15-779: Lecture 10

Graph-Level Optimizations

Zhihao Jia

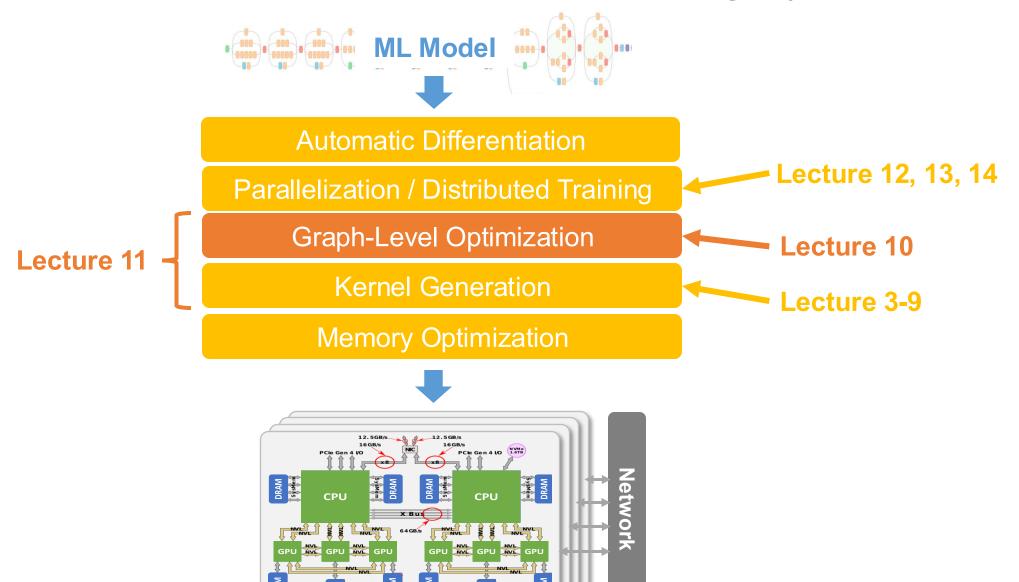
Carnegie Mellon University

Final Project Logistics

- Each group has 1-3 members
 - More group members -> larger project scope
- Potential project ideas & past project reports available on the website
- Proposal (1 page): due 10/10 (next Friday)
- Intermediate check-in meeting: 11/7 during the regular lecture slot
- Final presentations: 12/3 and 12/5
- Final project report: due 12/15

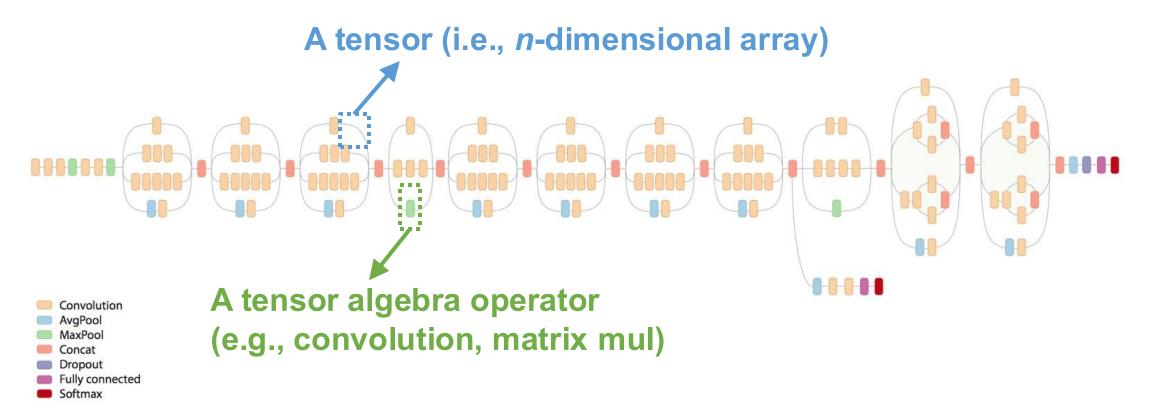


Recap: An Overview of Deep Learning Systems

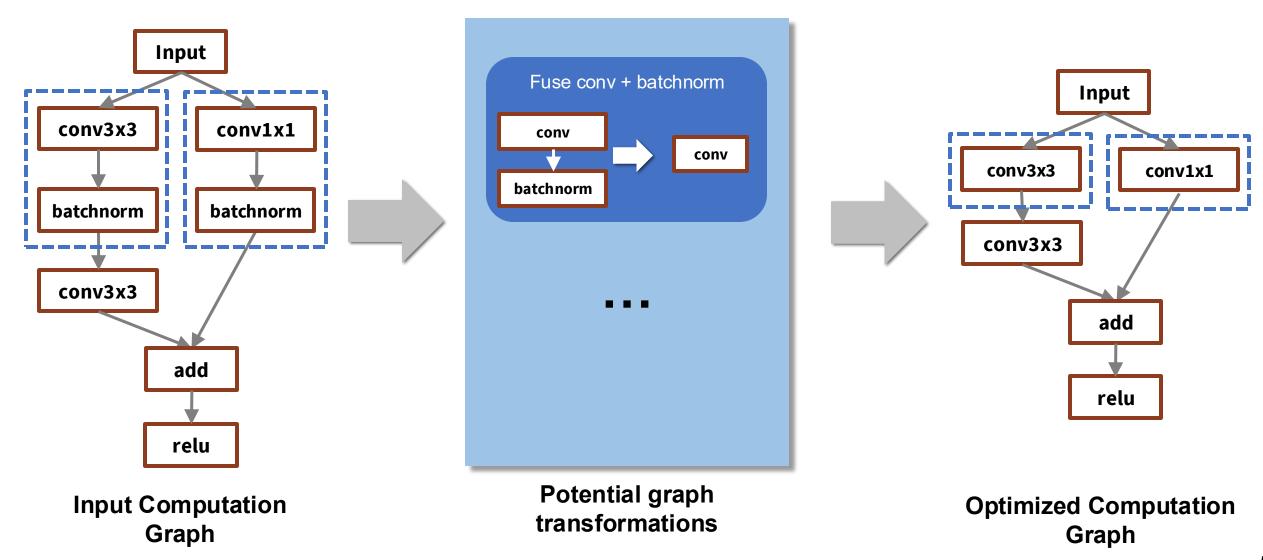


Recap: Deep Neural Network

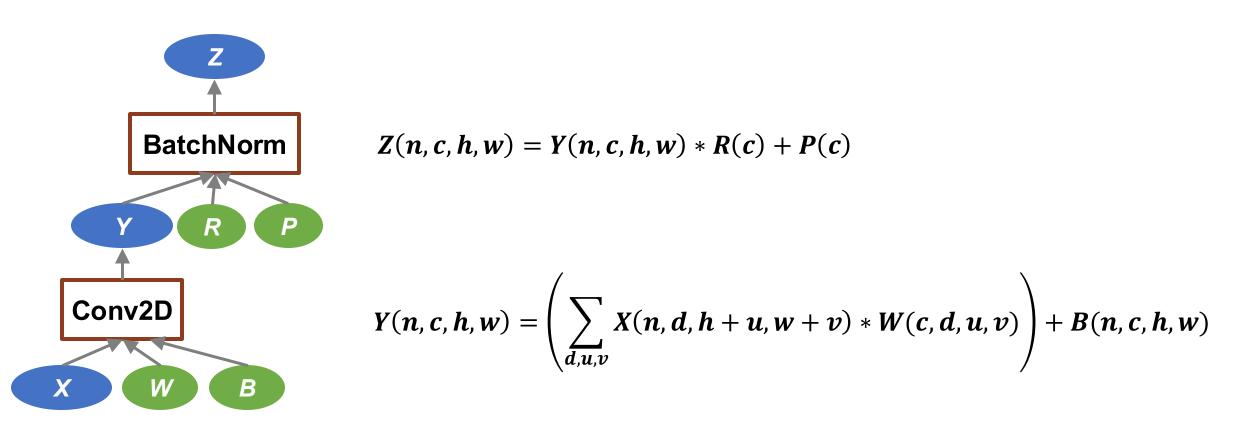
 Collection of simple trainable mathematical units that work together to solve complicated tasks



