

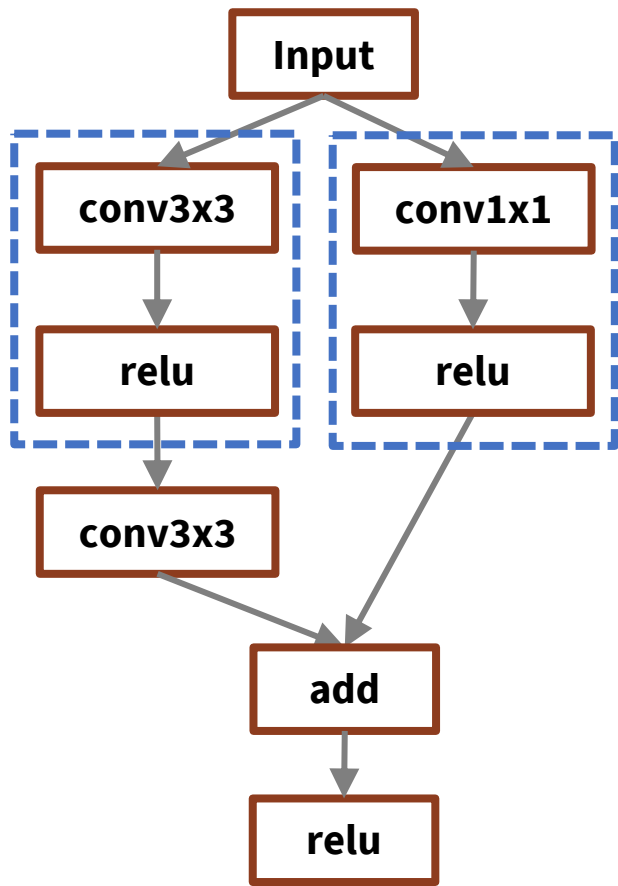
# TASO: Optimizing Deep Learning with Automatic Generation of Graph Substitutions

