

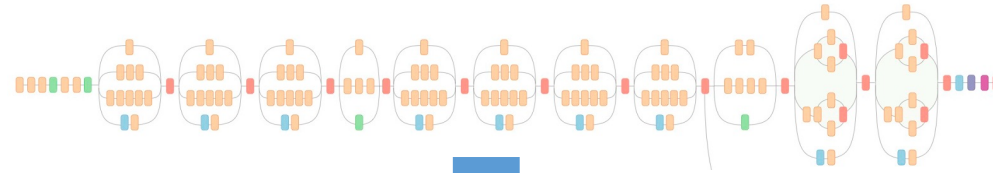
15-442/15-642: Machine Learning Systems

Three Lessons Learned from ML Systems

Tianqi Chen and Zhihao Jia

Carnegie Mellon University

ML Systems are a Key Ingredient in ML

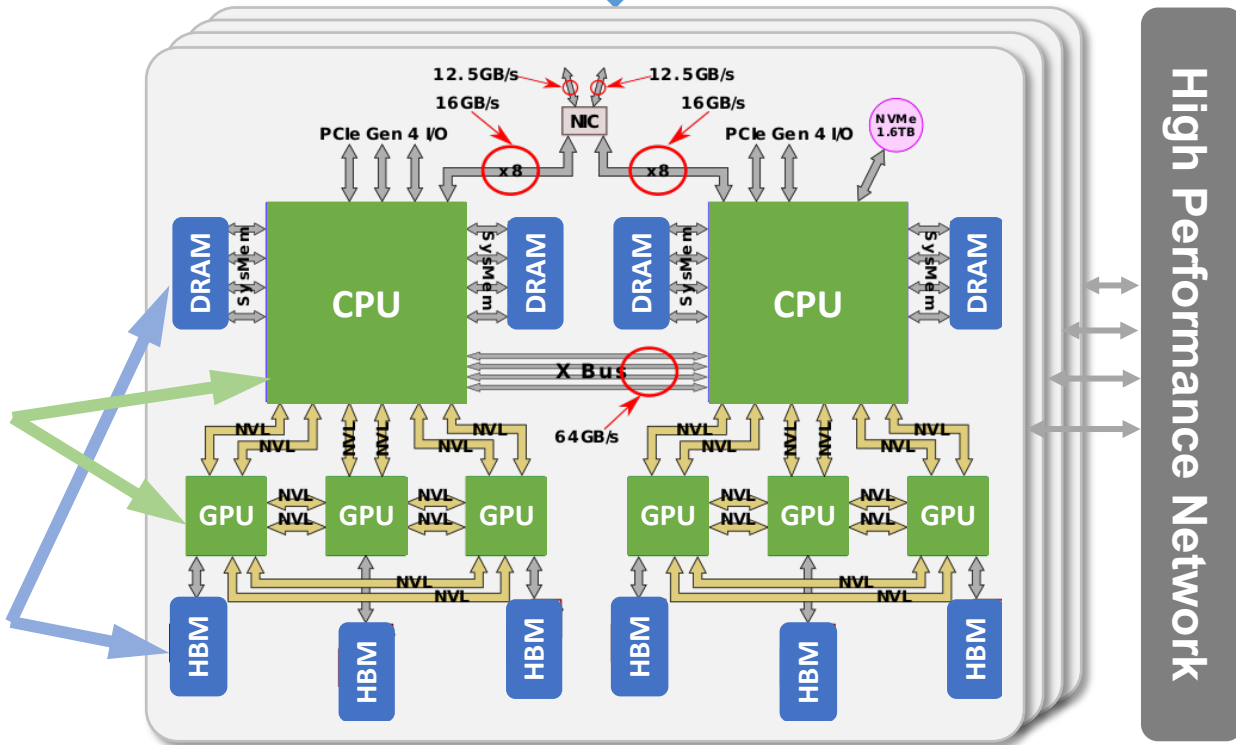


ML Model

ML Systems

Heterogenous Processors

Heterogenous Memories

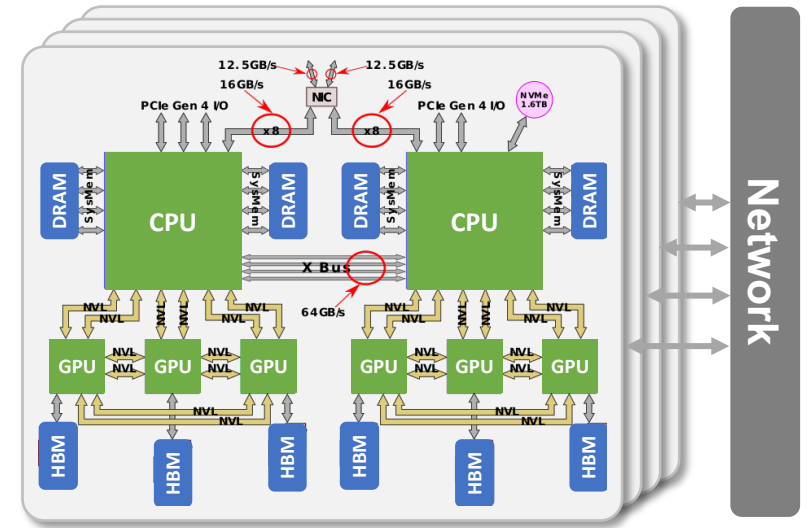
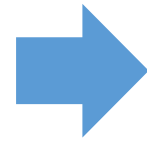
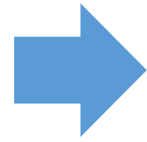


High Performance Network

Distributed Heterogenous Hardware Architectures

~4,600 compute nodes

Challenges of Building ML Systems



Increasingly diverse models

Large Language Models
Transformers,
Vision Language Models,
Graph Neural Networks,
Mixture of Experts,
Sparse NN,
Dynamic NN,
...

Increasingly heterogeneous hardware

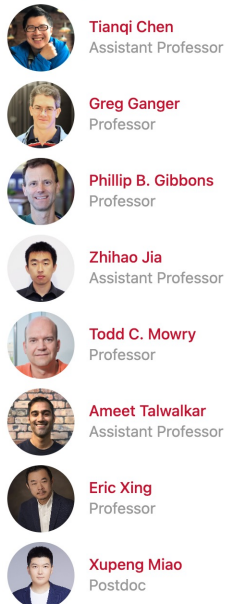
CPUs, GPUs, TPUs,
AI accelerators,
FPGAs, CGRAs,
Programmable networks,
and their combinations
...

CMU Automated Learning Systems Lab

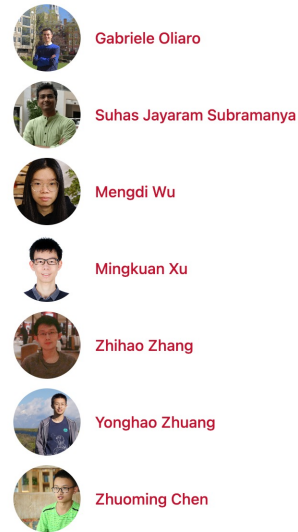
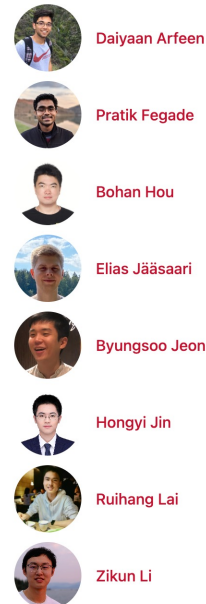
Mission: Automate the design and optimization of ML systems by leveraging

