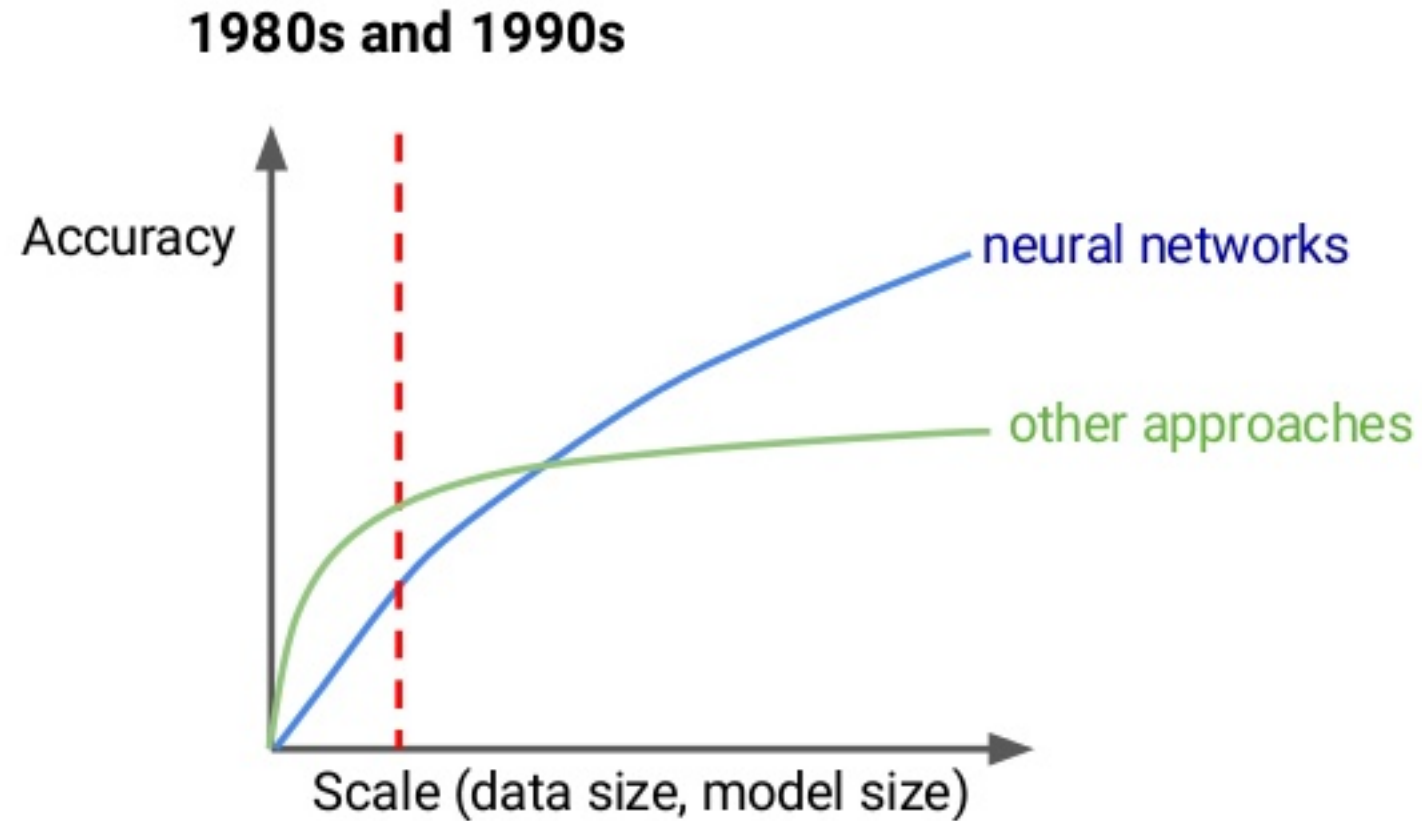


Automatically Discovering Systems Optimizations for Deep Learning

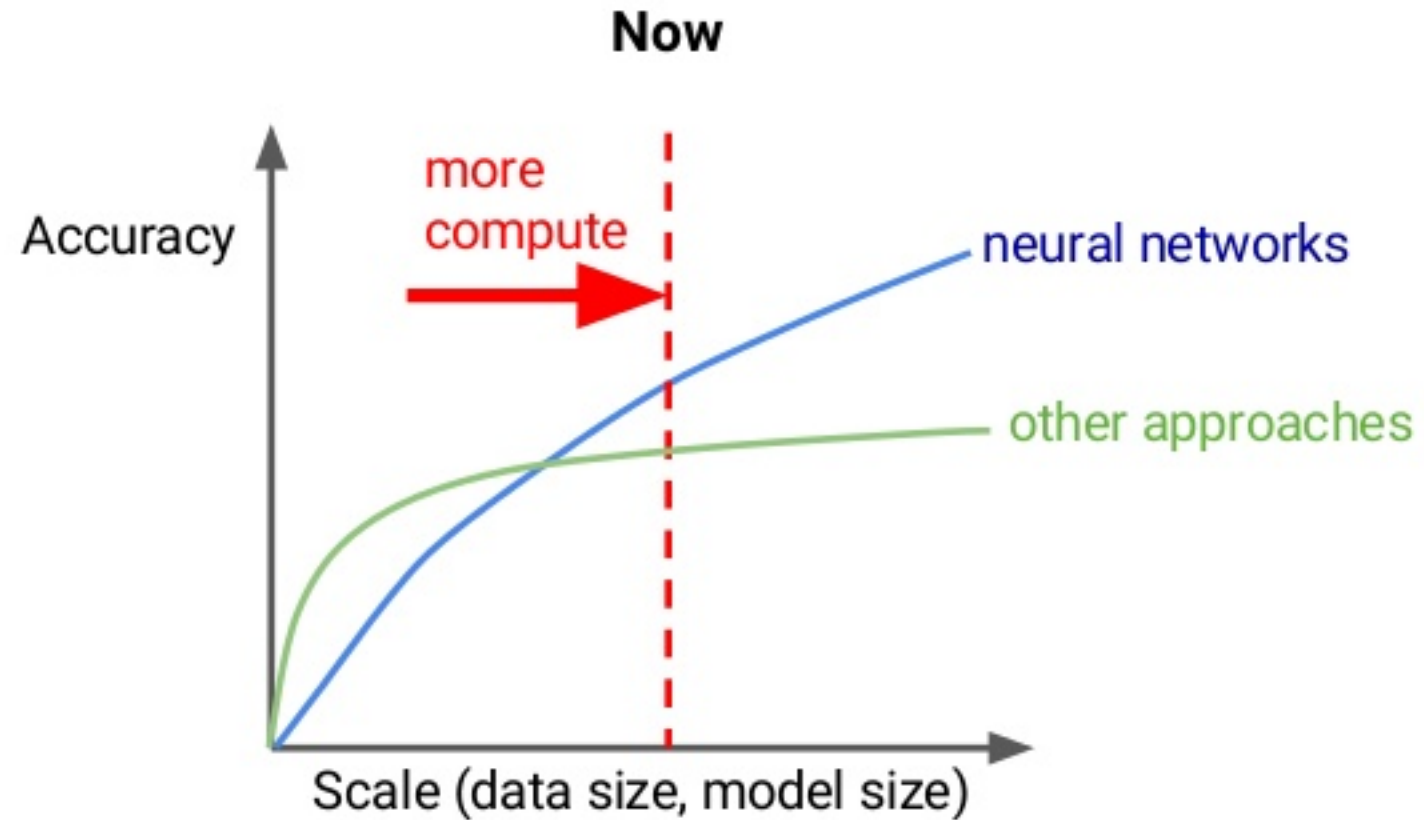
Zhihao Jia

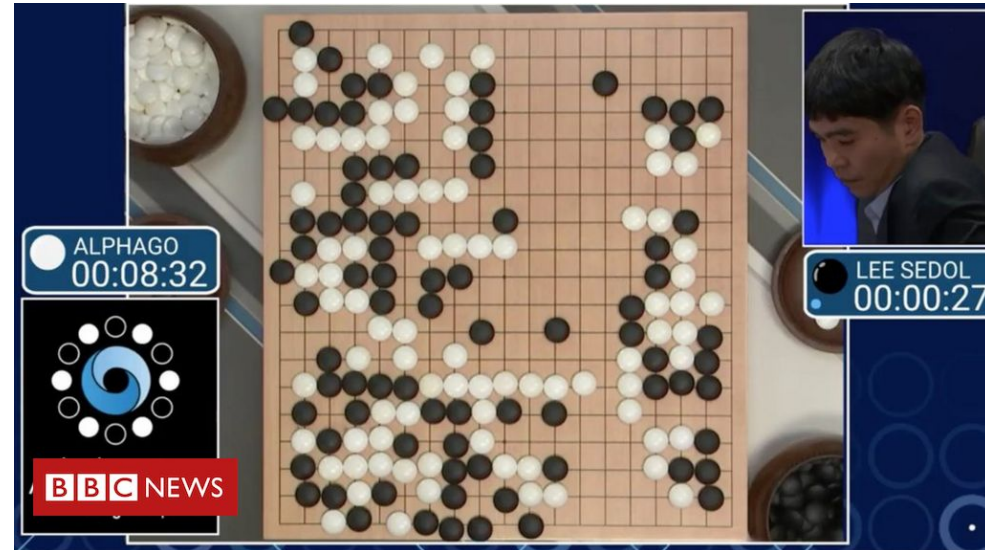
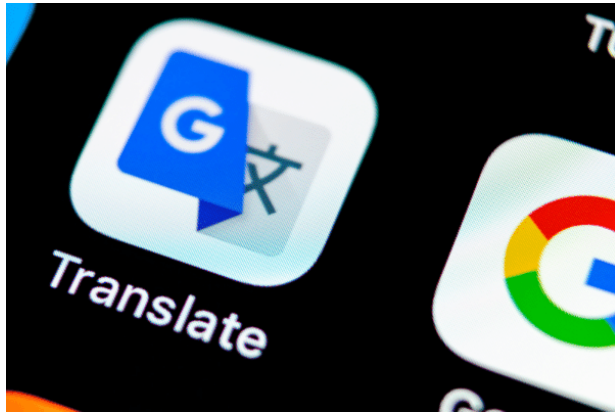
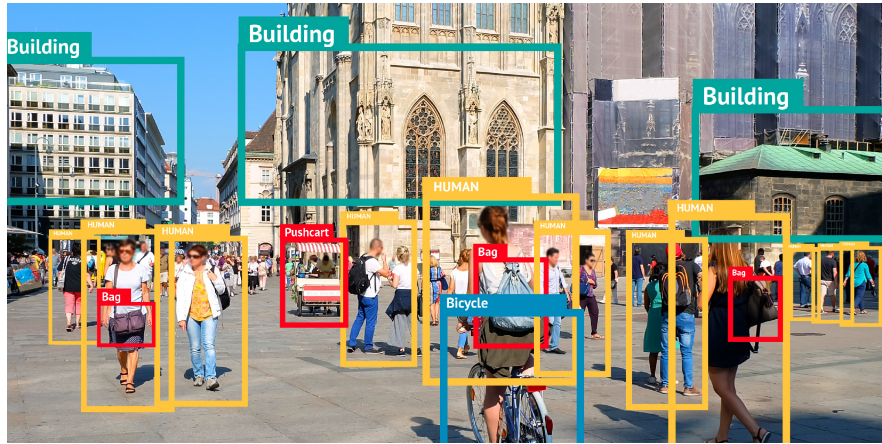
CMU and Facebook

The Rise of ML and Neural Networks

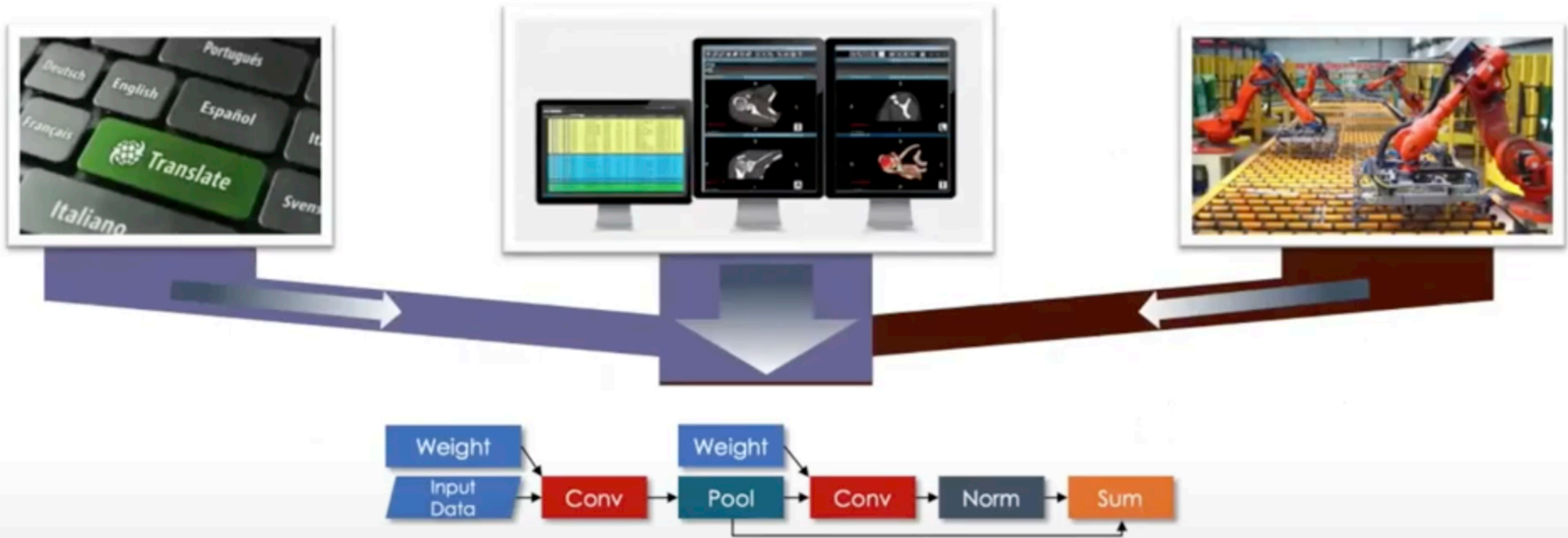


The Rise of ML and Neural Networks





Deep Neural Networks for Machine Translation

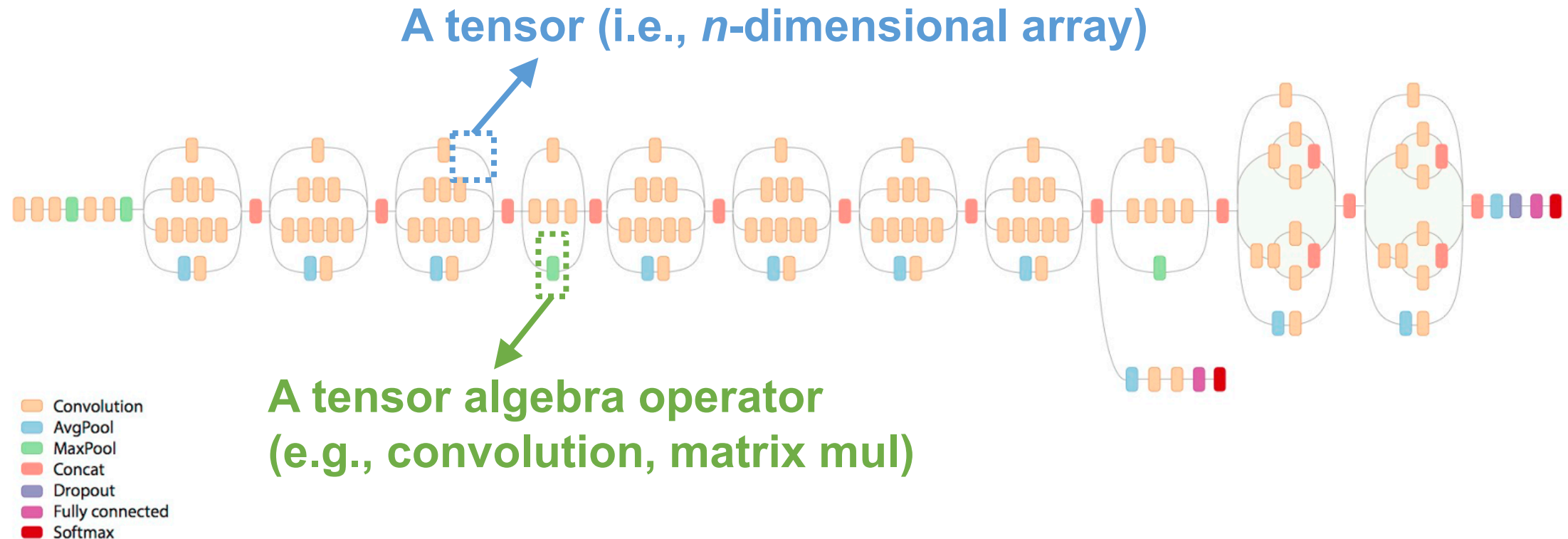


1000x Productivity

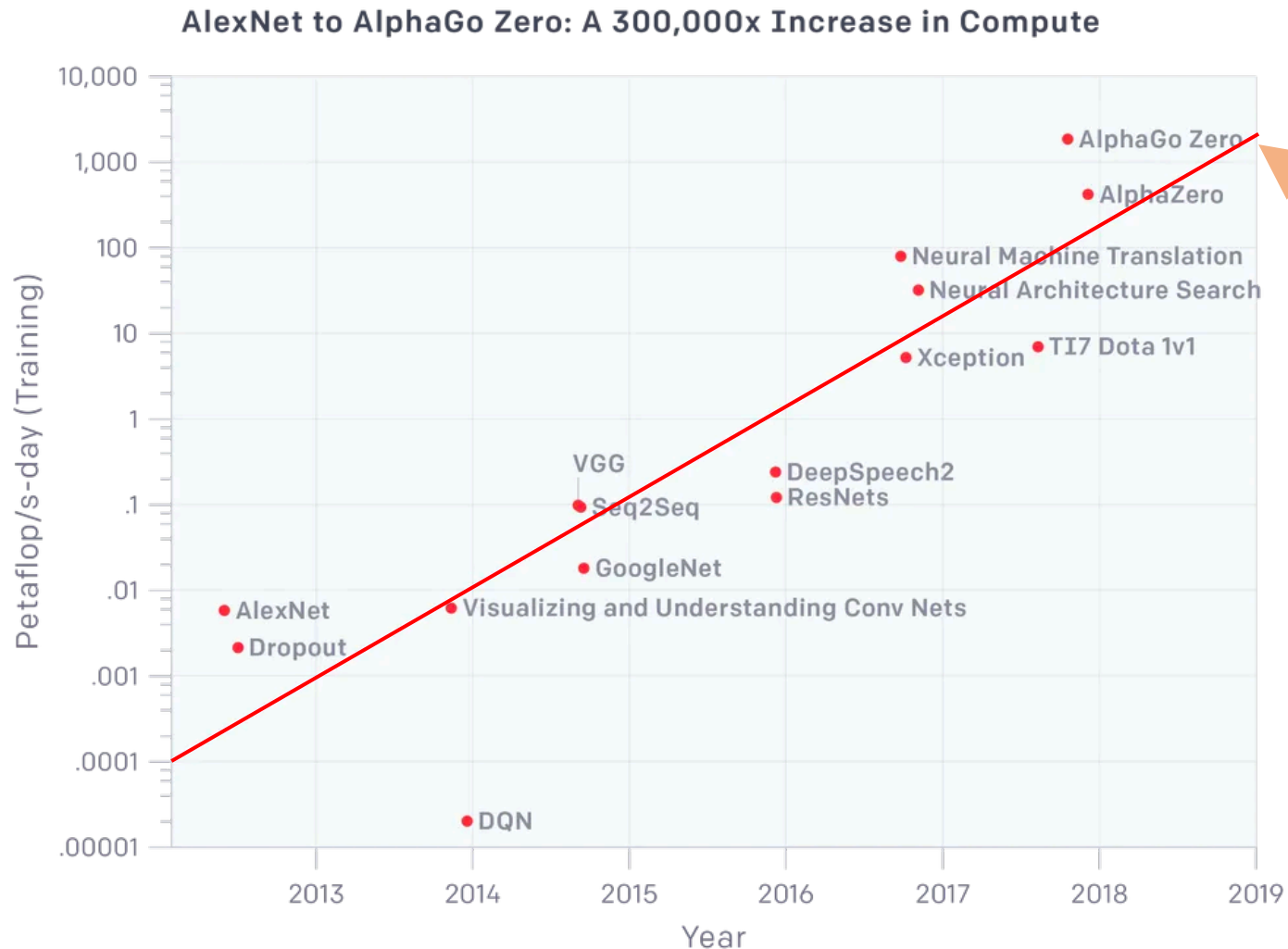
Google shrinks language translation code from 500k imperative LoC to **500 lines of dataflow**

What is a Deep Neural Network?