**Zhihao Jia**, Oded Padon, James Thomas, Todd Warszawski,  
Matei Zaharia, and Alex Aiken

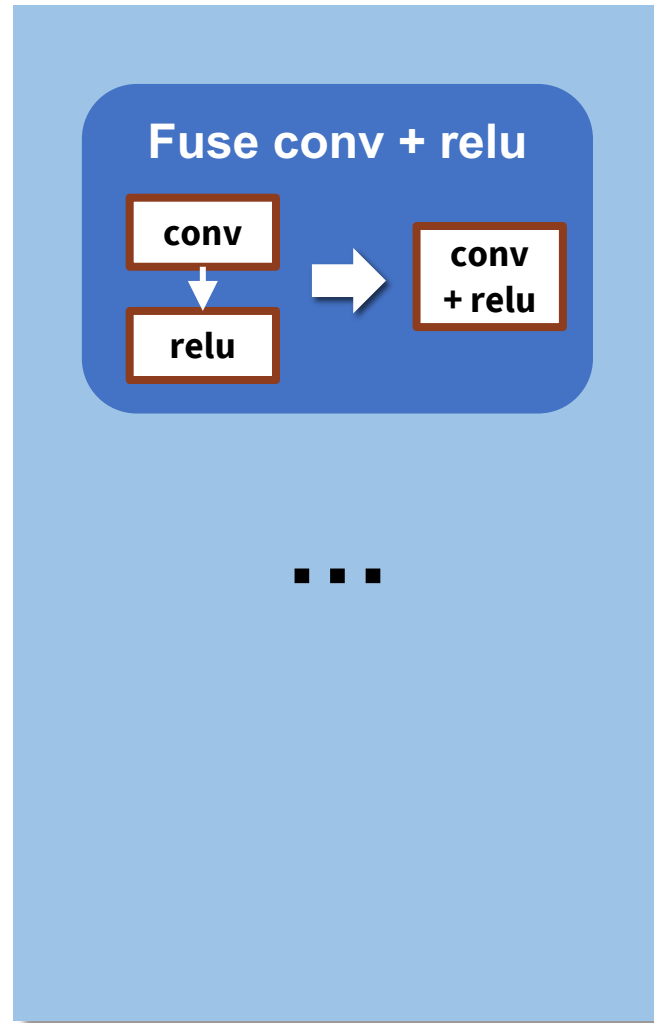
Stanford University



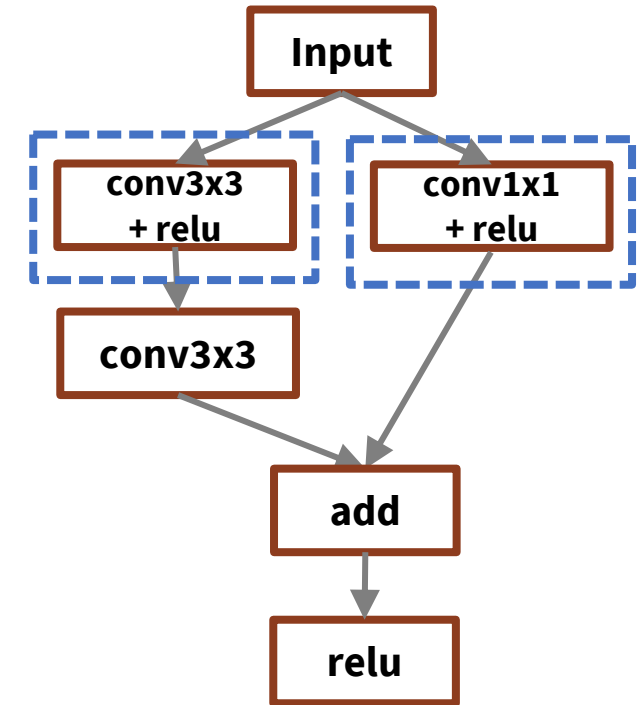
# Current Rule-based DNN Optimizations



Computation Graph



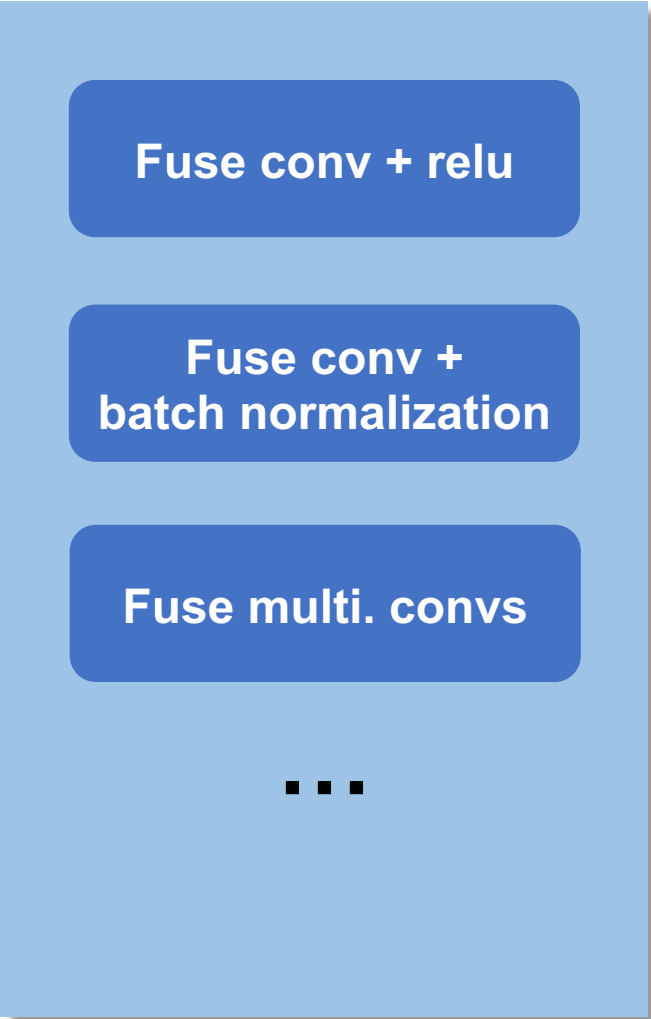
Rule-based Optimizer



Optimized Graph

# Current Rule-based DNN Optimizations

TensorFlow currently includes ~200 rules (~53,000 LOC)



Rule-based Optimizer

```
26 namespace tensorflow {
27 namespace graph_transforms {
28
29 // Converts Conv2D or MatMul ops followed by column-wise Muls into equivalent
30 // ops with the Mul baked into the convolution weights, to save computation
31 // during inference.
32 Status FoldBatchNorms(const GraphDef& input_graph_def,
33                      const TransformFuncContext& context,
34                      GraphDef* output_graph_def) {
35   GraphDef replaced_graph_def;
36   TF_RETURN_IF_ERROR(ReplaceMatchingOpTypes(
37     input_graph_def, // clang-format off
38     {"Mul", // mul_node
39      {
40        {"Conv2D|MatMul|DepthwiseConv2dNative", // conv_node
41         {
42           {"*"}, // input_node
43           {"Const"}, // weights_node
44         },
45         {"Const"}, // mul_values_node
46       },
47     }, // clang-format on
48     [(const NodeMatch& match, const std::set<string>& input_nodes,
49      const std::set<string>& output_nodes,
50      std::vector<NodeDef*> new_nodes) {
51       // Find all the nodes we expect in the subgraph.
52       const NodeDef& mul_node = match.node;
53       const NodeDef& conv_node = match.inputs[0].node;
54       const NodeDef& input_node = match.inputs[0].inputs[0].node;
55       const NodeDef& weights_node = match.inputs[0].inputs[1].node;
56       const NodeDef& mul_values_node = match.inputs[1].node;
57
58       // Check that nodes that we use are not used somewhere else.
59       for (const auto& node : {conv_node, weights_node, mul_values_node}) {
60         if (output_nodes.count(node.name())) {
61           // Return original nodes.
62           new_nodes->insert(new_nodes->end(),
63                            {mul_node, conv_node, input_node, weights_node,
64                             mul_values_node});
65         }
66         return Status::OK();
67       }
68     }
69
70   Tensor weights = GetNodeTensorAttr(weights_node, "value");
71   Tensor mul_values = GetNodeTensorAttr(mul_values_node, "value");
72
73   // Make sure all the inputs really are vectors, with as many entries as
74   // there are columns in the weights.
75   int64 weights_cols;
76   if (conv_node.op() == "Conv2D") {
77     weights_cols = weights.shape().dim_size(3);
78   } else if (conv_node.op() == "DepthwiseConv2dNative") {
79     weights_cols =
80       weights.shape().dim_size(2) * weights.shape().dim_size(3);
81   } else {
82     weights_cols = weights.shape().dim_size(1);
83   }
84   if ((mul_values.shape().dims() != 1) ||
85       (mul_values.shape().dim_size(0) != weights_cols)) {
86     return errors::InvalidArgument(
87       "Mul constant input to batch norm has bad shape: ",
88       mul_values.shape().DebugString());
89   }
90
91   // Multiply the original weights by the scale vector.
92   auto weights_vector = weights.flat<float>();
93   Tensor scaled_weights(DT_FLOAT, weights.shape());
94   auto scaled_weights_vector = scaled_weights.flat<float>();
95   for (int64 row = 0; row < weights_vector.dimension(0); ++row) {
96     scaled_weights_vector(row) =
97       weights_vector(row) *
98       mul_values.flat<float>()(row % weights_cols);
99   }
100
101   // Construct the new nodes.
102   NodeDef scaled_weights_node;
103   scaled_weights_node.set_op("Const");
104   scaled_weights_node.set_name(weights_node.name());
105   SetNodeAttr("dtype", DT_FLOAT, &scaled_weights_node);
106   SetNodeTensorAttr("value", scaled_weights, &scaled_weights_node);
107   new_nodes->push_back(scaled_weights_node);
108
109   new_nodes->push_back(input_node);
110
111   NodeDef new_conv_node;
112   new_conv_node = conv_node;
113   new_conv_node.set_name(mul_node.name());
114   new_nodes->push_back(new_conv_node);
115
116   return Status::OK();
117 }, &replaced_graph_def);
118 *output_graph_def = replaced_graph_def;
119 return Status::OK();
120 }
121
122 REGISTER_GRAPH_TRANSFORM("fold_batch_norms", FoldBatchNorms);
123 } // namespace graph_transforms
124 } // namespace tensorflow
```

# Limitations of Rule-based Optimizations

## Robustness

Experts' heuristics do not apply to all DNNs/hardware

The screenshot shows a GitHub issue page for 'Horovod with XLA is slower than without XLA (Tensorflow 1.12) #713'. The issue is marked as 'Closed' and was opened by LiweiPeng on Dec 19, 2018. A comment from LiweiPeng, dated Dec 19, 2018, describes a performance issue: 'I have a distributed nmt model (Transformer-based, AdamOptimizer) using Horovod (0.15.1). When I turned on XLA under tensorflow 1.12, the training speed is about 20% slower instead of faster. This result is sampled after training 1.5-hours and 4000 steps. I am using 4 V100 GPUs for the training. Because my current software is tightly coupled with Horovod, I couldn't test whether this is Horovod related or not. Does anyone have experience on whether this is expected?'. The issue also has a 'question' label and a 'Subscribe' button.

