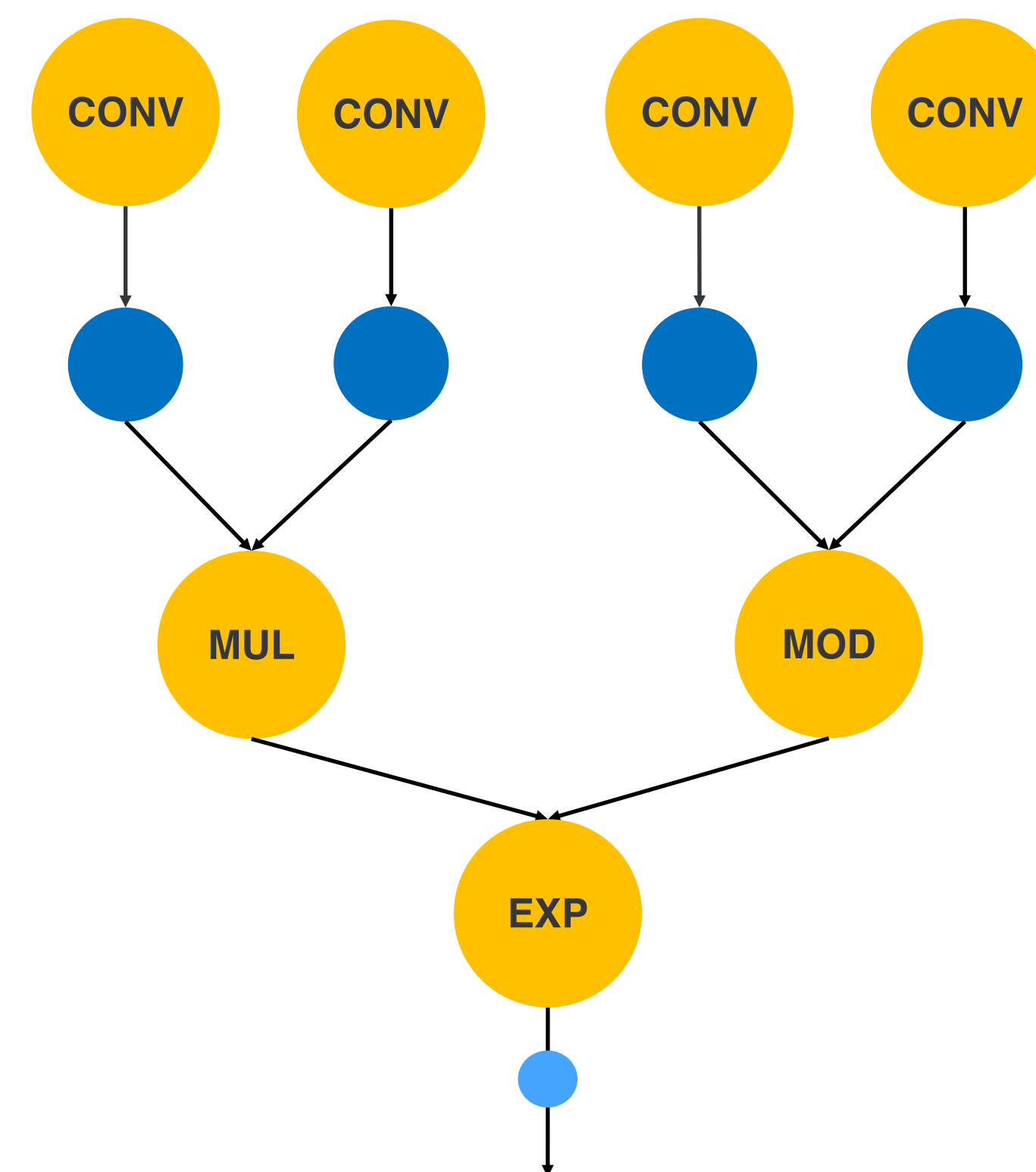
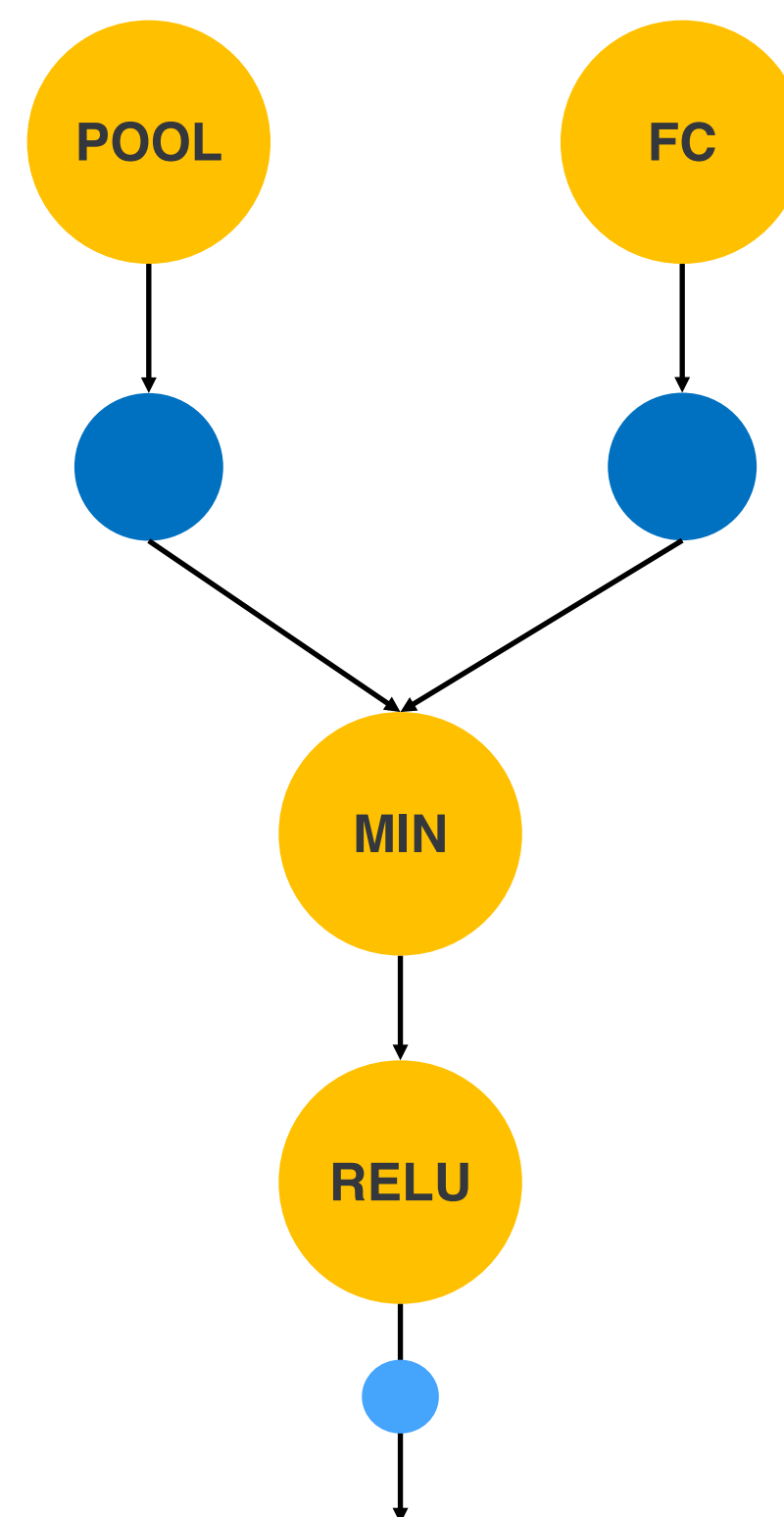
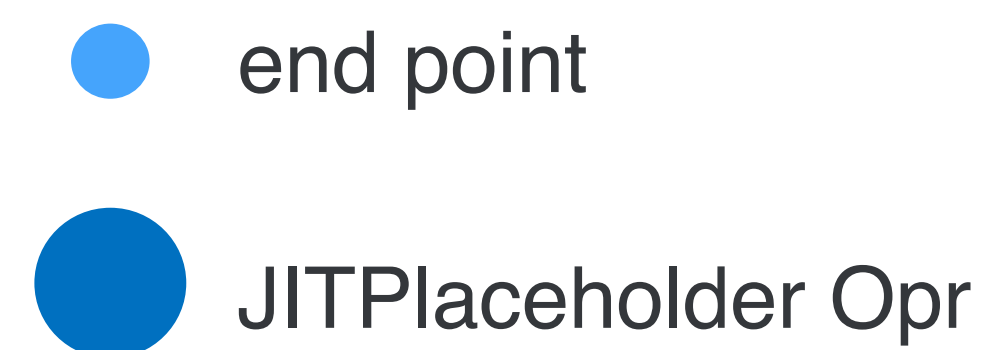
Algorithm : fusion – 将检测出的子图 **fuse** 成一个 **opr**

