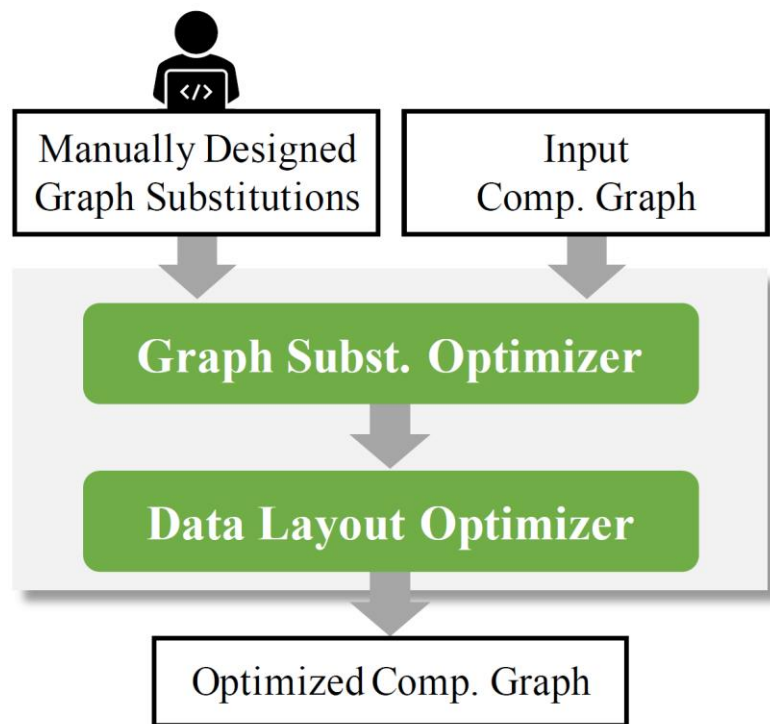# UNIVERSITY OF CAMBRIDGE

## TASO: Optimizing Deep Learning Computation with Automatic Generation of Graph Substitutions
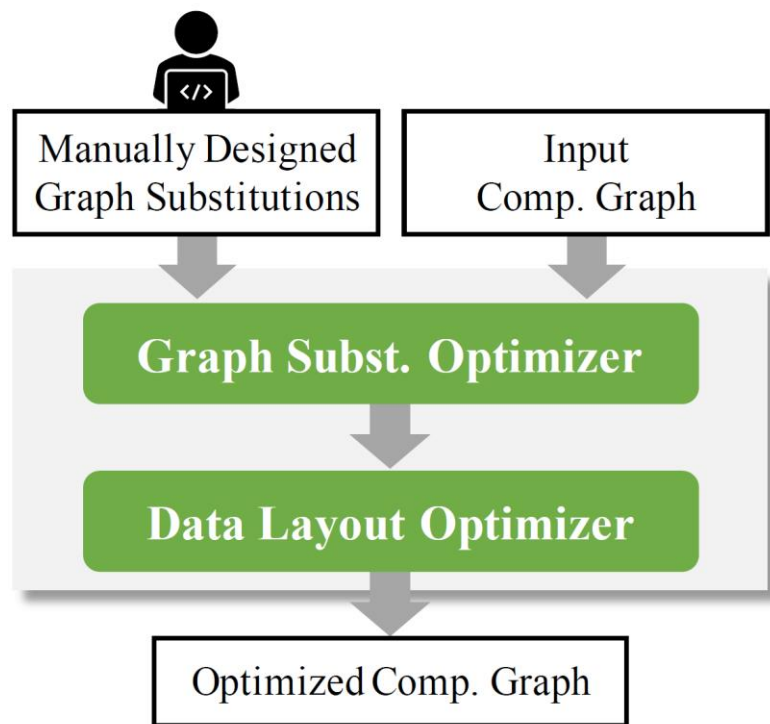
Jia, Z., Padon, O., Thomas, J., Warszawski, T., Zaharia, M., & Aiken, A. (2019, October). In Proceedings of the 27th ACM Symposium on Operating Systems Principles (pp. 47-62).

