

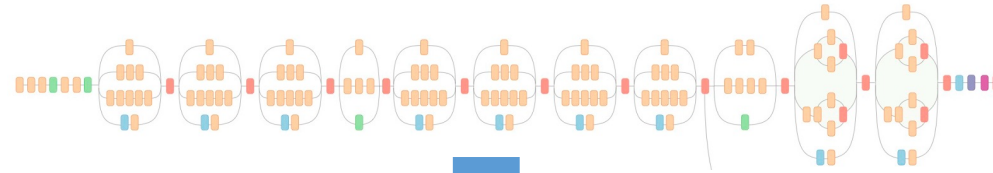
# Automatically Discovering ML Optimizations

**Zhihao Jia**

Computer Science Department  
Carnegie Mellon University

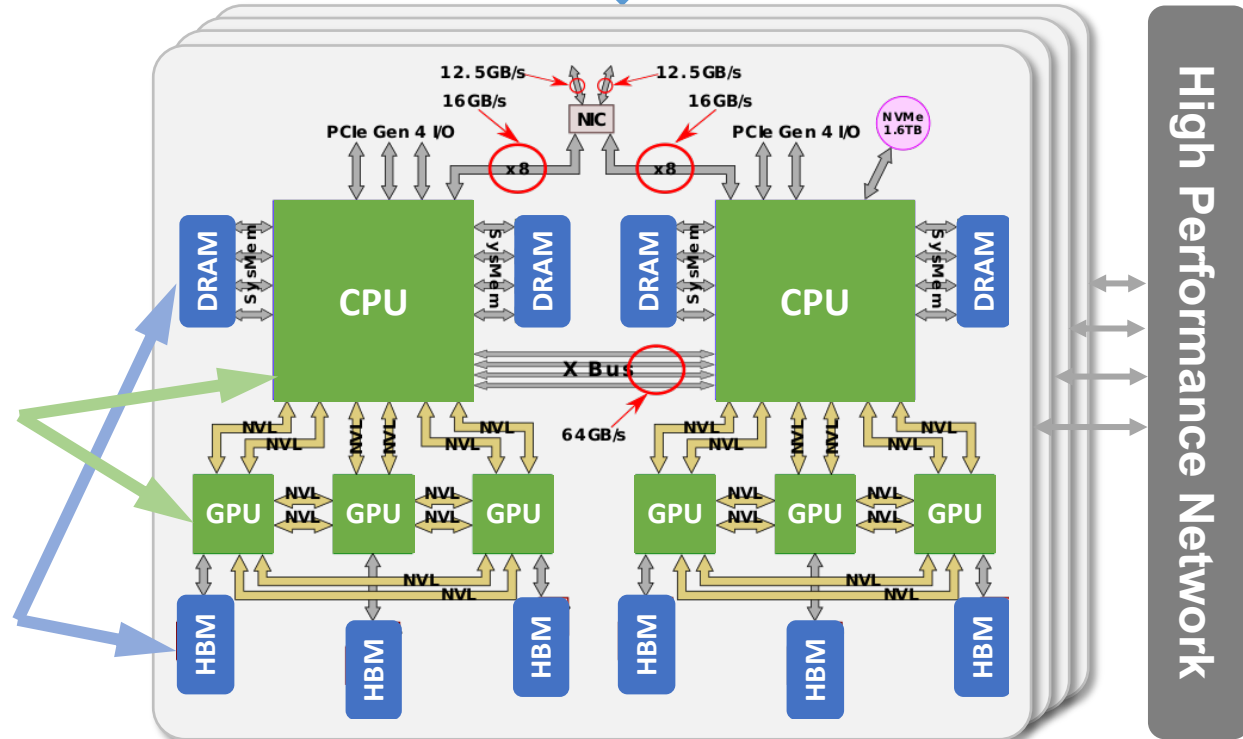


# My Research: Systems for ML



ML Model

ML Systems



Heterogenous Processors

Heterogenous Memories

High Performance Network

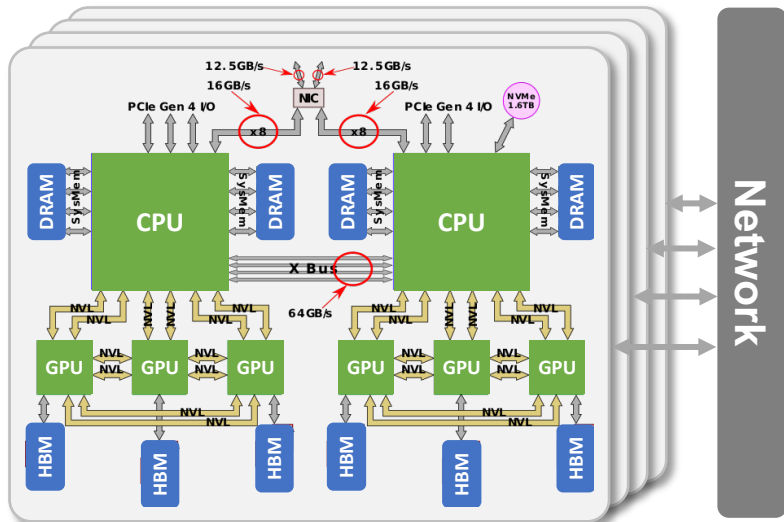
Distributed Heterogenous Hardware Architectures

~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

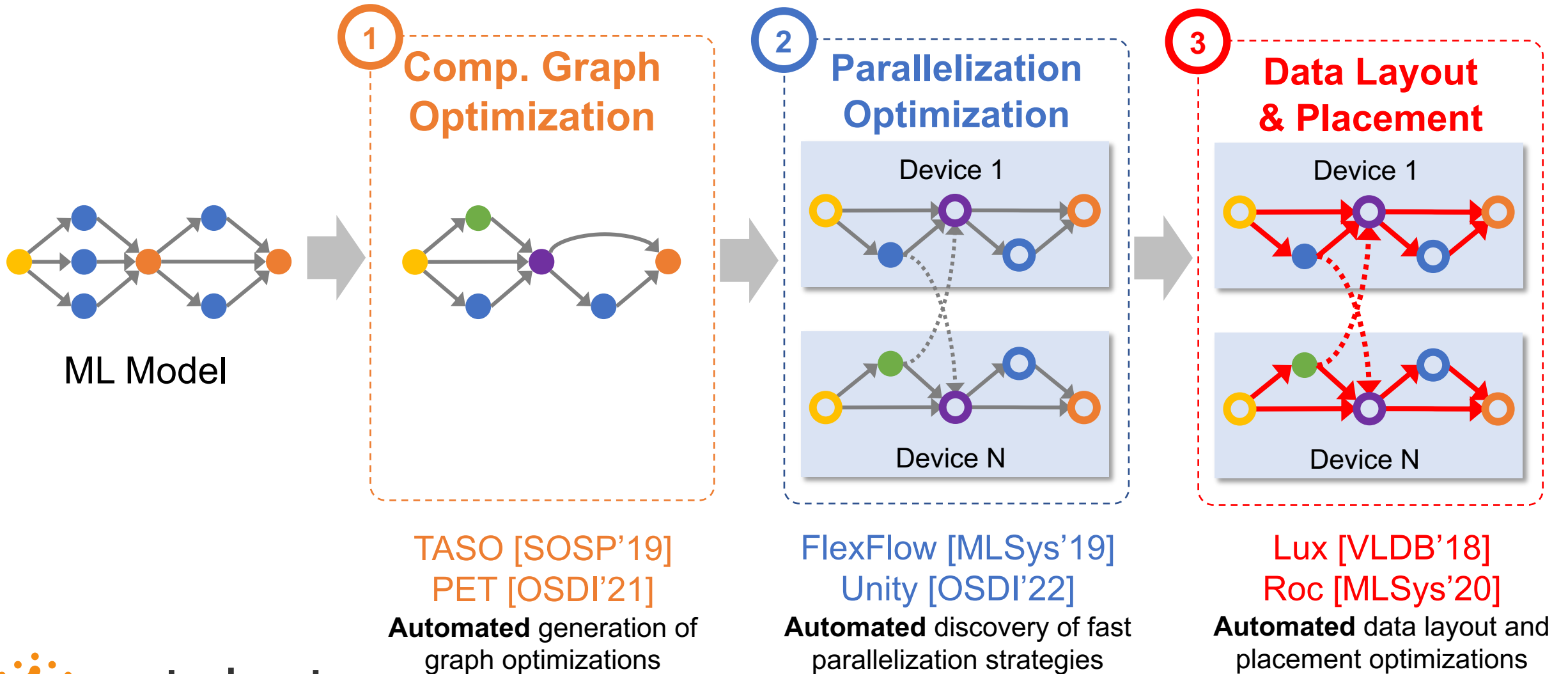
**Heterogenous Hardware**  
Different processor kinds and complex memory hierarchy

# CMU Automated Learning Systems Lab

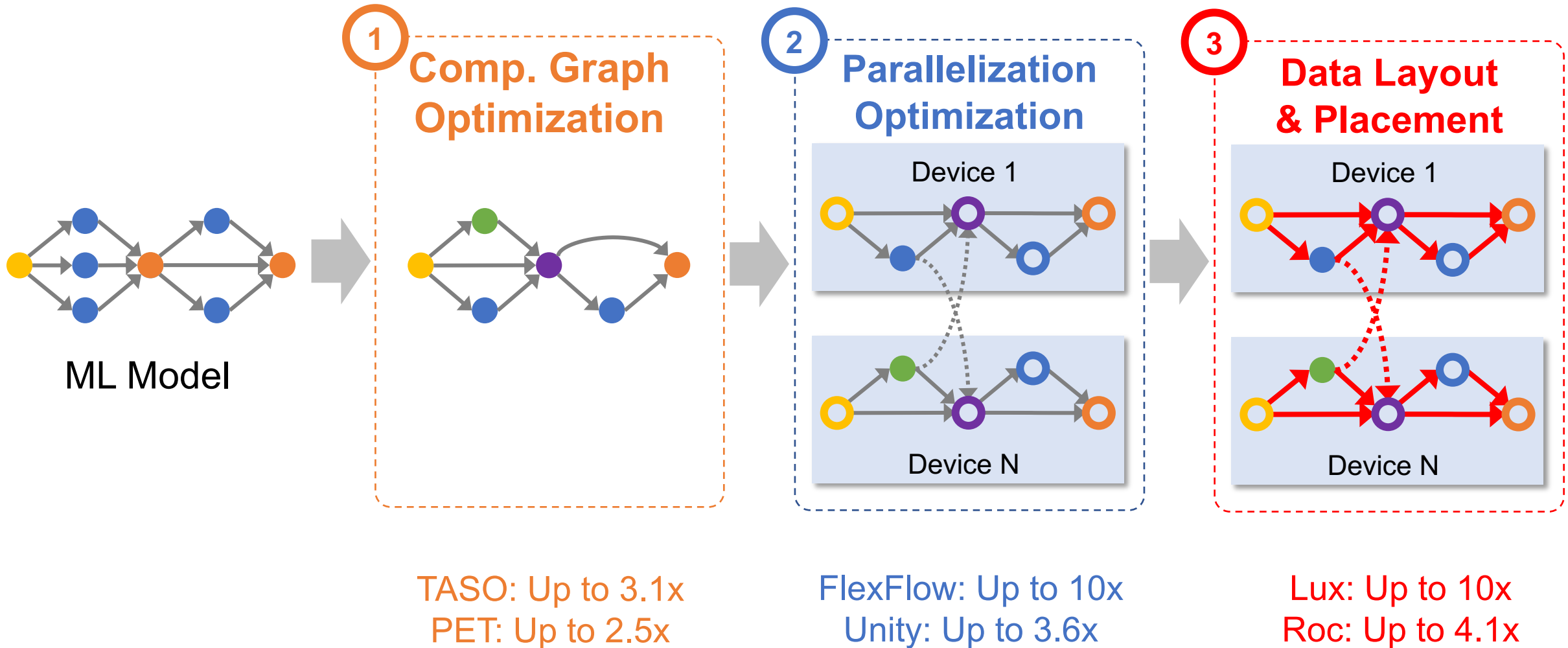
Mission: Automate the design and optimization of ML systems by leveraging

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

# Our Research: Automated Discovery of ML Optimizations



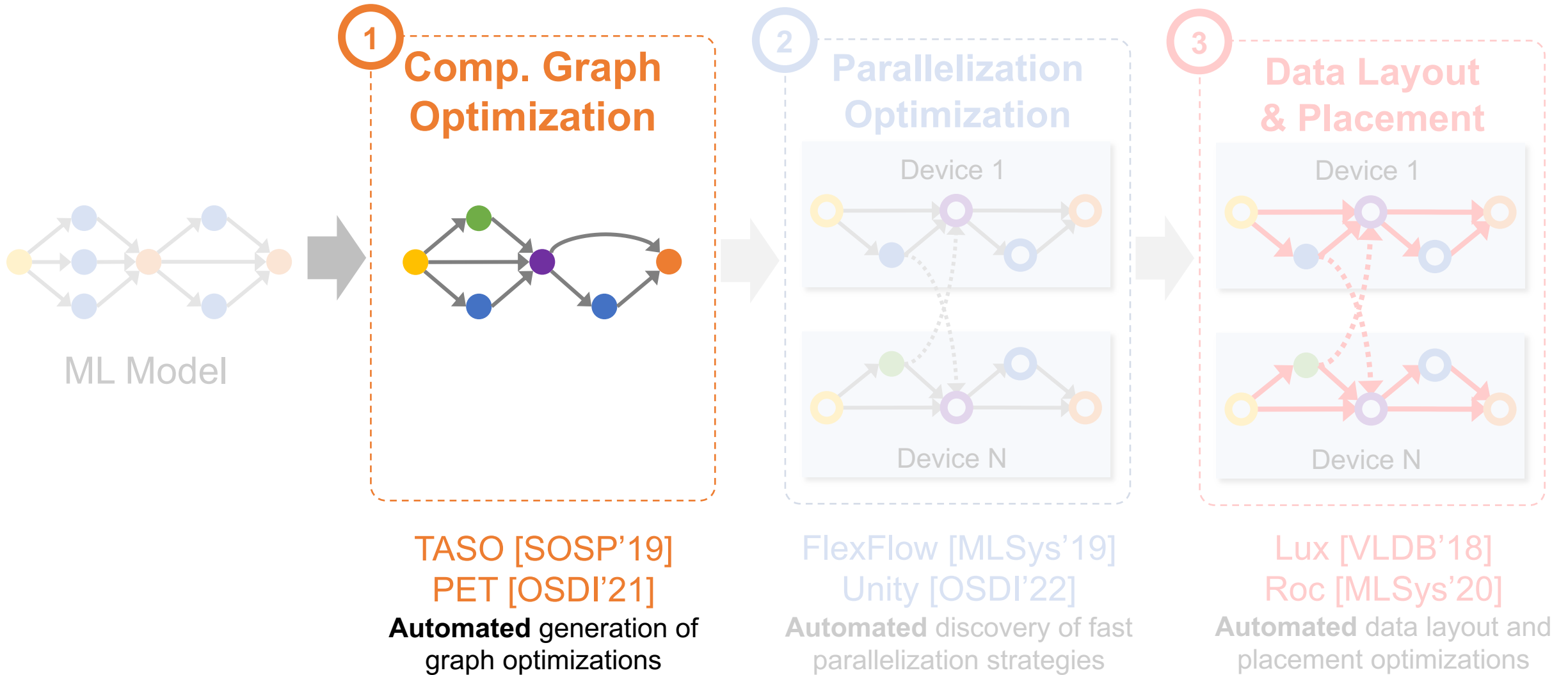
# Lesson 1: Automated Approaches Offer 3-10x Improvement



