

Automated GPU Kernel Fusion with XLA

EuroLLVM'19, April 8 2019



Thomas Joerg, Google
Presenting work done by the XLA team

Outline

- TensorFlow
- Kernel fusion
- XLA compiler
- Automated kernel fusion



Fit to screen

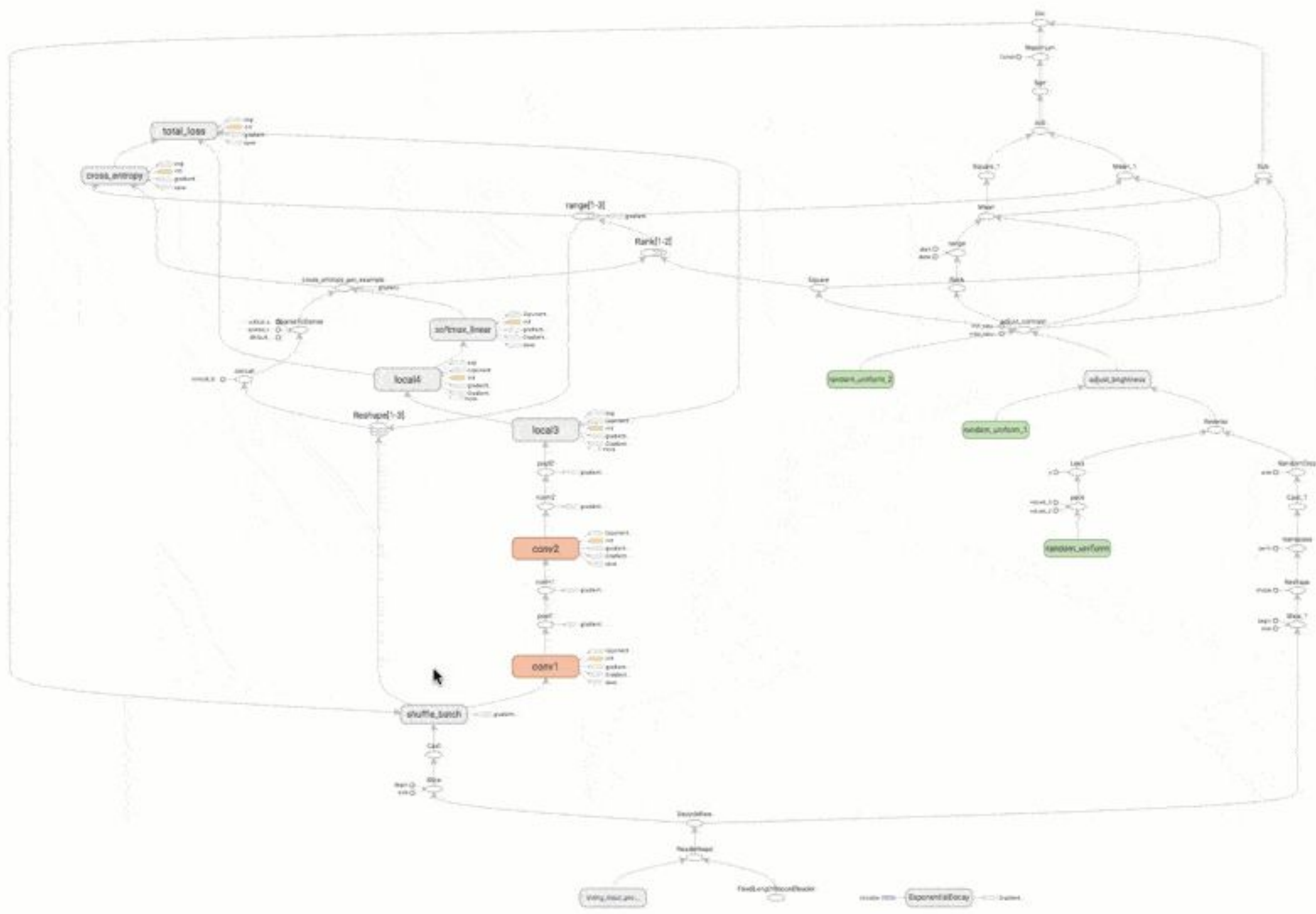
Run cifar-train

Upload Choose File

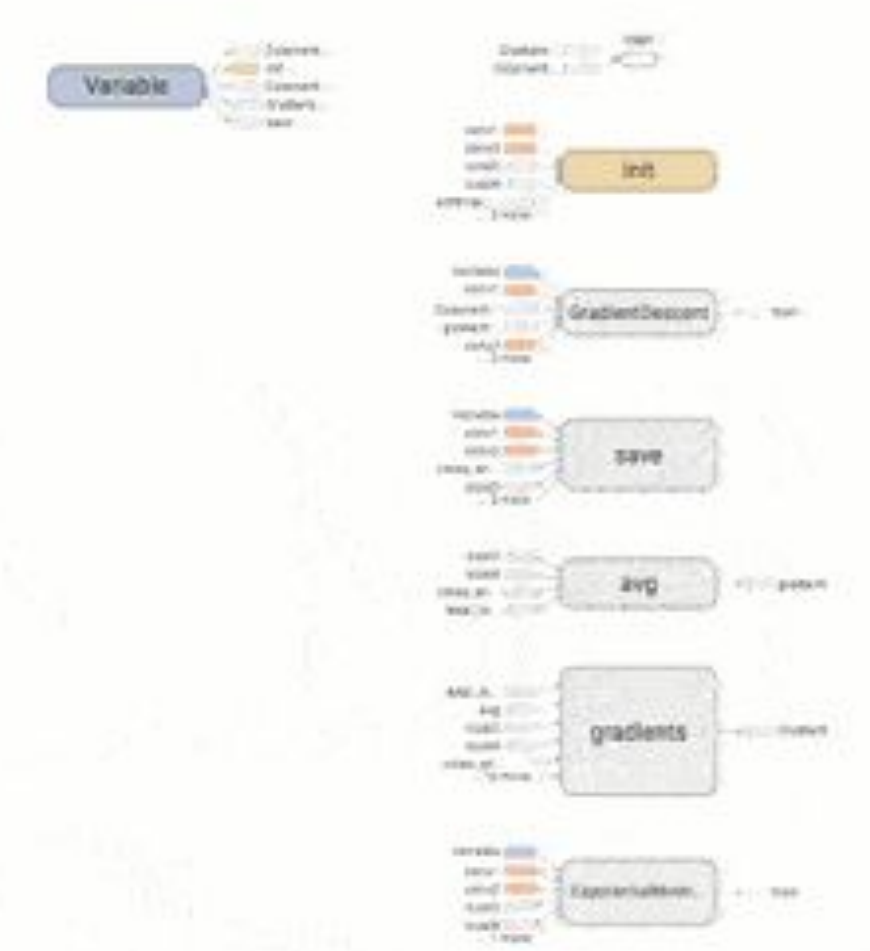
Color Structure

color: same substructure
gray: unique substructure

Main Graph



Auxiliary nodes



Graph (* = expandable)

- Namespace*
- OpNode
- Unconnected series*
- Connected series*
- Constant
- Summary
- Dataflow edge
- Control dependency edge
- Reference edge

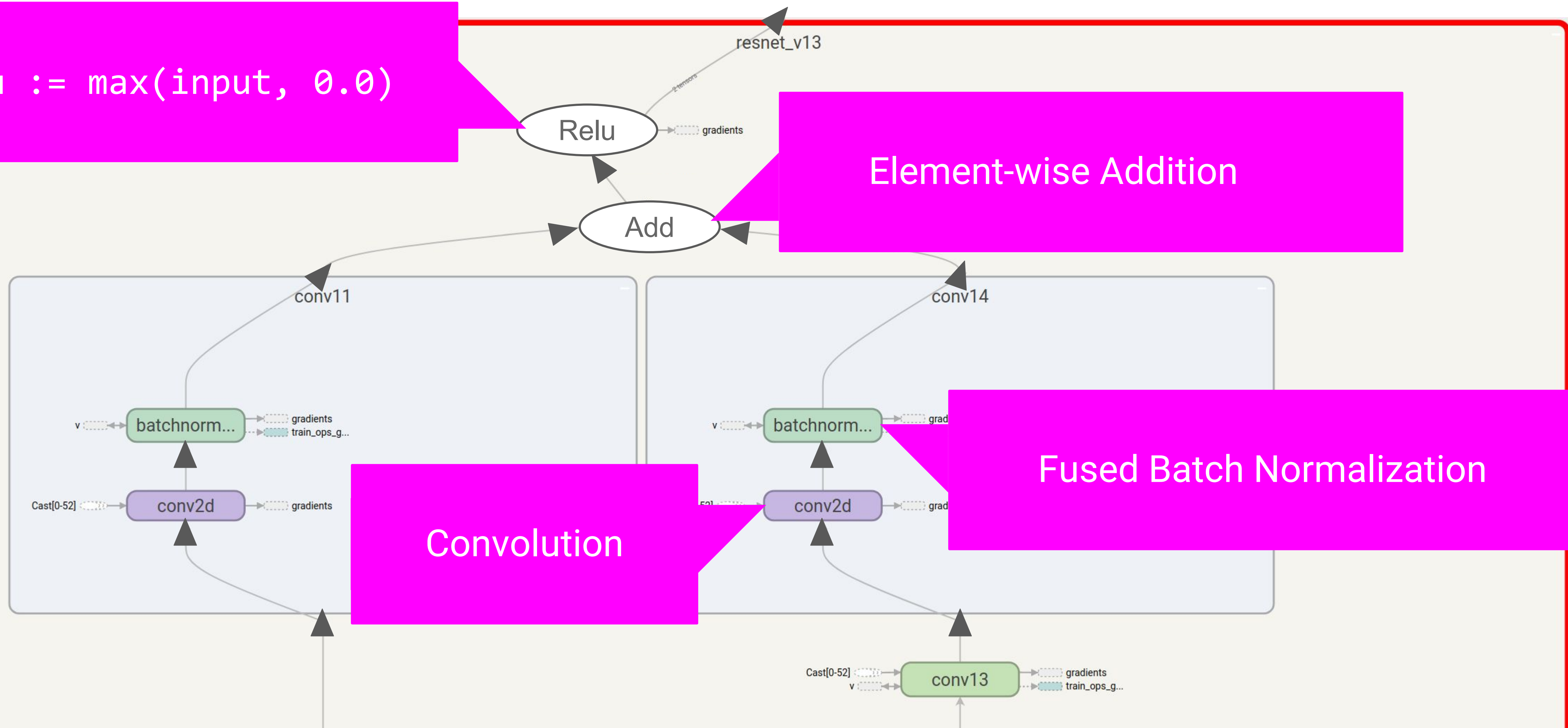
Example: ResNet block

ReLU := $\max(\text{input}, 0.0)$

Element-wise Addition

Fused Batch Normalization

Convolution



Fused Kernels

- Convenient
- Performant

```
// Compute  $a * x + y$ .  
// a is a scalar, x and y are tensors.  
tmp = tf.multiply(a, x)  
out = tf.add(tmp, y)
```

```
// Compute  $a * x + y$ .  
// a is a scalar, x and y are tensors.  
tmp = tf.multiply(a, x)  
out = tf.add(tmp, y)
```

```
__global__ void Multiply(int n, float a, float* x) {  
    int i = blockIdx.x * blockDim.x + threadIdx.x;  
    if (i < n) x[i] = a * x[i];  
}
```

```
// Compute a * x + y.  
// a is a scalar, x and y are tensors.  
tmp = tf.multiply(a, x)  
out = tf.add(tmp, y)
```