Graph-Level Optimizations

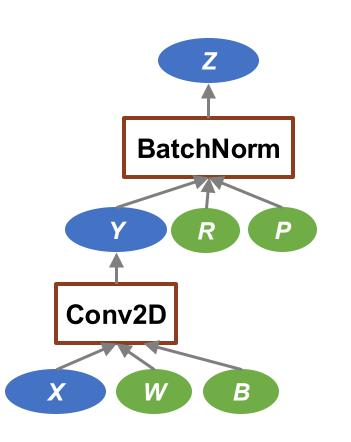


Example: Fusing Convolution and Batch Normalization

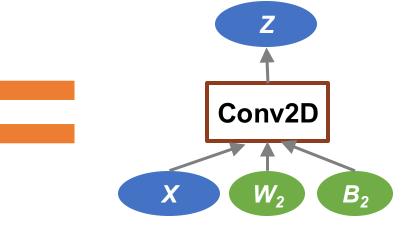


W, B, R, P are constant pre-trained weights

Fusing Conv and BatchNorm



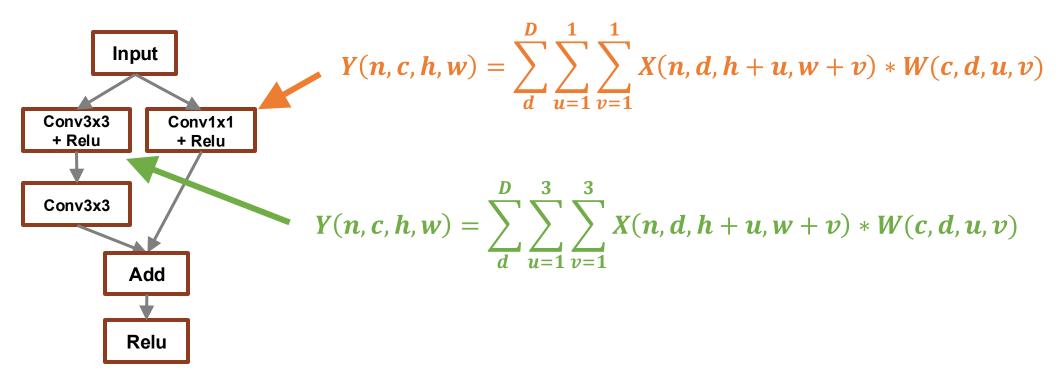
$$Z(n, c, h, w) = \left(\sum_{d,u,v} X(n, d, h + u, w + v) * W_2(c, d, u, v)\right) + B_2(n, c, h, w)$$



$$W_2(n,c,h,w) = W(n,c,h,w) * R(c)$$

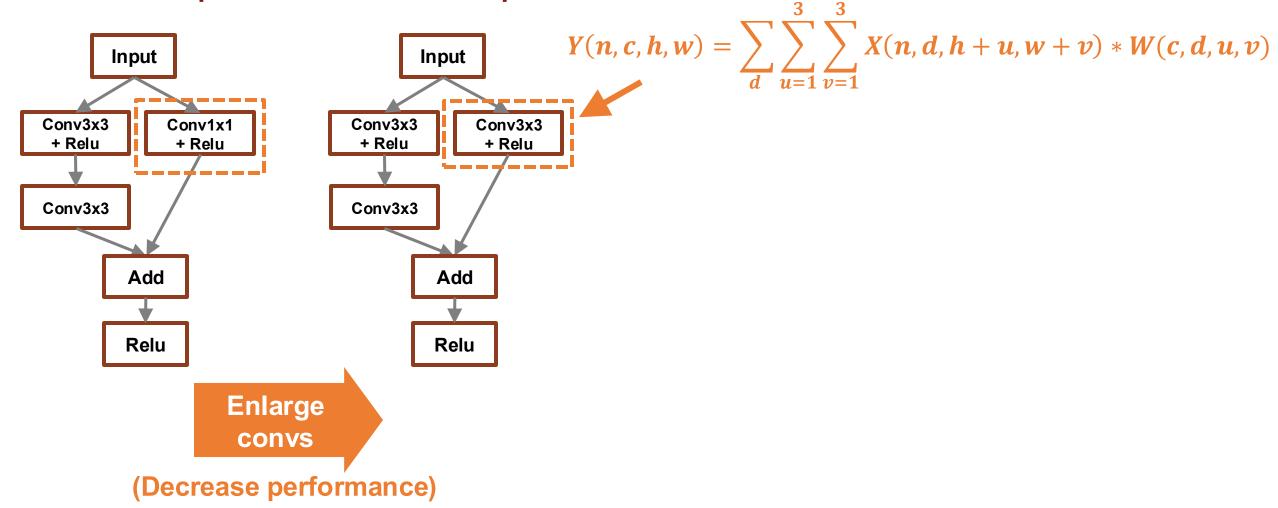
$$B_2(n, c, h, w) = B(n, c, h, w) * R(c) + P(c)$$

Recap: Resnet Example

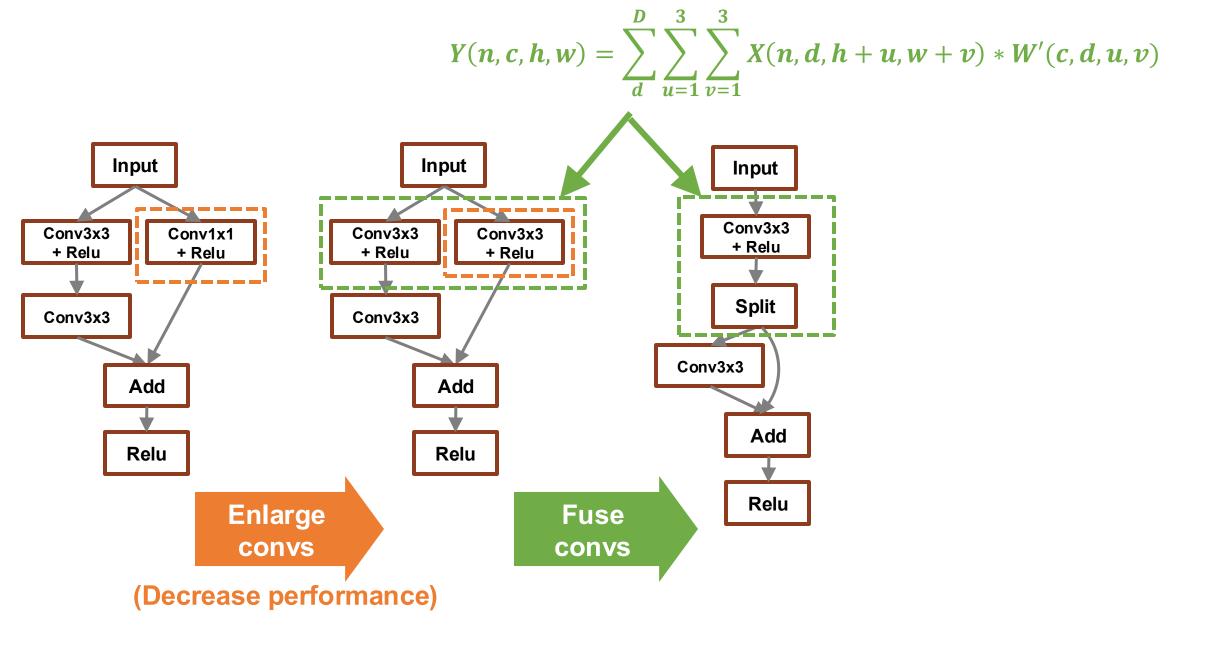


^{*} Kaiming He. et al. Deep Residual Learning for Image Recognition, 2015

Recap: Resnet Example

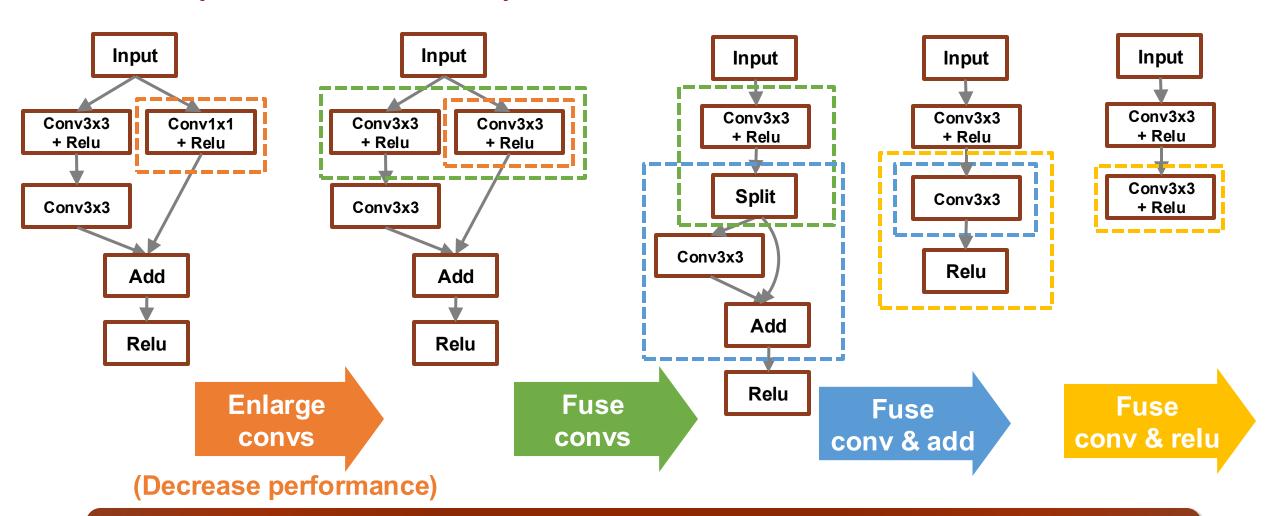


^{*} Kaiming He. et al. Deep Residual Learning for Image Recognition, 2015



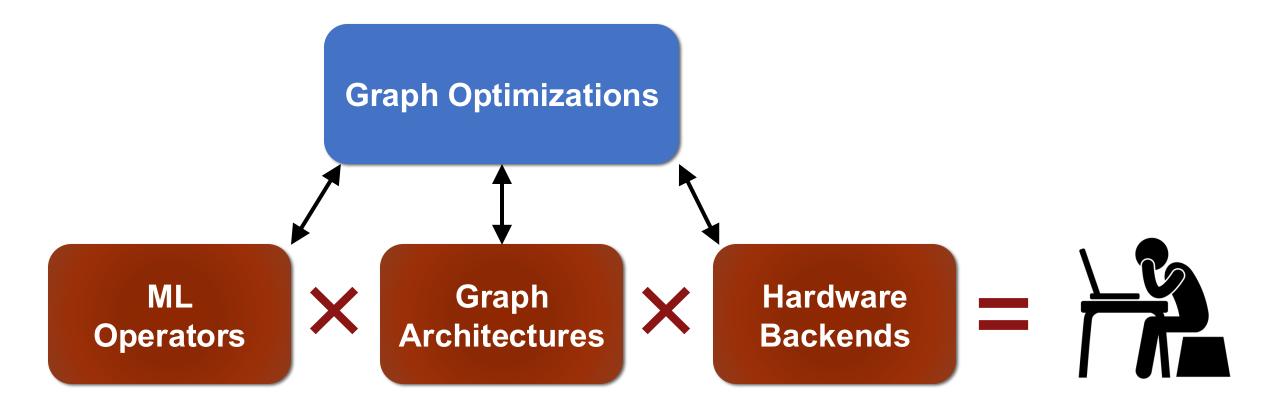
^{*} Kaiming He. et al. Deep Residual Learning for Image Recognition, 2015

