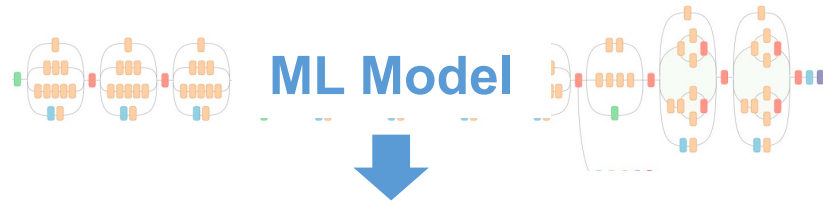# 15-442/15-642: Machine Learning Systems

# Graph-Level Optimizations

**Tianqi Chen and Zhihao Jia**

Carnegie Mellon University

# Recap: An Overview of Deep Learning Systems

**ML Model**

Automatic Differentiation

Graph-Level Optimization ← **Lecture 10**

Parallelization / Distributed Training ← **Lecture 11, 12**
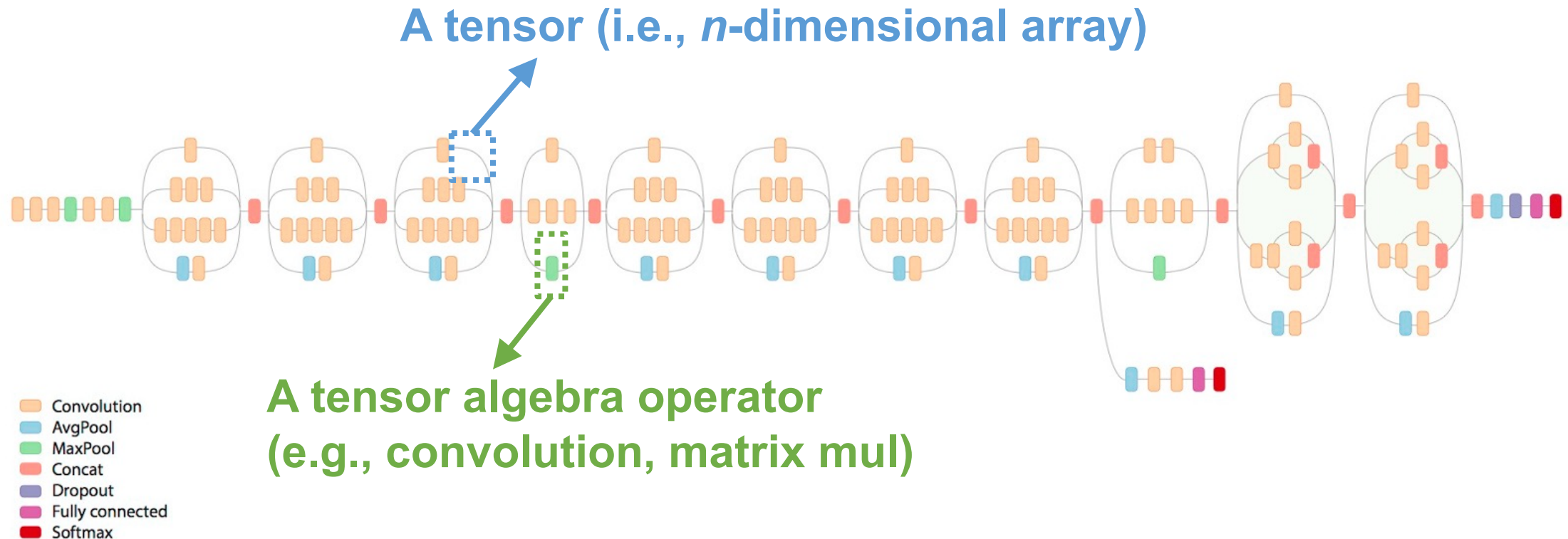
Kernel Generation

Memory Optimization ← **Lecture 13**



2

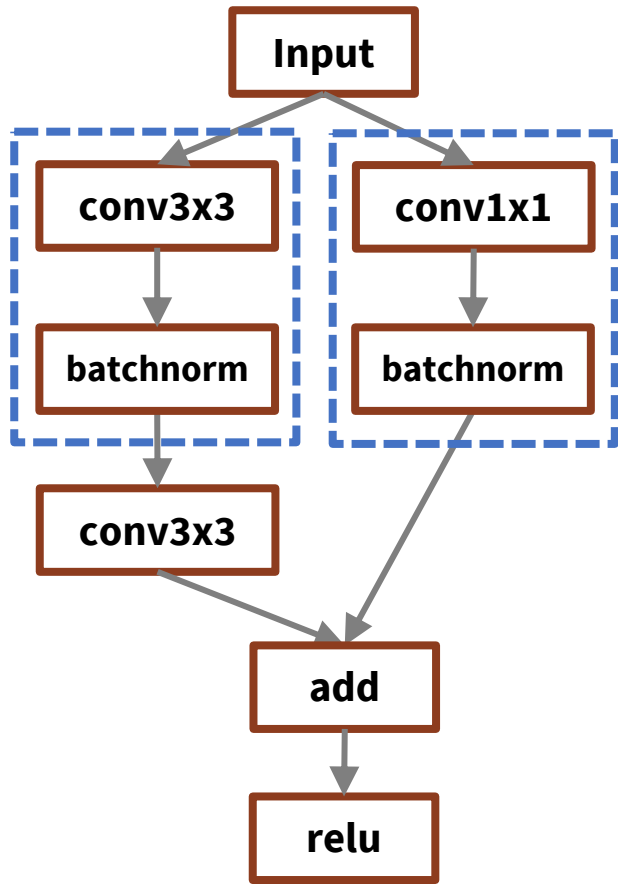# Recap: Deep Neural Network

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

**A tensor (i.e., *n*-dimensional array)**

**A tensor algebra operator (e.g., convolution, matrix mul)**

Convolution
AvgPool
MaxPool
Concat
Dropout
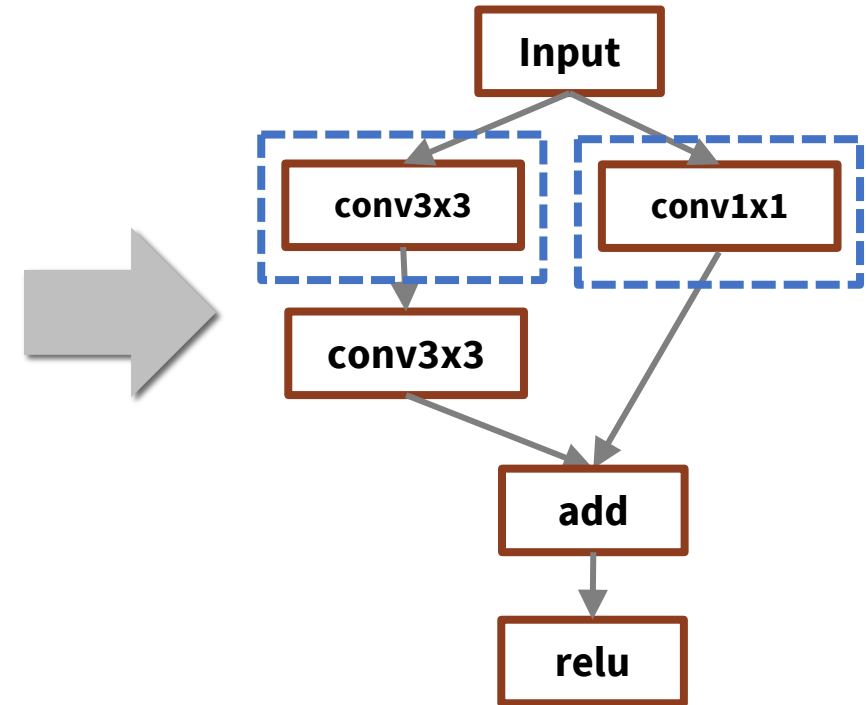Fully connected
Softmax