Tensors read + written: 4

```
__global__ void Multiply(int n, float a, float* x) {  
    int i = blockIdx.x * blockDim.x + threadIdx.x;  
    if (i < n) x[i] = a * x[i];  
}  
__global__ void Add(int n, float* x, float* y) {  
    int i = blockIdx.x * blockDim.x + threadIdx.x;  
    if (i < n) x[i] = x[i] + y[i];  
}
```



```
// Compute  $a * x + y$ .  
// a is a scalar, x and y are tensors.  
tmp = tf.multiply(a, x)  
out = tf.add(tmp, y)
```

```
__global__ void FusedMulAdd(int n, float a, float* x, float* y) {  
    int i = blockIdx.x * blockDim.x + threadIdx.x;  
    if (i < n) x[i] = a * x[i] + y[i];  
}
```

```
// Compute a * x + y.  
// a is a scalar, x and y are tensors.  
out = tf.fused_multiply_add(a, x, y)
```

Tensors read + written: 3
25% reduction!

```
__global__ void FusedMulAdd(int n, float a, float* x, float* y) {  
    int i = blockIdx.x * blockDim.x + threadIdx.x;  
    if (i < n) x[i] = a * x[i] + y[i];  
}
```

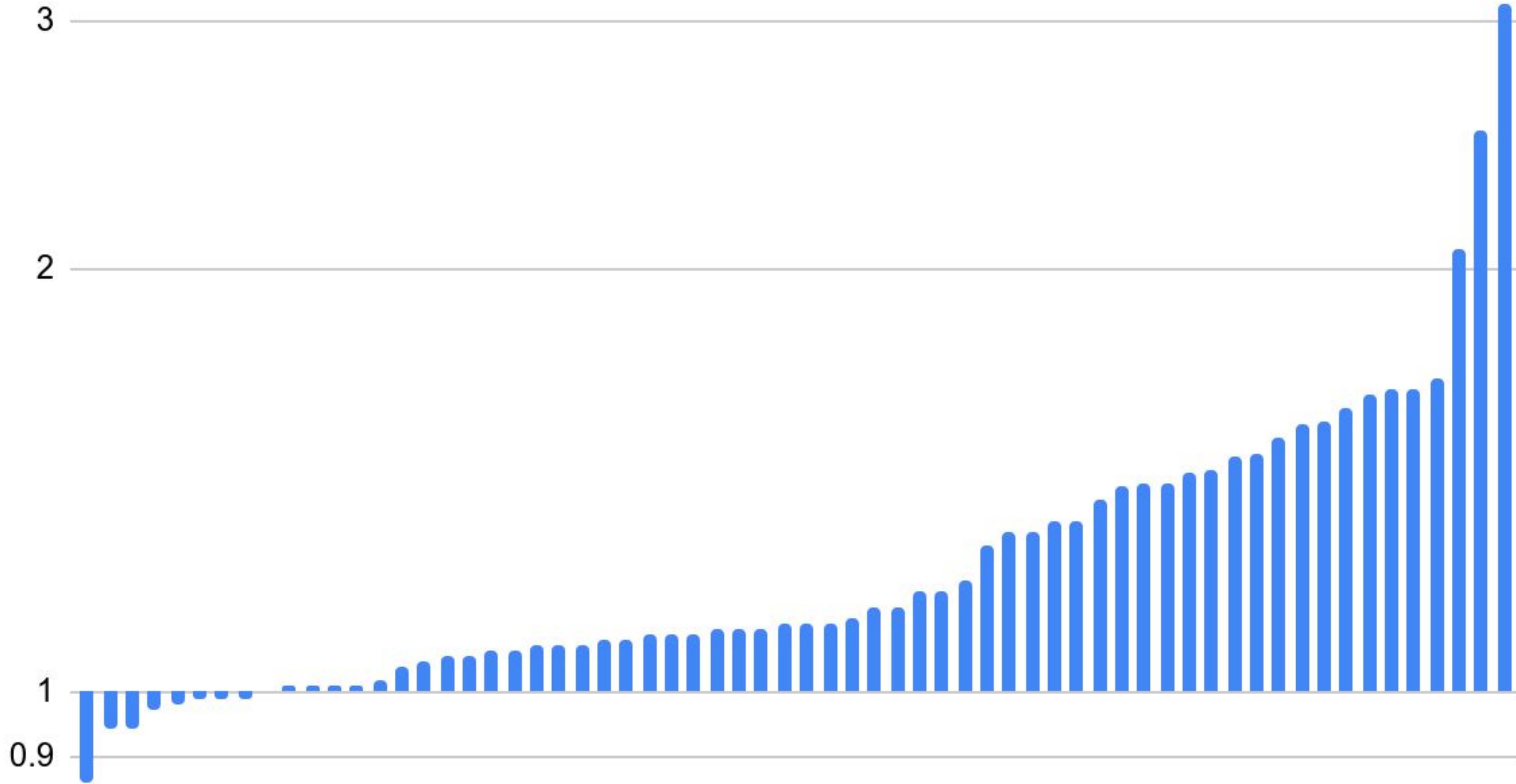
Fused Kernels

- Convenient
- Performant

But

- Development cost
- Inflexibel
- Hard to optimize

Speedup of TF with XLA vs TF without XLA





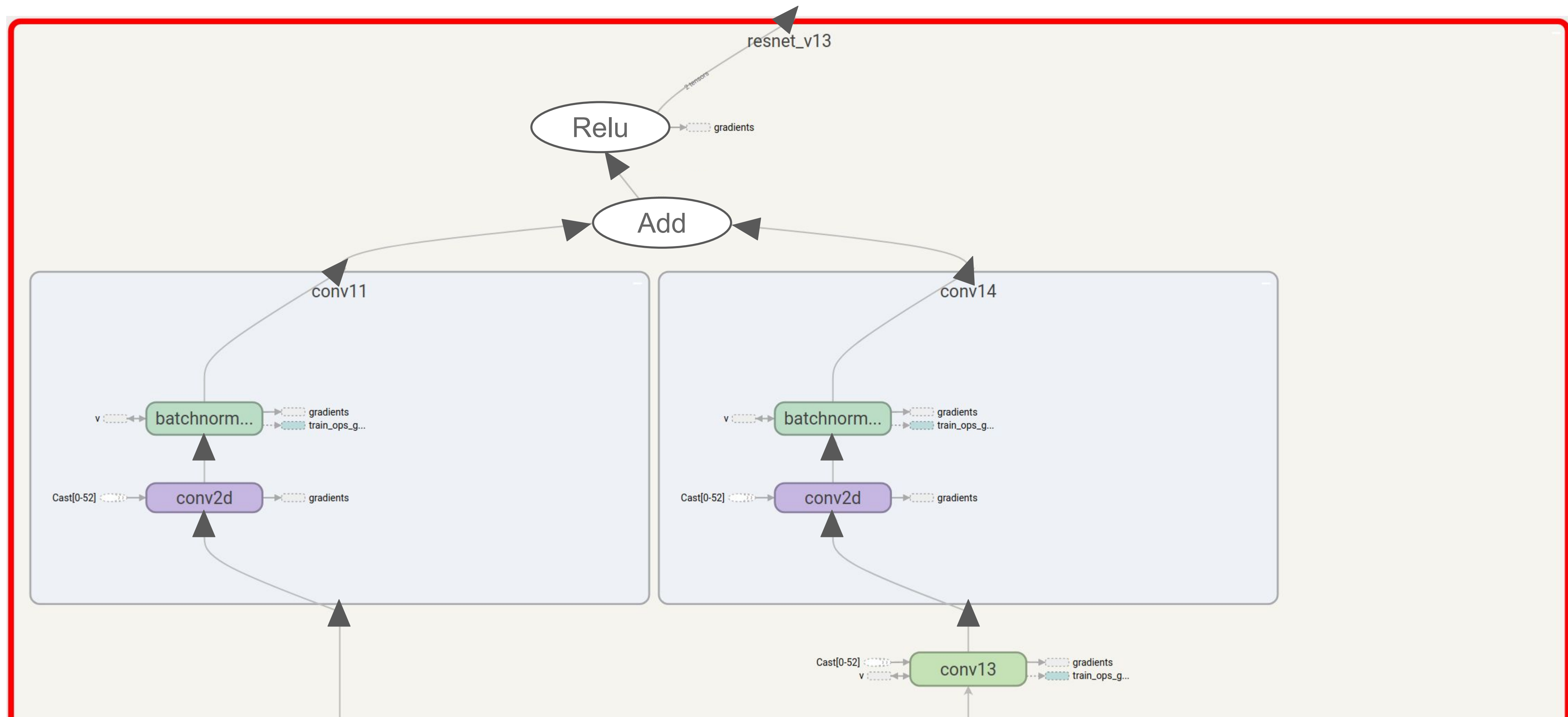
MLPerf

Submitter	Hardware	Chip count	Software	ResNet-50 v1.5 *
NVIDIA	DGX-1 (on premise)	8	ngc18.11_MXNet, cuDNN 7.4	65.6
Google	8x Volta V100 (Cloud)	8	TF 1.12, cuDNN 7.4	64.1