# Advantages of Automated Approaches

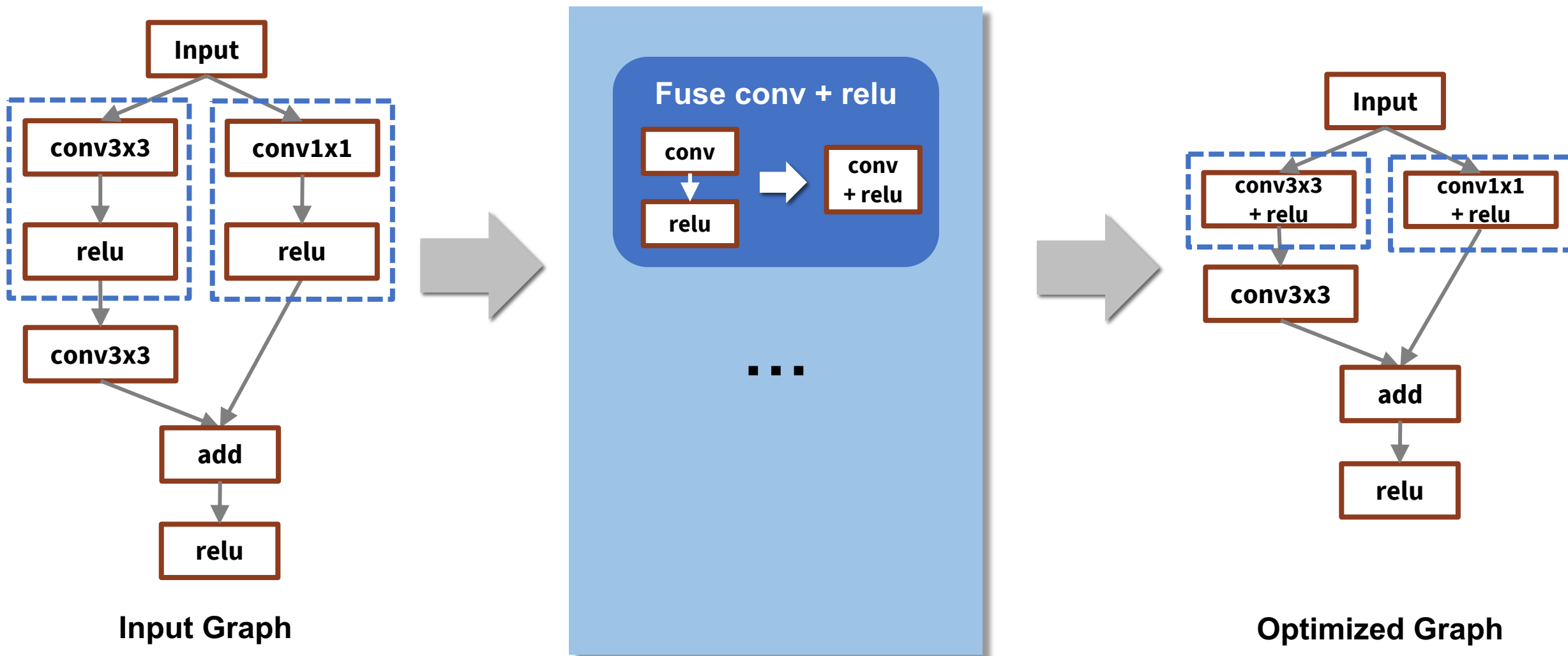
- **Better runtime performance:** discovering novel optimizations hard to manually designed, 3-10x speedup over manual optimizations
- **Less engineering effort:** code for discovering optimizations is generally much less than manual implementation of these optimizations
- **Stronger correctness guarantees:** using formal verification techniques

# Our Research: Automated Discovery of ML Optimizations





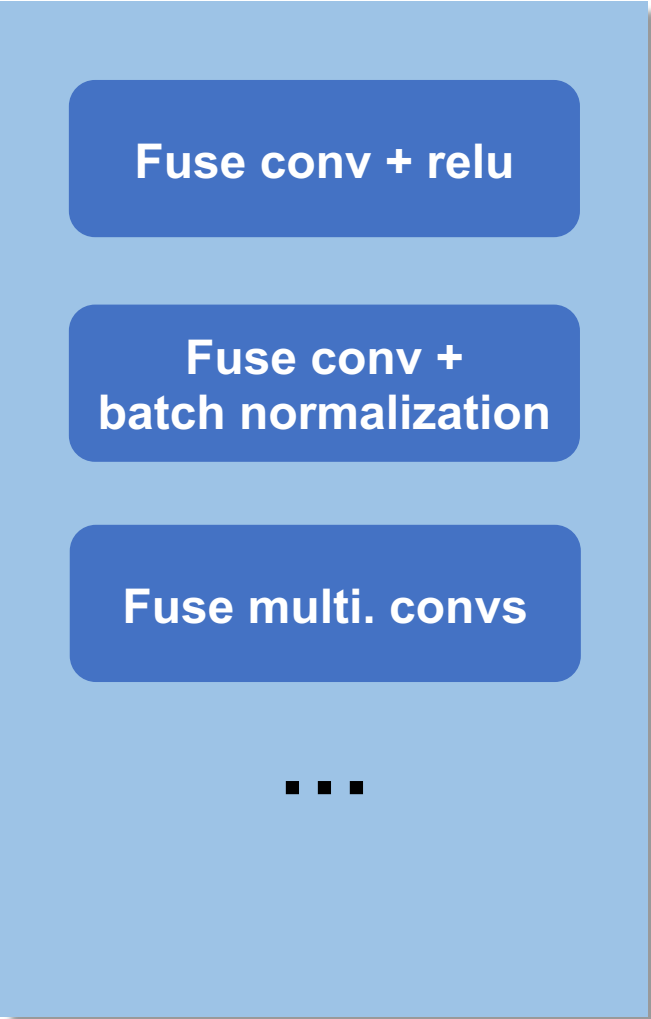
# Current Rule-based Graph Optimizations



Rule-based Optimizer

# Current Rule-based Graph Optimizations

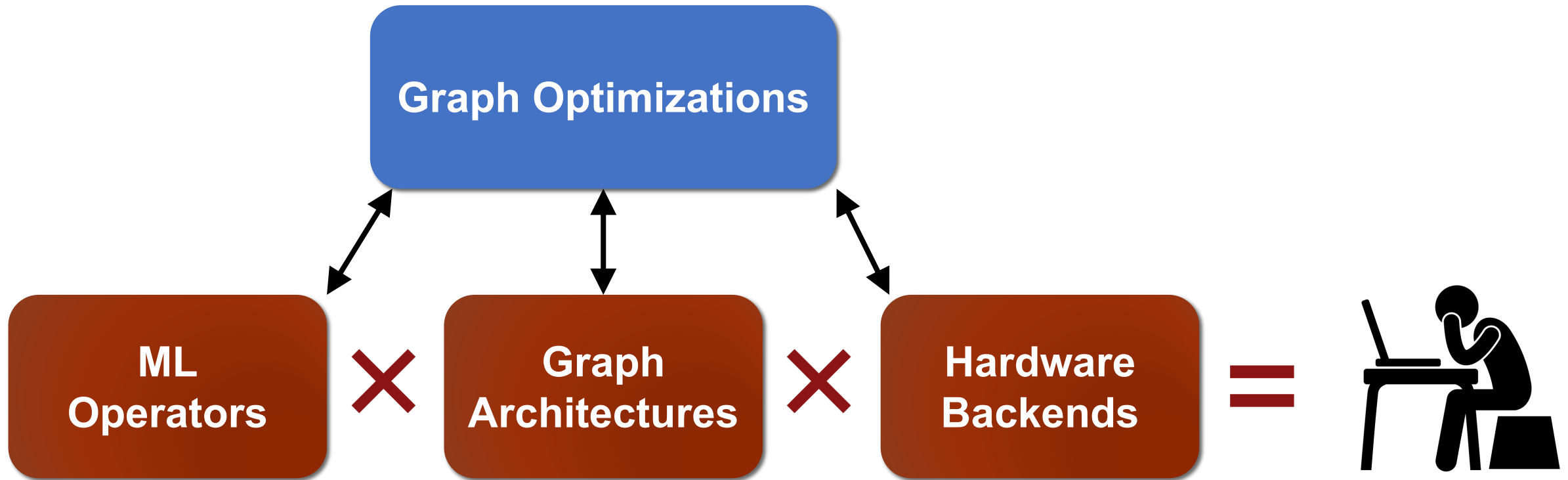
TensorFlow currently includes ~200 rules (~53,000 LOC)



