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

- Submitted conference/journel name/years: Published as a conference paper at OSDI, 2021.
- Author: Haojie Wang, Jidong Zhai, Mingyu Gao et al.
- Presented date: 2022.03.29.

Presented by: JunYeong Park Student ID: 72220216





- I. Introduction
- II. Design Overview
- III. Mutation Generator
- IV. Mutation Corrector
- V. Program Optimizer
- VI. Conclusion



- 1. Limitation of existing frameworks
- Optimize the tensor programs applying 'Fully Equivalent Transformation'
- Tensor Program: direct acyclic computation graphs describing the operations applied to a set of tensors
- <u>Fully Equivalent Transformation</u>: the new subprogram is mathematically equivalent to the original subprogram for arbitrary inputs



 Miss opportunities to optimize tensor programs using transformation method





#### 2. Partially Equivalent Transformation

- Transformation that do not preserve full equivalence on all elements of the output tensors.
- Example:

Changing the shape or linearization ordering of input tensors

Replacing less efficient operators with more optimized operators with similar mathematical behavior

Transforming the graph structure of a program

- Improve computational efficiency and optimize the performance









Figure 1: A partially equivalent transformation that improves the performance of convolution by manipulating tensor shape and linearization. The shaded boxes in (b) highlight non-equivalent elements between two programs in the transformation. The correction kernel in (c) is applied to these elements to recover the functional equivalence of the input program.

### 3. Challenges

- Directly applying partially equivalent transformation makes <u>different output tensors</u> from output tensors of original input program(applying only fully equivalent transformation).
- Different output tensors may adversely affect performance.
- So applying correction kernels is necessary. But it is <u>difficult to examine the equivalence and generate correction</u> <u>kernels</u>.



#### 4. Contribution

- First attempt in tensor program optimization to exploit partially equivalent transformation with automated corrections.
- Theoretical foundations that simplify the equivalence examination and correction kernel generation.
- Efficient generation and optimization approaches to explore the large design space automatically with <u>both fully and</u> <u>partially equivalent transformations</u>.
- Implementation of new <u>framework PET</u>. (2.5x speedup)

#### 5. PET architecture component

- Mutation Generator: generate *mutants*
- Mutation Corrector: examine the equivalence and automatically generates correction kernel
- <u>Program Optimizer</u>: identify mutant candidates with high performance and post-optimize them

## **Design Overview**

#### 1. MLTP (Multi-linear Tensor Program)

- A program *p*, in which all operators *op* are multi-linear.

```
op(I_1,...,I_{k-1},X,...,I_n) + op(I_1,...,I_{k-1},Y,...,I_n) 
= op(I_1,...,I_{k-1},X+Y,...,I_n) 
\alpha \cdot op(I_1,...,I_{k-1},X,...,I_n) = op(I_1,...,I_{k-1},\alpha \cdot X,...,I_n)
```

- DNN subprogram
  - = Multi-linear tensor operator + non-linear operator

#### 2. PET Overview

- Split input program to be optimized based on a non-linear operator.
- <u>Mutation generator</u> discovers partially equivalent transformation by generating mutants for MLTPs in the subprogram.
- <u>Mutation corrector</u> examines the equivalence between a mutant and its original MLTP, and generates correction kernels.
- Corrected mutants are sent to program optimizer.
- <u>Program optimizer</u> combines f.e.t with p.e.t and post-optimize it.





### **Mutation Generator**



#### Depth-first search algorithm & check the shape

- PET iteratively <u>adds a new operator</u> to current program *p* by enumerating the type of operator from
   *O* and the input tensors to the operator.
- Enumerates all potential MLTPs up to *depth*.
   For each mutant *p*, PET checks whether *p* and *p\_0* have the same number and shapes of inputs/outputs.
   *p* is a valid mutant if it passes this test.

