

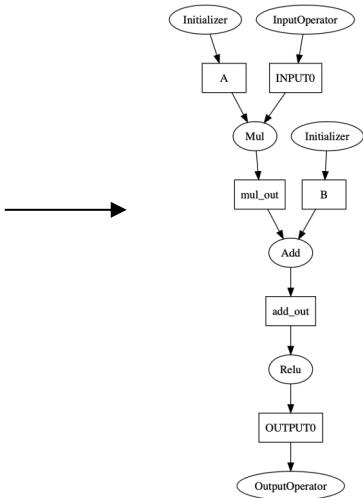
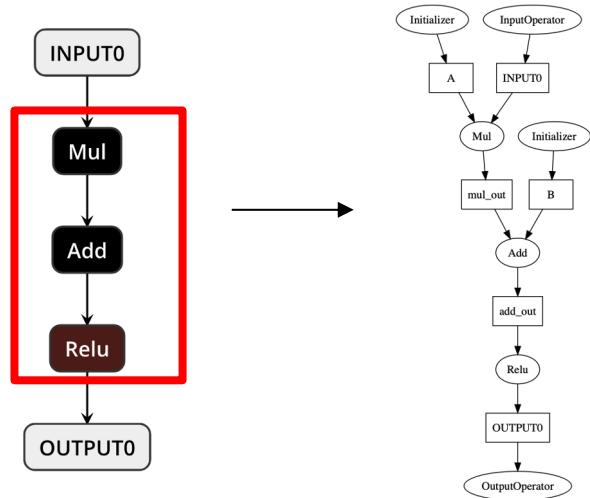
# **Hardware-Dependent Optimizations**

---

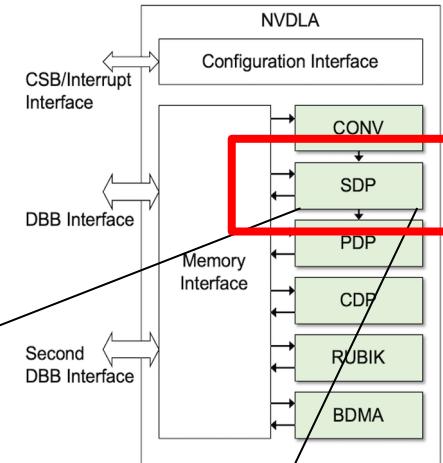
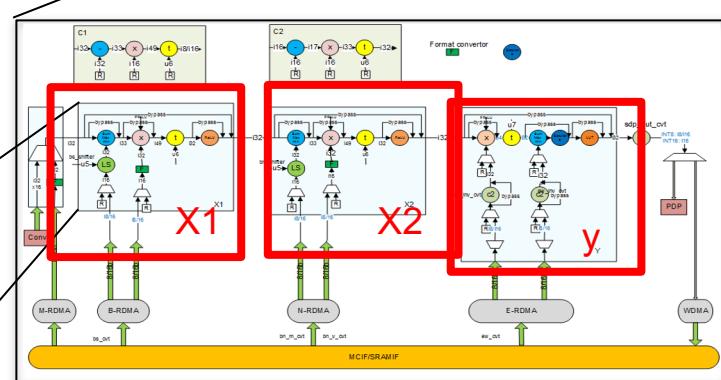
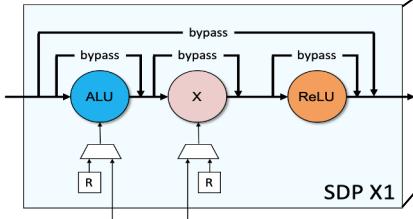
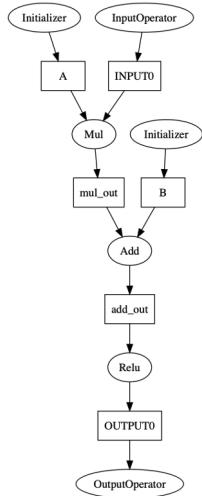
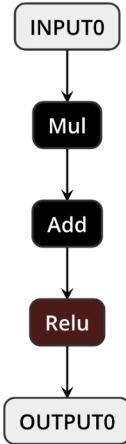
Cheng-Yi Chou, May 2020  
Skymizer Taiwan Inc.