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", "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
```

# Infeasible to Manually Design Graph Optimizations

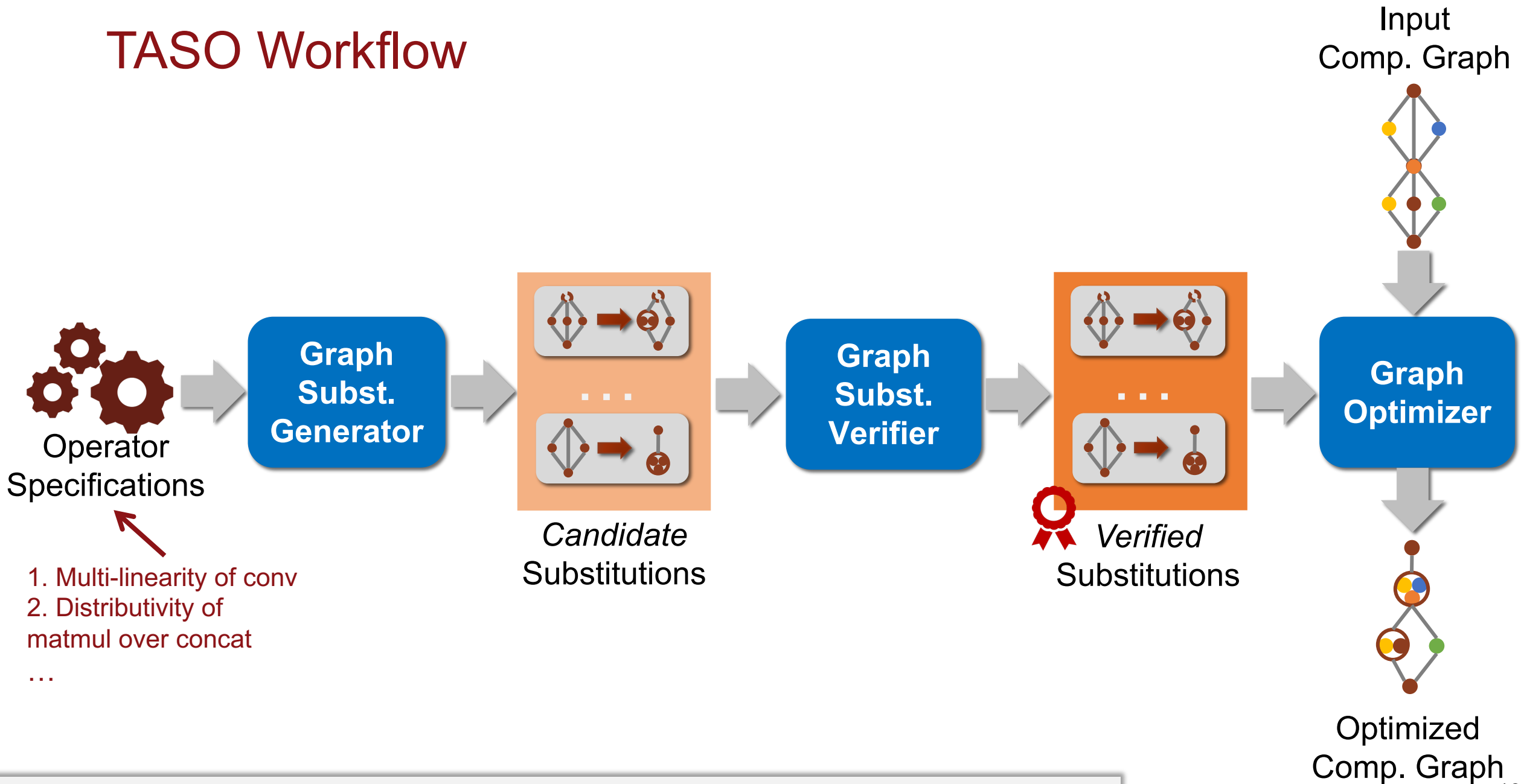


# TASO: Tensor Algebra SuperOptimizer

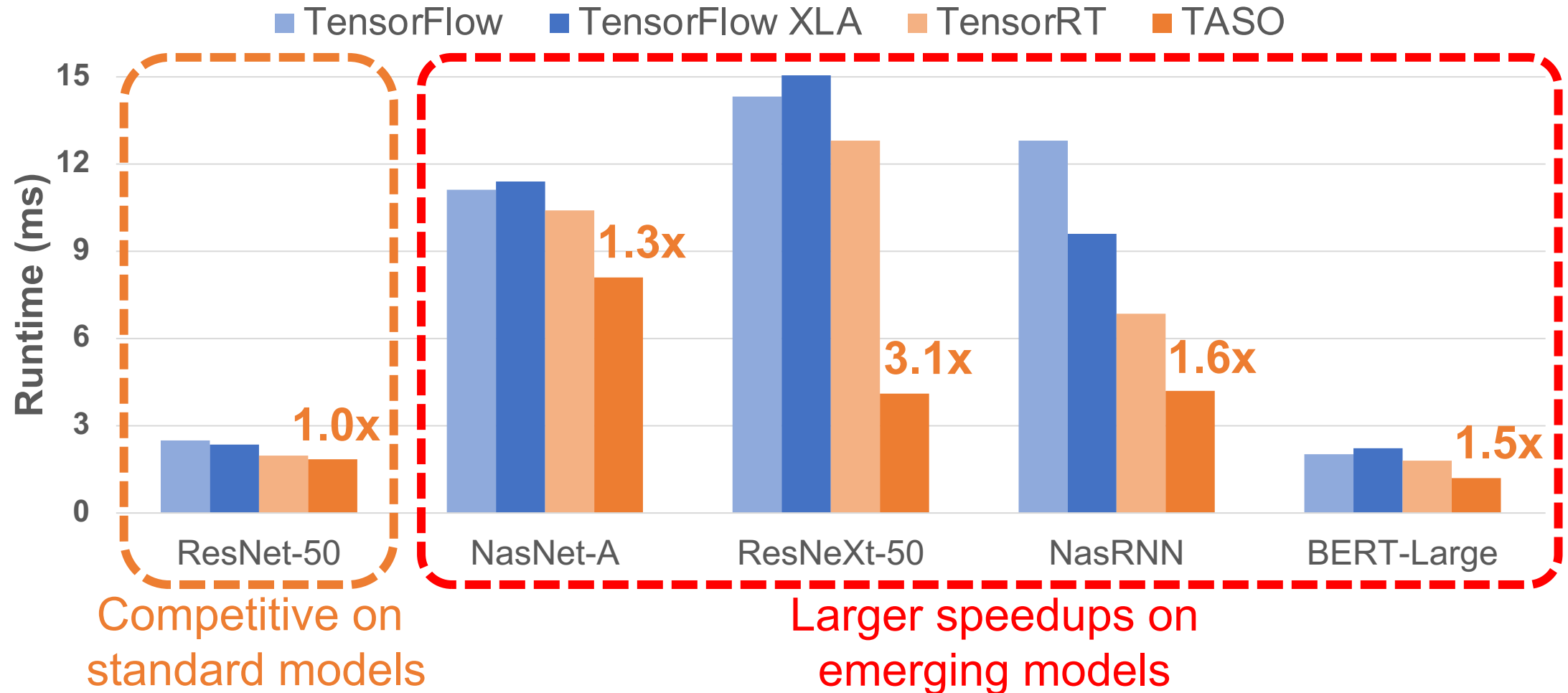
**Key idea:** replace manually-designed graph optimizations with *automated generation and verification* of graph substitutions for tensor algebra

- **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

# TASO Workflow



# 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



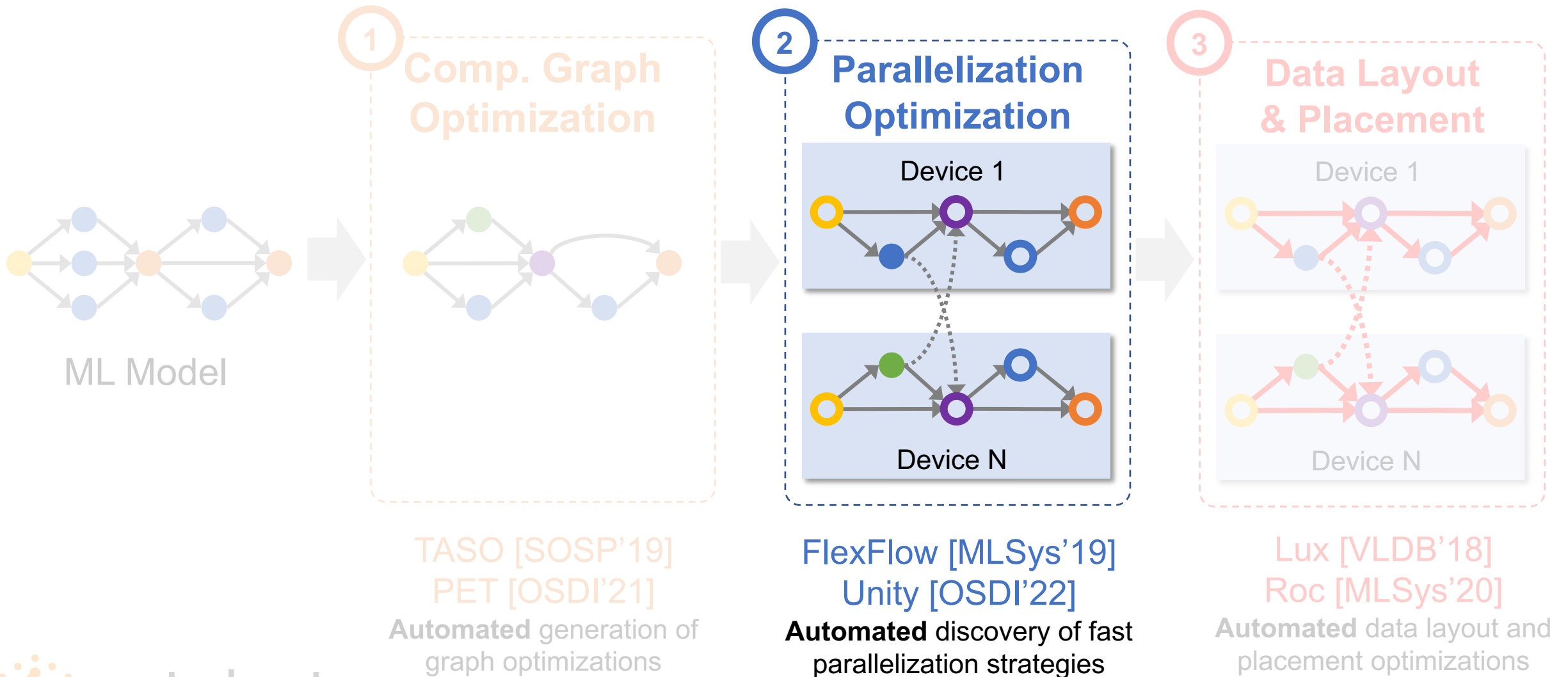
ONNX

ByteDance



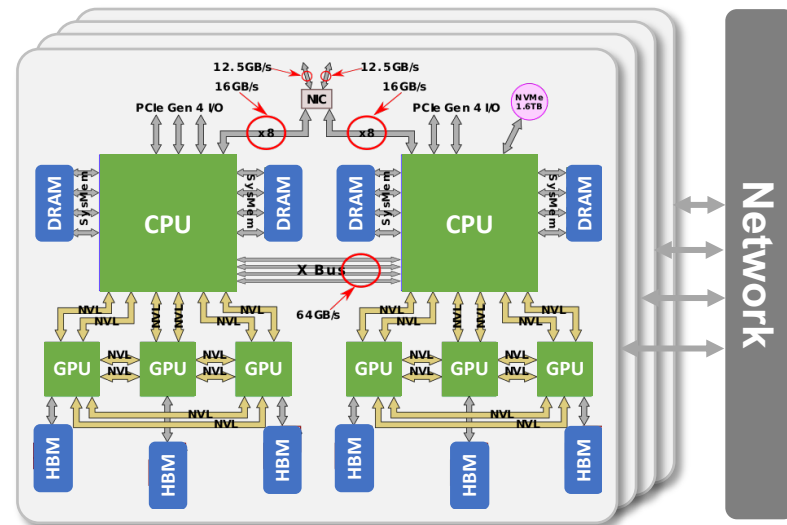
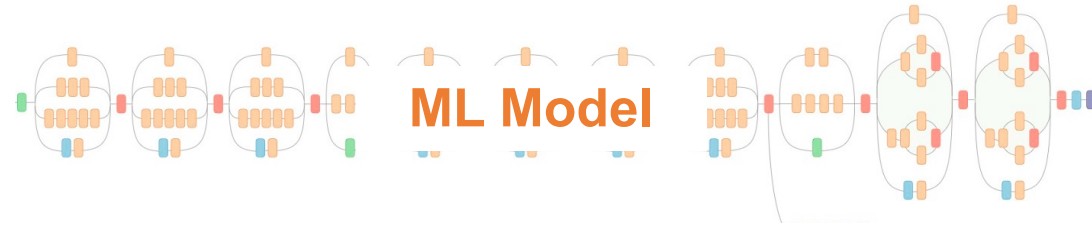
1. PET: Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections. OSDI'21.
2. TASO: Optimizing Deep Learning Computation with Automated Generation of Graph Substitutions. SOSP'19.
3. Optimizing DNN Computation with Relaxed Graph Substitutions. MLSys'19.
4. Exploring Hidden Dimensions in Parallelizing Convolutional Neural Networks. ICML'18.

# Our Research: Automated Discovery of ML Optimizations





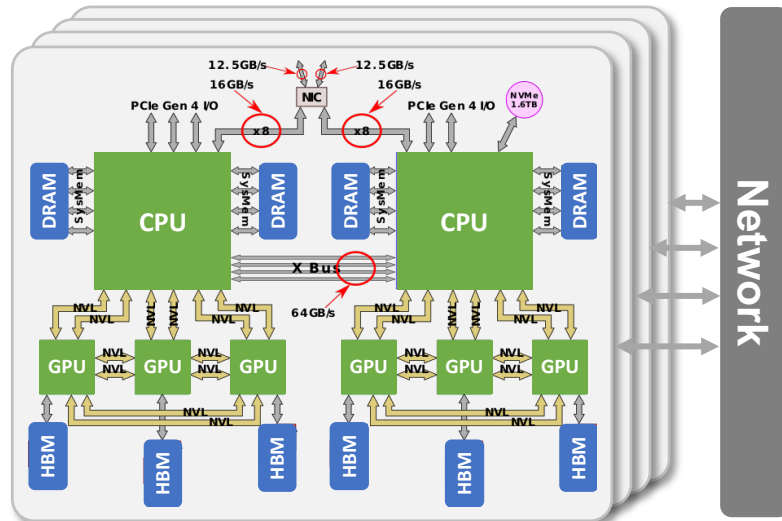
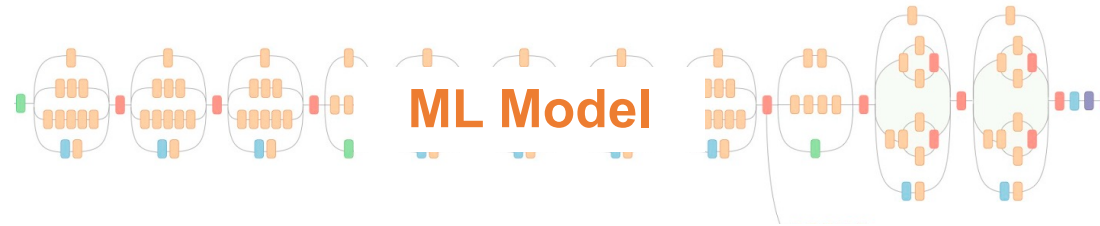
# Challenges of Parallelizing DNN Training



Need to simultaneously consider:

- Computation cost
- Communication overhead
- Resource usage