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

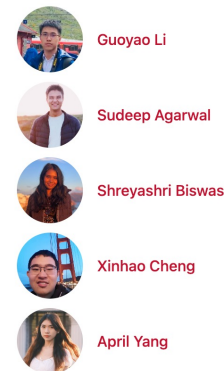
Faculty



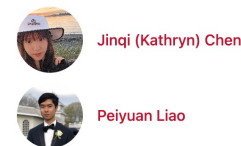
PhD Students



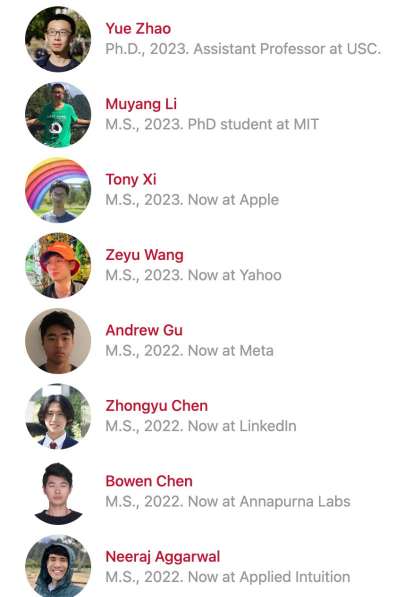
Masters Students



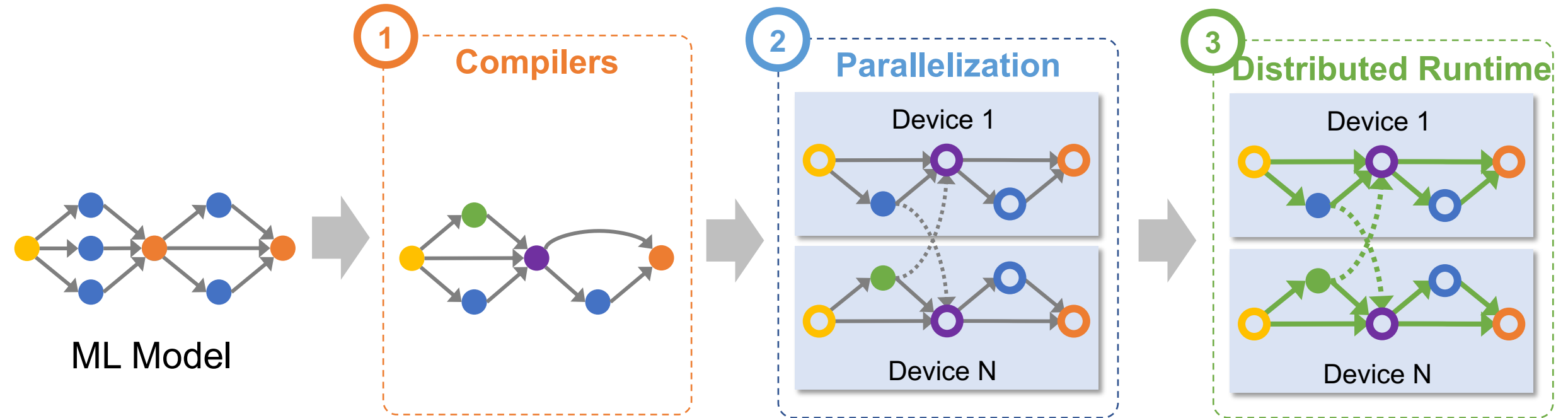
Undergraduate Students



Alumni

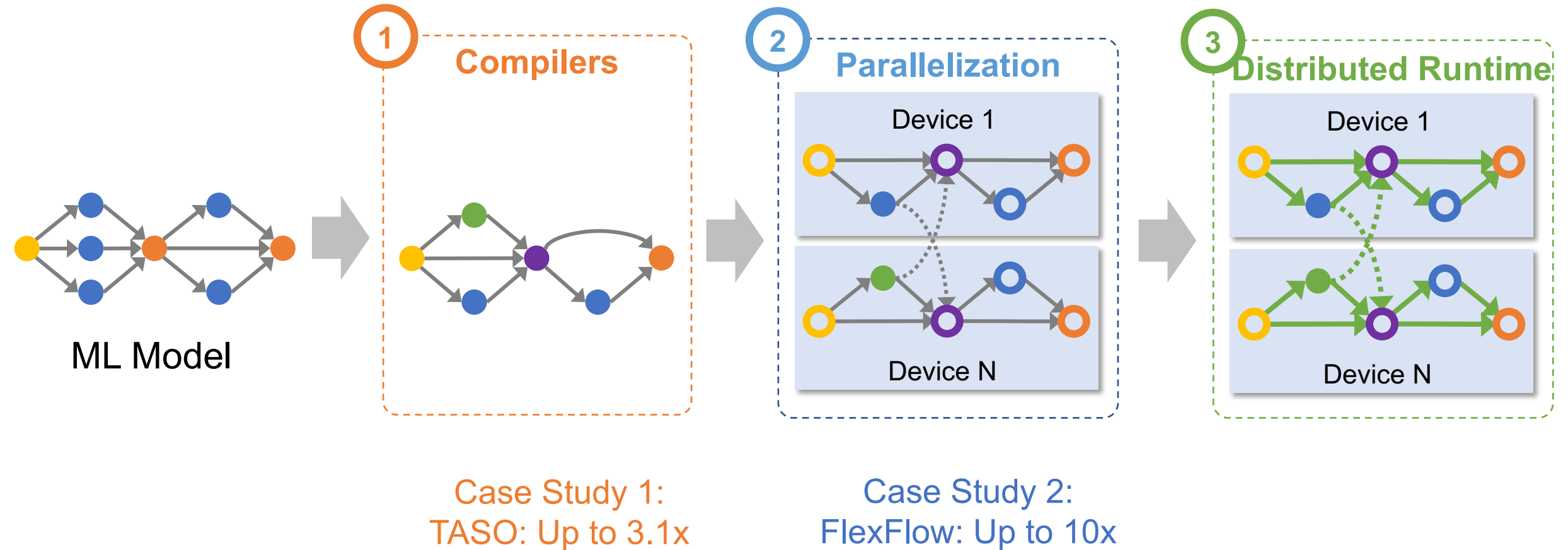


This Lecture: Three Lessons Learned from Our MLSys Research

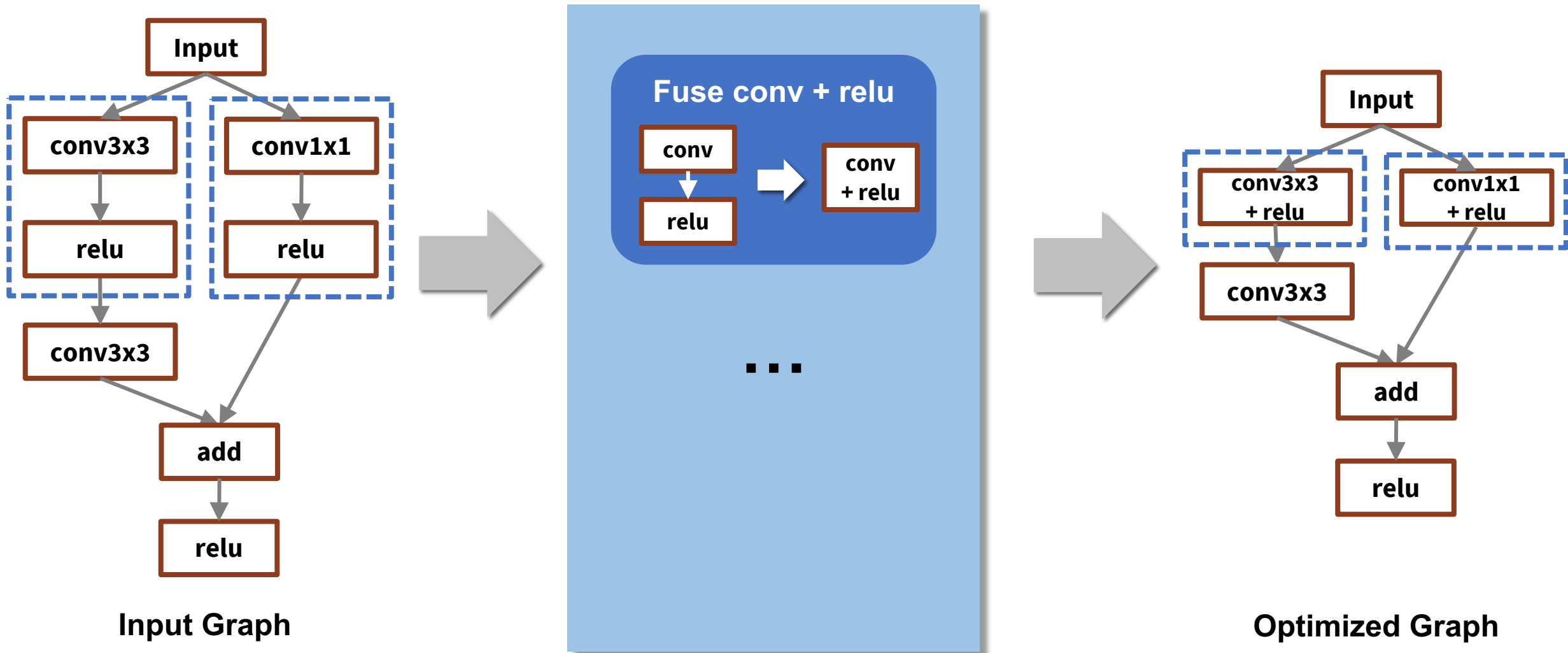


1. Automated approaches can offer **3-10x** improvement for most tasks
2. Joint optimization is **critical**
3. Combing systems and ML optimizations is **promising but challenging**