When I turned on XLA (TensorFlow's graph optimizer), the training speed is **about 20% slower**.

The screenshot shows a Stack Overflow question titled 'Tensorflow XLA makes it slower?'. The user asks: 'I am writing a very simple tensorflow program with XLA enabled. Basically it's something like:'. The code provided is as follows:

```
import tensorflow as tf

def ChainSoftMax(x, n)
    tensor = tf.nn.softmax(x)
    for i in range(n-1):
        tensor = tf.nn.softmax(tensor)
    return tensor

config = tf.ConfigProto()
config.graph_options.optimizer_options.global_jit_level = tf.OptimizerOptions.ON_1

input = tf.placeholder(tf.float32, [1000])
feed = np.random.rand(1000).astype('float32')

with tf.Session(config=config) as sess:
    res = sess.run(ChainSoftMax(input, 2000), feed_dict={input: feed})
```

The user concludes: 'Basically the idea is to see whether XLA can fuse the chain of softmax together to avoid multiple kernel launches. With XLA on, the above program is almost 2x slower than that without XLA on a machine with a GPU card. In my gpu profile, I saw XLA produces lots of kernels named as "reduce\_xxx" and "fusion\_xxx" which seem to overwhelm the overall runtime. Any one know what happened here?'

With XLA, my program is **almost 2x slower than** without XLA



# Limitations of Rule-based Optimizations

## **Robustness**

Experts' heuristics do not apply to all DNNs/hardware

## **Scalability**

New operators and graph structures require more rules

**TensorFlow currently uses ~4K LOC to optimize convolution**

# Limitations of Rule-based Optimizations

## **Robustness**

Experts' heuristics do not apply to all DNNs/hardware

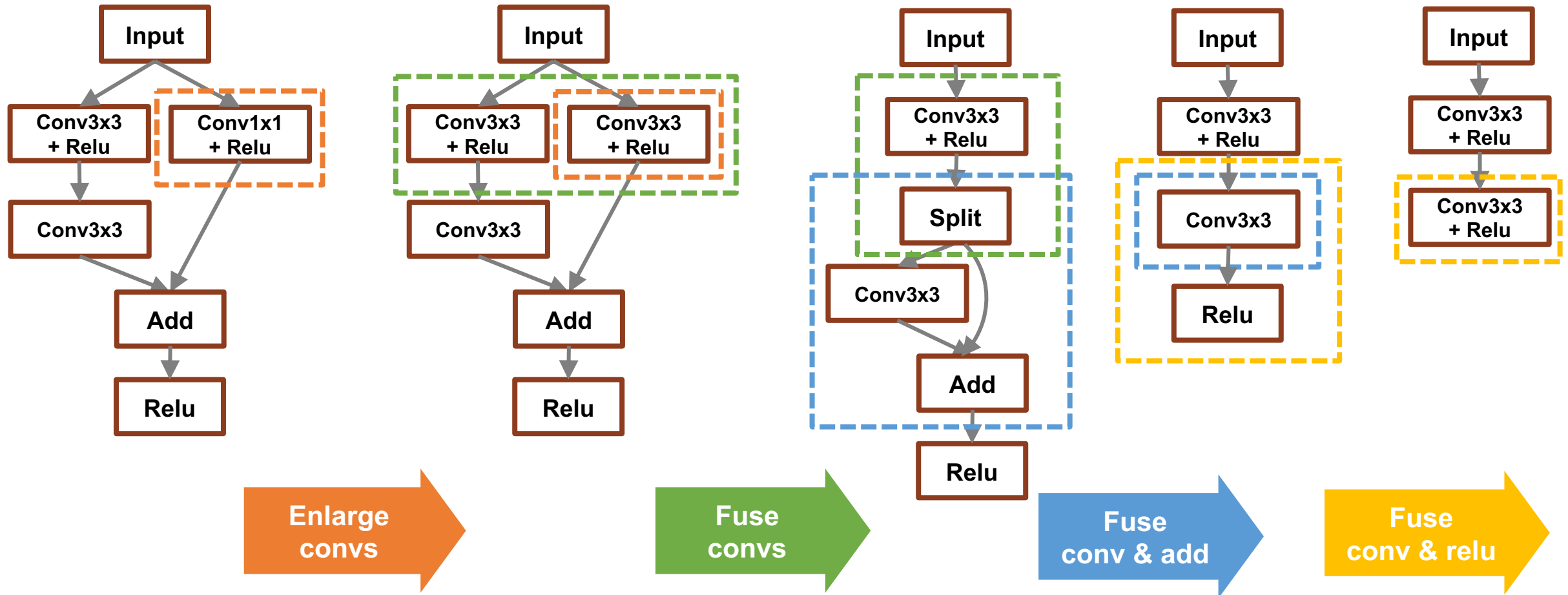
## **Scalability**

New operators and graph structures require more rules

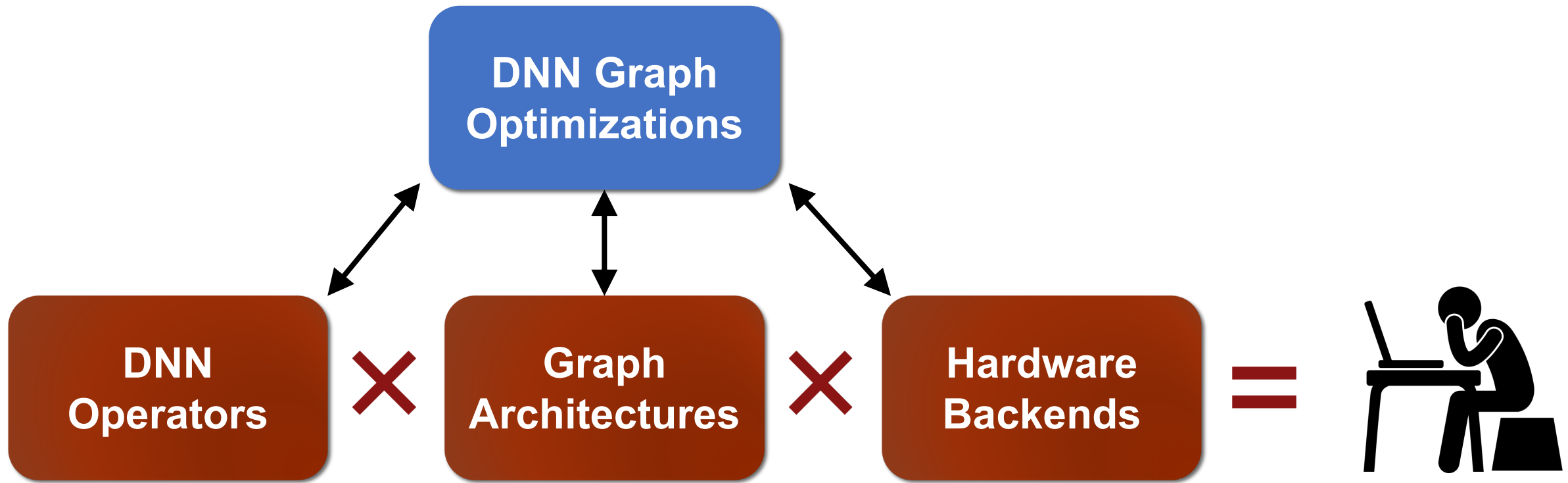
## **Performance**

Miss subtle optimizations for specific DNNs/hardware

# Motivating Example



The final graph is 30% faster on V100 but 10% slower on K80.