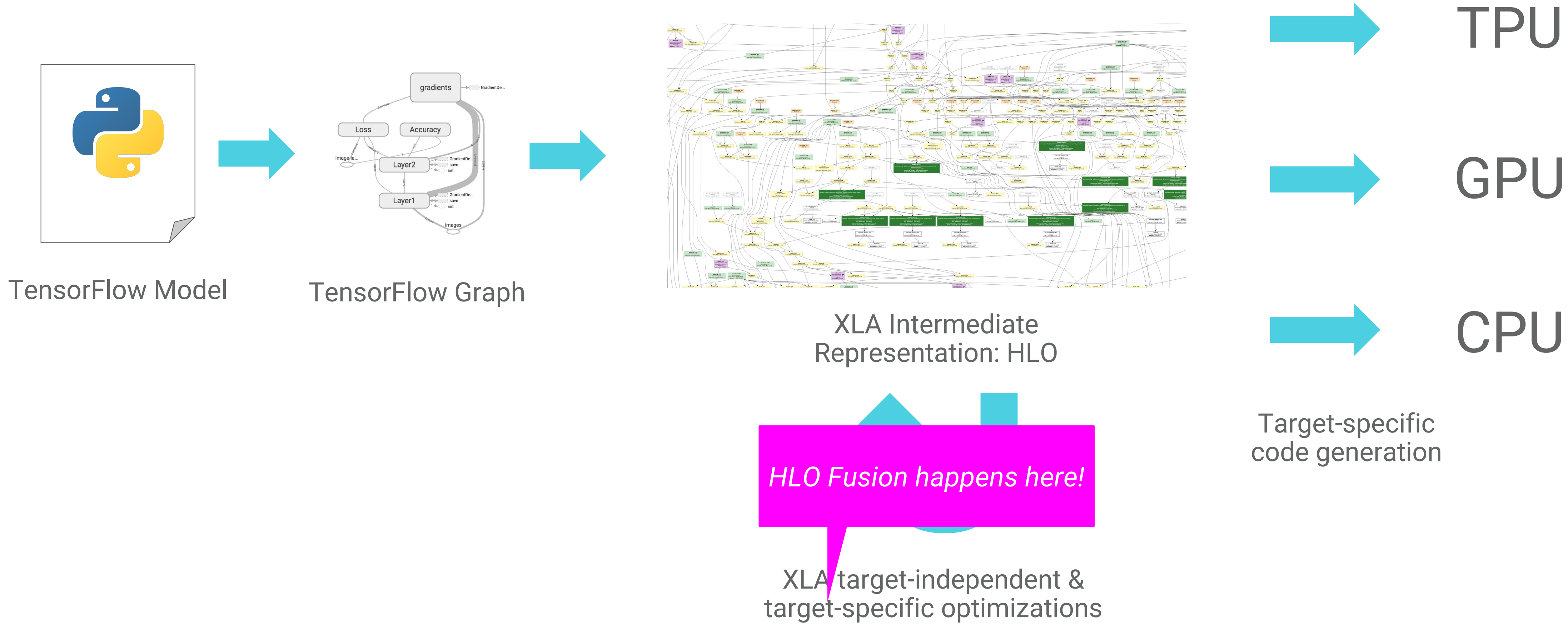
Full results: <https://mlperf.org/results/>

* speedup relative to reference implementation

Example: ResNet block



TensorFlow with XLA



HLO IR

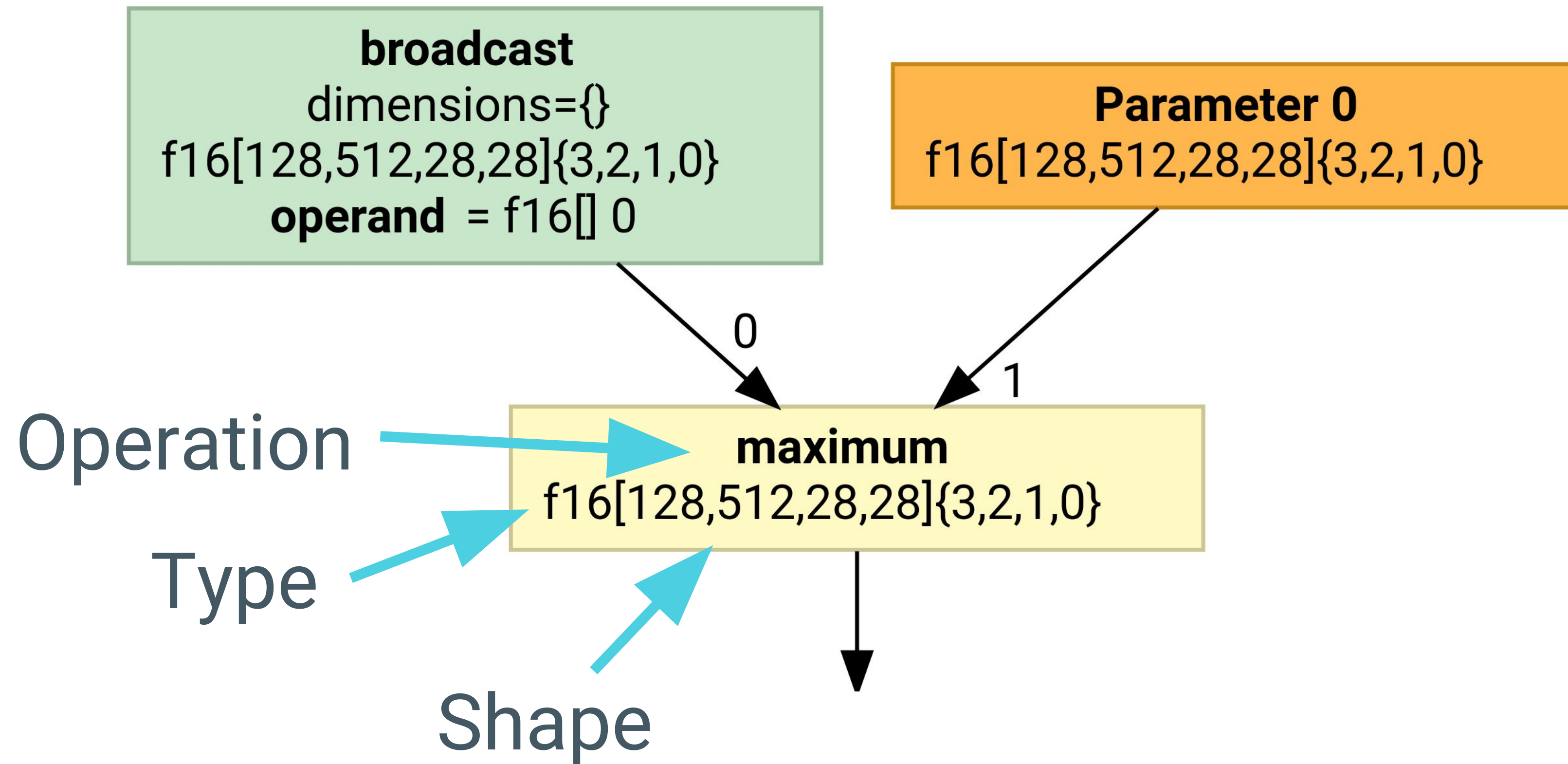
Sample HLO ops

- Elementwise math
 - Add, Tanh, Map
- Specialized math for neural nets
 - Dot, Convolution, Reduce
- Re-organize data
 - Reshape, Broadcast, Concat, Tuple
- Control flow
 - While, Call, CustomCall
- Data transfer
 - Parameter, Constant

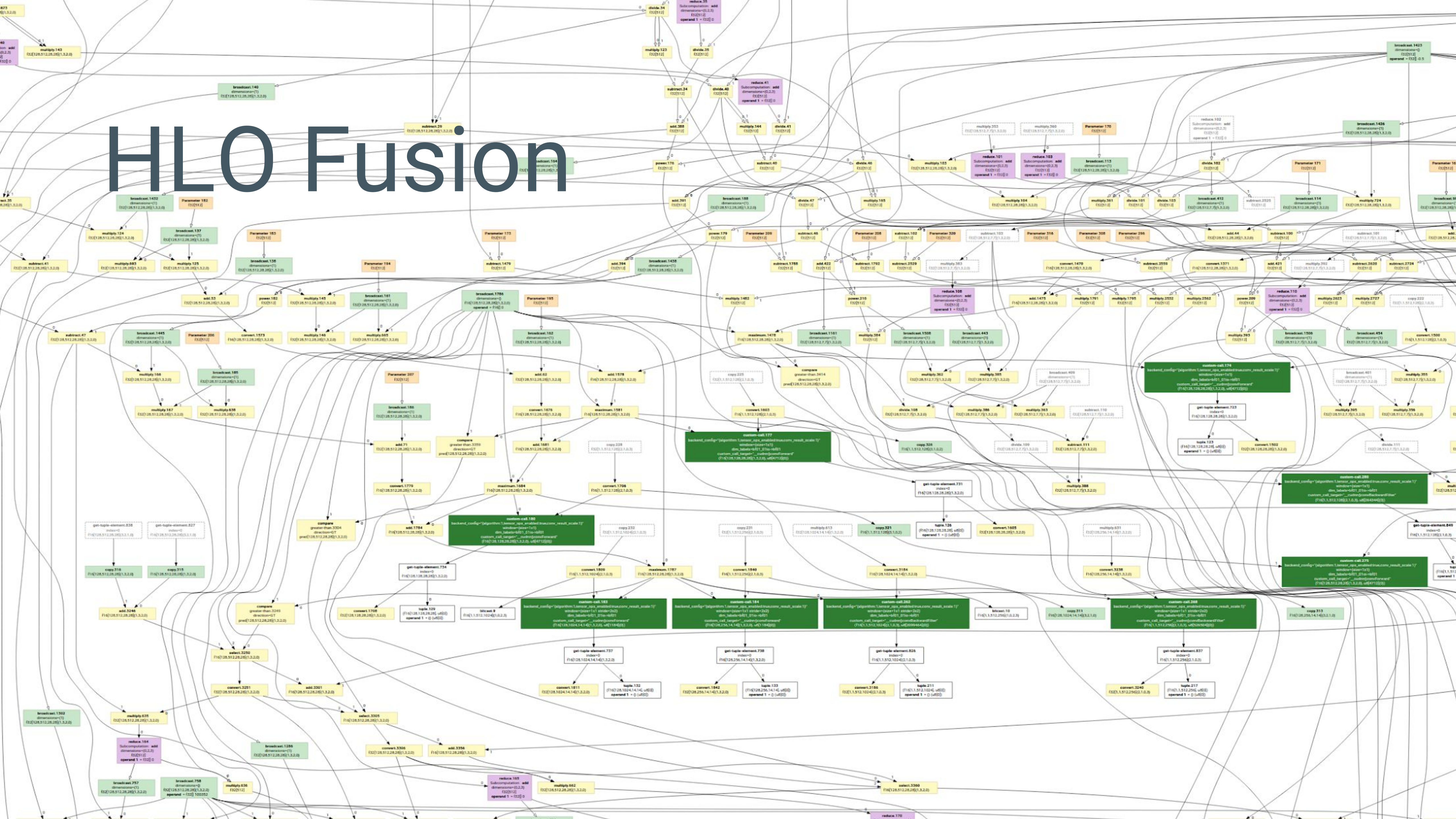
Sample data types

- Primitive types
 - PRED
 - F16
 - F32
- Composite types
 - array: F32[2,3], F16[]
 - tuple: TUPLE(F32[16], F16)