- Task scheduling

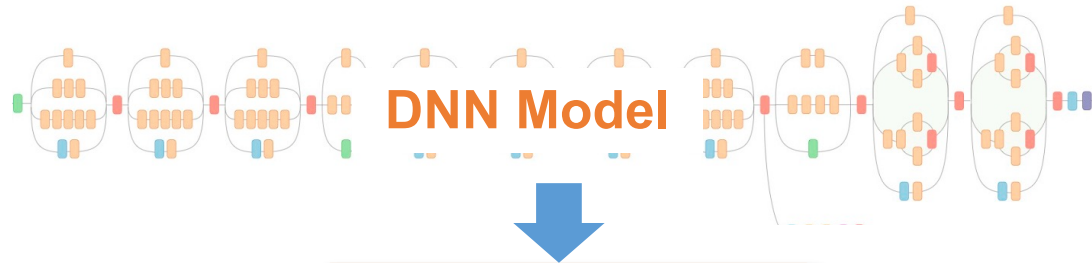
# Current Systems Rely on Manually Designed Strategies



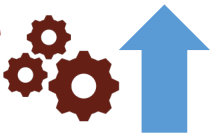
## Limitations:

- Hard to manually design
- Suboptimal performance
- Limited portability

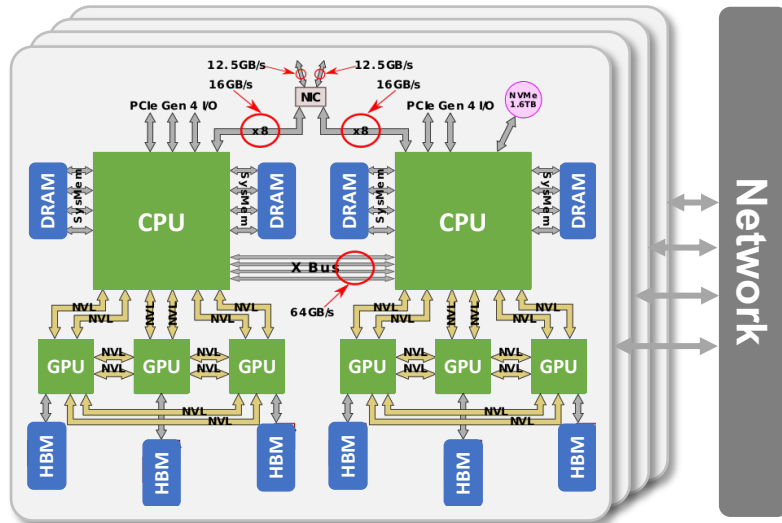
# FlexFlow: Automatically Optimizing DNN Parallelization



Hardware  
Topology



Fast Parallel  
Strategy



**Better Performance**

Up to 10x faster than manually designed strategies

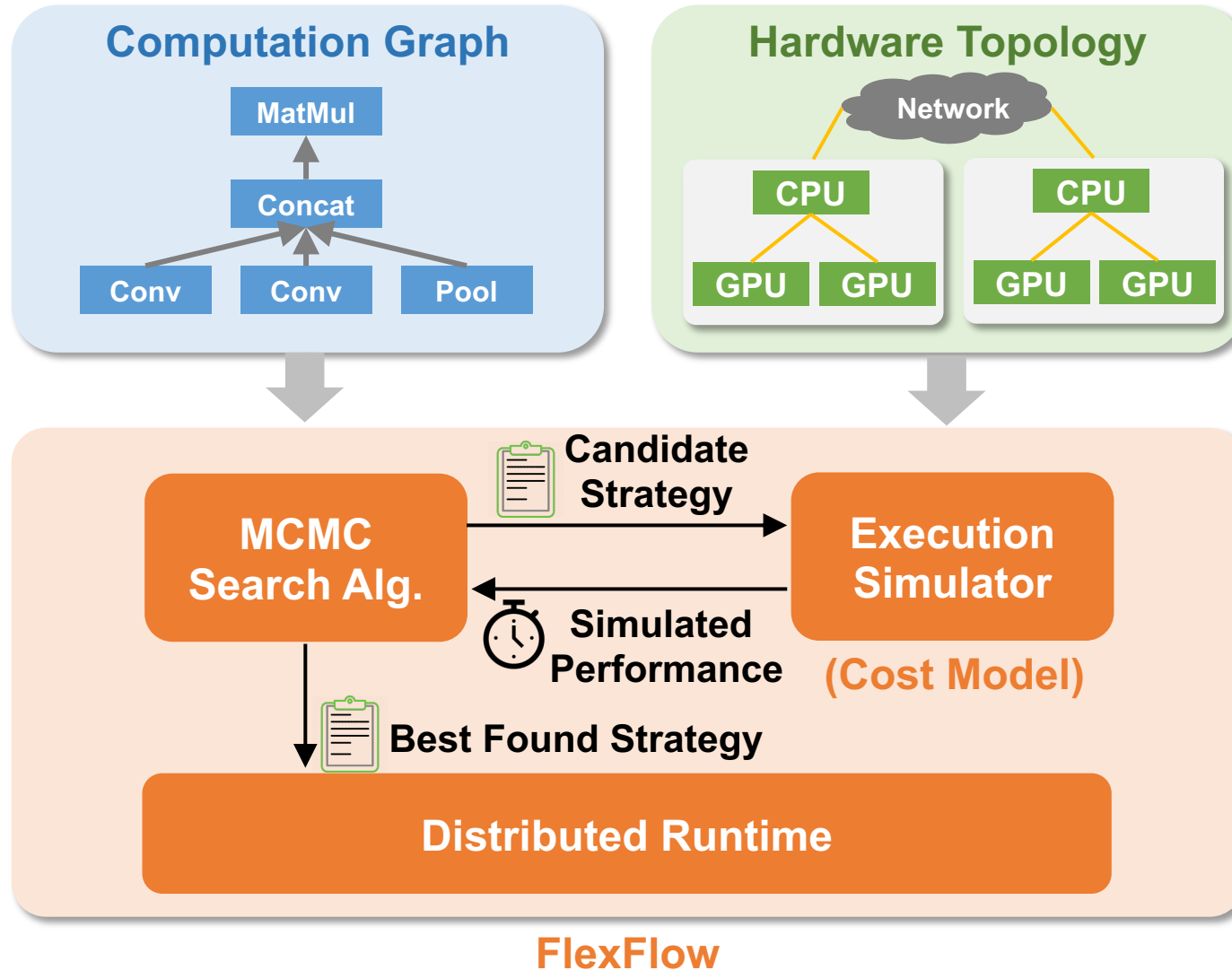
**Fast Deployment**

Minutes of automated search to discover performant strategies

**No Manual Effort**

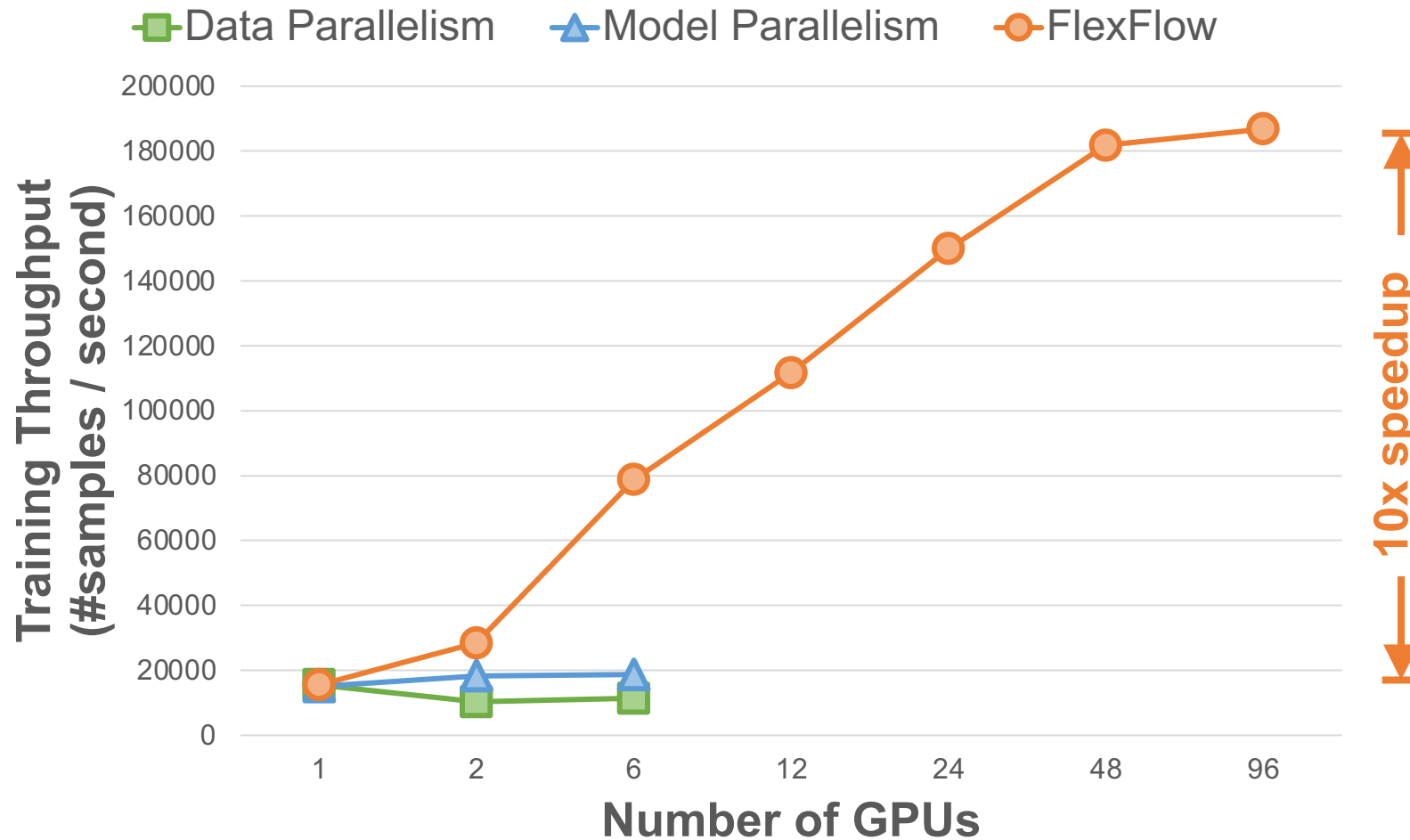
Automatically find strategies for new DNN models or hardware platforms

# FlexFlow Overview



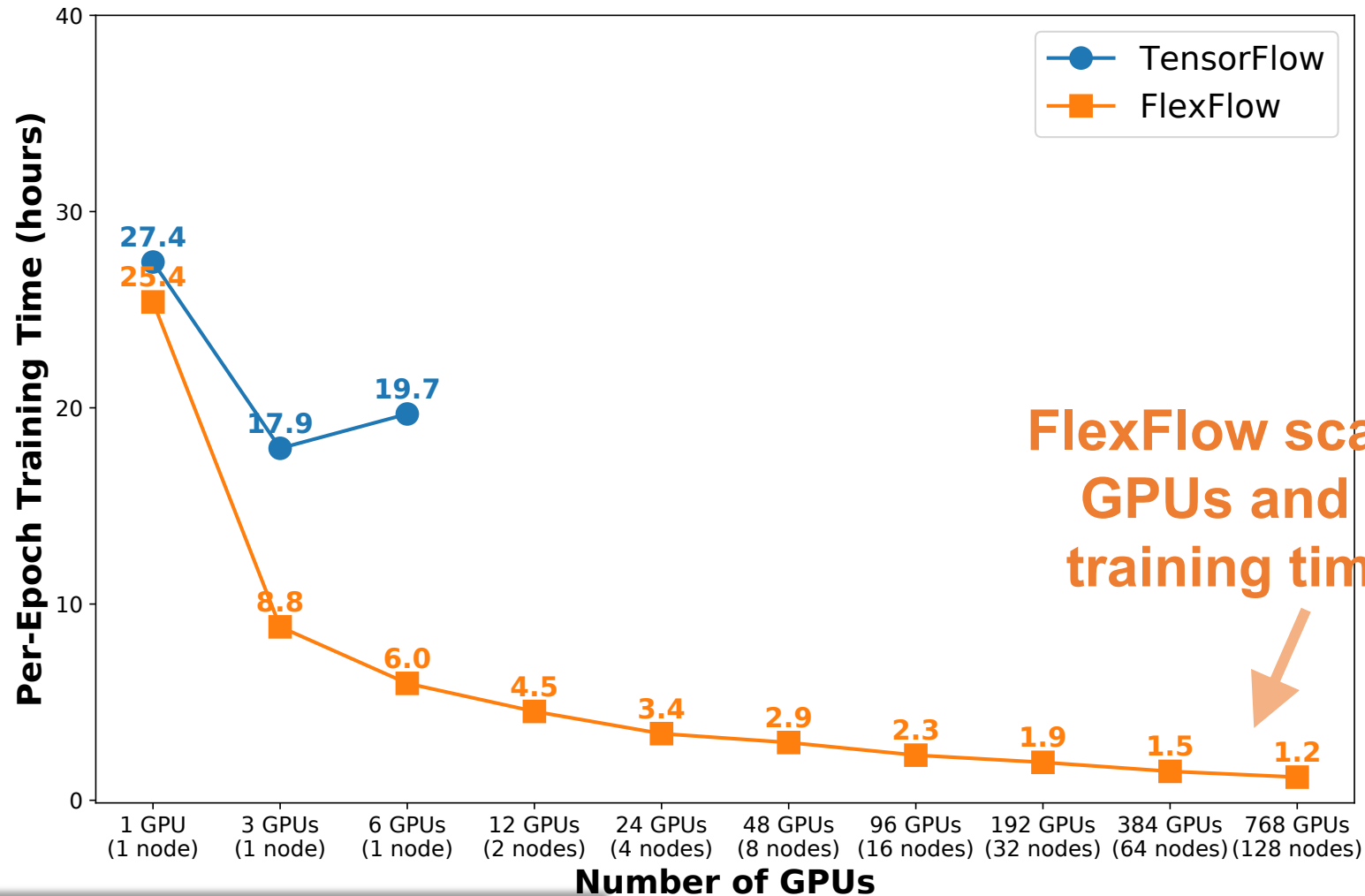
# Deep Learning Recommendation Model (DLRM)

## A deep learning model for ads recommendation



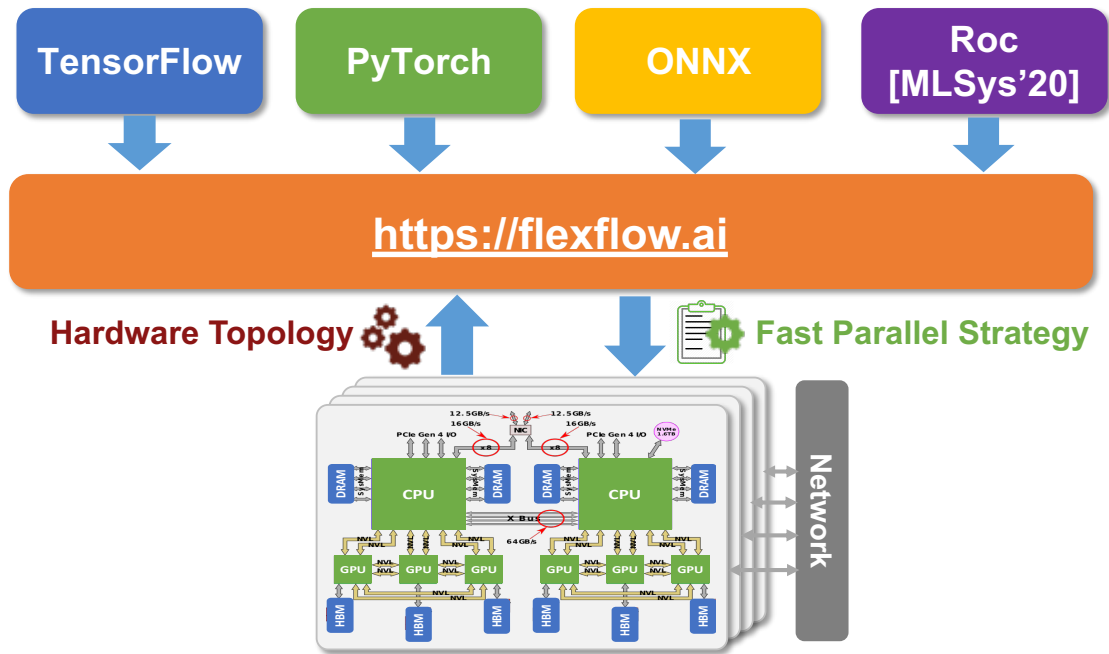
# ECP-CANDLE Training Performance

A deep learning model for precision medicine



FlexFlow scales to 768 GPUs and reduces training time by 15x

# FlexFlow: Automatically Discovering Fast and Scalable DNN Parallelization Strategies



<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.

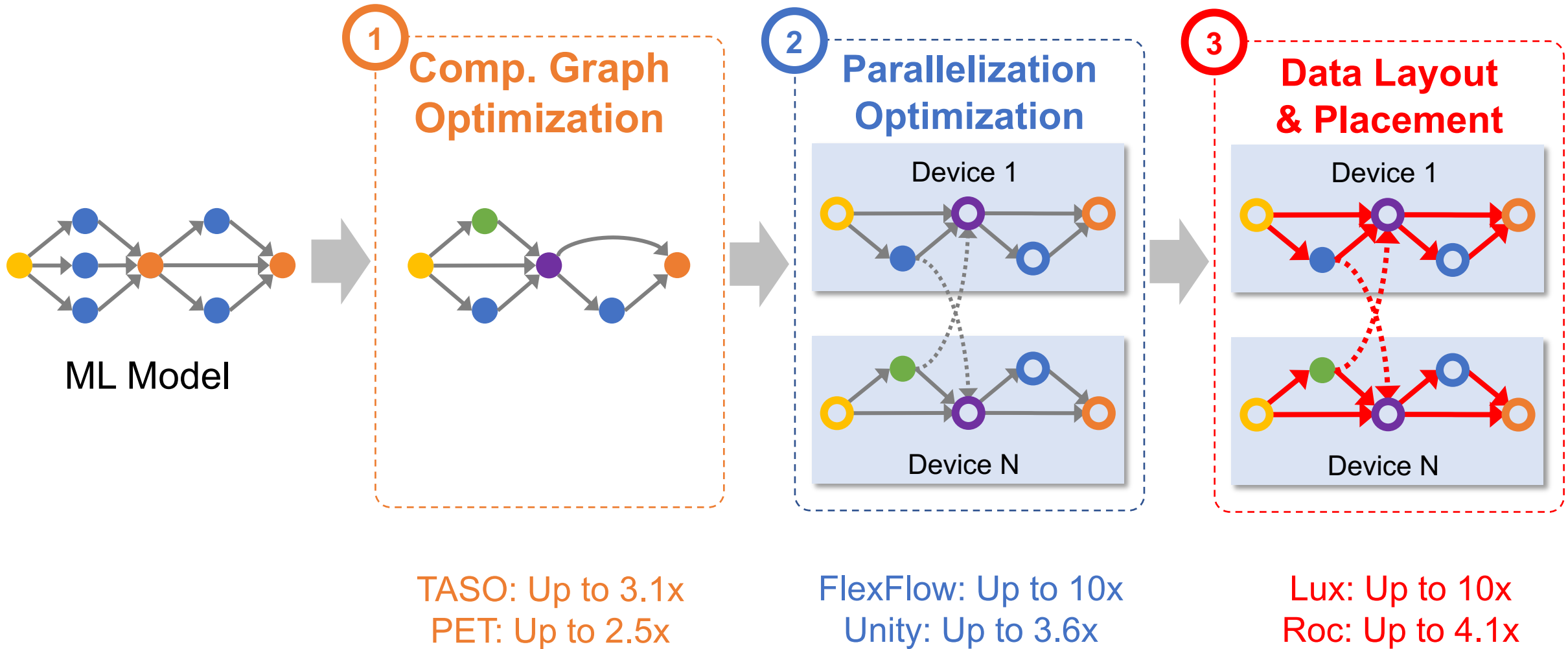
Learn more

### Large-Scale GNNs

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