input : G – 原始计算图; **M** – 子图的 **output var** 与对应的 **internal graph generator** ;

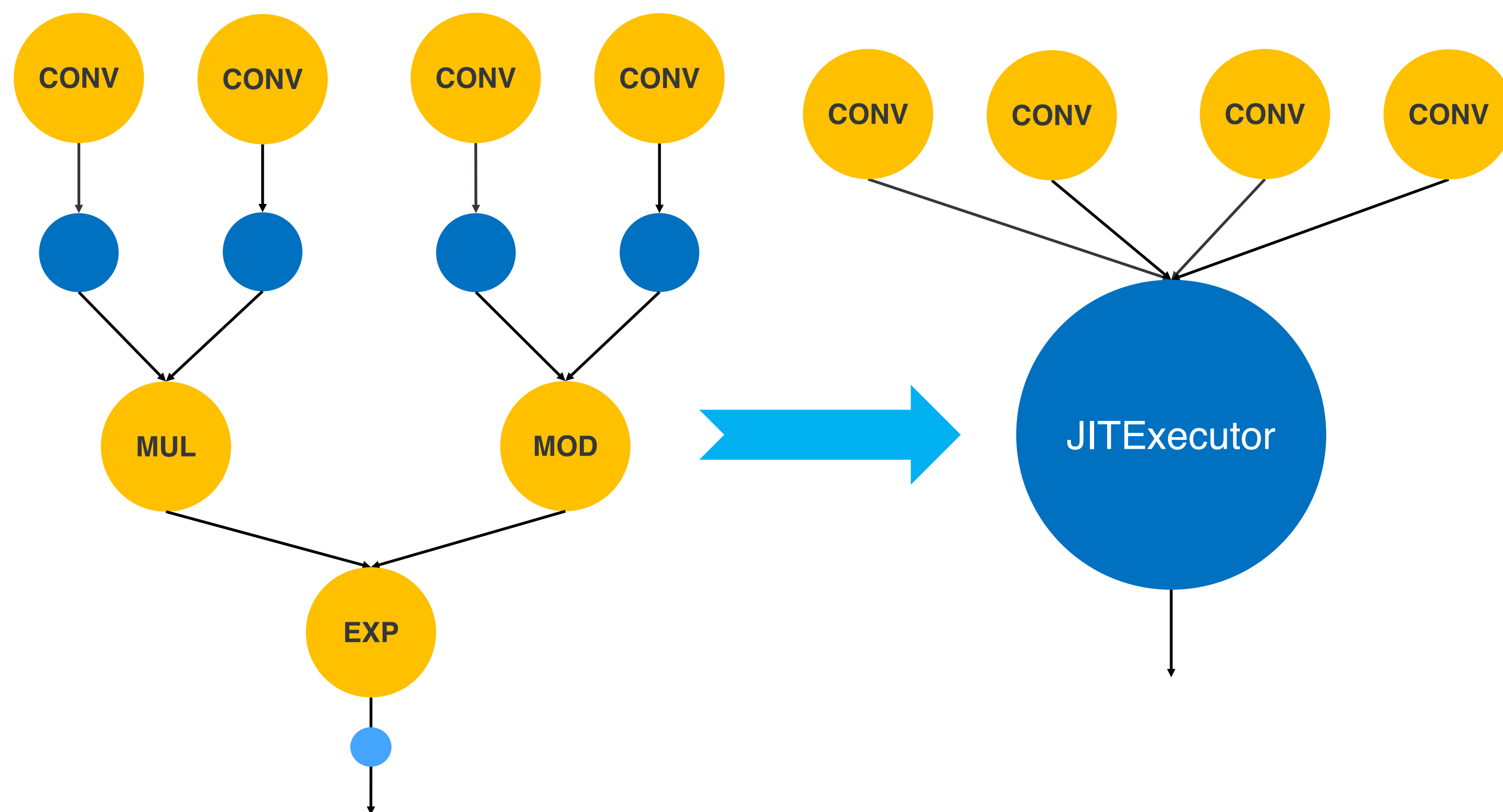
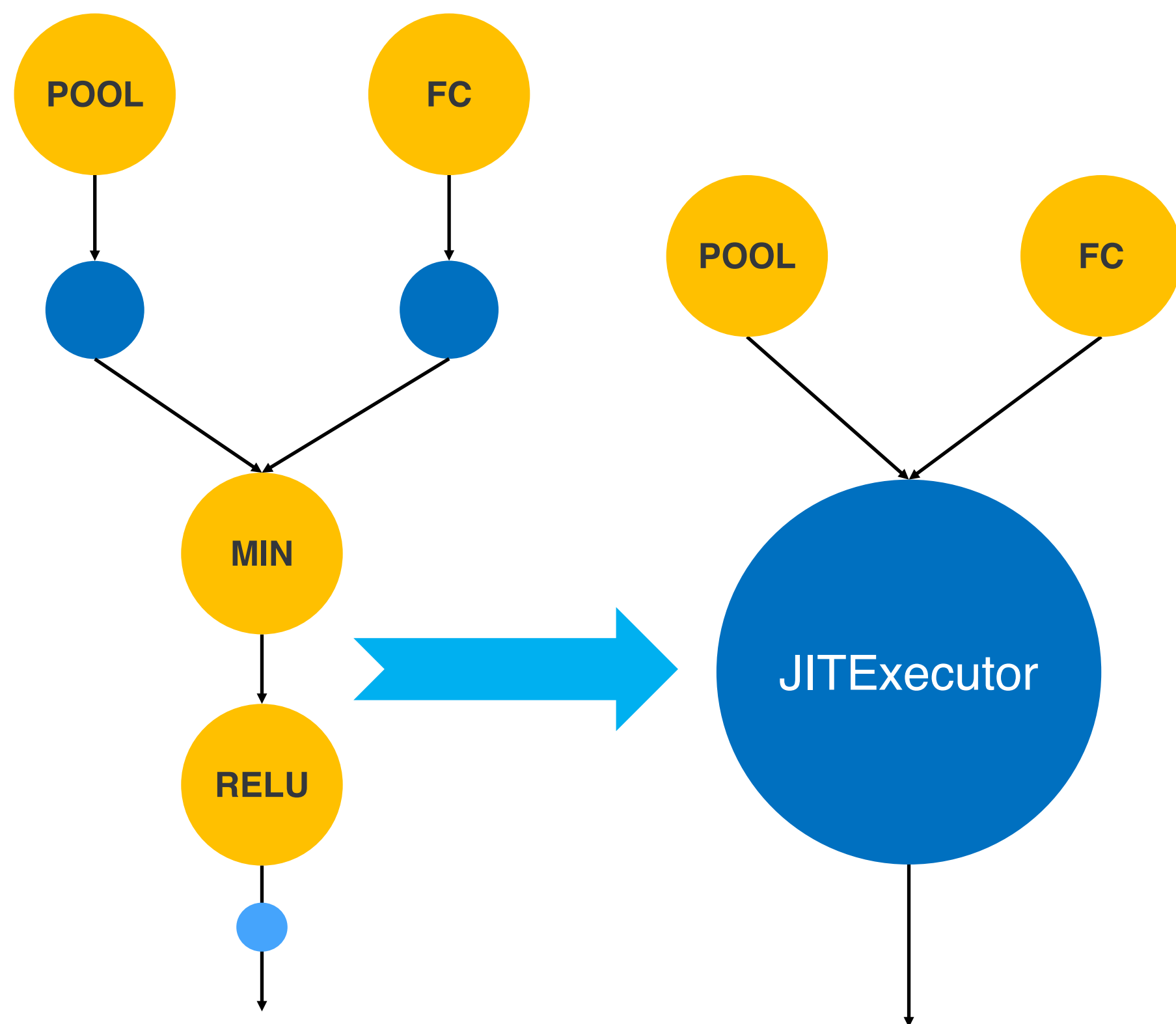
output : fuse 后的计算图 **G'**

1. 按照拓扑序列遍历图 **G** 中的算子 **opr**
2. 若 **opr** 的 **input var** 不是 **endpoint**, 返回步骤 1
3. 从 **M** 中拿到 **var** 对应的 **internal graph generator**, 生成 **internal graph**
4. 从 **internal graph** 创建 **JITExecutor**
5. 写回原始的计算图 **G**
6. 转到步骤 1