ReLu in HLO



HLO Fusion

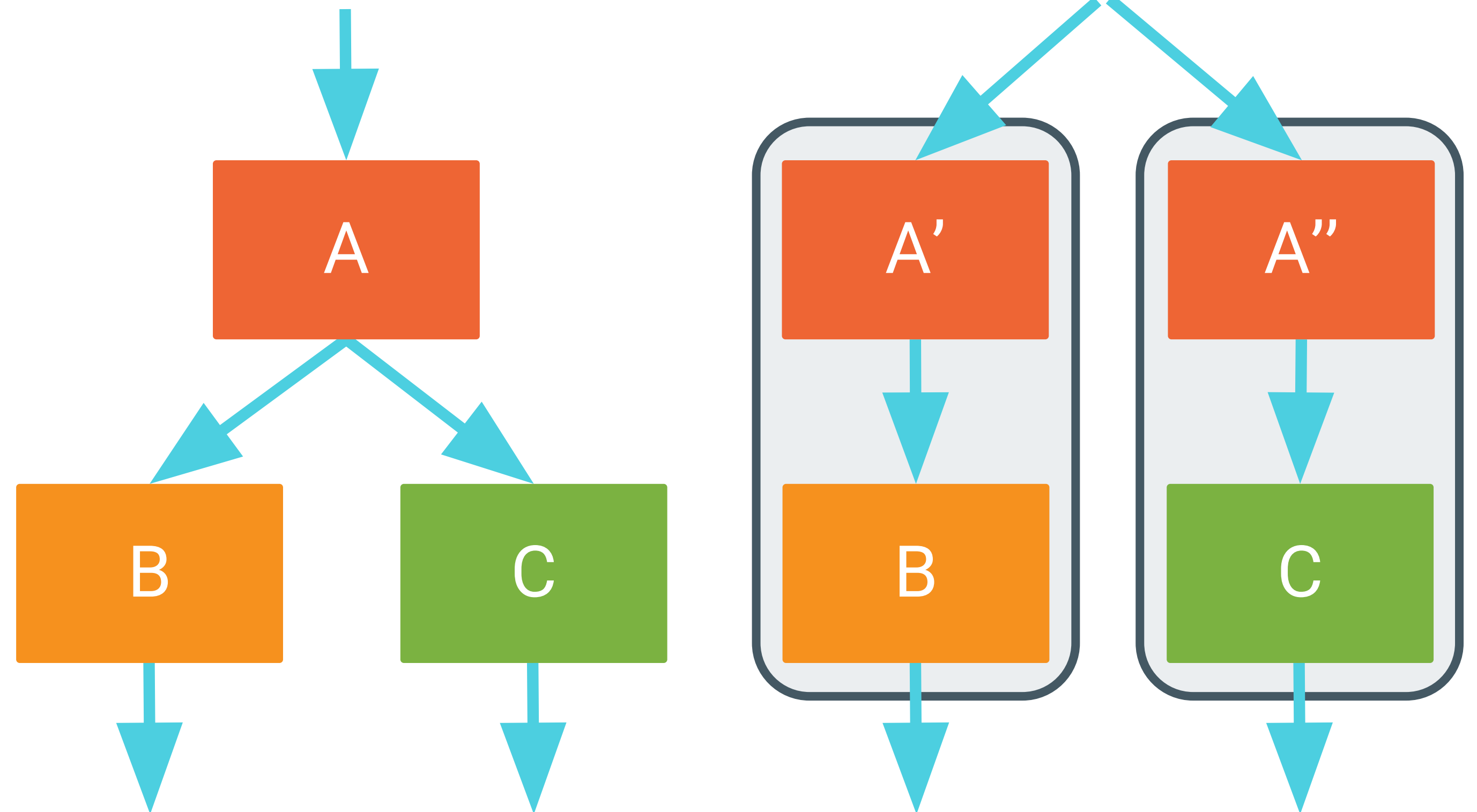


HLO Fusion

- Reduce memory bandwidth
- Compatible loop pattern
- Coalesced memory access

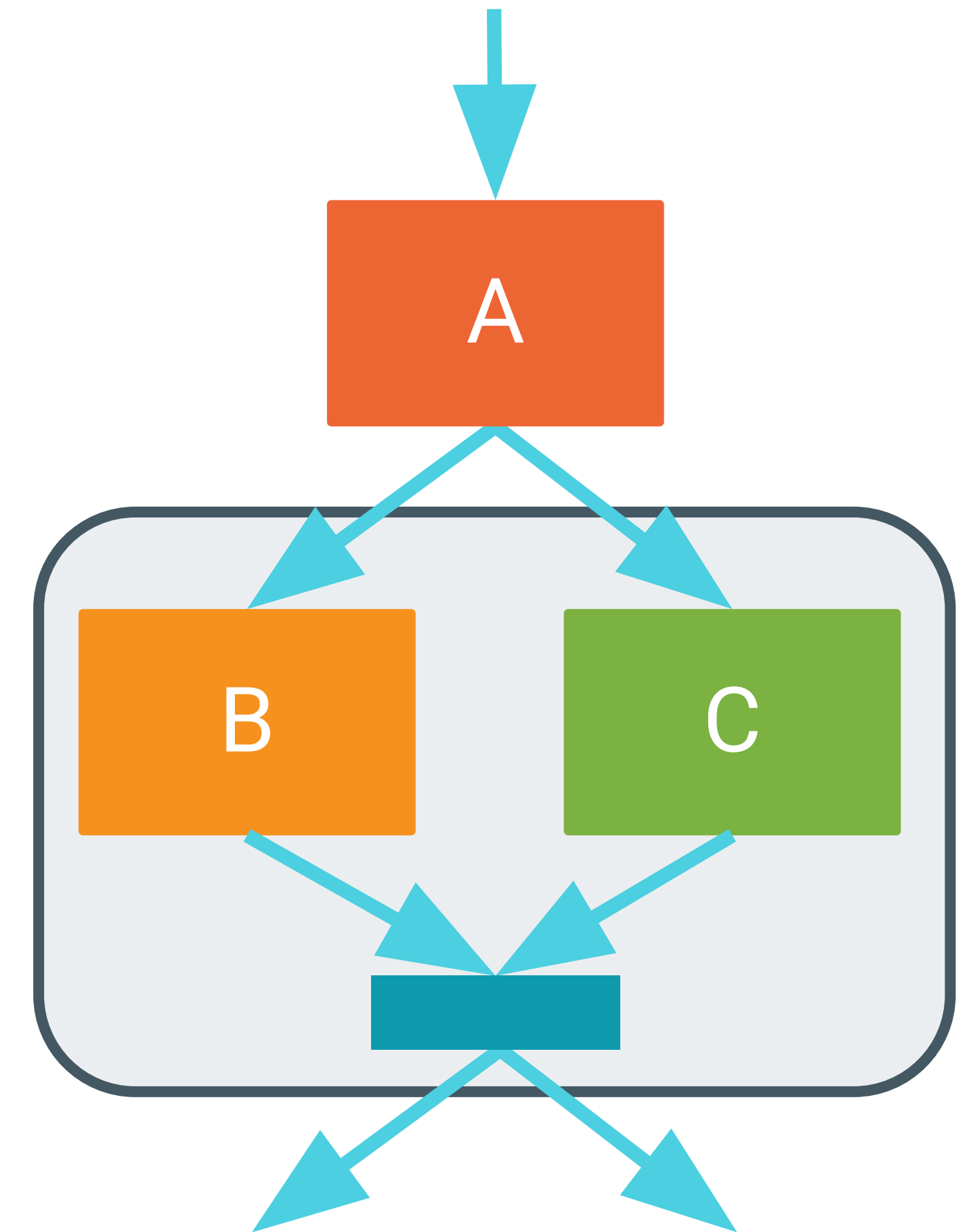
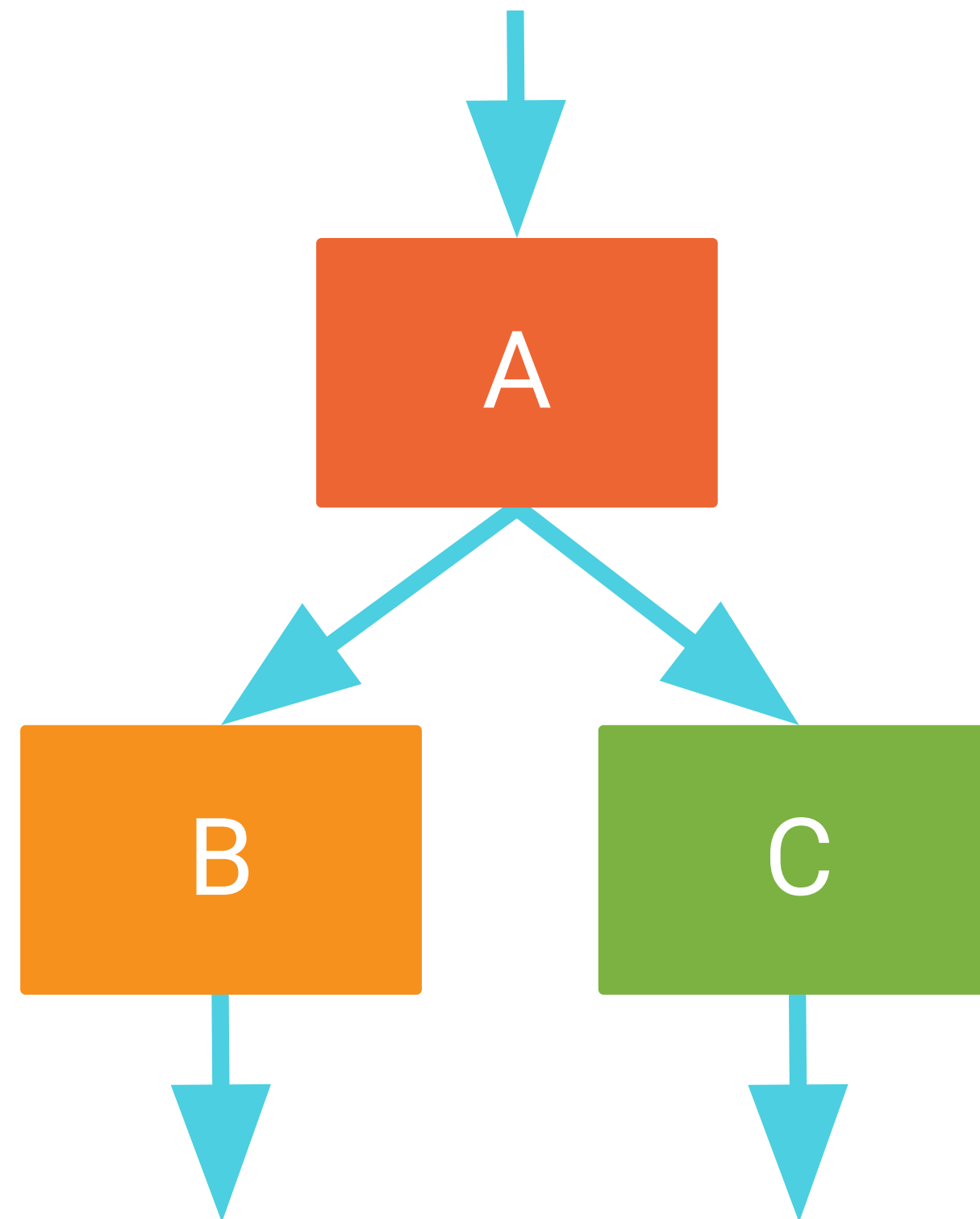
HLO Fusion