# Graph-Level Optimizations
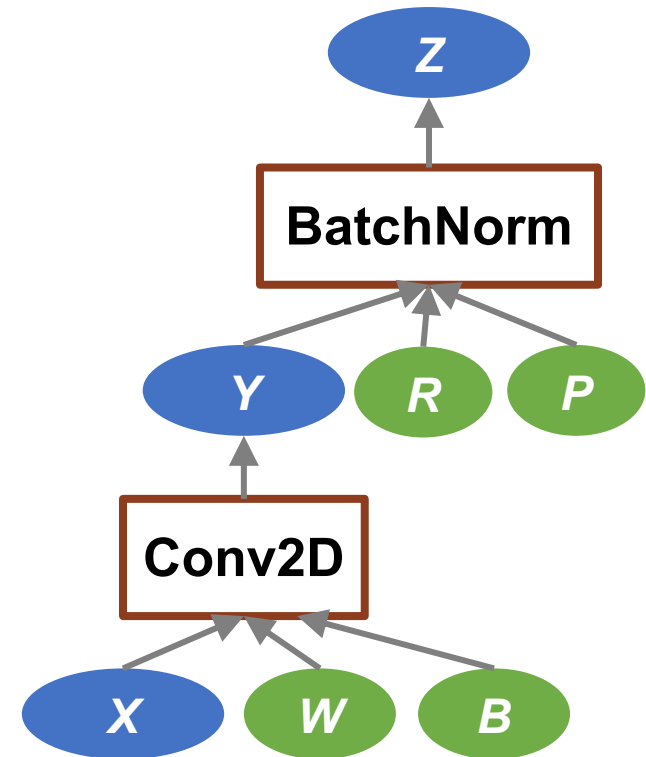


**Input Computation Graph**

**Potential graph transformations**

**Optimized Computation Graph**

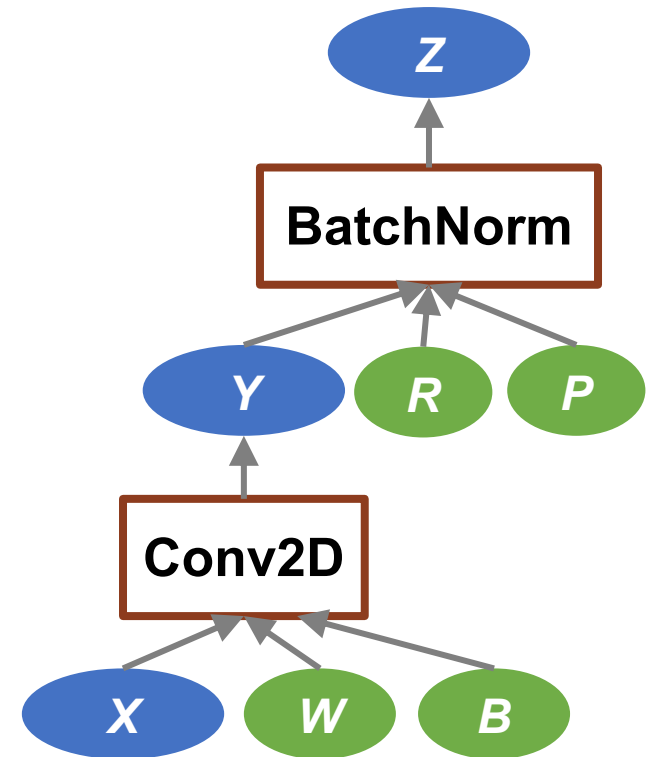# Example: Fusing Convolution and Batch Normalization



$$Z(n, c, h, w) = Y(n, c, h, w) * R(c) + P(c)$$

$$Y(n, c, h, w) = \left( \sum_{d,u,v} X(n, d, h + u, w + v) * W(c, d, u, v) \right) + B(n, c, h, w)$$
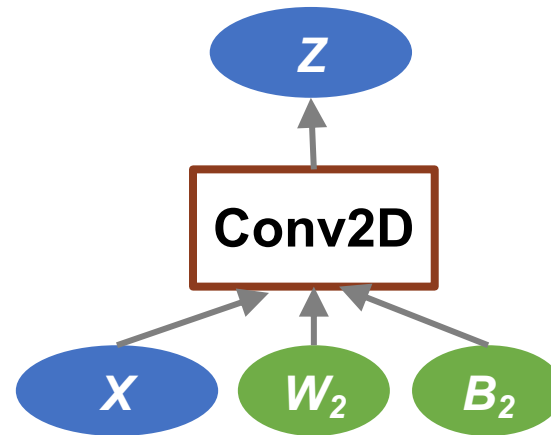
**W, B, R, P are constant pre-trained weights**

# Fusing Conv and BatchNorm



$$Z(n, c, h, w) = \left( \sum_{d,u,v} X(n, d, h + u, w + v) * W_2(c, d, u, v) \right) + B_2(n, c, h, w)$$
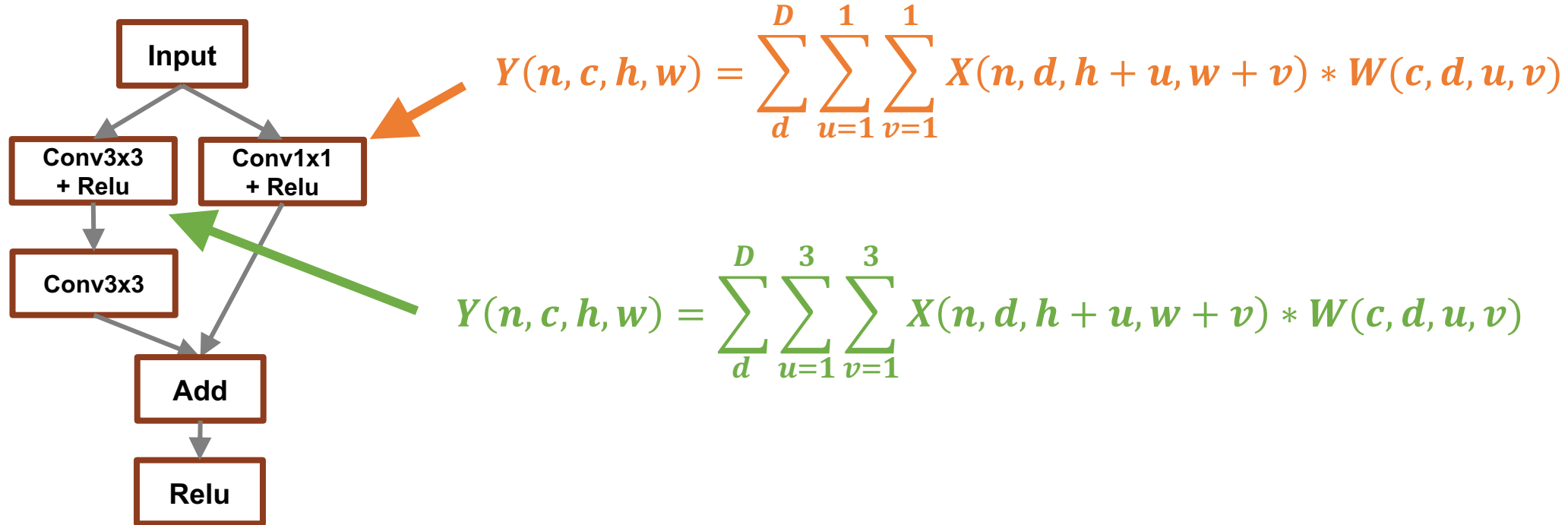
$$W_2(n, c, h, w) = W(n, c, h, w) * R(c)$$

$$B_2(n, c, h, w) = B(n, c, h, w) * R(c) + P(c)$$

# Recap: Resnet Example



$$Y(n, c, h, w) = \sum_{d}^{D} \sum_{u=1}^{1} \sum_{v=1}^{1} X(n, d, h+u, w+v) * W(c, d, u, v)$$

$$Y(n, c, h, w) = \sum_{d}^{D} \sum_{u=1}^{3} \sum_{v=1}^{3} X(n, d, h+u, w+v) * W(c, d, u, v)$$

* Kaiming He. et al. Deep Residual Learning for Image Recognition, 2015

# Recap: Resnet Example



$$Y(n, c, h, w) = \sum_{d} \sum_{u=1}^{3} \sum_{v=1}^{3} X(n, d, h+u, w+v) * W(c, d, u, v)$$

Enlarge convs

(Decrease performance)

* Kaiming He. et al. Deep Residual Learning for Image Recognition, 2015

8

$$Y(n, c, h, w) = \sum_{d}^{D} \sum_{u=1}^{3} \sum_{v=1}^{3} X(n, d, h+u, w+v) * W'(c, d, u, v)$$

