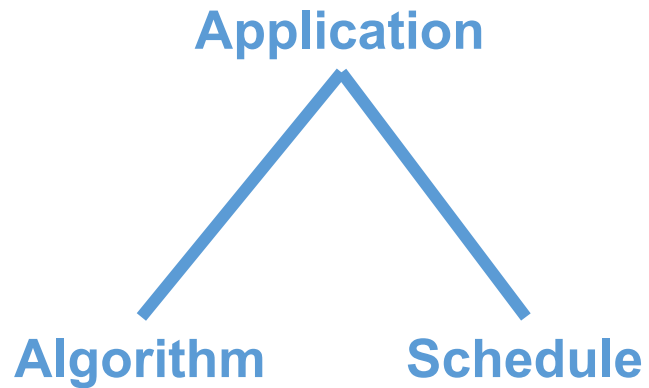


# Comparing Halide, TVM, and Anso



Halide: separate algorithm and schedule, require manual schedules

```
Func halide_blur(Func in) {  
  Func tmp, blurred;  
  Var x, y, xi, yi;  
  
  // The algorithm  
  tmp(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))/3;  
  blurred(x, y) = (tmp(x, y-1) + tmp(x, y) + tmp(x, y+1))/3;  
  
  // The schedule  
  blurred.tile(x, y, xi, yi, 256, 32)  
    .vectorize(xi, 8).parallel(y);  
  tmp.chunk(x).vectorize(x, 8);  
  
  return blurred;  
}
```

Halide: separate algorithm and schedule, require manual schedules

```

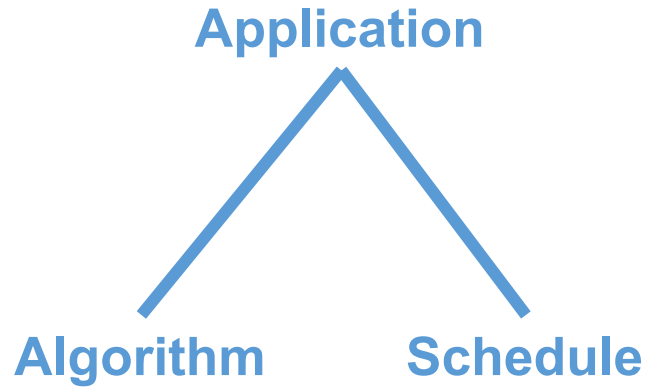
Func halide_blur(Func in) {
  Func tmp, blurred;
  Var x, y, xi, yi;

  // The algorithm
  tmp(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))/3;
  blurred(x, y) = (tmp(x, y-1) + tmp(x, y) + tmp(x, y+1))/3;

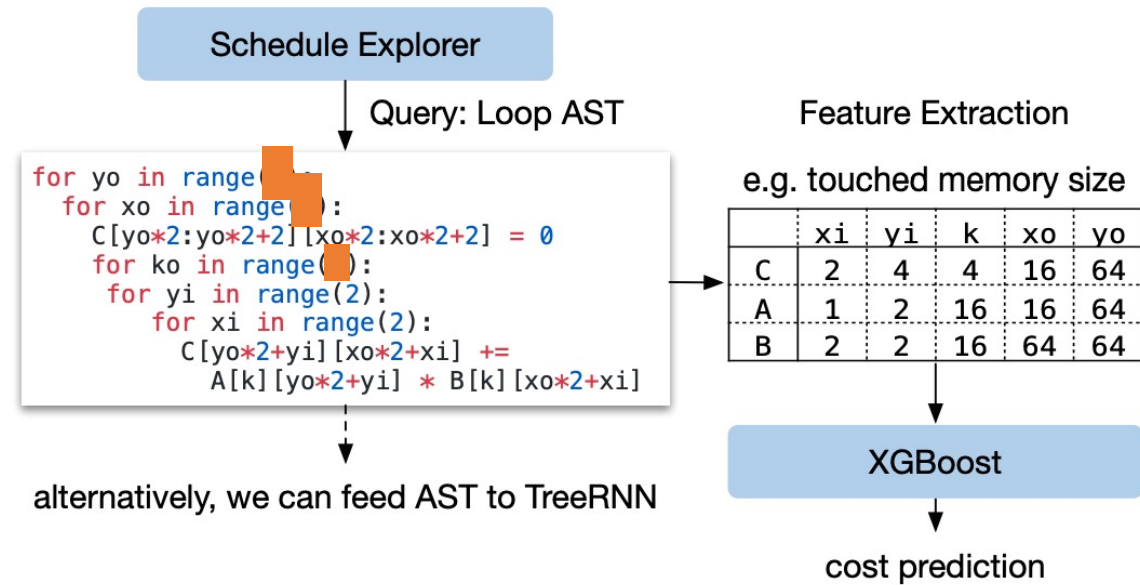
  // The schedule
  blurred.tile(x, y, xi, yi, 256, 32)
    .vectorize(xi, 8).parallel(y);
  tmp.chunk(x).vectorize(x, 8);

  return blurred;
}

```



TVM: require users to specify a schedule space, use ML to explore the space



Halide: separate algorithm and schedule, require manual schedules

```

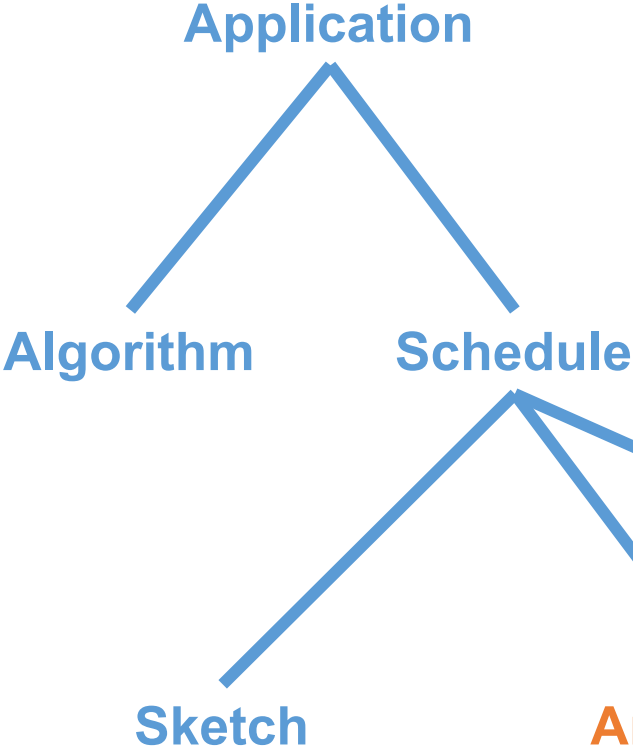
Func halide_blur(Func in) {
  Func tmp, blurred;
  Var x, y, xi, yi;

  // The algorithm
  tmp(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))/3;
  blurred(x, y) = (tmp(x, y-1) + tmp(x, y) + tmp(x, y+1))/3;

  // The schedule
  blurred.tile(x, y, xi, yi, 256, 32)
    .vectorize(xi, 8).parallel(y);
  tmp.chunk(x).vectorize(x, 8);

  return blurred;
}

```



TVM: require users to specify a schedule space, use ML to explore the space

Ansor: generate random sketches and annotations, auto-tune parameters

```

for i in range(8):
  for k in range(512):
    C[i, k] = max(A[i, k], 0.0) if k < 400 else 0
for i in range(8):
  for j in range(4):
    for k_o in range(TILE_K0):
      for k_i in range(TILE_KI):
        E.rf[...] += C[...] * D[...]
for i in range(8):
  for j in range(4):
    for k_i in range(TILE_KI):
      E[...] += E.rf[...]

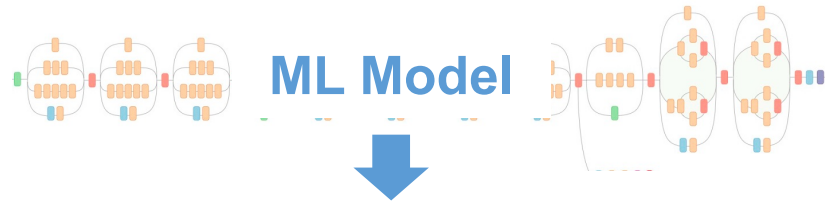
```

```

parallel i in range(8):
  for k in range(512):
    C[i, k] = ...
    for i in range(4):
      unroll k_o in range(32):
        vectorized k_i in range(16):
          E.rf[...] += C[...] * D[...]
parallel i in range(8):
  for i in range(4):
    unroll k_i in range(16):
      E[...] += E.rf[...]

```

# Recap: An Overview of Deep Learning Systems



Automatic Differentiation

Graph-Level Optimization

Parallelization / Distributed Training

Kernel Generation

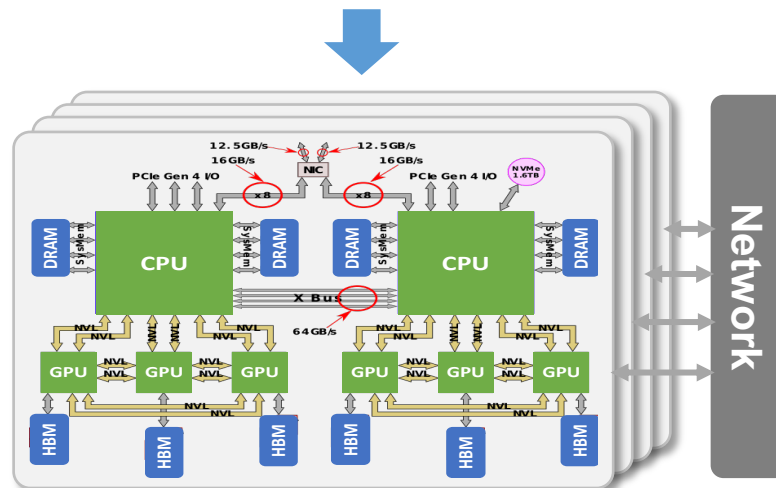
Memory Optimization

Week 5: this lecture

Week 6: Data, Model, Pipeline Parallelism

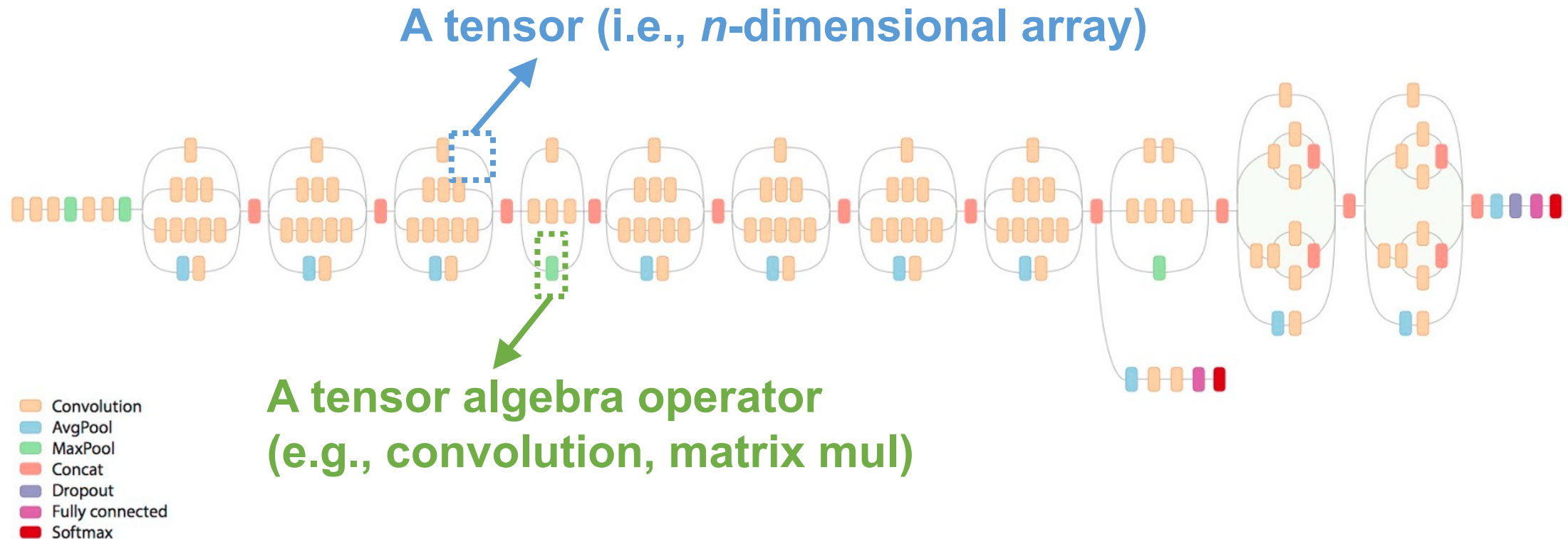
Week 4: Halide, TVM, Anson

Week 7: Zero-Redundancy, Tensor Rematerialization



# Recap: Deep Neural Network

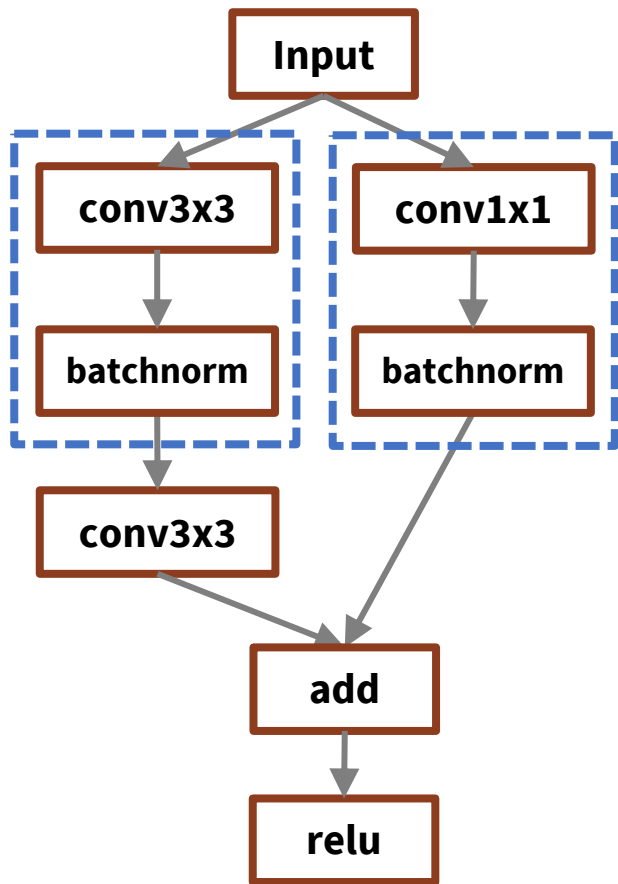
- Collection of simple trainable mathematical units that work together to solve complicated tasks