Recap: Resnet Example



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

Challenge of Graph Optimizations for ML



Infeasible to manually design graph optimizations for all cases

This Lecture

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

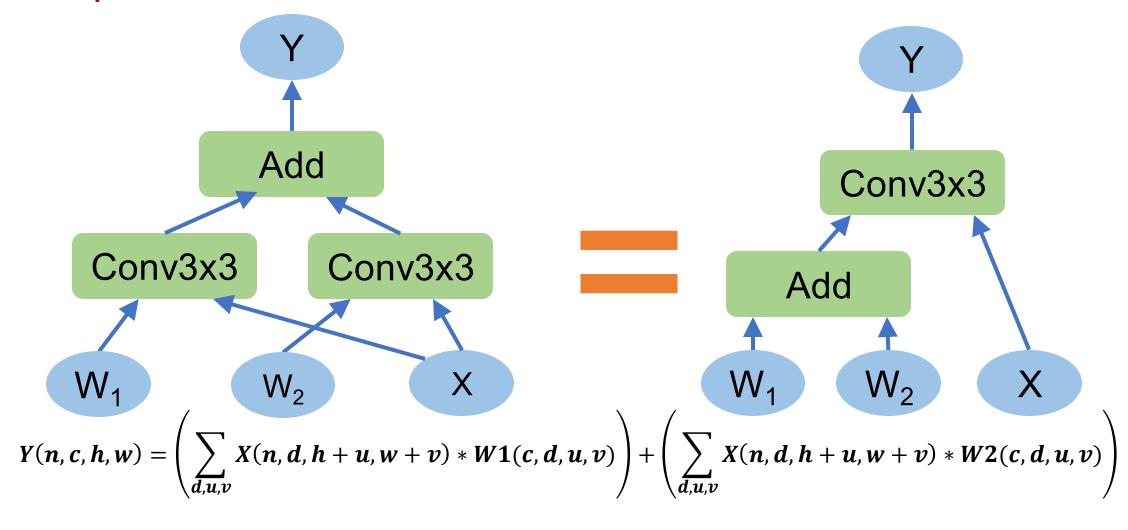
TASO: Optimizing Deep Learning with Automatic Generation of Graph Substitutions

TASO: Tensor Algebra SuperOptimizer

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

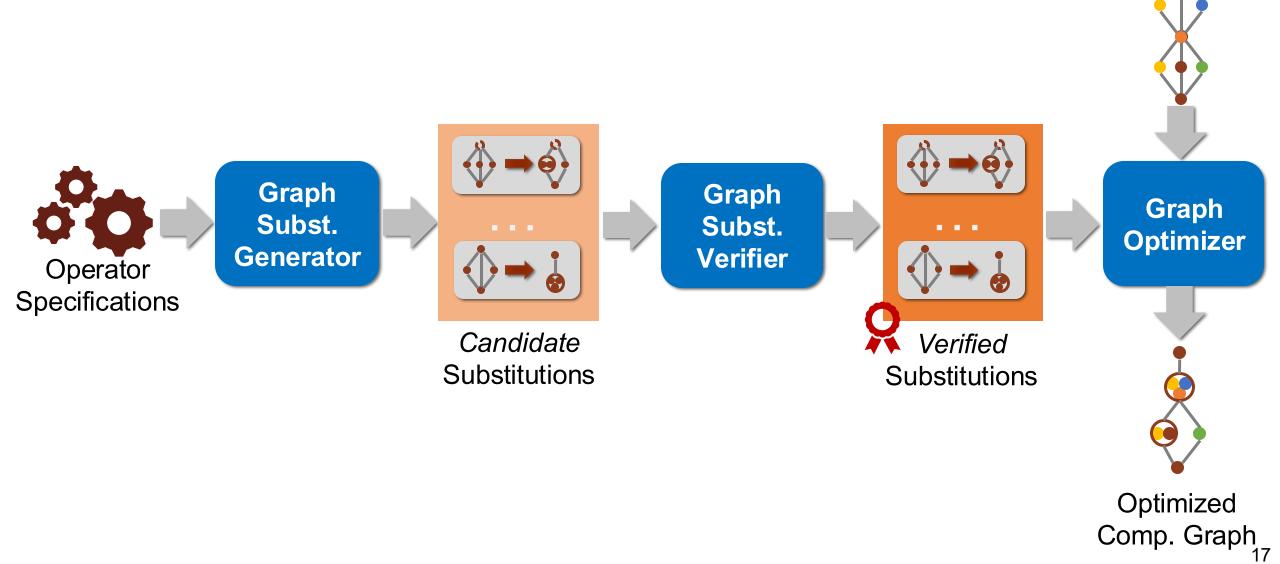
- Less engineering effort: <u>53,000</u> LOC for manual graph optimizations in TensorFlow → <u>1,400</u> LOC in TASO
- Better performance: outperform existing optimizers by up to 3x
- Stronger correctness: formally verify all generated substitutions

Graph Substitution



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

TASO Workflow

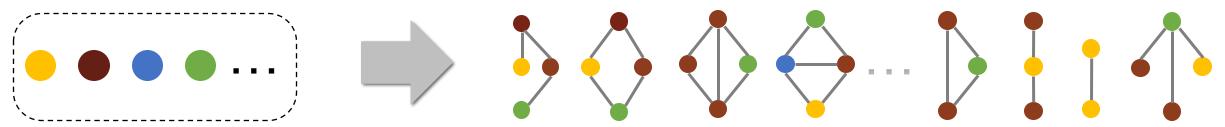


Input

Comp. Graph



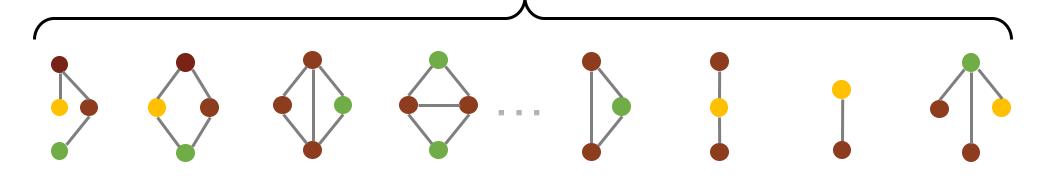
Enumerate <u>all possible</u> graphs up to a fixed size using available operators



