15-442/15-642: Machine Learning Systems

Graph-Level Optimizations

Tianqi Chen Carnegie Mellon University

2/17/2025



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







Input Computation Graph

transformations

Optimized Computation Graph

Example: Fusing Convolution and Batch Normalization



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

Fusing Conv and BatchNorm



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

Recap: Resnet Example



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



* 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



Challenge of Graph Optimizations for ML



Infeasible to manually design graph optimizations for all cases



- 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



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



Optimized Comp. Graph

Graph Substitution Generator



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



Operators supported by hardware backend



Graph Substitution Generator



A substitution = a pair of equivalent graphs

Explicitly considering all pairs does not scale



Graph Substitution Generator













Pruning Redundant Substitutions





Graph Substitution Verifier



Verification Workflow

 $\begin{aligned} \forall x, w_1, w_2 . \\ & \left(Conv(x, w_1), Conv(x, w_2) \right) \\ &= Split \left(Conv(x, Concat(w_1, w_2)) \right) \end{aligned}$



Operator Specifications







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 <u>optimal</u> program for an input program

PET:

Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections

Motivation: Fully v.s. Partially 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





Partially Equivalent Transformation



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



- 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









1. How to generate partially equivalent transformations? Superoptimization

2. How to correct them?

Multi-linearity of DNN computations

Mutant Generator



Superoptimization adopted from TASO¹



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

Mutant Generator



Superoptimization adopted from TASO¹



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



Challenges: Examine Transformations



Which part of the computation is not equivalent?
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

•
$$f(I_1, ..., X, ..., I_n) + f(I_1, ..., Y, ..., I_n) = f(I_1, ..., X + Y, ..., I_n)$$

- $\alpha \cdot \boldsymbol{f}(I_1, \dots, X, \dots, I_n) = \boldsymbol{f}(I_1, \dots, \alpha \cdot X, \dots, I_n)$
- 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: $O(m * n) \rightarrow O(n)$



*Proof details available in the paper

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 <u>1%</u> overhead





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

TensorFlow TensorRT TASO PET



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



- 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?