# Mapping IR Graph To NVDLA Hardware



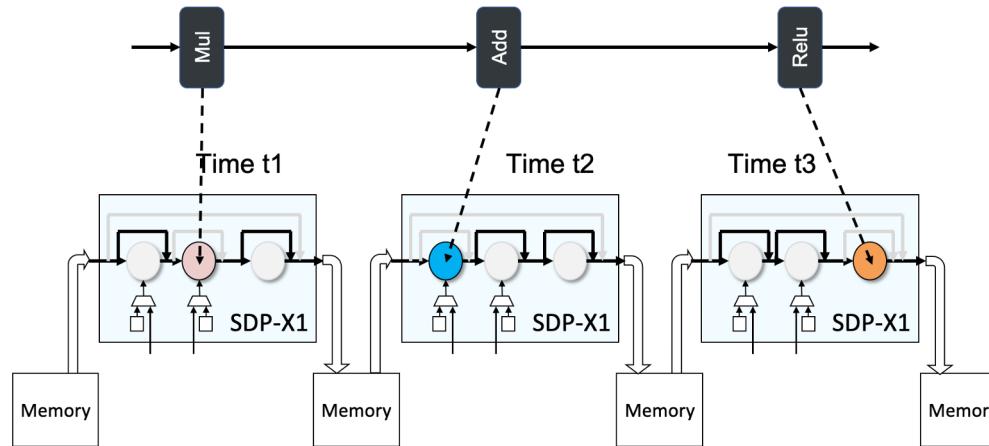
# Mapping IR Graph To NVDLA Hardware



# Mapping Options

If a model contains a series of **Mul-Add-Relu**, how to map operators into NVDLA hardware ?

Option 1: Execution time = T  
Memory request =  $\text{tensor\_size} * 6$

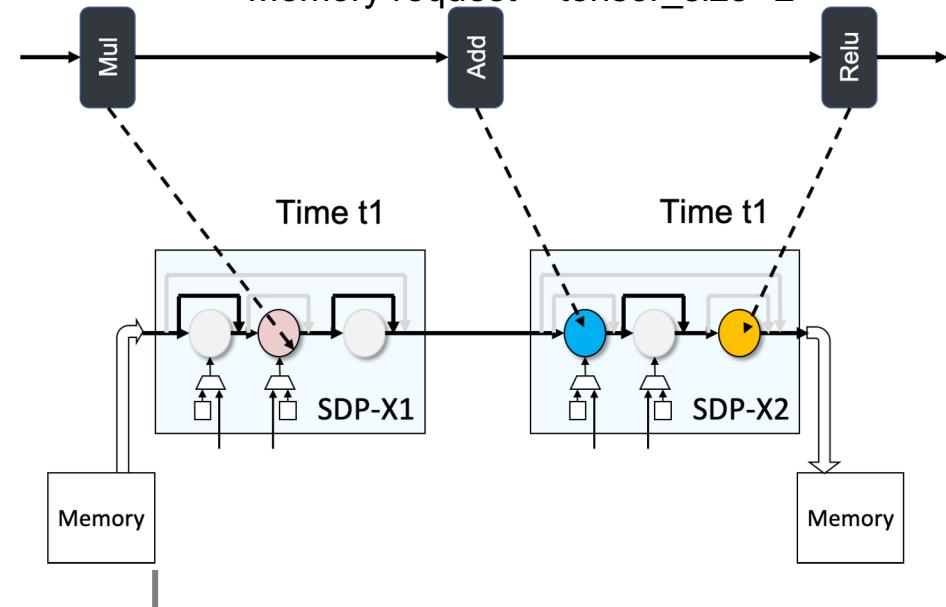


# Mapping Options

If a model contains a series of **Mul-Add-Relu**, how to map operators into NVDLA hardware ?