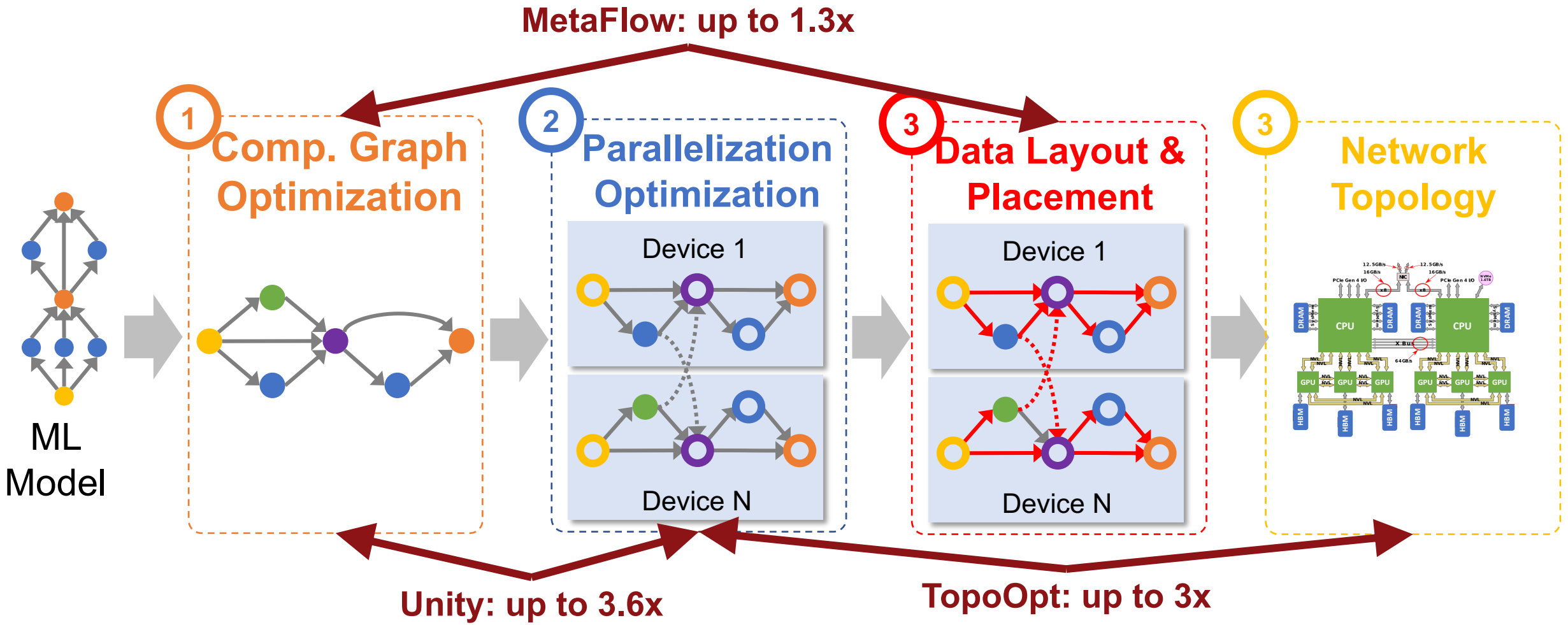
Learn more

# Lesson 1: Automated Approaches Offer 3-10x Improvement

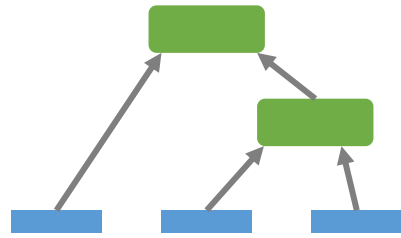




# Lesson 2: Joint Optimization is Critical to Performance



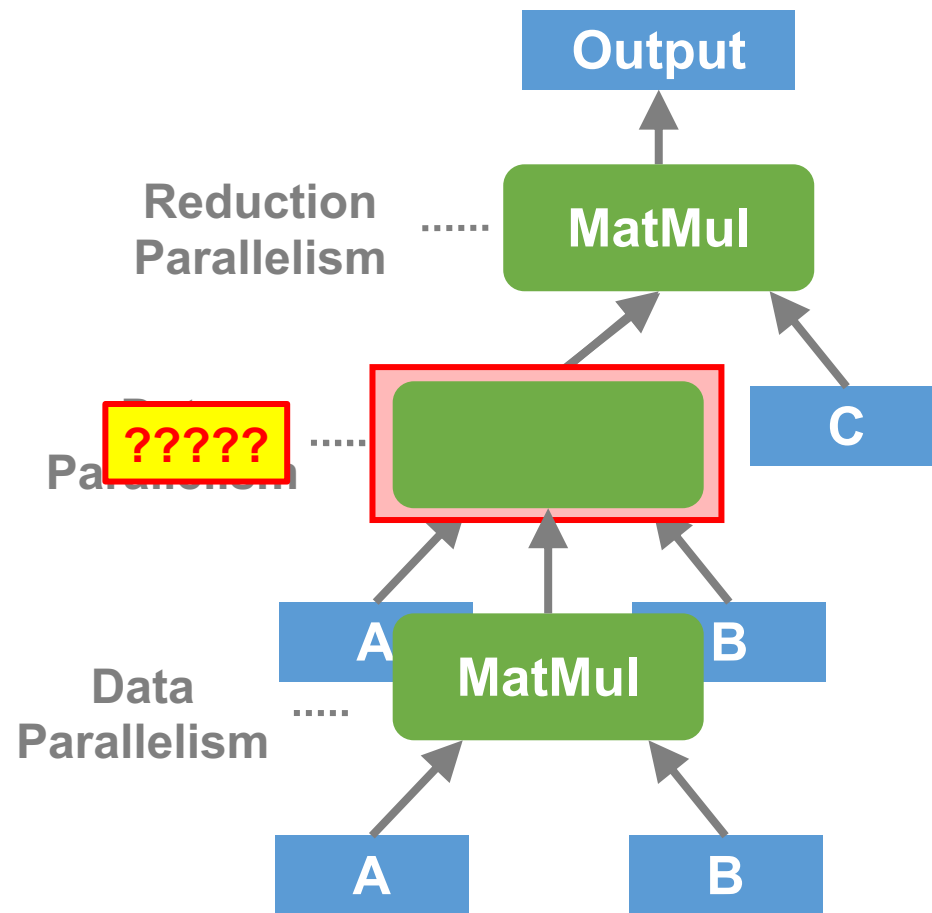
1. Unity: Accelerating DNN Training Through Joint Optimization of Algebraic Transformations and Parallelization. OSDI'22.
2. TopoOpt: Optimizing the Network Topology for Distributed DNN Training. NSDI'23.
3. MetaFlow: Optimizing DNN Computation with Relaxed Graph Substitutions. MLSys'19

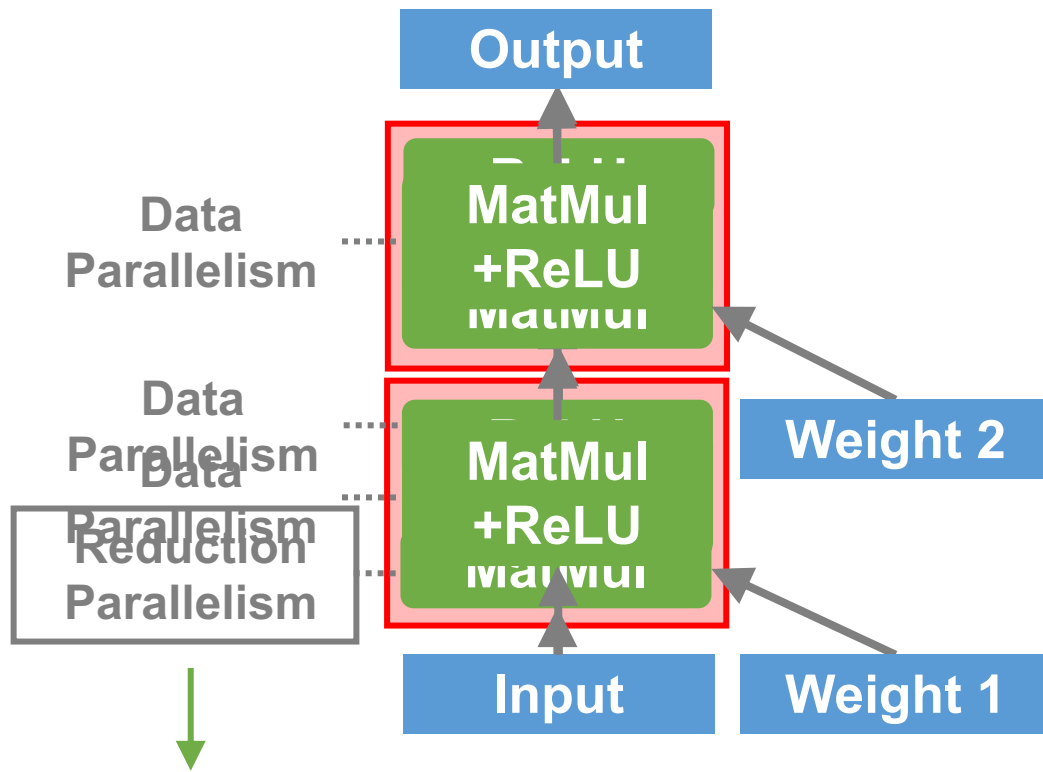


**Auto-Parallelization**

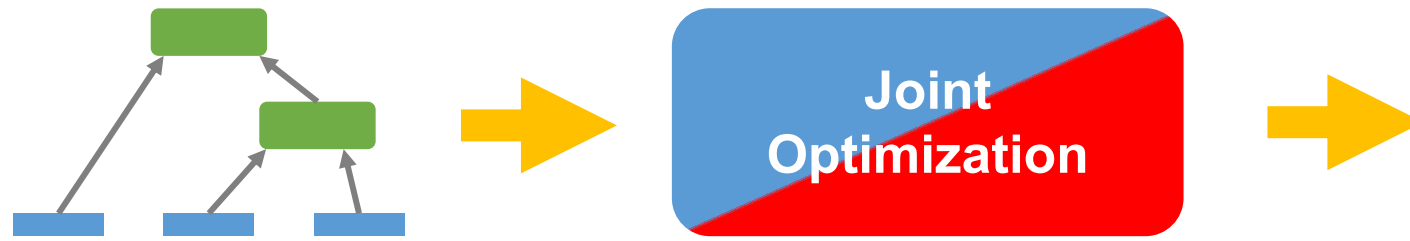
**Graph Optimization**







**≈ 6× less communication!**



1. Representation
2. Scalability

# Unity

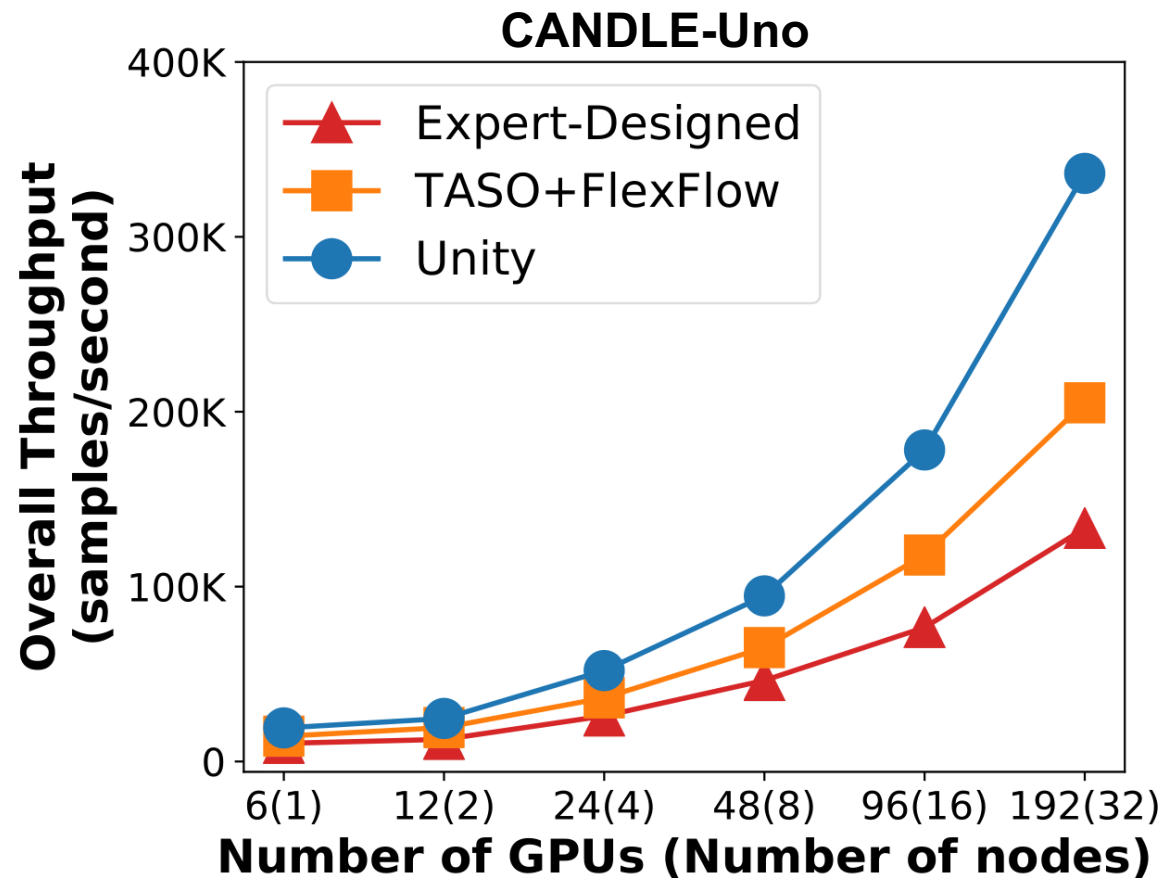
Representation

**Representation**  
Parallel Computation  
Graph (PCG)

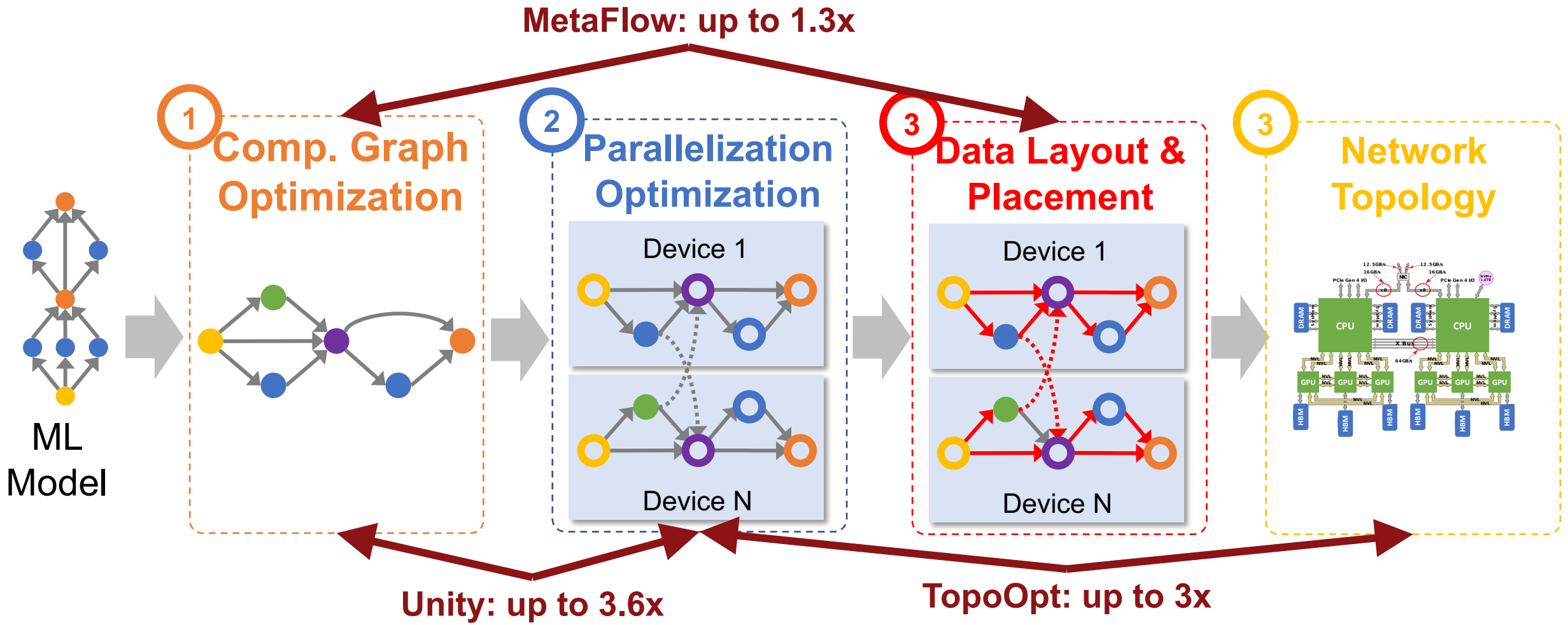
Scalability

**Scalability**  
Hierarchical Search  
Algorithm

# Joint Optimization Enables Better Performance and Scalability



# Lesson 2: Joint Optimization is Critical to Performance



1. Unity: Accelerating DNN Training Through Joint Optimization of Algebraic Transformations and Parallelization. OSDI'22.
2. TopoOpt: Optimizing the Network Topology for Distributed DNN Training. NSDI'23.