Operators supported by hardware backend



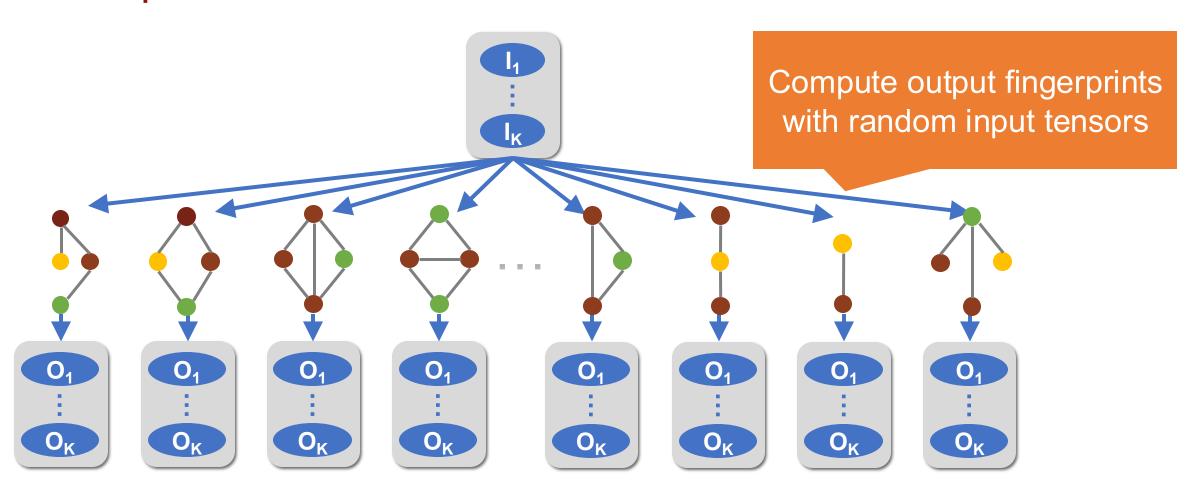
66M graphs with up to **4** operators



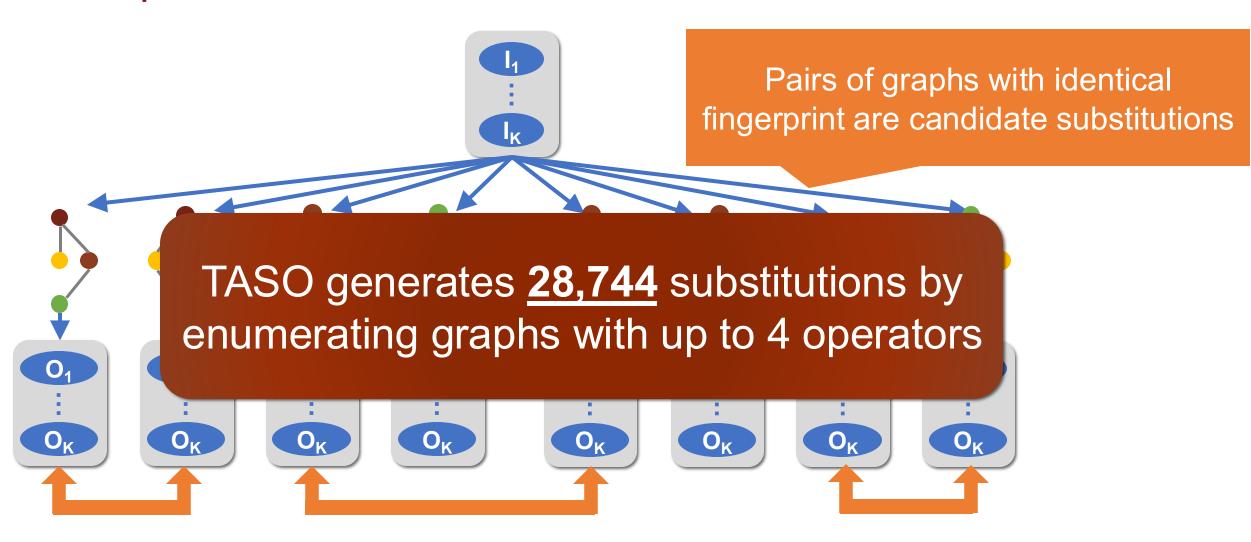
A substitution = a pair of equivalent graphs

Explicitly considering all pairs does not scale





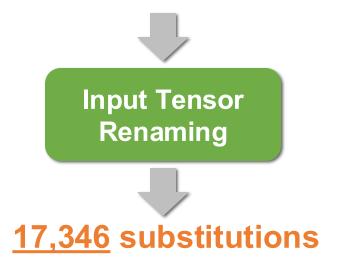


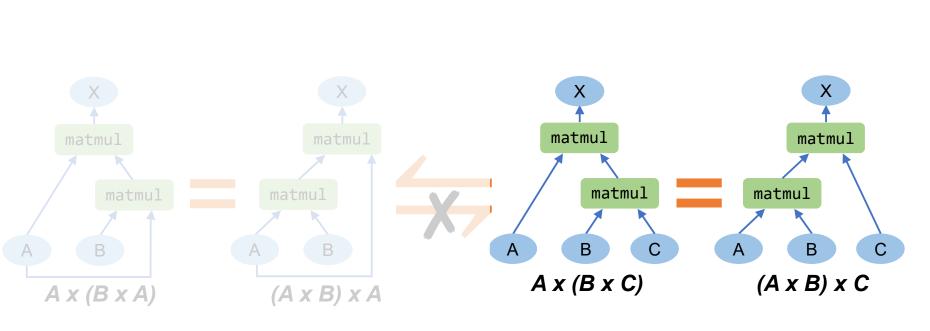




Pruning Redundant Substitutions

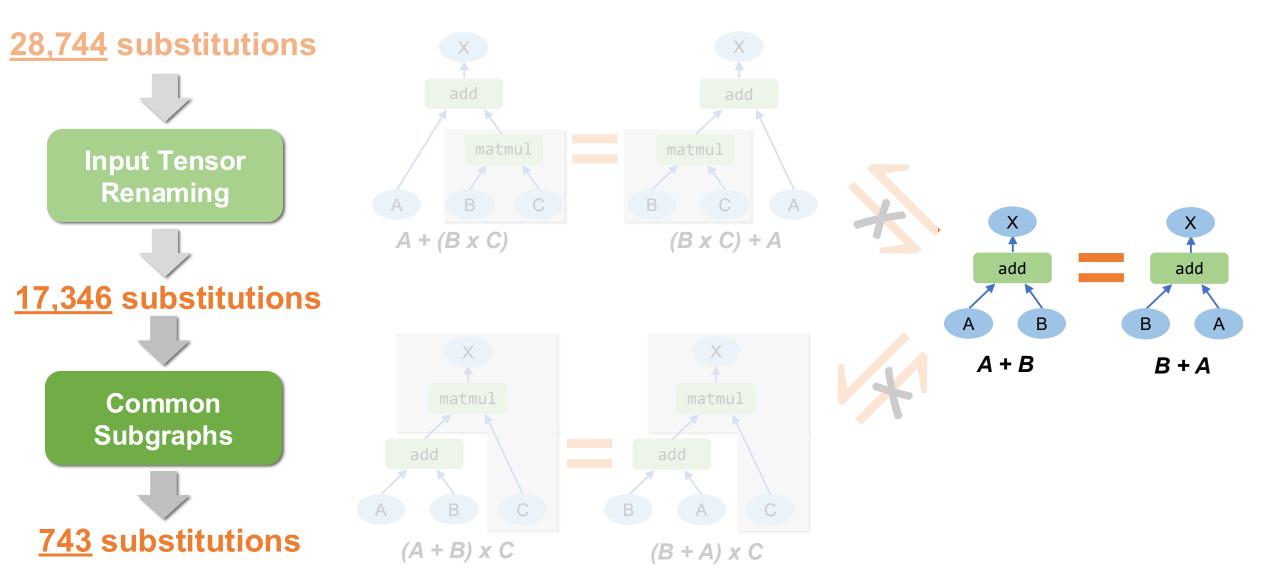
28,744 substitutions





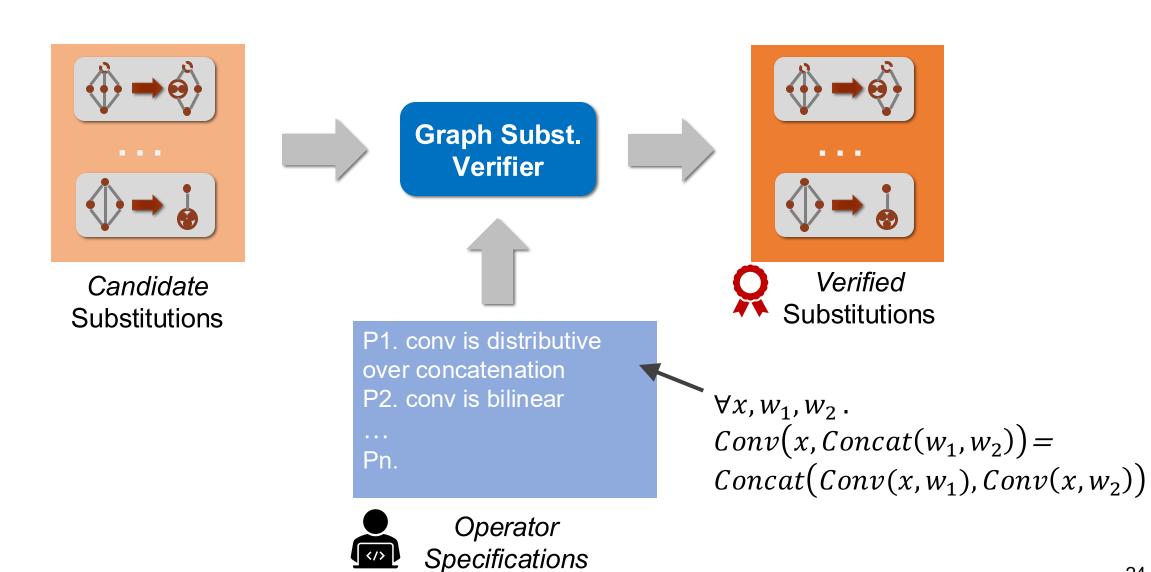


Pruning Redundant Substitutions





Graph Substitution Verifier



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)
```

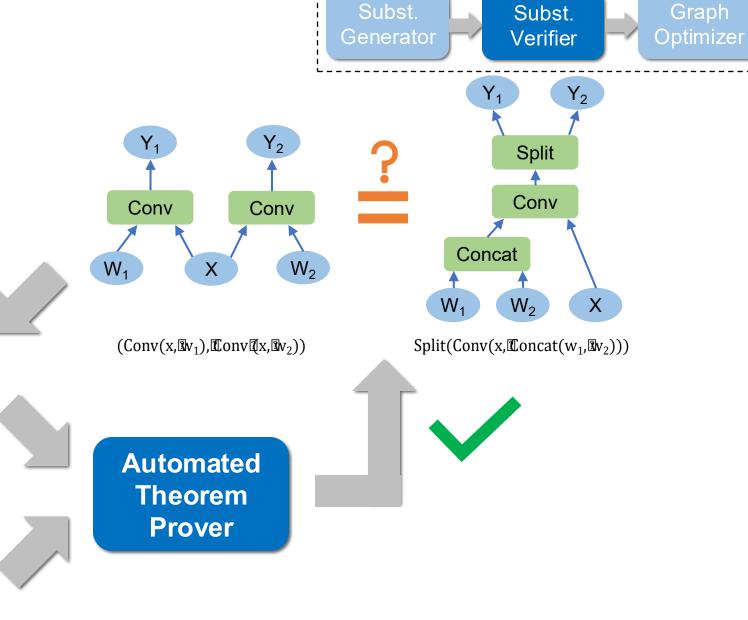
```
P1. \forall x, w_1, w_2.