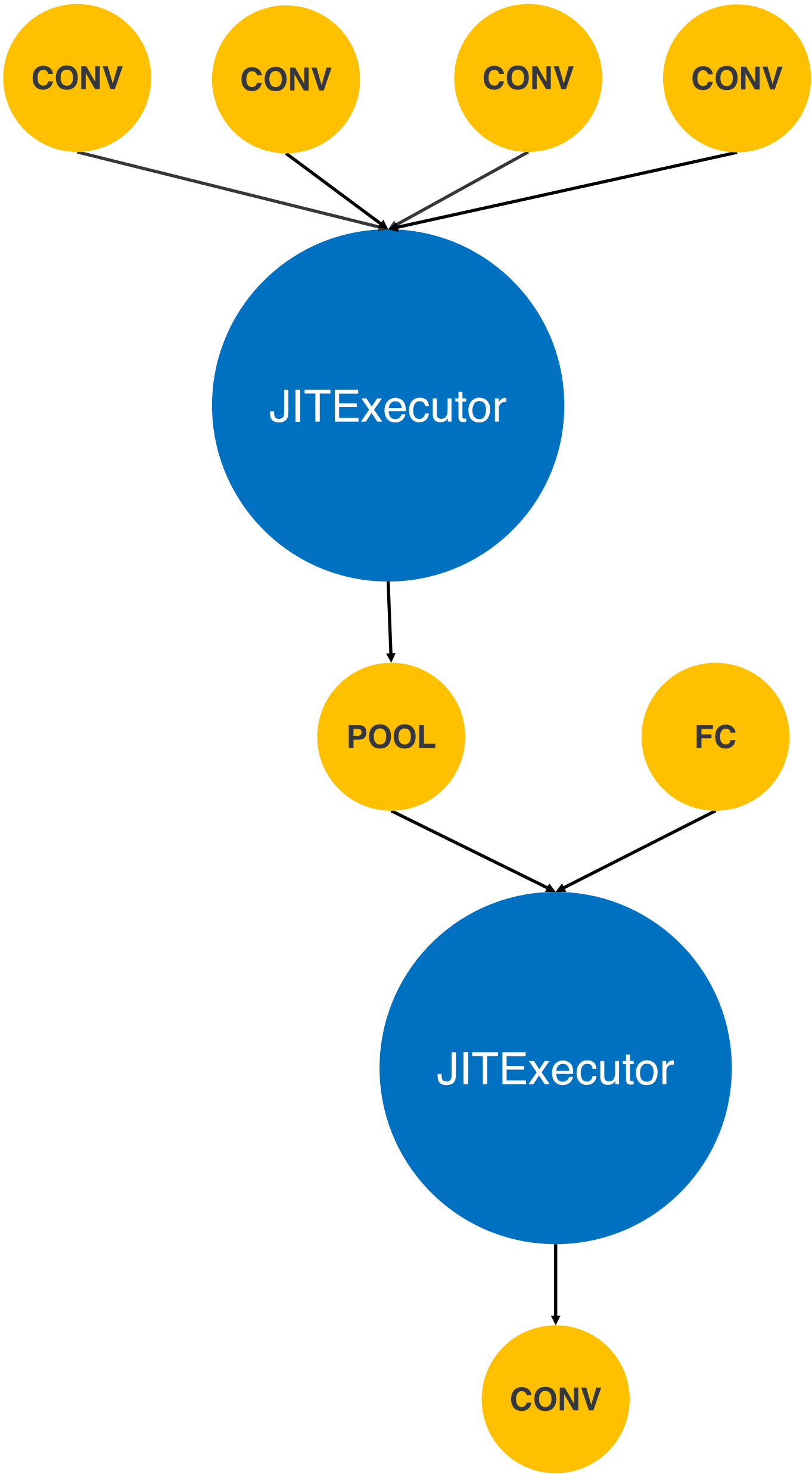
internal graph generator -> internal graph