Lesson 1: Automated Approaches Offer 3-10x Improvement



Case Study: Current Rule-based Graph Optimizations



Rule-based Optimizer

Case Study: Current Rule-based Graph Optimizations

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

Fuse conv + relu

Fuse conv +
batch normalization

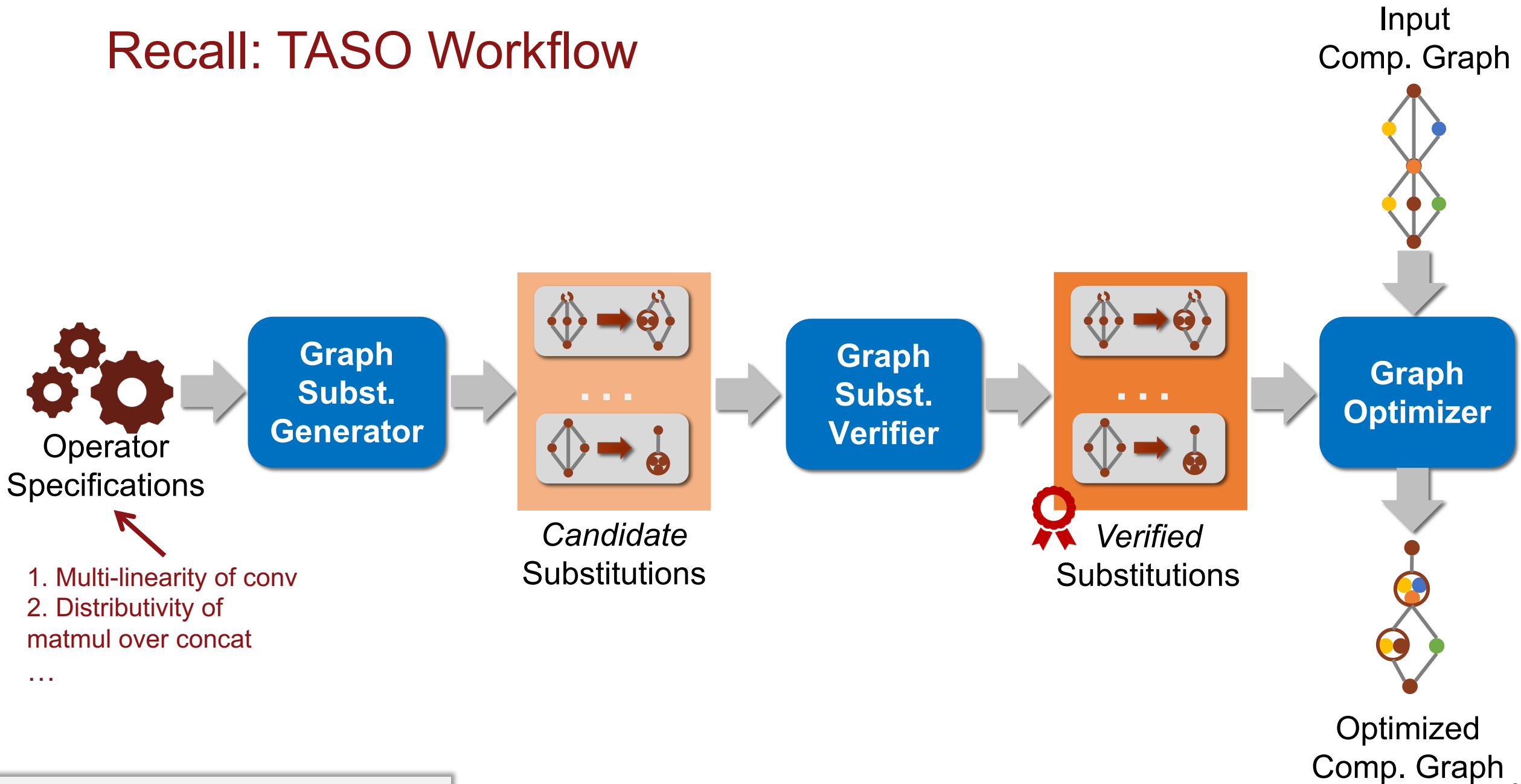
Fuse multi. convs

■ ■ ■

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
124 } // namespace graph_transforms
125 } // namespace tensorflow
```


Recall: TASO Workflow

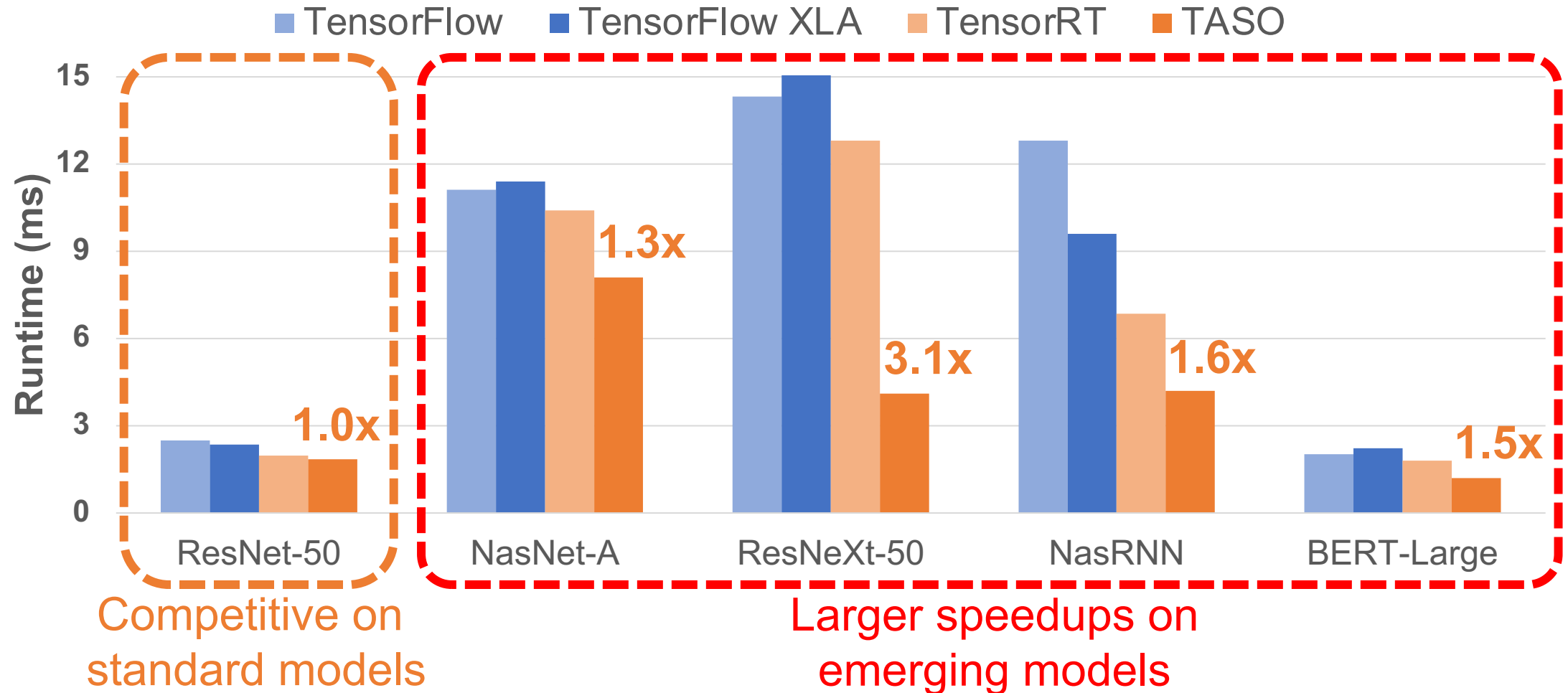


TASO: Tensor Algebra SuperOptimizer

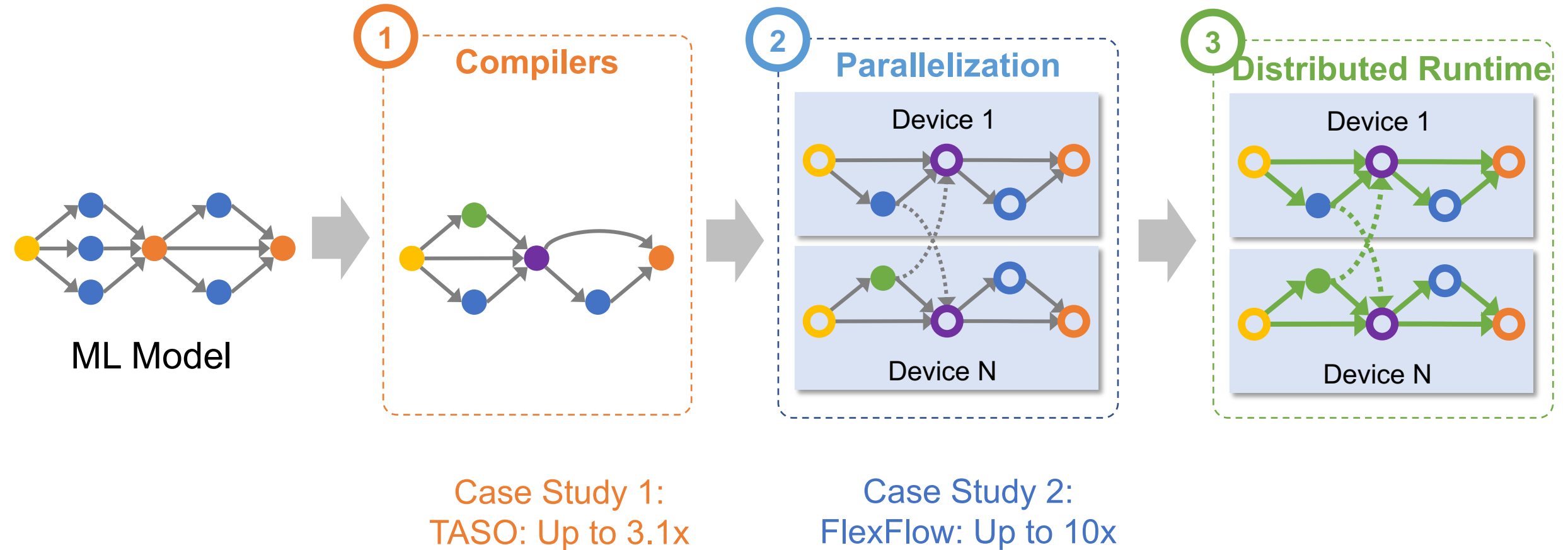
Key idea: replace manually-designed graph optimizations with *automated generation and verification* of graph substitutions for tensor algebra

- **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

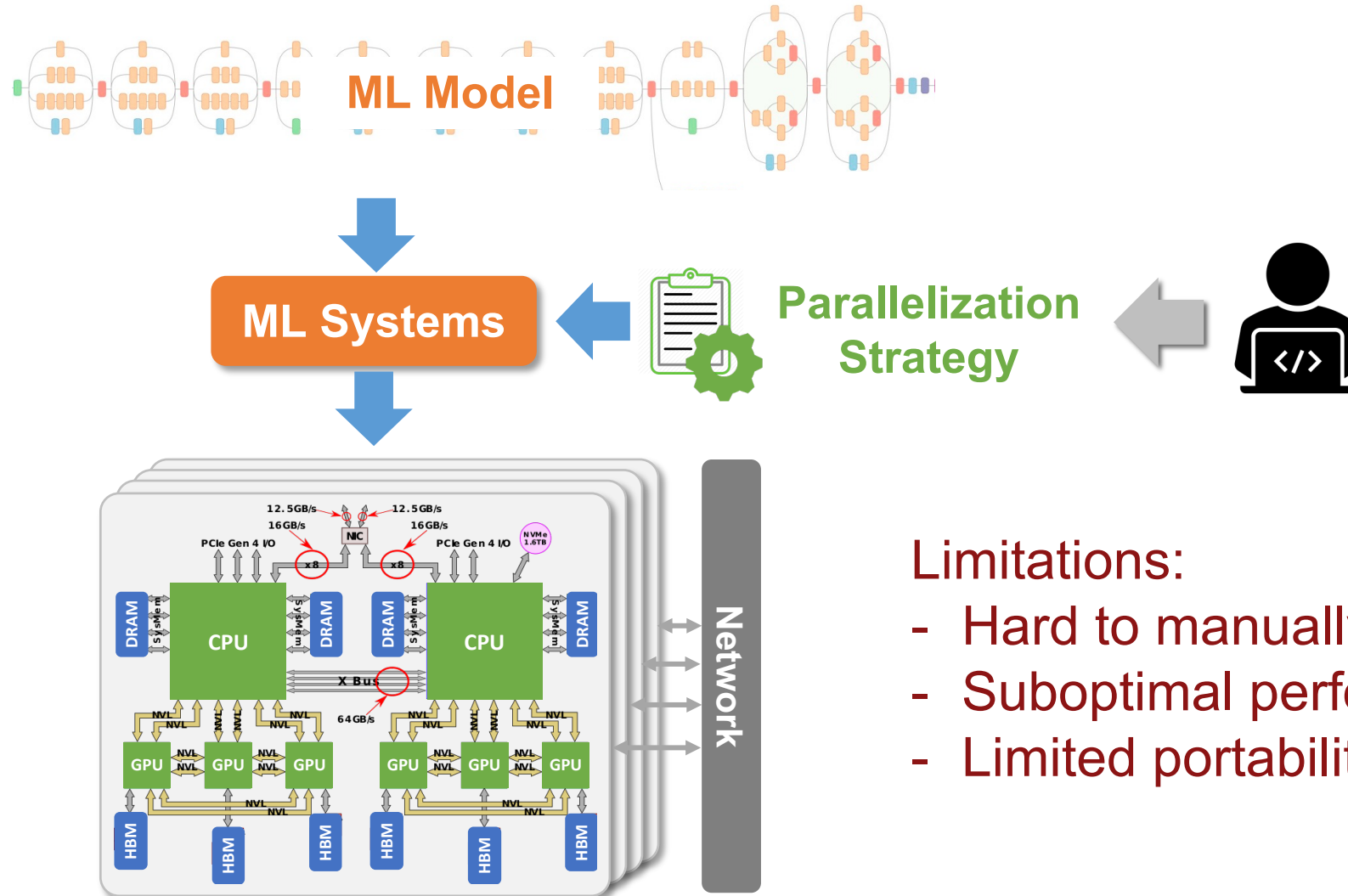
End-to-end Inference Performance (Nvidia V100 GPU)