- 1) Fusion (with duplication)
- 2) Sibling fusion
- 3) Fusion with multiple outputs



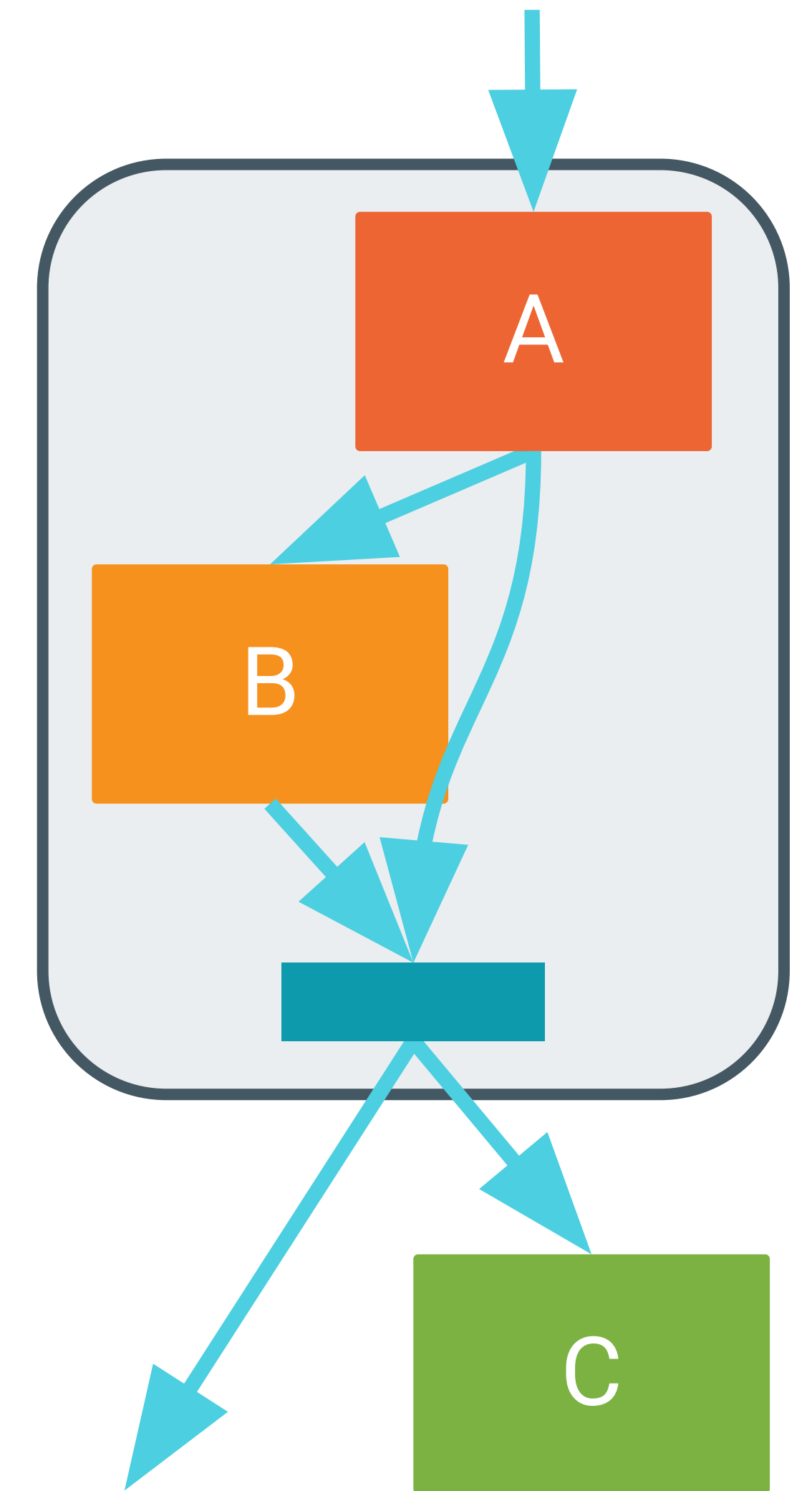
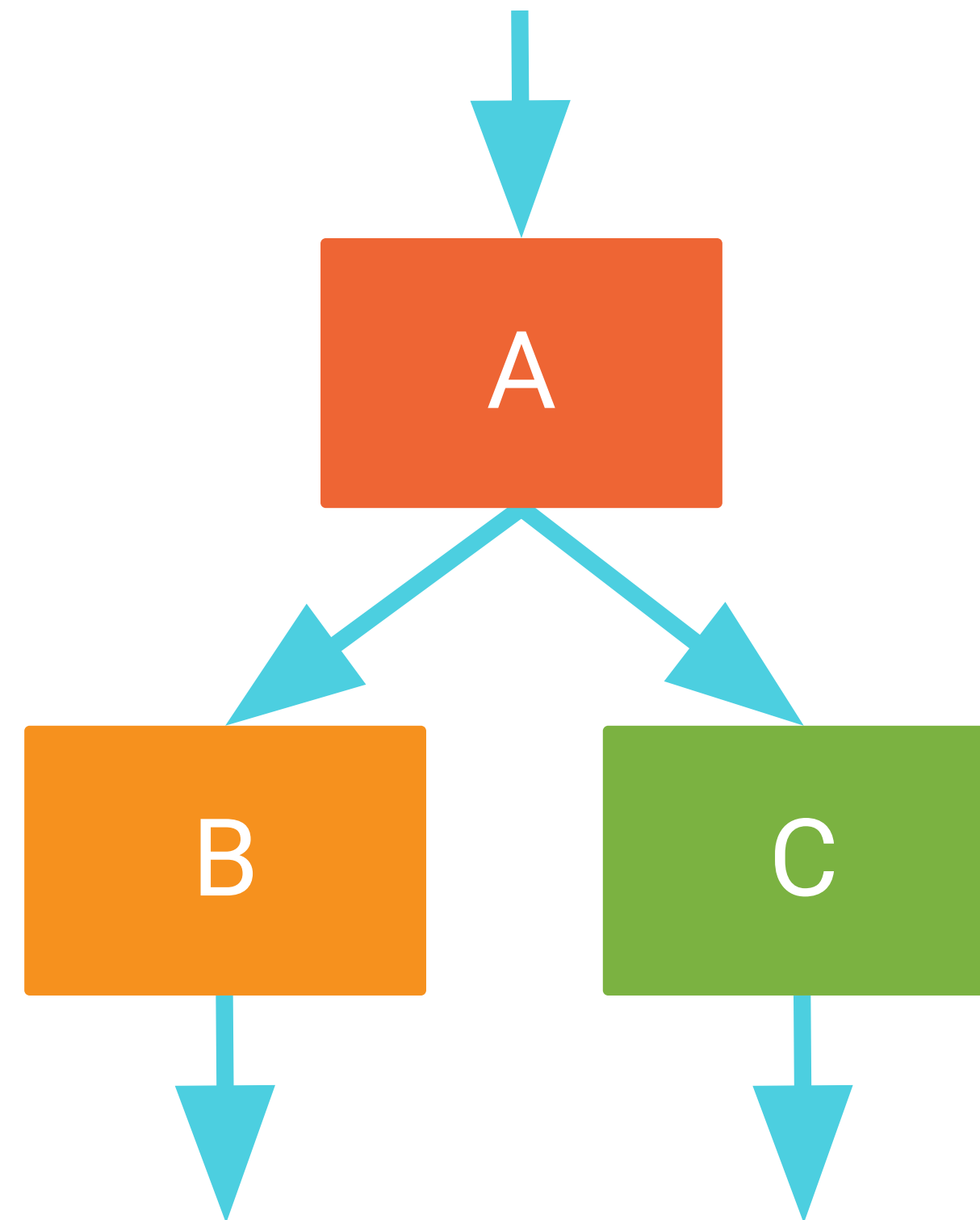
HLO Fusion

- 1) Fusion (with duplication)
- 2) Sibling fusion
- 3) Fusion with multiple outputs