3. MetaFlow: Optimizing DNN Computation with Relaxed Graph Substitutions. MLSys'19



# Lesson 3: Combining ML and Systems Optimizations is Promising but Challenging

## Systems Optimizations

- Graph Transformations
- Auto Parallelization
- Kernel Generation
- Data Layout and Placement

## ML Optimizations

- Quantization
- Pruning
- Distillation
- Neural Architecture Search

# Lesson 3: Combining ML and Systems Optimizations is Promising but Challenging

Systems  
Optimizations

👍 Pro: preserve functionality

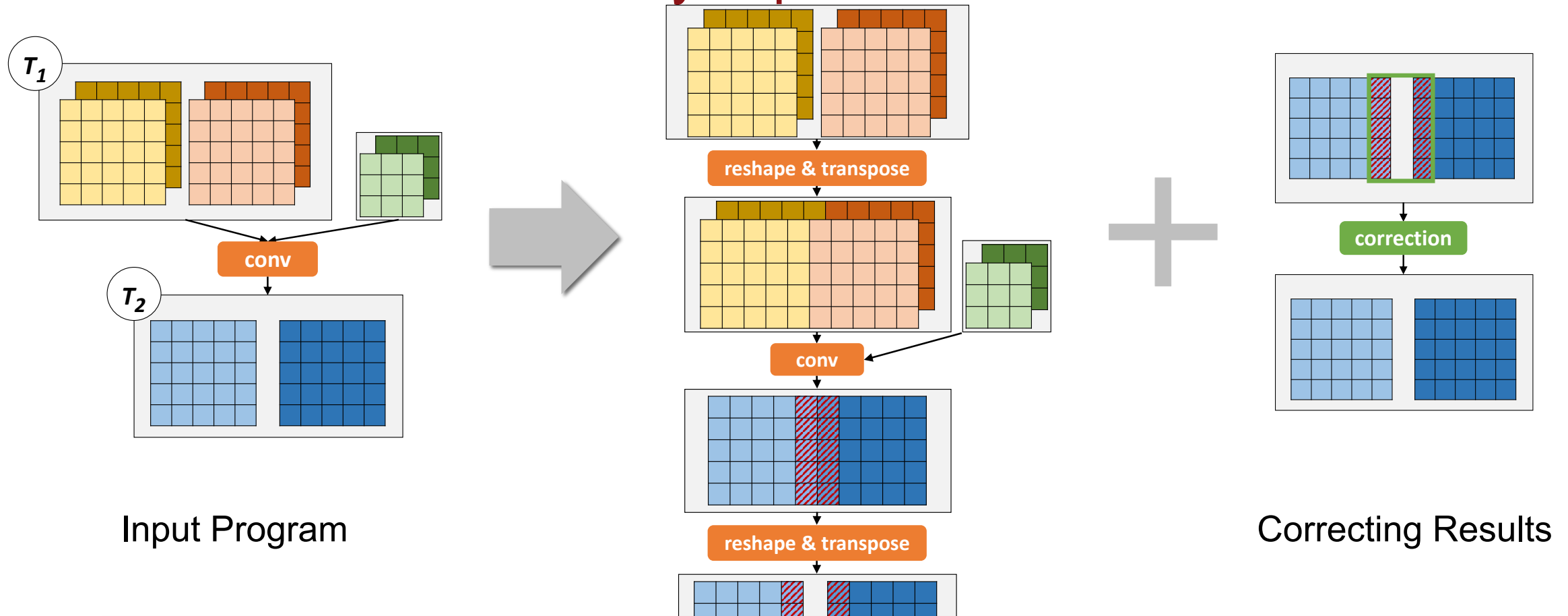
ML  
Optimizations

👍 Pro: better performance

- Faster ML operators
- Less Computation

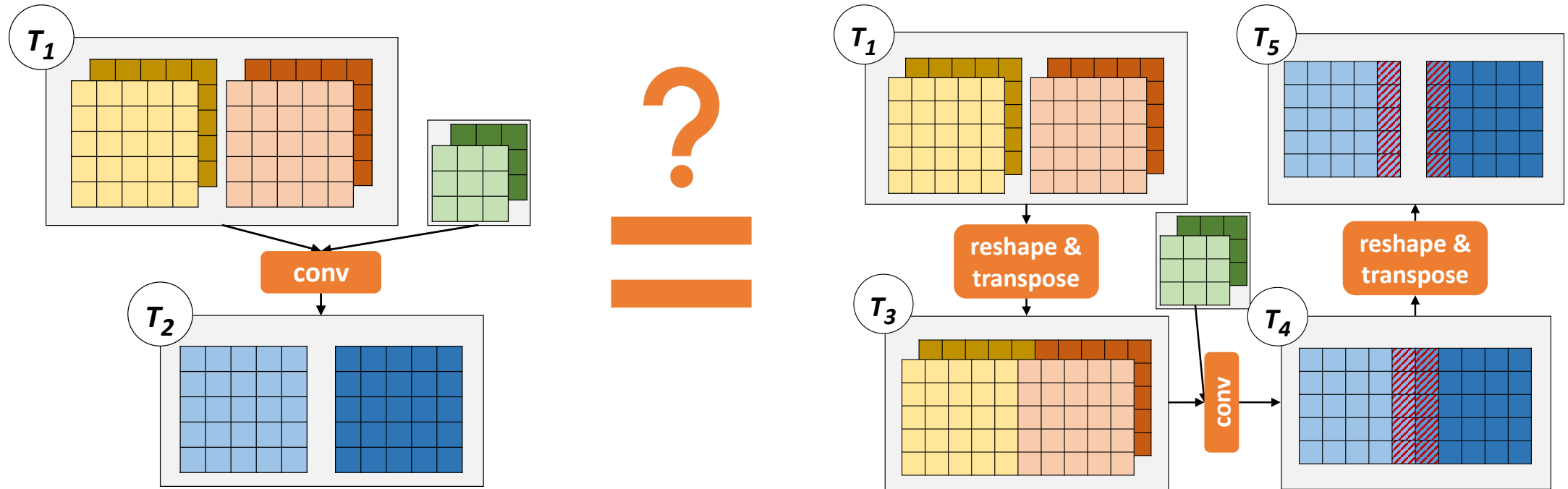
**Achieve the best of both worlds?**

# Hidden Treasure: Partially Equivalent Transformations



- Transformation and correction lead to **1.2x** speedup for ResNet-18
- Correction preserves end-to-end equivalence

# Hidden Treasure: Partially Equivalent Transformations



1. Which part of the computation is not equivalent?
2. How to correct the results?



# PET

- **Tensor program optimizer** with partially equivalent transformations and automated corrections
- **Larger optimization space** by combining fully and partially equivalent transformations
- **Better performance**: outperform existing optimizers by up to 2.5x
- **Correctness**: automated corrections to preserve end-to-end equivalence