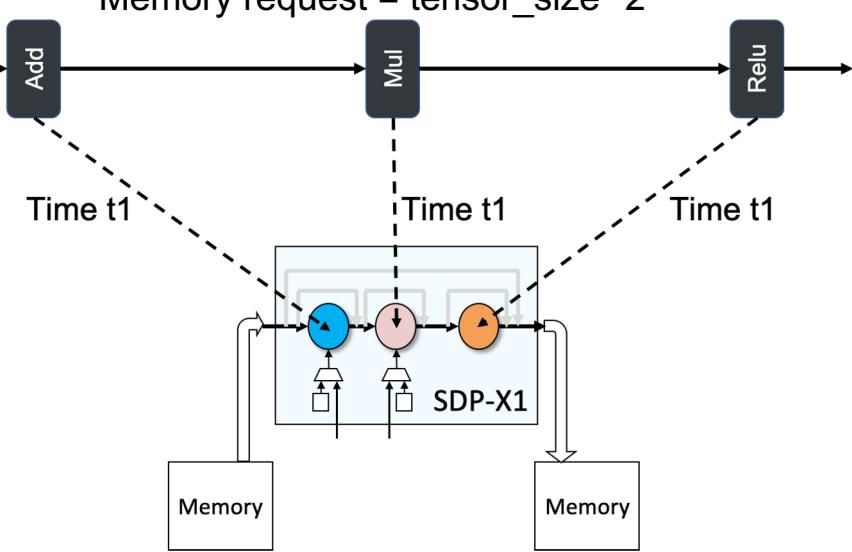
Option 2: Execution time  $\sim T/3$

Memory request =  $\text{tensor\_size} * 2$

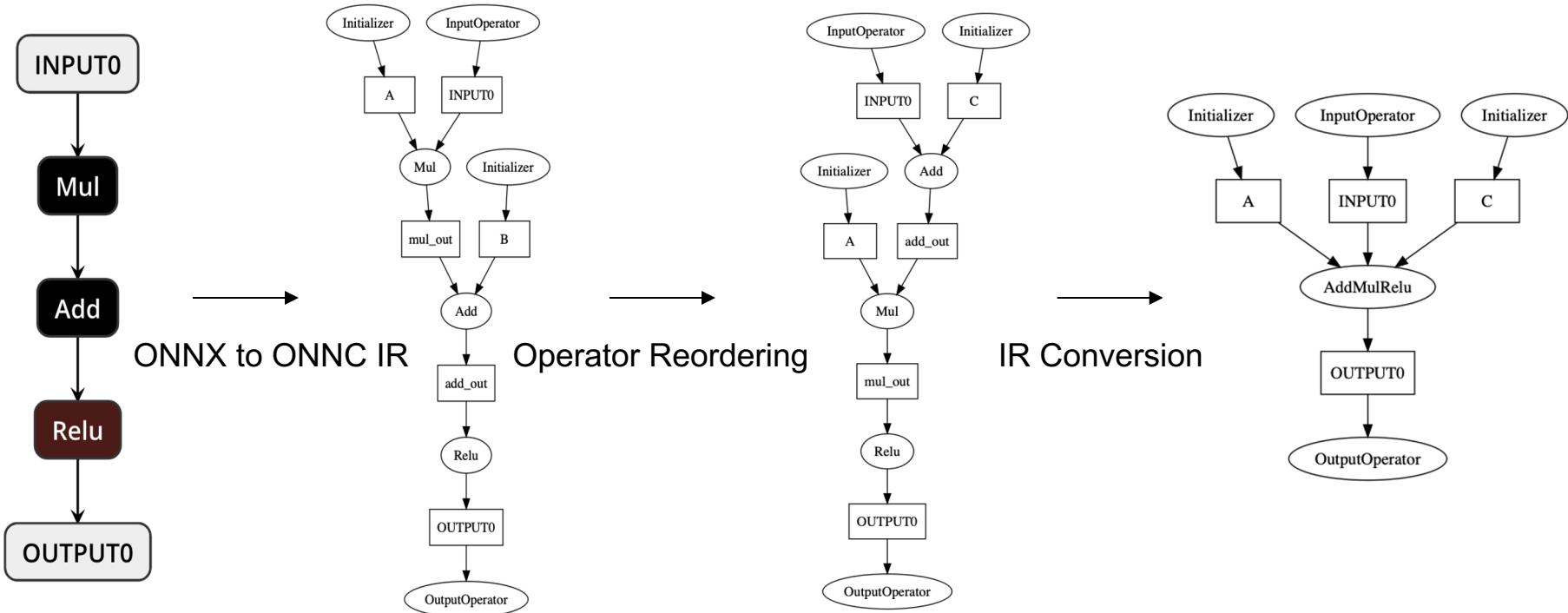


Option 3: Execution time  $\sim T/3$

Memory request =  $\text{tensor\_size} * 2$



# Hardware-Dependent Optimization For More Efficient Mapping



# Operator Reordering

$$Y = (X * \text{alpha}) + \text{beta}$$



$$Y = (X + \text{gamma}) * \text{alpha}$$
$$\text{gamma} = \text{beta} / \text{alpha}$$

# Operator Reordering

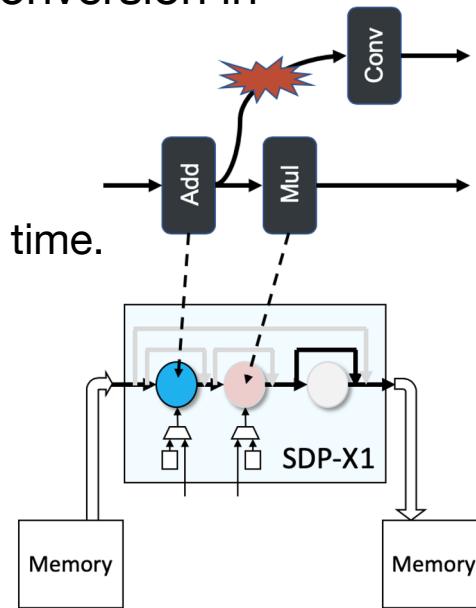
There are two pre-conditions to meet in order to have a valid conversion in this example.

