Tensor Algebra Superoptimizer (TASO)

Jia et al. (2019)

Timi Adeniran

Research Questions



• Can graph substitutions be automated?

• Will jointly optimizing graph substitutions with the data layout yield better performance?

Can we do better?





• Automatically generates substitutions.

• Verifies their correctness.

• Jointly optimizes substitutions with data layout on an input graph.



Generating Substitutions

Algorithm 1 Graph substitution generation algorithm.

```
1: Input: A set of operators \mathcal{P}, and a set of input tensors \mathcal{I}.
 2: Output: Candidate graph substitutions S.
 3:
 4: // Step 1: enumerating potential graphs.
 5: \mathcal{D} = \{\} // \mathcal{D} \text{ is a graph hash table indexed by their fingerprints.} \}
6: BUILD(1, \emptyset, I)
7: function BUILD(n, G, I)
        if G contains duplicated computation then
 8:
             return
 9:
        \mathcal{D} = \mathcal{D} + (\text{FingerPrint}(\mathcal{G}), \mathcal{G})
10:
        if n < threshold then
11:
            for op \in \mathcal{P} do
12:
                 for i \in \mathcal{I} and i is a valid input to op do
13:
                     Add operator op into graph G.
14:
                     Add the output tensors of op into I.
15:
                     BUILD(n+1, G, I)
16:
                     Remove operator op from G.
17:
                     Remove the output tensors of op from I.
18:
19:
20: // Step 2: testing graphs with identical fingerprint.
21: S = \{\}
22: for G_1, G_2 \in \mathcal{D} with the same FINGERPRINT(·) do
        if G_1 and G_2 are equivalent for all test cases then
23:
             S = S + (G_1, G_2)
24:
25: return S
```

- DFS algorithm to generate all possible graphs for a given set of operators and input tensors.
- Evaluates graphs on random inputs and hashes their result (fingerprint).
- Graphs with the same fingerprint are good candidates.

Verifying Substitutions

• Tests candidate substitutions against user-provided operator properties.

• Users provide a symbolic representation of properties in Python.

• Uses the Z3 theorem prover to verify properties holds for substitutions.

Pruning Substitutions

• Eliminates input tensor renaming redundancies.

• Prunes substitutions with a common subgraph between source and target.

Pruning	Remaining	Reduction
Techniques	Substitutions	v.s. Initial
Initial	28744	1×
Input tensor renaming	17346	1.7×
Common subgraph	743	39×

Joint Optimization

• New: TASO jointly optimizes graph substitution and data layout.

• TASO uses a cost-based backtracking algorithm and picks the best* layout for each substitution.

• Cost model sums execution time of each operator in a graph.

Joint Optimization - Backtracking Algorithm

Algorithm 2 Cost-Based Backtracking Search

```
1: Input: an input graph G_{in}, verified substitutions S, a cost
    model Cost(\cdot), and a hyper parameter \alpha.
 2: Output: an optimized graph.
 3:
 4: \mathcal{P} = \{\mathcal{G}_{in}\} // \mathcal{P} is a priority queue sorted by Cost.
 5: while \mathcal{P} \neq \{\} do
         \mathcal{G} = \mathcal{P}.dequeue()
 6:
         for substitution s \in S do
 7:
             // LAYOUT(G, s) returns possible layouts applying s on G.
 8:
             for layout l \in LAYOUT(\mathcal{G}, s) do
 9:
                  // APPLY(G, s, l) applies s on G with layout l.
10:
                  G' = \operatorname{Apply}(G, s, l)
11:
                  if G' is valid then
12:
                      if Cost(G') < Cost(G_{opt}) then
13:
                           G_{opt} = G'
14:
                      if Cost(G') < \alpha \times Cost(G_{opt}) then
15:
                           \mathcal{P}.engueue(G')
16:
17: return Gopt
```

Evaluation

• Applied 743 substitutions generated from up to four operators on ResNet-50, ResNeXt-50, NasNet-A, NasRNN and BERT.

• Compared inference performance with TensorFlow, TensorFlow XLA, TensorRT, TVM, and MetaFlow on a V100 GPU.

• Compared inference performance on BERT of joint optimization with independent optimizations.

Results - End-to-end Inference Performance Comparison



Results - Joint Optimization



Discussion - Strengths

- Novel approach to tensor graph optimization:
 - Automatically generating graph substitutions.
 - Use of formal methods to verify them.
 - Jointly optimizing graph substitutions with data layout.

• Produces good results.

Discussion - Limitations (1)

• Only supports a maximum of four operators.

• Manual approach for adding operators.

• Too many redundant substitutions.

• Evaluation does not consider time taken to optimize.

Discussion - Limitations (2)

Table 1. Discrepancy between TASO's cost model estimates and TASO's end-to-end inference latency on some unoptimised DNNs. E2E stands for end-to-end inference latency. Time is measured in milliseconds.

Oversimplified cost model

DNNs	COST MODEL	E2E	DIFF (%)
DALL-E	1.8269	1.7324	5.2%
INCEPTIONV3	8.3650	9.2098	10.1%
BERT	1.0453	1.1264	7.8%
SQUEEZENET	1.3082	1.4006	7.1%
RESNEXT-50	6.1545	7.6498	24%
T-T	2.4828	2.7281	9.9%

From X-RLflow by He et al.

Impact

- Inspired follow-up work into automatic graph substitutions:
 - PET (2021): Wang et al. "PET: Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections," OSDI 2021.
 - TENSAT (2021): Yang et al. "Equality Saturation for Tensor Graph Superoptimization," MLSys 2021.
 - Unity (2022): Unger et al. "Accelerating DNN Training Through Joint Optimization of Algebraic Transformations and Parallelization". OSDI 2022.
 - X-RLFlow (2023): He at al. "X-RLflow: Graph Reinforcement Learning for Neural Network Subgraphs Transformation," MLSys 2023.

Questions?