Equivalent  
Optimizations



Partially-Equivalent  
Optimizations



Non-Equivalent  
Optimizations



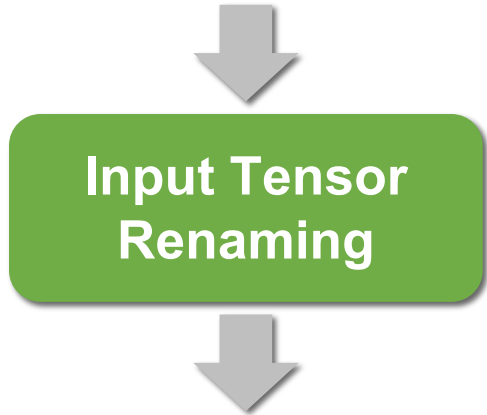
## Three Lessons

1. Automated approaches can offer **3-10x** improvement on most tasks
2. Joint optimization is **critical**
3. Combing systems and ML optimizations is **promising but challenging**

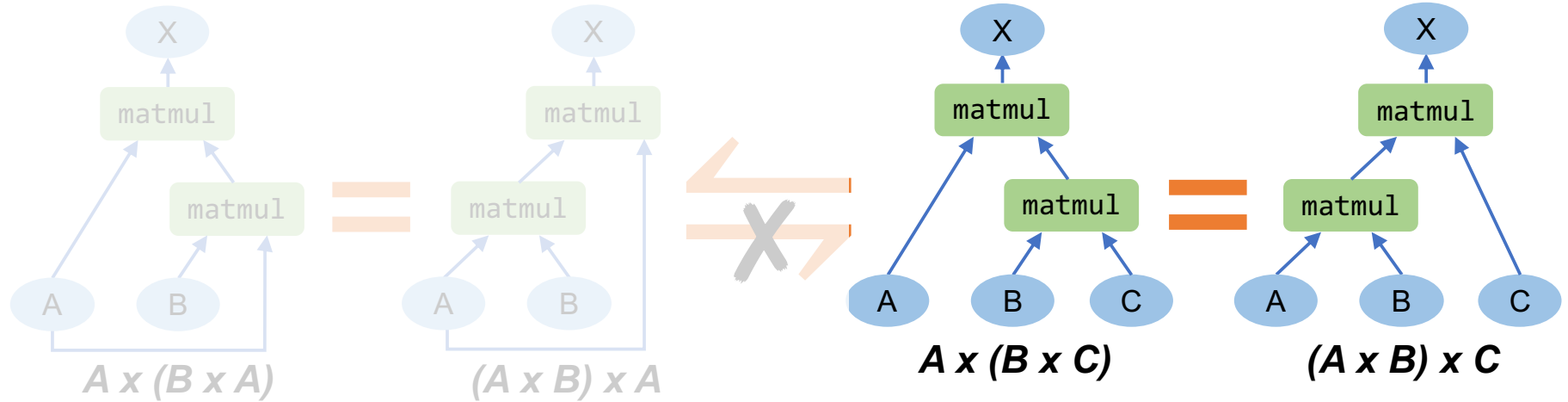
# Pruning Redundant Substitutions



28,744 substitutions



17,346 substitutions







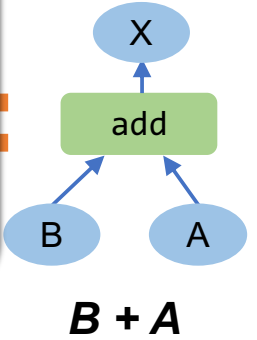
# Pruning Redundant Substitutions

28,744 substitutions

Input Tensor Renaming

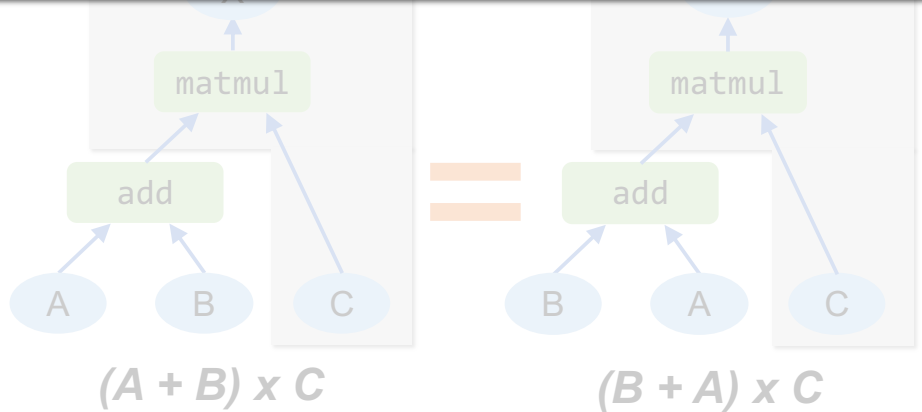


Pruning techniques reduce the number of candidate substitutions by 39x



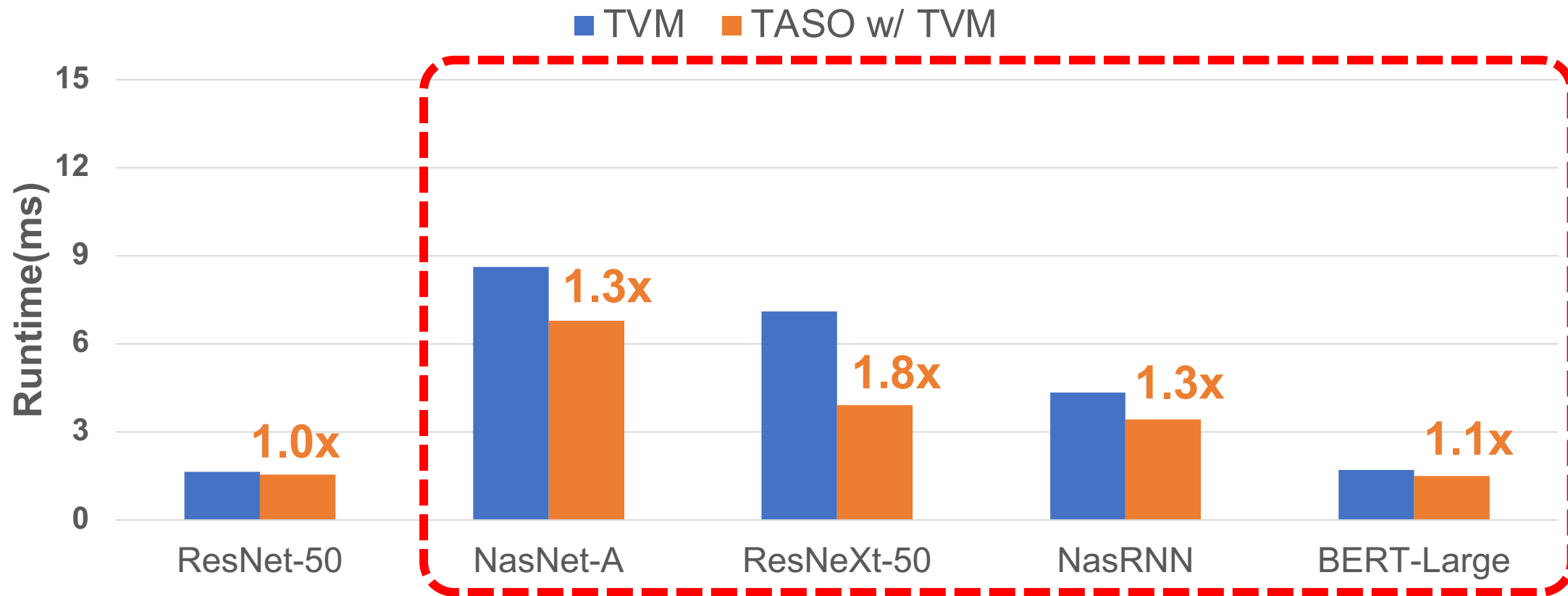
17,346 substitutions

Common Subgraphs



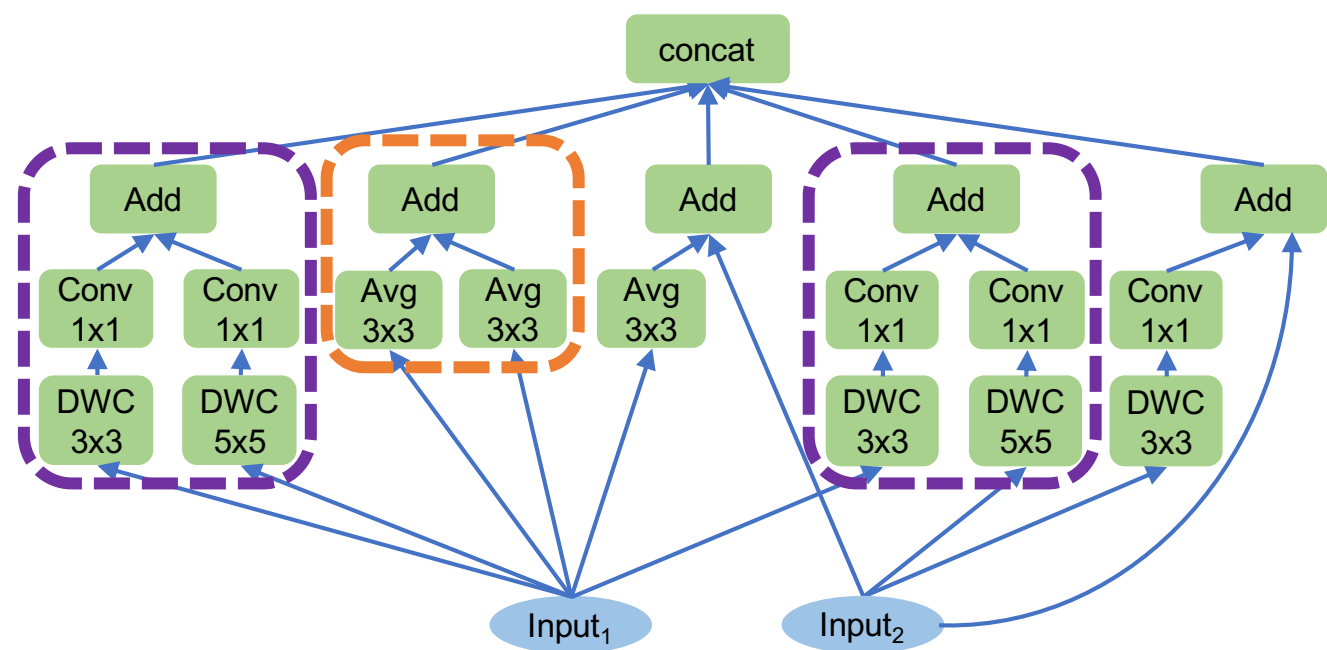
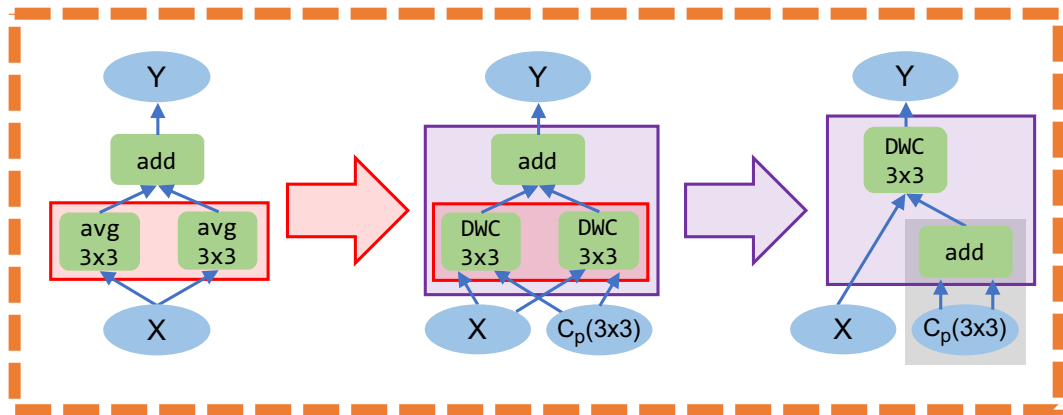
743 substitutions

# End-to-end Inference Performance (TVM)

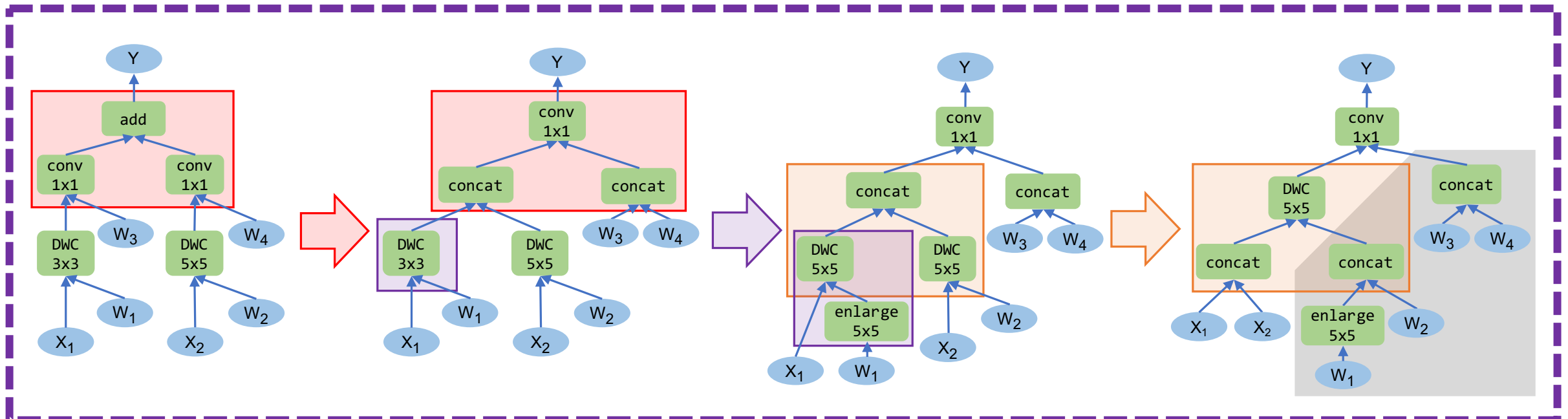


**Larger speedups on emerging models**

# Case Study: NASNet



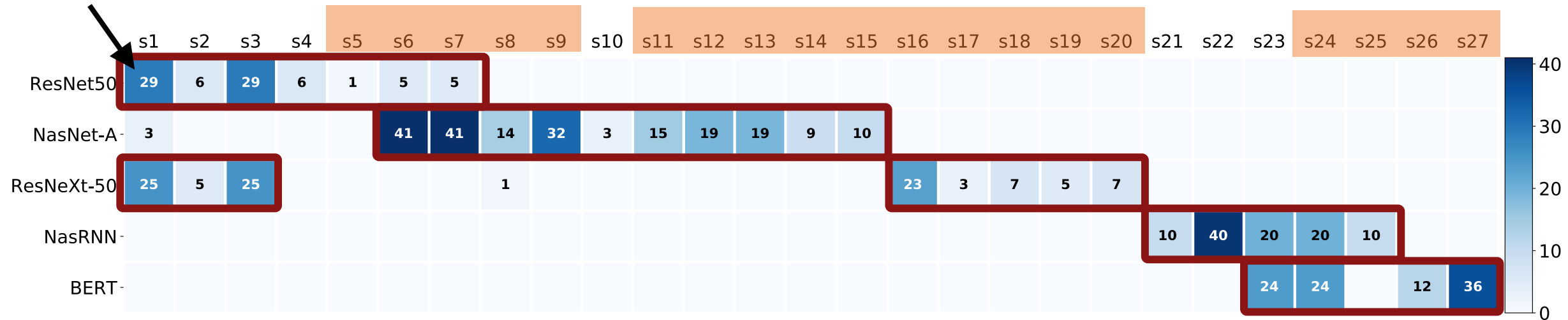
\*DWC: depth-wise convolution



# Heatmap of Used Substitutions

Not covered in TensorFlow

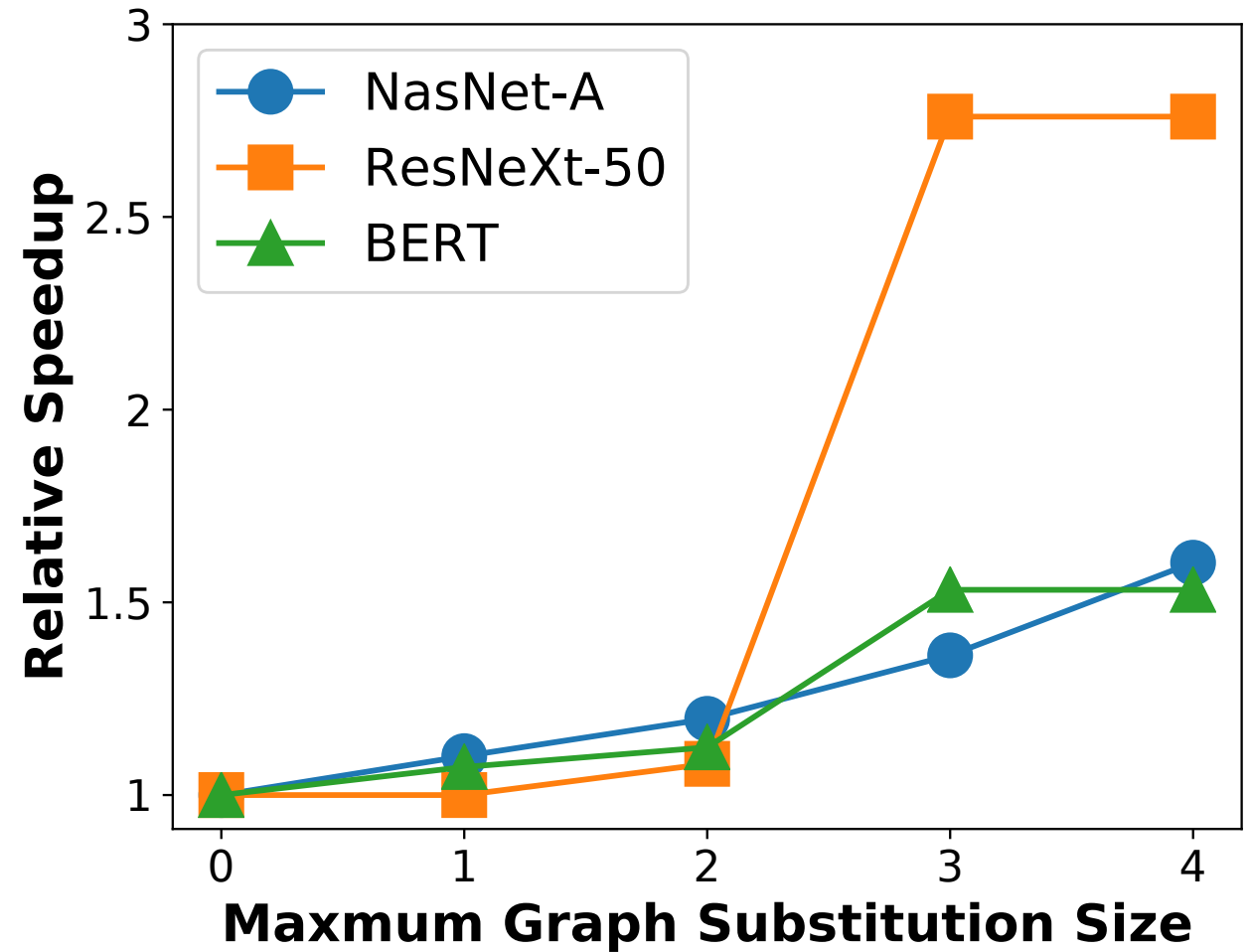
How many times a subst. is used to optimize a DNN



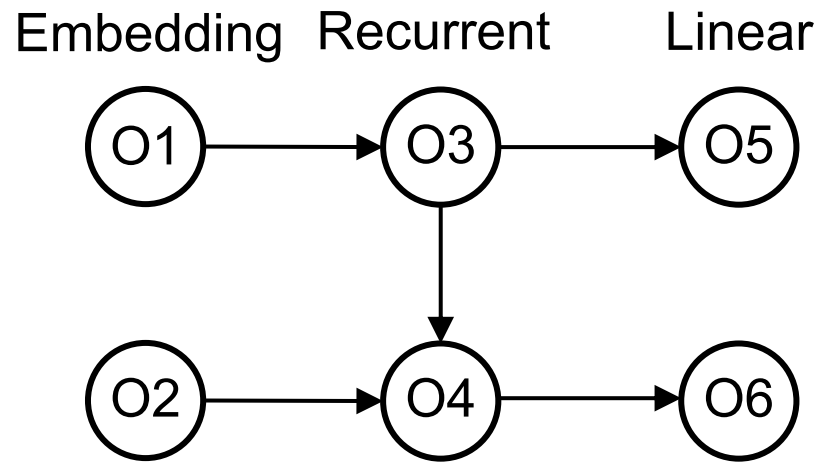