Algorithm I MLTP mutation generation algorithm.		
1: <b>Input:</b> A set of operators $O$ ; an input MLTP $\mathcal{P}_0$		
2: <b>Output:</b> A set of valid program mutants $\mathcal{M}$ for $\mathcal{P}_0$		
3: $I_0$ = the set of input tensors in $\mathcal{P}_0$		
4: $\mathcal{M} = arnothing$		
5: $BUILD(1, \emptyset, I_0)$		
6: // Depth-first search to construct mutants		
7: function $\operatorname{BUILD}(n, \mathcal{P}, I)$		
8: <b>if</b> $\mathcal{P}$ and $\mathcal{P}_0$ have the same input/output shapes <b>then</b>		
9: $\mathcal{M} = \mathcal{M} + \{\mathcal{P}\}$		
10: <b>if</b> $n < depth$ <b>then</b>		
11: <b>for</b> $op \in O$ <b>do</b>		
12: <b>for</b> $i \in I$ and $i$ is a valid input to <i>op</i> <b>do</b>		
13: Add operator $op$ into program $\mathcal{P}$		
14: Add the output tensors of <i>op</i> into <i>I</i>		
15: $\operatorname{BUILD}(n+1, \mathcal{P}, I)$		
16: Remove operator $op$ from $\mathcal{P}$		
17: Remove the output tensors of <i>op</i> from <i>I</i>		
18: return $\mathcal{M}$		

#### Table 1: Multi-linear tensor operators used in PET.

Operator	Description
add	Element-wise addition
mul	Element-wise multiplication
conv	Convolution
groupconv	Grouped convolution
dilatedconv	Dilated convolution
batchnorm	Batch normalization
avgpool	Average pooling
matmul	Matrix multiplication
batchmatmul	Batch matrix multiplication
concat	Concatenate multiple tensors
split	Split a tensor into multiple tensors
transpose	Transpose a tensor's dimensions
reshape	Decouple/combine a tensor's dimensions



#### 1. Why needs?

- Mutations resulted from partially equivalent transformations potentially leads to accuracy loss.
- Therefore, PET preserve the statistical behavior and guarantees the same model accuracy.

#### 2. Specifically what does it do?

- Input: an MLTP *p\_0*, mutant *p*
- Automatically generates correction kernels and applies them to the output tensors to maintain equivalence.

#### 3. Challenges

- The output tensors may be very large.  $\rightarrow$  infeasible to verify every single elements of the outputs.
- Numerically enumerating all possible values for many input variables is impractical.



#### 4. Proposal

- Only needs to verify a few representative output positions with a <u>small number of randomly generated</u> <u>input values</u>.

#### 5. Theoretical Foundations

- Assumption: input MLTP  $p_0$  and its mutant p each has one output.

$$\mathcal{P}(I_1,\ldots,I_n)[\vec{v}] = \sum_{\vec{r}\in\mathcal{R}(\vec{v})} \prod_{j=1}^n I_j[\mathbf{L}_j(\vec{v},\vec{r})]$$

- $\circ \hspace{0.2cm} I = (I_1,...,I_n)$ : n input tensors
- $\circ \hspace{0.2cm} I_{j}[ec{u}]$ : input value at position  $ec{u}$  of  $I_{j}$
- $\circ \ p(I)$ : output tensor of running p on n input tensors
- $\circ \ p(I)[ec{v}]$ : output value at position  $ec{v}$
- $R(\vec{v})$ : the summation interval of  $\vec{v}$  $(p(I)[\vec{v}]$  이거 계산할때 반복될 간격)
- $\circ~~ec{u}={f L}_j(ec{v},ec{r})$ : linear mapping from  $(ec{v},ec{r})$  to a position  $ec{u}$  of the j-th iput



#### 5-1. Example (3x3 convolution, zero padding)

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

(D: the number of channels, H: height, W: width, I\_1: input image, I\_2: kernel)

- Group the output positions with identical summation interval into a box
- → A box is a region of an output tensor whose elements all have the same summation intervals (9 boxes)
- All output positions in the same box have an identical summation interval and share similar mathematical properties.



Figure 4: The nine boxes of a convolution with a  $3 \times 3$  kernel and zero padding, as well as their summation intervals. A convolution has three summation dimensions (i.e., d, x, and y in Equation (1)). The channel dimension (i.e., d) has the same internal in all boxes and is thus omitted.

**Theorem 1** For two MLTPs  $\mathcal{P}_1$  and  $\mathcal{P}_2$  with an m-dimension output tensor, let  $\vec{e}_1, \ldots, \vec{e}_m$  be a set of m-dimension base vectors. That is,  $\vec{e}_i = (0, \ldots, 0, 1, 0, \ldots, 0)$  is an m-tuple with all coordinates equal to 0 except the i-th.