*R244 Large-scale data processing and optimisation*
*Presentation by Martin Graf on 16/11/2022*

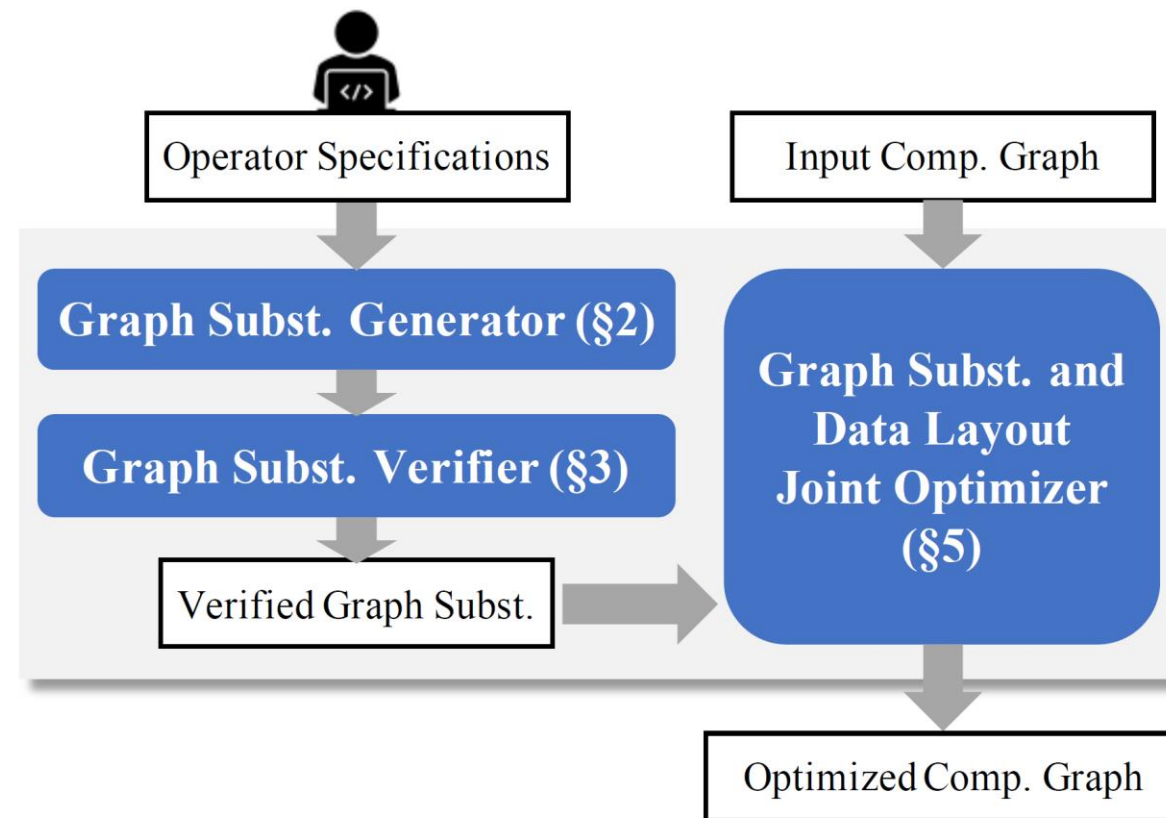