1. The output of the **Add** operator can only have one consumer.
2. The values of **alpha** and **beta** need to be determined at compile time.

$$Y = (X * \text{alpha}) + \text{beta}$$



$$Y = (X + \text{gamma}) * \text{alpha}$$
$$\text{gamma} = \text{beta} / \text{alpha}$$



# Lab - Mul And Add Re-Ordering And Fusion

In this lab, we will show how to create hardware-dependent optimization passes described in earlier slides

1. NvDlaReorderMulAddPass
2. NvDlaFuseAddMulReluPass

# Setup Environment

```
1 # Skip this command if you have already had the Docker image.  
2 $ docker pull onnc/onnc-community  
3  
4 # Prepare ONNC and tutorial source code.  
5 $ git clone https://github.com/ONNC/onnc.git  
6 $ cd onnc; git checkout tags/1.3.0; cd ..  
7 $ git clone https://github.com/ONNC/onnc-tutorial.git  
8  
9 # Start the onnc/onnc-community Docker.  
10 $ docker run -ti --rm \  
11 -v <absolute/path/to/onnc>:/onnc/onnc \  
12 -v <absolute/path/to/tutorial>:/tutorial \  
13 onnc/onnc-community
```

# Create New Backend - FooNvdla

```
1 $ cd /onnc/onnc
2 $ ./scripts/create-new-backend.sh FooNvdla
3
4 # Install the pre-built FooNvdla backend.
5 $ tar -zxvf <path/to/tutorial>/lab_4_Code_Emitting/src/FooNvdla.tar.gz \
6 -C <path/to/onnc>/lib/Target
7
8 $ cd /onnc/onnc-umbrella/build-normal
9 # Rebuild ONNC.
10 $ smake -j8 install
11
12 # Run ONNC to compile a DNN model.
13 $ onnc -mquadruple foonvdla \
14 /tutorial/models/test_group_Conv/test_group_Conv.onnx
15 FooNvdla is invoked
```

# NvDlaReorderMulAddPass - Operator Reordering

- The re-ordering optimization is done in a pass named **NvDlaReorderMulAddPass**.
- We provide all of the related files in  
***<path/to/tutorial>/lab\_8\_Mul\_Add\_Reordering\_and\_Fusion/src/*** directory.

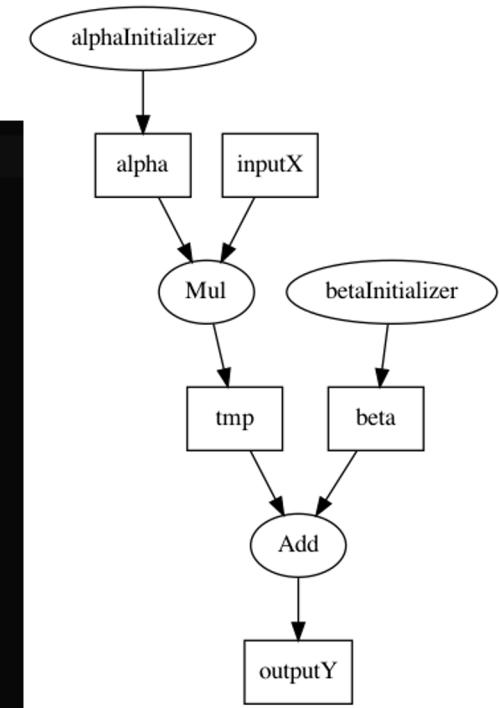
```
1 | $ cp \
2 | /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/NvDlaReorderMulAddPass.* \
3 | /onnc/onnc/lib/Target/FooNvdla
```

# 1. Search For The Matched Pattern

```
40 // File: NvDlaReorderMulAddPass.cpp
41 Pass::ReturnType NvDlaReorderMulAddPass::runOnComputeGraph(ComputeGraph& pCG)
42 {
43     Pass::ReturnType ret = Pass::kModuleNoChanged;
44     //-----
45     // Search for the Mul-Add patterns that can be reordered.
46     //-----
47     std::vector<ComputeOperator*> mulList;
48     for (ComputeOperator& node : pCG) {
49         if (canBeReordered(&node)) {
50             mulList.emplace_back(&node);
51             ret |= Pass::kModuleChanged;
52         }
53     }
54 }
```