Let  $\mathcal{B}$  be a box for  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , and let  $\vec{v}_0$  be an arbitrary position in  $\mathcal{B}$ . Define  $\vec{v}_j = \vec{v}_0 + \vec{e}_j$ ,  $1 \le j \le m$ . If  $\forall I, 0 \le i \le m$ ,  $\mathcal{P}_1(I)[\vec{v}_i] = \mathcal{P}_2(I)[\vec{v}_i]$ , then  $\forall I, \vec{v} \in \mathcal{B}$ ,  $\mathcal{P}_1(I)[\vec{v}] = \mathcal{P}_2(I)[\vec{v}]$ .

- We want to check the equivalence of two output tensors

 $\rightarrow$  choose one random position for each box (v0)  $\rightarrow$  compare p\_0(v0~v(m+1)) == p(v0~v(m+1))

- If all position in the If the values in the m+1 positions are all the same, the values in the overlapping box ensure equivalence.







Figure 4: The nine boxes of a convolution with a  $3 \times 3$  kernel and zero padding, as well as their summation intervals. A convolution has three summation dimensions (i.e., d, x, and yin Equation (1)). The channel dimension (i.e., d) has the same

internal in all boxes and is thus omitted.





**Theorem 2** For two MLTPs  $\mathcal{P}_1$  and  $\mathcal{P}_2$  with n input tensors, let  $\vec{v}$  be a position where  $\mathcal{P}_1$  and  $\mathcal{P}_2$  are not equivalent, i.e.,  $\exists I, \mathcal{P}_1(I)[\vec{v}] \neq \mathcal{P}_2(I)[\vec{v}]$ . Let I' be a randomly generated input uniformly sampled from a finite field  $\mathbb{F}$ . The probability that  $\mathcal{P}_1(I')[\vec{v}] = \mathcal{P}_2(I')[\vec{v}]$  is at most  $\frac{n}{p}$ , where p is the number of possible values in  $\mathbb{F}$ .

- n/p =  $(\frac{1}{2^{31}})^t$ 

(in t random integer)

- If program *p* and *p\_0* are not equivalent,

the probability that the output values of the two programs are the same when random inputs are inserted is very low.

Examine the equivalence for each pair of <u>overlapped boxes</u> from the two MLTPs, <u>using a small number of random</u>
 <u>tests</u>



Figure 4: The nine boxes of a convolution with a  $3 \times 3$  kernel and zero padding, as well as their summation intervals. A convolution has three summation dimensions (i.e., d, x, and yin Equation (1)). The channel dimension (i.e., d) has the same internal in all boxes and is thus omitted.

#### 6. Mutation Correction Algorithm

#### (1) Box Propagation

- Infer the split points of its output tensors based on the split points of tis input tensors and the operator type and hyper-parameter.
- a set of split points for each dimension of a tensor
  - $\rightarrow$  identify the boundaries of its boxes.

#### (2) Random testing for each box pair

- If two boxes do not have any overlapped region, the can be skipped. (not equivalent)
- Examines the equivalence of the two programs on <u>m+1 positions</u> identified by Theroem1
   (4 dimensional output tensor → v0 ~ v5)
- For each of these m+1 positions, PET runs a set of random tests by assigning input tensors with values uniformly sampled from a finite field F.
- <u>caching optimization</u>: all random inputs for the tests are the same.



Figure 5: Box propagation for the example in Figure 1. The red arrows indicate the split points of each tensor dimension.





#### 6. Mutation Correction Algorithm

#### (3) Correction kernel generation



- For each box failing the random tests, PET generates correction kernels to maintain original statistical behavior (mathematical equivalence)
- Correction kernel performs the <u>same set of operations as the original MLTP</u> but only on those boxes where the two input programs are <u>not equivalent</u>.
- Kernel generation leverage existing libraries and techniques.
- To reduce the correction overhead, PET <u>opportunistically fuses</u> the correction kernels with existing tensor operators.



#### 7. Fusing Correction Kernels

- Launching the correction kernels  $\rightarrow$  substantial overhead
  - $\rightarrow$  eliminate the advantage of partially equivalent transformation (performance gain)
- Opportunistically <u>fuses correction kernels</u> with other tensor operators.



Figure 6: Fusing correction kernels with DNN kernels.