- Collection of simple trainable mathematical units that work together to solve complicated tasks

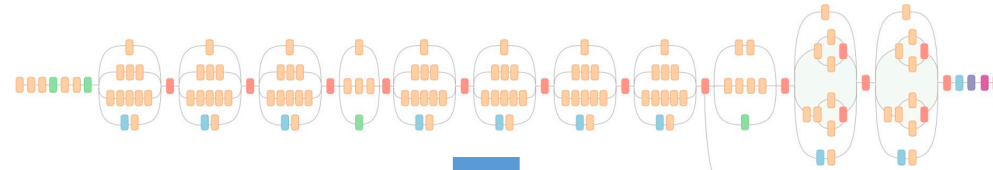


ML Computation is Increasing Exponentially



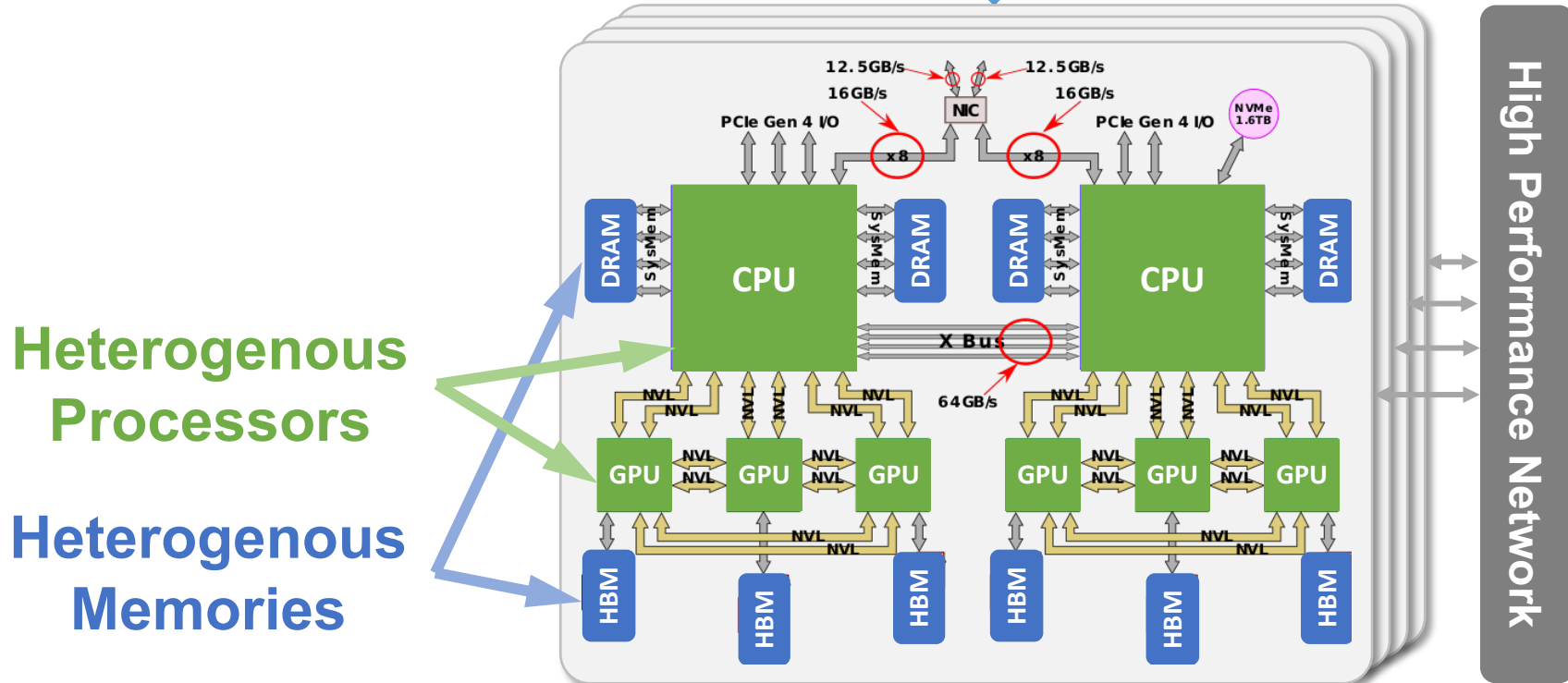
Amount of computation in largest ML training doubles every 3.4 months

ML Systems are a Key Ingredient in ML



ML Model

ML Systems

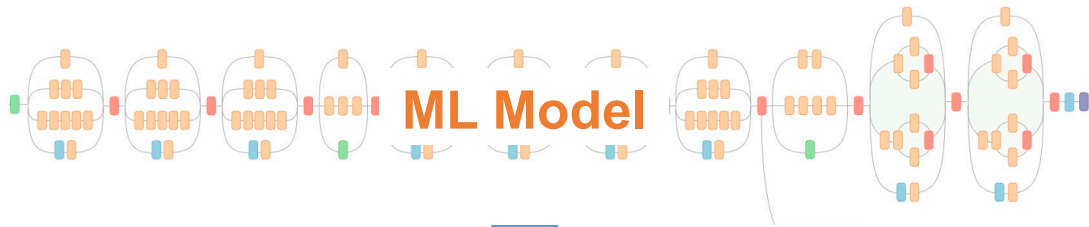


High Performance Network

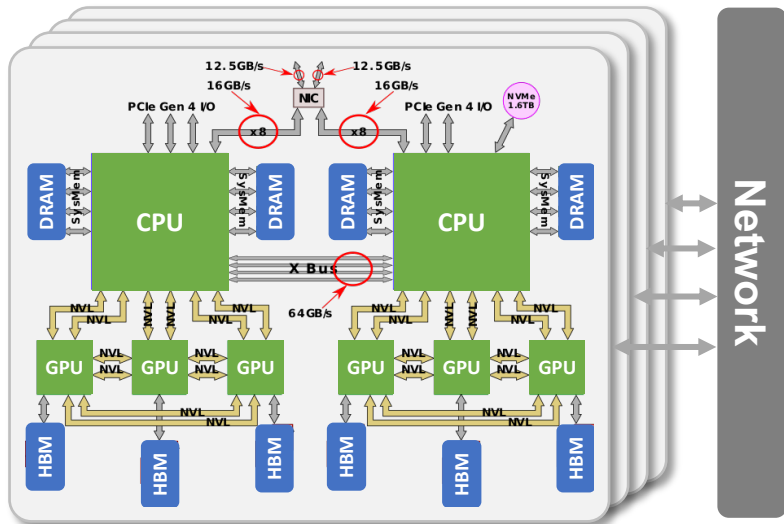
Distributed Heterogenous Hardware Architectures

~4,600 compute nodes

Challenges of Building ML Systems



ML Systems



Massively Parallel
Tensor algebra is parallelizable in many dimensions

New ML Operators
Continuously introduced into ML systems

Heterogenous Hardware
Different processor kinds and complex memory hierarchy

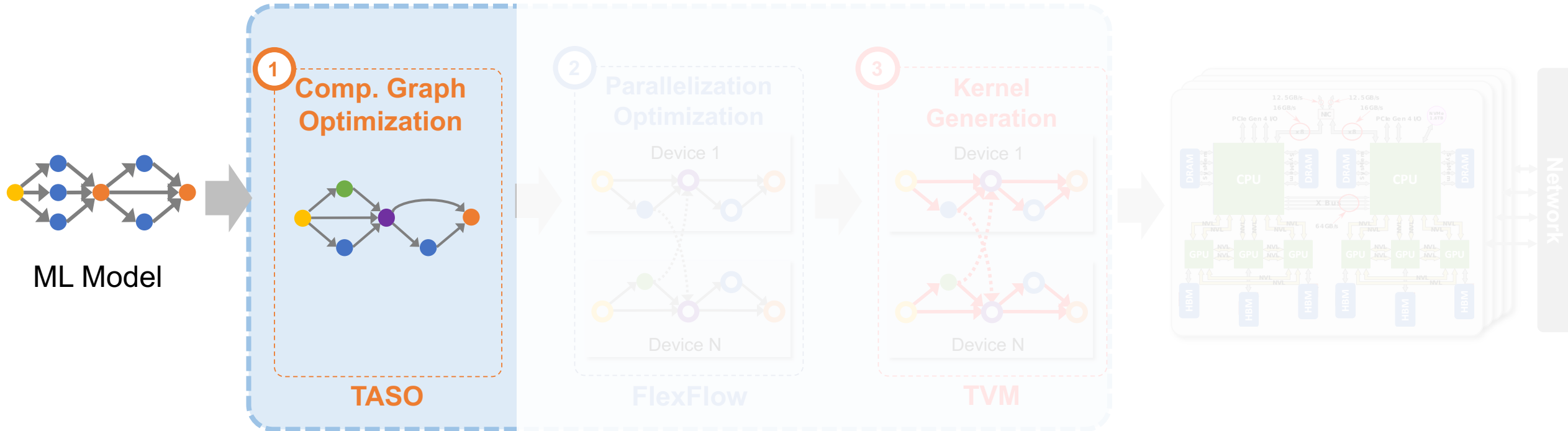
CMU Automated Learning Systems Lab

Mission: Automate the design and optimization of ML systems by leveraging

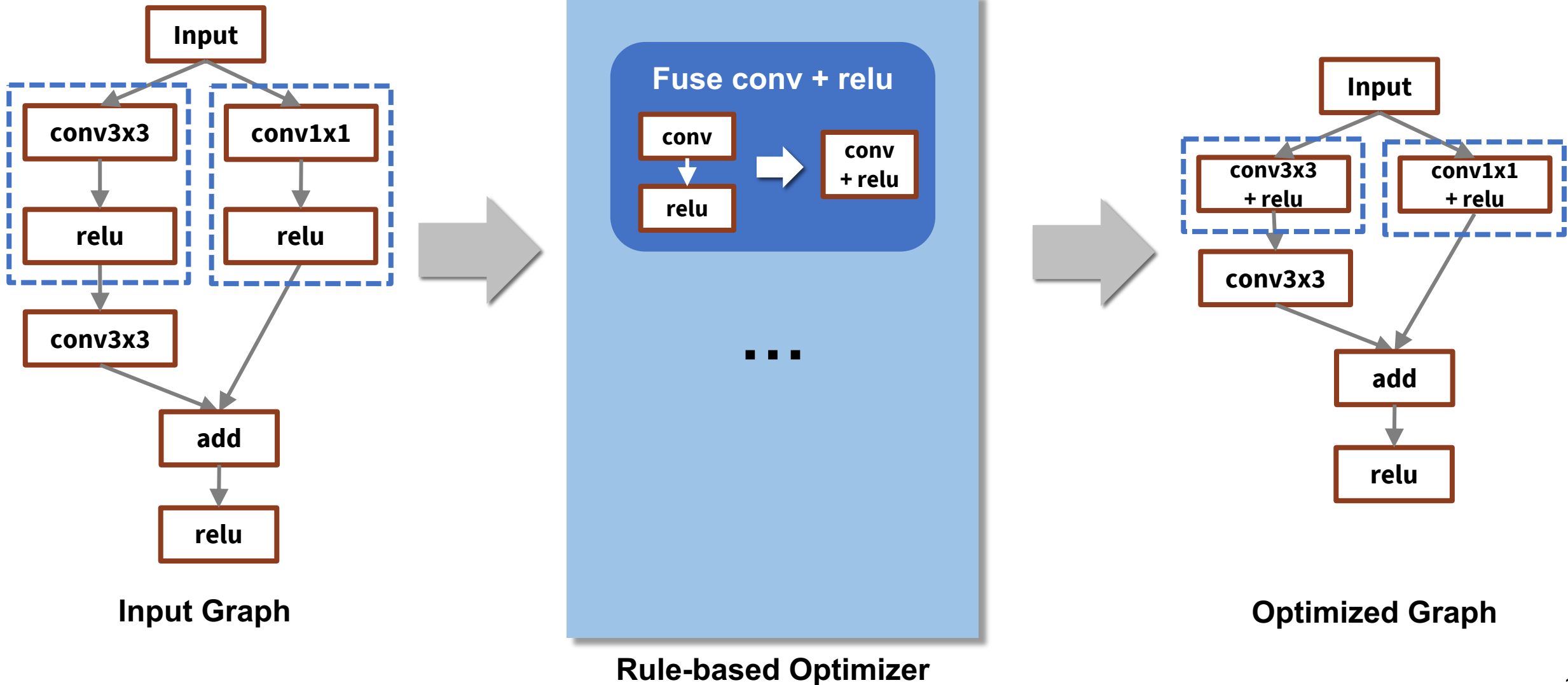
1. Statistical and mathematical properties of ML algorithms
2. Domain knowledge of modern hardware platforms

CMU Automated Learning Systems Lab

Automated ML Systems

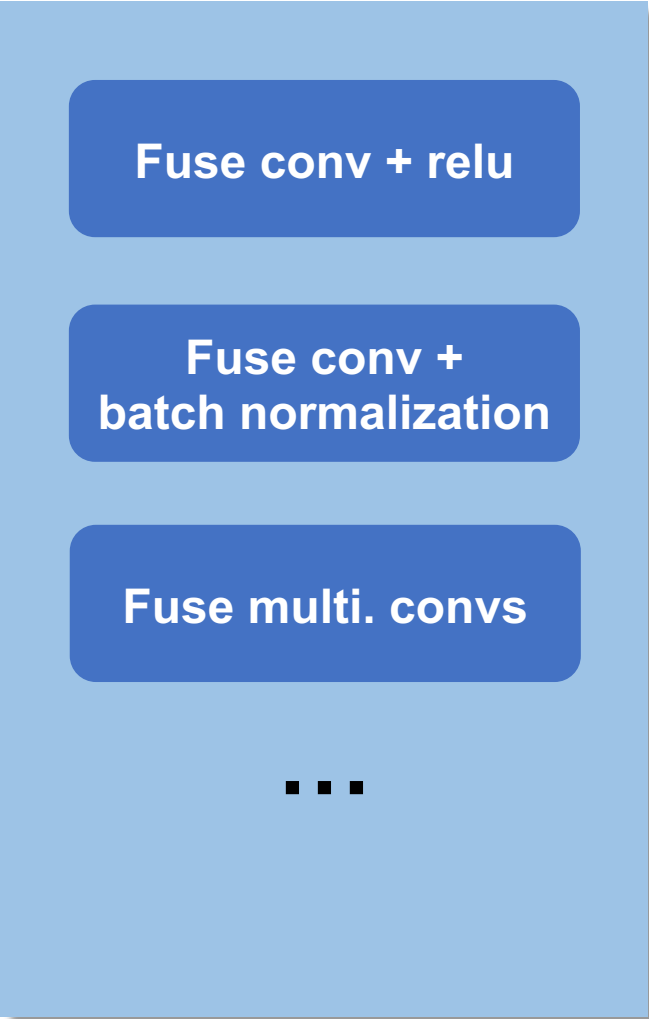


Current Rule-based Graph Optimizations



Current Rule-based Graph 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       },
46       {"Const"}, // mul_values_node
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 models/hardware

The screenshot shows a GitHub issue page for 'Horovod with XLA is slower than without XLA (Tensorflow 1.12) #713'. The issue is closed and was opened by LiweiPeng on Dec 19, 2018. A comment from LiweiPeng, dated Dec 19, 2018, describes the problem: '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 is labeled as a 'question' and has no assignees or milestones.