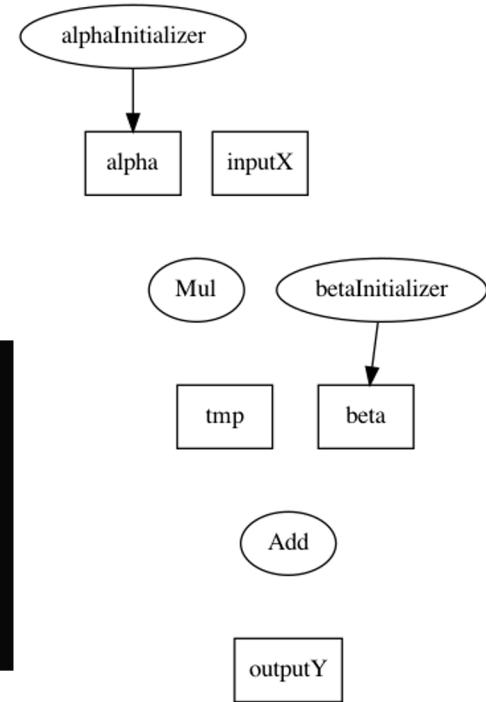
## 2. Iterate Through The Matched Patterns

```
54 // File: NvDlaReorderMulAddPass.cpp
55 //-----
56 // Perform re-ordering on the found patterns.
57 //-----
58 // The original pattern is:
59 //   outputY = (inputX * alpha) + beta
60 //
61 // We will re-arrange the above pattern by:
62 //   outputY = (inputX + gamma) * alpha, where
63 //   gamma = beta / alpha
64
65 for (ComputeOperator* node : mulList) {
66     Mul* mul = dyn_cast<Mul>(node);
67     Add* add = dyn_cast<Add>(node->getOutput(0)->getUses()[0].getUser());
68     ...
69 }
```



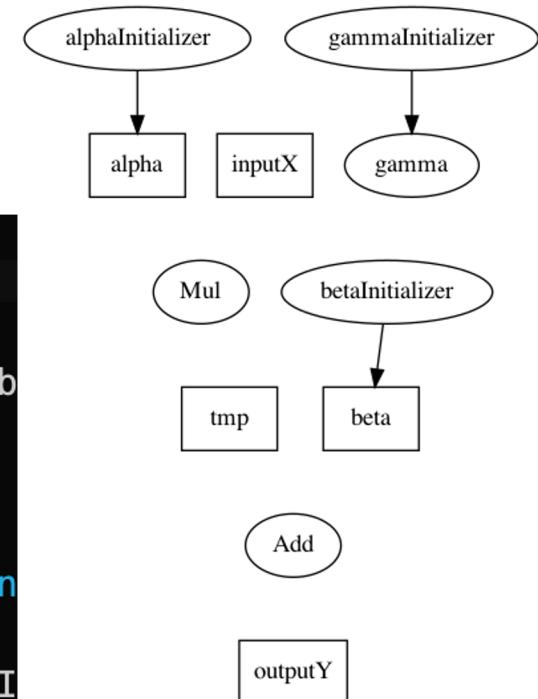
### 3. Remove Edges

```
120 // File: NvDlaReorderMulAddPass.cpp  
121     mul->removeAllInputs();  
122     mul->removeAllOutputs();  
123     add->removeAllInputs();  
124     add->removeAllOutputs();
```



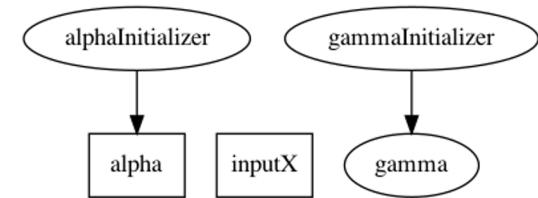
## 4. Create Gamma Tensor

```
143 // File: NvDlaReorderMulAddPass.cpp
144     // Create a new tensor gamma.
145     FloatTensor* gamma = dynamic_cast<FloatTensor*>(b
146     ...
147     // Add gamma into the ONNC IR graph.
148     gamma = pCG.addValue<FloatTensor>(gamma);
149     assert((gamma != nullptr) && "The name must be un
150     ...
151     Initializer* gammaInitializer = pCG.addOperator<I
152     gammaInitializer->setTensor(*gamma);
```