Lesson 1: Automated Approaches Offer 3-10x Improvement



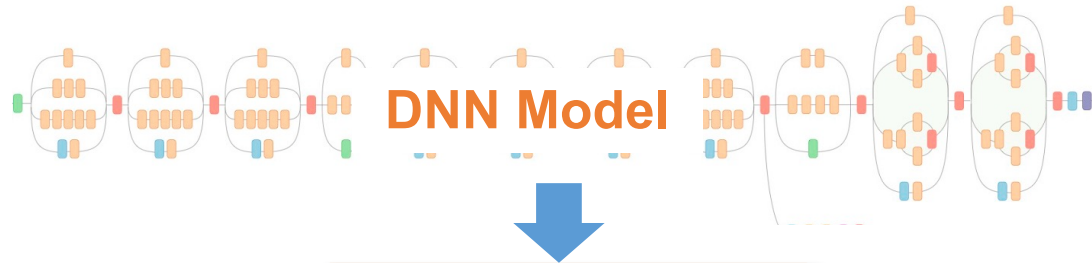
Challenges of Parallelizing DNN Training



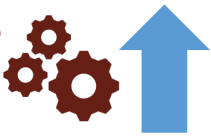
Limitations:

- Hard to manually design and implement
- Suboptimal performance
- Limited portability

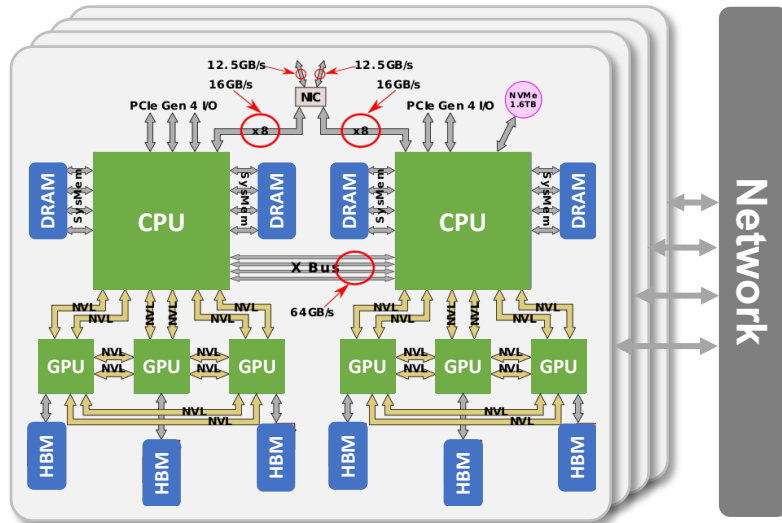
FlexFlow: Automatically Optimizing DNN Parallelization



Hardware
Topology



Fast Parallel
Strategy



Better Performance

Up to 10x faster than manually designed strategies

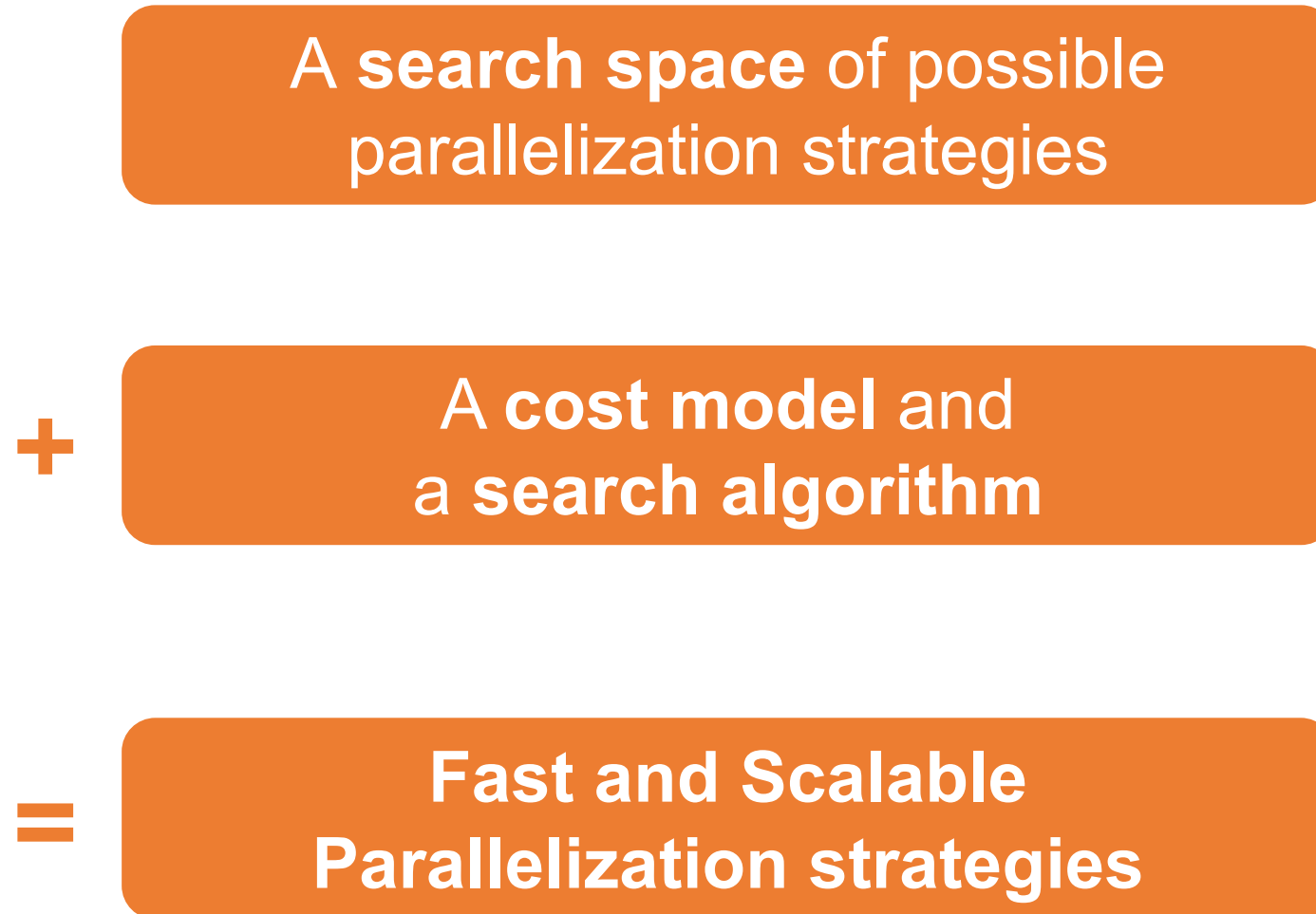
Fast Deployment

Minutes of automated search to discover performant strategies

No Manual Effort

Automatically find strategies for new DNN models or hardware platforms

FlexFlow: Searching for Efficient Parallelization Strategies

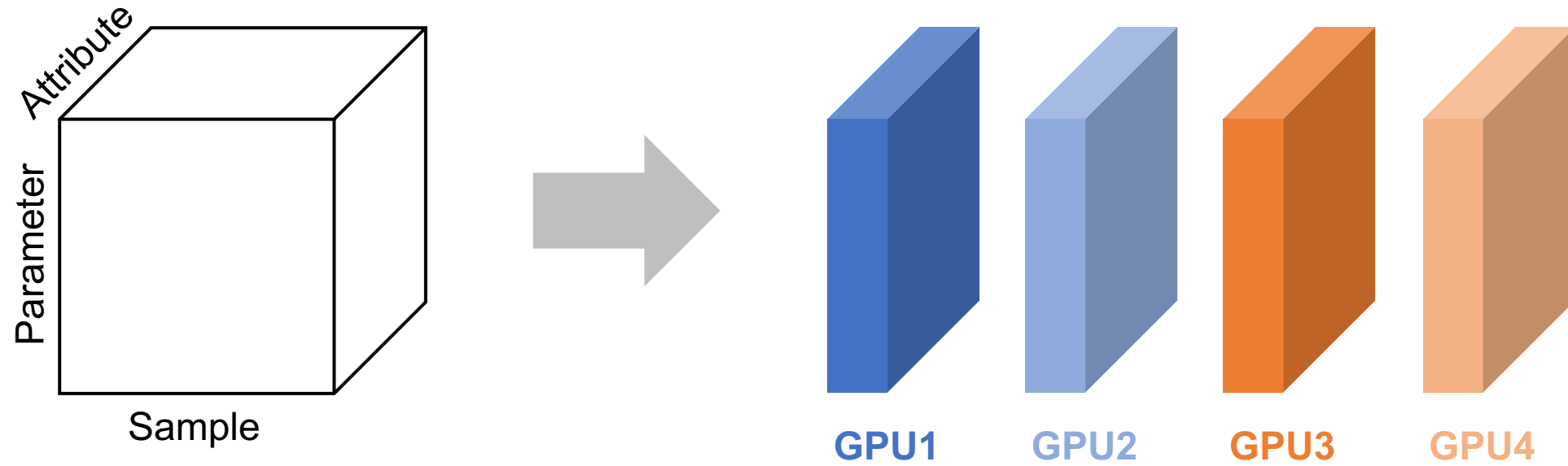


The SOAP Search Space

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