Input

Conv3x3 + Relu | Conv1x1 + Relu

Conv3x3

Add

Relu

**Enlarge convs**

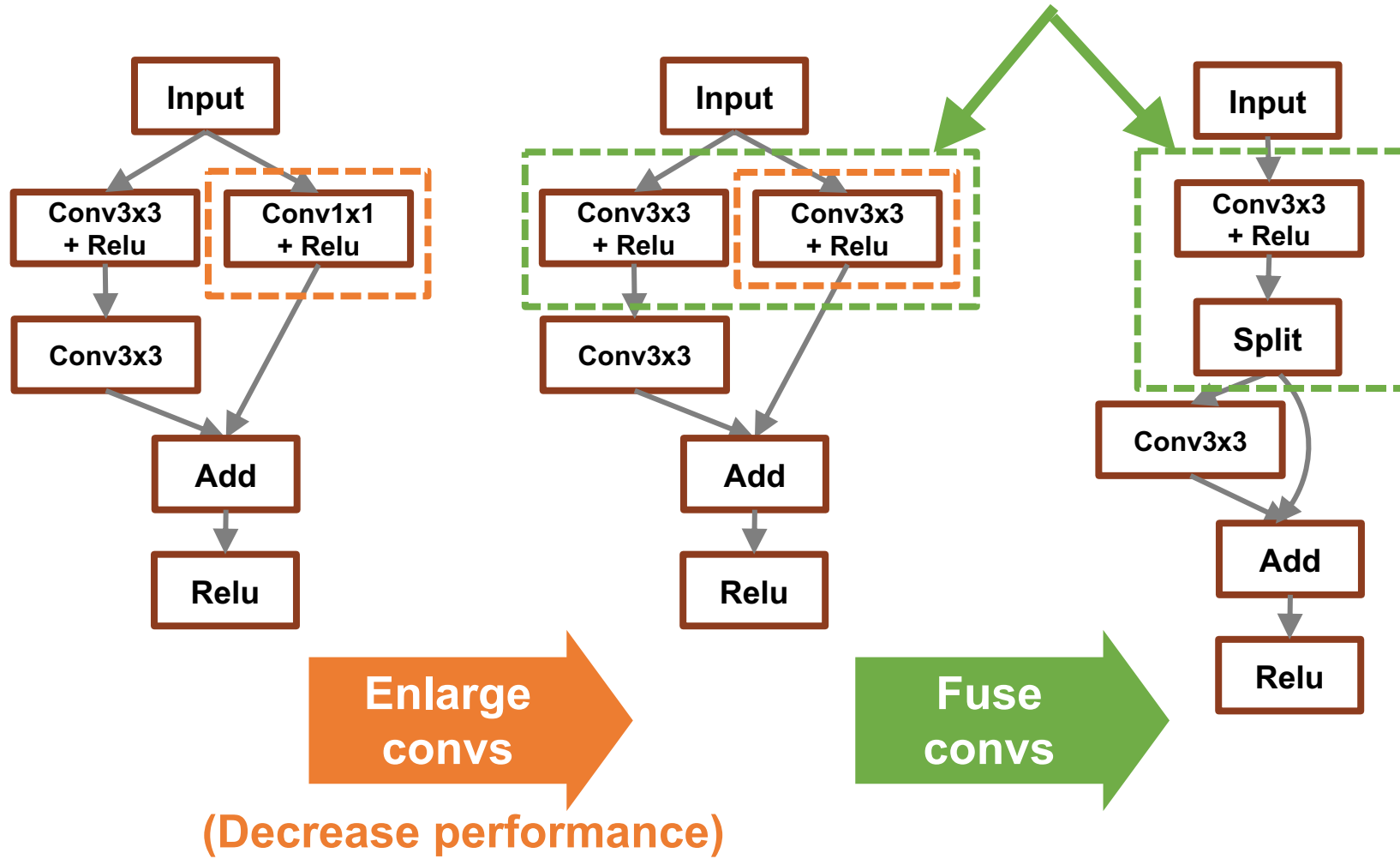**(Decrease performance)**

Input

Conv3x3 + Relu | Conv3x3 + Relu

Conv3x3

Add

Relu

**Fuse convs**

Input

Conv3x3 + Relu

Split

Conv3x3

Add

Relu

\* Kaiming He. et al. Deep Residual Learning for Image Recognition, 2015

# Recap: Resnet Example



**Enlarge convs** (Decrease performance)

**Fuse convs**

**Fuse conv & add**

**Fuse conv & relu**

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

* Kaiming He. et al. Deep Residual Learning for Image Recognition, 2016

10

# Challenge of Graph Optimizations for ML



**Graph Optimizations**

**ML Operators** × **Graph Architectures** × **Hardware Backends** =

Infeasible to manually design graph optimizations for all cases

# This Lecture

- TASO: Automatically Generate Graph Transformations
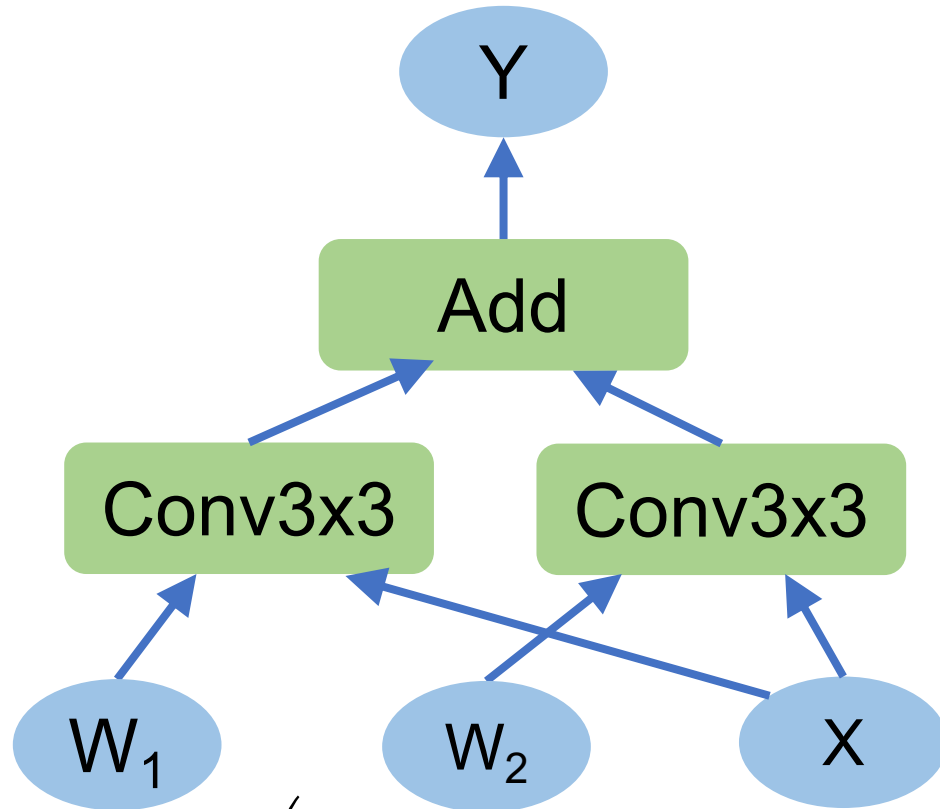- PET: Discover Partially-Equivalent Graph Transformations

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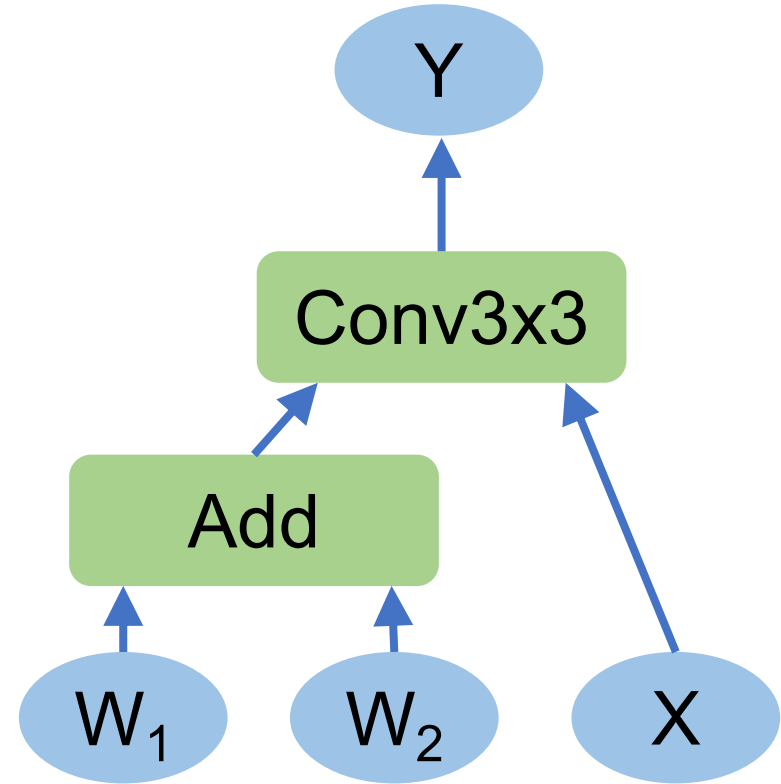
# TASO: Tensor Algebra SuperOptimizer

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

- **Less engineering effort:** <u>53,000</u> LOC for manual graph optimizations in TensorFlow → <u>1,400</u> LOC in TASO

- **Better performance:** outperform existing optimizers by up to <u>3x</u>

- **Stronger correctness:** formally verify all generated substitutions
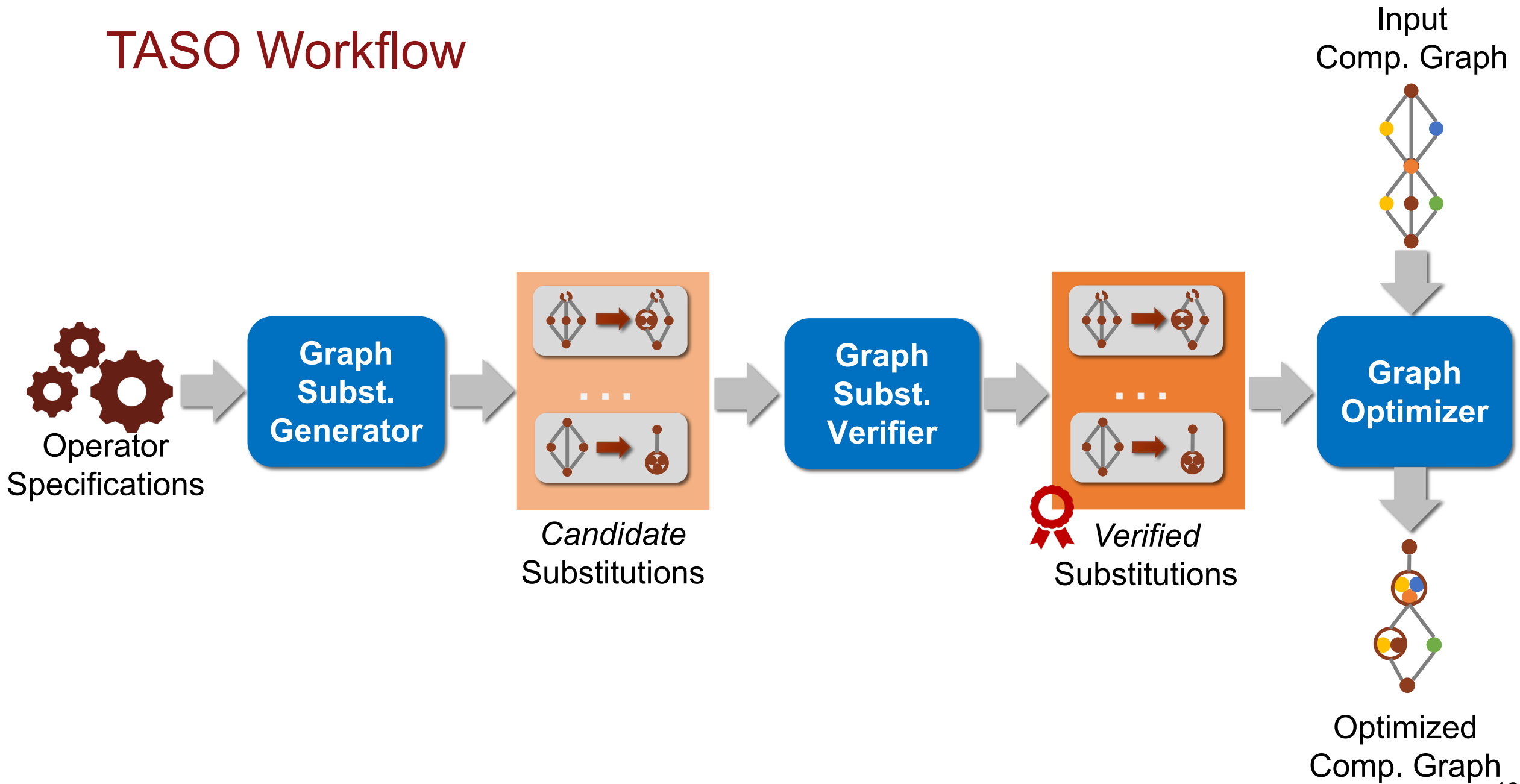
# Graph Substitution