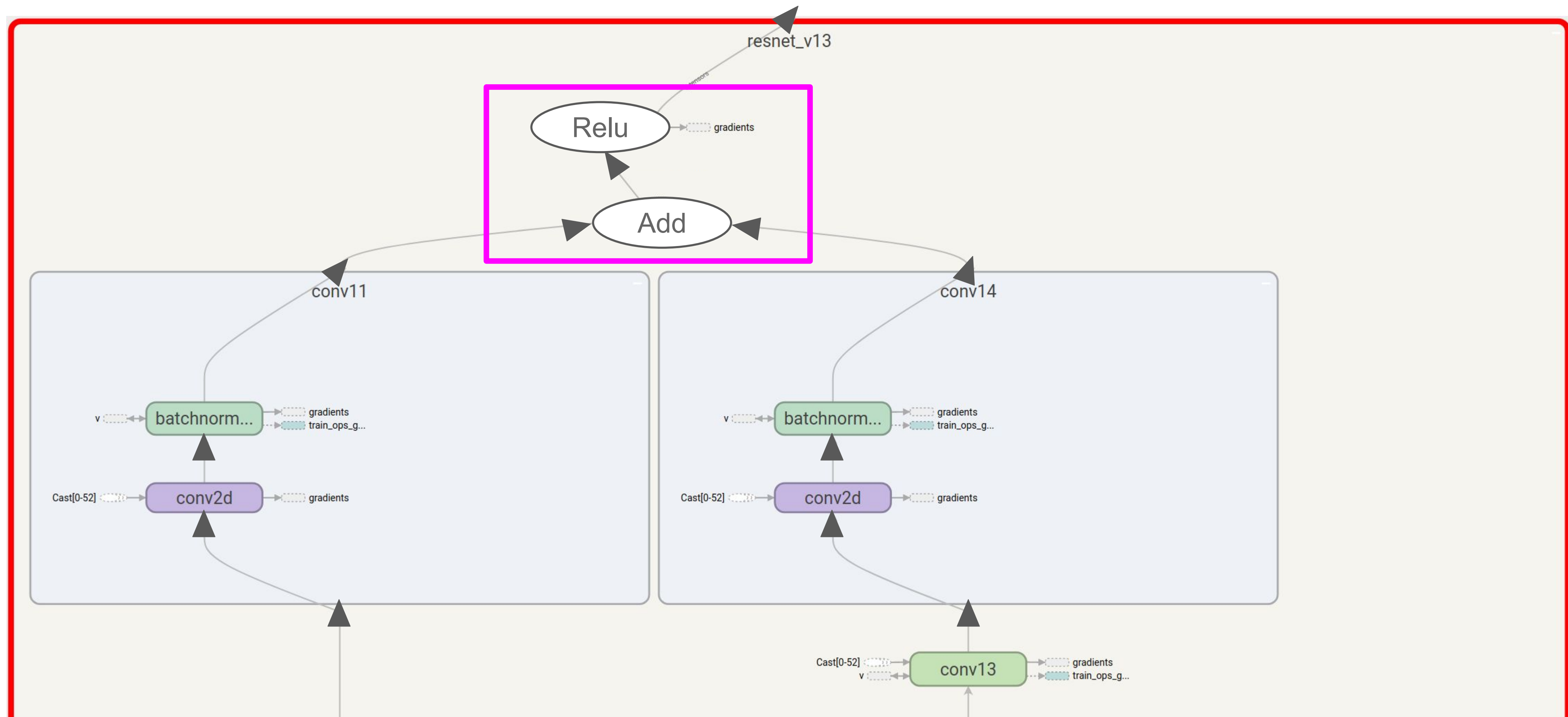


HLO Fusion

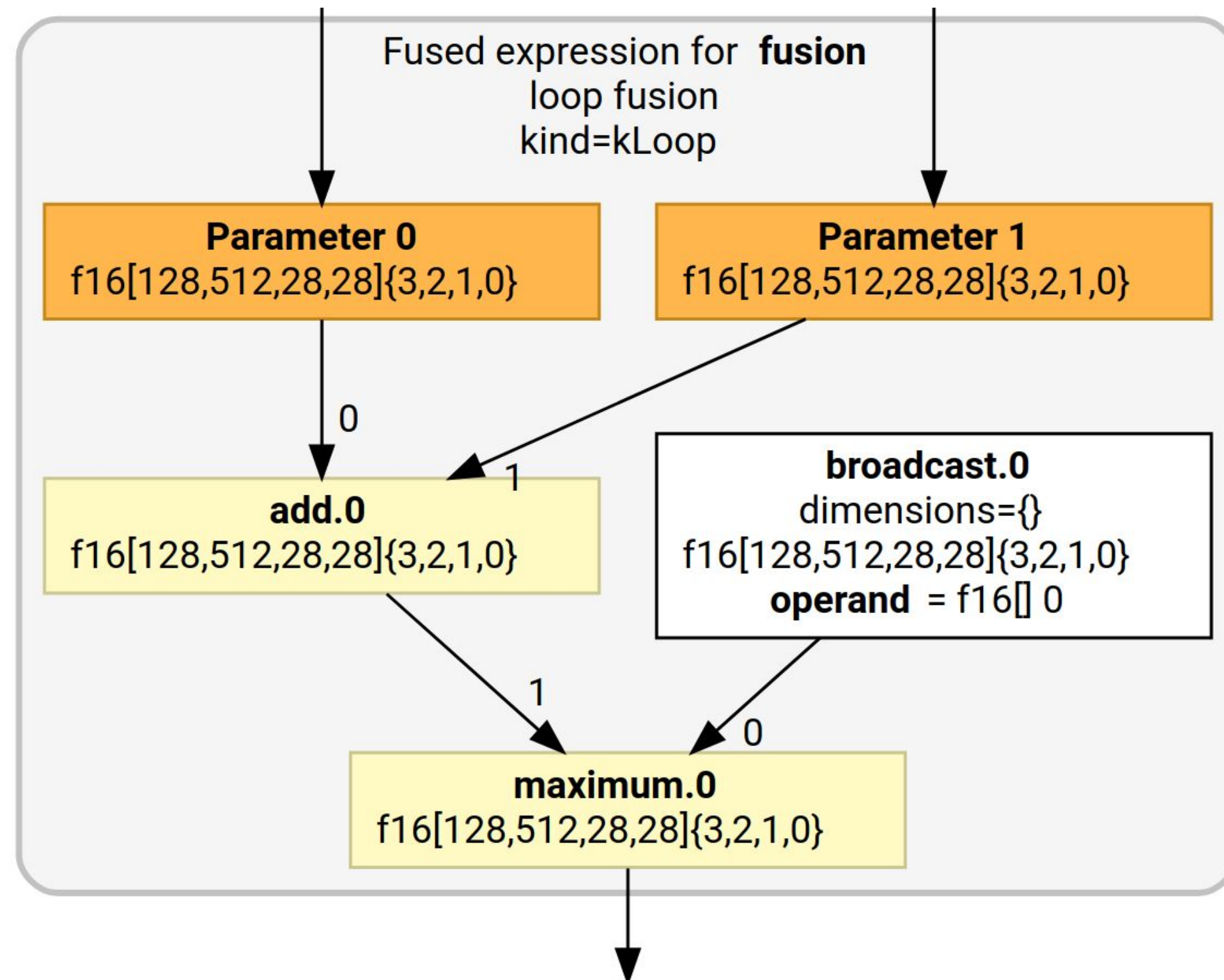
- 1) Fusion (with duplication)
- 2) Sibling fusion
- 3) Fusion with multiple outputs



Example: ResNet block



Fused Add + ReLu

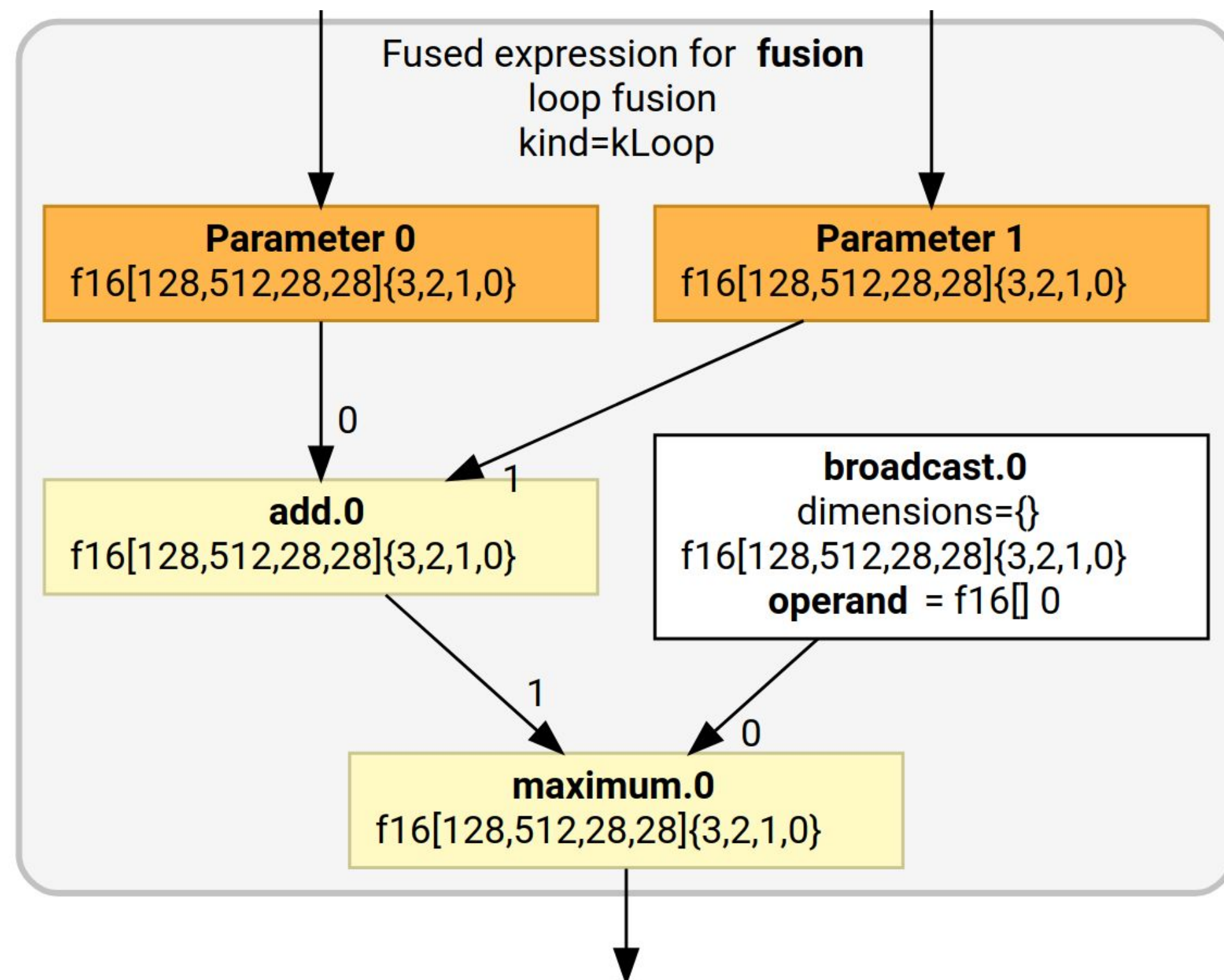


```
__global__ void fusion(float *lhs,  
float *rhs, float* output) {  
  
int i = blockIdx.x * blockDim.x +  
threadIdx.x;  
  
if (i < 128*512*28*28) {  
output[i] =  
}  
}
```