The SOAP Search Space

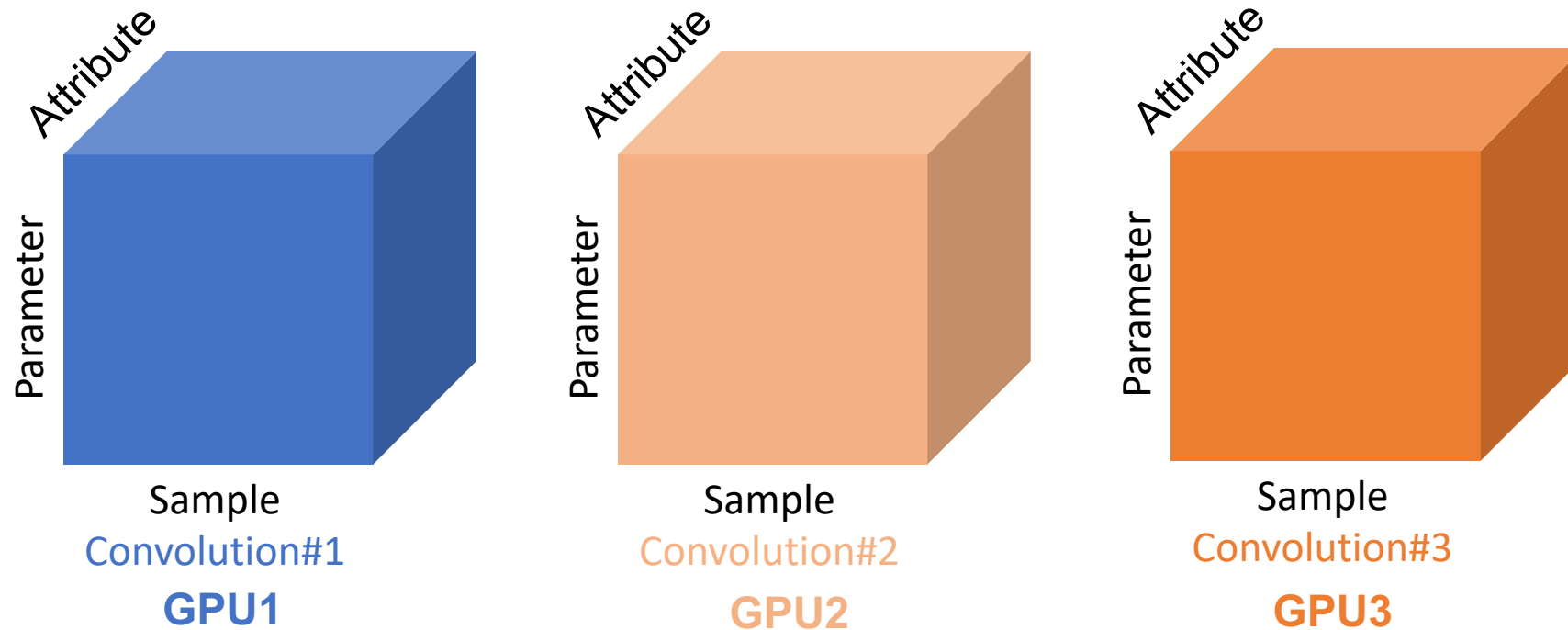
- **S**amples: partitioning training samples (Data Parallelism)
- **O**perators
- **A**tributes
- **P**arameters



Parallelizing a 1D convolution in *Sample*

The SOAP Search Space

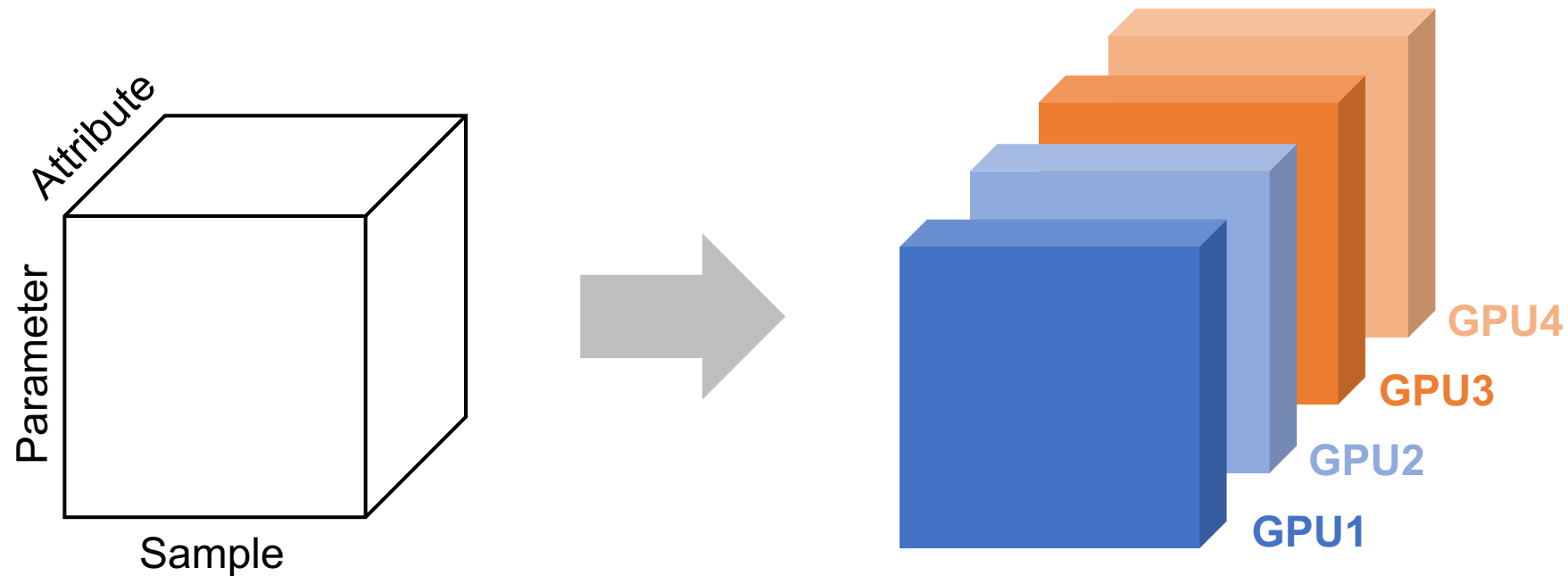
- **S**amples: partitioning training samples (Data Parallelism)
- **O**perators: partitioning ML operators (Model Parallelism)
- **A**tributes
- **P**arameters



Parallelizing multiple convolutions in *Operator*

The SOAP Search Space

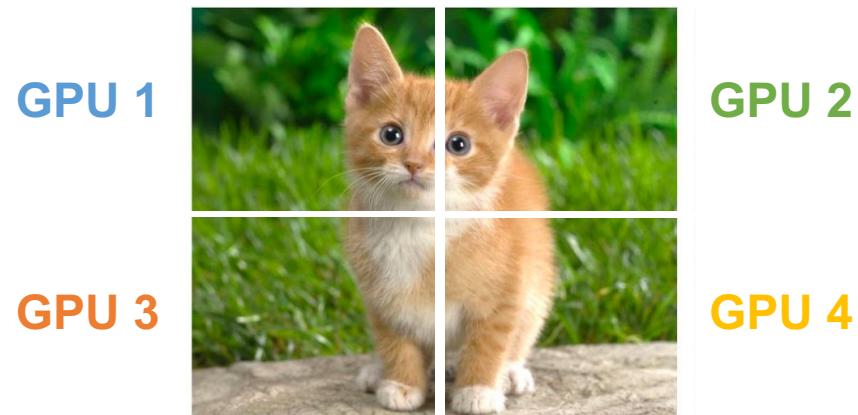
- **S**amples: partitioning training samples (Data Parallelism)
- **O**perators: partitioning ML operators (Model Parallelism)
- **A**tributes: partitioning attributes in a sample (e.g., pixels)
- **P**arameters



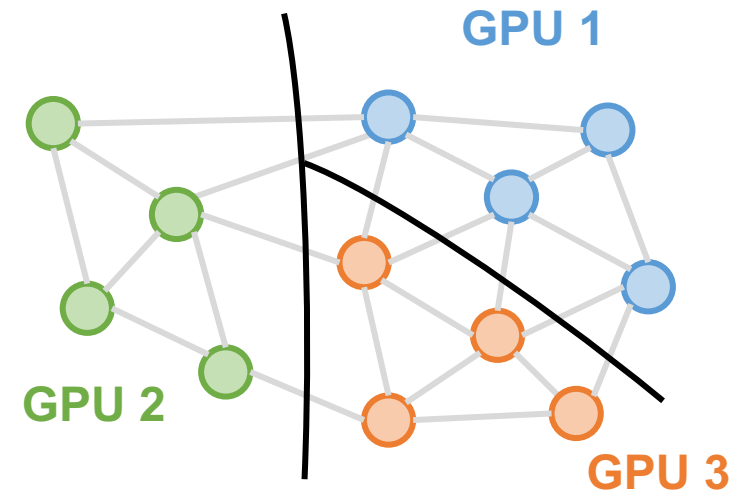
Parallelizing a 1D convolution in *Attribute*

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