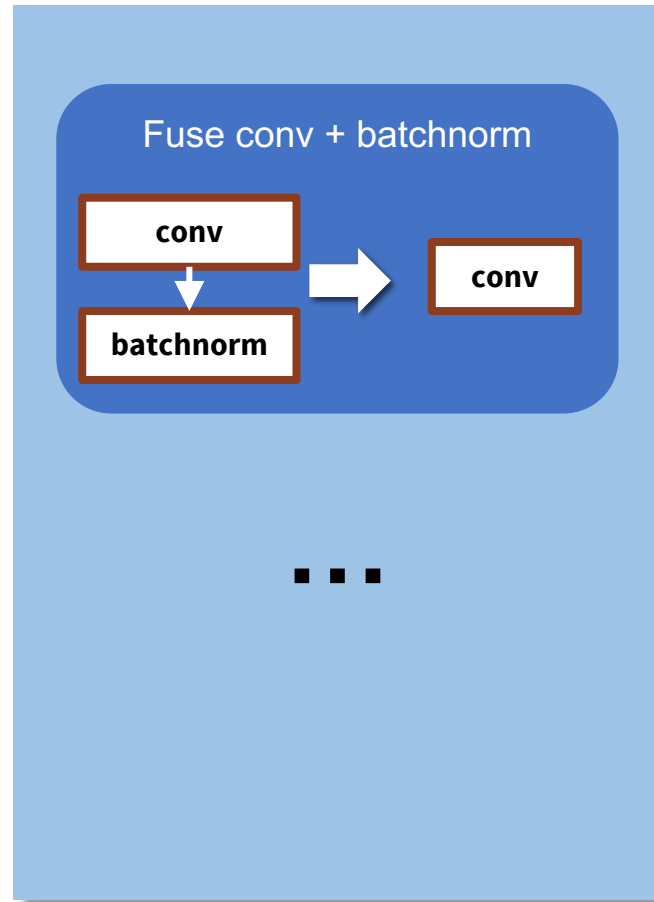
# TASO: Optimizing Deep Learning with Automatic Generation of Graph Substitutions

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

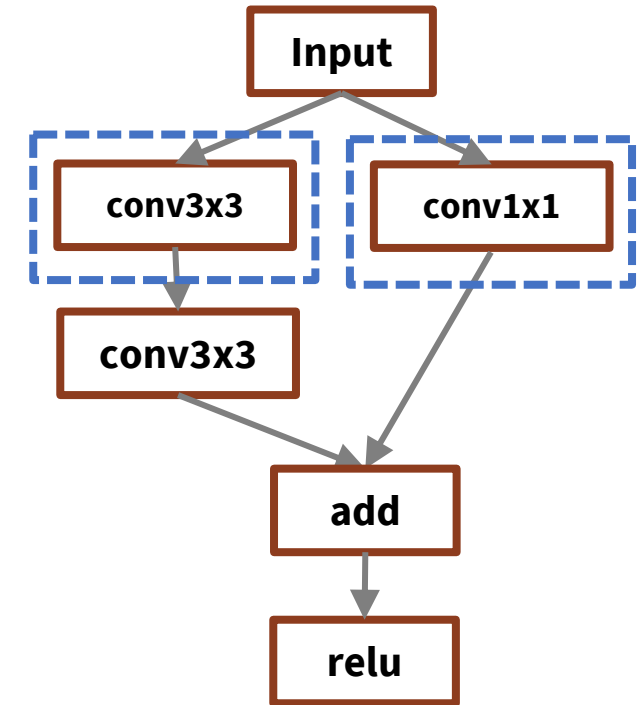
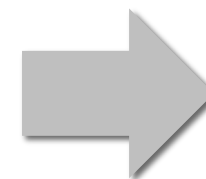
# Graph-Level Optimizations



Input Computation Graph

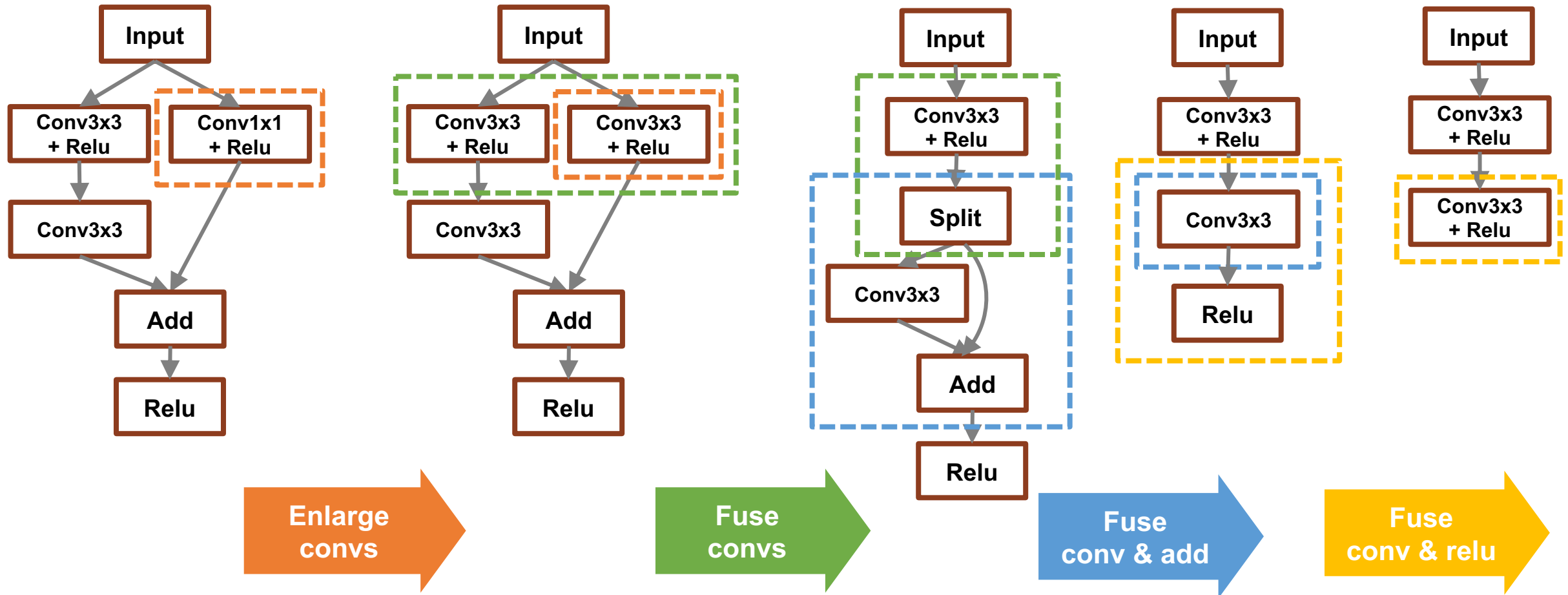


Potential graph transformations



Optimized Computation Graph

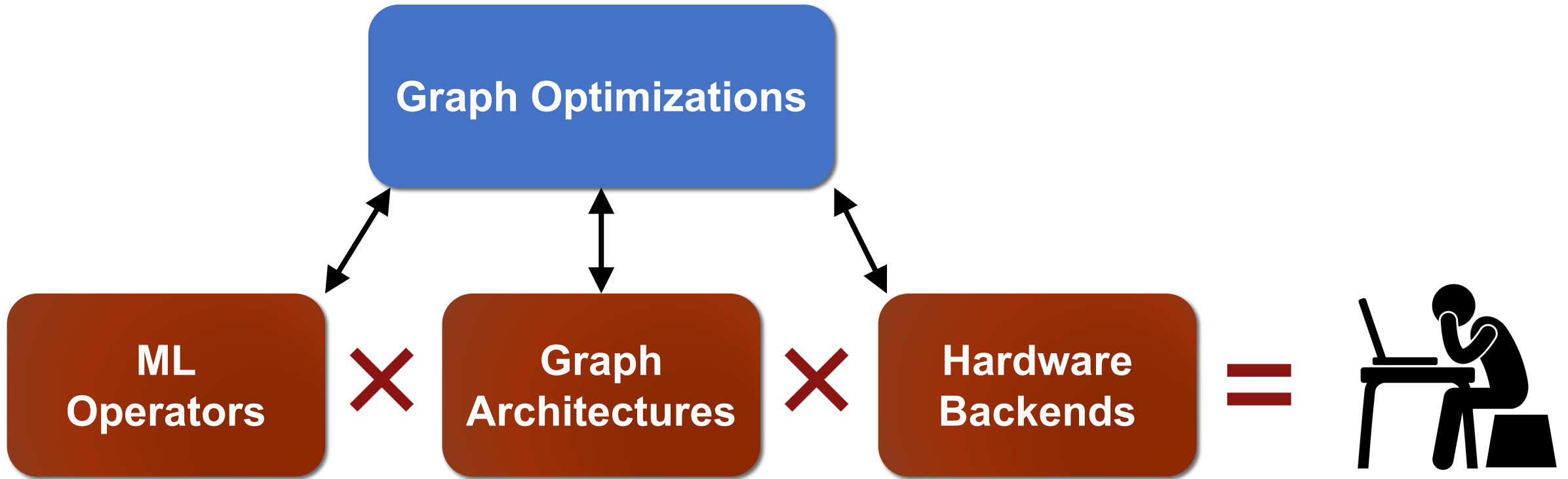
# Recap: ResNet Example



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



# Challenge of Graph Optimizations for ML



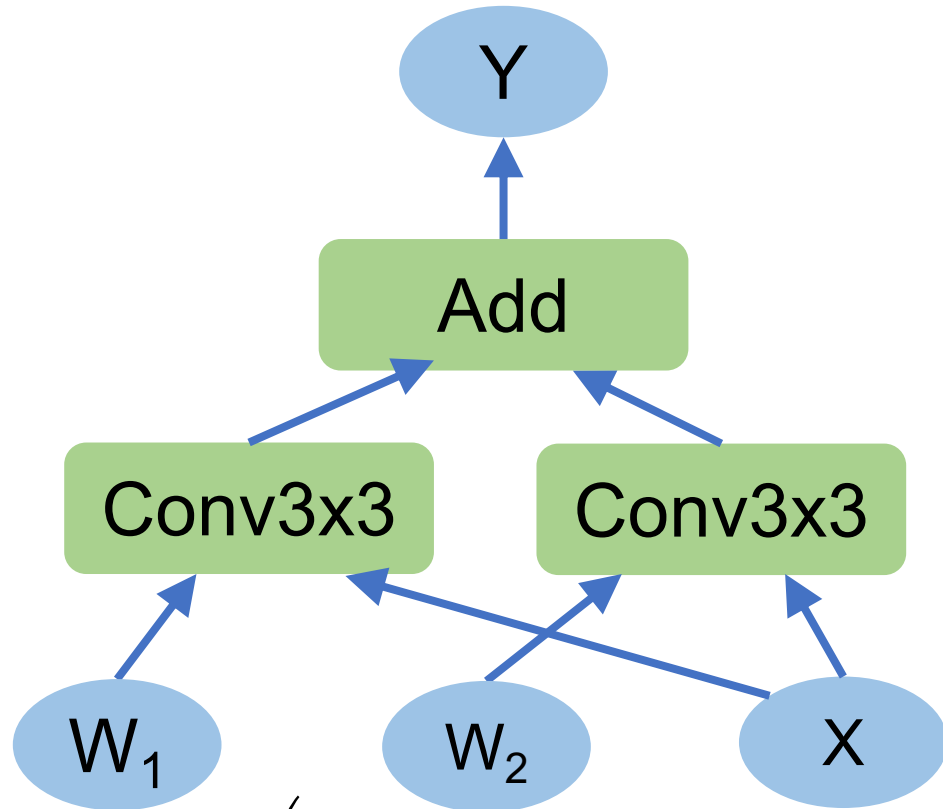
Infeasible to manually design graph optimizations for all cases

# TASO: Tensor Algebra SuperOptimizer

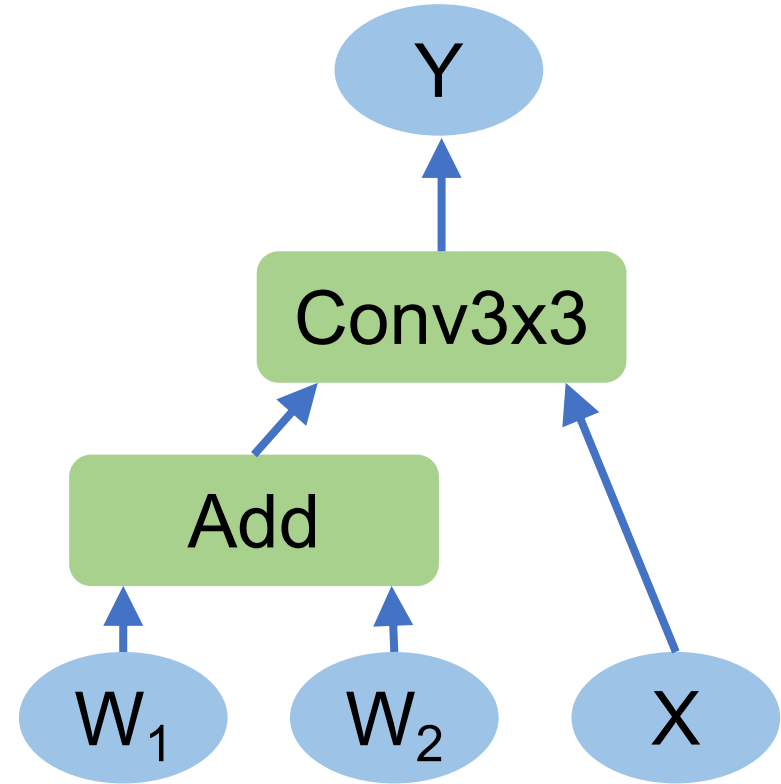
**Key idea:** replace manually-designed graph optimizations with *automated generation and verification* of graph substitutions for tensor algebra