Conv(x, Concat(w_1, w_2)) =

Concat(Conv(x, w_1), Conv(x, w_2))

P2. ...
```

Operator Specifications



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. ewadd(x, y) = 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 <u>743</u> substitutions in 5 minutes, and verifies them against <u>43</u> 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)

\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))
```

Supporting a new operator requires <u>a few hours</u> of human effort to specify its properties

 $\forall a, x, y.$ split $_0(a,$ concat(a, x, y)) = x

```
Operator specifications in TASO ≈ <u>1,400</u> LOC
```

Manual graph optimizations in TensorFlow ≈ <u>53,000</u> LOC

```
 \forall s, p, x, y, z, w. \ \operatorname{conv}(s, p, c, x, y), \operatorname{conv}(s, p, c, x, z)) = \operatorname{conv}(s, p, c, x, c) \\ \forall s, p, x, y, z, w. \ \operatorname{conv}(s, p, \mathsf{A}_{\mathsf{none}}, \mathsf{concat}(1, x, z), \mathsf{concat}(1, y, w)) = \\ & \qquad \qquad \operatorname{ewadd}(\operatorname{conv}(s, p, \mathsf{A}_{\mathsf{none}}, x, y), \operatorname{conv}(s, p, \mathsf{A}_{\mathsf{none}}, z, w)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{avg}}(k, s, p, x), \operatorname{pool}_{\mathsf{avg}}(k, s, p, y)) = \operatorname{pool}_{\mathsf{avg}}(k, s, p, \mathsf{concat}(1, x, y)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, y)) = \operatorname{pool}_{\mathsf{max}}(k, s, p, \mathsf{concat}(1, x, y)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, y)) = \operatorname{pool}_{\mathsf{max}}(k, s, p, \mathsf{concat}(1, x, y)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, y)) = \operatorname{pool}_{\mathsf{max}}(k, s, p, \mathsf{concat}(1, x, y)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, y)) = \operatorname{pool}_{\mathsf{max}}(k, s, p, \mathsf{concat}(1, x, y)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, y)) = \operatorname{pool}_{\mathsf{max}}(k, s, p, \mathsf{concat}(1, x, y)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x), \operatorname{pool}_{\mathsf{max}}(k, s, p, x)) \\ \forall k, s, p, x, y. \ \operatorname{concat}(1, \operatorname{pool}_{\mathsf{max}}(k, s, p, x),
```

ommutativity
is its own inverse
ommutativity
ommutativity
ommutativity
associative

inear
inear
d transpose
dinear
conv is bilinear

conv is bilinear
inear
ivolution kernel
Arelu applies relu
mmutativity
conv. with Cpool
rnel

split definition

atrix

of concatenation ommutativity ommutativity ommutativity

ommutativity ion and transpose ion and matrix mul. ion and matrix mul.

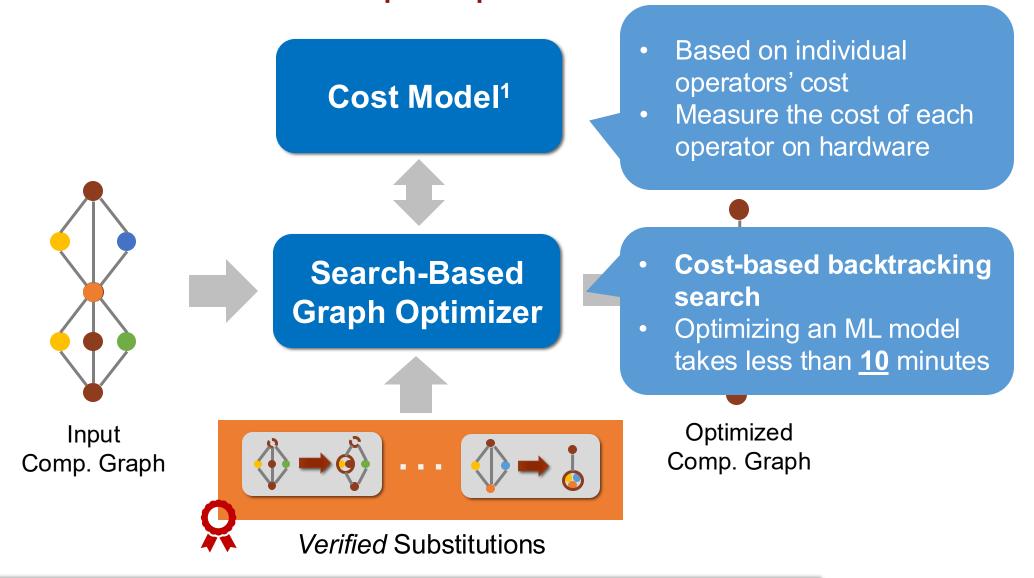
ion and conv.

concarenation and conv.

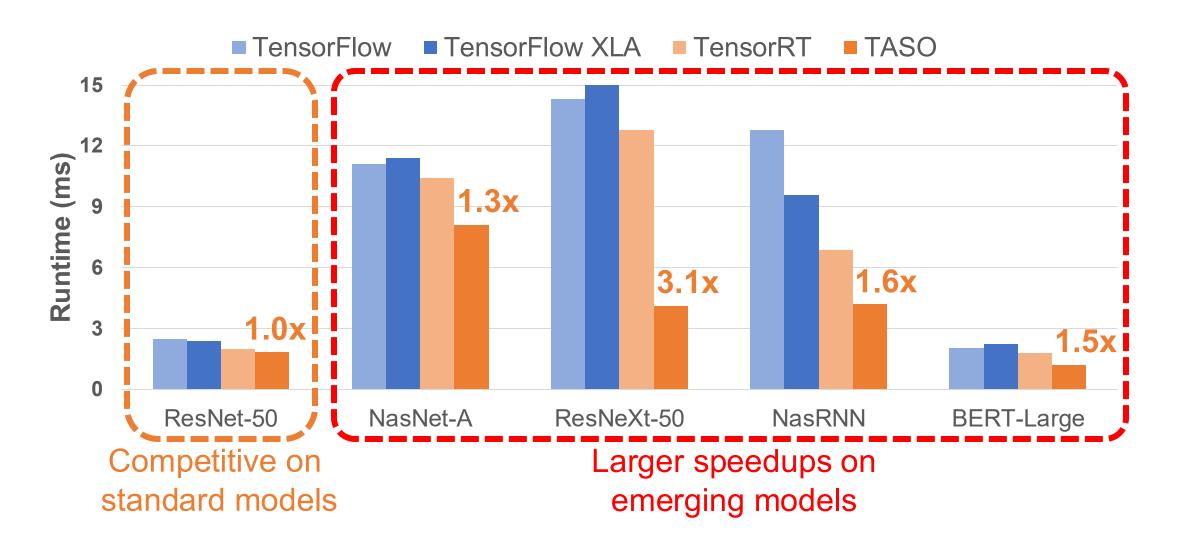
concatenation and conv.

concatenation and pooling concatenation and pooling concatenation and pooling 26

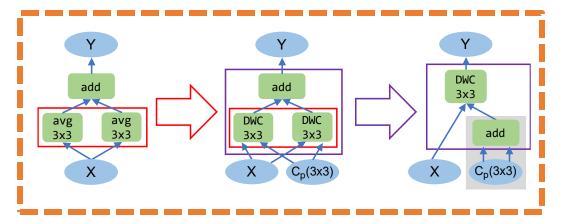
Search-Based Graph Optimizer

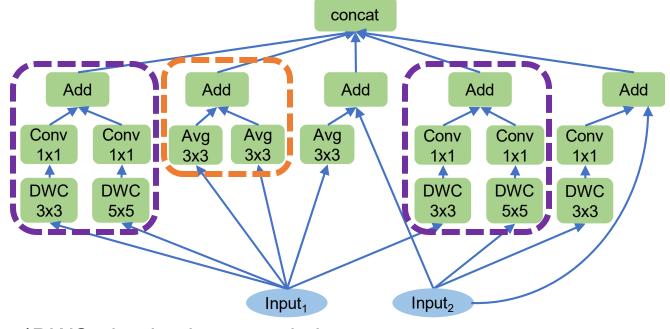


End-to-end Inference Performance (Nvidia V100 GPU)

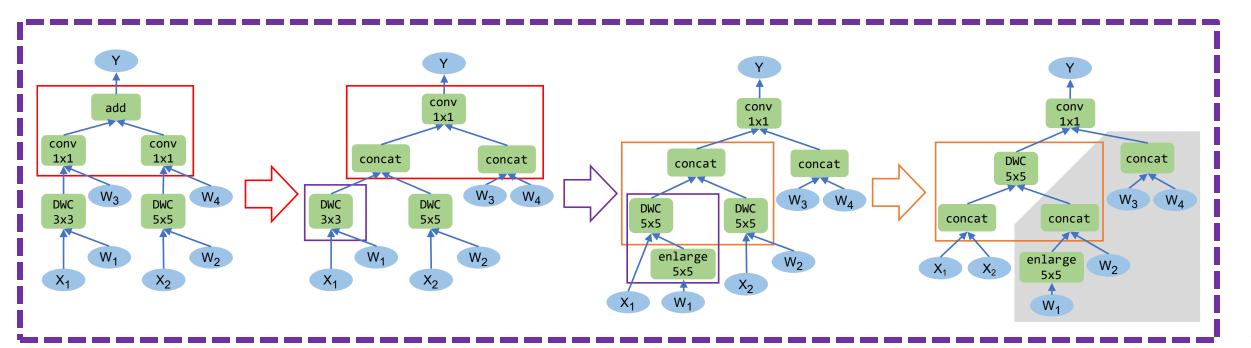


Case Study: NASNet





*DWC: depth-wise convolution



Why TASO is a SuperOptimizer?