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:

```
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 explains: '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 models/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 models/hardware

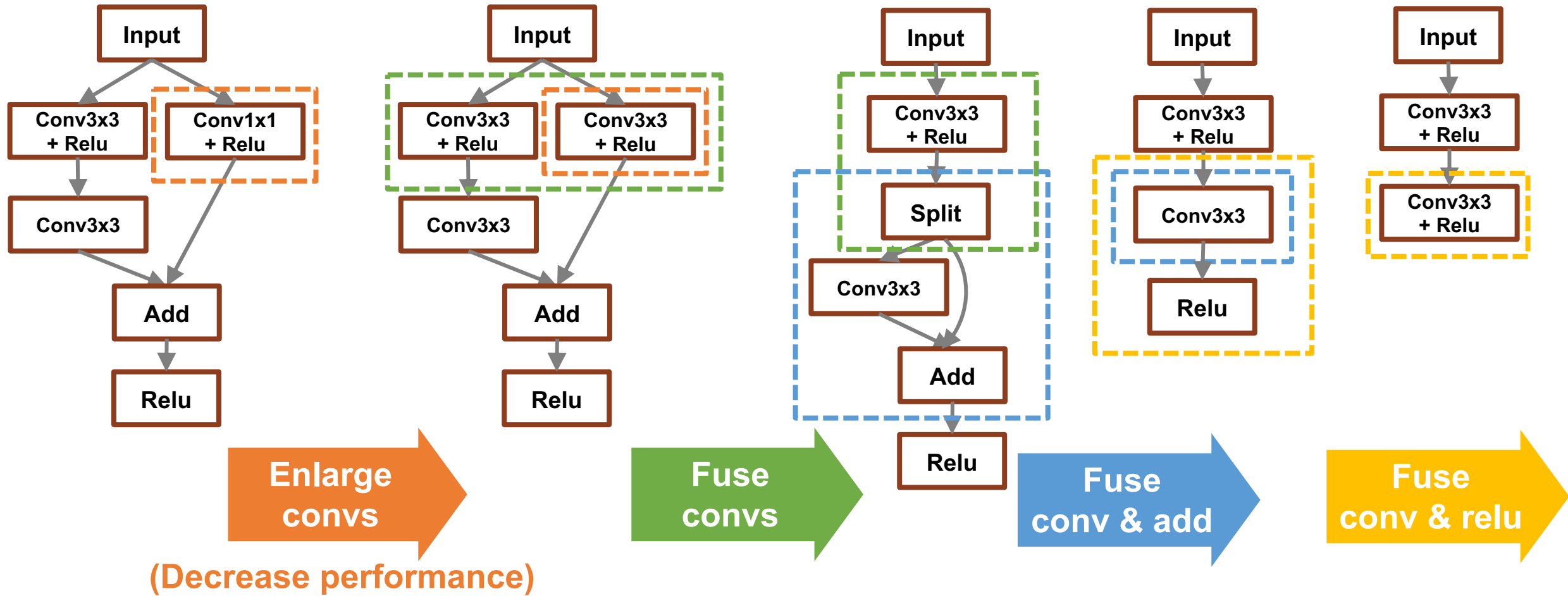
Scalability

New operators and graph structures require more rules

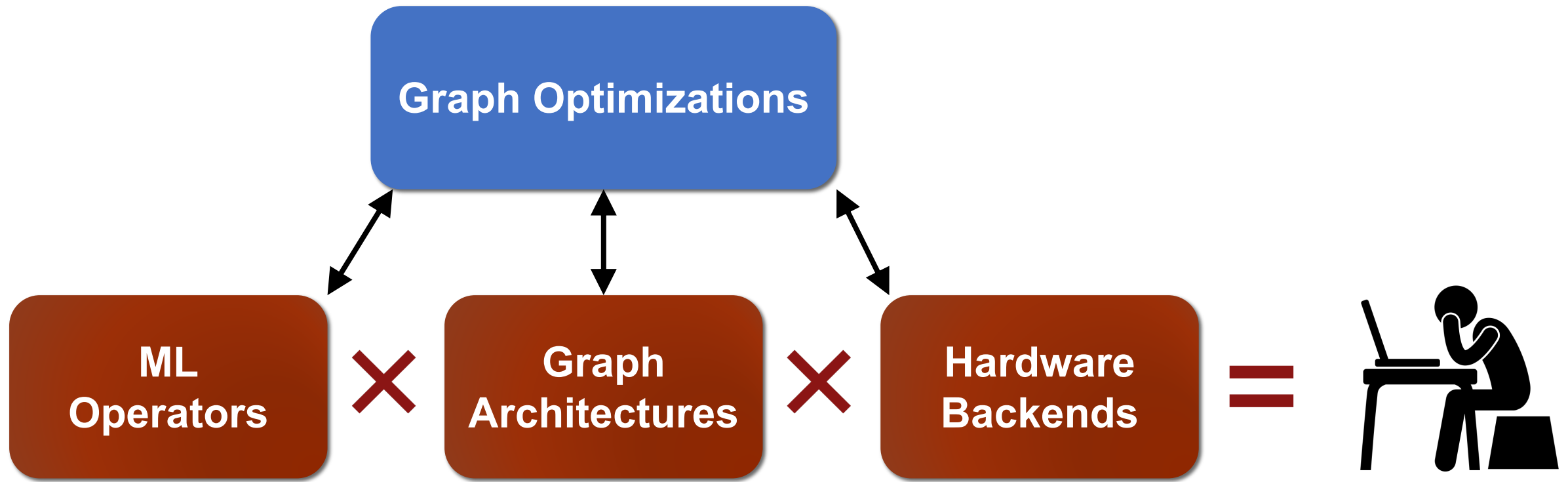
Performance

Miss subtle optimizations for specific models/hardware

Motivating Example (ResNet)



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



Infeasible to manually design graph optimizations for all cases

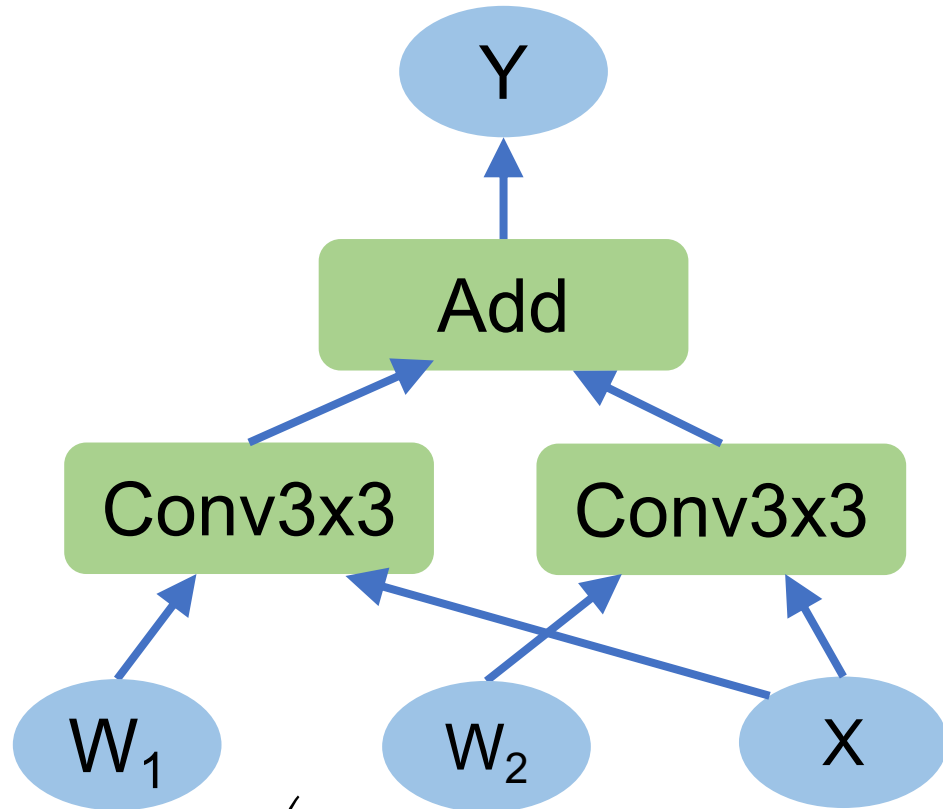
Is it possible to generate them automatically?

TASO: Tensor Algebra SuperOptimizer

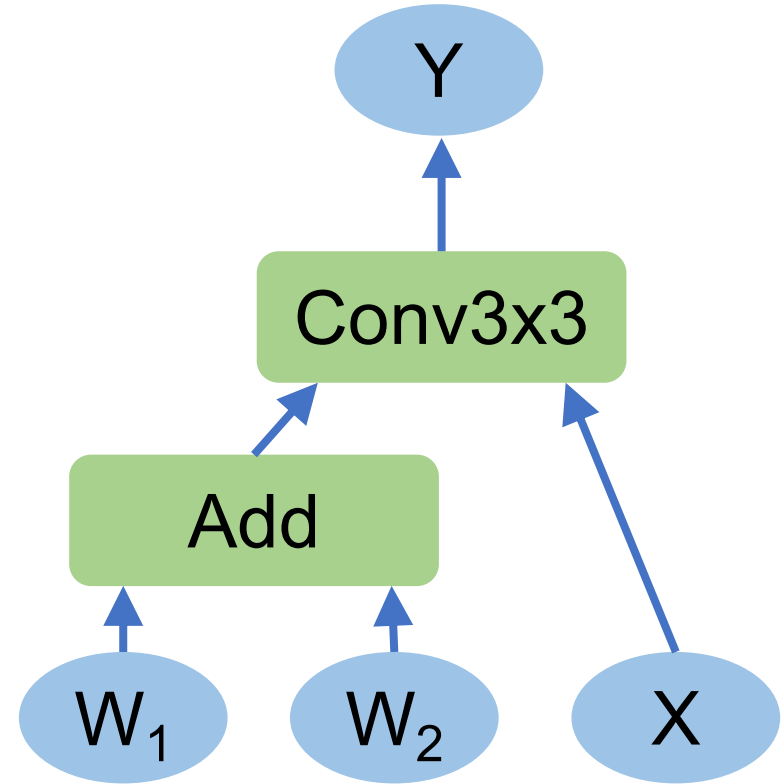
Key idea: replace manually-designed graph optimizations with *automated generation and verification* of graph substitutions for tensor algebra

- **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 3x
- **Stronger correctness:** formally verify all generated substitutions

Graph Substitution



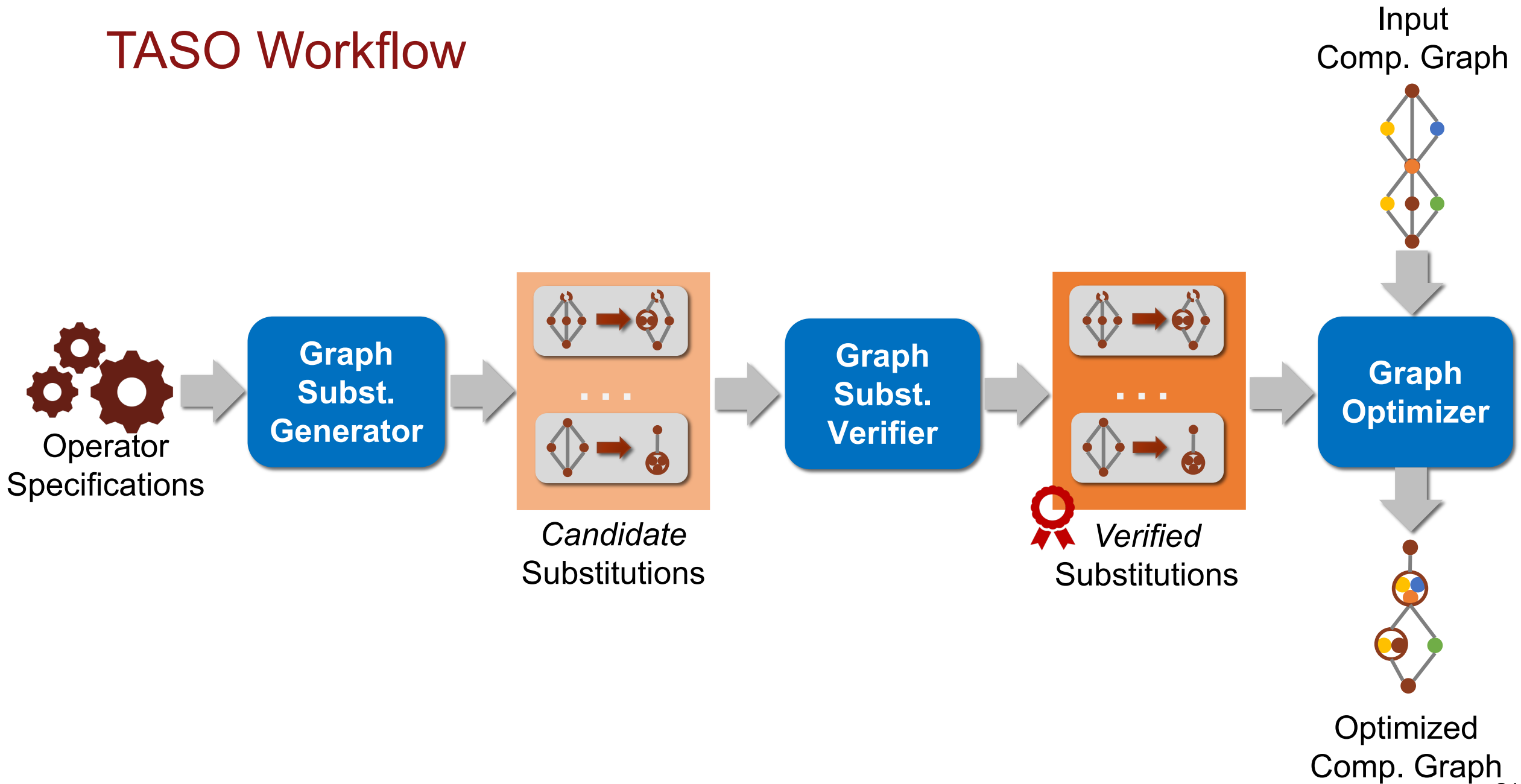
=



$$Y(n, c, h, w) = \left(\sum_{d,u,v} X(n, d, h + u, w + v) * W1(c, d, u, v) \right) + \left(\sum_{d,u,v} X(n, d, h + u, w + v) * W2(c, d, u, v) \right)$$

$$\Leftrightarrow Y(n, c, h, w) = \sum_{d,u,v} X(n, d, h + u, w + v) * ((W_1(c, d, u, v) + W_2(c, d, u, v)))$$

TASO Workflow



Key Challenges

1. How to generate potential substitutions?

Graph fingerprints

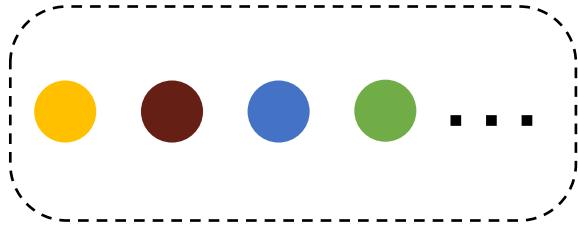
2. How to verify their correctness?

Operator specifications + theorem prover

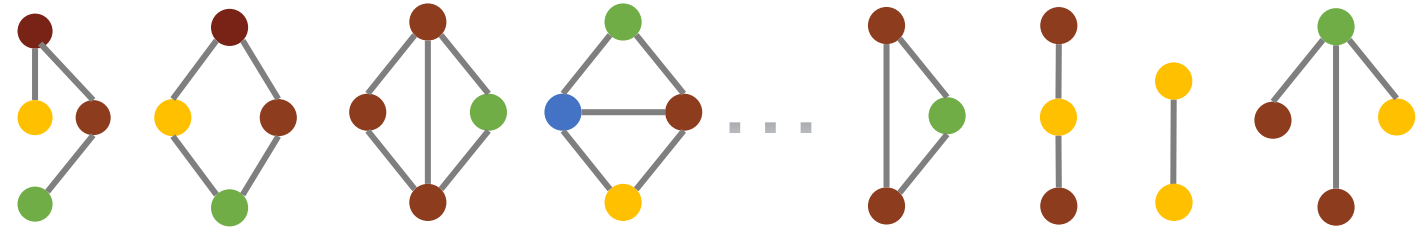
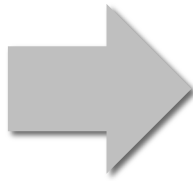
Graph Substitution Generator



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



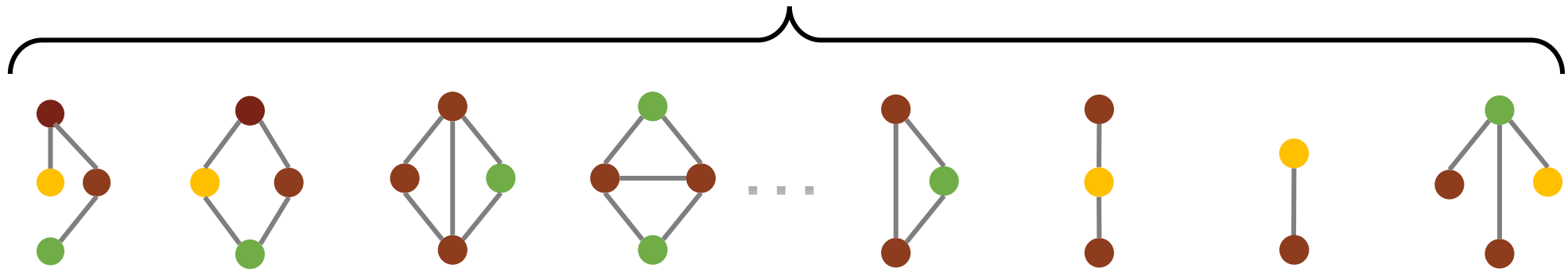
Operators supported by hardware backend





Graph Substitution Generator

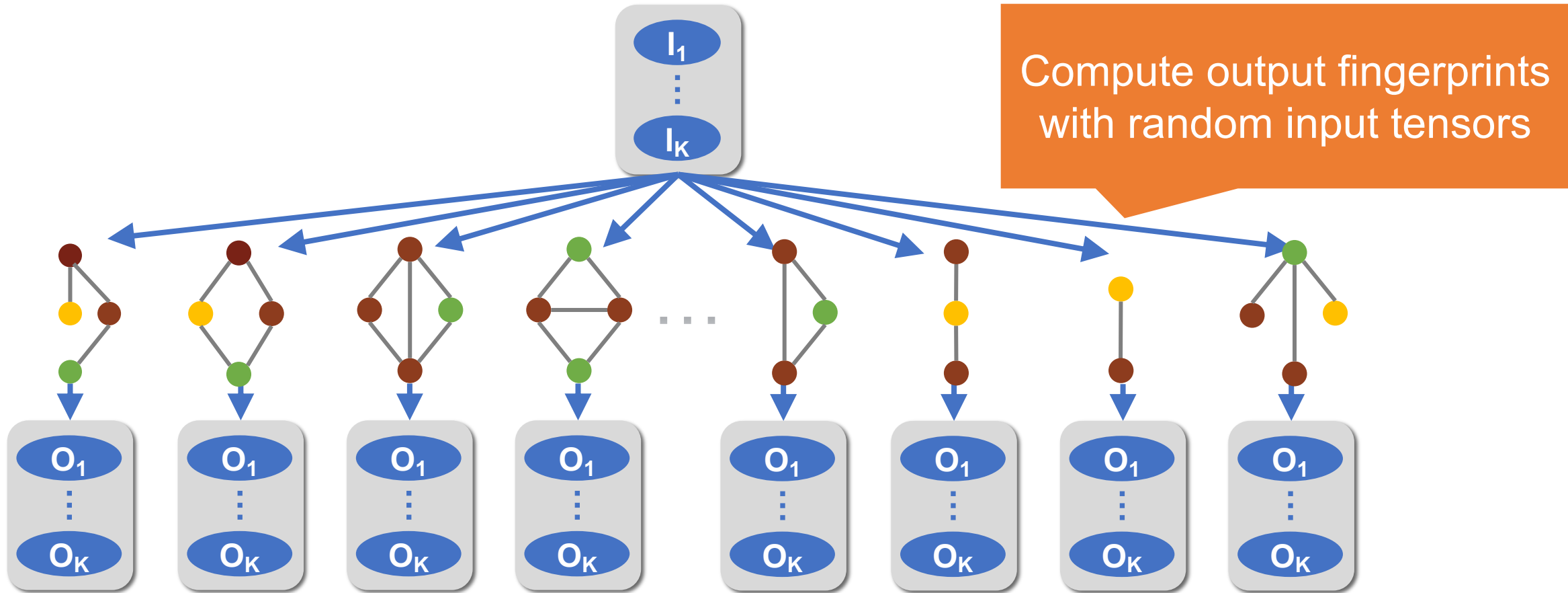
66M graphs with up to **4** operators



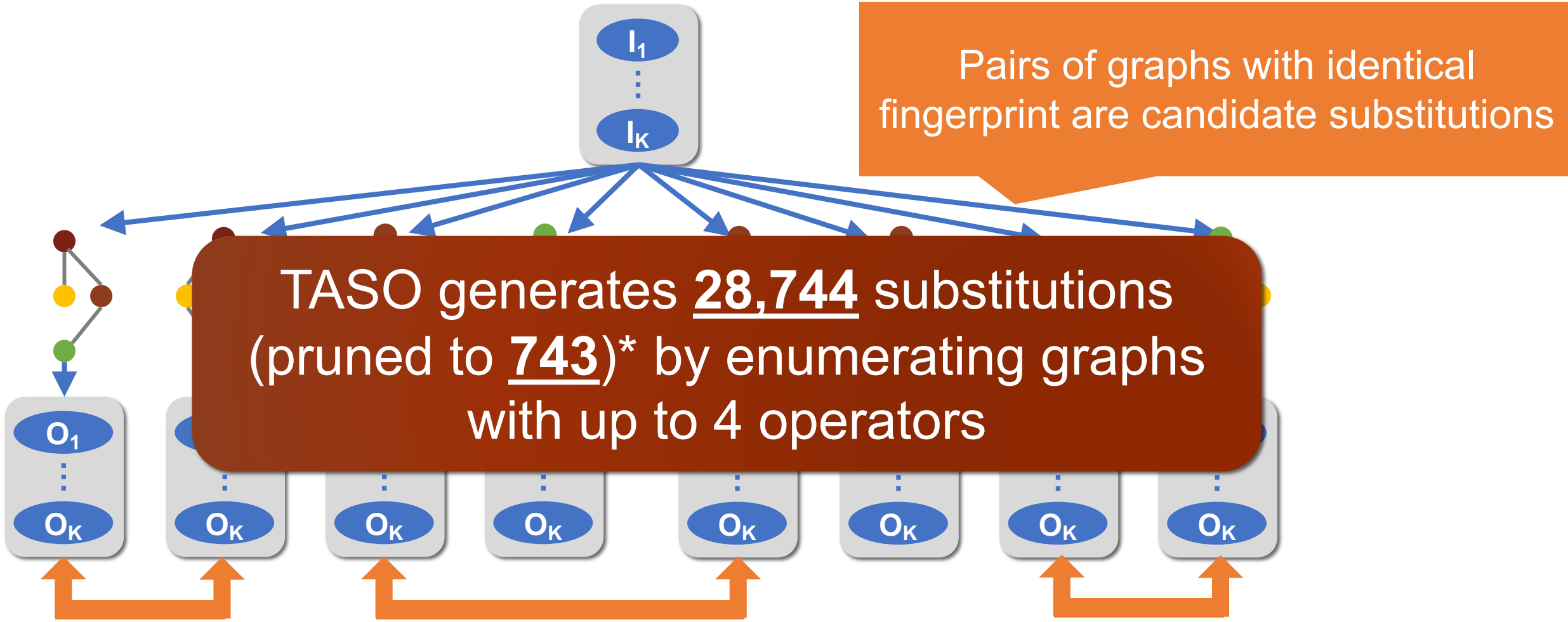
A substitution = a pair of equivalent graphs

Explicitly considering all pairs does not scale

Graph Substitution Generator

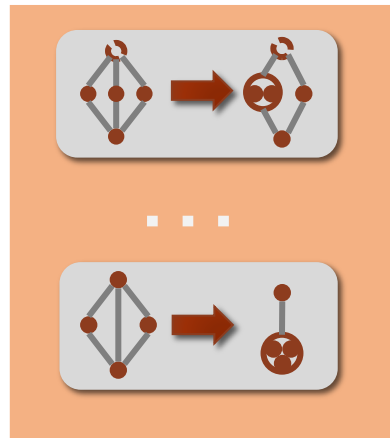


Graph Substitution Generator