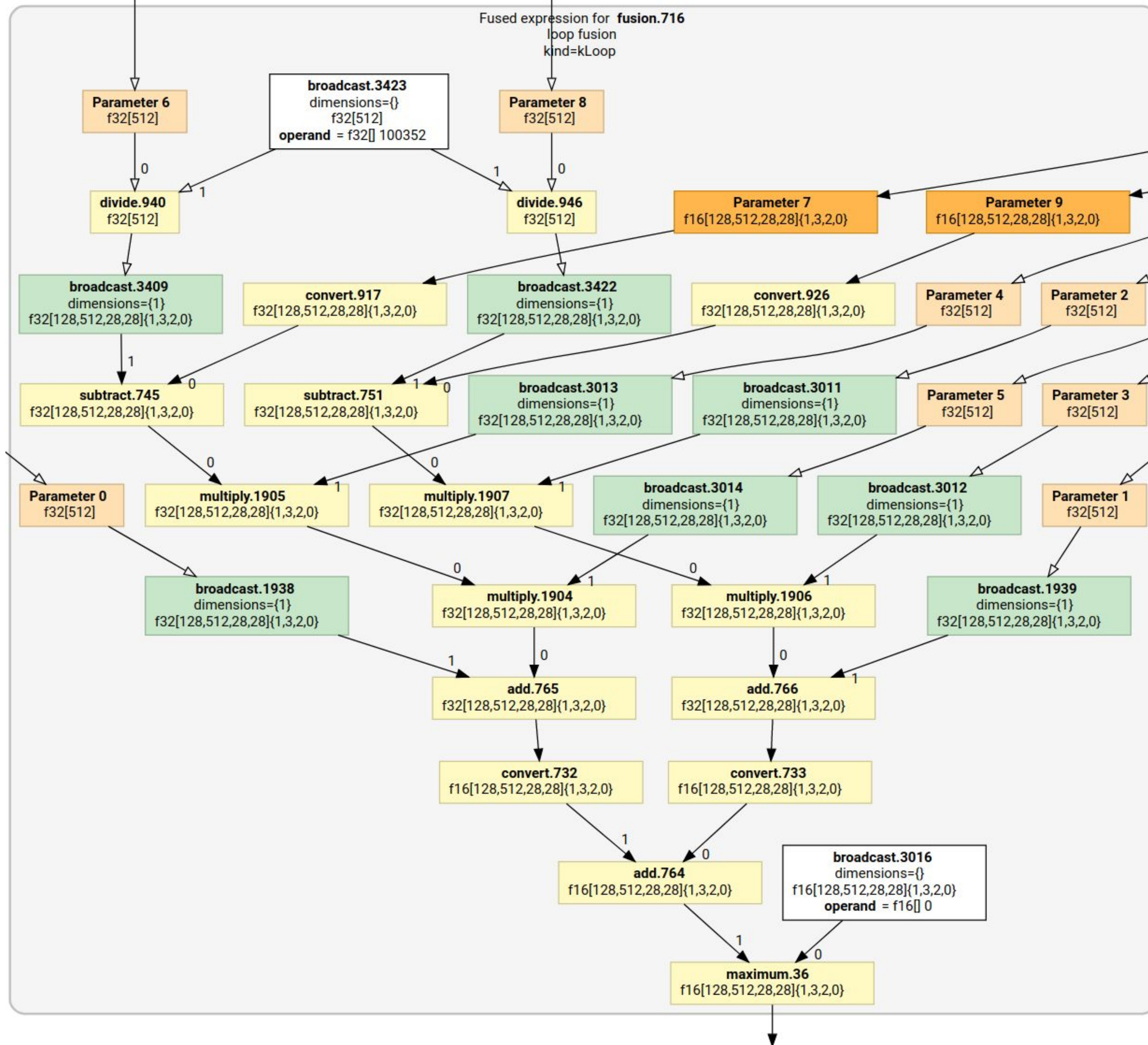


```
std::function<llvm::Value*>(const IrArray::Index& index)
    MakeElementGenerator(const HloInstruction* hlo,
                        HloToElementGeneratorMap& operand_to_generator) {
switch (hlo->opcode()) {
case HloOpcode::kMaximum:
    return [...](const IrArray::Index& index) {
        llvm::Value* lhs =
            operand_to_generator.at(hlo->operand(0))(index);
        llvm::Value* rhs =
            operand_to_generator.at(hlo->operand(1))(index);
        auto cmp = b->CreateFCmpUGE(lhs, rhs);
        return ir_builder_->CreateSelect(cmp, lhs, rhs);
    };
    ...
}
```

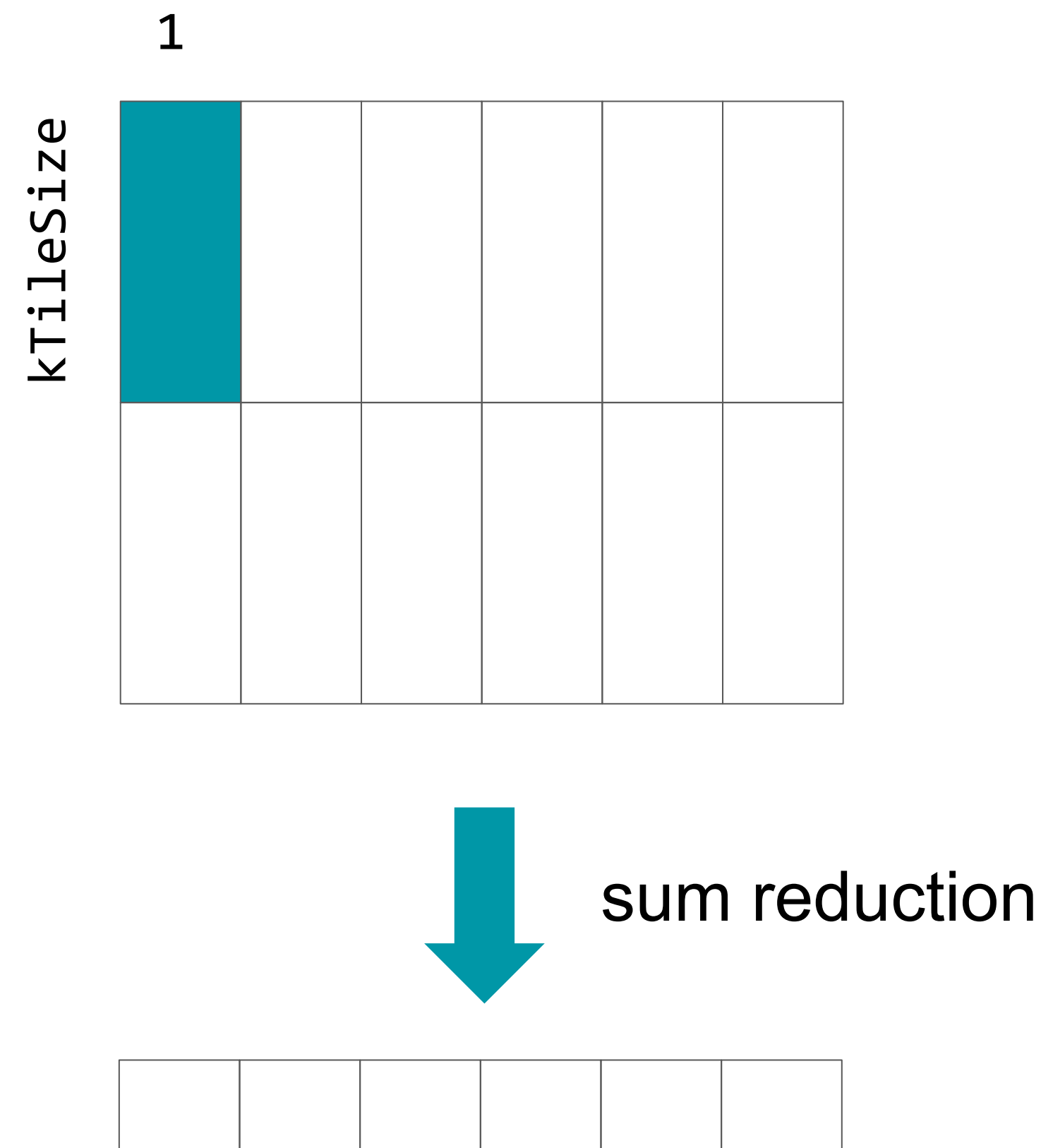
Fused Add + ReLu



```
__global__ void fusion(float *lhs,  
float *rhs, float* output) {  
  
int i = blockIdx.x * blockDim.x +  
threadIdx.x;  
  
if (i < 128*512*28*28) {  
output[i] = max(0.0, lhs[i] + rhs[i]);  
}  
  
}
```



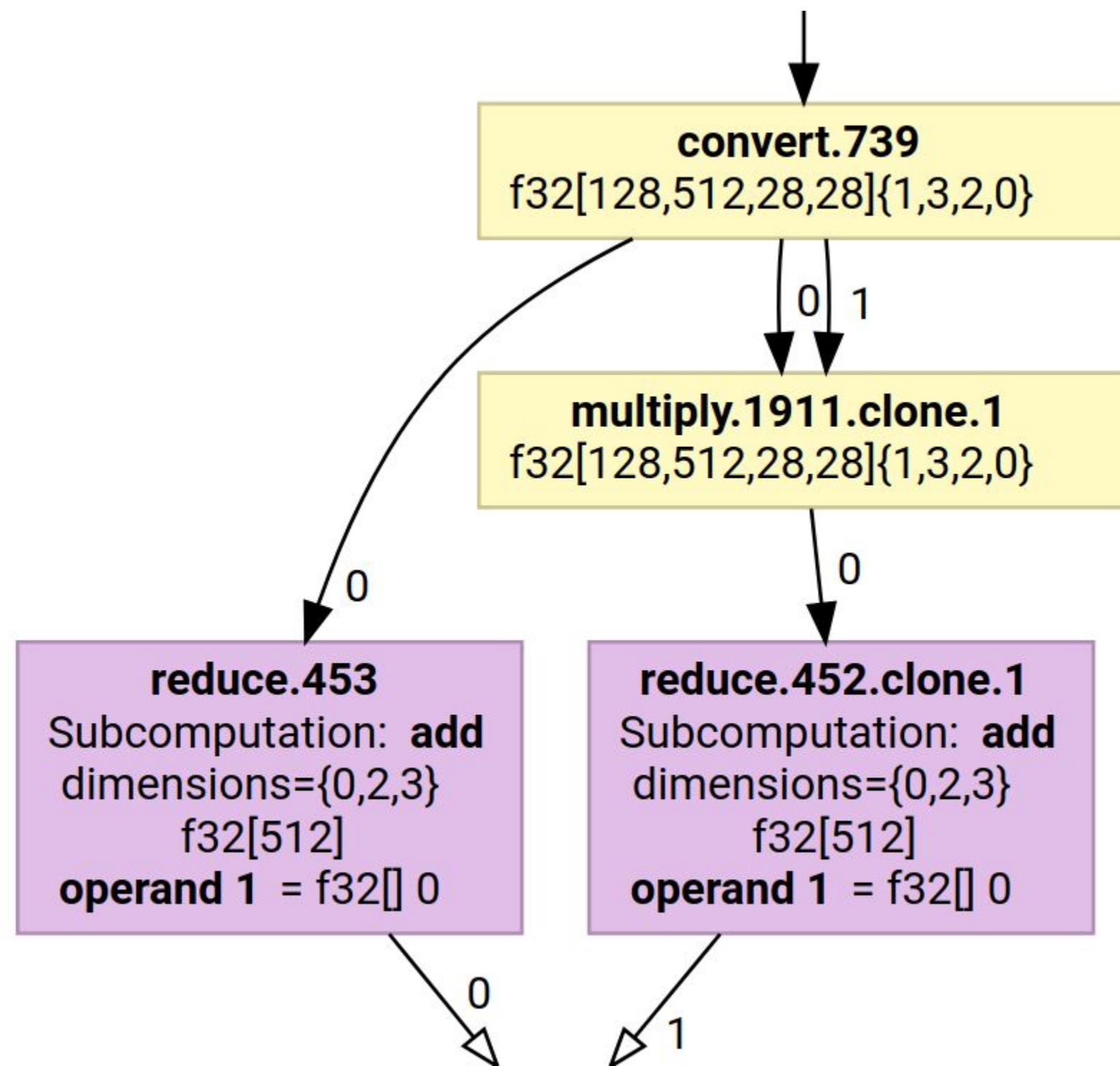
Reduction



```
i = blockIdx.x * blockDim.x + threadIdx.x;  
y_in_tiles = i / width;  
x = i % width;
```

```
for (int j = 0; j < kTileSize: ++j) {  
    y = y_in_tiles * kTileSize + j;  
    if (y < height) {  
        partial_sum += generator(y, x);  
    }  
}  
atomicAdd(&output[x], partial_sum);
```