What is the difference between optimizer and super-optimizer?

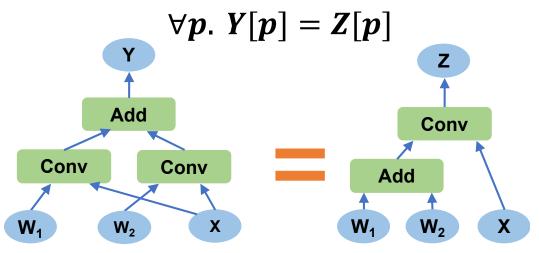
Goal: gradually <u>improve</u> an input program by greedily applying optimizations

Goal: automatically find an optimal program for an input program

PET:

Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections

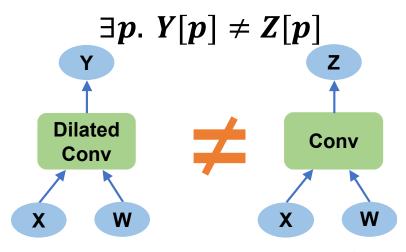
Motivation: Fully v.s. Partially Equivalent Transformations



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

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 improve performance while preserving equivalence?

W₁

N₂

 W_2

7

X



Partially Equivalent Transformations



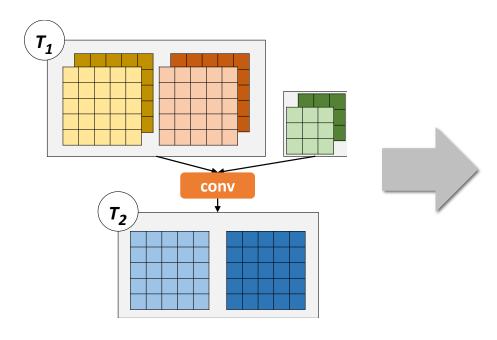
- Faster ML operators
- More efficient tensor layouts
- Hardware-specific optimizations
- Con: potential accuracy loss

Fully Equivalent Transformations

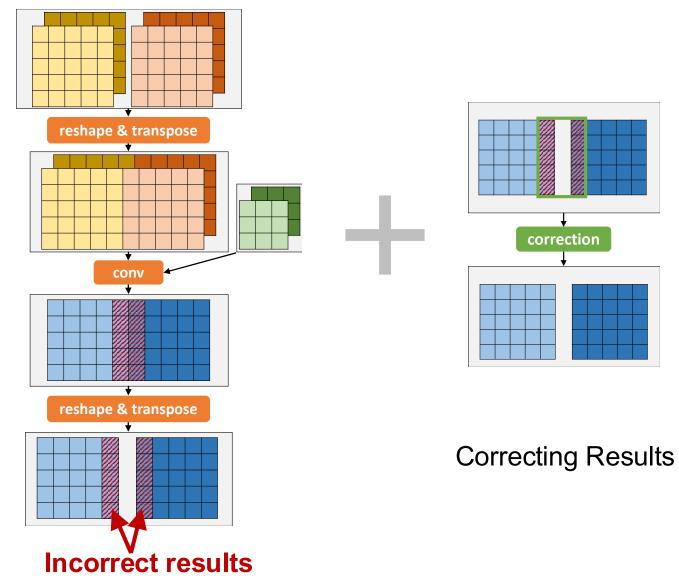
Pro: preserve functionality

Con: miss optimization opportunities

Motivating Example

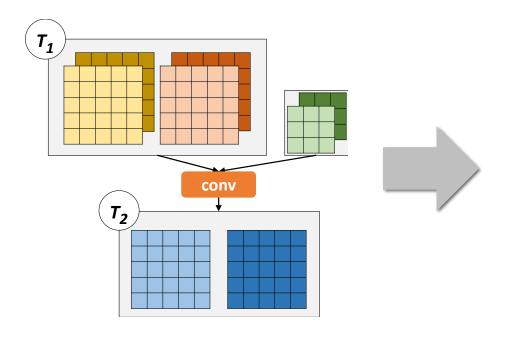


Input Program

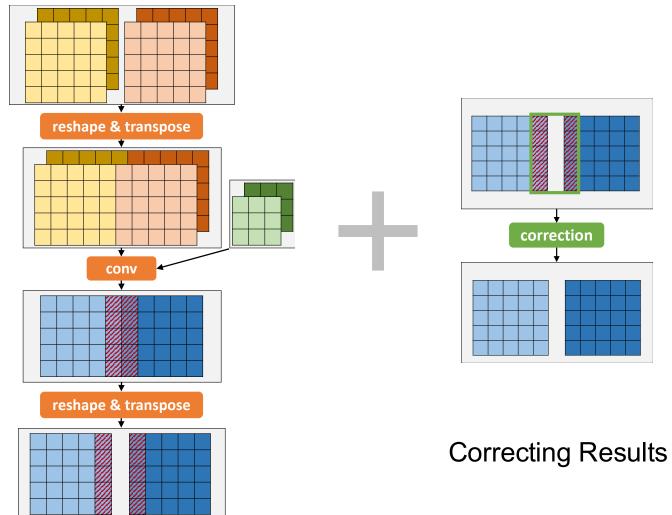


Partially Equivalent Transformation

Motivating Example



Input Program

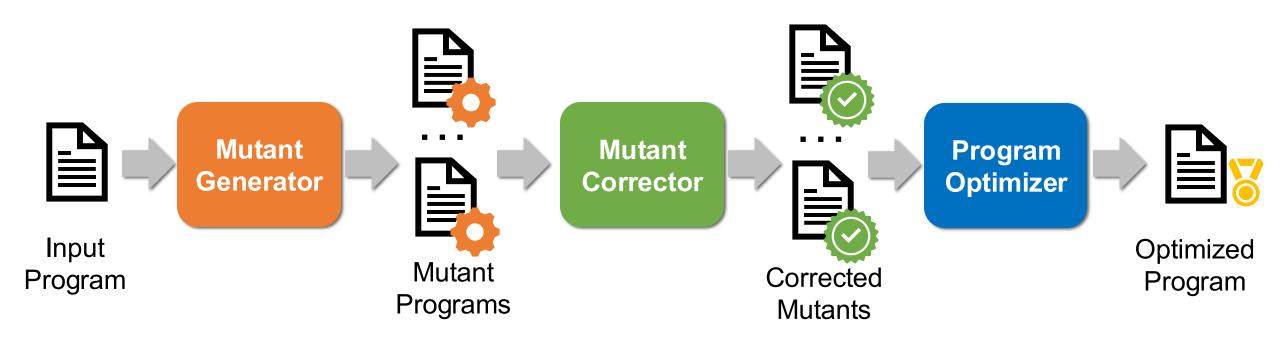


- Transformation and correction lead to 1.2x speedup for ResNet-18
- Correction preserves end-to-end equivalence