- **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



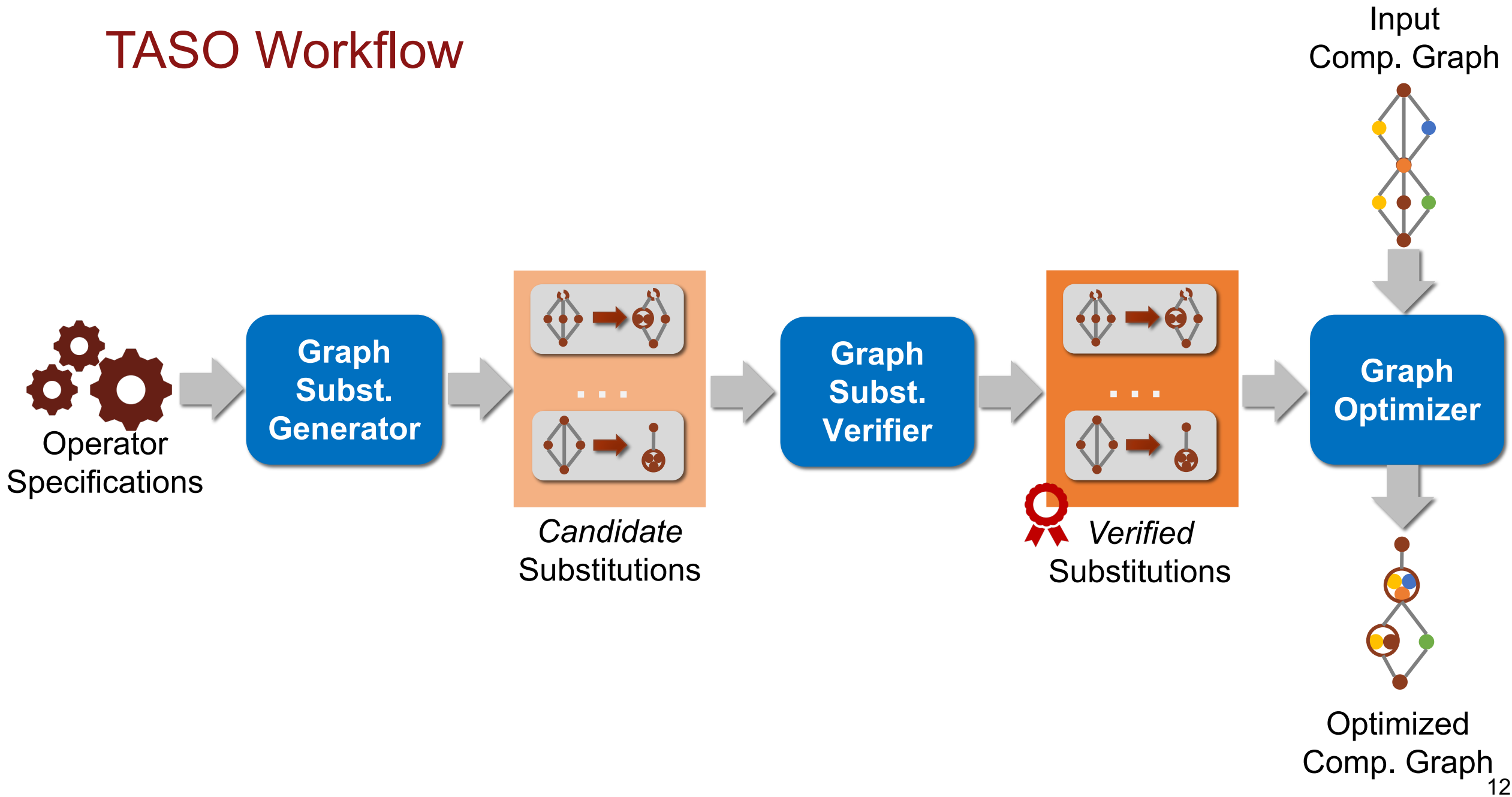
=



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

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

# TASO Workflow



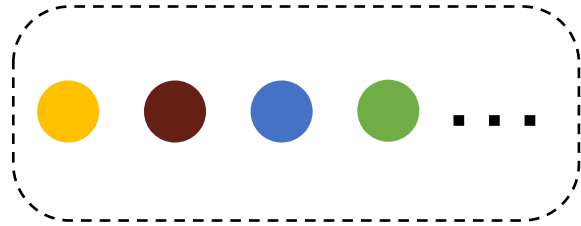
# Graph Substitution Generator

Subst.  
Generator

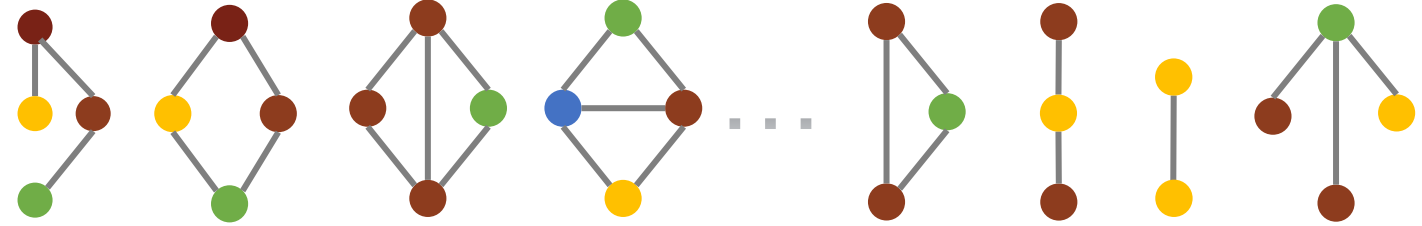
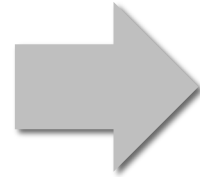
Subst.  
Verifier

Graph  
Optimizer

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



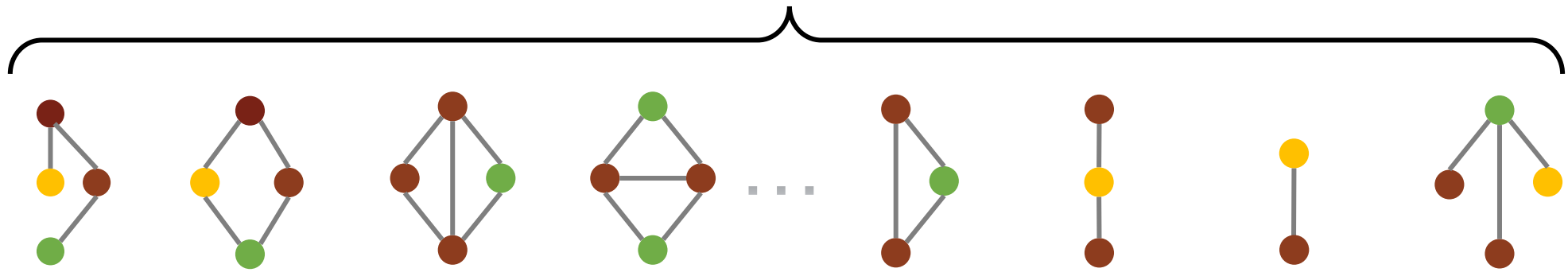
Operators supported by hardware backend





# Graph Substitution Generator

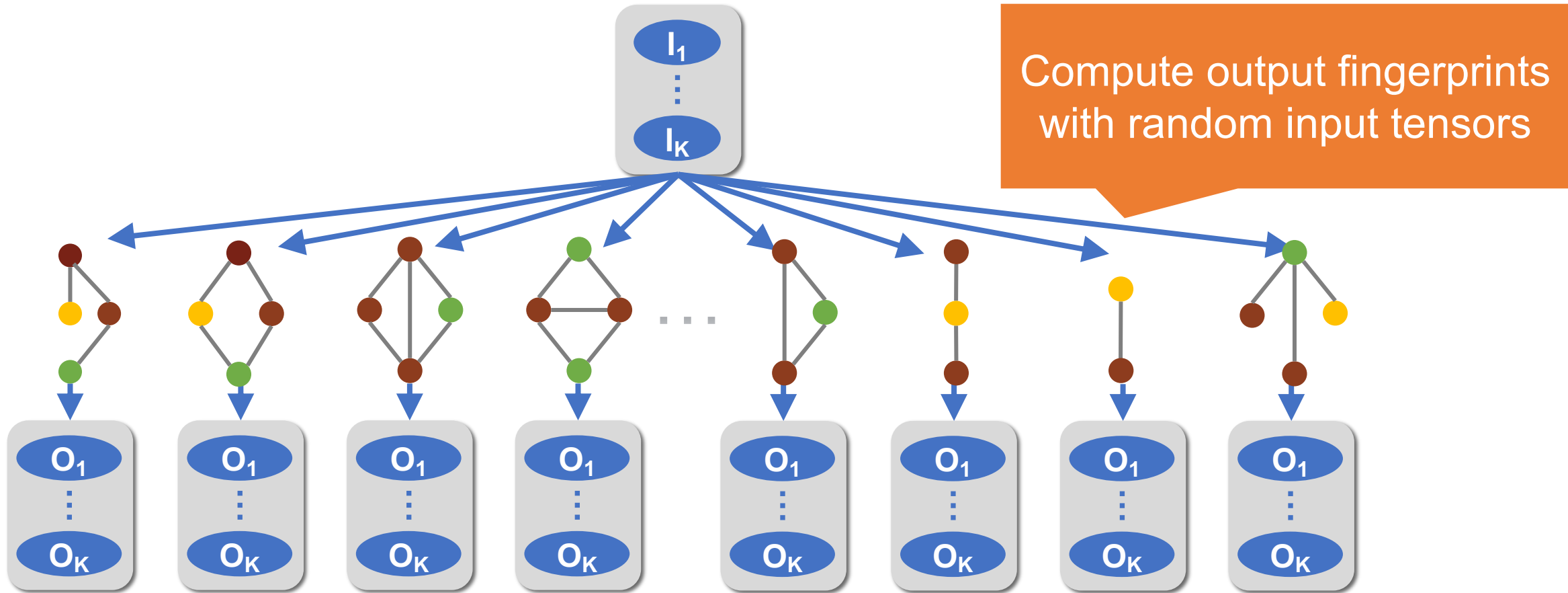
**66M** graphs with up to **4** operators



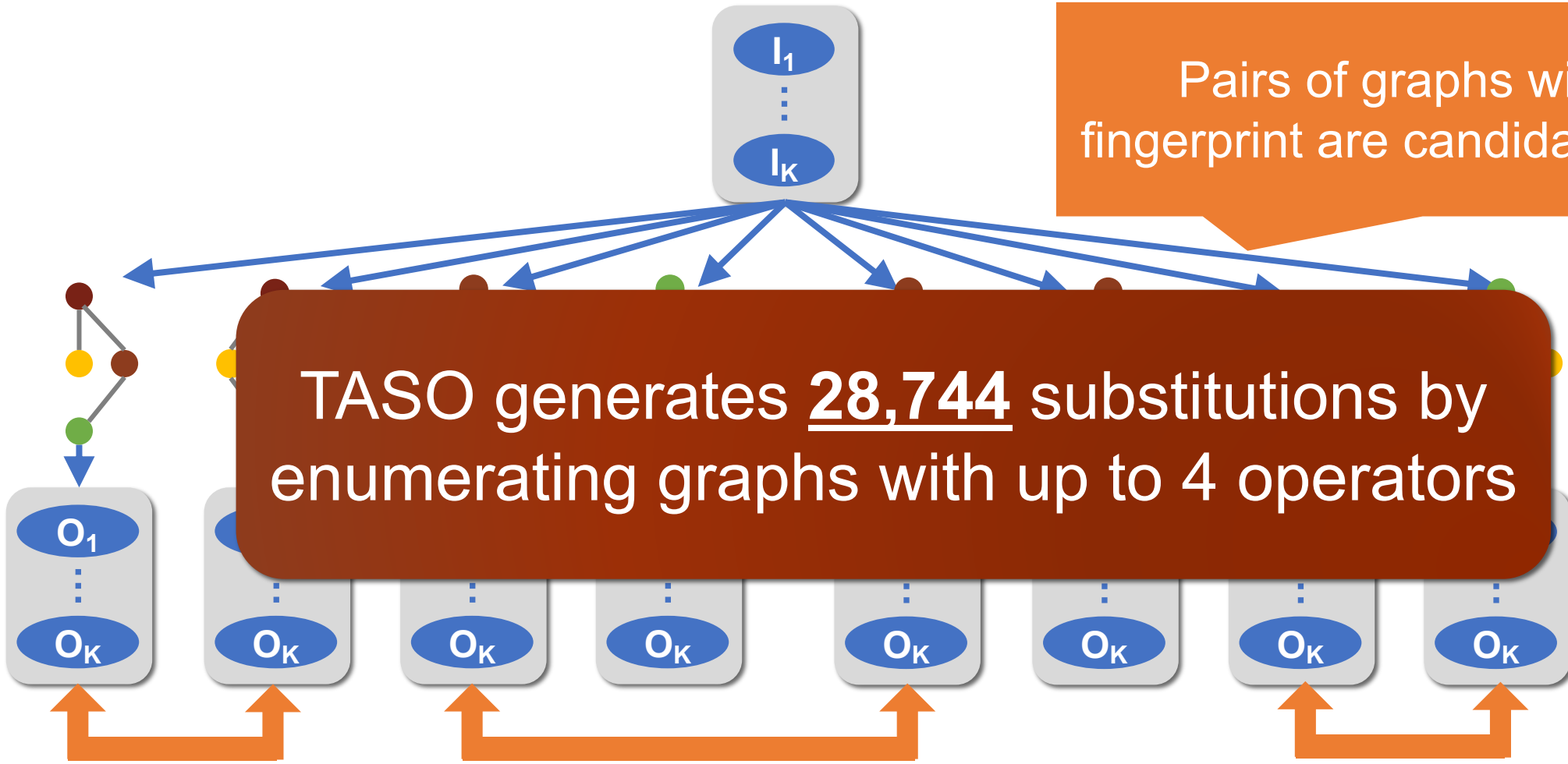
A substitution = a pair of equivalent graphs