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

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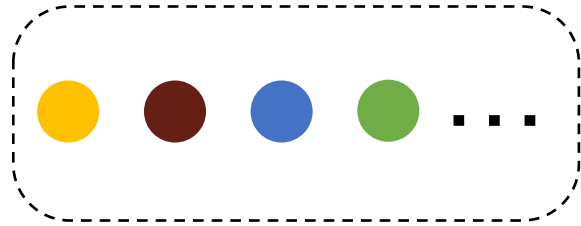
TASO Workflow

Operator Specifications → Graph Subst. Generator → Candidate Substitutions → Graph Subst. Verifier → Verified Substitutions → Graph Optimizer

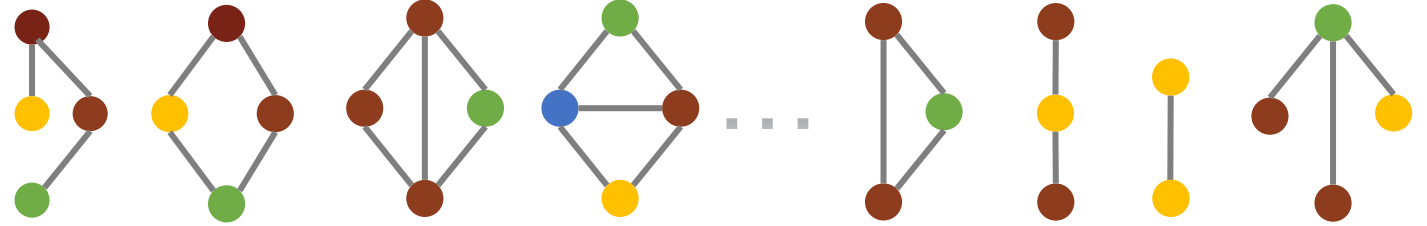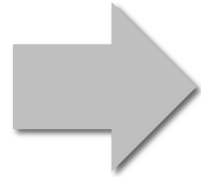Input Comp. Graph → Graph Optimizer → Optimized Comp. Graph

# Graph Substitution Generator

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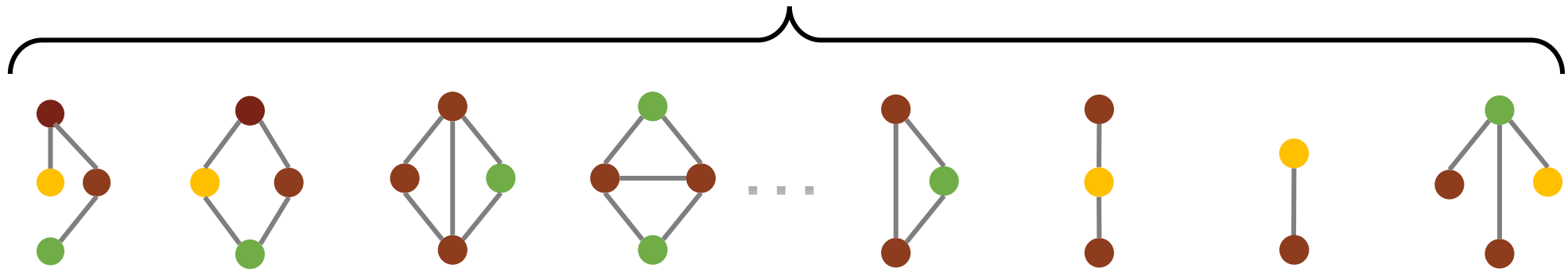


Operators supported by hardware backend

# Graph Substitution Generator

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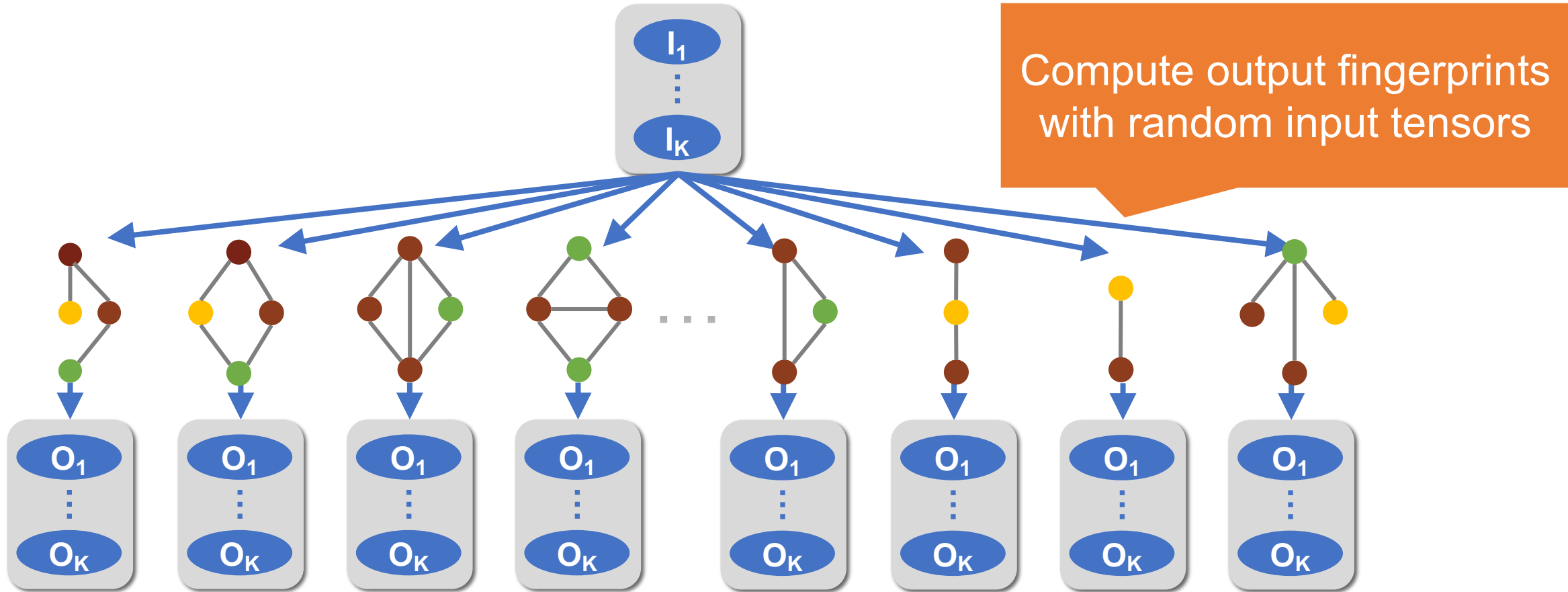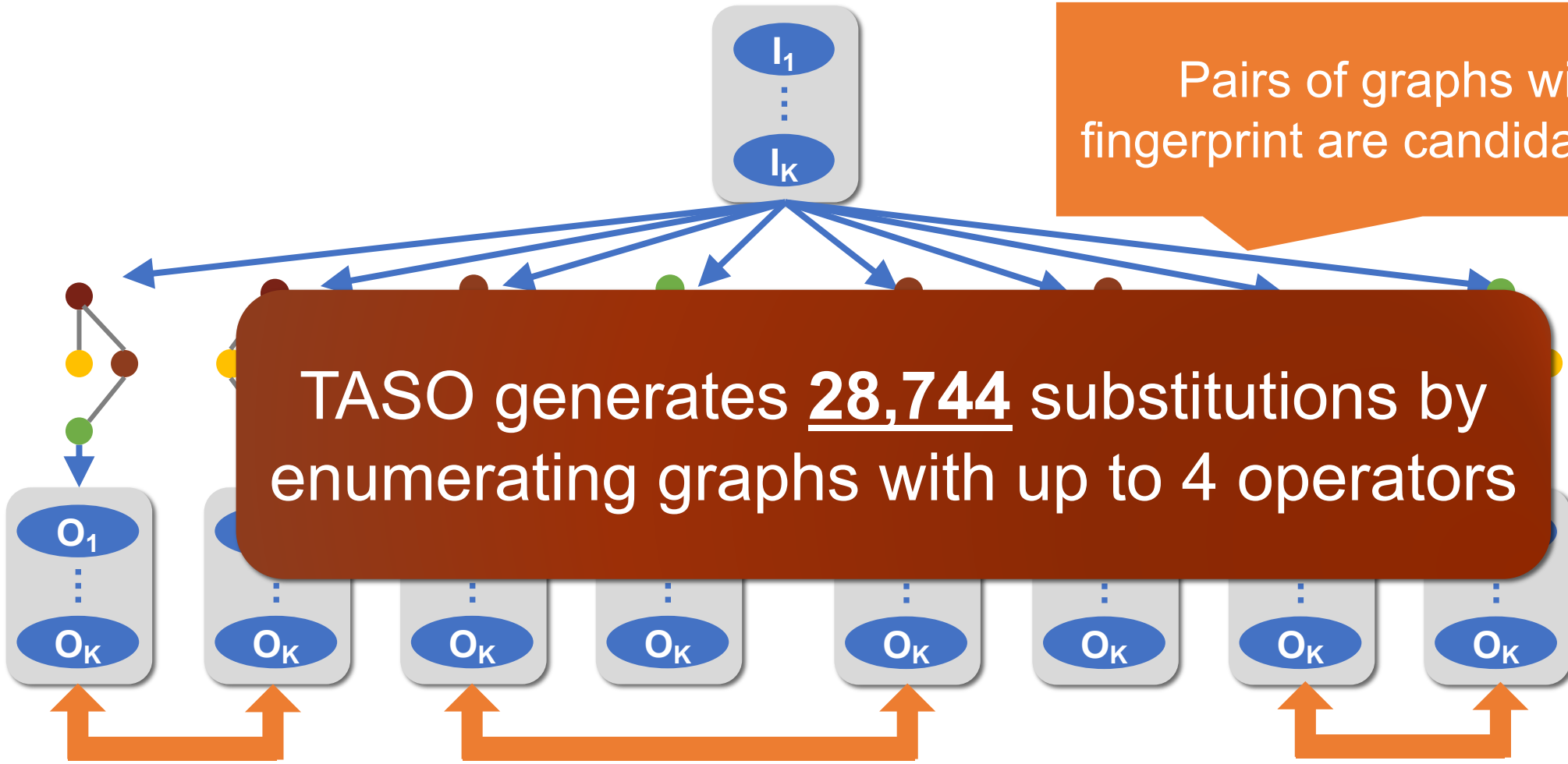