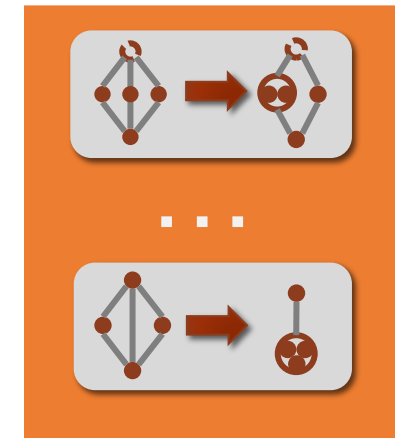


*Pruning details available in **Z. Jia et al. SOSP'19**

Graph Substitution Verifier



Candidate Substitutions



Verified Substitutions

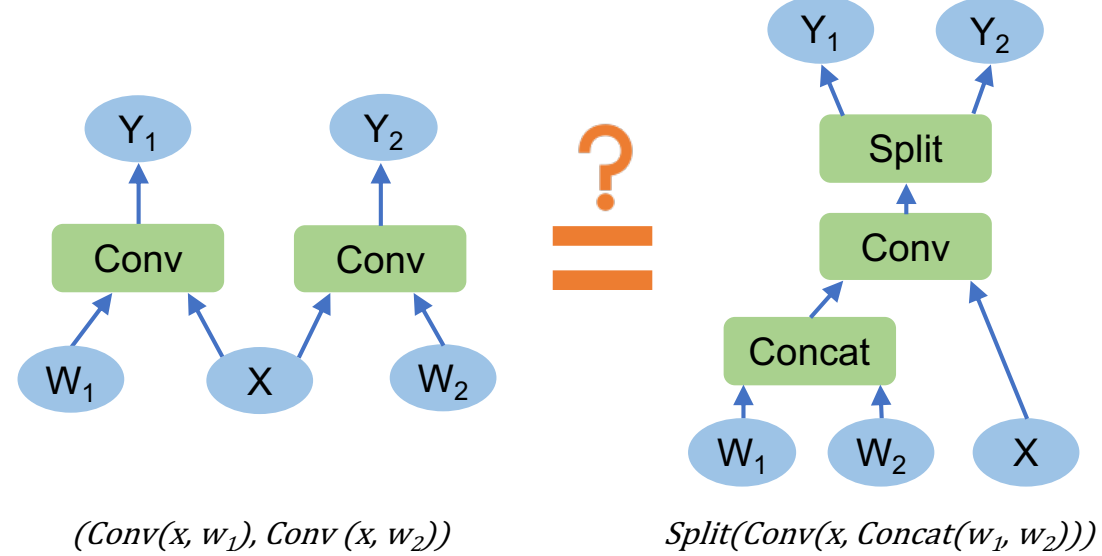
P1. conv is distributive over concatenation
P2. conv is bilinear
...
Pn.



Operator Specifications

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

Verification Workflow

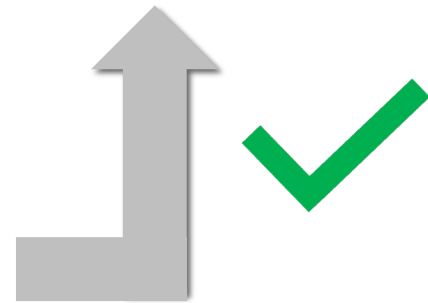


$\forall x, w_1, w_2 .$
 $(Conv(x, w_1), Conv(x, w_2))$
 $= Split(Conv(x, Concat(w_1, w_2)))$

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

Automated Theorem Prover



Verification Effort

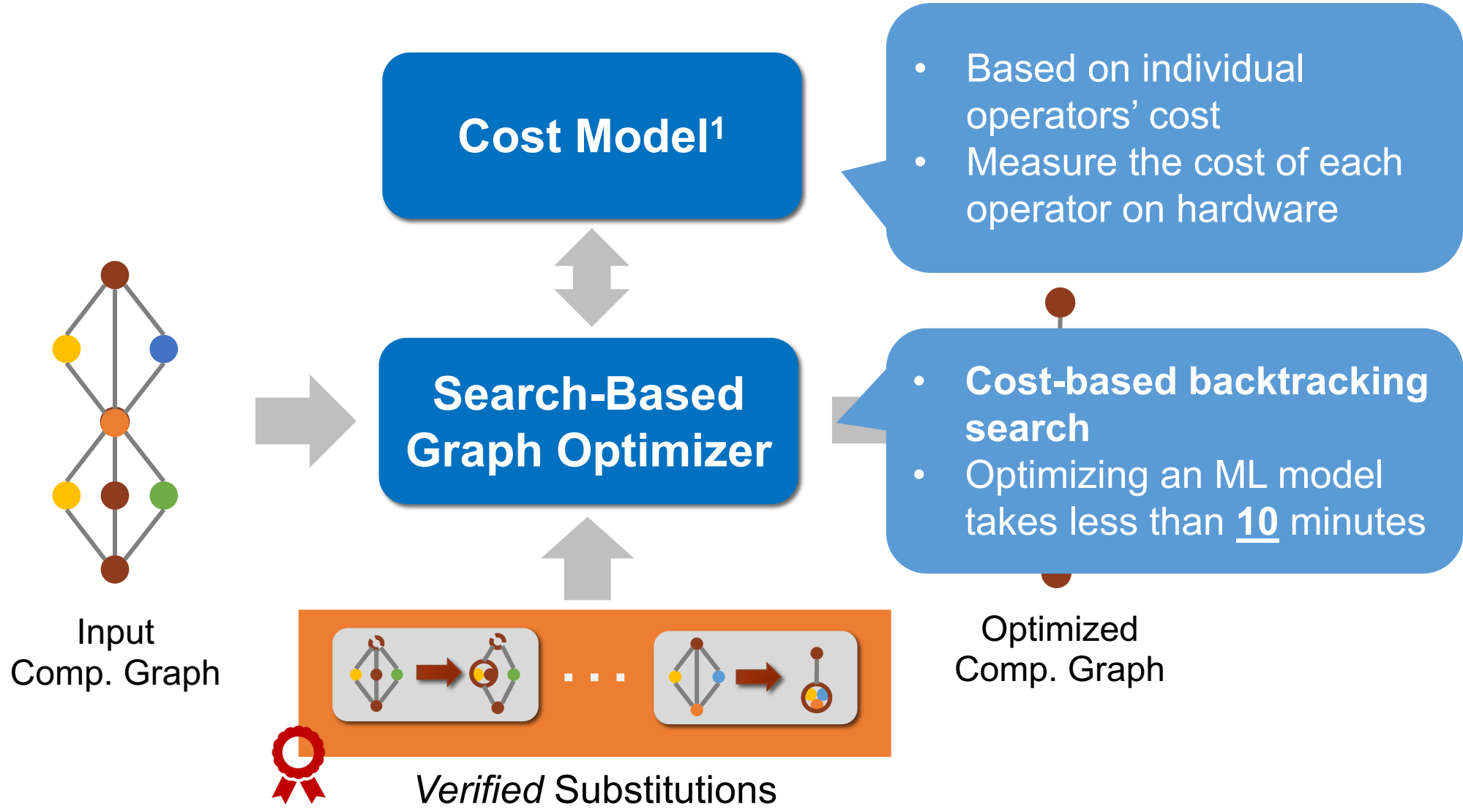
Operator Property	Comment
$\forall x, y, z. \text{ewadd}(x, \text{ewadd}(y, z)) = \text{ewadd}(\text{ewadd}(x, y), z)$	ewadd is associative
$\forall x, y. \text{ewadd}(x, y) = \text{ewadd}(y, x)$	ewadd is commutative
$\forall x, y, z. \text{ewmul}(x, \text{ewmul}(y, z)) = \text{ewmul}(\text{ewmul}(x, y), z)$	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
$\forall x, u, w. \text{smul}(\text{ewadd}(x, u), w) = \text{ewadd}(\text{smul}(x, w), \text{smul}(u, w))$	distributivity
	commutativity
	is its own inverse
	commutativity
	commutativity
	commutativity
	associative
	linear
	linear
	and transpose
	linear
$\forall s, p, x, y, w. \text{smul}(\text{conv}(s, p, A_{\text{none}}, x, y), w) = \text{conv}(s, p, A_{\text{none}}, \text{smul}(x, w), y)$	conv is bilinear
$\forall s, p, x, y, z. \text{conv}(s, p, A_{\text{none}}, x, \text{ewadd}(y, z)) = \text{ewadd}(\text{conv}(s, p, A_{\text{none}}, x, y), \text{conv}(s, p, A_{\text{none}}, x, z))$	conv is bilinear
	linear
	convolution kernel
	A _{relu} applies relu
	commutativity
	conv. with C _{pool}
	kernel
	matrix
	identity
$\forall a, x, y. \text{split}_0(a, \text{concat}(a, x, y)) = x$	split definition
	definition
	of concatenation
	commutativity
	commutativity
	commutativity
	commutativity
	tion and transpose
	tion and matrix mul.
	tion and matrix mul.
	tion and conv.
$\forall s, p, c, x, y, z. \text{concat}(1, \text{conv}(s, p, c, x, y), \text{conv}(s, p, c, x, z)) = \text{conv}(s, p, c, x, \text{concat}(0, y, z))$	concatenation and conv.