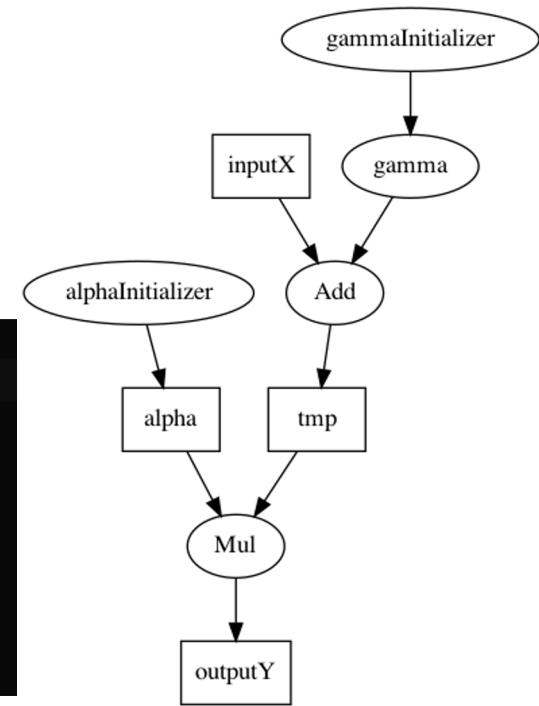
## 5. Remove Beta Tensor

```
184 // File: NvDlaReorderMulAddPass.cpp
185     // Calculate the constant data of gamma.
186     int tensorSize = beta->getValues().size();
187     for (int i = 0; i < tensorSize; i++) {
188         gamma->getValues().push_back( betaData[i] / alpha);
189     }
190     // Remove beta from the ONNC IR graph. We don't need it.
191     Initializer* betaInitializer = static_cast<Initializer*>(beta);
192     pCG.erase(*betaInitializer);
193     pCG.erase(*beta);
```



## 6. Re-Connect The Operators

```
213 // File: NvDlaReorderMulAddPass.cpp  
214     add->addInput(*inputX);  
215     add->addInput(*gamma);  
216     add->addOutput(*tmp);  
217     mul->addInput(*tmp);  
218     mul->addOutput(*outputY);
```



# NvDlaFuseAddMulReluPass: Fuse Add-Mul-Relu Operators

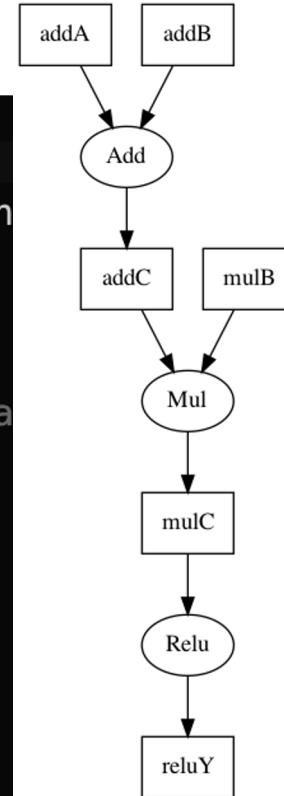
- The fuse optimization is done in a pass named **NvDlaReorderMulAddPass**.
- We provide all of the related files in  
*<path/to/tutorial>/lab\_8\_Mul\_Add\_Reordering\_and\_Fusion/src/*  
directory.

# NvDlaFuseAddMulReluPass: Fuse Add-Mul-Relu Operators

```
1 $ mkdir -p /onnc/onnc/lib/Target/FooNvdla/Compute
2 # These files are about the definition of new IR.
3 $ cp /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/NvDlaAddMulRelu.* \
4 /onnc/onnc/lib/Target/FooNvdla/Compute
5
6 # These files are about deploying the new IR into the model graph.
7 $ cp \
8 /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/NvDlaFuseAddMulReluPass.* \
9 /onnc/onnc/lib/Target/FooNvdla
10
11 # These files are about the code emitting functions for the new IR.
12 $ cp /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/CodeEmitVisitor.* \
13 /onnc/onnc/lib/Target/FooNvdla
14
15 # These files are about visualizing the optimization effect.
16 $ cp /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/PrintONNCIRPass.* \
17 /onnc/onnc/lib/Target/FooNvdla
```