**Explicitly considering all pairs does not scale**

# Graph Substitution Generator



# Graph Substitution Generator



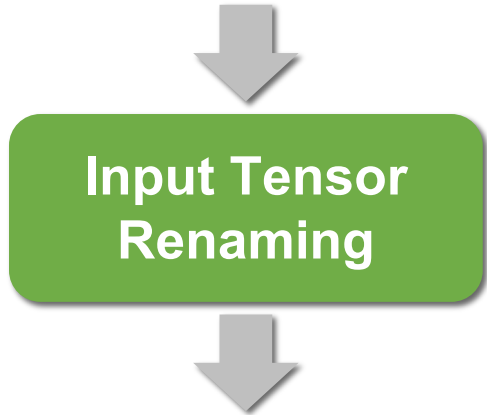
Pairs of graphs with identical fingerprint are candidate substitutions



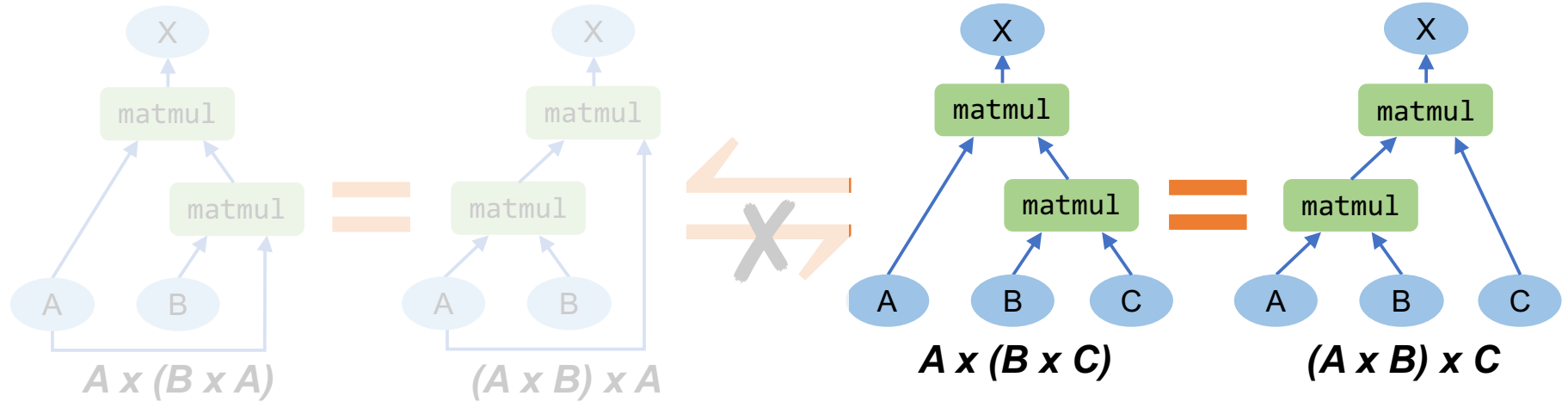
# Pruning Redundant Substitutions



28,744 substitutions



17,346 substitutions



# Pruning Redundant Substitutions

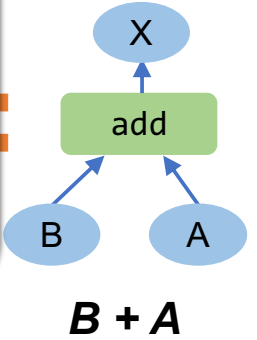


28,744 substitutions

Input Tensor Renaming

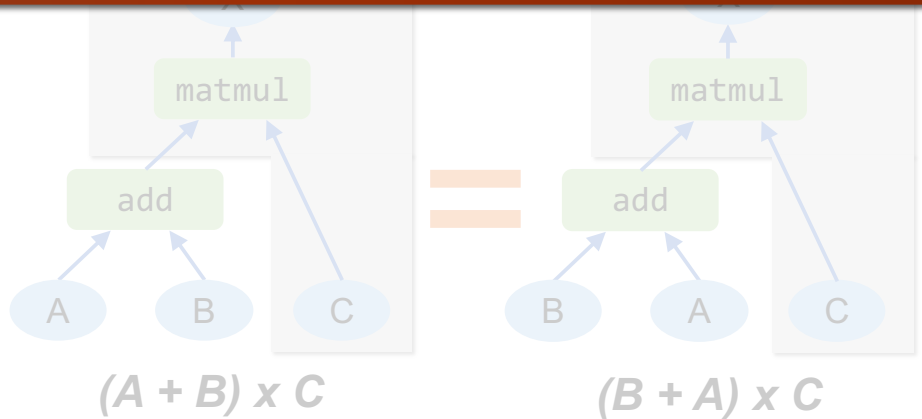


Pruning techniques reduce the number of candidate substitutions by 39x



17,346 substitutions

Common Subgraphs

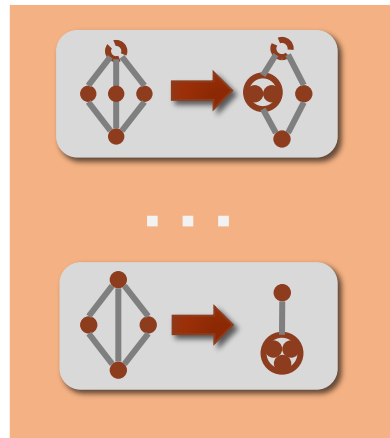


$A + B$

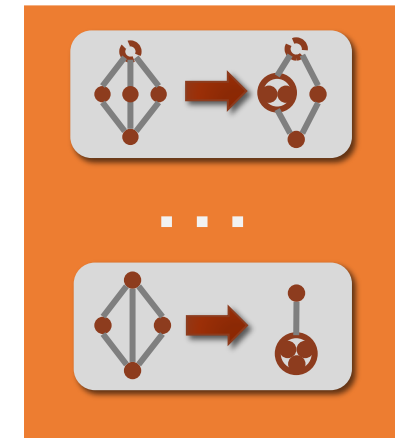
$B + A$

743 substitutions

# Graph Substitution Verifier



Candidate Substitutions



Verified Substitutions

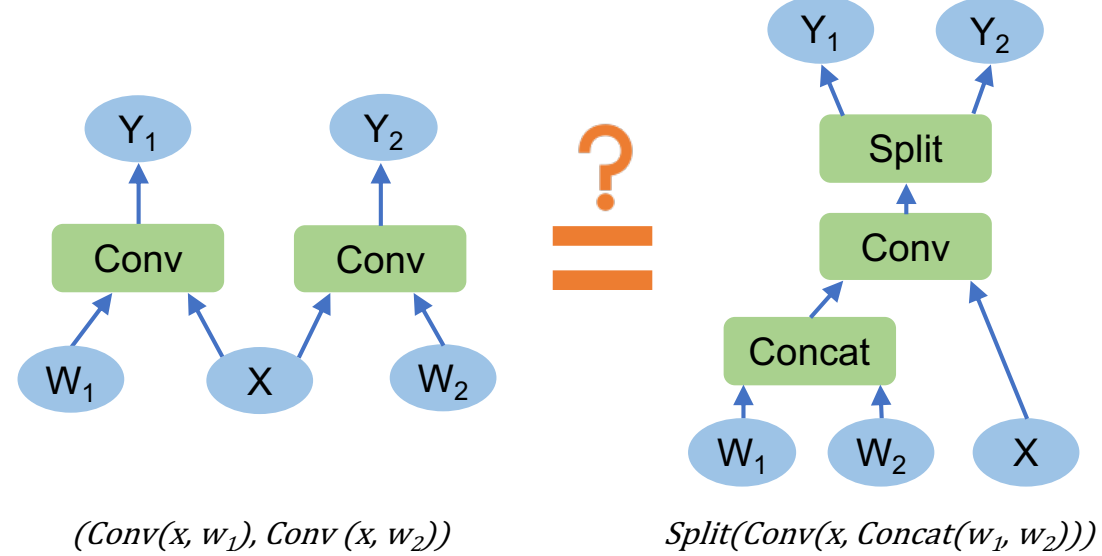
P1. conv is distributive over concatenation  
P2. conv is bilinear  
...  
Pn.



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

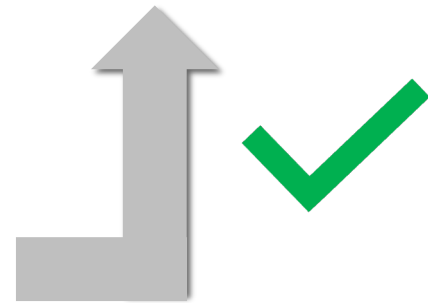


$\forall x, w_1, w_2 .$   
 $(Conv(x, w_1), Conv(x, w_2))$   
 $= Split(Conv(x, Concat(w_1, w_2)))$

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

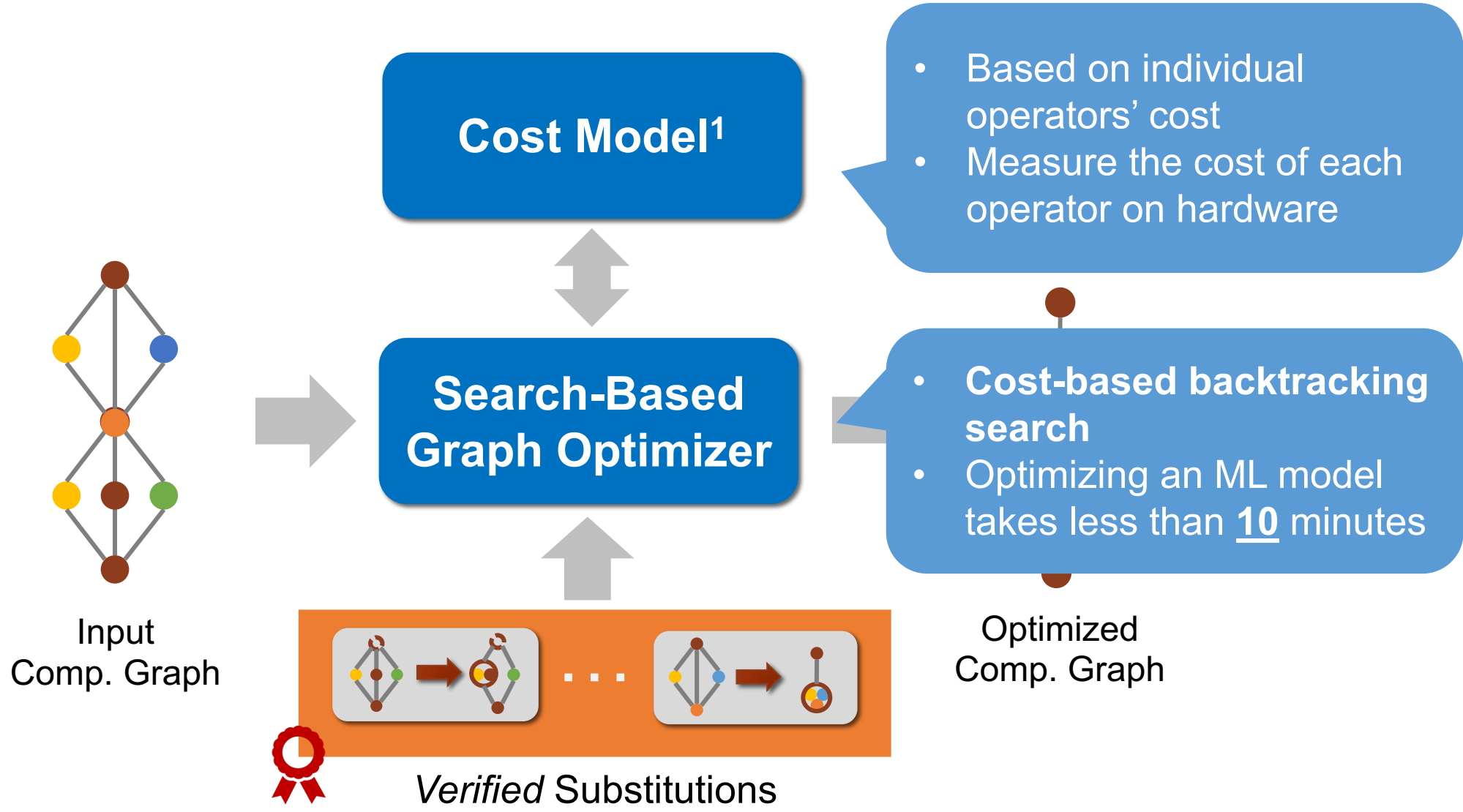
Operator Property	Comment
$\forall x, y, z. \text{ewadd}(x, \text{ewadd}(y, z)) = \text{ewadd}(\text{ewadd}(x, y), z)$	ewadd is associative
$\forall x, y. \text{ewadd}(x, y) = \text{ewadd}(y, x)$	ewadd is commutative
$\forall x, y, z. \text{ewmul}(x, \text{ewmul}(y, z)) = \text{ewmul}(\text{ewmul}(x, y), z)$	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
$\forall x, y, w. \text{smul}(\text{ewadd}(x, y), w) = \text{ewadd}(\text{smul}(x, w), \text{smul}(y, w))$	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
	ReLU applies relu
	commutativity
	conv. with C <sub>pool</sub>
	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