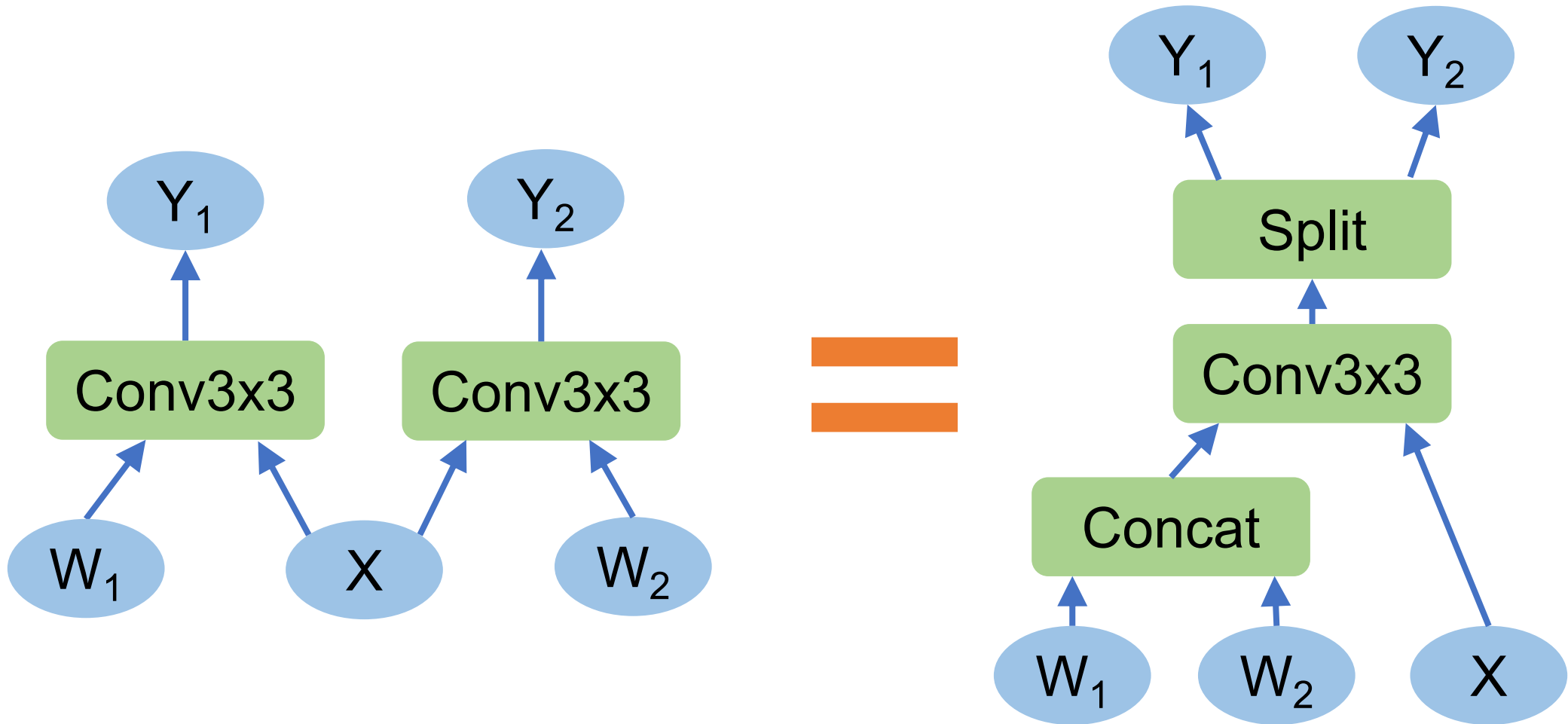


How should we address the complexity of designing DNN graph optimizations?

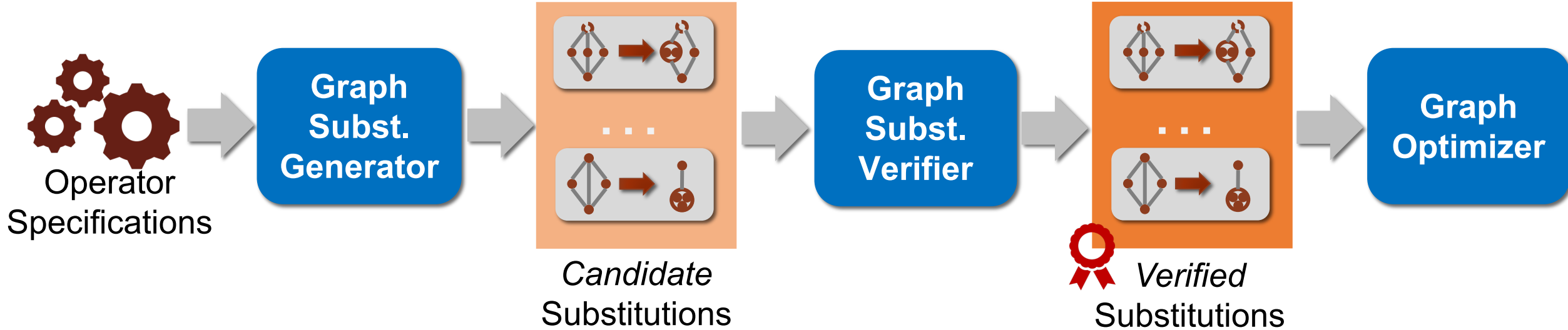
# TASO: Tensor Algebra SuperOptimizer

- Key idea: replace manually-designed graph optimizations with ***automated generation and verification*** of graph substitutions for deep learning
- **Less engineering effort:** 53,000 LOC for manual graph optimizations in TensorFlow → 1,400 LOC in TASO
- **Better performance:** outperform existing optimizers by up to 2.8x