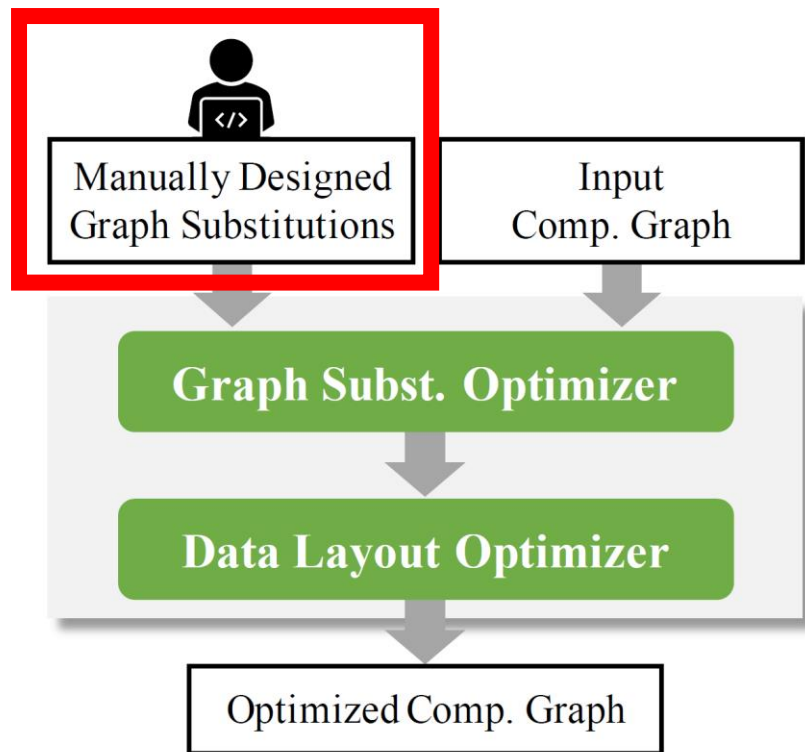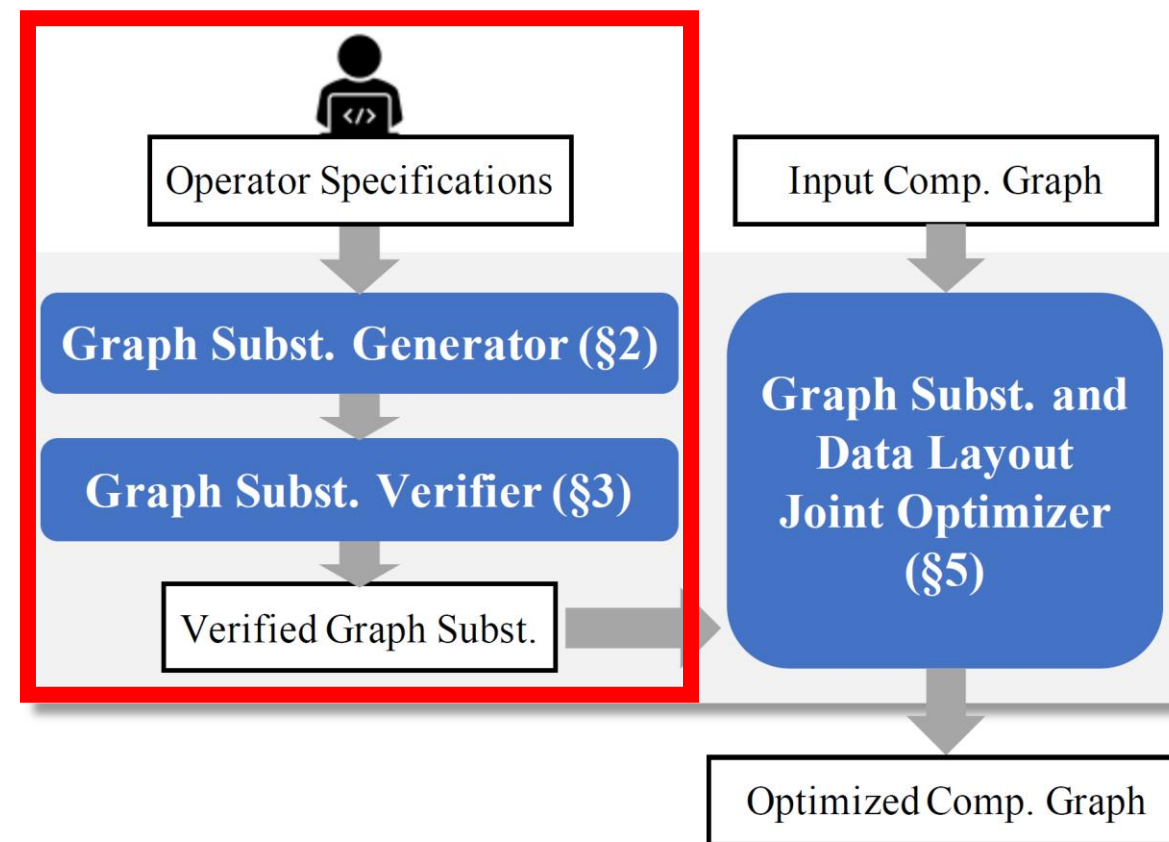**(a)** Existing DNN frameworks.