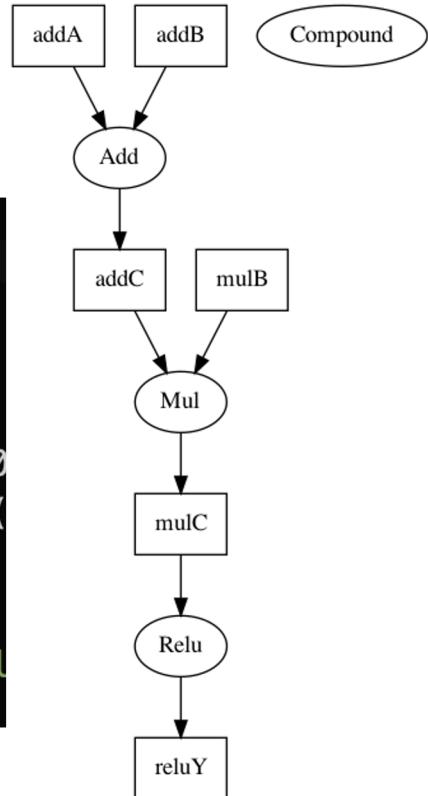
# 1. Find The Matched Pattern

```
37 // File: NvDlaFuseAddMulReluPass.cpp
38 Pass::ReturnType NvDlaFuseAddMulReluPass::runOnComputeGraph
39 {
40     Pass::ReturnType ret = Pass::kModuleNoChanged;
41
42     // Search for the Add-Mul-Relu patterns that can be replaced
43     std::vector<ComputeOperator*> patternList;
44     for (ComputeOperator& node : pCG) {
45         if (isAddMulRelu(&node)) {
46             patternList.emplace_back(&node);
47             ret |= Pass::kModuleChanged;
48         }
49     }
```



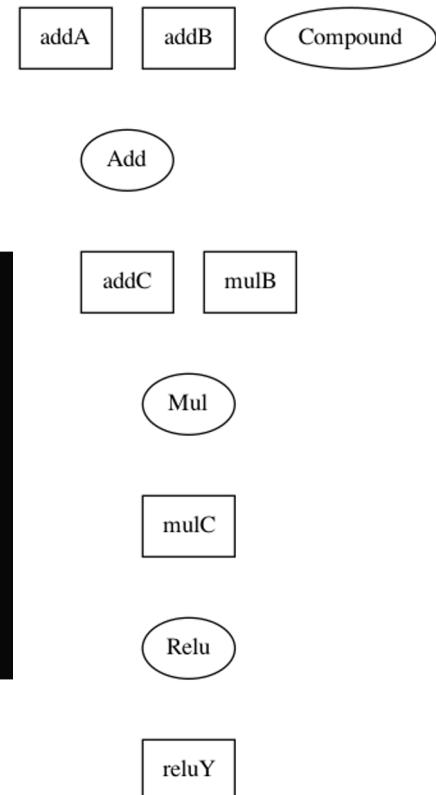
## 2. Create A New Compound IR

```
50 // File: NvDlaFuseAddMulReluPass.cpp
51     for (ComputeOperator* node : patternList) {
52         // Derive original IRs.
53         Add* add = dyn_cast<Add>(node);
54         Mul* mul = dyn_cast<Mul>(add->getOutput(0)->getUses()[0]);
55         Relu* relu = dyn_cast<Relu>(mul->getOutput(0)->getUses());
56         ...
57         // Create a new AddMulRelu IR.
58         NvDlaAddMulRelu* compound = pCG.addOperator<NvDlaAddMulRelu>(add, mul, relu);
```



### 3. Remove Edges

```
109 // File: NvDlaFuseAddMulReluPass.cpp
110     add->removeAllInputs();
111     add->removeAllOutputs();
112     mul->removeAllInputs();
113     mul->removeAllOutputs();
114     relu->removeAllInputs();
115     relu->removeAllOutputs();
```



## 4. Remove Original IRs

addA

addB

Compound

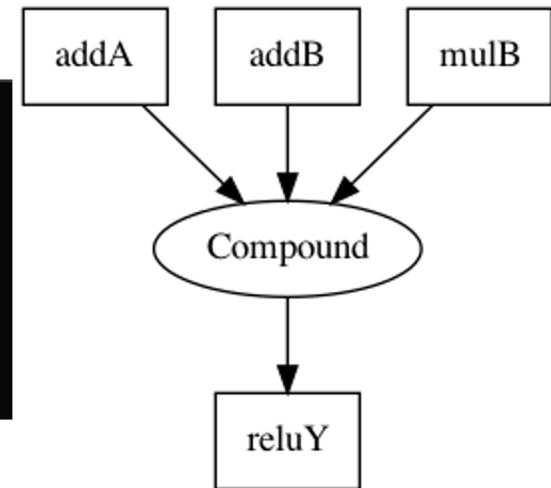
```
135 // File: NvDlaFuseAddMulReluPass.cpp
136     pCG.erase(*add);
137     pCG.erase(*mul);
138     pCG.erase(*relu);
139     pCG.erase(*addC);
140     pCG.erase(*mulC);
```