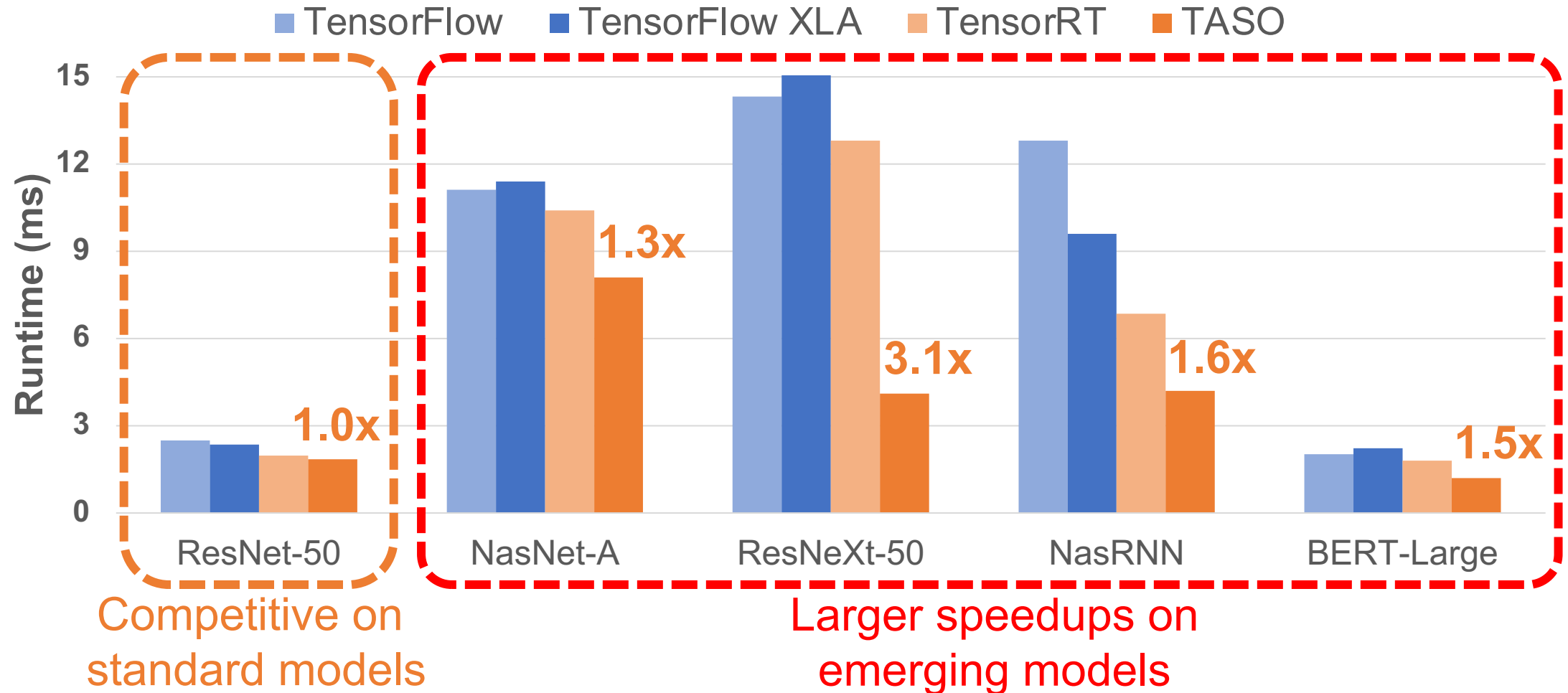
Supporting a new operator requires a few hours of human effort to specify its properties

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

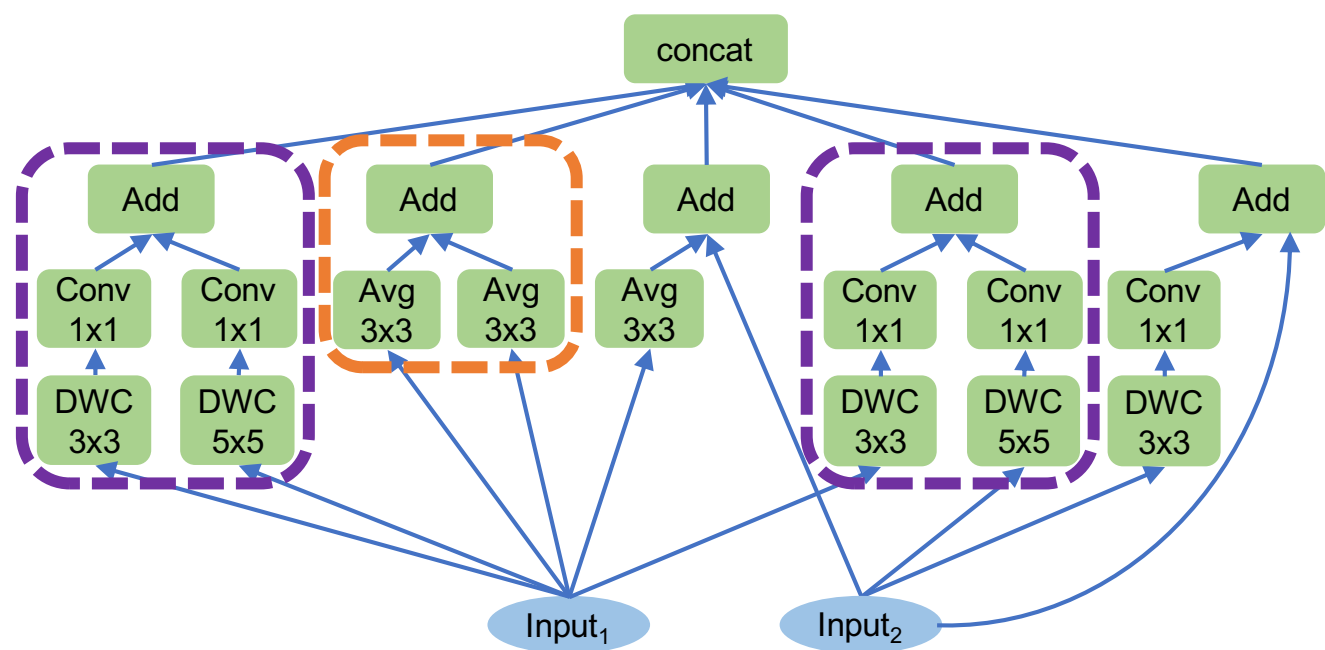
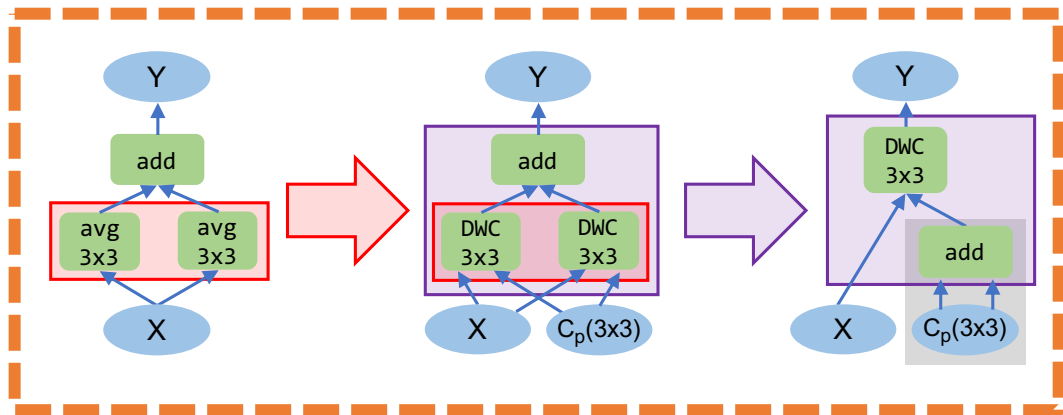
# Search-Based Graph Optimizer



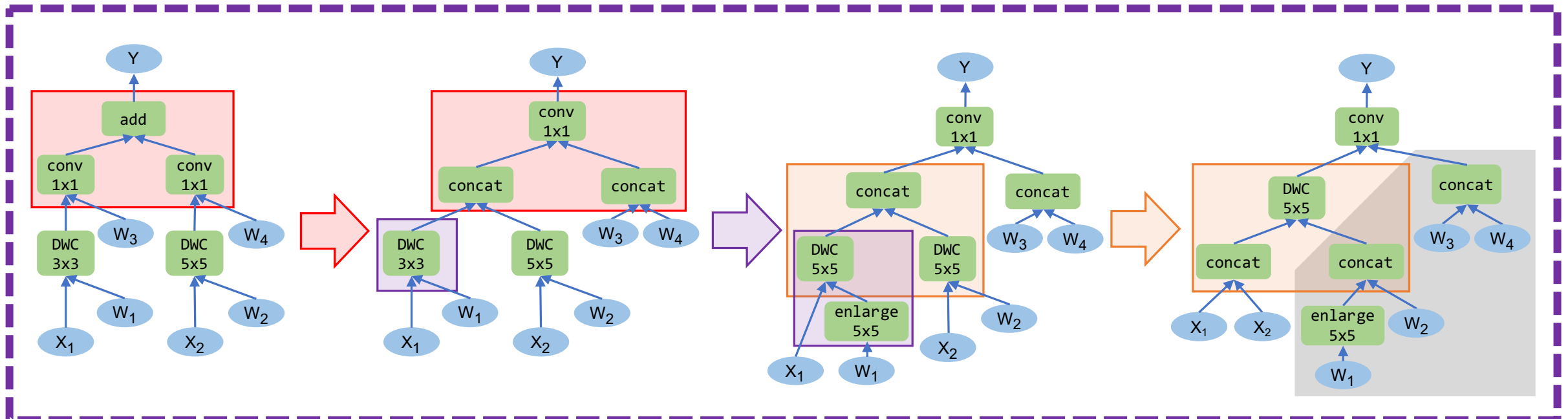
# End-to-end Inference Performance (Nvidia V100 GPU)



# Case Study: NASNet



\*DWC: depth-wise convolution





# Why TASO is a SuperOptimizer?

*What is the difference between optimizer and super-optimizer?*

Goal: gradually *improve* an input program by greedily applying optimizations

Goal: automatically find an *optimal* program for an input program

**PET:**

# Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections

Haojie Wang, Jidong Zhai, Mingyu Gao, Zixuan Ma, Shizhi Tang,  
Liyan Zheng, Yuanzhi Li, Kaiyuan Rong, Yuanyong Chen, Zhihao Jia

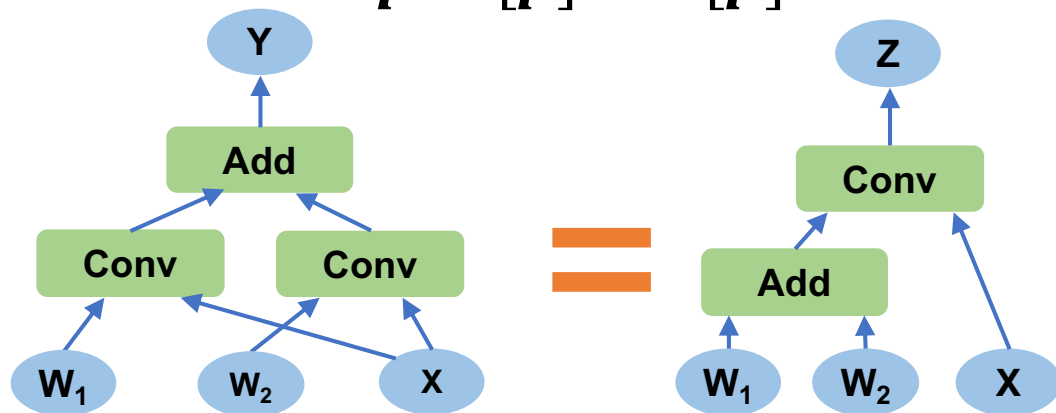
Tsinghua University

Carnegie Mellon University

Facebook

# Motivation: Current Systems Consider only Fully Equivalent Transformations

$$\forall p. Y[p] = Z[p]$$

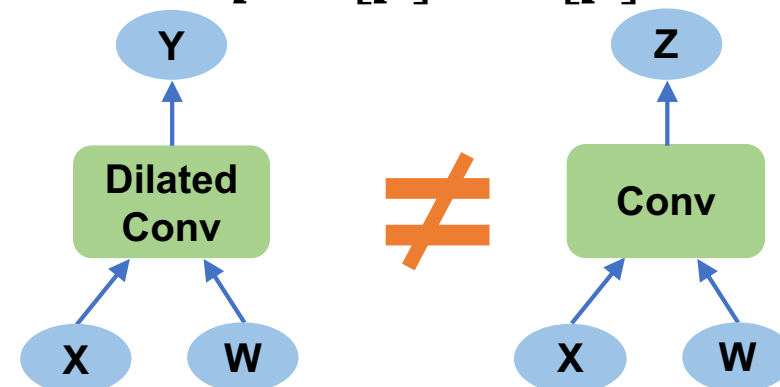


Fully Equivalent Transformations

👍 Pro: preserve functionality

👎 Con: miss optimization opportunities

$$\exists p. Y[p] \neq Z[p]$$



Partially Equivalent Transformations

👍 Pro: better performance

- Faster ML operators
- More efficient tensor layouts
- Hardware-specific optimizations