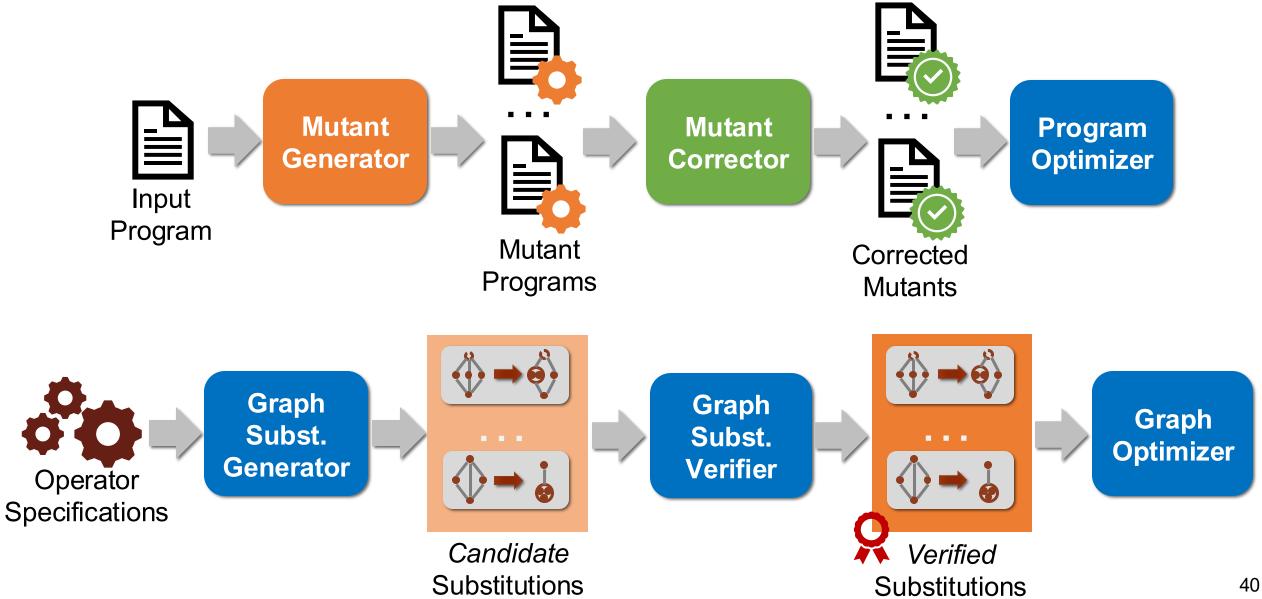
PET

- 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



PET vs TASO



Key Challenges

1. How to generate partially equivalent transformations?

Superoptimization

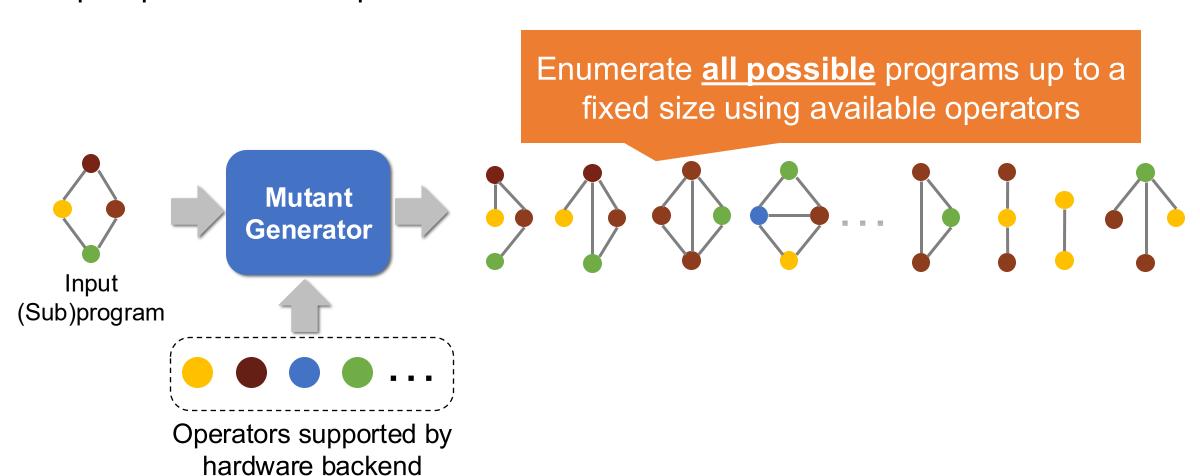
2. How to correct them?

Multi-linearity of DNN computations



Mutant Generator

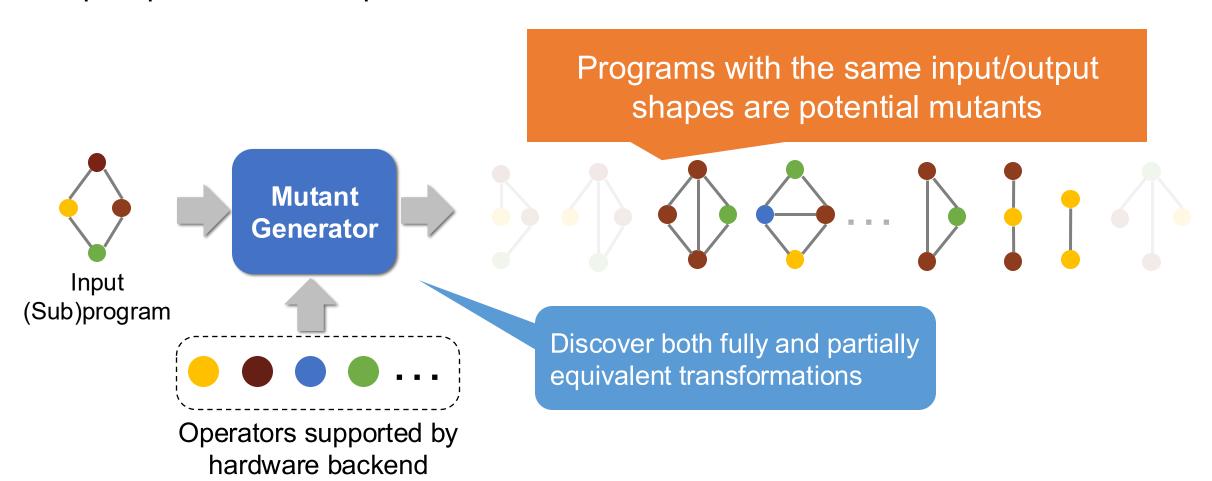
Superoptimization adopted from TASO¹





Mutant Generator

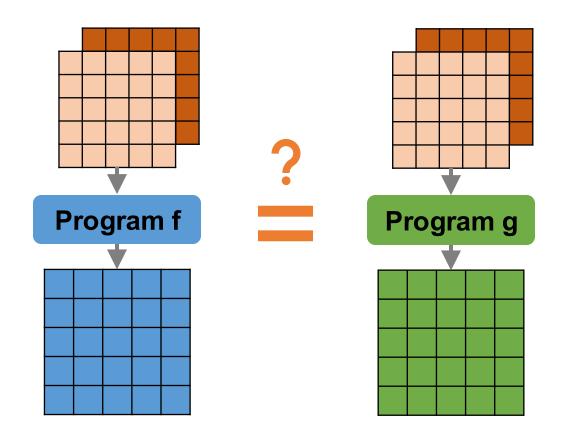
Superoptimization adopted from TASO¹



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



Challenges: Examine Transformations

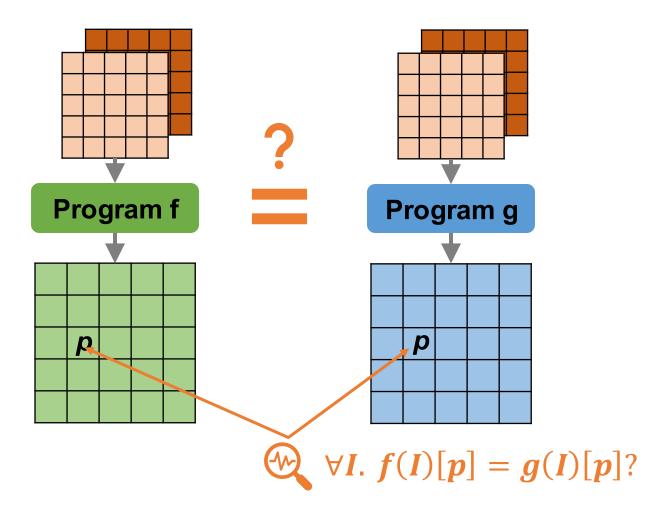


- 1. Which part of the computation is not equivalent?
- 2. How to correct the results?

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
 - $f(I_1, ..., X, ..., I_n) + f(I_1, ..., Y, ..., I_n) = f(I_1, ..., X + Y, ..., I_n)$
 - $\alpha \cdot f(I_1, \dots, X, \dots, I_n) = f(I_1, \dots, \alpha \cdot X, \dots, I_n)$
- DNN computation = MLTP + non-linear activations

Majority of the computation

O(m * n) examinations in strawman approach



O(1) examinations in PET's approach

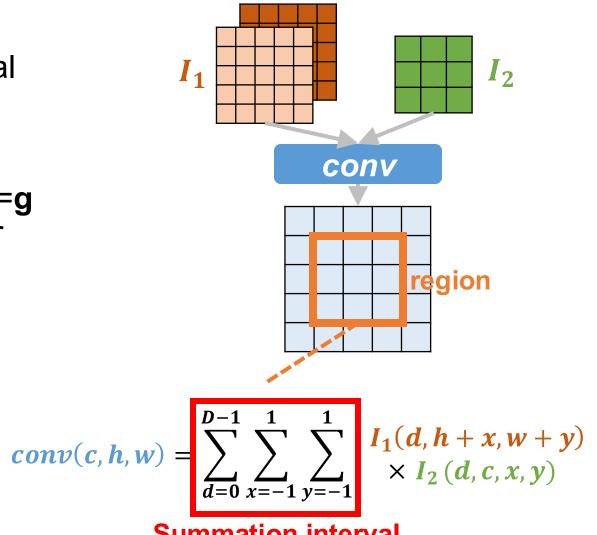
Insight #1: No Need to Enumerate All Output Positions