internal graph -> JITExecutor

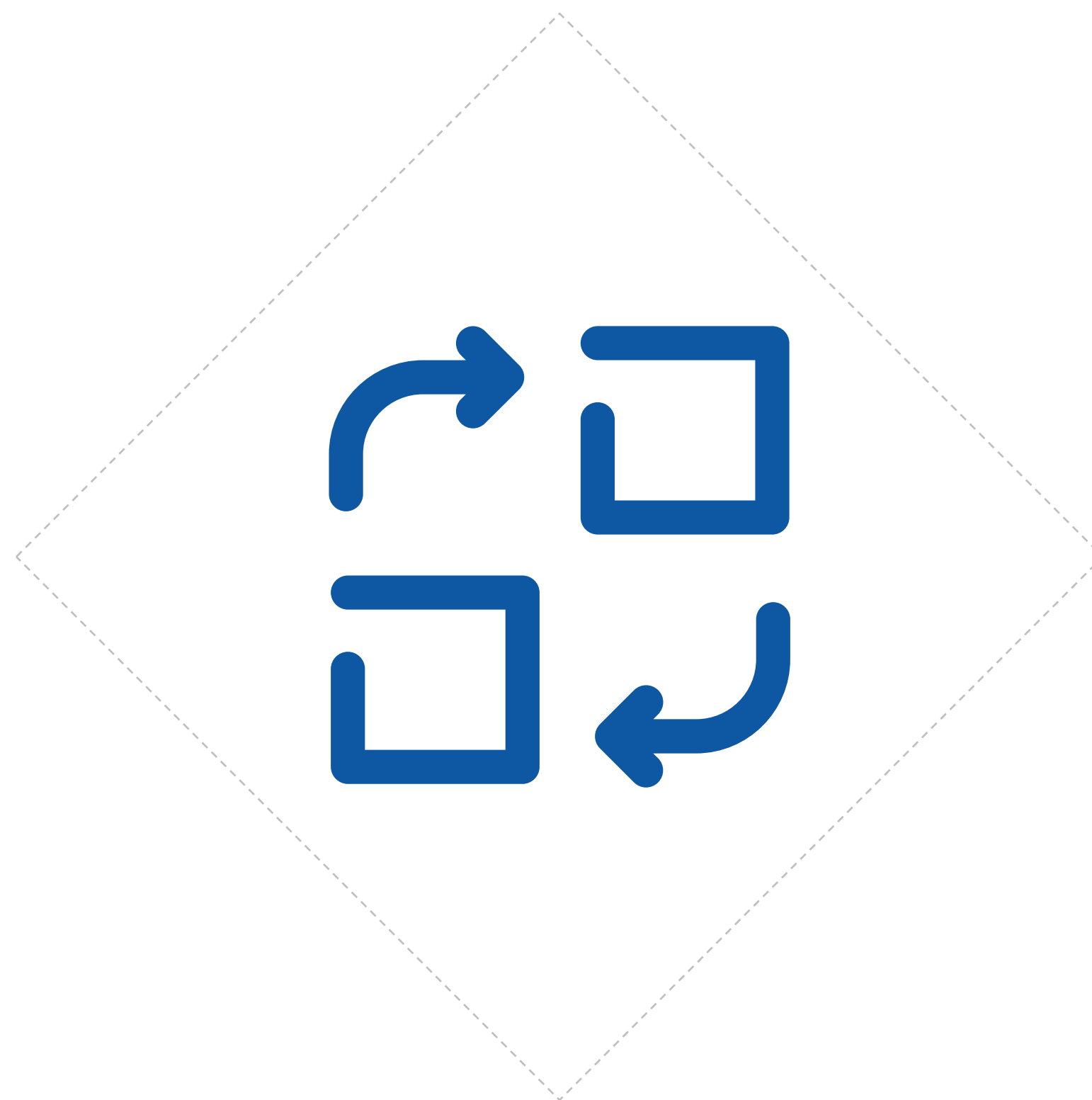


写回计算图



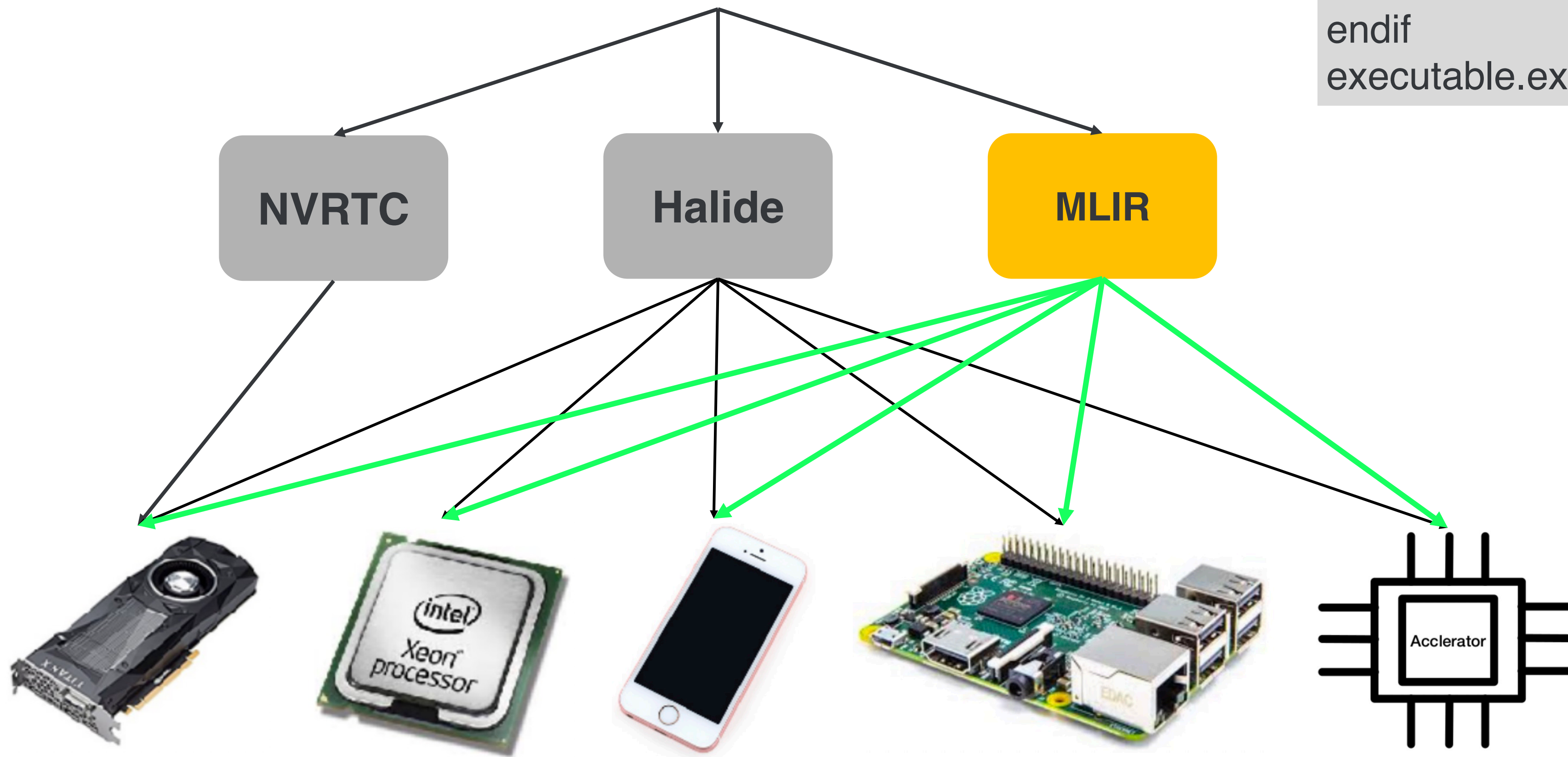
① Motivation

④ Evaluation



② Detect Fusion

③ JIT Compilation



Algorithm : JITExecutor::execute

```
if executable is not valid:  
    executable = compiler.compile()  
endif  
executable.execute()
```

MLIR(Multi-Level Intermediate Representation)

1. Operator level

```
%C = "mgb.matmul"(%A, %B)
      : (memref<2048x2048xf32>, memref<2048x2048xf32>) ->
      tensor<2048x2048xf32>
```