👎 Con: potential accuracy loss

# Motivation: Current Systems Consider only Fully Equivalent Transformations

$$\forall p. Y[p] = Z[p]$$



$$\exists p. Y[p] \neq Z[p]$$



**Is it possible to exploit partially equivalent transformations to improve performance while preserving equivalence?**



Fully Equivalent Transformations

Pro: preserve functionality

Con: miss optimization opportunities



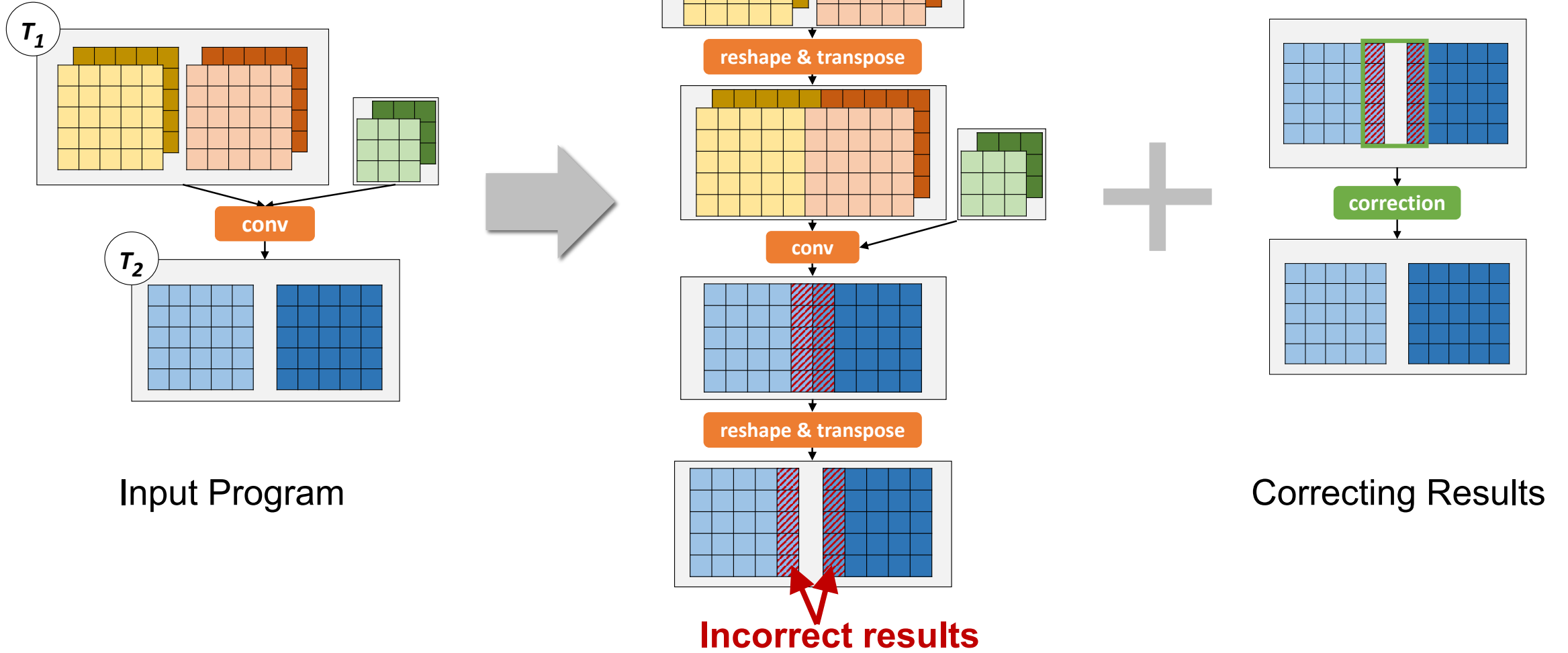
Partially Equivalent Transformations

Pro: better performance

- Faster ML operators
- More efficient tensor layouts
- Hardware-specific optimizations

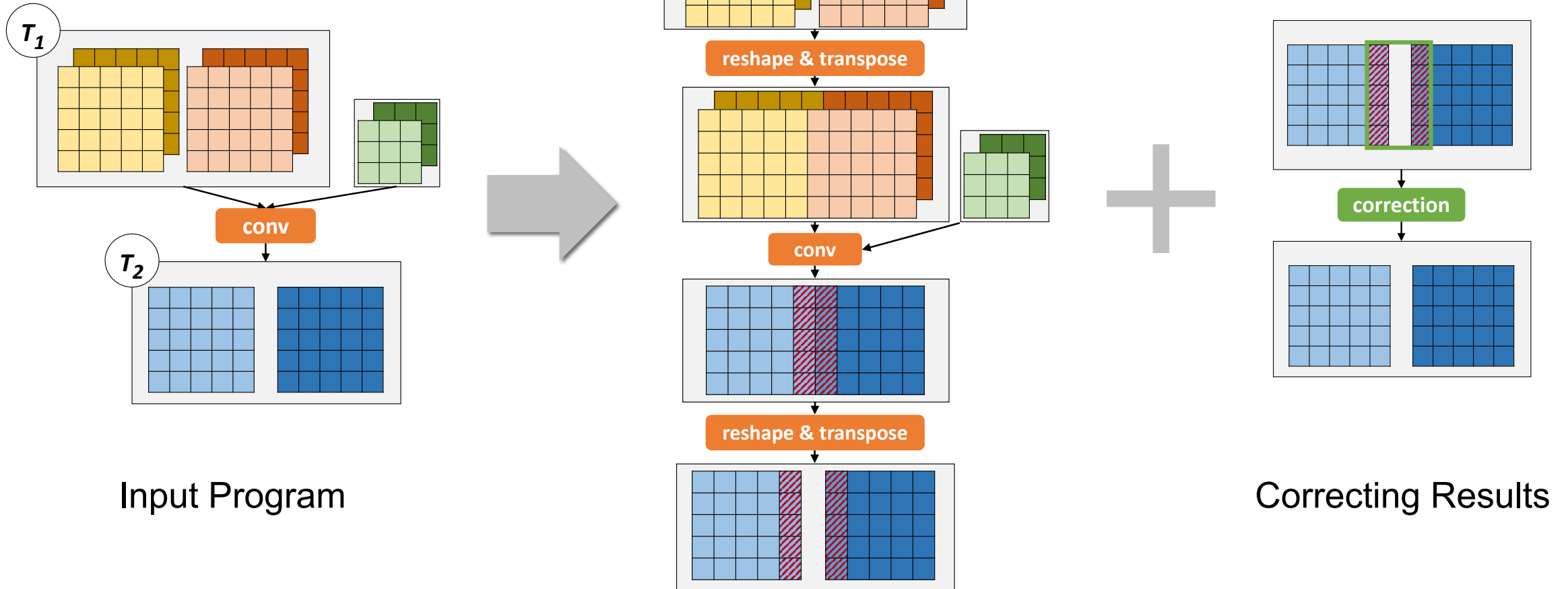
Con: potential accuracy loss

# Motivating Example



Partially Equivalent Transformation

# Motivating Example

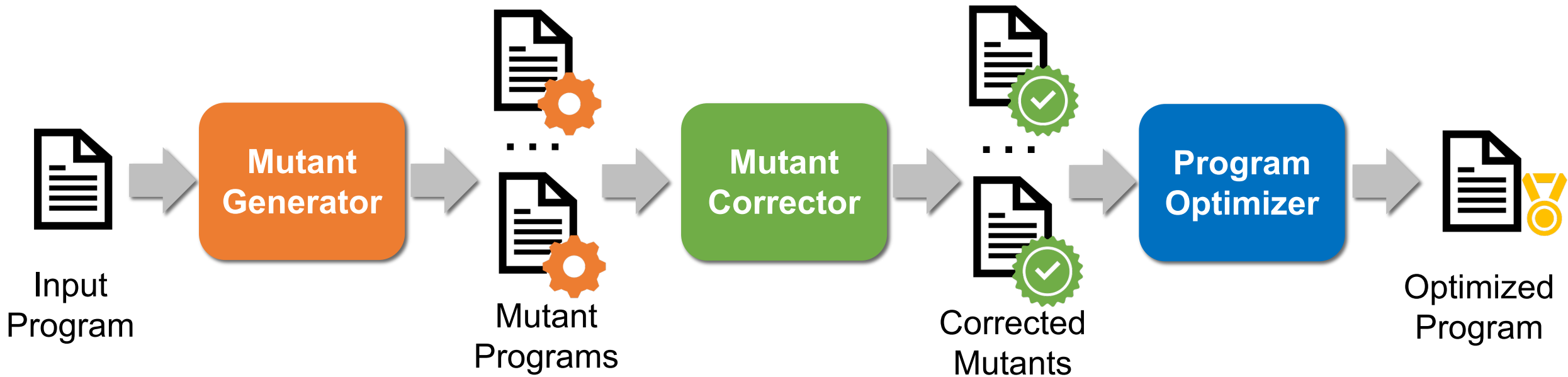


- Transformation and correction lead to **1.2x** speedup for ResNet-18
- Correction preserves end-to-end equivalence

# PET

- **First tensor program optimizer** with partially equivalent transformations
- **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

# PET Overview





# Key Challenges

1. How to generate partially equivalent transformations?

Superoptimization

2. How to correct them?

Multi-linearity of DNN computations



# Mutant Generator

Superoptimization adapted from TASO<sup>1</sup>

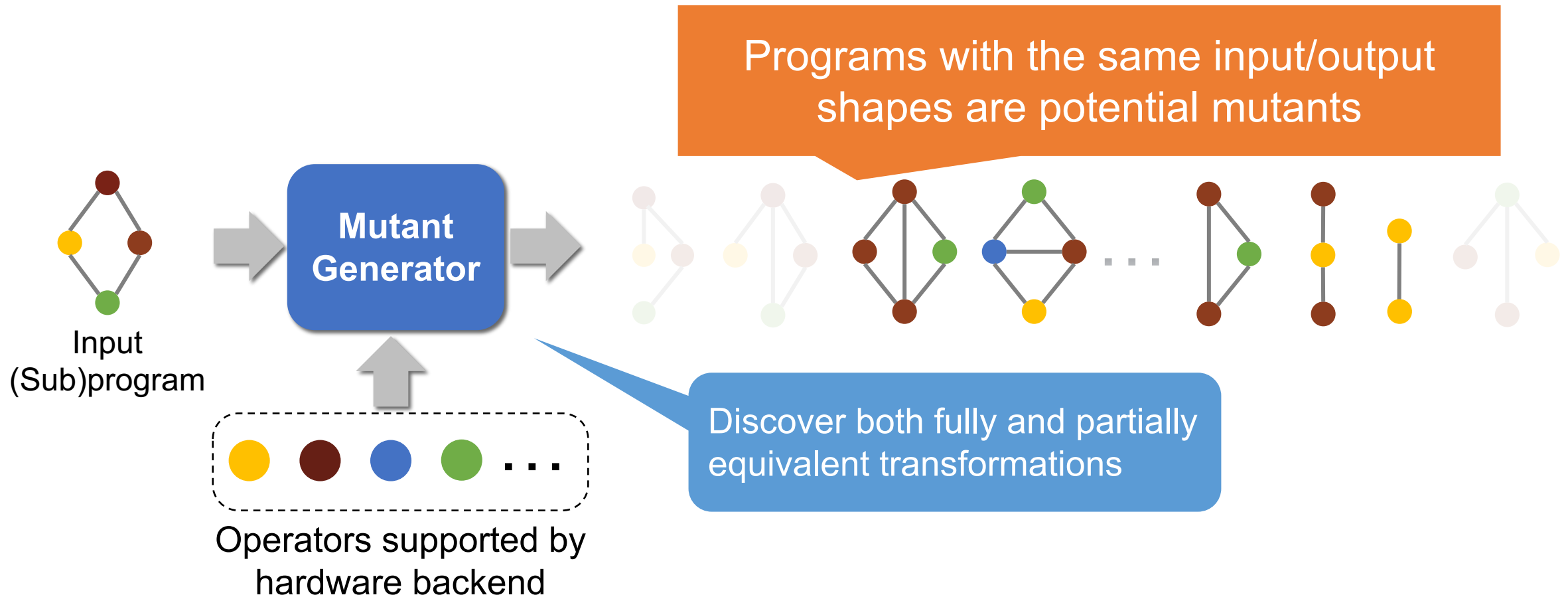
Enumerate all possible programs up to a fixed size using available operators