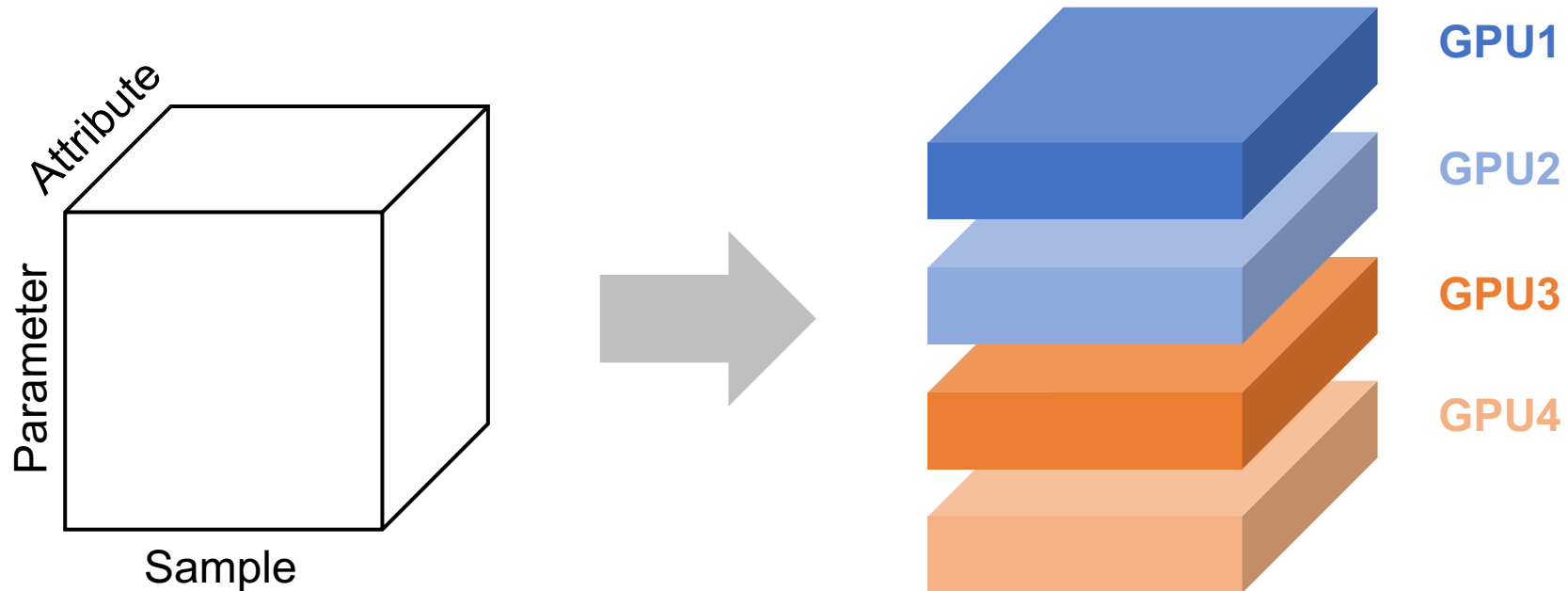
Attribute parallelism for CNNs
(attribute = pixel)



Attribute parallelism for GNNs
(attribute = vertex)

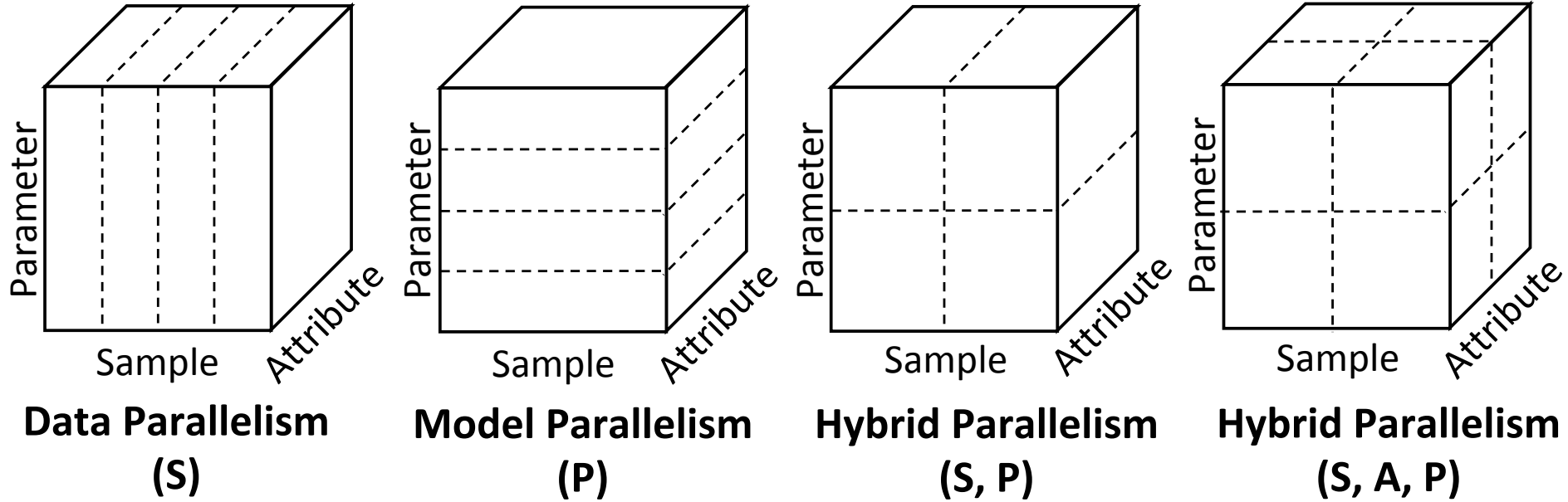
The SOAP Search Space

- **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 (Tensor Model Parallelism)



Parallelizing a 1D convolution in *Parameter*

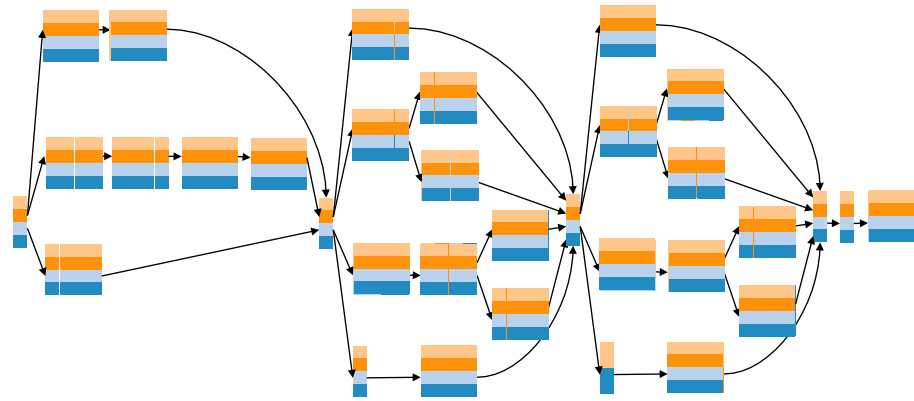
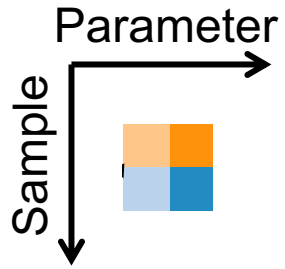
Hybrid Parallelism in SOAP



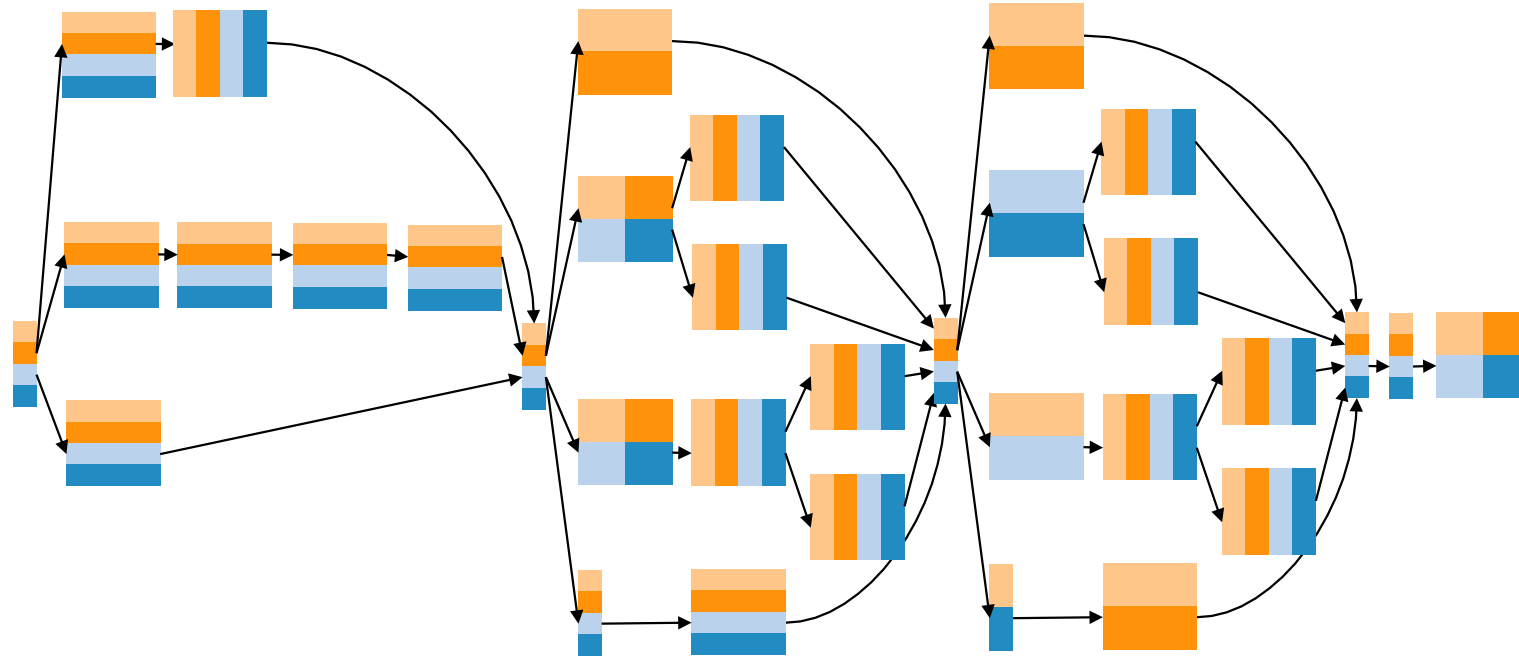
Example parallelization strategies for 1D convolution

Different strategies perform the same computation.

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

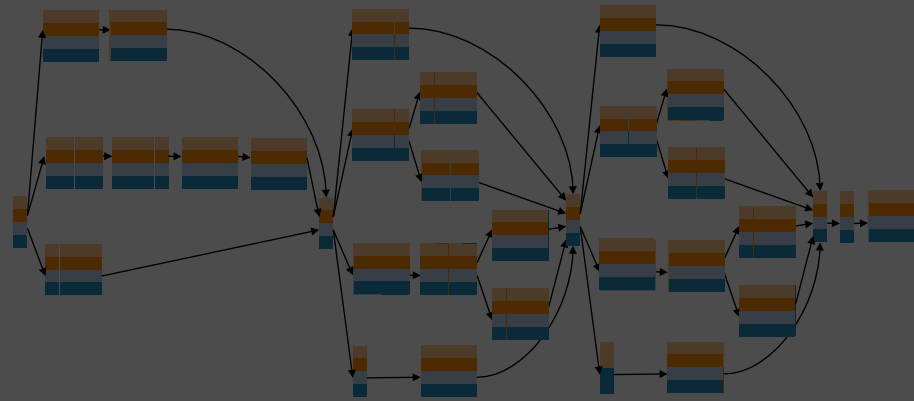
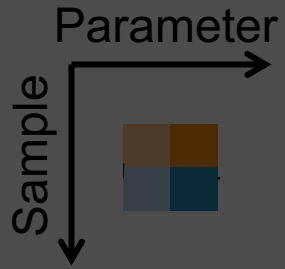


Data parallelism

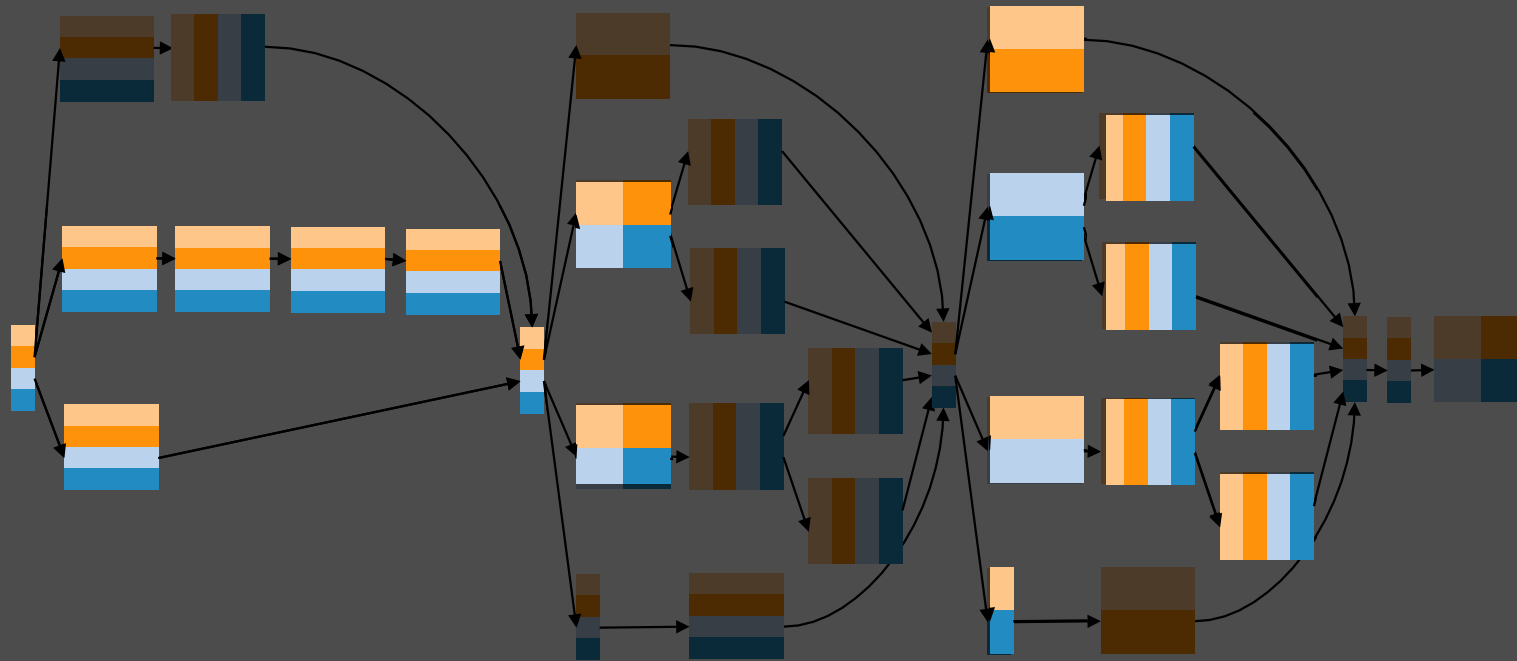


A parallelization strategy in SOAP **(1.2x faster)**

- 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