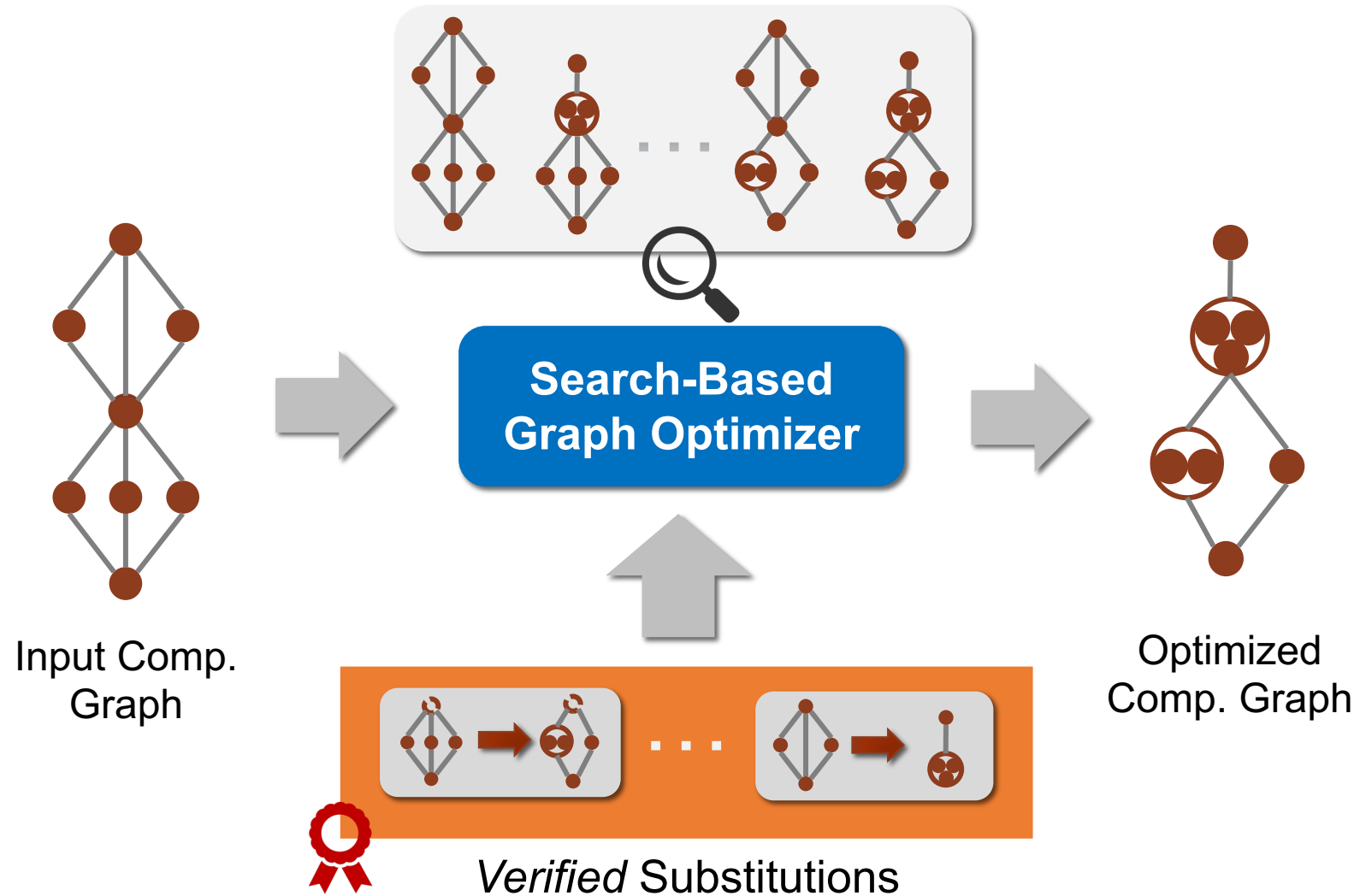
# Graph Substitution



# TASO Workflow



# TASO Workflow





# Key Challenges

1. How to generate potential substitutions?

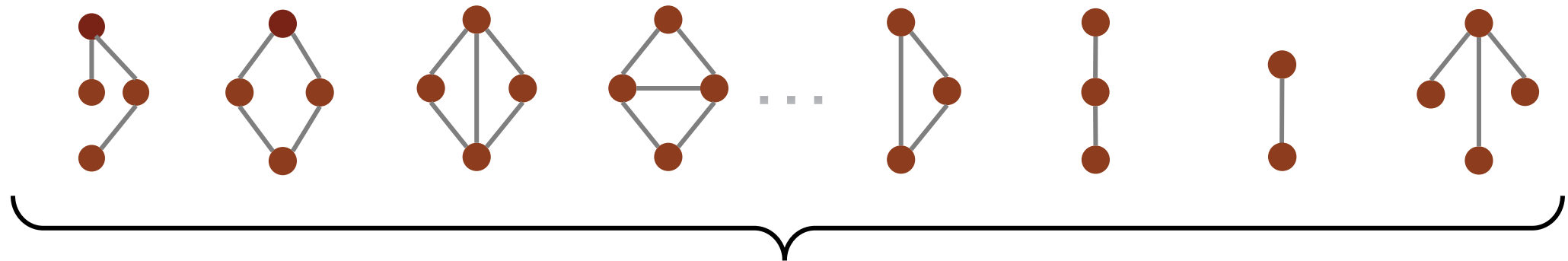
Graph fingerprints

2. How to verify their correctness?

Operator specifications + theorem prover

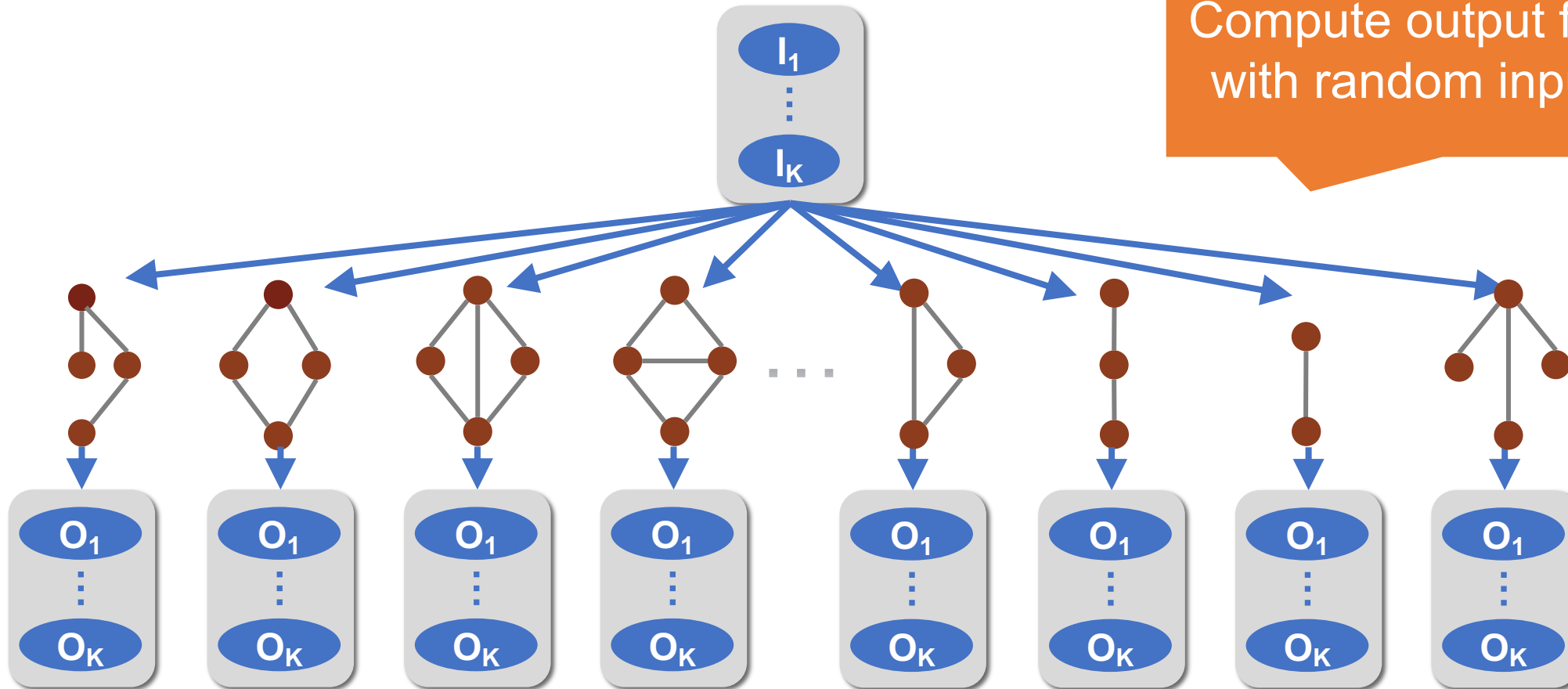