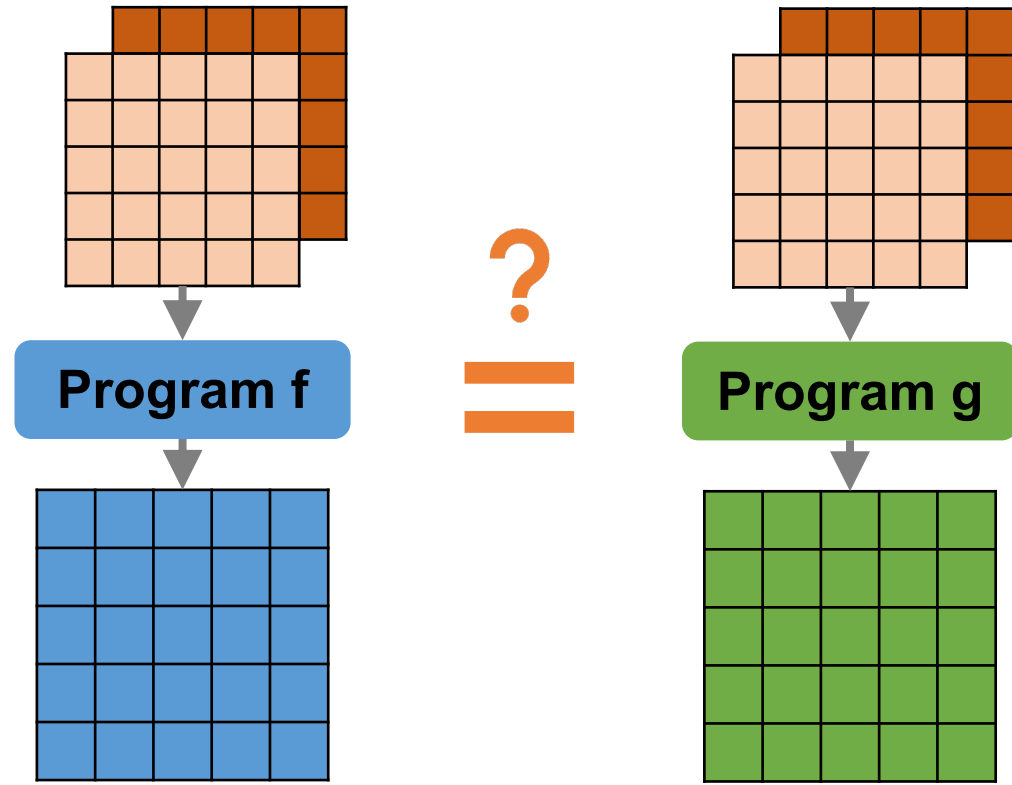
# Mutant Generator

Superoptimization adapted from TASO<sup>1</sup>





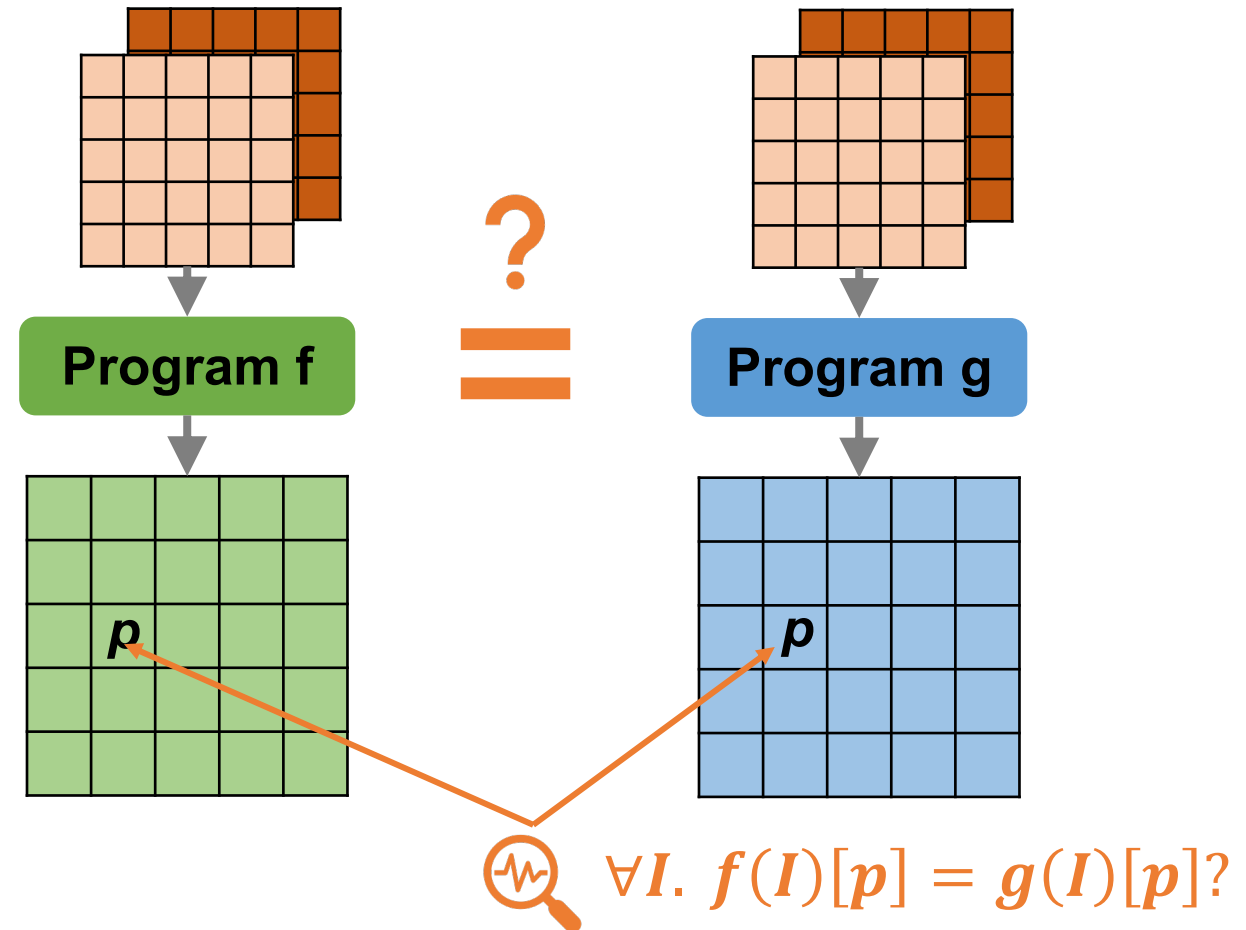
# Challenges: Examine Transformations



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

# A Strawman Approach

- **Step 1:** Explicitly consider all output positions (m positions)
- **Step 2:** For each position  $p$ , examine all possible inputs (n inputs)



**Require  $O(m * n)$  examinations, but both m and n are too large to explicitly enumerate**

# Multi-Linear Tensor Program (MLTP)

- A program  $f$  is multi-linear if the output is linear to all inputs
  - $f(I_1, \dots, X, \dots, I_n) + f(I_1, \dots, Y, \dots, I_n) = f(I_1, \dots, X + Y, \dots, I_n)$
  - $\alpha \cdot f(I_1, \dots, X, \dots, I_n) = f(I_1, \dots, \alpha \cdot X, \dots, I_n)$
- DNN computation = MLTP + non-linear activations

Majority of the computation

**$O(m * n)$  examinations  
in strawman approach**

**MLTP**

**$O(1)$  examinations in  
PET's approach**

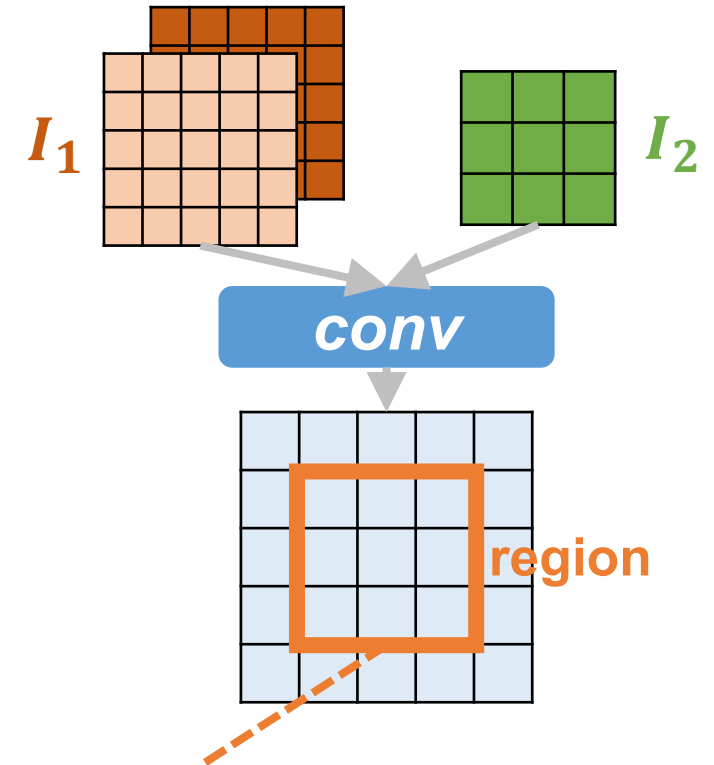
# No Need to Enumerate All Output Positions

Group all output positions with an identical **summation interval** into a **region**

**\*Theorem 1:** For two MLTPs **f** and **g**, if **f=g** for **O(1)** positions in a **region**, then **f=g** for all positions in the **region**

Only need to examine **O(1)** positions for each region.

**Complexity:**  $O(m * n) \rightarrow O(n)$



$$conv(c, h, w) = \sum_{d=0}^{D-1} \sum_{x=-1}^1 \sum_{y=-1}^1 I_1(d, h+x, w+y) \times I_2(d, c, x, y)$$

**Summation interval**

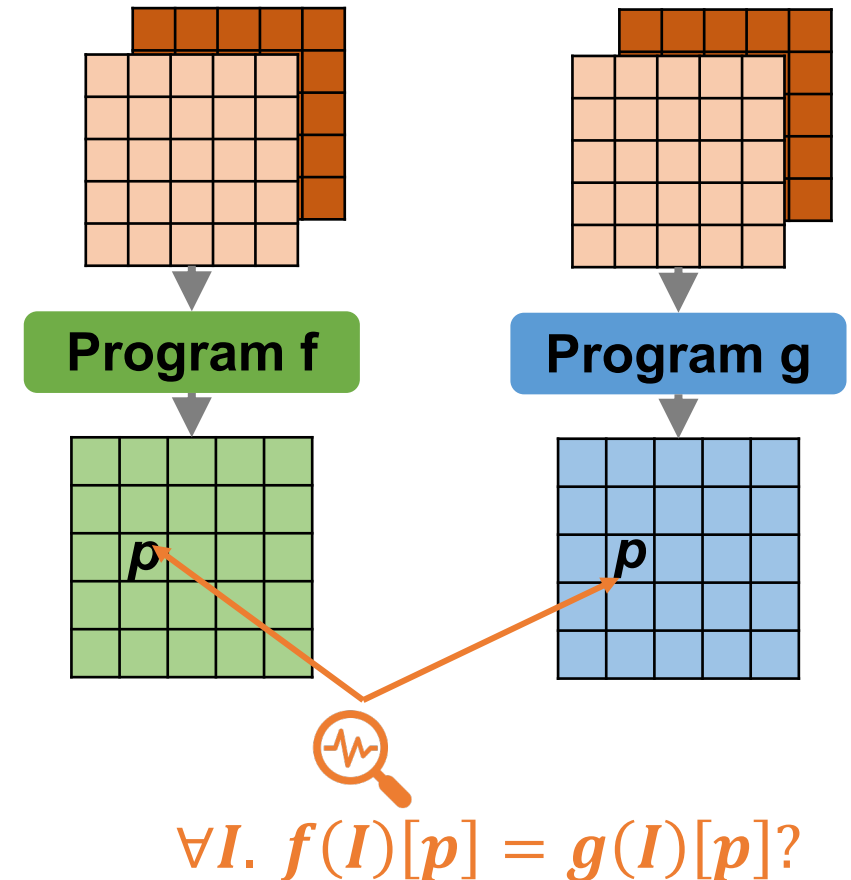
# No Need to Consider All Possible Inputs

Examining equivalence for a single position is still challenging

**\*Theorem 2:** If  $\exists I. f(I)[p] \neq g(I)[p]$ , then the probability that **f** and **g** give identical results on  $t$  random integer inputs is  $(\frac{1}{2^{31}})^t$

Run  $t$  random tests for each position  $p$

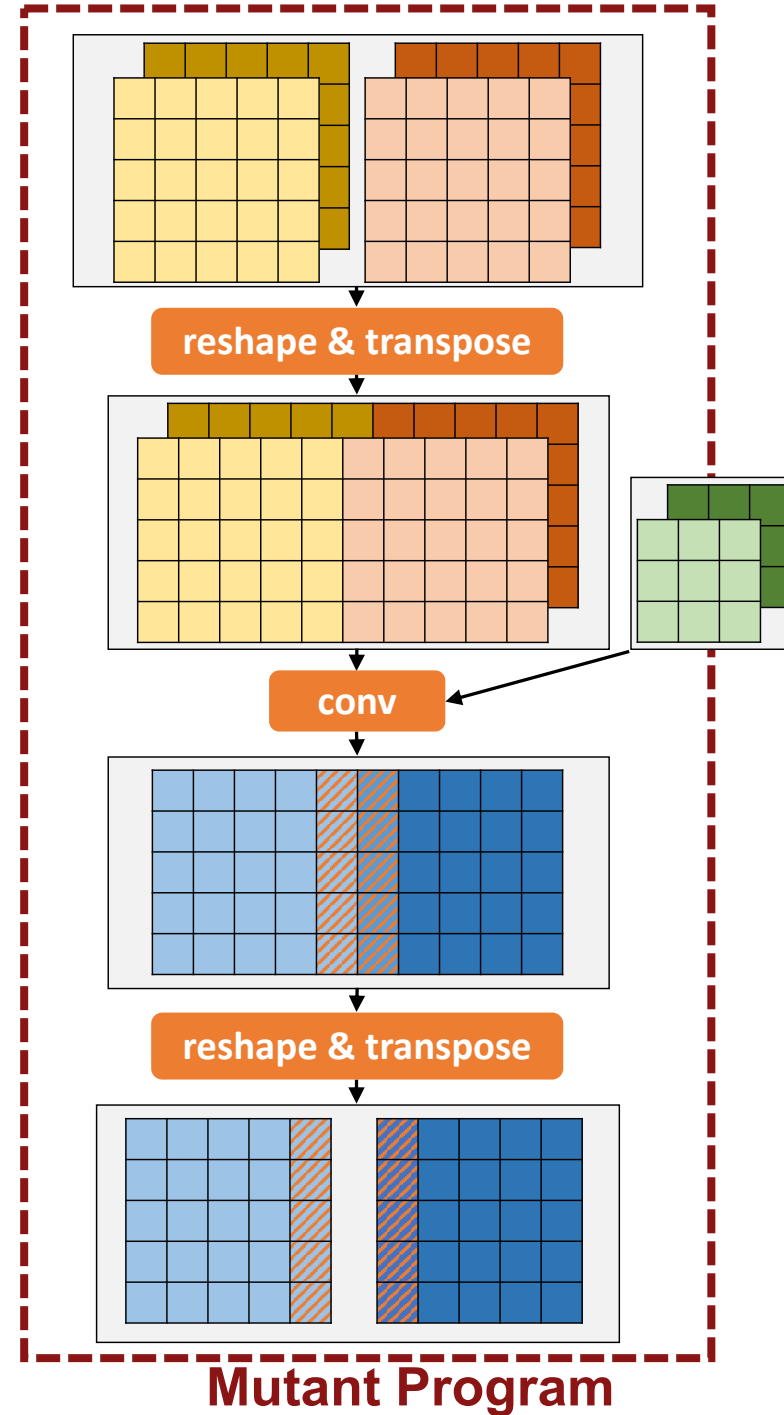
**Complexity:**  $O(n) \rightarrow O(t) = O(1)$