A substitution = a pair of equivalent graphs

**Explicitly considering all pairs does not scale**

# Graph Substitution Generator

Compute output fingerprints with random input tensors

# Graph Substitution Generator

Pairs of graphs with identical fingerprint are candidate substitutions

TASO generates **28,744** substitutions by enumerating graphs with up to 4 operators

# Pruning Redundant Substitutions

**28,744** substitutions

↓

**Input Tensor Renaming**

↓

**17,346** substitutions



X
matmul
matmul
A  B

*A x (B x A)*

X
matmul
matmul
A  B

*(A x B) x A*

X
matmul
matmul
A  B  C

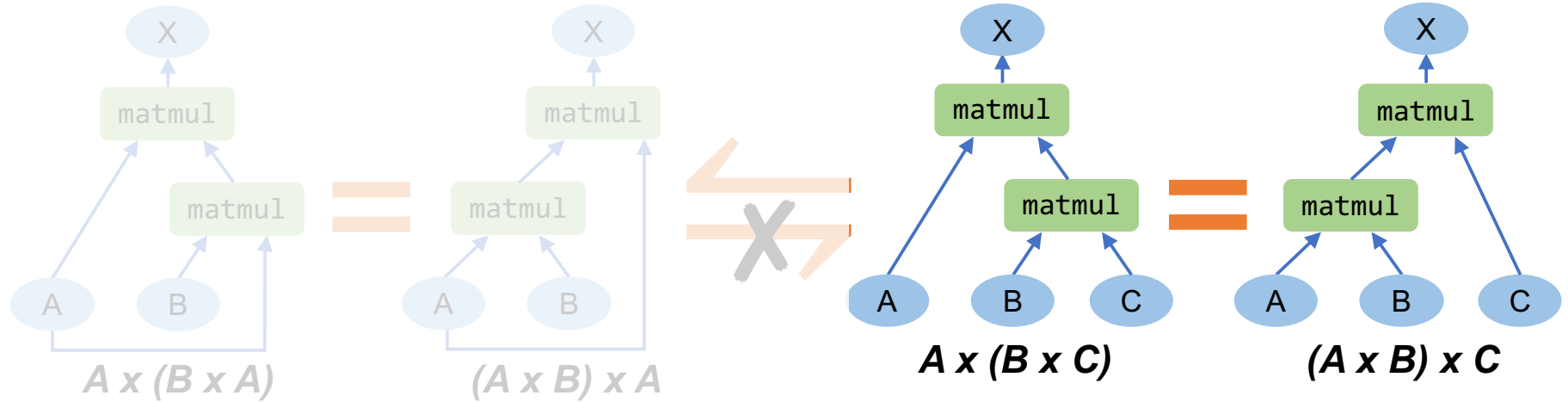*A x (B x C)*

X
matmul
matmul
A  B  C

*(A x B) x C*

# Pruning Redundant Substitutions

**28,744 substitutions**

**Input Tensor Renaming**

**17,346 substitutions**

**Common Subgraphs**

**743 substitutions**



*A + (B x C)*   =   *(B x C) + A*

*(A + B) x C*   =   *(B + A) x C*

*A + B*   =   *B + A*

# Graph Substitution Verifier

*Candidate* Substitutions

**Graph Subst. Verifier**

*Verified* Substitutions

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

*Operator Specifications*

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

23

# Verification Workflow

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

*(Conv(x, w₁), Conv (x, w₂))*

*Split(Conv(x, Concat(w₁, w₂)))*

**Automated Theorem Prover**

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*

24

# Verification Effort

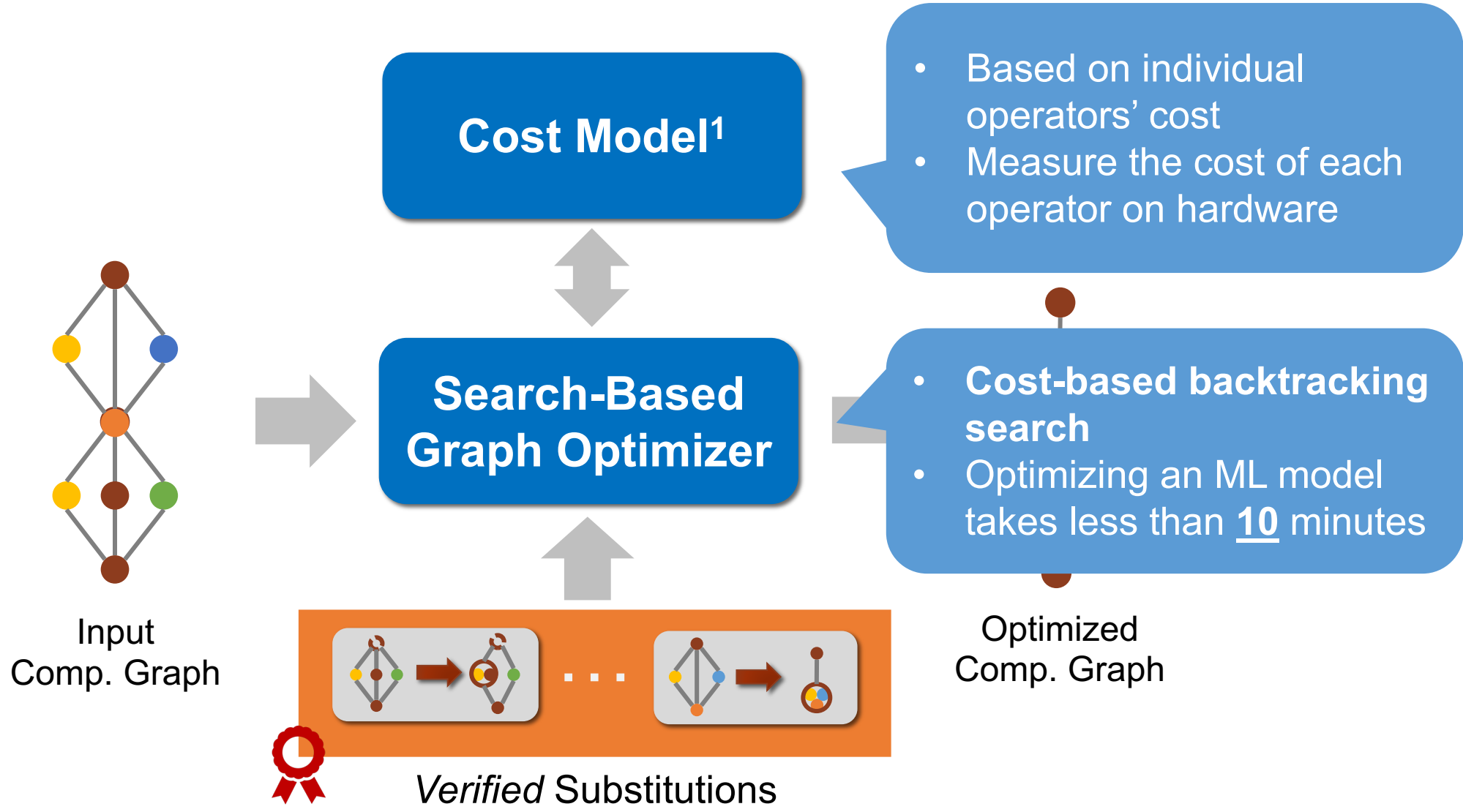| Operator Property | Comment |
|---|---|
| $\forall x, y, z.\ \text{ewadd}(x, \text{ewadd}(y, z)) = \text{ewadd}(\text{ewadd}(x, y), z)$ | ewadd is associative |
| $\forall x, y.\ \text{ewadd}(x, y) = \text{ewadd}(y, x)$ | ewadd is commutative |
| $\forall x, y, z.\ \text{ewmul}(x, \text{ewmul}(y, z)) = \text{ewmul}(\text{ewmul}(x, y), z)$ | 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 |

**TASO generates all _743_ substitutions in 5 minutes, and verifies them against _43_ operator properties in 10 minutes**

| | |
|---|---|
| $\forall s, p, x, y, w.\ \text{smul}(\text{conv}(s, p, A_{\text{none}}, x, y), w) = \text{conv}(s, p, A_{\text{none}}, \text{smul}(x, w), y)$ | 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 |