# Graph Substitution Generator

Enumerate all possible subgraphs up to a fixed size using available operators



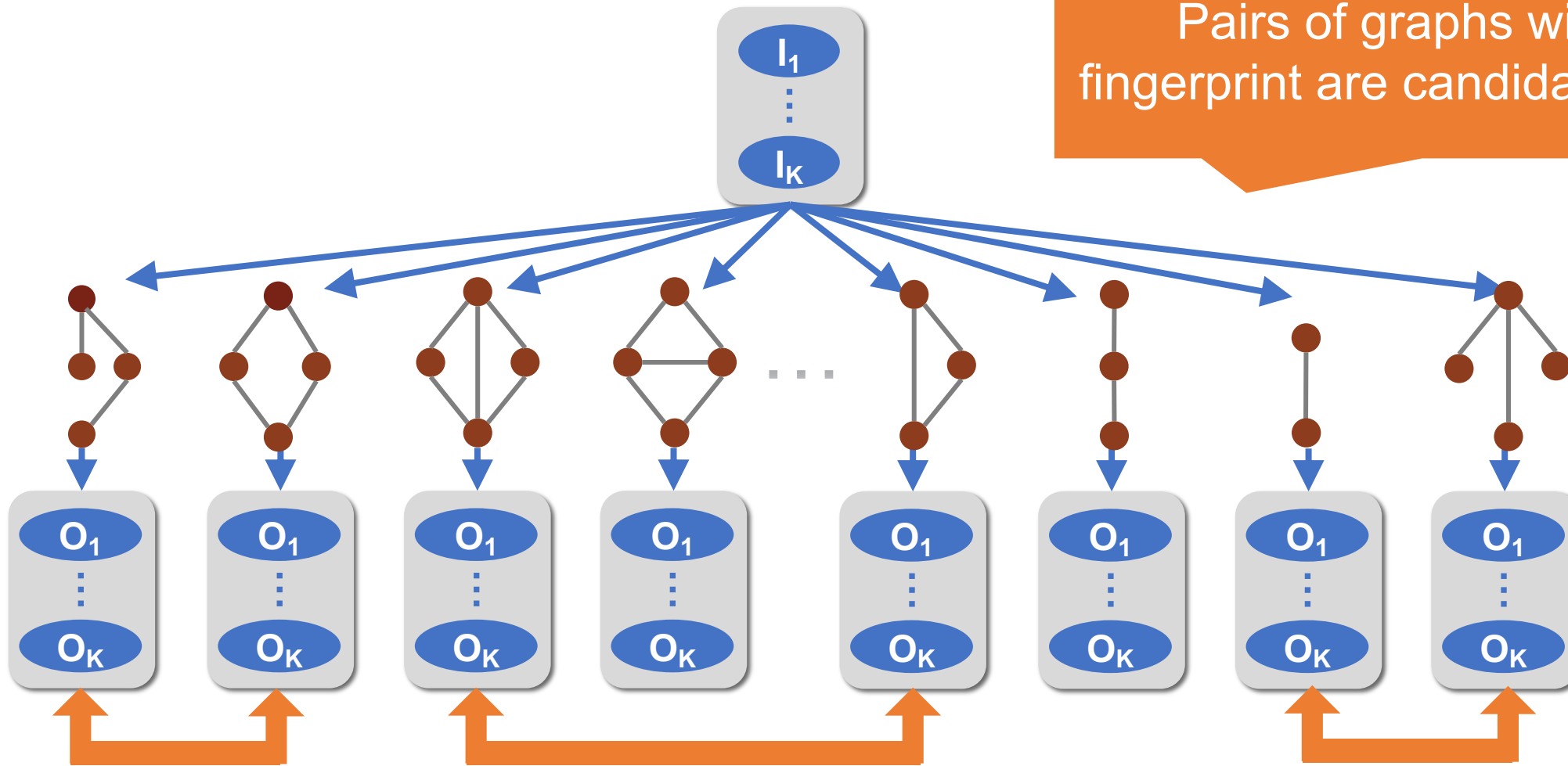
~**66M** graphs with up to 4 operators

# Graph Substitution Generator



Compute output fingerprints  
with random input tensors

# Graph Substitution Generator



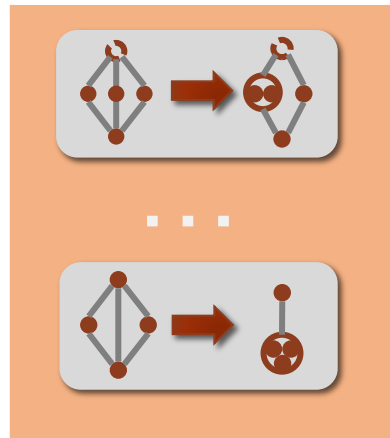
Pairs of graphs with identical fingerprint are candidate substitutions