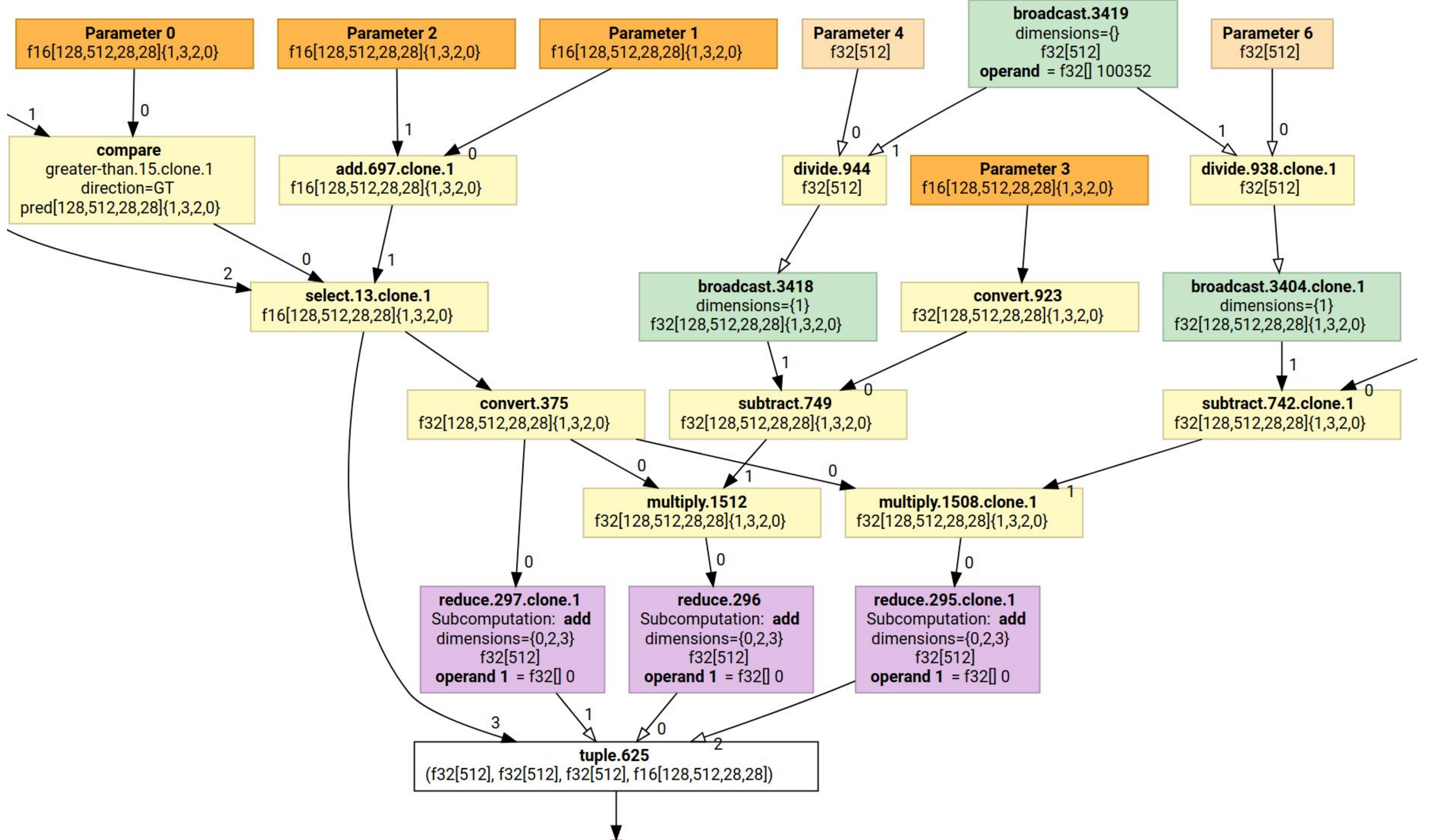
Multi-output fusion



```
i = blockIdx.x * blockDim.x + threadIdx.x;
y_in_tiles = i / width;
x = i % width;
```

```
for (int j = 0; j < kTileSize: ++j) {
    y = y_in_tiles * kTileSize + j;
    if (y < height) {
        partial_sum[0] += generator[0](y, x);
        partial_sum[1] += generator[1](y, x);
    }
}
```

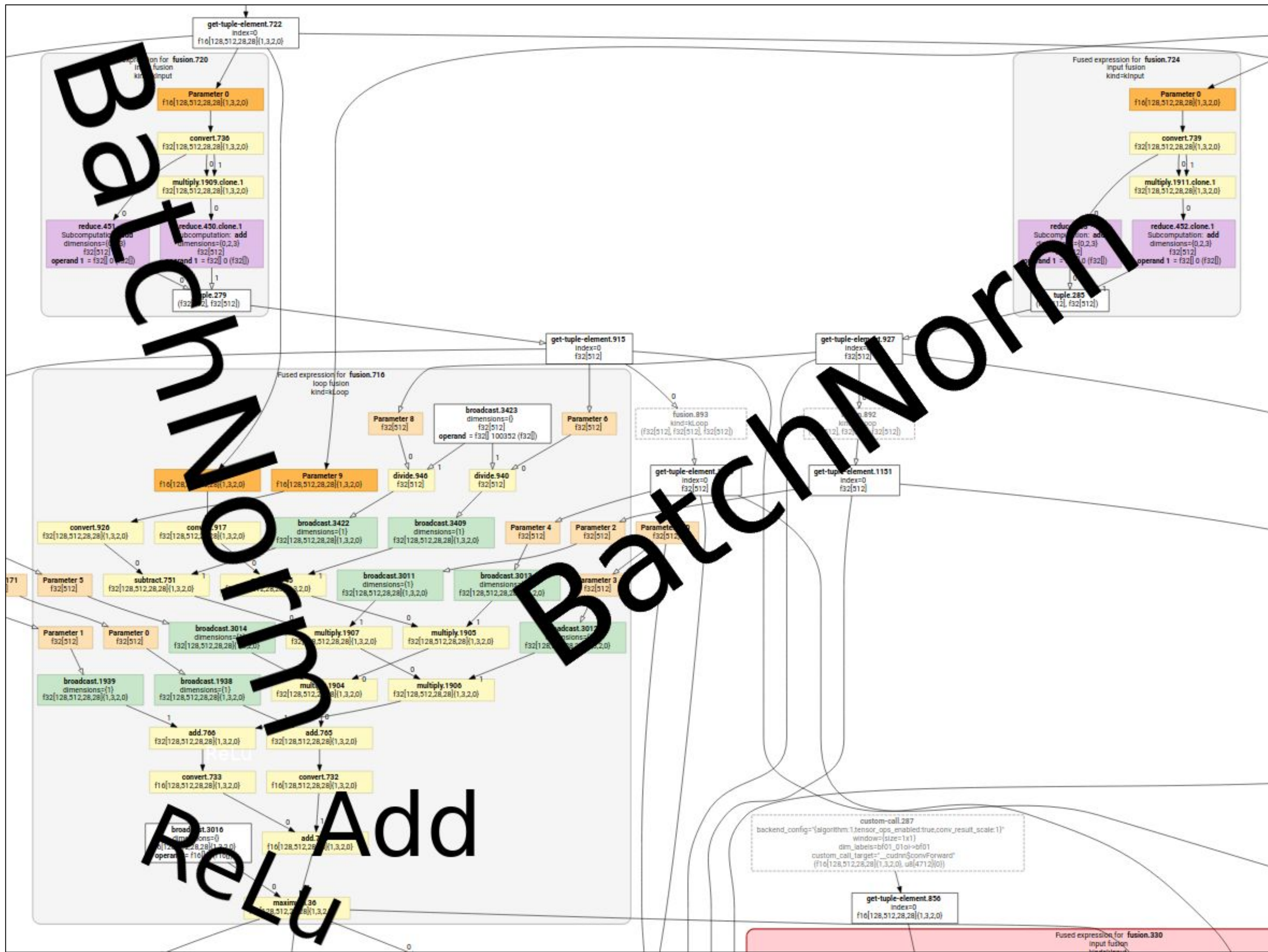
```
atomicAdd(&output[0][x], partial_sum[0]);
atomicAdd(&output[1][x], partial_sum[1]);
```



```
i = blockIdx.x * blockDim.x + threadIdx.x;
y_in_tiles = i / width;
x = i % width;

for (int j = 0; j < kTileSize: ++j) {
    y = y_in_tiles * kTileSize + j;
    if (y < height) {
        partial_sum[0] += generator[0](y, x);
        partial_sum[1] += generator[1](y, x);
        partial_sum[2] += generator[2](y, x);
        output[3][y, x] = generator[3](y, x);
    }
}

atomicAdd(&output[0][x], partial_sum[0]);
atomicAdd(&output[1][x], partial_sum[1]);
atomicAdd(&output[2][x], partial_sum[2]);
```



Thank you! Questions?

XLA documentation

<https://www.tensorflow.org/xla/overview>

Public XLA mailing list

xla-dev@googlegroups.com

XLA on Github

<https://github.com/tensorflow/tensorflow/tree/master/tensorflow/compiler>