**Supporting a new operator requires _a few hours_ of human effort to specify its properties**

| | |
|---|---|
| $\forall a, x, y.\ \text{split}_0(a, \text{concat}(a, x, y)) = x$ | split definition |

**Operator specifications in TASO ≈ _1,400_ LOC**
**Manual graph optimizations in TensorFlow ≈ _53,000_ LOC**

| | |
|---|---|
| $\forall s, p, x, y, z, w.\ \text{conv}(s, p, A_{\text{none}}, \text{concat}(1, x, z), \text{concat}(1, y, w)) =$ | concatenation and conv. |
| $\qquad\qquad\qquad \text{ewadd}(\text{conv}(s, p, A_{\text{none}}, x, y), \text{conv}(s, p, A_{\text{none}}, z, w))$ | |
| $\forall k, s, p, x, y.\ \text{concat}(1, \text{pool}_{\text{avg}}(k, s, p, x), \text{pool}_{\text{avg}}(k, s, p, y)) = \text{pool}_{\text{avg}}(k, s, p, \text{concat}(1, x, y))$ | concatenation and pooling |
| $\forall k, s, p, x, y.\ \text{concat}(0, \text{pool}_{\text{max}}(k, s, p, x), \text{pool}_{\text{max}}(k, s, p, y)) = \text{pool}_{\text{max}}(k, s, p, \text{concat}(0, x, y))$ | concatenation and pooling |
| $\forall k, s, p, x, y.\ \text{concat}(1, \text{pool}_{\text{max}}(k, s, p, x), \text{pool}_{\text{max}}(k, s, p, y)) = \text{pool}_{\text{max}}(k, s, p, \text{concat}(1, x, y))$ | concatenation and pooling |

# Search-Based Graph Optimizer

**Cost Model[1]**

- Based on individual operators' cost
- Measure the cost of each operator on hardware

**Search-Based Graph Optimizer**

- **Cost-based backtracking search**
- Optimizing an ML model takes less than **10** minutes

Input Comp. Graph

Optimized Comp. Graph

*Verified* Substitutions

1. Exploring Hidden Dimensions in Parallelizing Convolutional Neural Networks. ICML'18.

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

# Case Study: NASNet



*DWC: depth-wise convolution

28

# Why TASO is a SuperOptimizer?

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

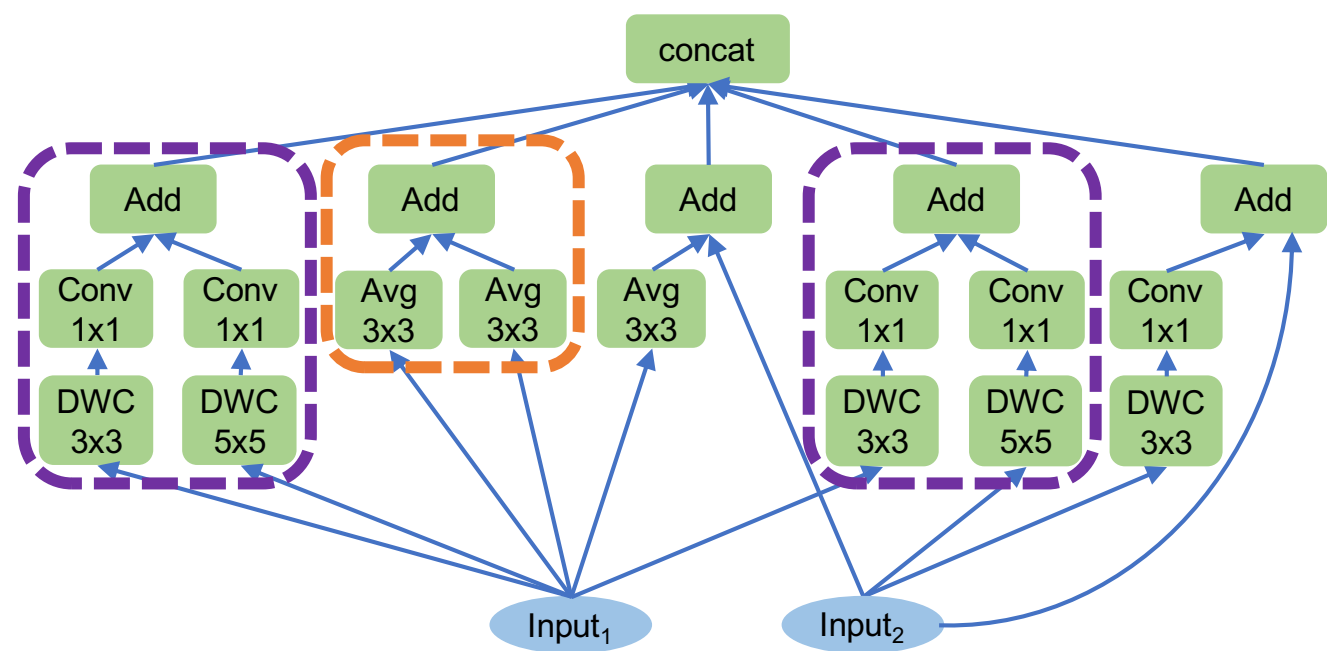Goal: gradually *improve* an input program by greedily applying optimizations

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

# PET:
# Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections

PET: Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections. OSDI'21

30

# Motivation: Fully v.s. Partially 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: Fully v.s. Partially Equivalent Transformations

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

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

**Is it possible to exploit partially equivalent transformations to <u>improve performance</u> while <u>preserving equivalence</u>?**

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



Input Program

Partially Equivalent Transformation

**Incorrect results**

Correcting Results

# Motivating Example



Input Program

Correcting Results

- 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



Input Program → **Mutant Generator** → Mutant Programs → **Mutant Corrector** → Corrected Mutants → **Program Optimizer** → Optimized Program

38

# PET vs TASO

# Key Challenges

1. How to generate partially equivalent transformations?

Superoptimization

2. How to correct them?

Multi-linearity of DNN computations

# Mutant Generator

Superoptimization adopted from TASO[1]



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

Input (Sub)program

**Mutant Generator**

Operators supported by hardware backend

# Mutant Generator

Superoptimization adopted from TASO[1]



Programs with the same input/output shapes are potential mutants

Mutant Generator

Input (Sub)program

Operators supported by hardware backend

Discover both fully and partially equivalent transformations

1. TASO: Optimizing Deep Learning Computation with Automated Generation of Graph Substitutions. SOSP'19.

42

# Challenges: Examine Transformations

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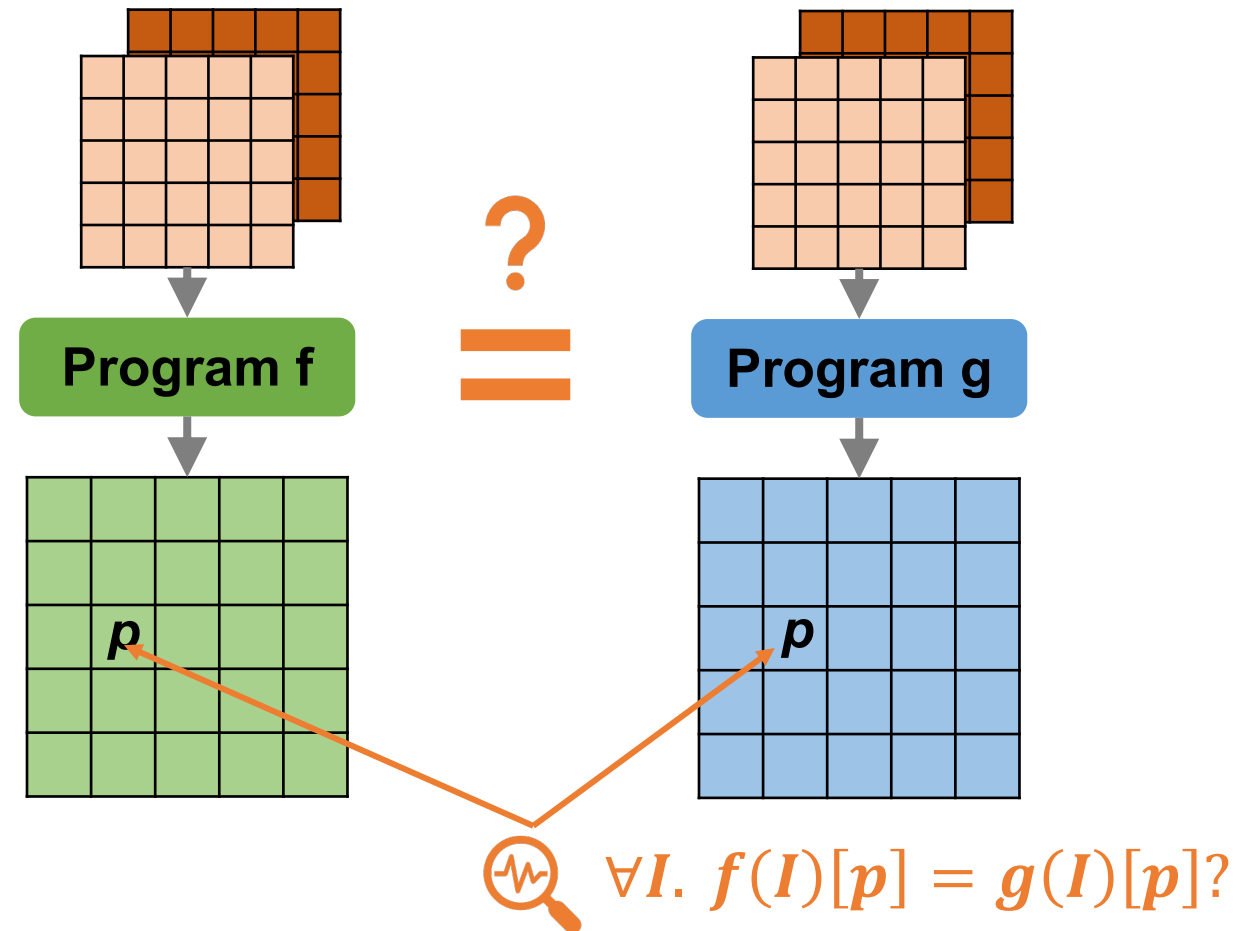2. How to correct the results?