# Graph Substitution Generator

TASO generates ~29,000 substitutions by enumerating graphs w/ up to 4 operators

743 substitutions remain after applying pruning techniques to eliminate redundancy

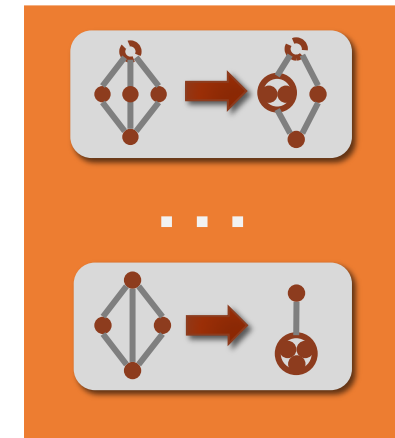
# Graph Substitution Verifier



Candidate Substitutions



Graph Subst.  
Verifier



Verified Substitutions

P1. conv is distributive  
over concatenation  
P2. conv is bilinear  
...  
Pn.

$\forall x, w_1, w_2 .$

$$\text{Conv}(x, \text{Concat}(w_1, w_2)) = \text{Concat}(\text{Conv}(x, w_1), \text{Conv}(x, w_2))$$



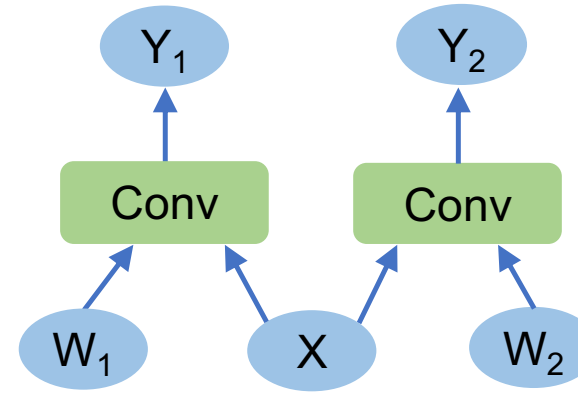
Operator Specifications

# Verification Workflow

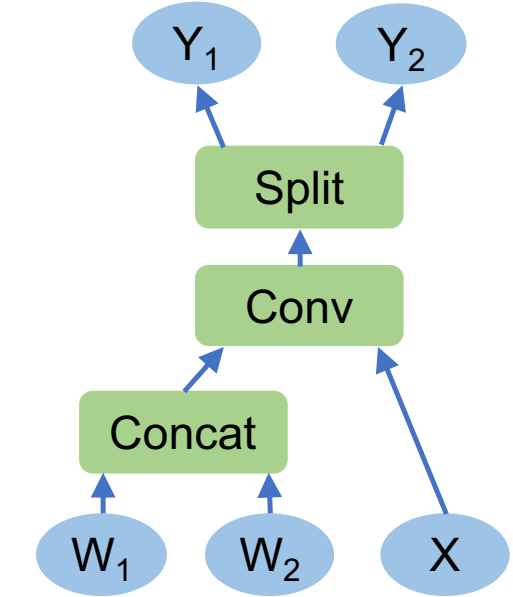
$\exists x, w_1, w_2 .$   
 $(Conv(x, w_1), Conv(x, w_2))$   
 $\neq Split(Conv(x, Concat(w_1, w_2)))$

P1.  $\forall x, w_1, w_2 .$   
 $Conv(x, Concat(w_1, w_2)) =$   
 $Concat(Conv(x, w_1), Conv(x, w_2))$   
 P2. ...

Operator Specifications



$(Conv(x, w_1), Conv(x, w_2))$



$Split(Conv(x, Concat(w_1, w_2)))$

**Theorem Prover**

UNSAT



# Verification Effort

TASO generates all **743** substitutions in 5 minutes, and verifies them against **43** operator properties in 10 minutes

Supporting a new operator requires a few hours of human effort to discover its properties

Operator specifications in TASO  $\approx$  **1,400** LOC  
Manual graph optimizations in TensorFlow  $\approx$  **53,000** LOC