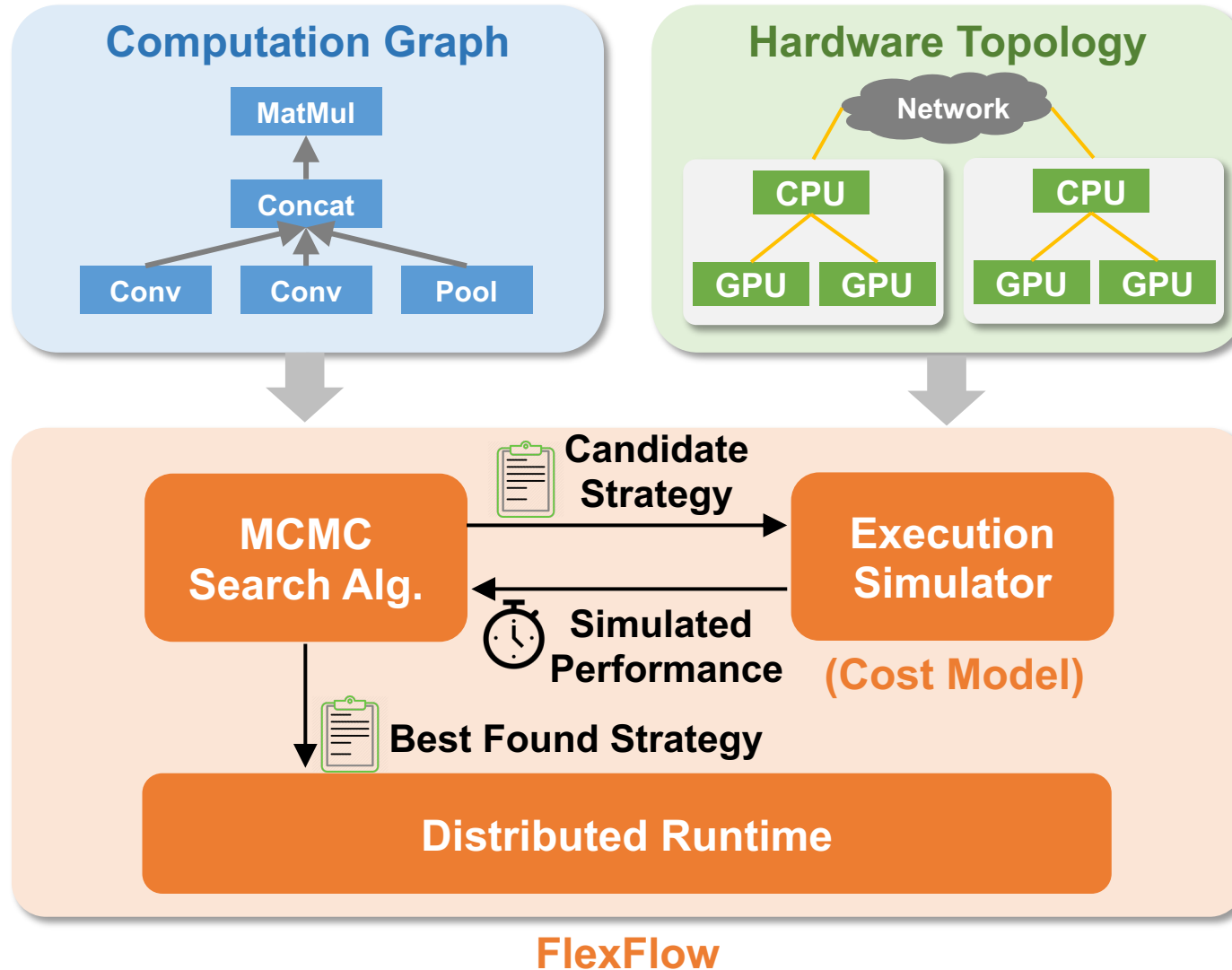
1. SOAP contains billions or more possible strategies

MCMC search algorithm

2. Evaluating a strategy on hardware is too slow

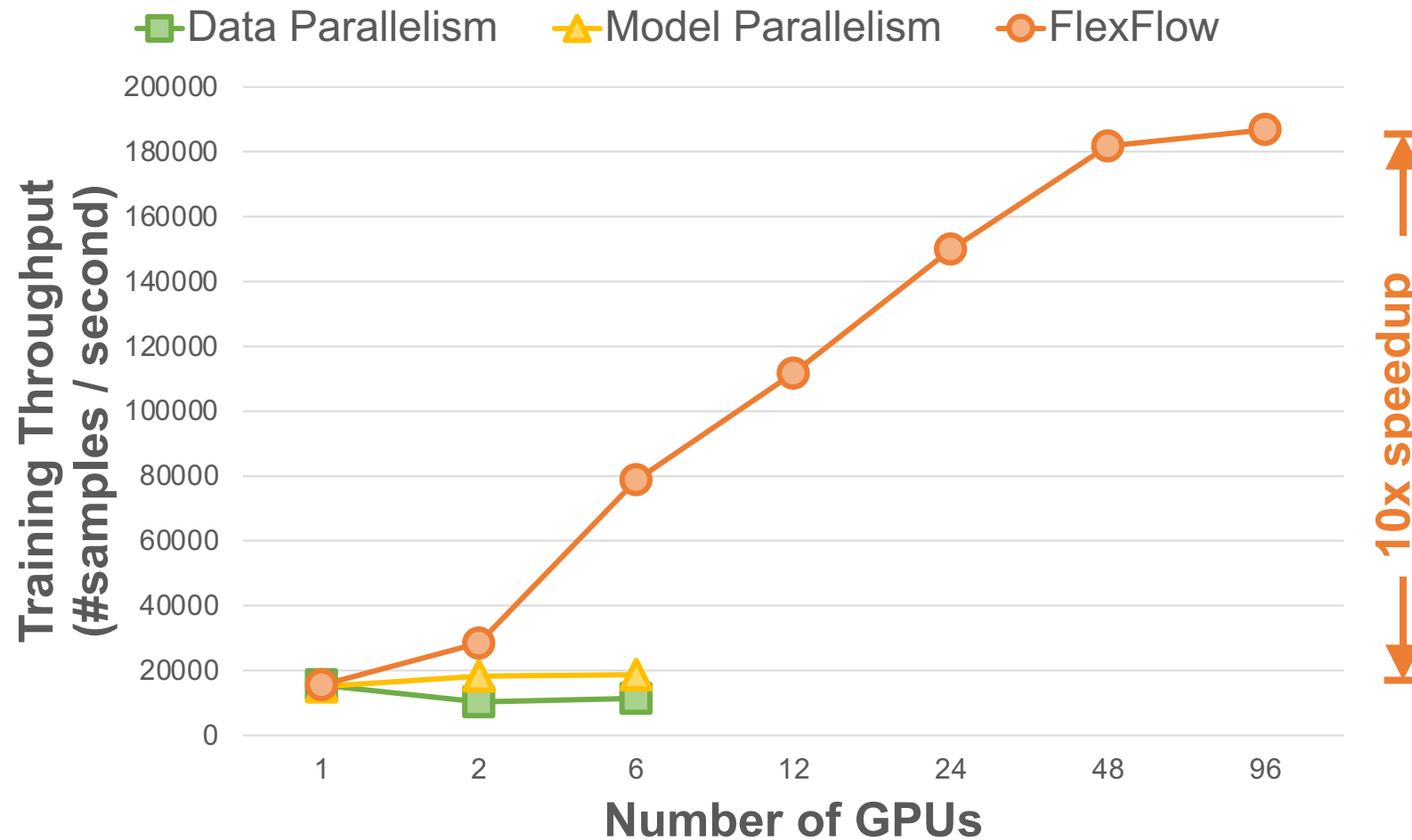
Execution simulator

FlexFlow Overview

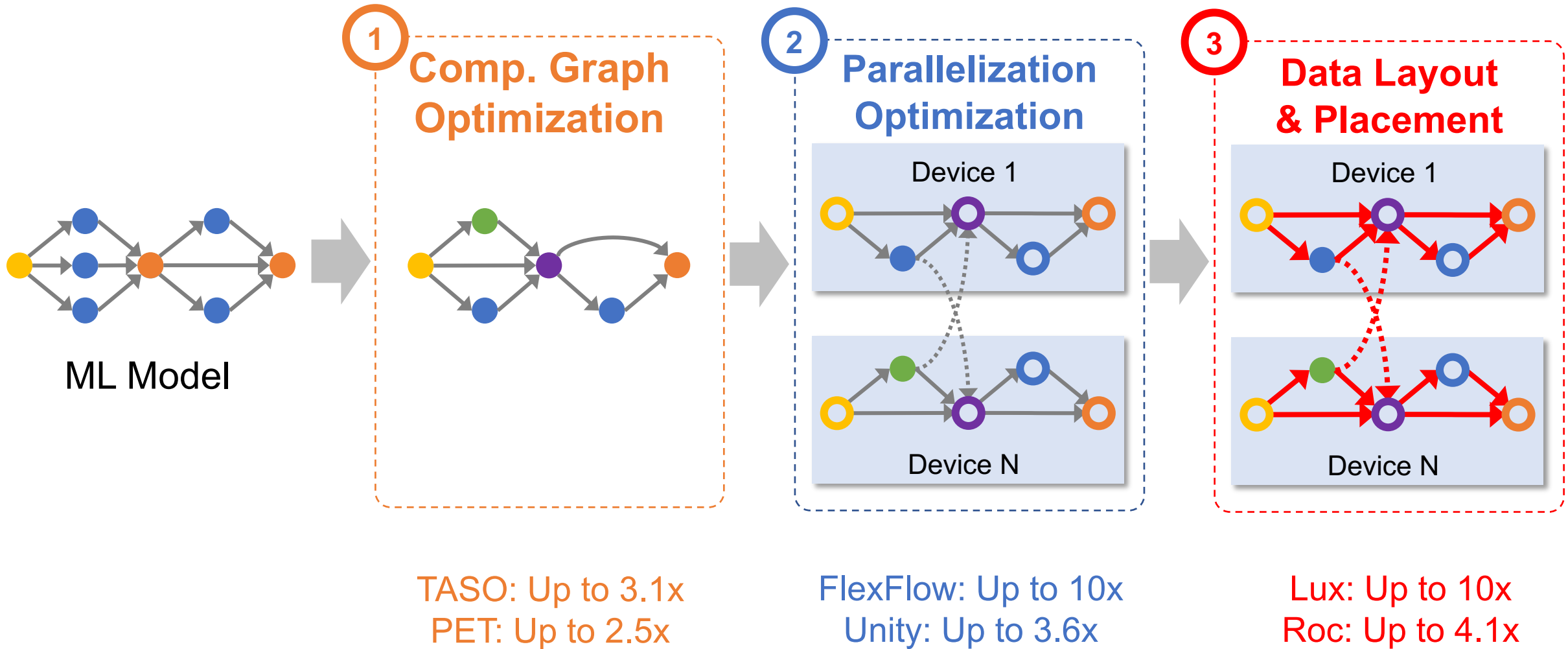


Deep Learning Recommendation Model (DLRM)

A deep learning model for ads recommendation



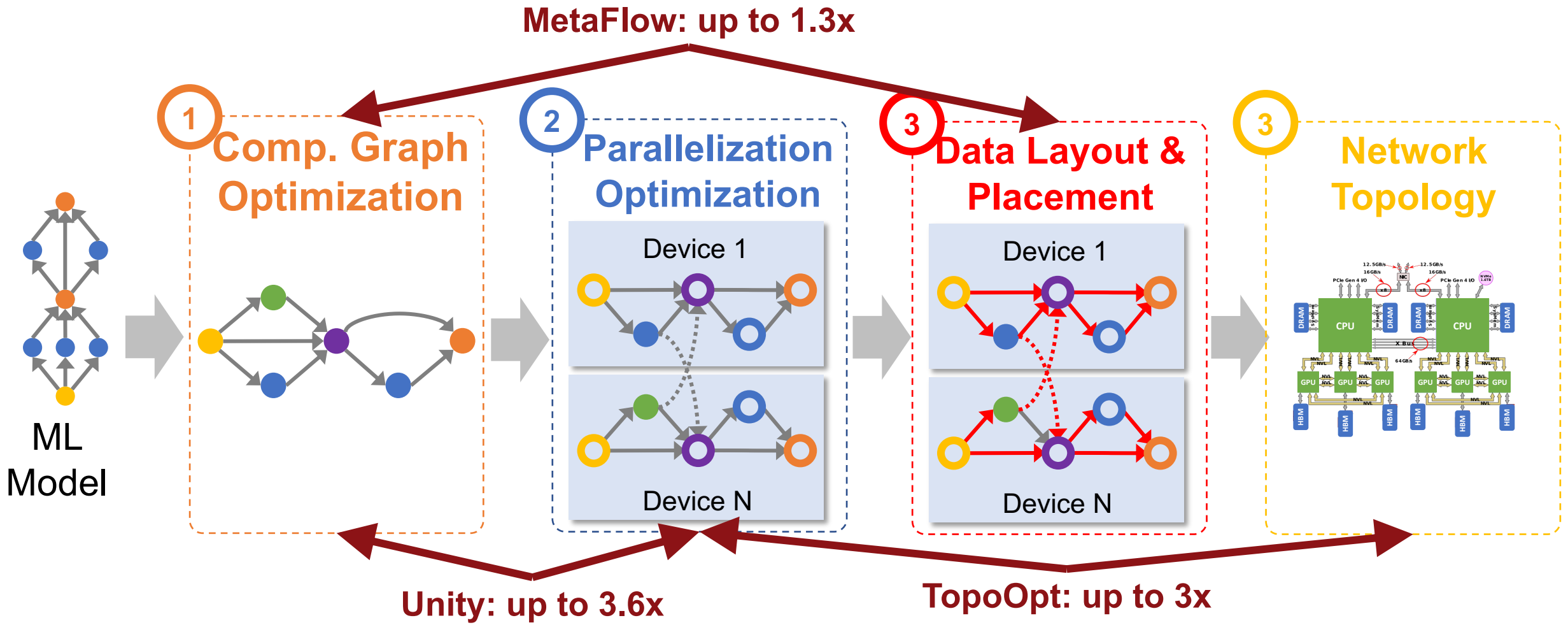
Lesson 1: Automated Approaches Offer 3-10x Improvement



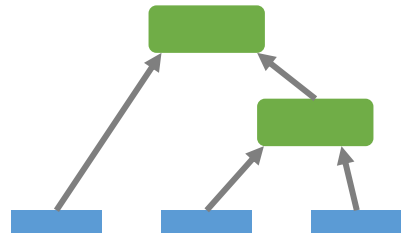
Common Advantages of Automated Approaches

- **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

Lesson 2: Joint Optimization is Critical to Performance



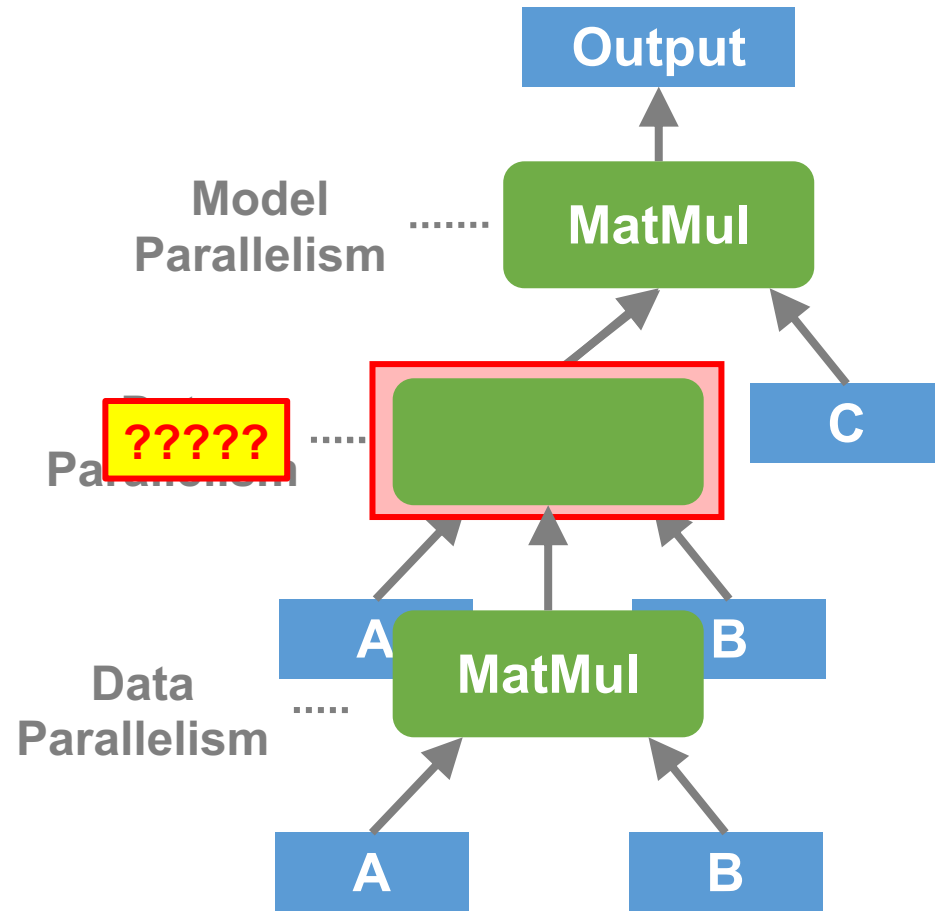
1. Unity: Accelerating DNN Training Through Joint Optimization of Algebraic Transformations and Parallelization. OSDI'22.
2. TopoOpt: Optimizing the Network Topology for Distributed DNN Training. NSDI'23.