Different DNN models require **different** substitutions.

# Scalability Analysis

Max Num. of Operators	Mem. to Cache Fingerprints
1	0.9 KB
2	35.8 KB
3	6.9 MB
4	5.35 GB

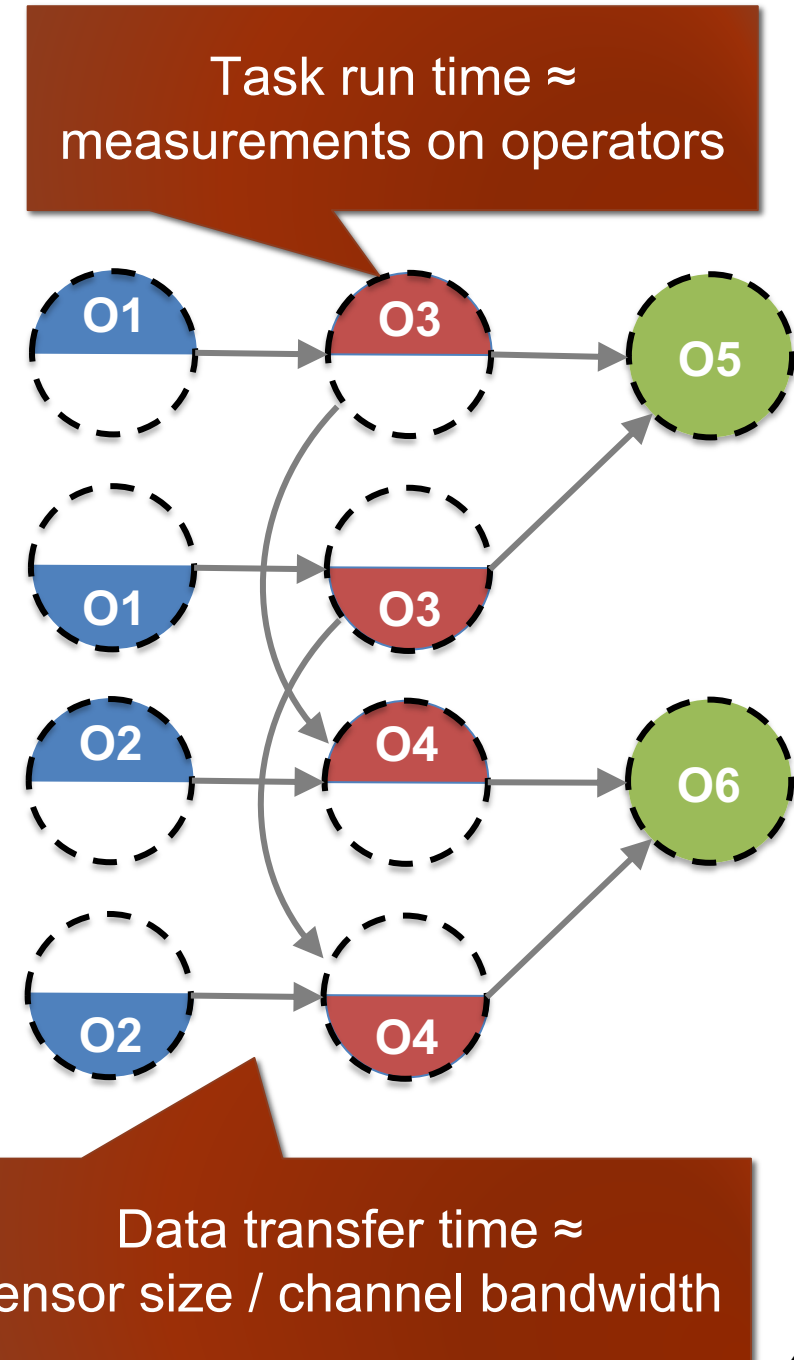


# Execution Simulator

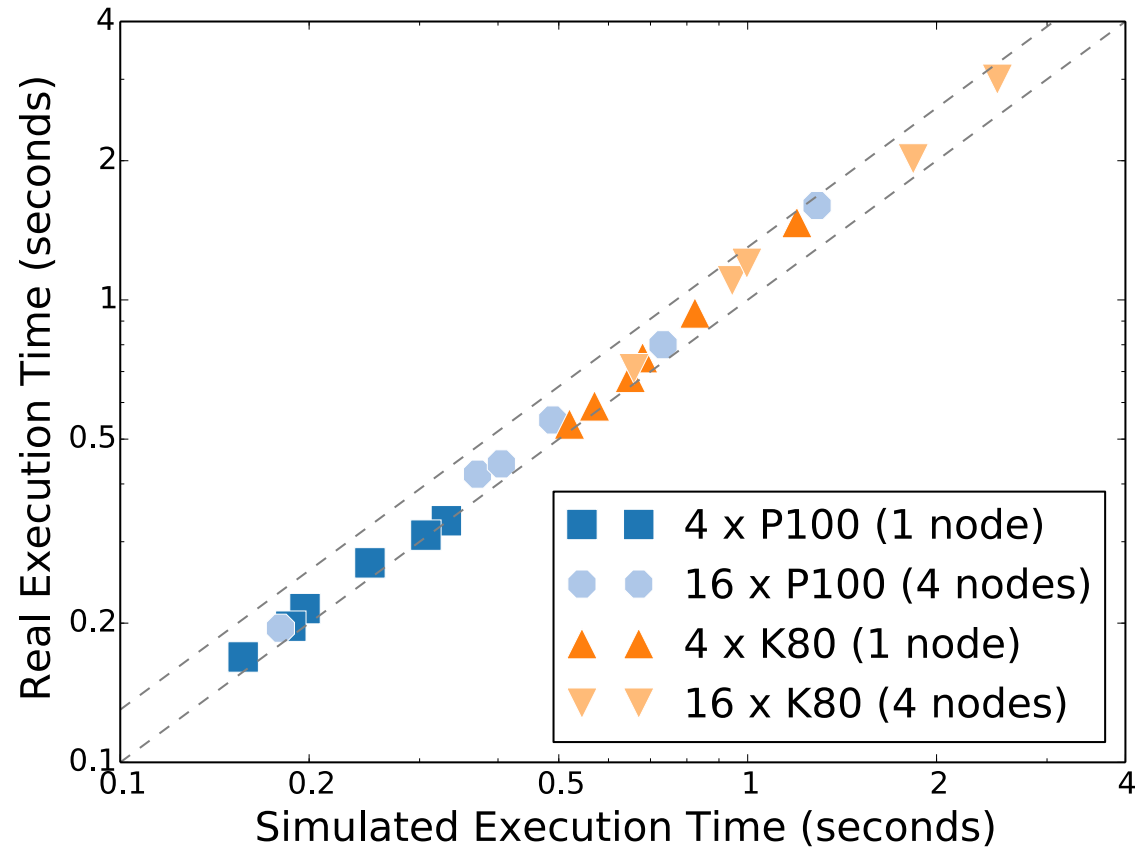


**ML Architecture**

Parallelization Strategy	
O1, O2	Degree(sample) = 2 GPU1
O3, O4	Degree(sample) = 2 GPU2
O5, O6	Degree(sample) = 1 GPU3



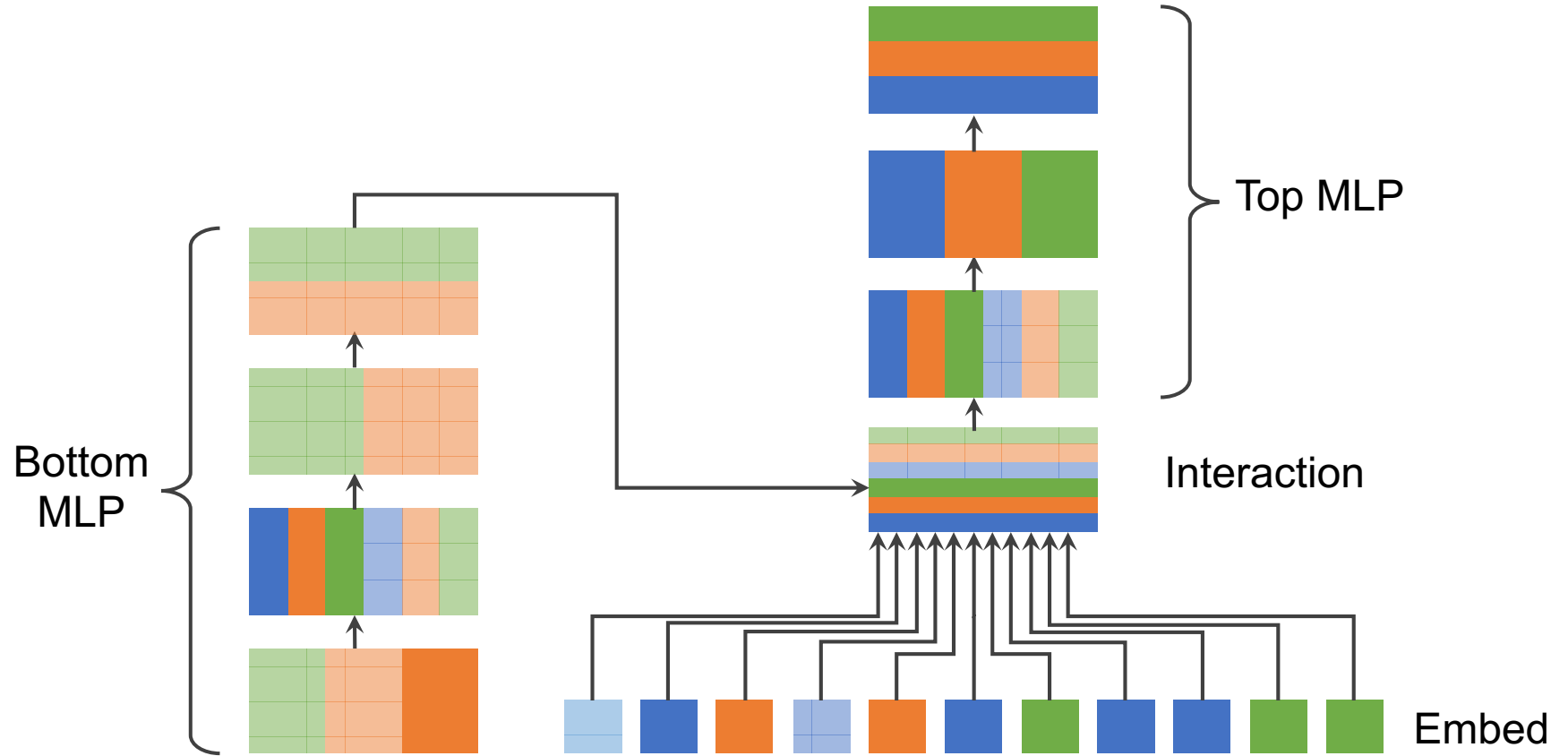
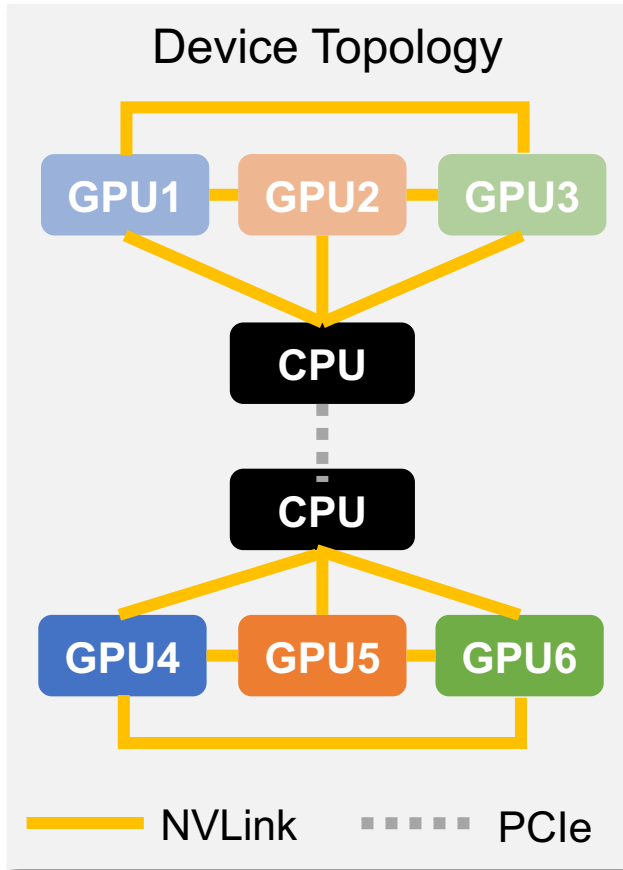
# Execution Simulator



Relative difference between simulated and actual execution time is less than 30%

Simulated execution time preserves real execution time ordering

# Case Study



Layers	Communication	Compute	Strategy
Embed	Heavy	Light	Parallelism across <b>Operators</b>
Interaction	Light	Heavy	Parallelism across <b>Samples</b>
MLP	Heavy	Heavy	Parallelism across <b>Samples and Parameters</b>



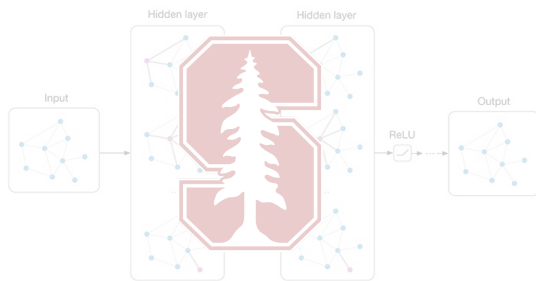
# FlexFlow Impact

The Facebook logo, consisting of the word "facebook" in white lowercase letters on a dark blue rectangular background.

Facebook uses FlexFlow to train production ML models. Increase training throughput by 10x.



Used by LANL to train ML models for precision medicine. Reduce training time from days to hours.



Improve the accuracy, scalability, and performance of graph neural networks.