# A Strawman Approach

- **Step 1**: Explicitly consider all output positions (m positions)

- **Step 2**: For each position **p**, examine all possible inputs (n inputs)



$$\forall I. \; f(I)[p] = g(I)[p]?$$

**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
  - $f(I_1, \ldots, X, \ldots, I_n) + f(I_1, \ldots, Y, \ldots, I_n) = f(I_1, \ldots, X + Y, \ldots, I_n)$
  - $\alpha \cdot f(I_1, \ldots, X, \ldots, I_n) = f(I_1, \ldots, \alpha \cdot X, \ldots, I_n)$

- DNN computation = MLTP + non-linear activations

Majority of the computation

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

**MLTP**

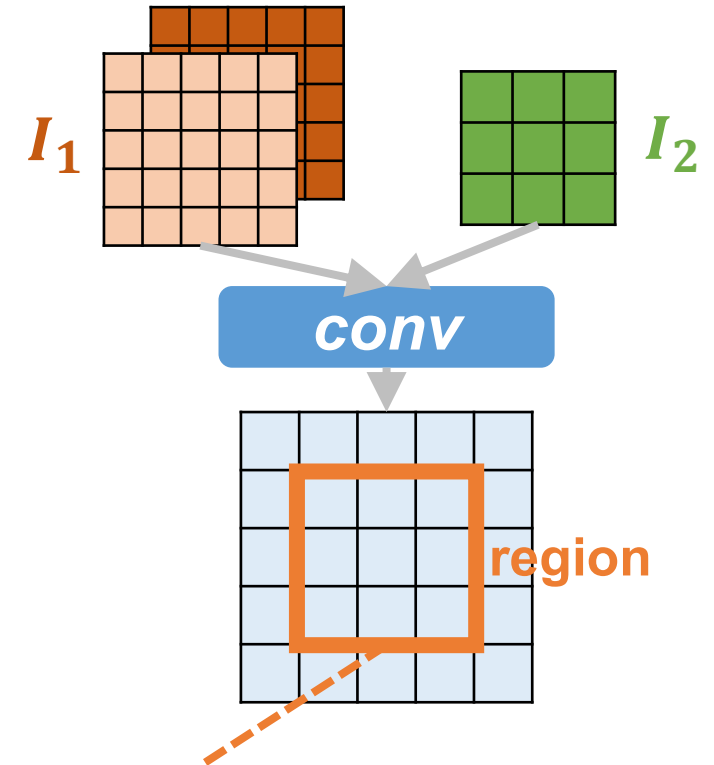**O(1) examinations in PET's approach**

# Insight #1: 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) → O(n)**



$I_1$

$I_2$

*conv*

region

$$conv(c, h, w) = \sum_{d=0}^{D-1} \sum_{x=-1}^{1} \sum_{y=-1}^{1} \begin{array}{l} I_1(d, h+x, w+y) \\ \times I_2(d, c, x, y) \end{array}$$
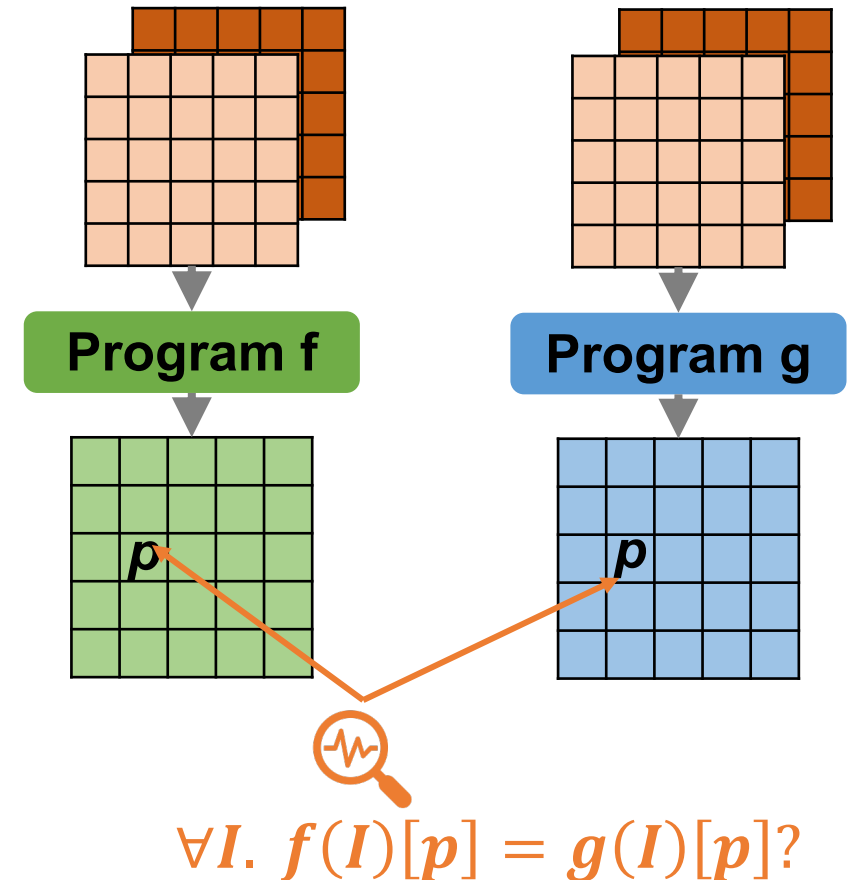
**Summation interval**

# Insight #2: 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 $(\frac{1}{2^{31}})^t$
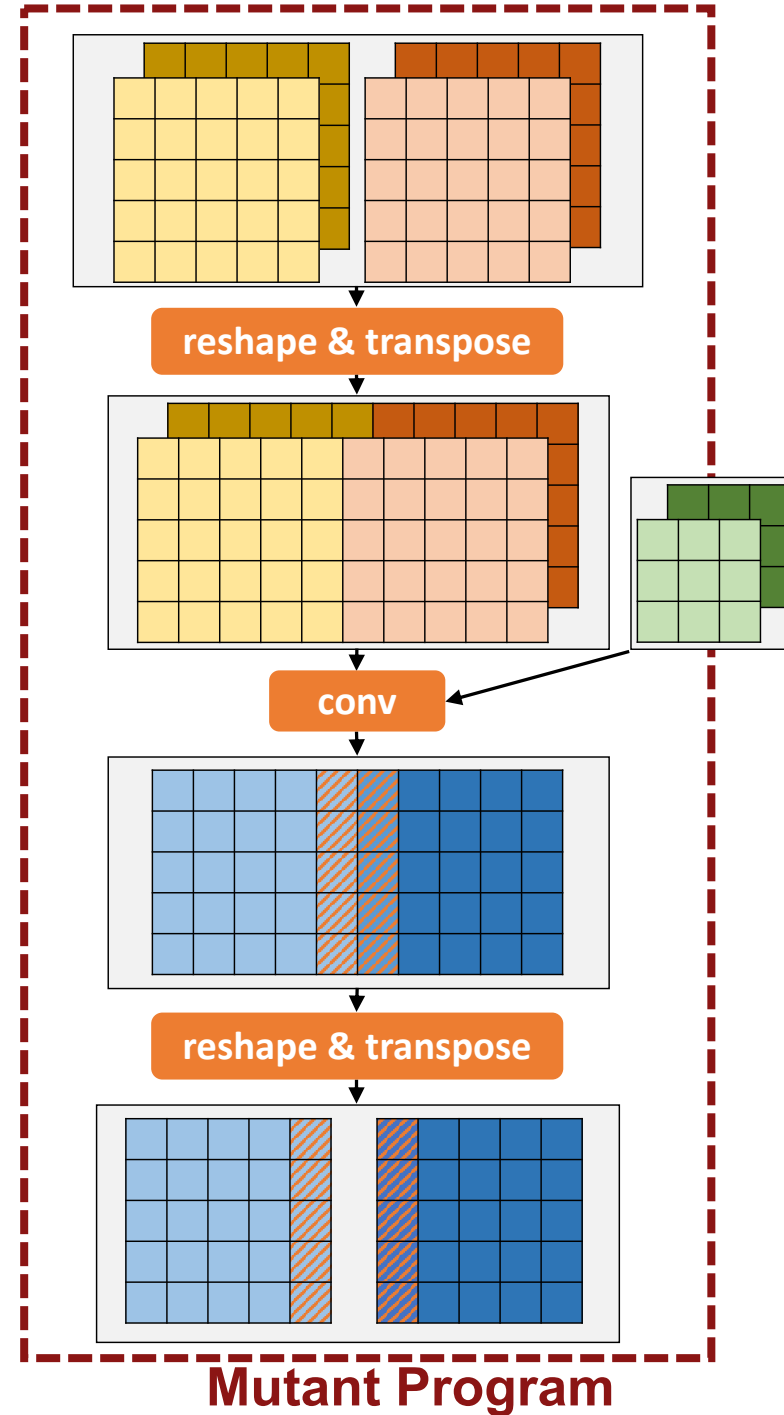
Run $t$ random tests for each position $p$