2. Loop level

```
affine.for %i = 0 to 2048 {
  affine.for %j = 0 to 2048 {
    affine.for %k = 0 to 2048 {
      %a = affine.load %A[%i, %k] : memref<2048x2048xf32>
      %b = affine.load %B[%k, %j] : memref<2048x2048xf32>
      %c = affine.load %C[%i, %j] : memref<2048x2048xf32>
      %p = mul %a, %b : f32
      %co = addf %c, %p : f32
      affine.store %co, %C[%i, %j] : memref<2048x2048xf32>
    }
  }
}
```

3. Low-level : 类似汇编语言

```
%a = load %A[%i2, %i3] : memref<2048x2048xf32>
%b = load %B[%i2, %i3] : memref<2048x2048xf32>
%c = addf %a, %b : f32
store %c, %C[%i2, %i3] : memref<2048x2048xf32>
```

MLIR 可以做：

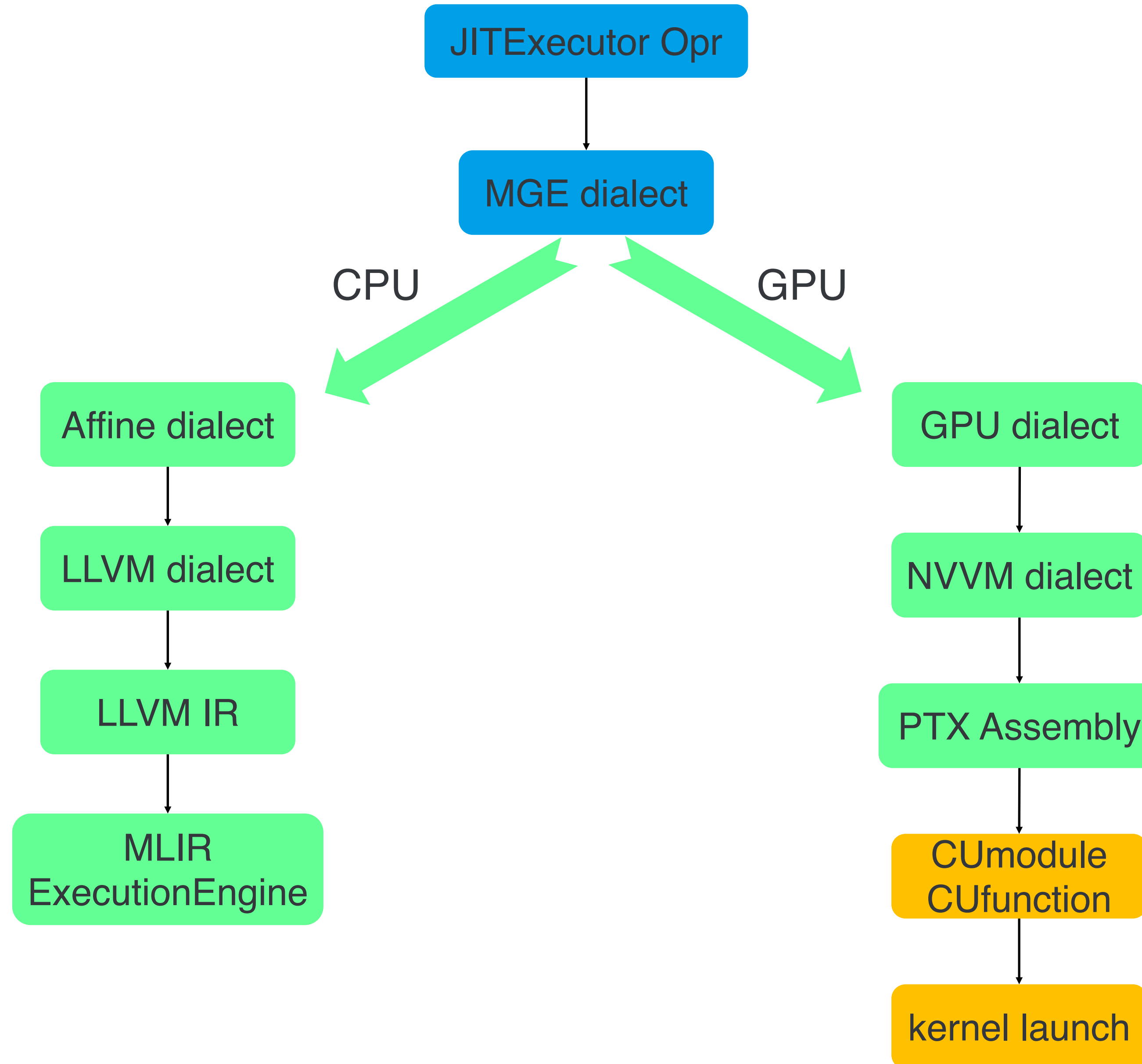
- 表达数据流图（如静态图模式下的 MegEngine Graph）
- 进行各种优化如算子融合（**kernel fusion**）、循环融合、分块和内存格式（**memory layout**）转换等
- 自动代码生成、自动向量化
- 表达适用于专用加速器的语言