3. MetaFlow: Optimizing DNN Computation with Relaxed Graph Substitutions. MLSys'19

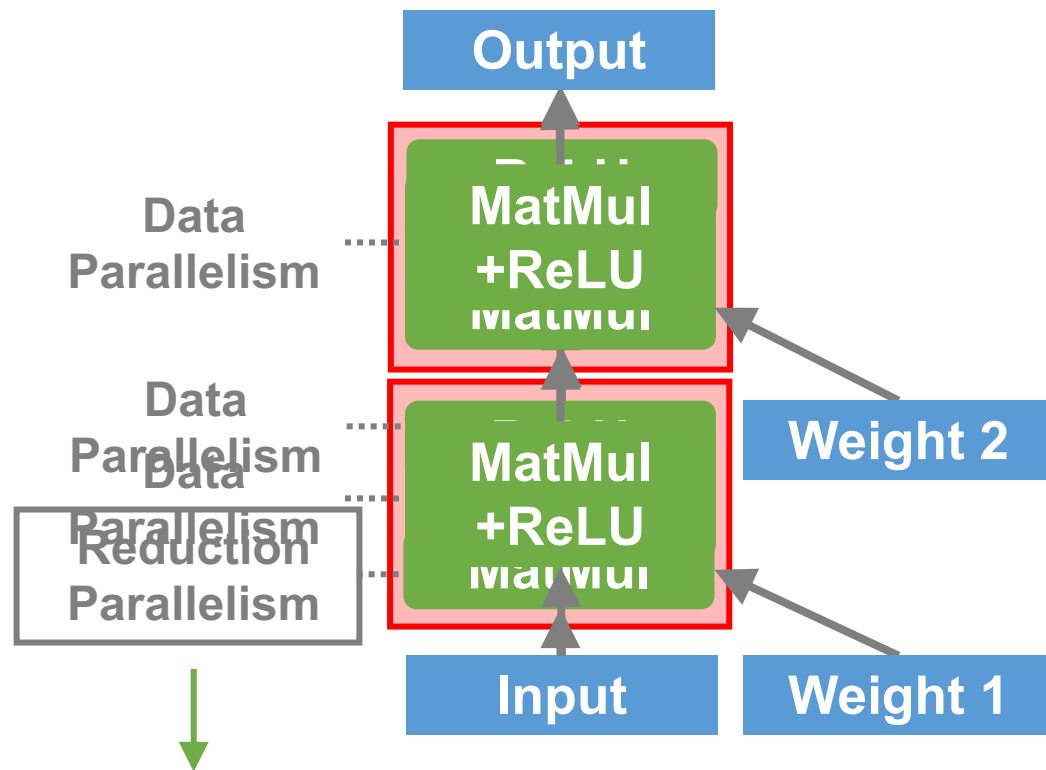


Auto-Parallelization

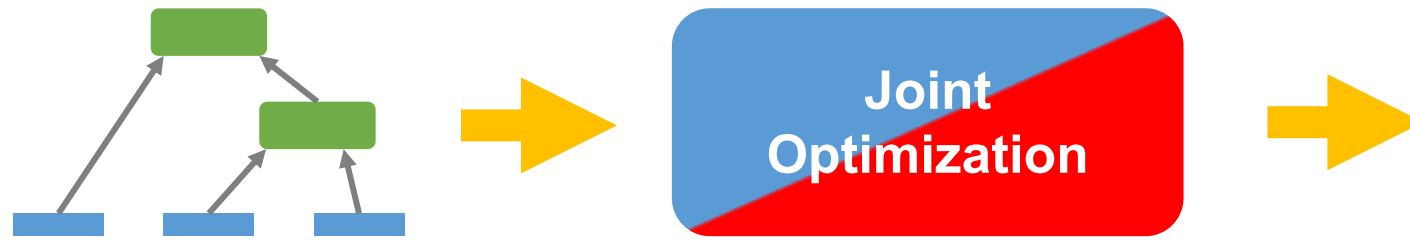
Graph Optimization







≈ 6× less communication!



1. Representation
2. Scalability

Unity

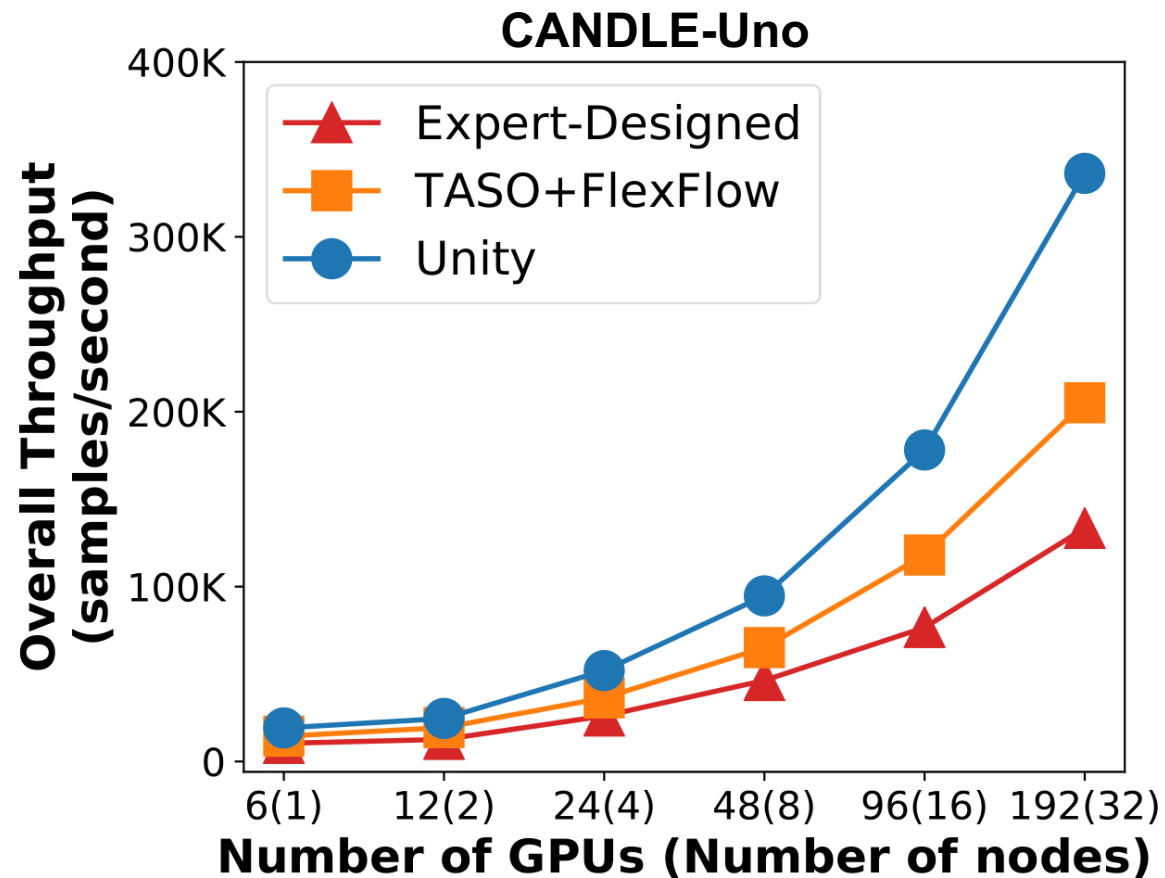
Representation

Parallel Computation
Graph (PCG)

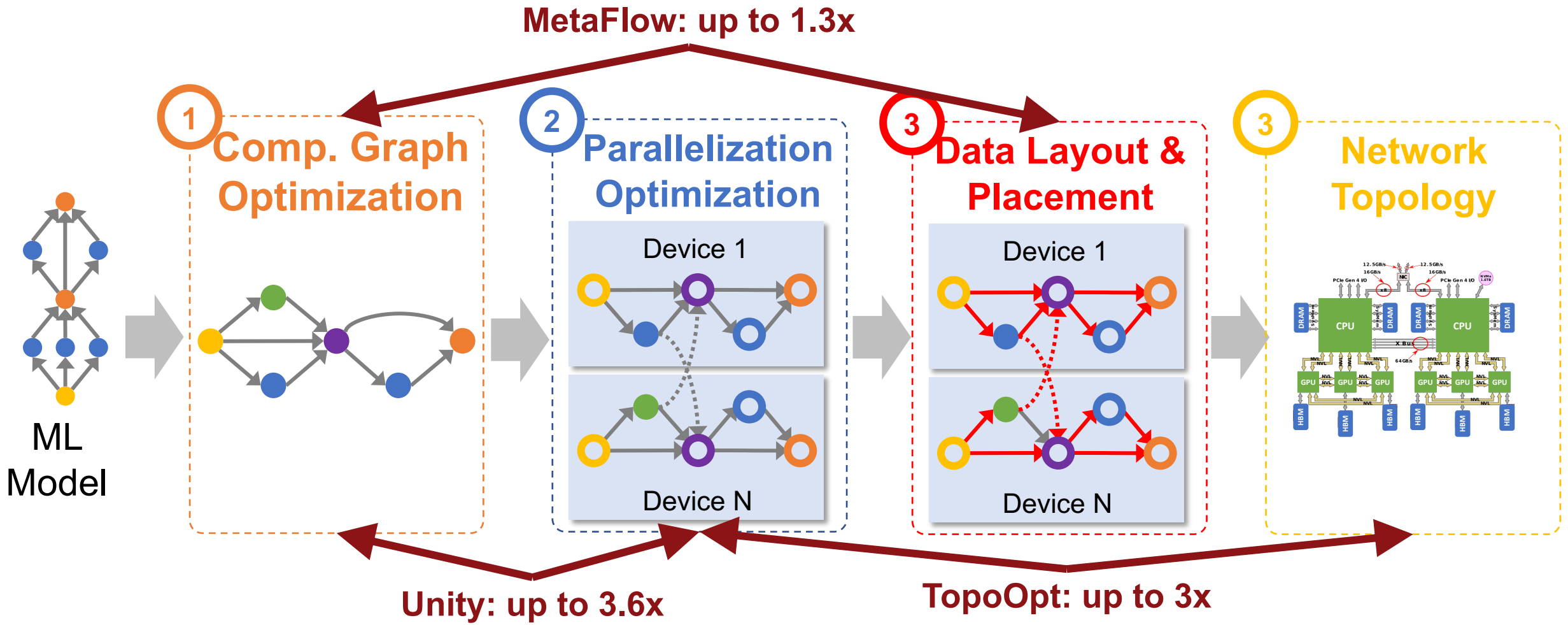
Scalability

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
- Low-Rank Adaptation
- Distillation
- Neural Architecture Search

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

Systems
Optimizations

👍 Pro: preserve equivalence

ML
Optimizations

👍 Pro: better performance

- Faster ML operators
- Less Computation

Achieve the best of both worlds?

Equivalent
Optimizations



Approximate
Optimizations



Three Lessons

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