$\forall s, p, x, y, z, w. \text{conv}(s, p, A_{\text{none}}, \text{concat}(1, x, z), \text{concat}(1, y, w)) =$ $\text{ewadd}(\text{conv}(s, p, A_{\text{none}}, x, y), \text{conv}(s, p, A_{\text{none}}, z, w))$	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

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

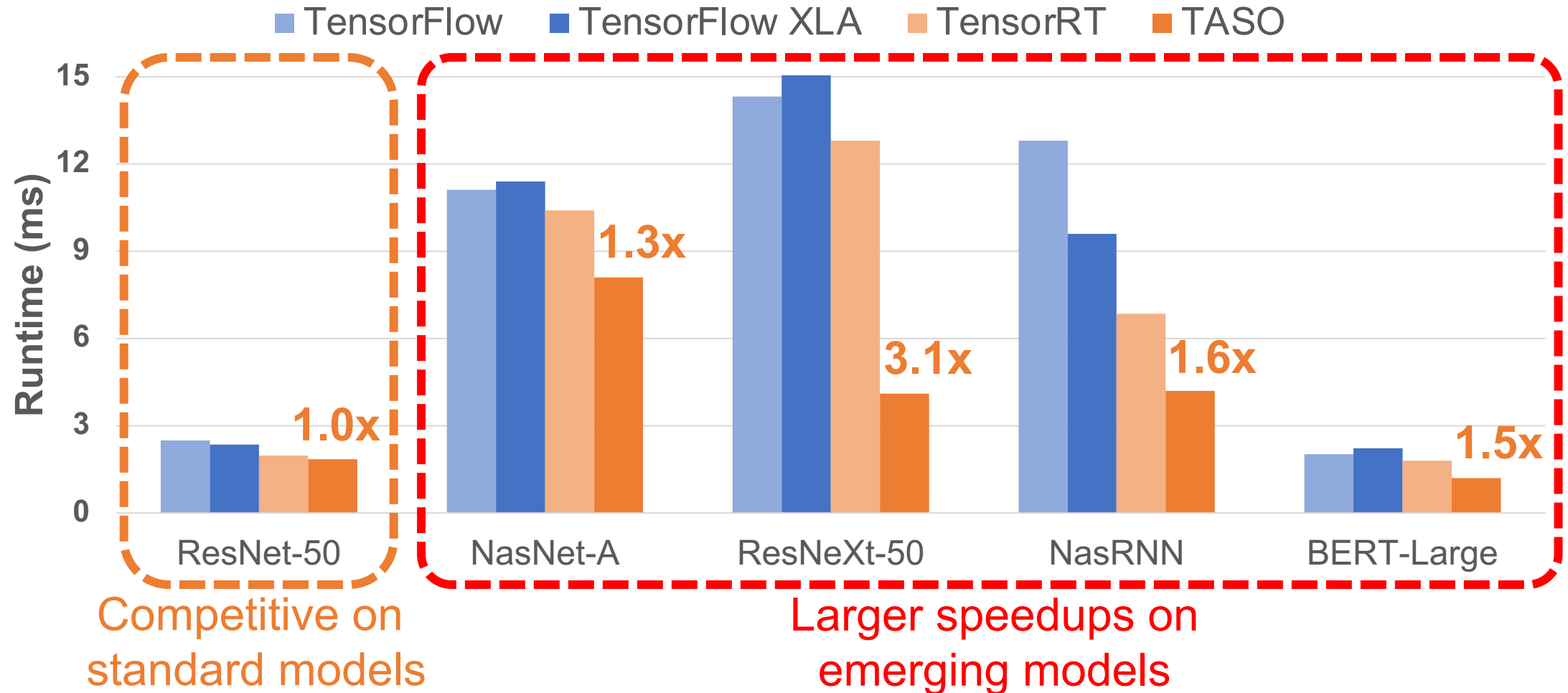
Supporting a new operator requires a few hours of human effort to specify its properties

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

Search-Based Graph Optimizer



End-to-end Inference Performance (Nvidia V100 GPU)



TASO



First DNN graph optimizer that automatically generates substitutions

- Less engineering effort
- Better runtime performance
- Stronger correctness guarantee



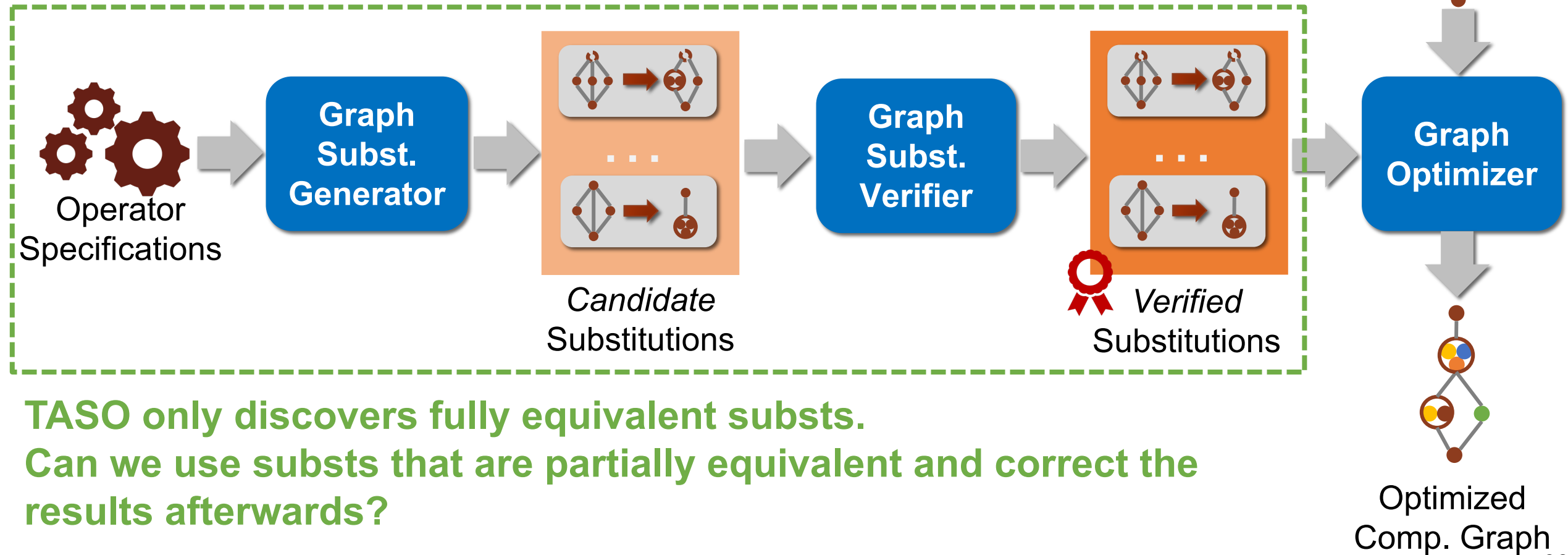
ONNX

ByteDance



1. TASO: Optimizing Deep Learning Computation with Automated Generation of Graph Substitutions. SOSP'19.
2. Optimizing DNN Computation with Relaxed Graph Substitutions. MLSys'19.
3. Exploring Hidden Dimensions in Parallelizing Convolutional Neural Networks. ICML'18.

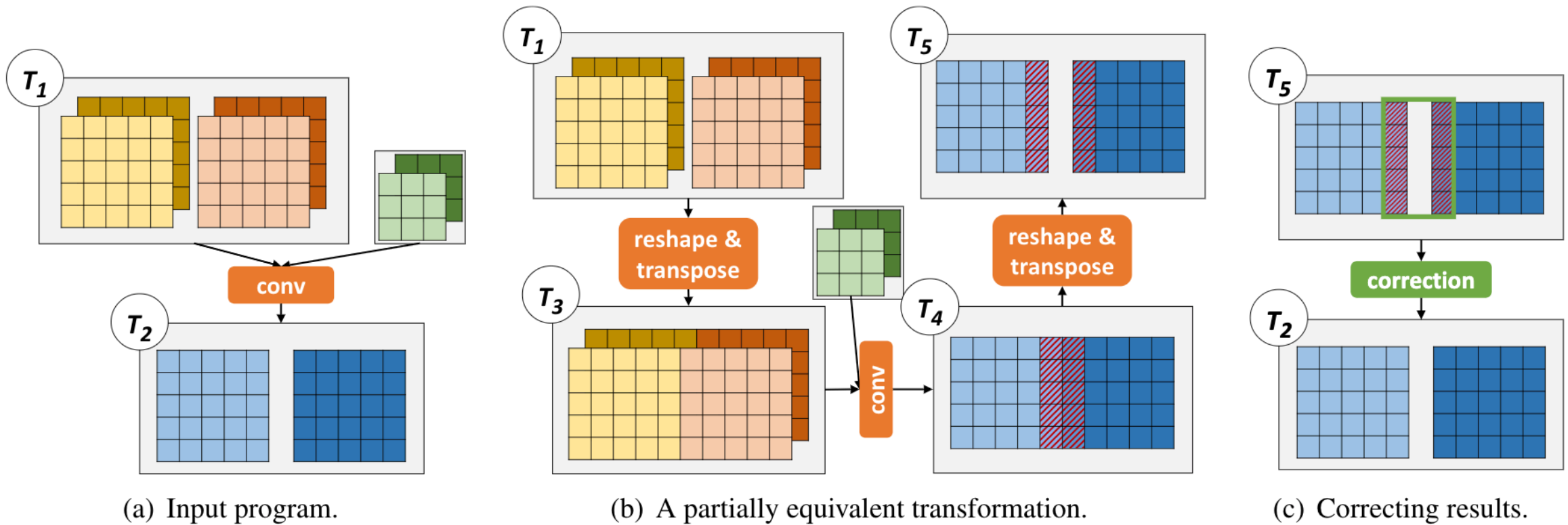
Can we improve TASO?



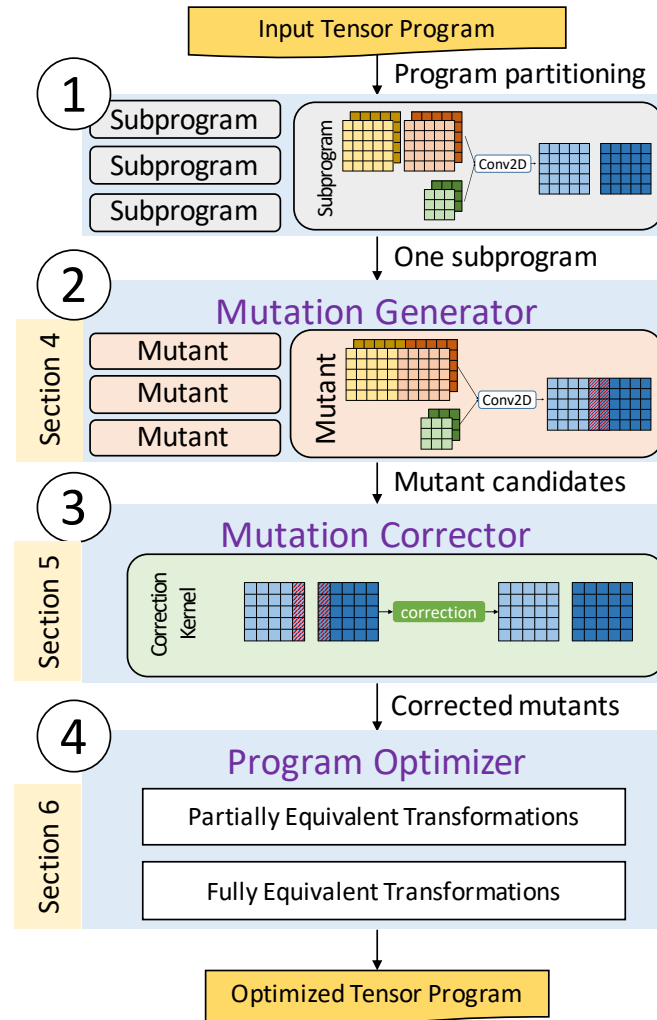
TASO only discovers fully equivalent substs.

Can we use substs that are partially equivalent and correct the results afterwards?

Motivating Example: Partially Equivalent Substs



PET: Partially Equivalent Substs and Auto Corrections



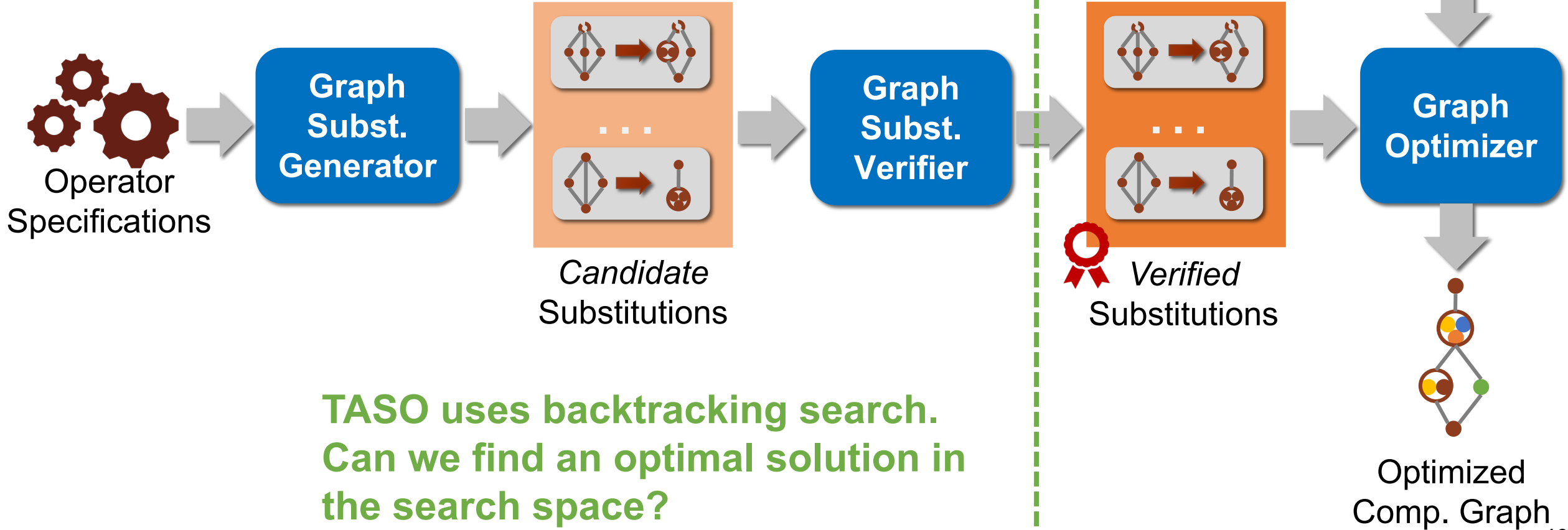
Generate both fully and partially equivalent substs

Auto-correct substs to maintain end-to-end equivalence

Optimize DNNs using both fully and partially equivalent substs

Up to 2.5x faster than TASO

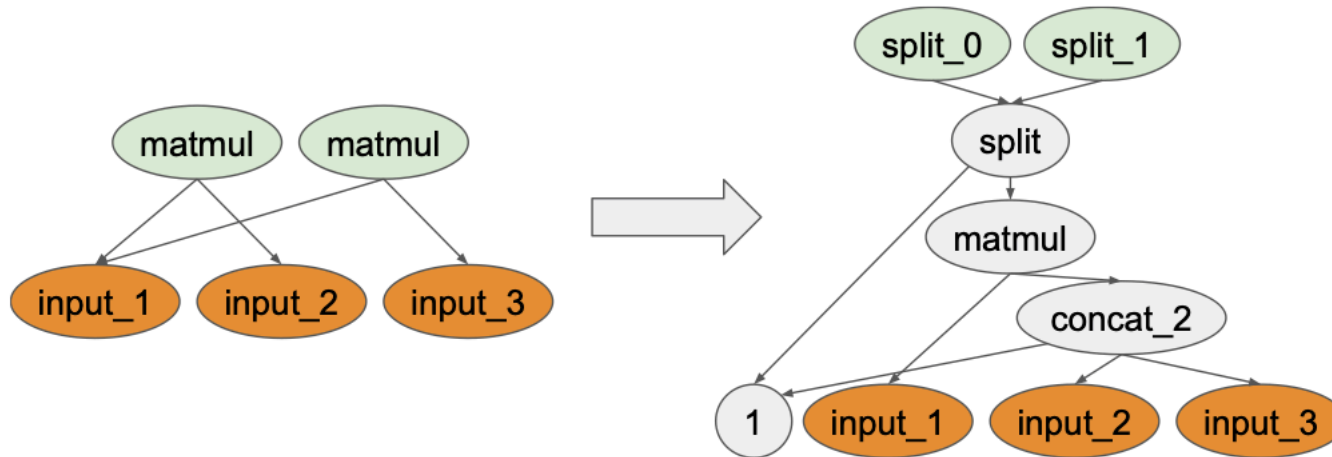
Can we improve TASO?



**TASO uses backtracking search.
Can we find an optimal solution in
the search space?**

Equality Saturation for TASO

- Key idea: a new representation that can express all possible computation graphs at once



Source: (matmul ?input₁ ?input₂), (matmul ?input₁ ?input₃)
Target: (split₀ (split 1 (matmul ?input₁ (concat₂ 1 ?input₂ ?input₃))))),
(split₁ (split 1 (matmul ?input₁ (concat₂ 1 ?input₂ ?input₃))))

The optimal graph is **16%** better.

The search is **48x** faster.

Can we apply TASO to other problem domains?