**Complexity**: **O(n) → O($t$) = O(1)**



$\forall I.\ f(I)[p] = g(I)[p]?$

# Mutant Corrector

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



**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 Program**

**Correction Kernel**

49

# Mutant Corrector

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

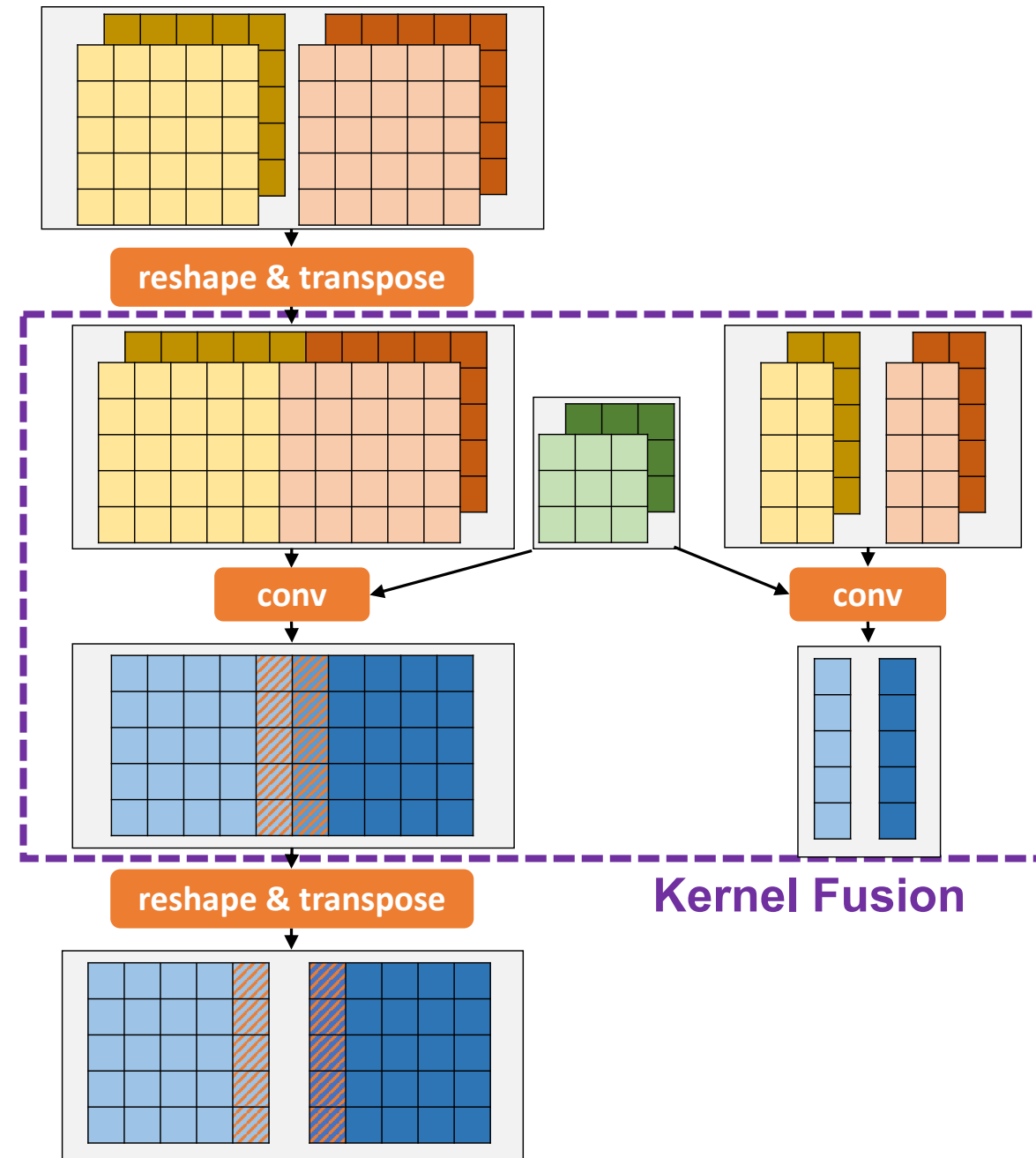**Step 1**: recompute the incorrect outputs using the original program

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

Correction introduces less than **1%** overhead

# Program Optimizer

Input Program

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

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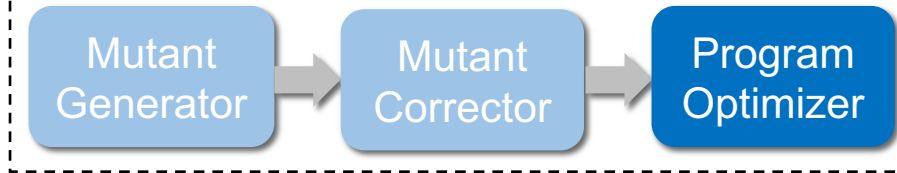- Constant folding
- Redundancy elimination

**Search-Based Program Optimizer**

MLTP

**Mutant Generator & Corrector**
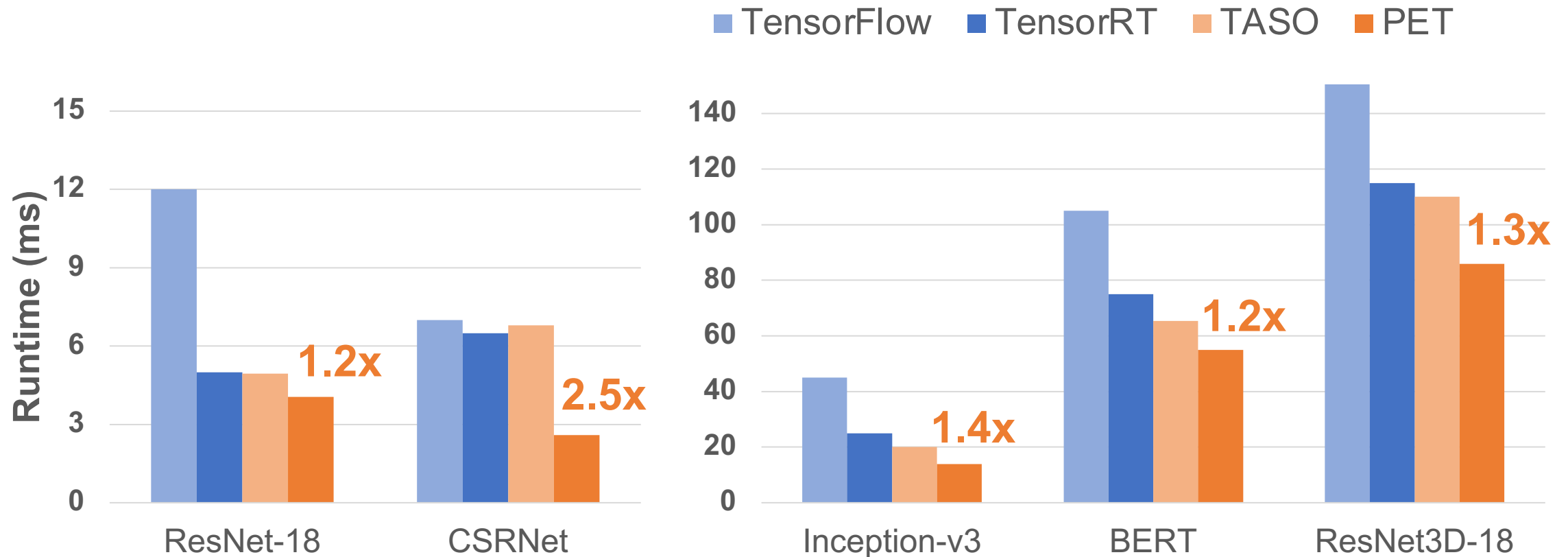
Mutants w/ Corrections

Optimized Program

51

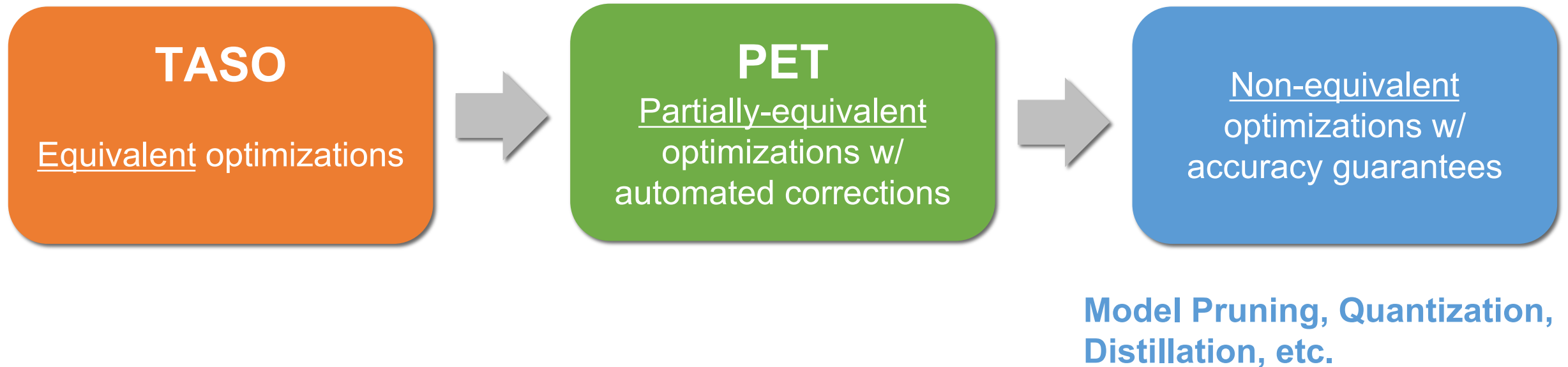# End-to-end Inference Performance (Nvidia V100 GPU)



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

# Recap: 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

**TASO**
Equivalent optimizations

**PET**
Partially-equivalent
optimizations w/
automated corrections

Non-equivalent
optimizations w/
accuracy guarantees

**Model Pruning, Quantization,
Distillation, etc.**

55

# 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?