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



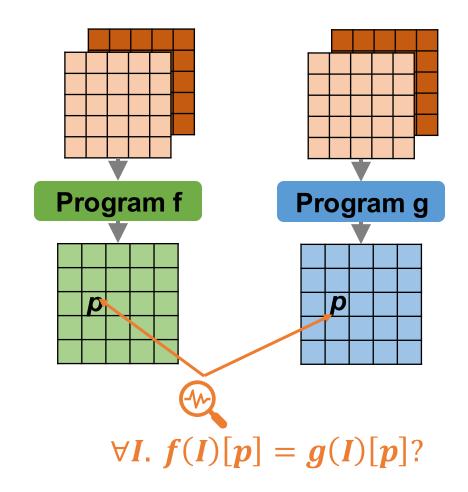
Insight #2: No Need to Consider All Possible Inputs

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 $(\frac{1}{2^{31}})^t$

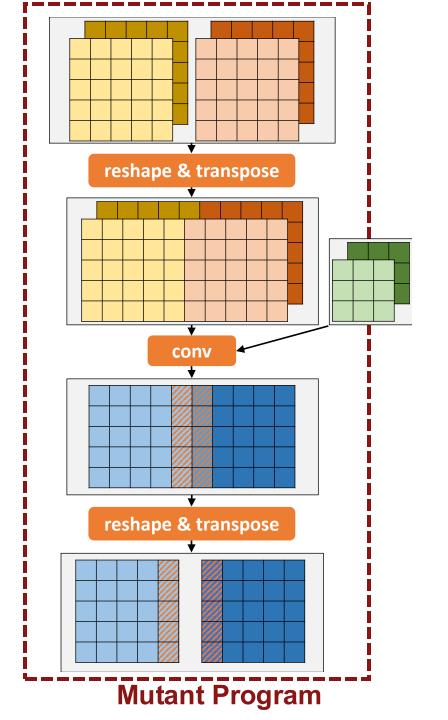
Run *t* random tests for each position *p*

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



Mutant Corrector

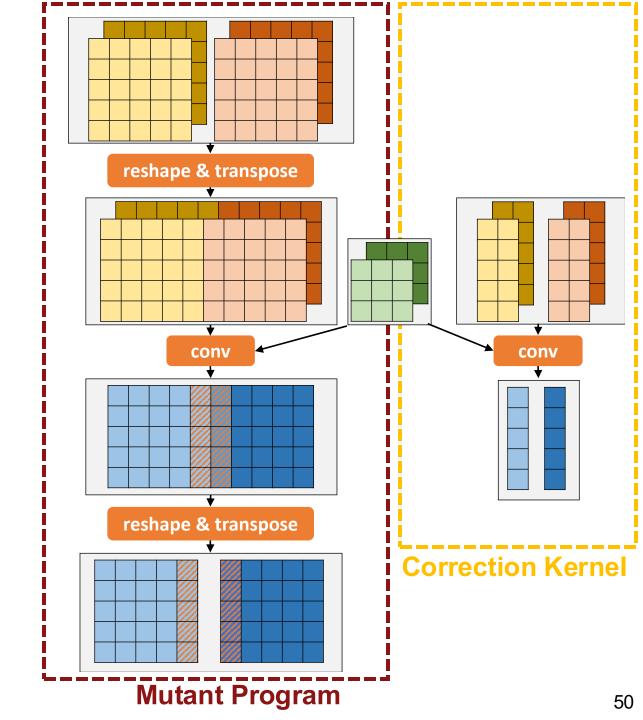
Goal: quickly and efficiently correcting the outputs of a mutant program



Mutant Corrector

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

Step 1: recompute the incorrect outputs using the original program



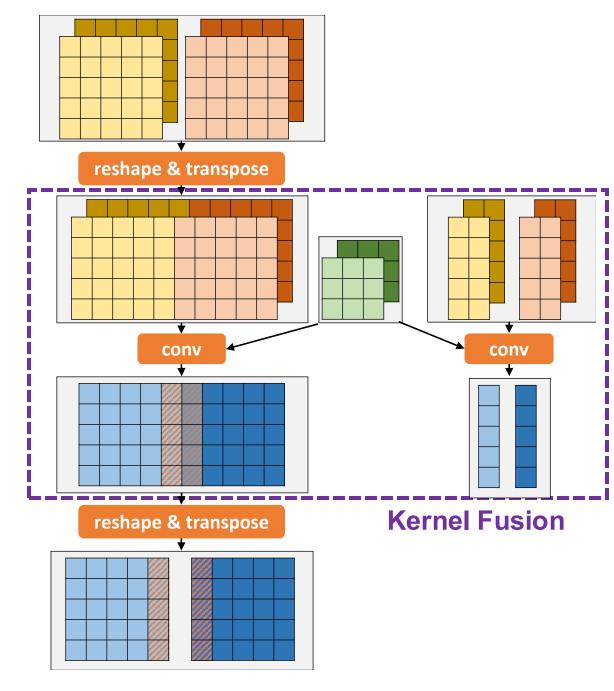
Mutant Corrector

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

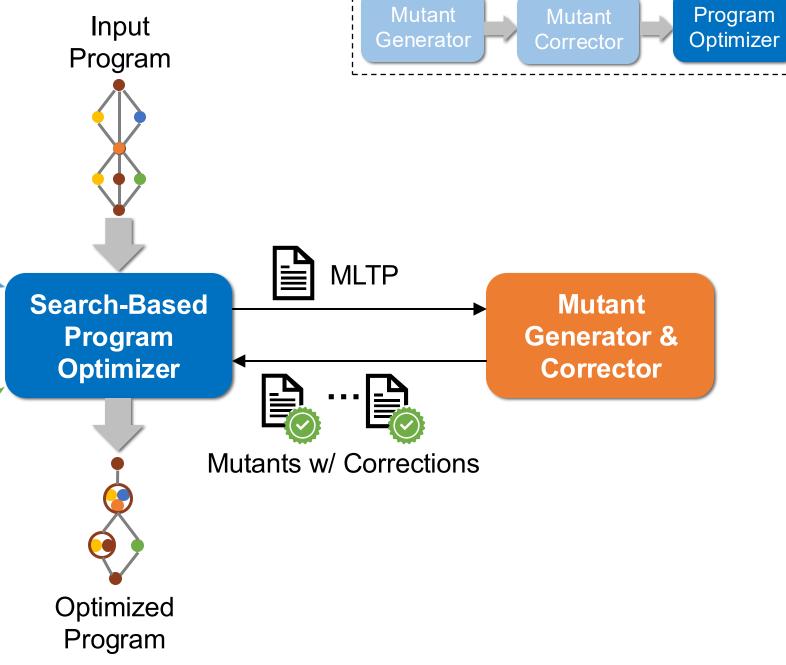




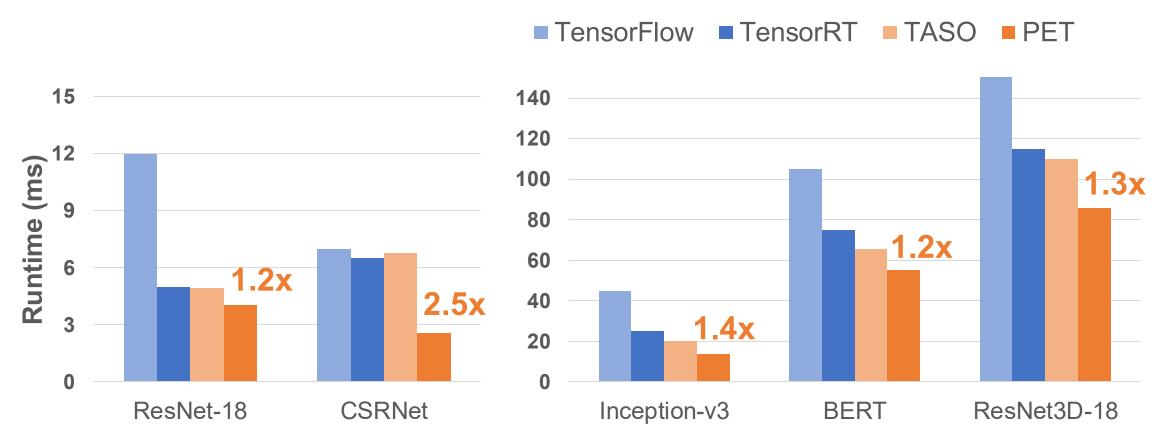
- Beam search
- Optimizing a DNN architecture takes less than <u>30</u> minutes

Other optimizations:

- Operator fusion
- Constant folding
- Redundancy elimination



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

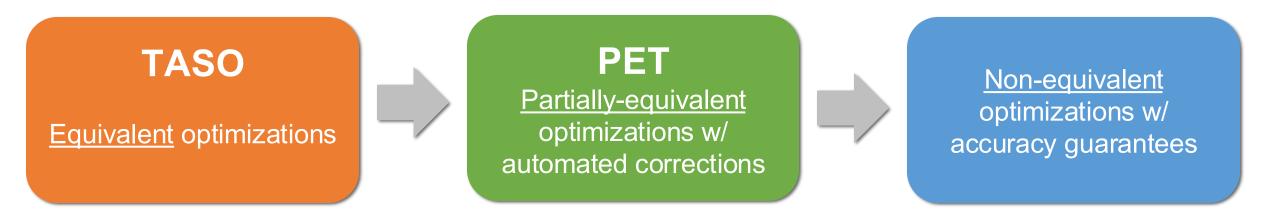


PET outperforms existing optimizers by 1.2-2.5x by combining fully and partially equivalent transformations

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



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?