mulB

reluY

## 5. Connect The New Compound IR

```
155 // File: NvDlaFuseAddMulReluPass.cpp  
156 compound->addInput(*addA);  
157 compound->addInput(*addB);  
158 compound->addInput(*mulB);  
159 compound->addOutput(*reluY);
```



# Enable Passes In The Backend

```
1 $ cp /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/FooNvdlaBackend.cpp \
2 /onnc/onnc/lib/Target/FooNvdla
3
4 $ cp /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/CMakeLists.txt \
5 /onnc/onnc/lib/Target/FooNvdla
6 $ cp /tutorial/lab_8_Mul_Add_Reordering_and_Fusion/src/Makefile.am \
7 /onnc/onnc/lib/Target/FooNvdla
```

# Rebuild ONNC & Compile Example Model

```
1 $ cd /onnc/onnc-umbrella/build-normal
2
3 # Rebuild ONNC.
4 $ smake -j8 install
5
6 # Execute ONNC to compile the model.
7 $ onnc -mquadruple foonvdla \
8 /tutorial/models/test_Mul_Add_ReLU/test_Mul_Add_ReLU.onnx
```

```
1 FooNvdla is invoked
2 === PrintONNCIRPass =====
3 %A<float>[1, 1, 5, 5] = Initializer<unimplemented>()
4 %B<float>[1, 1, 5, 5] = Initializer<unimplemented>()
5 %INPUT0<float>[1, 1, 5, 5] = InputOperator<unimplemented>()
6 %mul_out<float>[1, 1, 5, 5] = Mul(%INPUT0<float>[1, 1, 5, 5], %A<float>[1, 1,
7 %add_out<float>[1, 1, 5, 5] = Add(%mul_out<float>[1, 1, 5, 5], %B<float>[1, 1,
8 %OUTPUT0<float>[1, 1, 5, 5] = Relu(%add_out<float>[1, 1, 5, 5])
9   = OutputOperator<unimplemented>(%OUTPUT0<float>[1, 1, 5, 5])
10 =====
11 NvDlaReorderMulAddPass is called...
12 === PrintONNCIRPass =====
13 %A<float>[1, 1, 5, 5] = Initializer<unimplemented>()
14 %INPUT0<float>[1, 1, 5, 5] = InputOperator<unimplemented>()
15 %B_gamma_0<float>[1, 1, 5, 5] = Initializer<unimplemented>()
16 %add_out<float>[1, 1, 5, 5] = Add(%INPUT0<float>[1, 1, 5, 5], %B_gamma_0<float>
17 %mul_out<float>[1, 1, 5, 5] = Mul(%add_out<float>[1, 1, 5, 5], %A<float>[1, 1,
18 %OUTPUT0<float>[1, 1, 5, 5] = Relu(%mul_out<float>[1, 1, 5, 5])
19   = OutputOperator<unimplemented>(%OUTPUT0<float>[1, 1, 5, 5])
20 =====
21 NvDlaFuseAddMulReluPass is called...
22 === PrintONNCIRPass =====
23 %A<float>[1, 1, 5, 5] = Initializer<unimplemented>()
24 %INPUT0<float>[1, 1, 5, 5] = InputOperator<unimplemented>()
25 %B_gamma_0<float>[1, 1, 5, 5] = Initializer<unimplemented>()
26 %OUTPUT0<float>[1, 1, 5, 5] = AddMulRelu<>(%INPUT0<float>[1, 1, 5, 5], %B_gamma_0<float>
27   = OutputOperator<unimplemented>(%OUTPUT0<float>[1, 1, 5, 5])
28 =====
29 visit(NvDlaAddMulRelu) is called
```



## Skymizer Taiwan Inc.

### CONTACT US

**E-mail** [sales@skymizer.com](mailto:sales@skymizer.com)    **Tel** +886 2 8797 8337

**HQ** 12F-2, No.408, Ruiguang Rd., Neihu Dist., Taipei City 11492, Taiwan

**BR** Center of Innovative Incubator, National Tsing Hua University, Hsinchu Taiwan



**skymizer**

Boost deep learning accelerator  
with compiler technology



<https://skymizer.com>