Dialect：一个算子、类型、属性的抽象集合（ MGE dialect, Affine dialect, GPU dialect, LLVM dialect）

Lowering：dialect 之间的转换，例如 MGE dialect -> Affine dialect

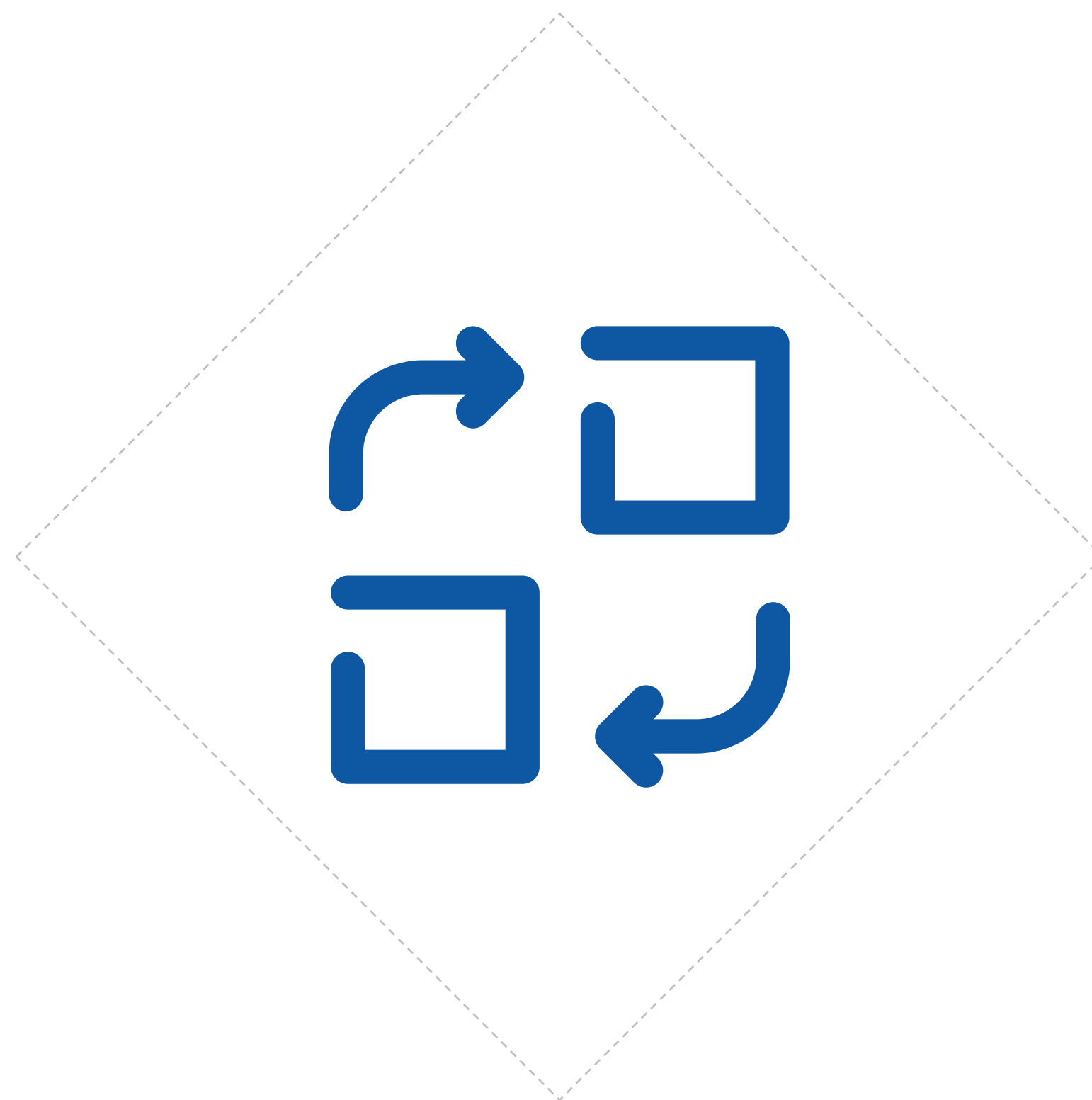
Lowering is all we need !