#### 1. Program splitting

- Splits an input program into multiple disjoint subprograms
   with smaller sizes. (larger program → high complexity)
- Split point: non-linear operators (activation layers)

#### 2. Subprogram optimization for best mutant

- Call mutation generator  $\rightarrow$  keeping the top-K candidates
- Larger K  $\rightarrow$  more memory and higher computation cost
- At each step, each of the obtained mutants <u>replace its</u>
   <u>corresponding subprogram</u> in each of the current candidates(*p*)
   to generate a new candidate(*p\_new*)
- Apply a series of post-optimization

```
Algorithm 2 Program optimization algorithm.
 1: Input: An input tensor program \mathcal{P}_0
 2: Output: An optimized tensor program \mathcal{P}_{opt}
 3:
 4: Split \mathcal{P}_0 into a list of subprograms
 5: Initialize a heap \mathcal{H} to record the top-K programs
 6: \mathcal{H}.insert(\mathcal{P}_0)
 7: // Greedily mutate each subprogram
 8: for each subprogram S \in \mathcal{P}_0 do
          mutants = GETMUTANTS(S)
 9:
          Initialize a new heap \mathcal{H}_{new}
10:
          for \mathcal{P} \in \mathcal{H} do
11:
               for \mathcal{M} \in mutants do
12:
                    \mathcal{P}_{new} = replace \mathcal{S} with \mathcal{M} in \mathcal{P}
13:
                    Apply post-optimizations on \mathcal{P}_{new}
14:
15:
                    \mathcal{H}_{new}.insert (\mathcal{P}_{new})
          \mathcal{H} = \mathcal{H}_{new}
16:
17: \mathcal{P}_{opt} = the program with the best performance in \mathcal{H}
18: return \mathcal{P}_{opt}
19:
20: function GETMUTANTS(S_0)
           O = the set of mutant operators for S_0
21:
          Q = \{S_0\}, mutants = \{S_0\}
22:
          for r rounds do
23:
               Q_{\text{new}} = \{\}
24:
25:
               for S \in Q do
26:
                    for each subset of operators S' \in S do
                         for \mathcal{M}' \in MUTATIONGENERATOR(\mathcal{O}, \mathcal{S}') do
27:
                              \mathcal{M} = replace \mathcal{S}' with \mathcal{M}' in \mathcal{S}
28:
                              Add \mathcal{M} to Q_{new} and mutants
29:
30:
               Q = Q_{\text{new}}
```

```
31: return mutants
```

## Program Optimizer

#### 3. Post Optimizations

- The optimized mutant for all subprograms need to be stitched together.
- Connecting their input and output tensors
   → fuse R/T operators across subprograms
   → fuse the non-linear operators
- [fuse R/T operators] Figure 7(b): Group reshape/transpose operators between subprograms by reordering the R/T operators with non-linear activations.
- [fuse non-linear operators] Figure 7(c): Conv+ReLU = Conv-ReLU

#### [3 post optimization]

- Inverse elimination: eliminate any pairs of R/T operators(called inverse group) that can be cancel out each other.
- Operator fusion:

Fuses remaining R/T operators into a single operator to reduce the kernel launch cost. non-linear activations in a tensor program are fused with R/T or other linear operations.

- <u>Preprocessing</u>: preprocess any operator if all its input tensors are statically known.



Figure 7: Post-optimizations applied when stitching two subprograms SG-1 and SG-2. R/T refers to a reshape followed by a transpose. Conv and ReLU denote a convolution and a ReLU operator, respectively.



### Conclusion



- 5 DNN Architectures: ResNet-18 (Image Classification), CSRNet (semantic segmentation), Inception-v3, BERT (language representation), Resnet3D-18 (video processing)
- Default mutation generation depth = 4 (Algorithm 1. DFS) search rounds = 4 (Algorithm 2, Greedy search)
- 1. End-to-end inference performance
- Batch size under 1 and 16
- Larger batch
  - $\rightarrow$  relatively higher performance



Figure 8: End-to-end performance comparison between PET and existing frameworks. For each DNN, the numbers above the PET bars show the speedups over the best baseline. TASO does not support the 3D convolution operators in Resnet3D-18.

#### 2. Partial equivalent transformations discovered by PET

- Manually modify TASO as <u>adding partial equivalent transformations and corresponding</u> <u>correction kernels</u>.
- Evaluate the end-to-end inference performance of TASO compared to original TASO and PET.
- 'TASO + PET' transformations and corrections' is relatively <u>faster</u> than original TASO.



Figure 9: Performance benefits after adding PET's partially equivalent transformations into TASO.



### Thank You For Listening