UNIVERSITY OF
CAMBRIDGE

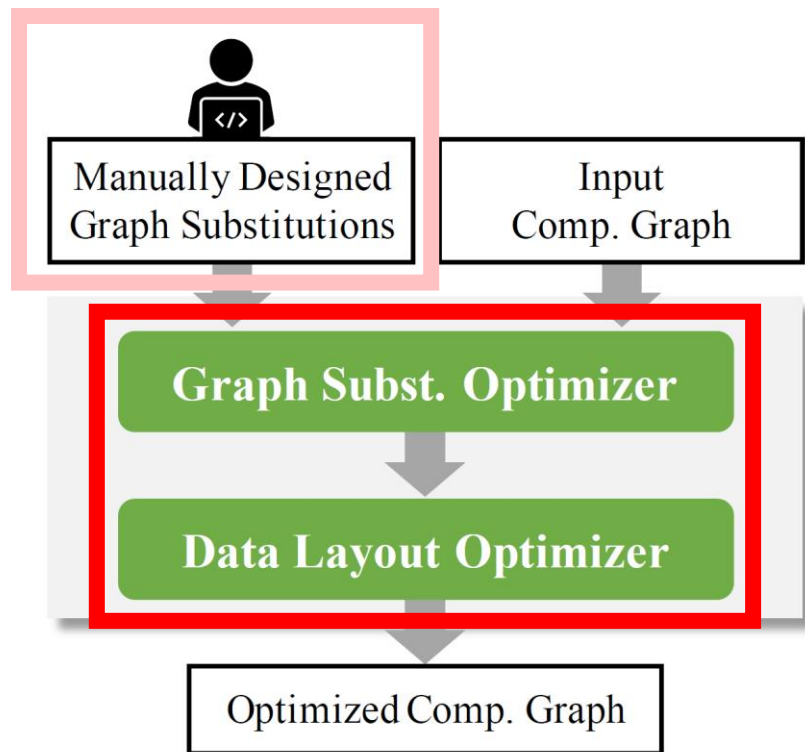(a) Existing DNN frameworks.

(b) TASO.