# Mutant Corrector

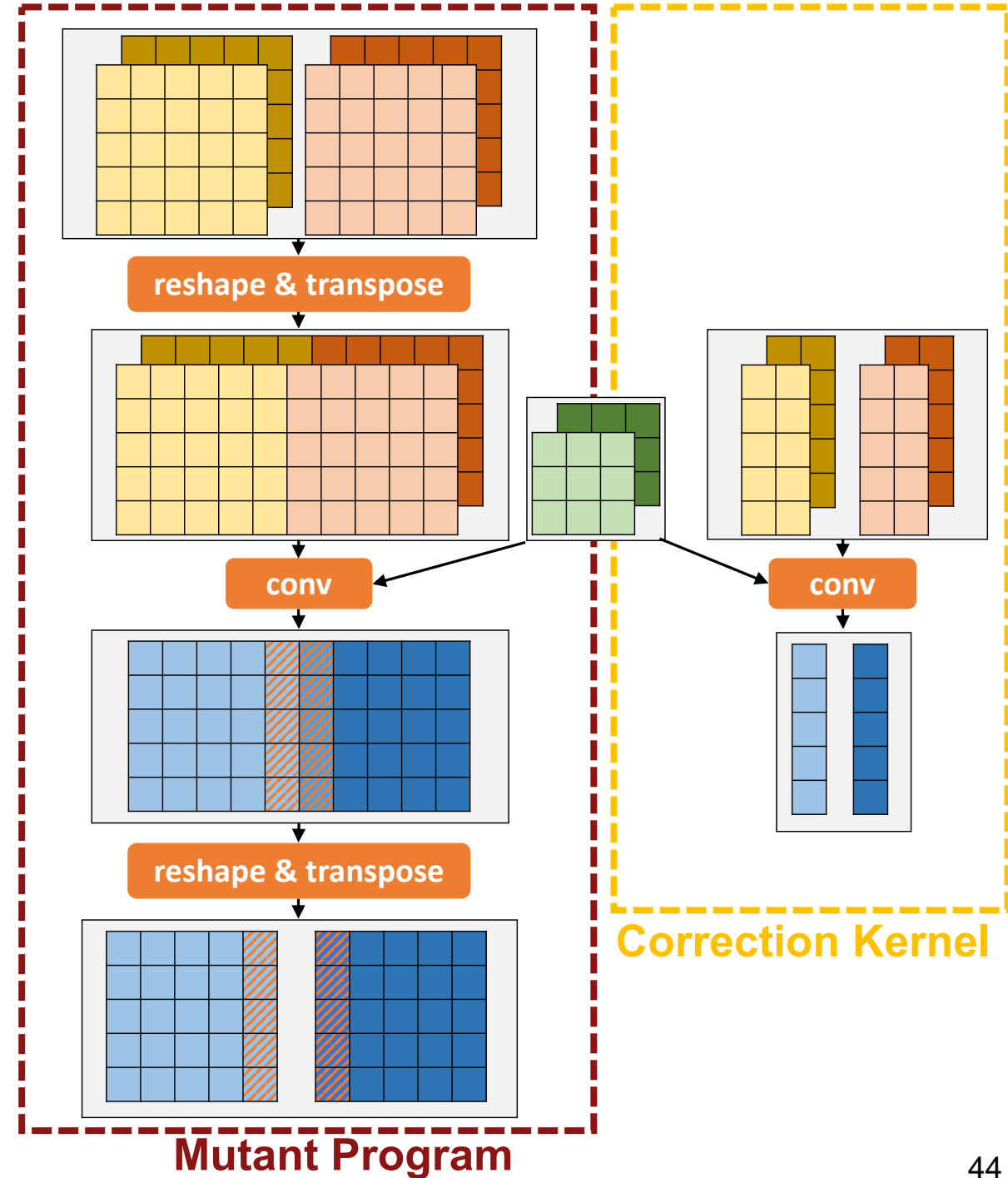
**Goal:** quickly and efficiently correcting the outputs of a mutant program



# Mutant Corrector

**Goal:** quickly and efficiently correcting the outputs of a mutant program

**Step 1:** recompute the incorrect outputs using the original program



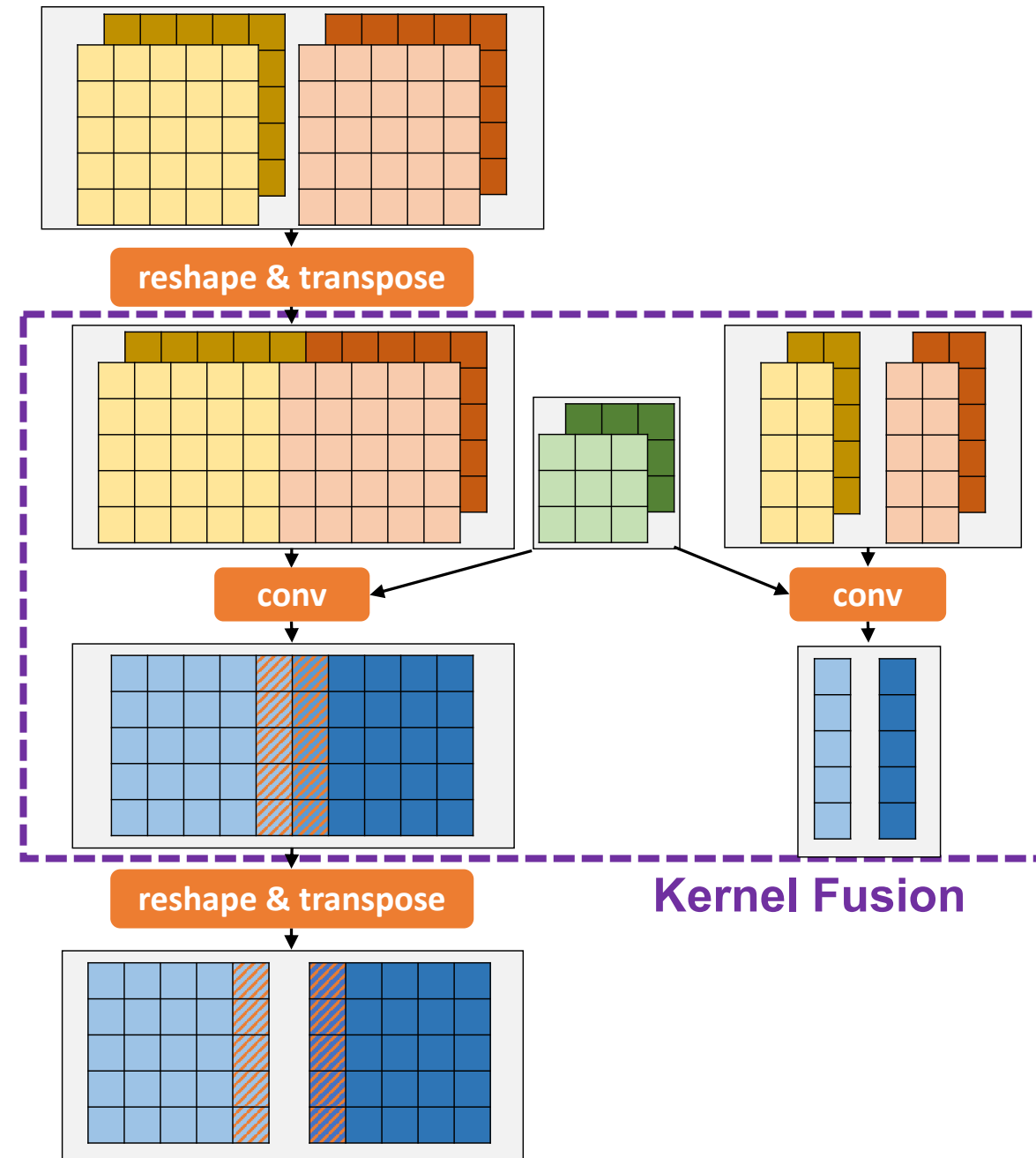
# Mutant Corrector

**Goal:** quickly and efficiently correcting the outputs of a mutant program

**Step 1:** recompute the incorrect outputs using the original program

**Step 2:** opportunistically fuse correction kernels with other operators

Correction introduces less than 1% overhead



# Program Optimizer

- **Beam search**
- Optimizing a DNN architecture takes less than 30 minutes

- Other optimizations:
- Operator fusion
  - Constant folding
  - Redundancy elimination

Input Program



**Search-Based Program Optimizer**



Optimized Program



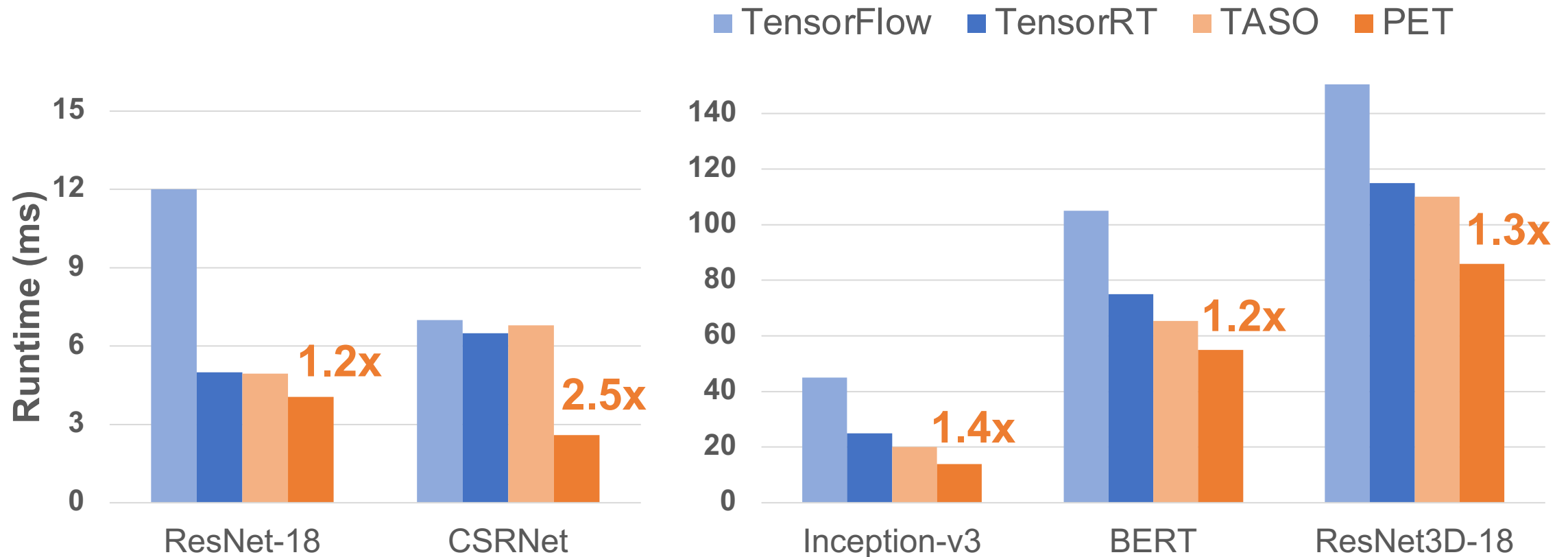
MLTP

**Mutant Generator & Corrector**



Mutants w/ Corrections

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

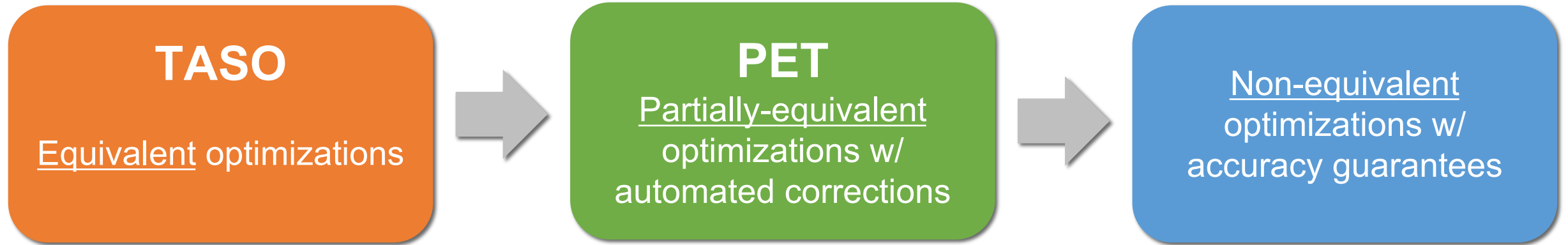


**PET outperforms existing optimizers by 1.2-2.5x by combining fully and partially equivalent transformations**

# PET

- A **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

# From Equivalent to Non-Equivalent Optimizations for ML



**Week 11: Model Pruning,  
Quantization, Distillation, etc.**

# Questions to Discuss

1. How does PET differ from TASO in generating graph transformations?
2. How does PET differ from TASO in verifying/correcting transformations?
3. How can we combine graph optimizations with kernel optimizations?