- **Cross-optimizations between ML and DB operations**

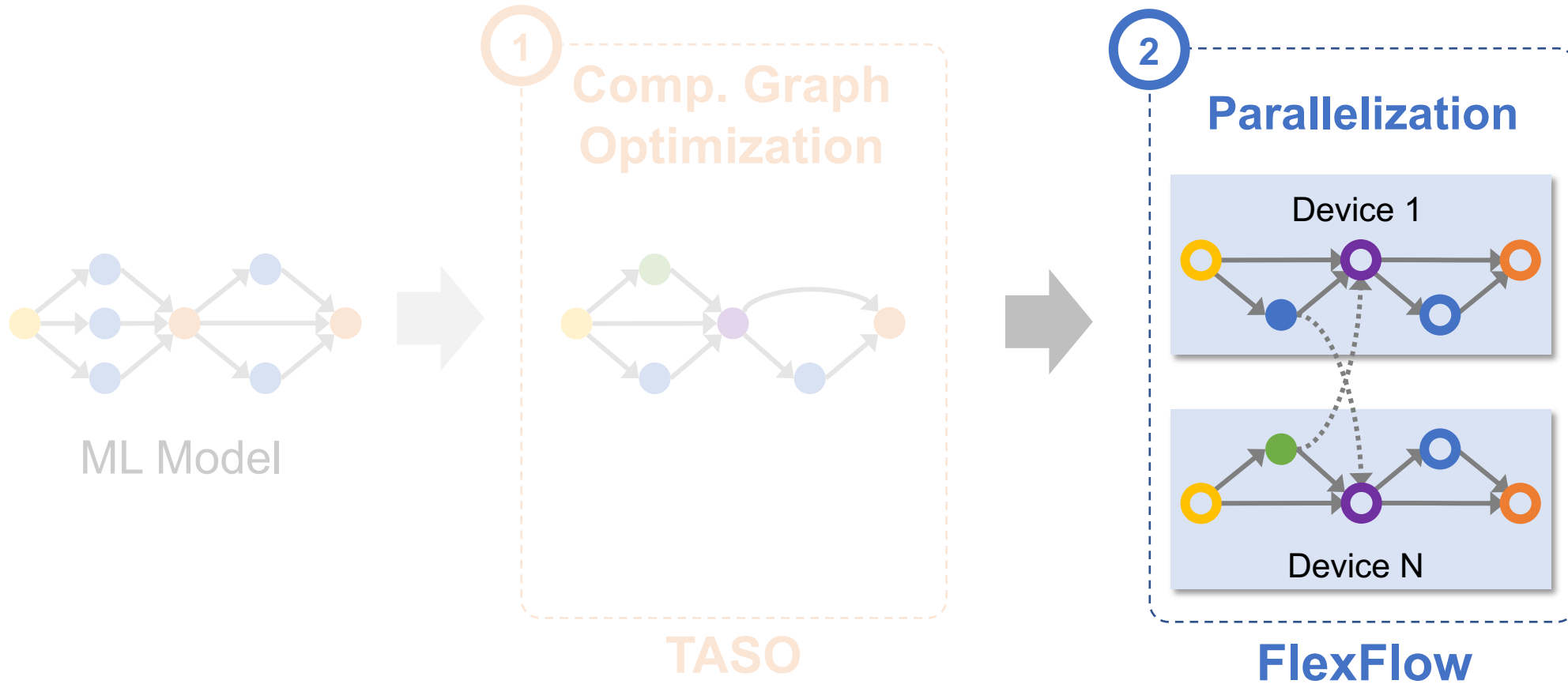
- Why: ML and DB operations are optimized separately in today's DB systems
- How: automatically generate co-optimizations of linear algebra and relational algebra operations

- **Optimizing Compilers for Quantum Computing**

- Why: today's quantum machines support different sets of instructions -> impossible to manually design optimizations for all quantum architectures
- How: automatically generate quantum program transformations given a set of instructions

- **Others?**

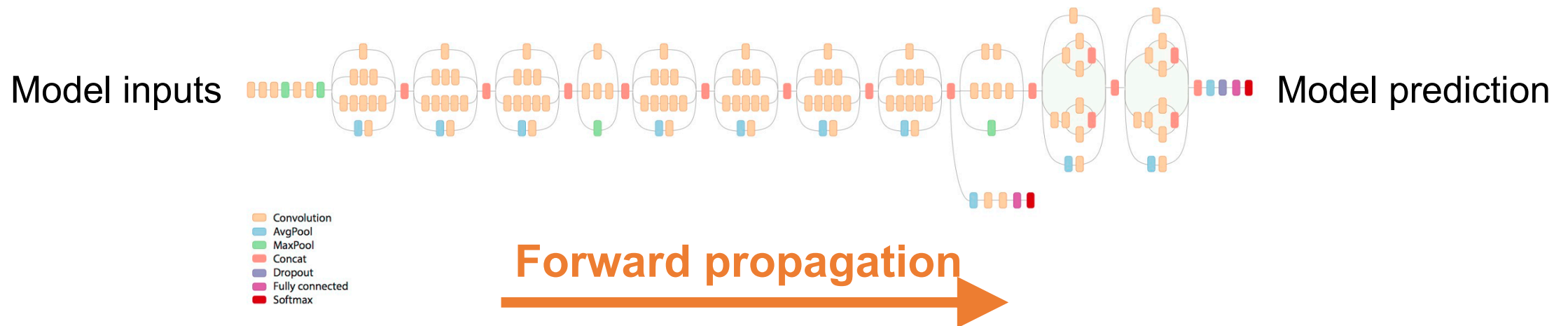
Automated Discovery of ML Optimizations



Stochastic Gradient Descent (SGD)

Train ML models through many iterations of 3 stages

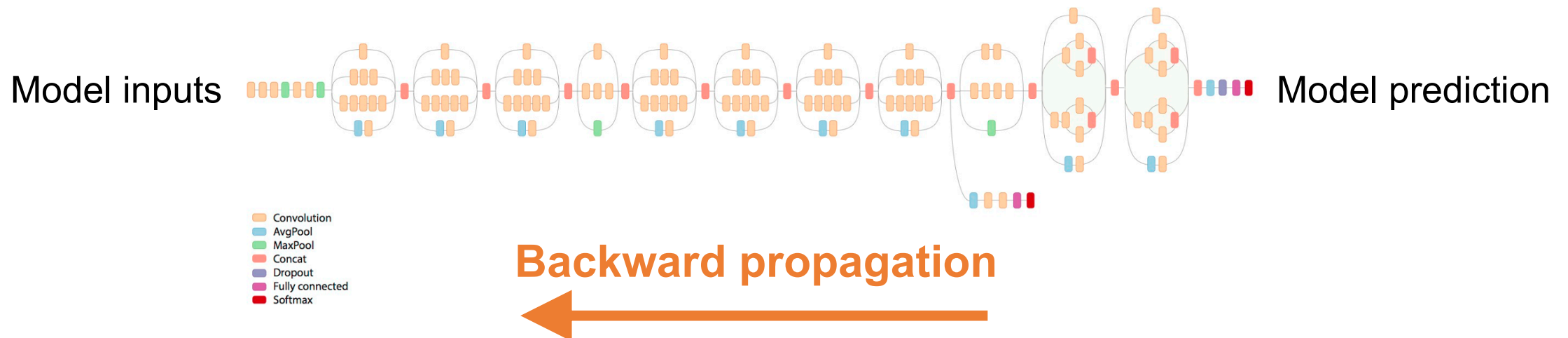
1. **Forward propagation**: apply model to a batch of input samples and run calculation through operators to produce a prediction
2. **Backward propagation**: run the model in reverse to produce error for each trainable weight
3. **Weight update**: use the loss value to update model weights



Stochastic Gradient Descent (SGD)

Train ML models through many iterations of 3 stages

1. **Forward propagation**: apply model to a batch of input samples and run calculation through operators to produce a prediction
2. **Backward propagation**: run the model in reverse to produce error for each trainable weight
3. **Weight update**: use the loss value to update model weights



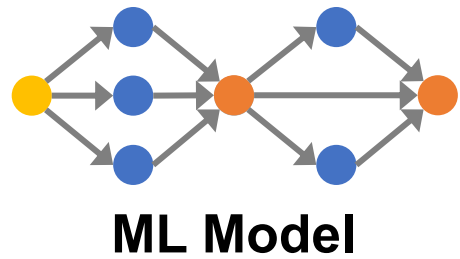
Stochastic Gradient Descent (SGD)

Train ML models through many iterations of 3 stages

1. **Forward propagation**: apply model to a batch of input samples and run calculation through operators to produce a prediction
2. **Backward propagation**: run the model in reverse to produce error for each trainable weight
3. **Weight update**: use the loss value to update model weights

$$w_i := w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^n \nabla L_j(w_i)$$

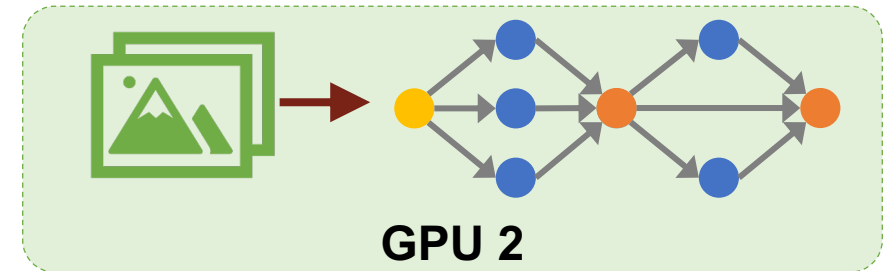
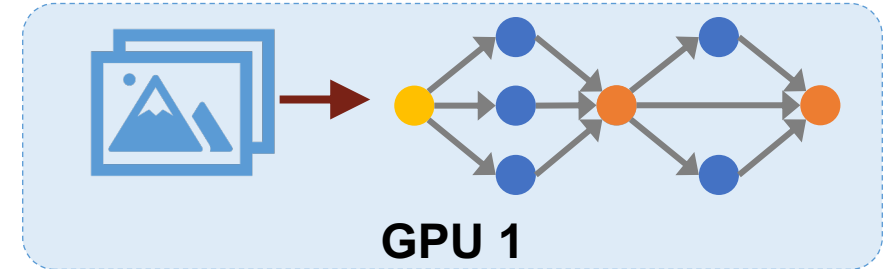
Current Strategies to Parallelize ML Training: Data and Model Parallelism



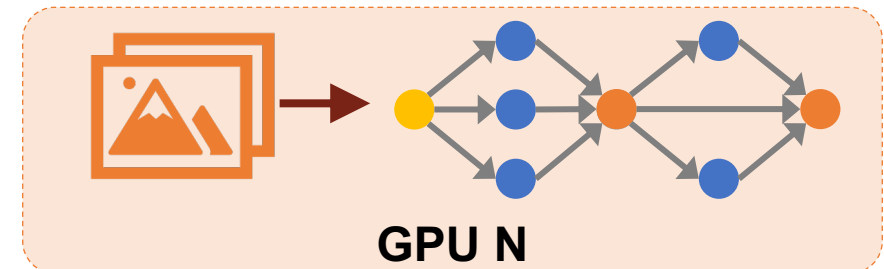
Training Dataset

Data
Parallelism

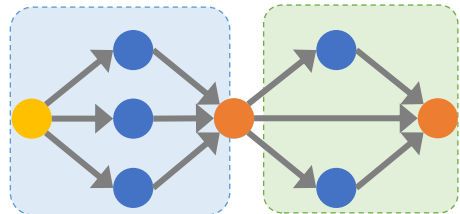