(a) Existing DNN frameworks.
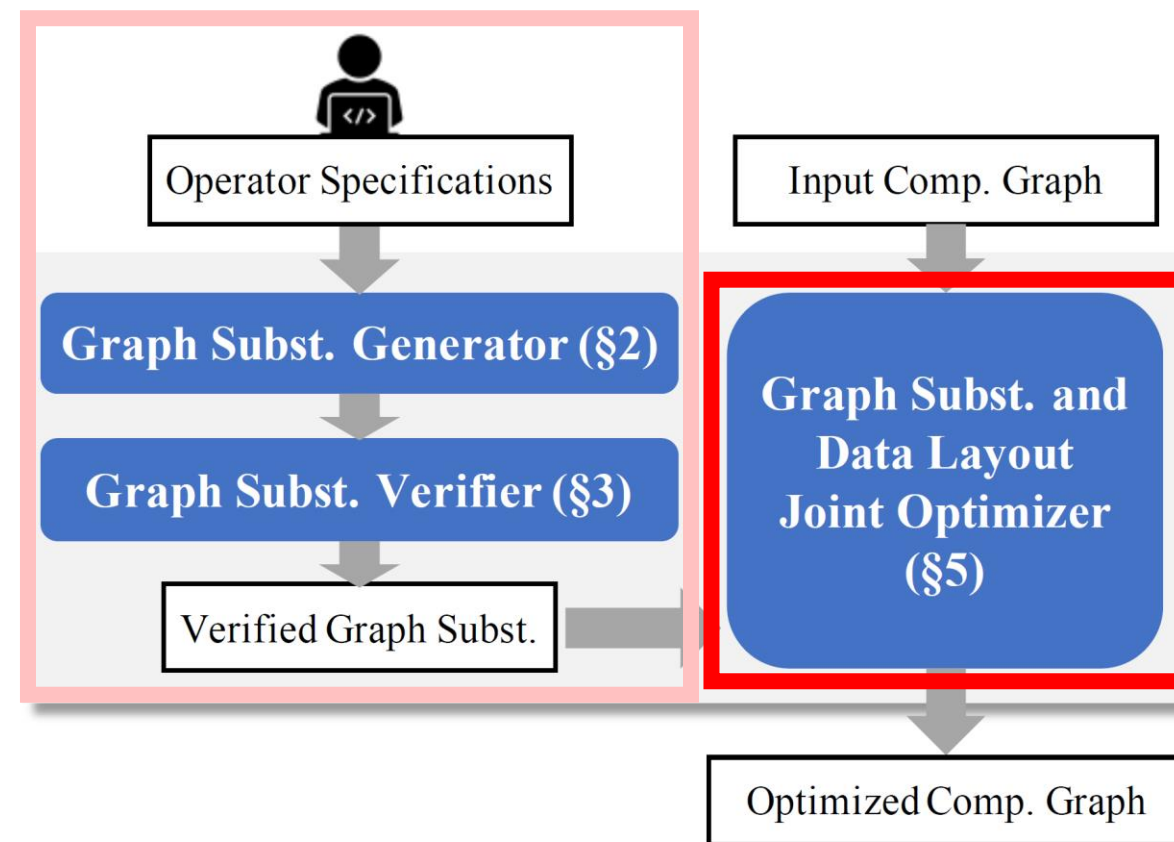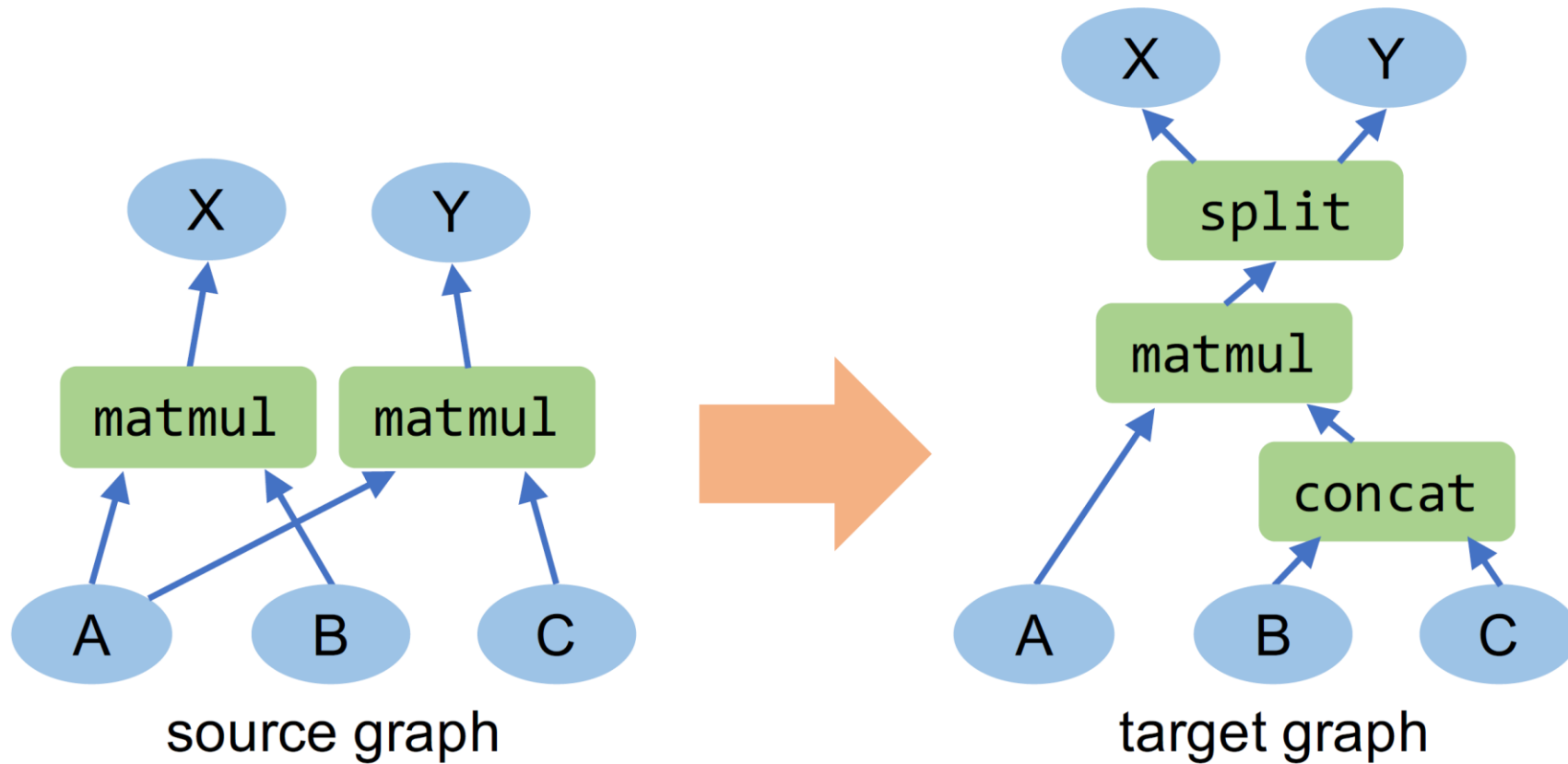
(b) TASO.