JIT Compilation



① Motivation

④ Evaluation



② Detect Fusion

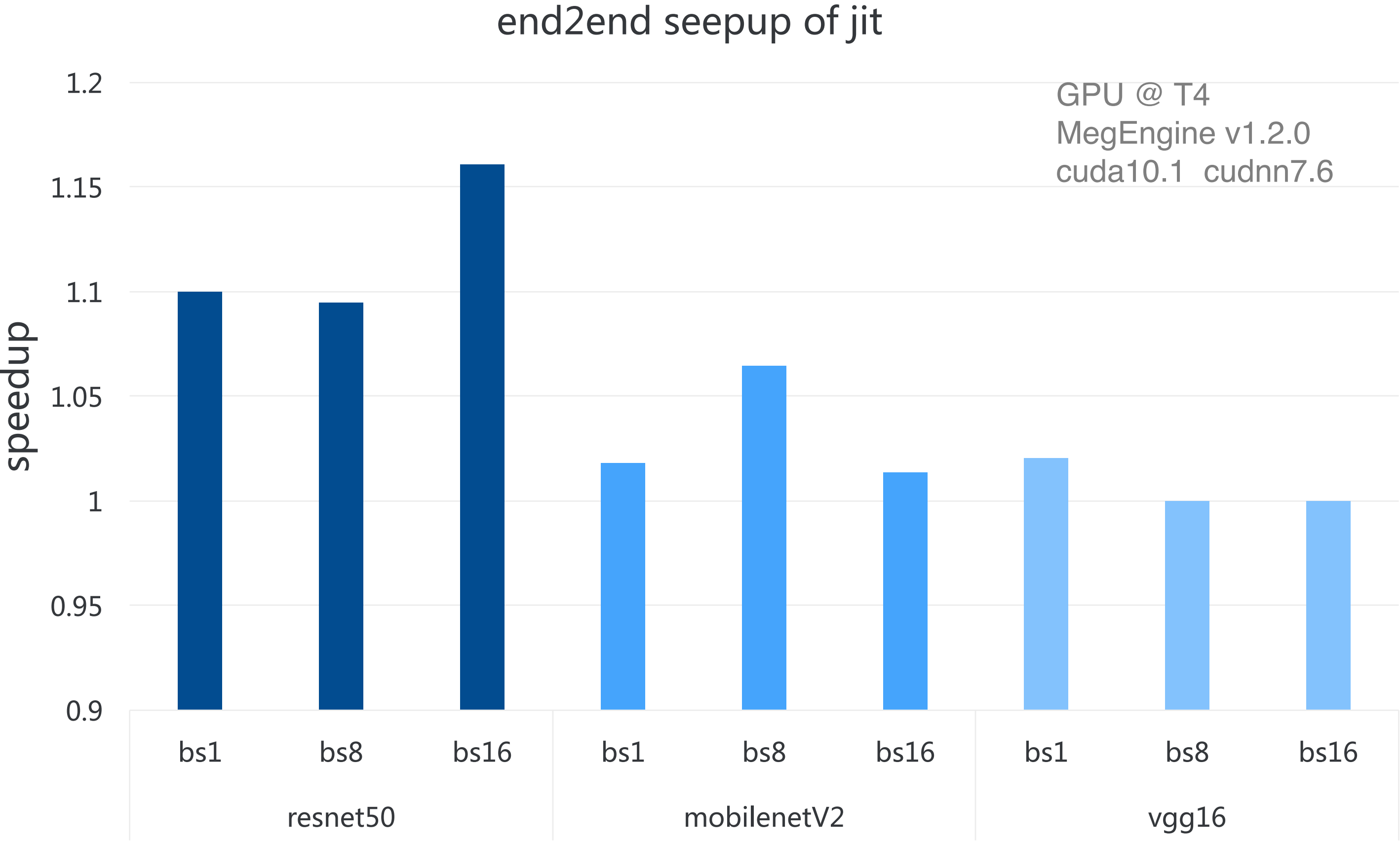
③ JIT Compilation

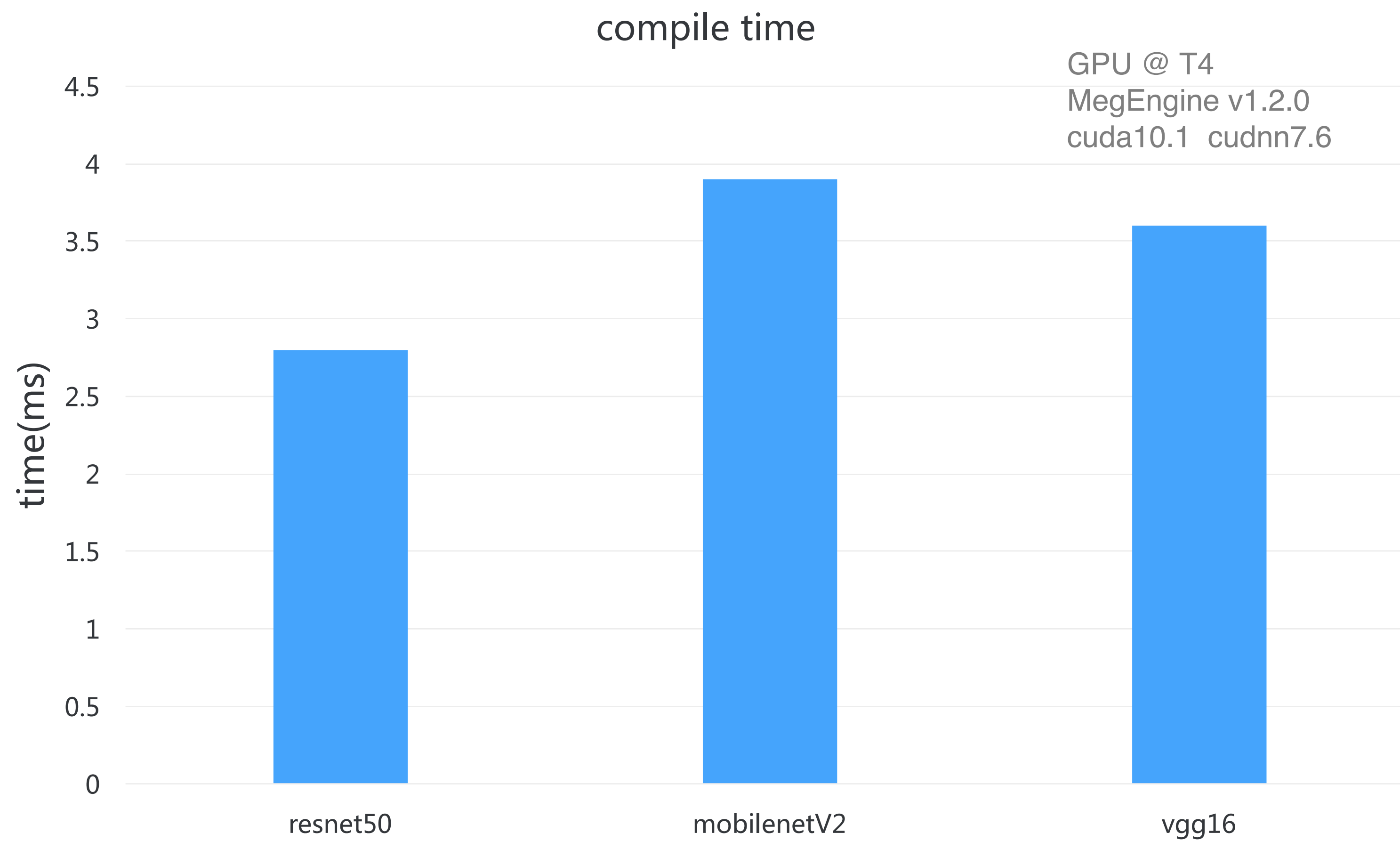
1. cmake 源码编译 MegEngine (仅 MLIR 需要)

2. 改 python 代码

```
if __name__ == '__main__':
    gm = ad.GradManager().attach(model.parameters())
    opt = optim.SGD(model.parameters(), lr=0.0125, momentum=0.9, weight_decay=1e-4,)
    # 通过 trace 转换为静态图
    @trace(symbolic=True, opt_level=3)
    def train():
        with gm:
            logits = model(image)
            loss = F.loss.cross_entropy(logits, label)
            gm.backward(loss)
            opt.step()
            opt.clear_grad()
            return loss
    loss = train()
    loss.numpy()
```

3. **export** MGB_JIT_BACKEND= "MLIR" / "NVRTC" / "HALIDE"





- 将 jit 编译的结果离线保存，解决线上首次运行编译慢的问题
- jit 支持更多的算子
- jit 支持更多数据类型
- 动态图 jit（真jit），即检测热点代码并编译

MEGVII 旷视



天元
MegEngine



天元开发者交流群

群号: 1029741705



扫一扫二维码，加入群聊。