$$w_i := w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^n \nabla L_j(w_i)$$



...



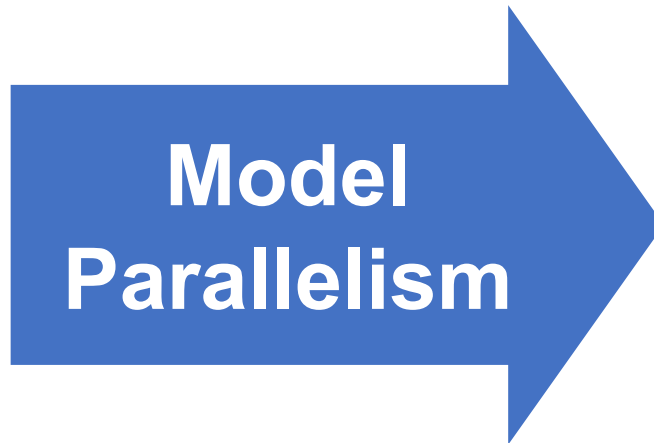
Current Strategies to Parallelize ML Training: Data and Model Parallelism



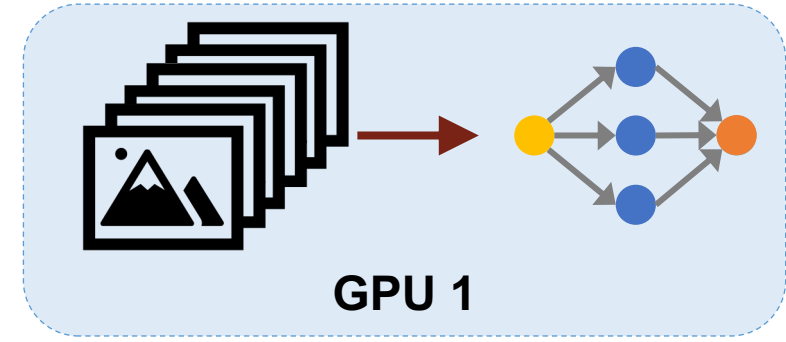
ML Model



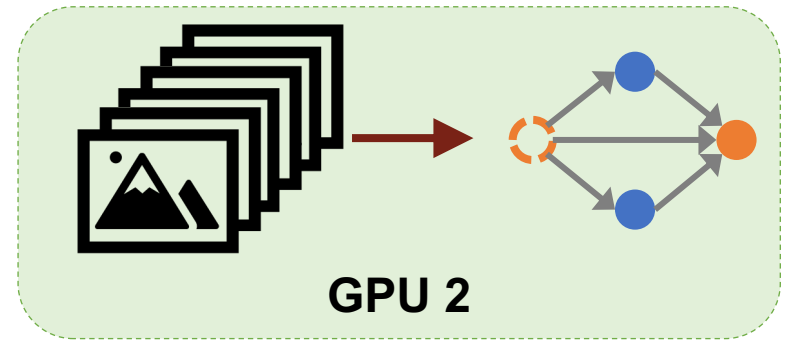
Training Dataset



Model
Parallelism



GPU 1



GPU 2

$$w_i := w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^n \nabla L_j(w_i)$$

Are there strategies beyond data/model parallelism?
Can we discover fast ones automatically?

FlexFlow: Automated Search for Fast Strategies

Define a **search space** of possible parallelization strategies

+

A **cost model** and a **search algorithm**

=

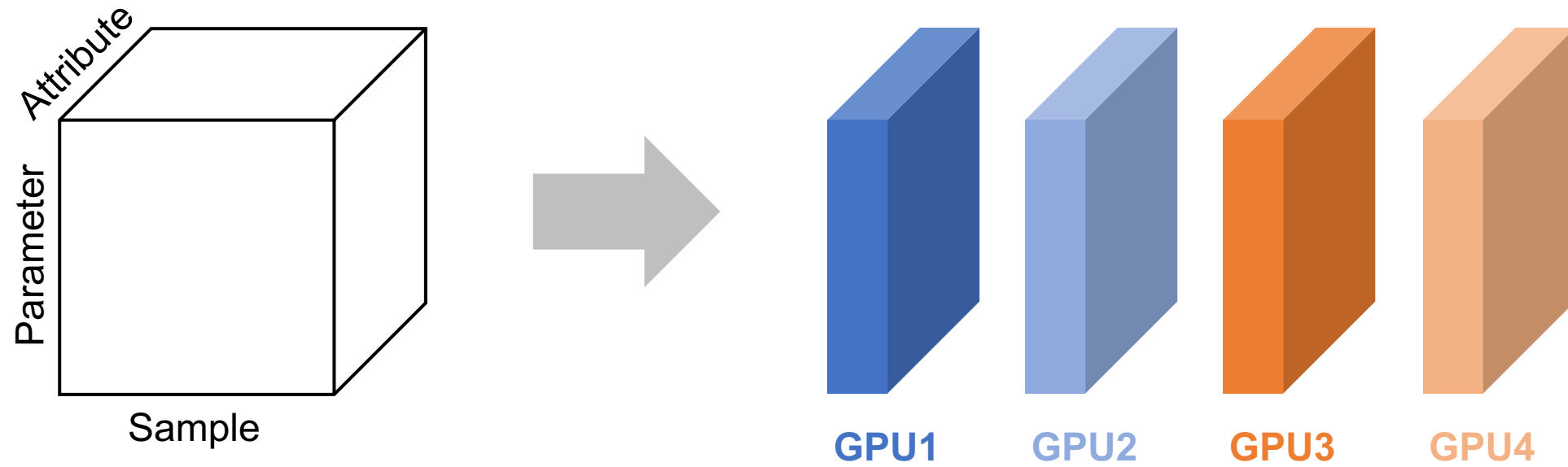
**Optimized
Parallelization strategies**

The SOAP Search Space

- **S**amples
- **O**perators
- **A**tributes
- **P**arameters

The SOAP Search Space

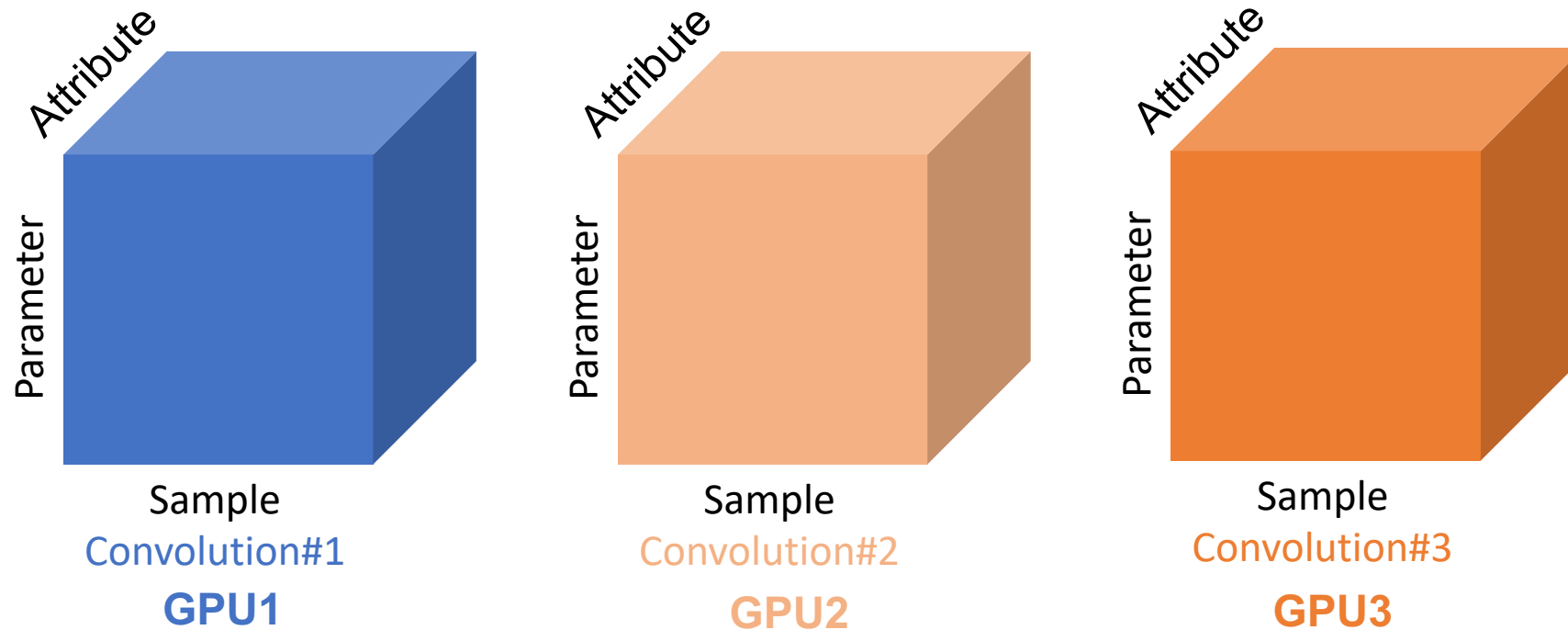
- **S**amples: partitioning training samples (Data Parallelism)
- **O**perators
- **A**tributes
- **P**arameters



Parallelizing a 1D convolution in *Sample*

The SOAP Search Space

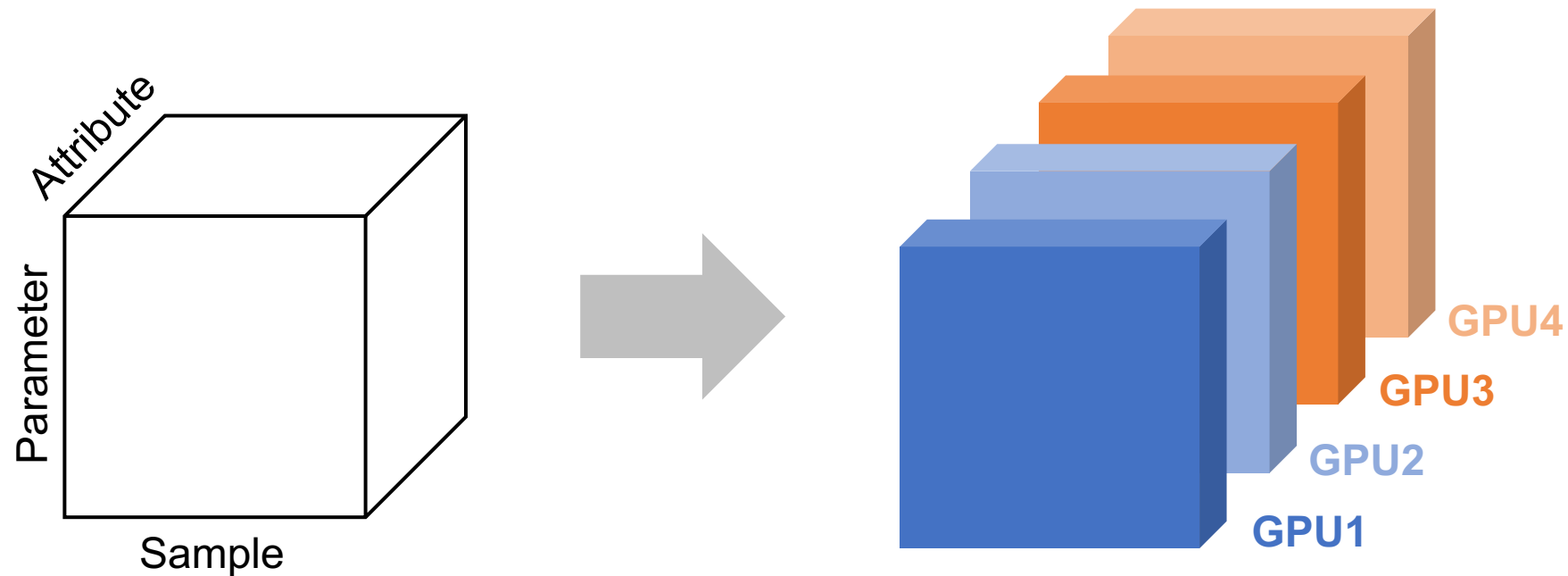
- **S**amples: partitioning training samples (Data Parallelism)
- **O**perators: partitioning ML operators (Model Parallelism)
- **A**tributes
- **P**arameters



Parallelizing multiple convolutions in *Operator*

The SOAP Search Space

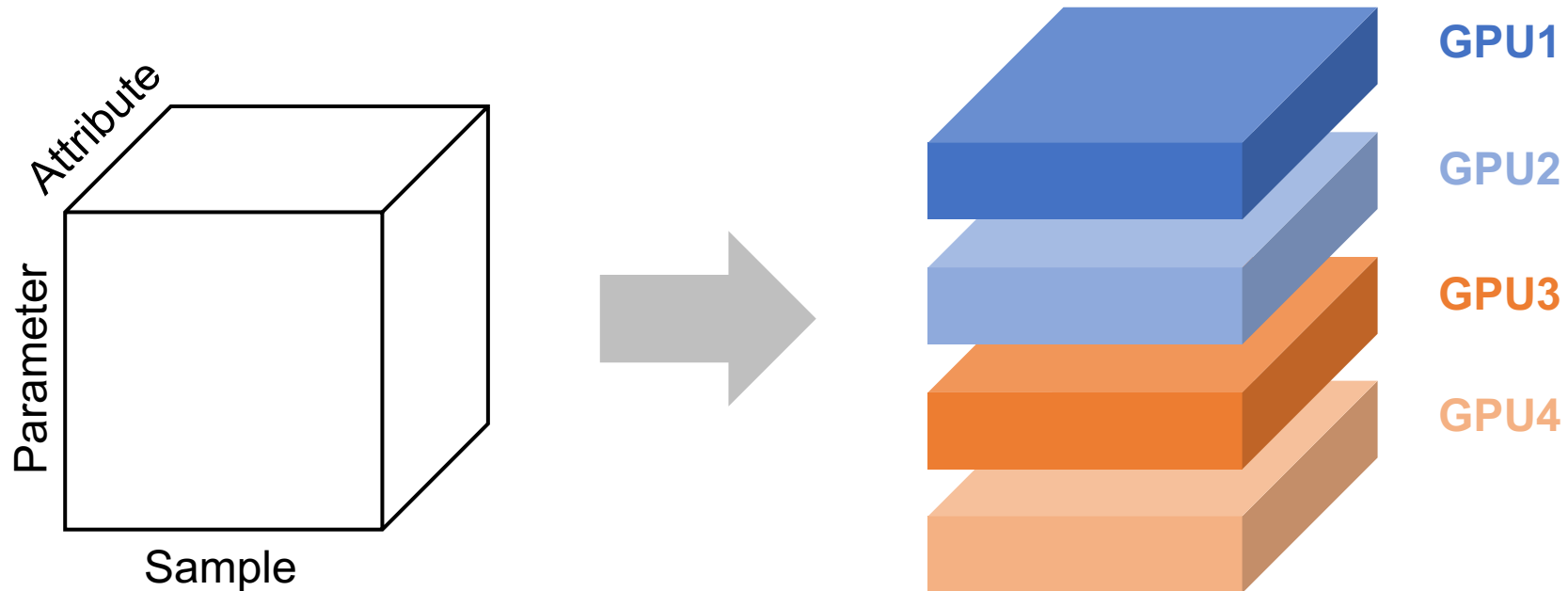
- **S**amples: partitioning training samples (Data Parallelism)
- **O**perators: partitioning ML operators (Model Parallelism)
- **A**tributes: partitioning attributes in a sample (e.g., pixels)
- **P**arameters



Parallelizing a 1D convolution in *Attribute*

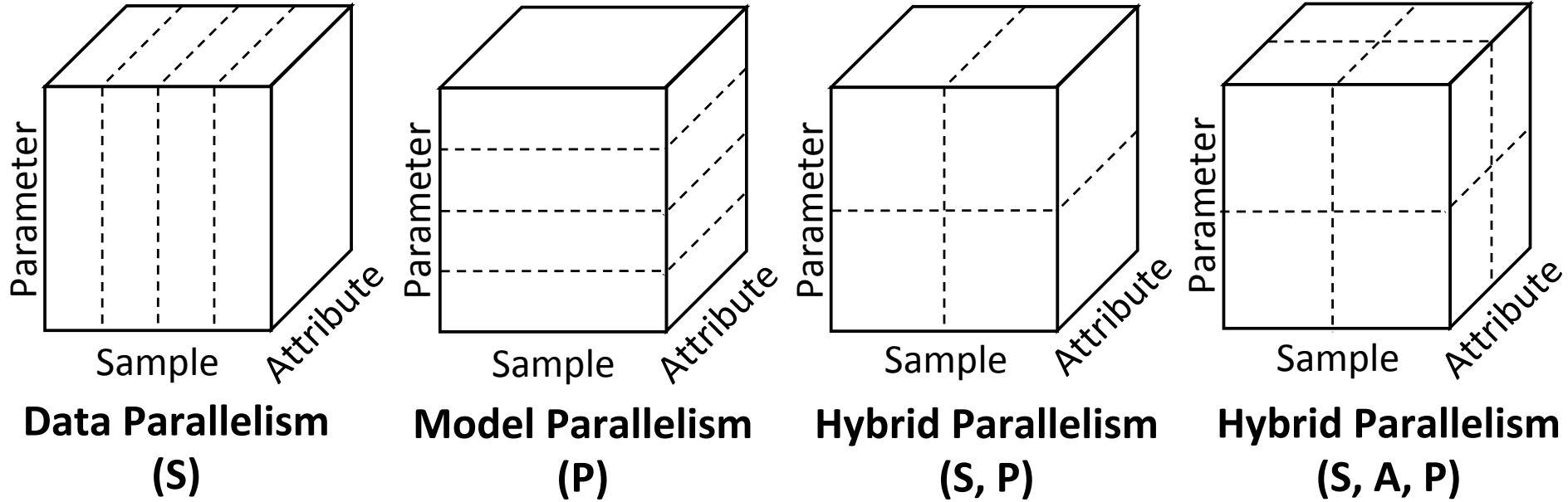
The SOAP Search Space

- **S**amples: partitioning training samples (Data Parallelism)
- **O**perators: partitioning ML operators (Model Parallelism)
- **A**tributes: partitioning attributes in a sample (e.g., pixels)
- **P**arameters: partitioning parameters in an operator



Parallelizing a 1D convolution in *Parameter*

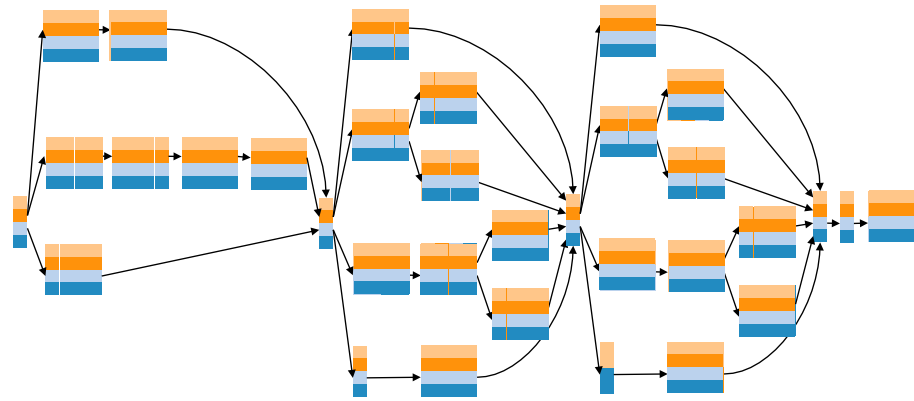
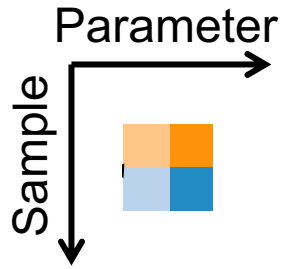
Hybrid Parallelism in SOAP



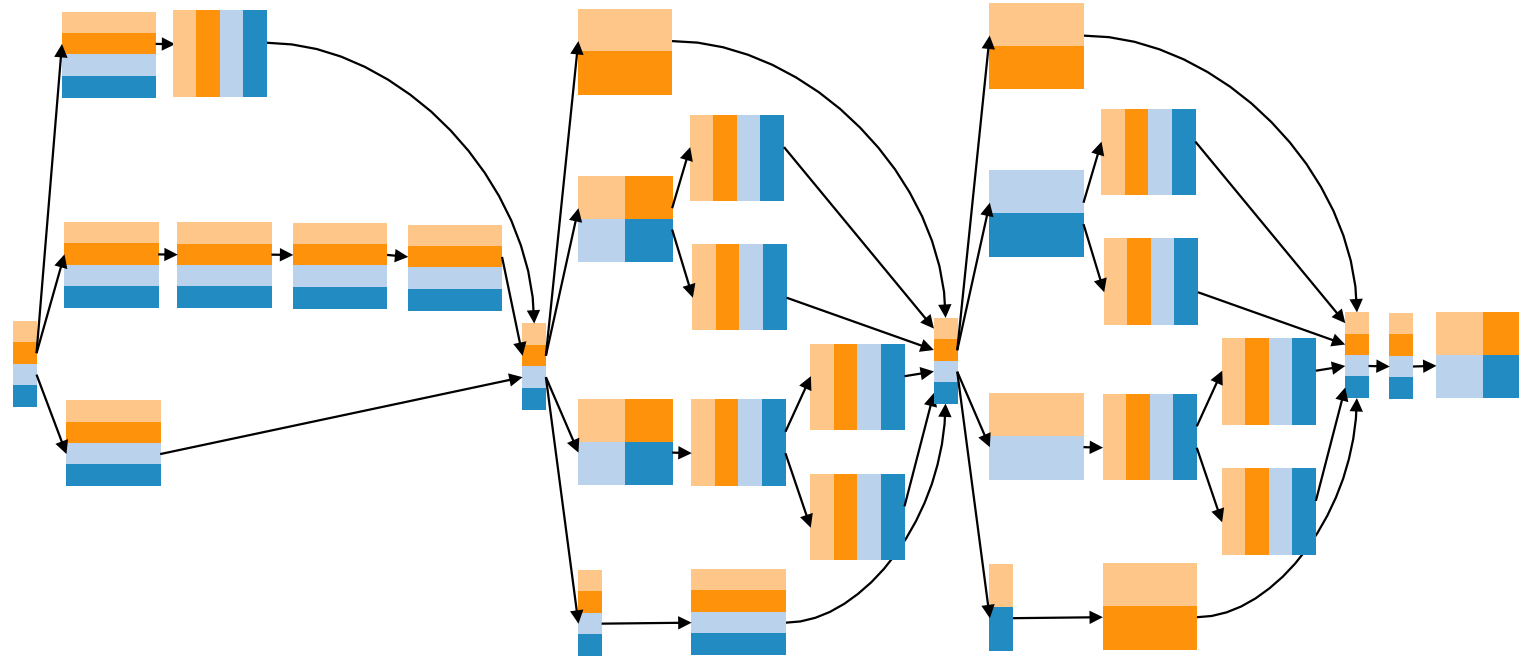