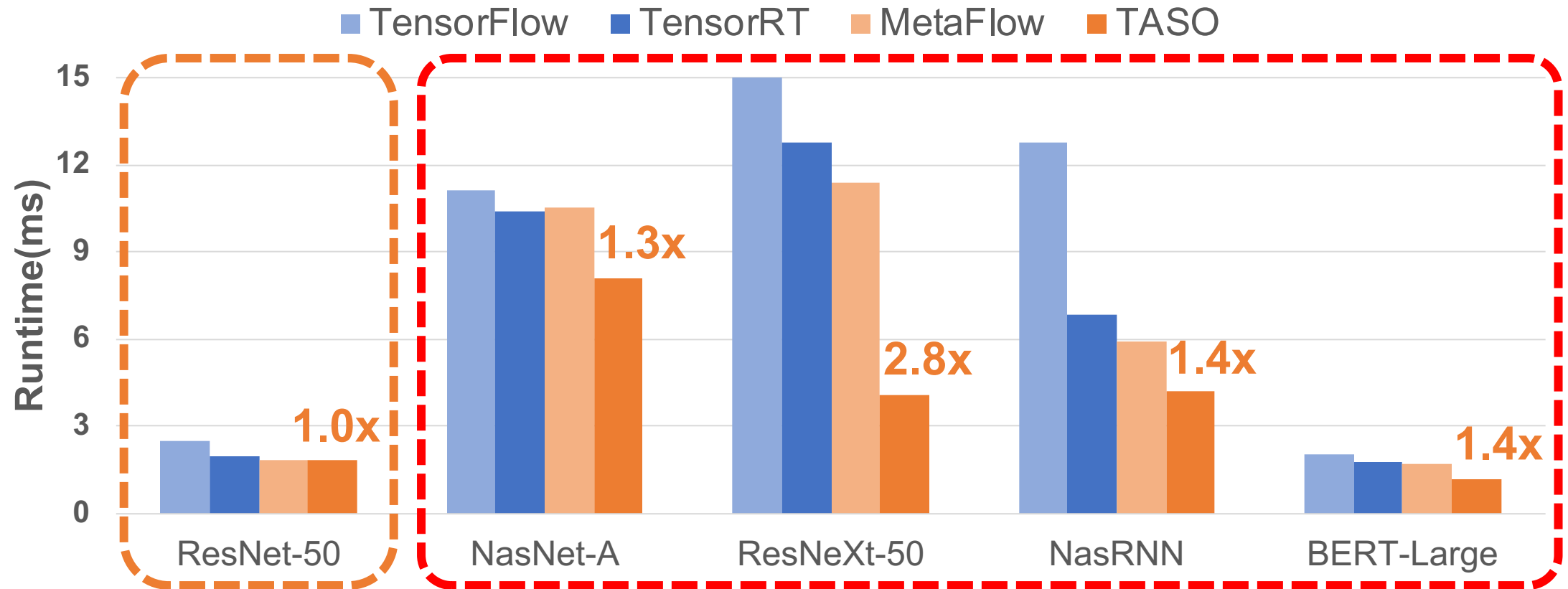
Operator Property	Comment
$\forall x, y, z. \text{ewadd}(x, \text{ewadd}(y, z)) = \text{ewadd}(\text{ewadd}(x, y), z)$	ewadd is associative
$\forall x, y. \text{ewadd}(x, y) = \text{ewadd}(y, x)$	ewadd is commutative
$\forall x, y, z. \text{ewmul}(x, \text{ewmul}(y, z)) = \text{ewmul}(\text{ewmul}(x, y), z)$	ewmul is associative
$\forall x, y. \text{ewmul}(x, y) = \text{ewmul}(y, x)$	ewmul is commutative
$\forall x, y, z. \text{ewmul}(\text{ewadd}(x, y), z) = \text{ewadd}(\text{ewmul}(x, z), \text{ewmul}(y, z))$	distributivity
$\forall x, y, w. \text{smul}(\text{smul}(x, y), w) = \text{smul}(x, \text{smul}(y, w))$	smul is associative
$\forall x, y, w. \text{smul}(\text{ewadd}(x, y), w) = \text{ewadd}(\text{smul}(x, w), \text{smul}(y, w))$	distributivity
$\forall s, p, x, y, w. \text{smul}(\text{conv}(s, p, \text{none}, x, y), w) = \text{conv}(s, p, \text{none}, \text{smul}(x, w), y)$	conv is bilinear
$\forall s, p, x, y, z. \text{conv}(s, p, \text{none}, x, \text{ewadd}(y, z)) = \text{ewadd}(\text{conv}(s, p, \text{none}, x, y), \text{conv}(s, p, \text{none}, x, z))$	conv is bilinear
$\forall a, x, y. \text{split}_0(a, \text{concat}(a, x, y)) = x$	split definition
$\forall s, p, c, x, y, z. \text{concat}(1, \text{conv}(s, p, c, x, y), \text{conv}(s, p, c, x, z)) = \text{conv}(s, p, c, \text{concat}(0, x, y))$	concatenation and conv.
$\forall s, p, x, y, z, w. \text{conv}(s, p, \text{none}, \text{concat}(1, x, z), \text{concat}(1, y, w)) = \text{ewadd}(\text{conv}(s, p, \text{none}, x, y), \text{conv}(s, p, \text{none}, z, w))$	concatenation and conv.
$\forall k, s, p, x, y. \text{concat}(1, \text{pool}_{\text{avg}}(k, s, p, x), \text{pool}_{\text{avg}}(k, s, p, y)) = \text{pool}_{\text{avg}}(k, s, p, \text{concat}(1, x, y))$	concatenation and pooling
$\forall k, s, p, x, y. \text{concat}(0, \text{pool}_{\text{max}}(k, s, p, x), \text{pool}_{\text{max}}(k, s, p, y)) = \text{pool}_{\text{max}}(k, s, p, \text{concat}(0, x, y))$	concatenation and pooling
$\forall k, s, p, x, y. \text{concat}(1, \text{pool}_{\text{max}}(k, s, p, x), \text{pool}_{\text{max}}(k, s, p, y)) = \text{pool}_{\text{max}}(k, s, p, \text{concat}(1, x, y))$	concatenation and pooling



# Search-Based Graph Optimizer (MetaFlow [SysML19])

- **Goal:** applying verified substitutions to obtain an optimized graph
- **Cost model**
  - Based on the sum of individual operators' cost
  - Measure the cost of each operator on hardware
- **Cost-based backtracking search**
  - Backtrack local optimal solutions
  - Optimizing a DNN model takes less than **10** minutes

# End-to-end Inference Performance (V100 GPU w/ cuDNN)



Competitive on standard models

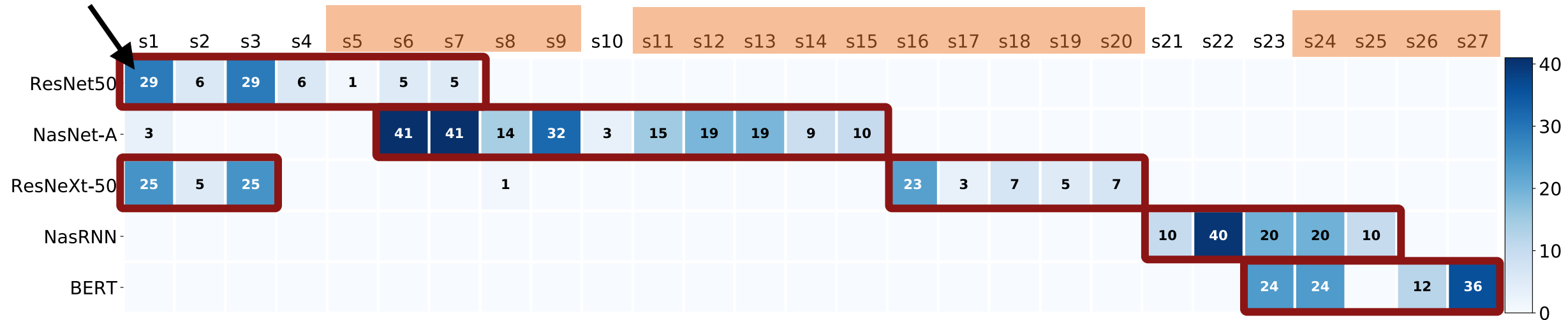
Larger speedups on emerging models

Similar speedups on the TVM backend

# Heatmap of Used Substitutions

Not covered in TensorFlow

How many times a subst. is used to optimize a DNN



Different DNN models require **different** substitutions.

# Conclusion

TASO is the first DNN optimizer that automatically generates substitutions

- Less engineering effort
- Better performance
- Formal verification

<https://github.com/jiazhihao/taso>

- Support DNN models in ONNX, TensorFlow, and PyTorch