Example parallelization strategies for 1D convolution

Different strategies perform the same computation.

- GPU 1
- GPU 2
- GPU 3
- GPU 4

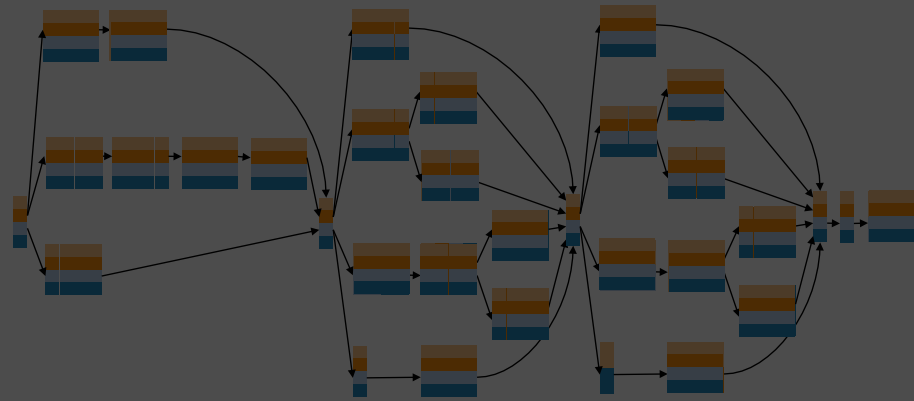
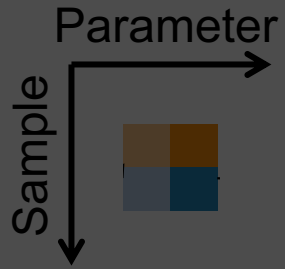


Data parallelism

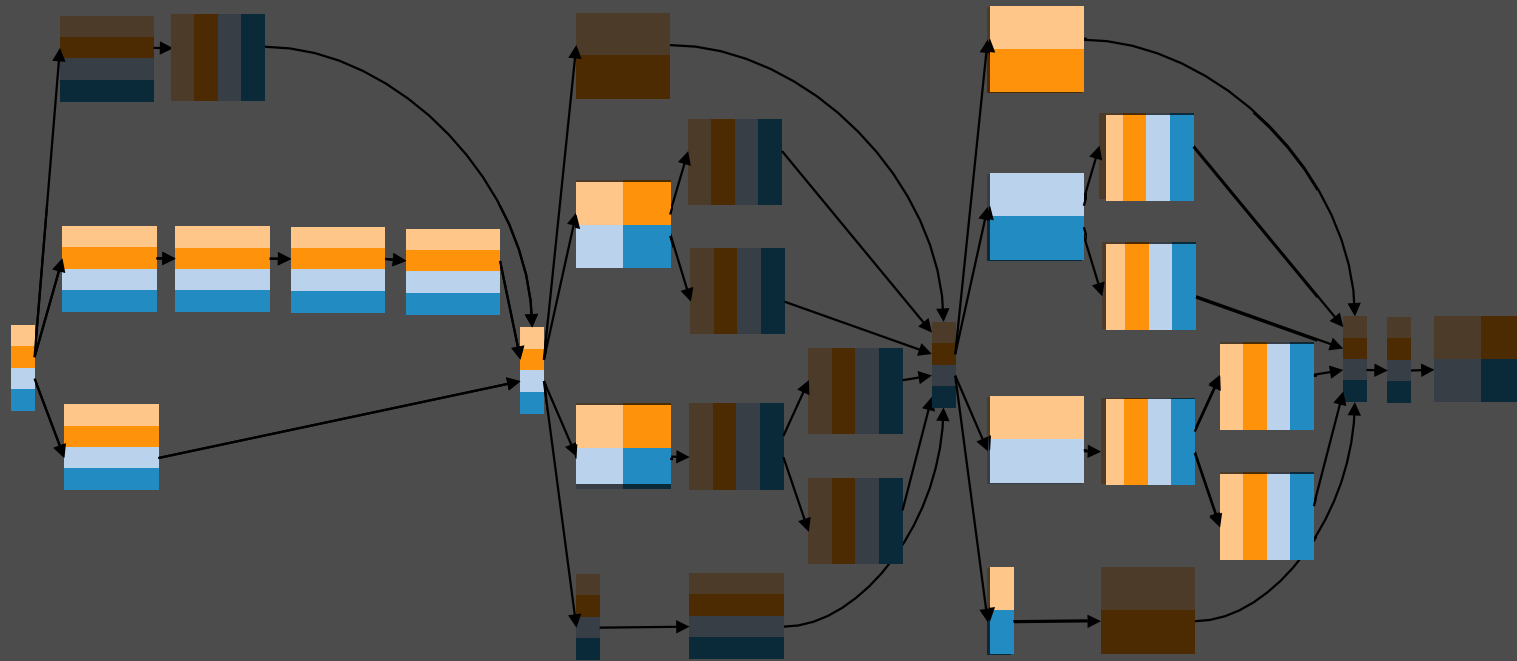


A parallelization strategy in SOAP **(1.2x faster)**

- GPU 1
- GPU 2
- GPU 3
- GPU 4



Data parallelism



A parallelization strategy in SOAP **(1.2x faster)**

Challenges of Discovering Fast Strategies in SOAP

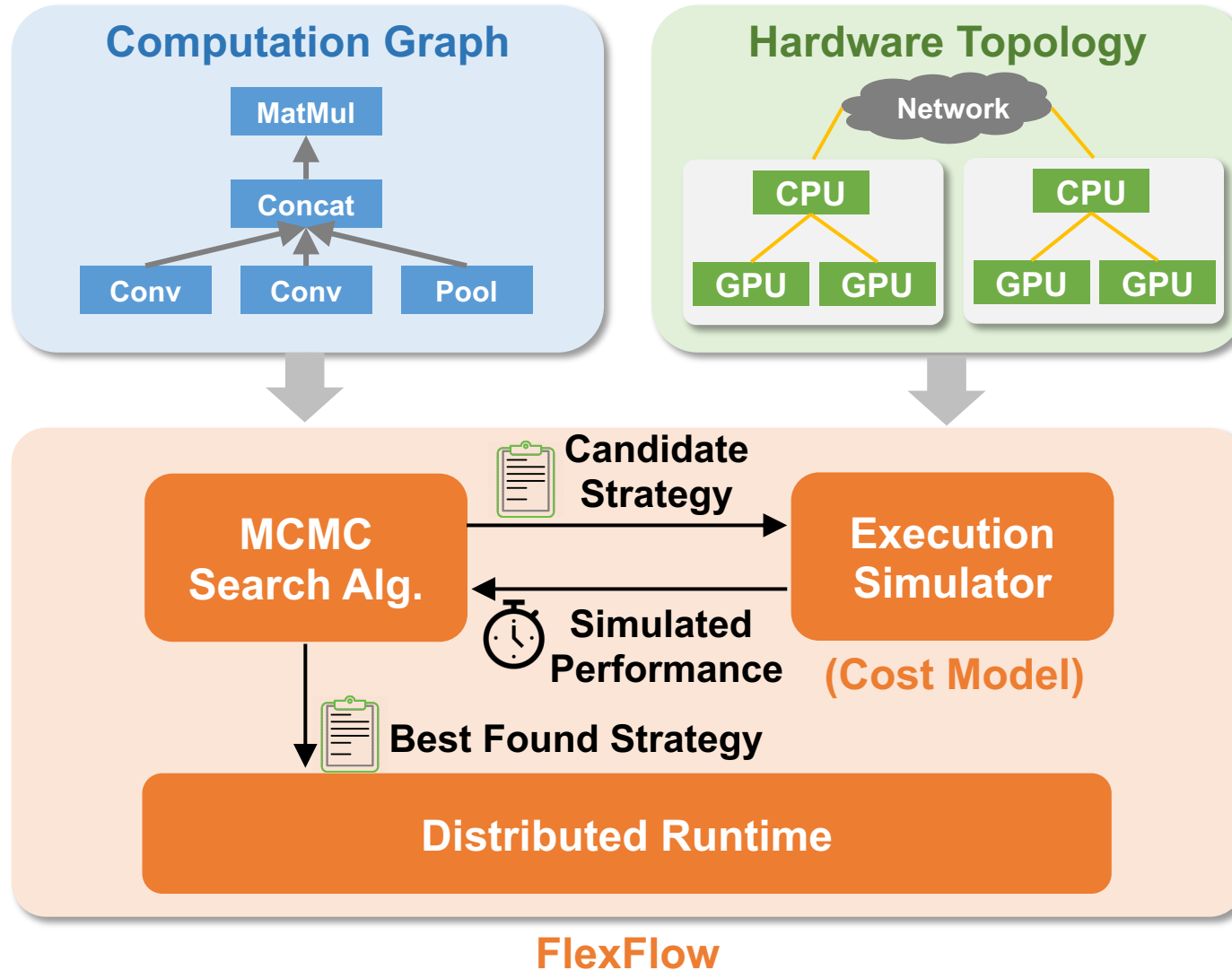
1. SOAP contains billions or more possible strategies

MCMC search algorithm

2. Evaluating a strategy on hardware is too slow

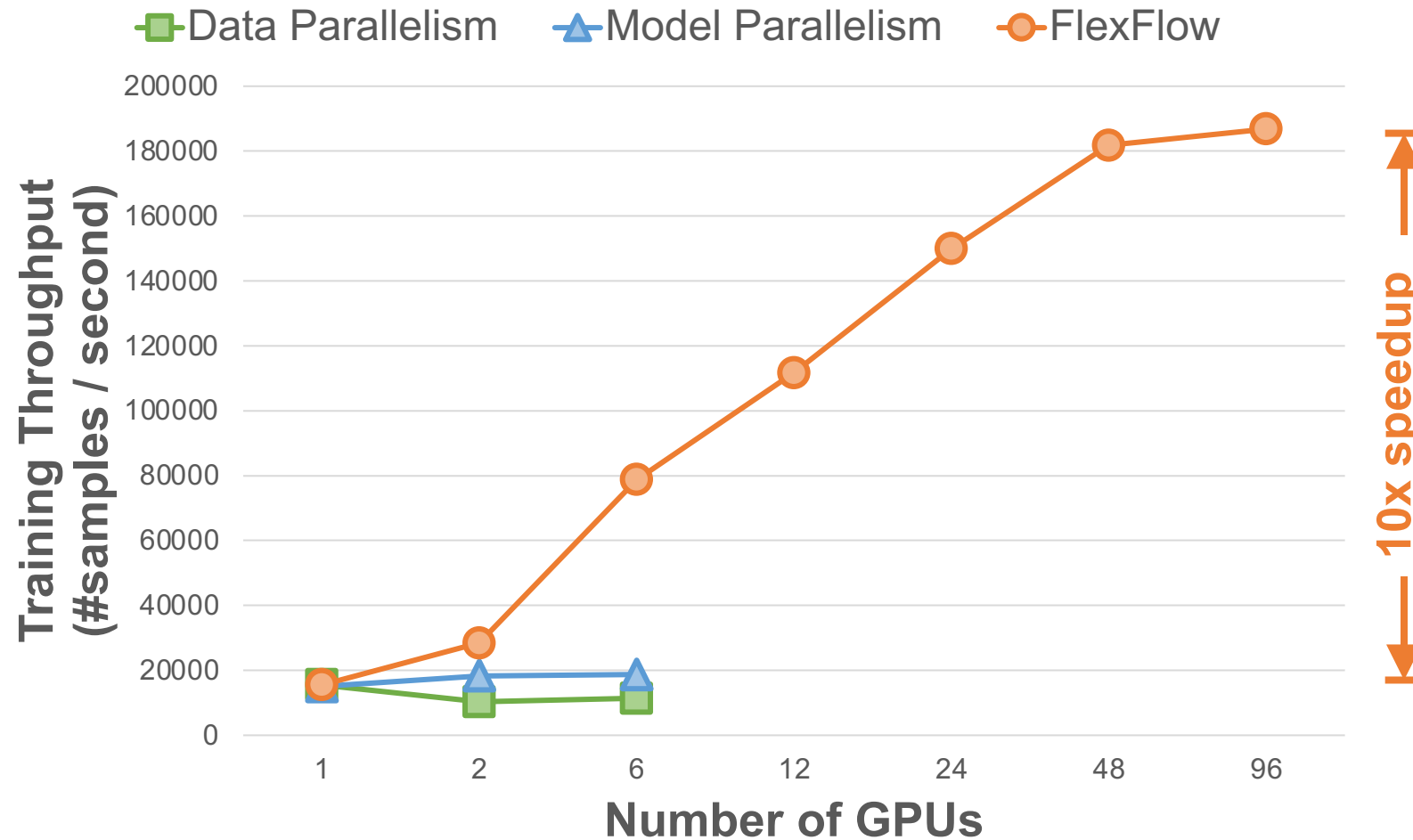
Execution simulator

FlexFlow Overview

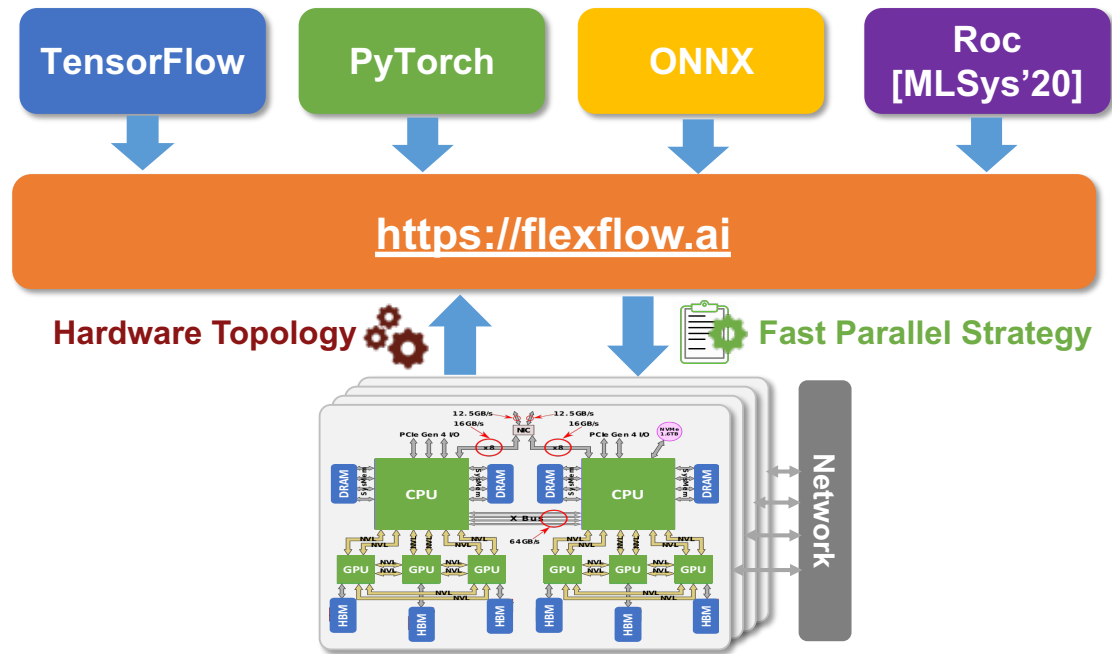


Deep Learning Recommendation Model (DLRM)

A deep learning model for ads recommendation



FlexFlow: Automatically Discover Fast Parallelization for DNNs



<https://flexflow.ai>

FlexFlow

Automatically Discovering Fast Parallelization Strategies for Distributed Deep Neural Network Training

[Latest release r20.08](#)

[Install now](#)

Performance Autotuning

FlexFlow accelerates DNN training by automatically discovering fast parallelization strategies for a specific parallel machine.

[Learn more](#)

Keras Support

FlexFlow provides a drop-in replacement for TensorFlow Keras and requires only a few lines of changes to existing Keras programs.

[Learn more](#)

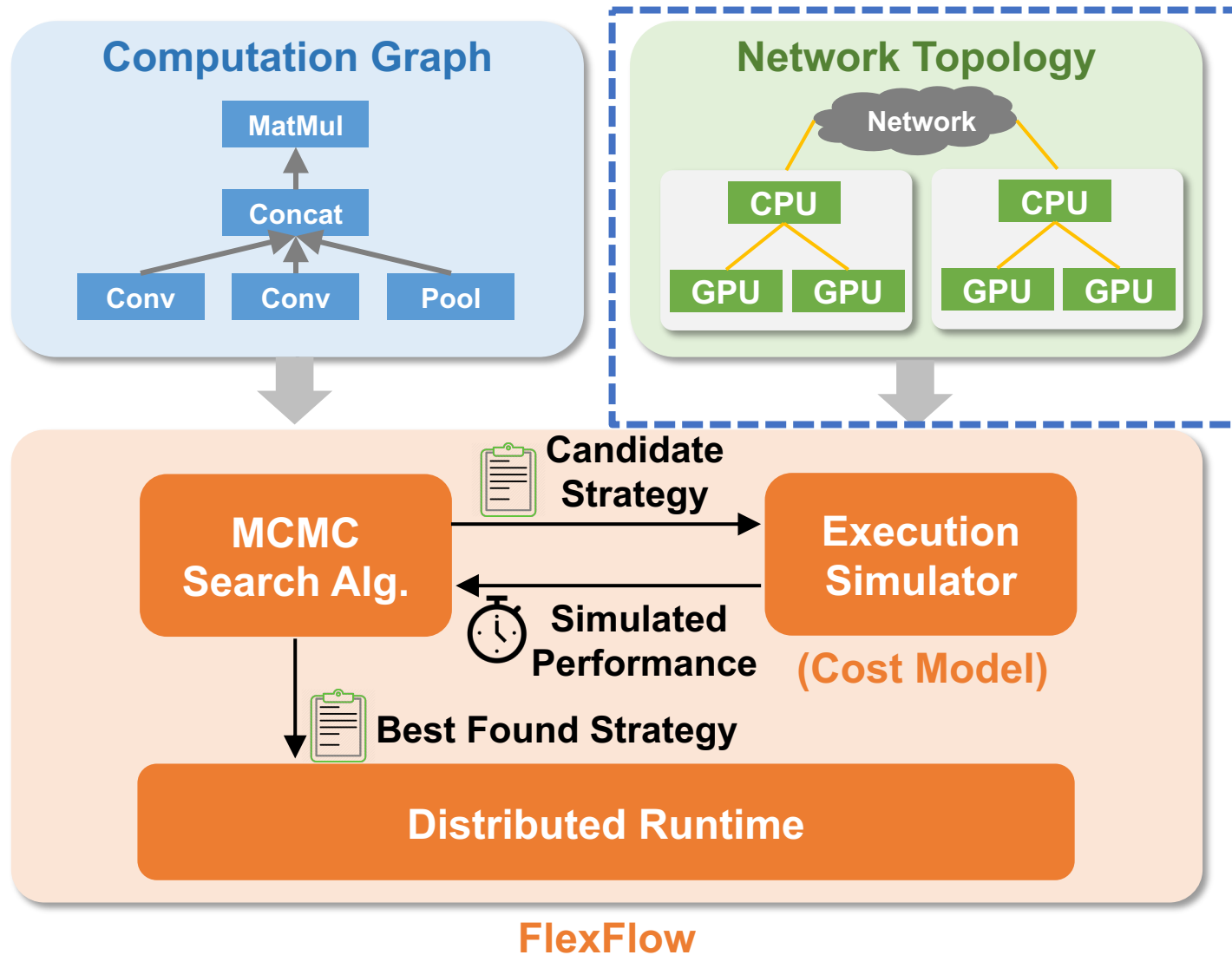
Large-Scale GNNs

FlexFlow enables fast graph neural network training and inference on large-scale graphs by exploring attribute parallelism.

[Learn more](#)



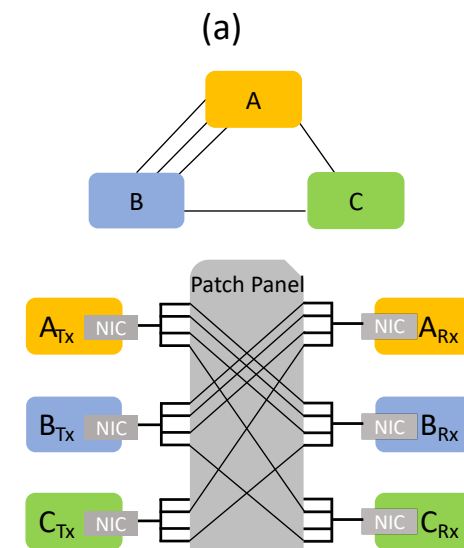
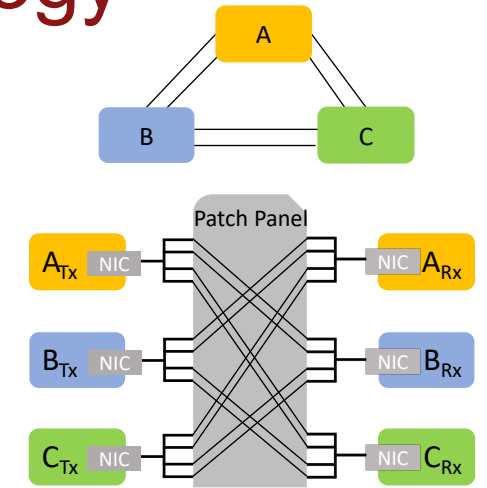
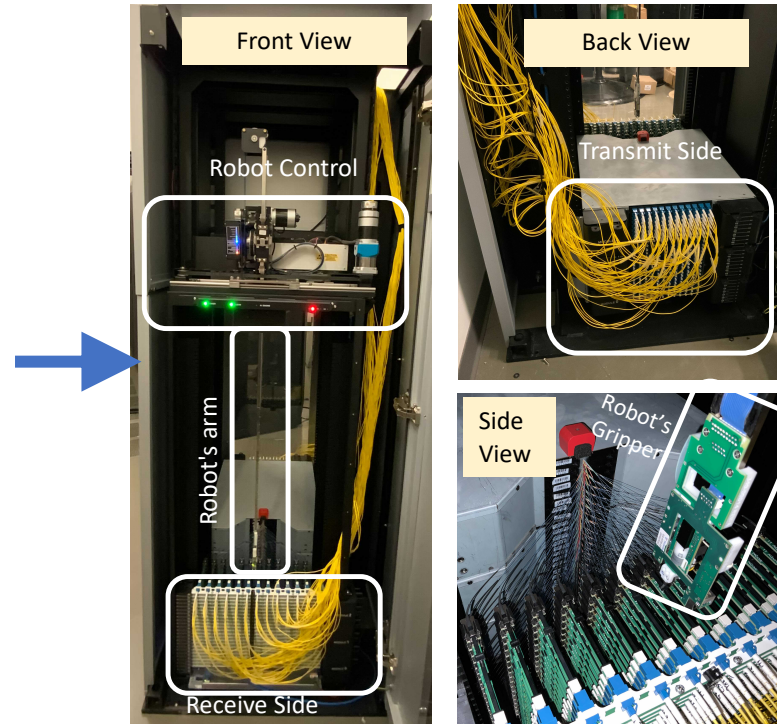
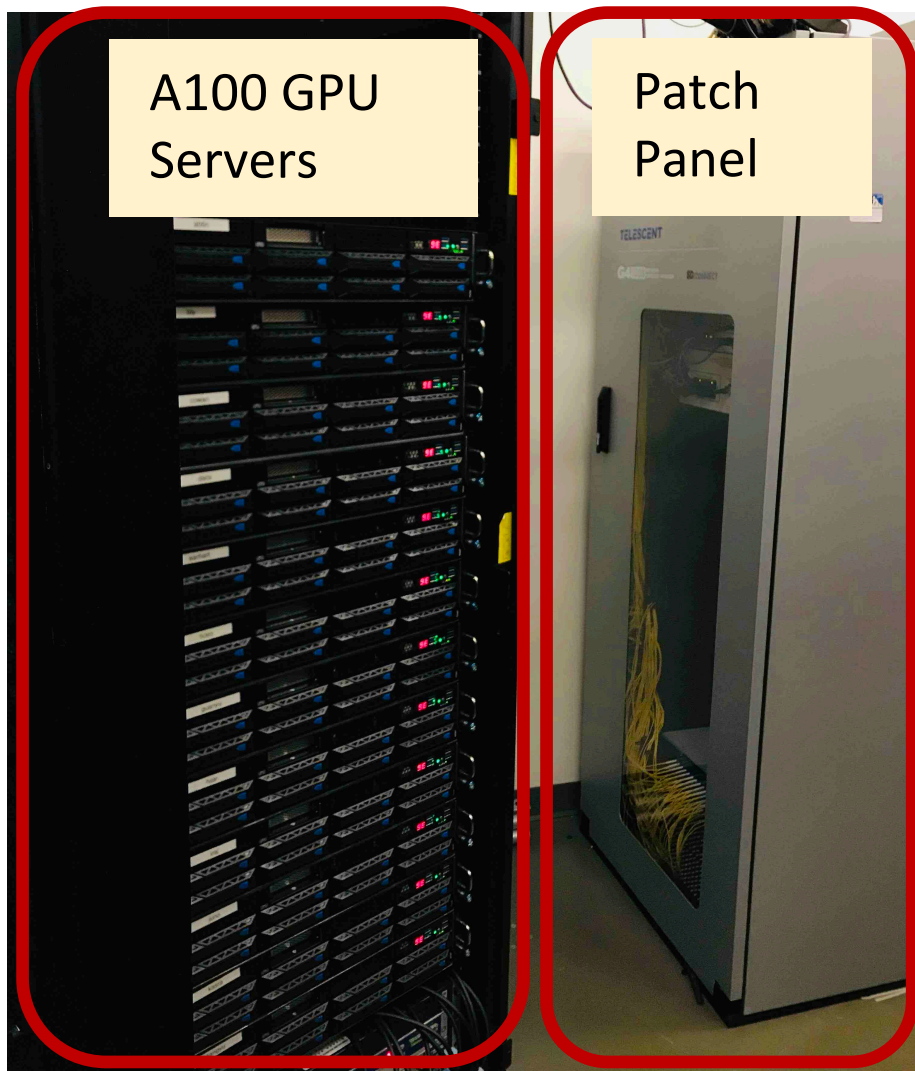
Can we improve FlexFlow?



FlexFlow takes network topology as an input.

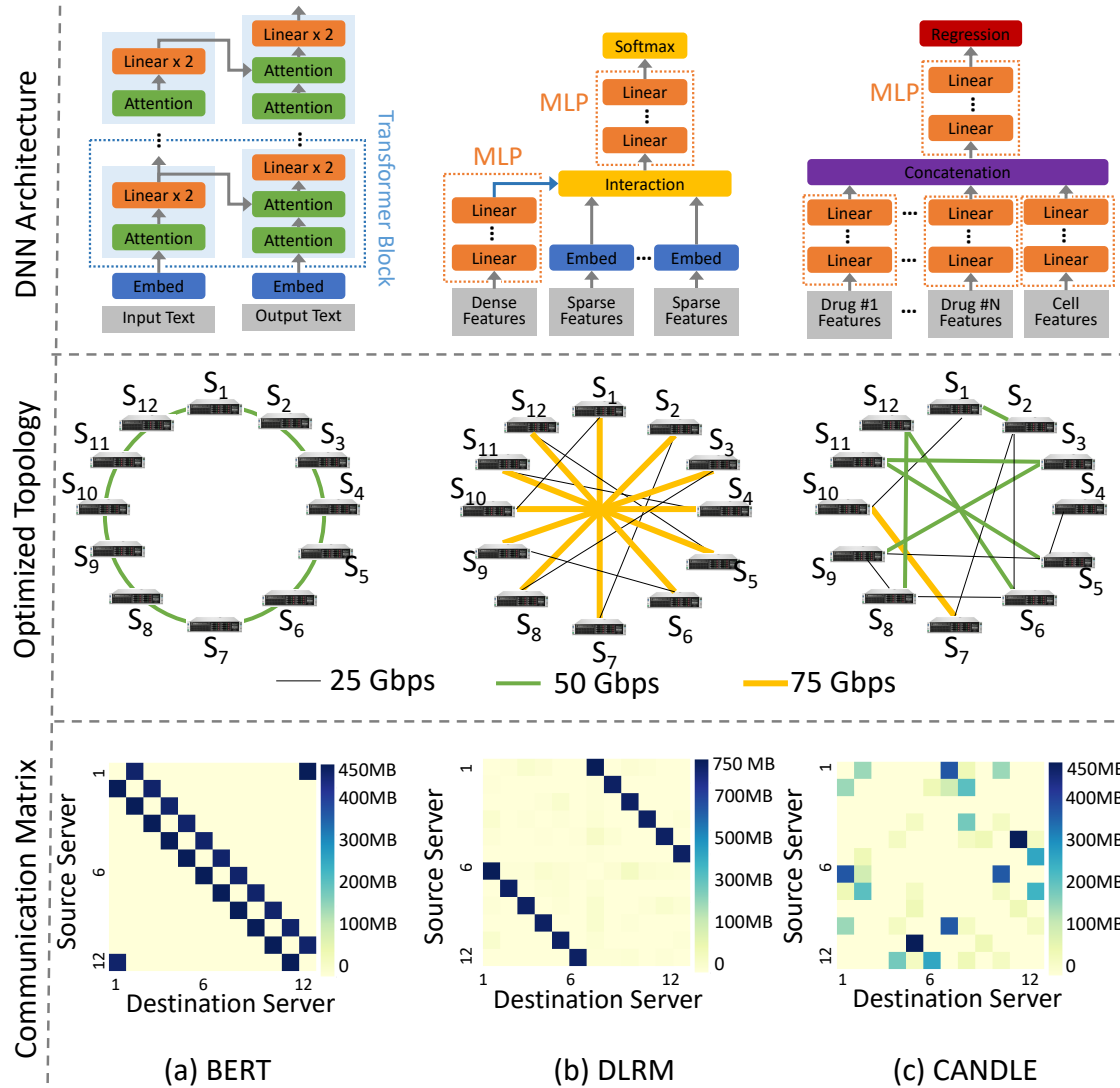
Can we co-optimize network topology and parallelization strategies?

Testbed with Reconfigurable Network Topology



Reconfigurable network topology

Joint Optimization of Parallelization Strategy and Network Topology



Up to 5.7x faster than FlexFlow w/ Fat-tree interconnect

Automated Machine Learning Systems