(a) Existing DNN frameworks.

(b) TASO.

**(b)** Fusing two matrix multiplications using concatenation and split.

**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 $\mathcal{S}$.
3:
4:   *// Step 1: enumerating potential graphs.*
5:   $\mathcal{D} = \{\}$  *// $\mathcal{D}$ is a graph hash table indexed by their fingerprints.*
6:   BUILD$(1, \emptyset, \mathcal{I})$
7:   **function** BUILD$(n, \mathcal{G}, \mathcal{I})$
8:      **if** $\mathcal{G}$ contains duplicated computation **then**
9:        **return**
10:      $\mathcal{D} = \mathcal{D} + (\text{FINGERPRINT}(\mathcal{G}), \mathcal{G})$
11:      **if** $n < threshold$ **then**
12:        **for** $op \in \mathcal{P}$ **do**
13:          **for** $i \in \mathcal{I}$ and $i$ is a valid input to $op$ **do**
14:            Add operator $op$ into graph $\mathcal{G}$.
15:            Add the output tensors of $op$ into $\mathcal{I}$.
16:            BUILD$(n + 1, \mathcal{G}, \mathcal{I})$
17:            Remove operator $op$ from $\mathcal{G}$.
18:            Remove the output tensors of $op$ from $\mathcal{I}$.
19:
20:   *// Step 2: testing graphs with identical fingerprint.*
21:   $\mathcal{S} = \{\}$
22:   **for** $\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{D}$ with the same FINGERPRINT$(\cdot)$ **do**
23:      **if** $\mathcal{G}_1$ and $\mathcal{G}_2$ are equivalent for all test cases **then**
24:        $\mathcal{S} = \mathcal{S} + (\mathcal{G}_1, \mathcal{G}_2)$
25:   **return** $\mathcal{S}$

Enumerate potential graphs

UNIVERSITY OF
CAMBRIDGE

**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 $\mathcal{S}$.
3:
4: // *Step 1: enumerating potential graphs.*
5: $\mathcal{D} = \{\}$ // $\mathcal{D}$ is a graph hash table indexed by their fingerprints.
6: BUILD$(1, \emptyset, \mathcal{I})$
7: **function** BUILD$(n, \mathcal{G}, \mathcal{I})$
8:     **if** $\mathcal{G}$ contains duplicated computation **then**
9:         **return**
10:     $\mathcal{D} = \mathcal{D} + (\text{FINGERPRINT}(\mathcal{G}), \mathcal{G})$
11:     **if** $n < threshold$ **then**
12:         **for** $op \in \mathcal{P}$ **do**
13:             **for** $i \in \mathcal{I}$ and $i$ is a valid input to $op$ **do**
14:                 Add operator $op$ into graph $\mathcal{G}$.
15:                 Add the output tensors of $op$ into $\mathcal{I}$.
16:                 BUILD$(n + 1, \mathcal{G}, \mathcal{I})$
17:                 Remove operator $op$ from $\mathcal{G}$.
18:                 Remove the output tensors of $op$ from $\mathcal{I}$.
19:
20: // *Step 2: testing graphs with identical fingerprint.*
21: $\mathcal{S} = \{\}$
22: **for** $\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{D}$ with the same FINGERPRINT$(\cdot)$ **do**
23:     **if** $\mathcal{G}_1$ and $\mathcal{G}_2$ are equivalent for all test cases **then**
24:         $\mathcal{S} = \mathcal{S} + (\mathcal{G}_1, \mathcal{G}_2)$
25: **return** $\mathcal{S}$

Enumerate potential graphs

Collect and test candidate substitutions

**UNIVERSITY OF CAMBRIDGE**

**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 $\mathcal{S}$.
3:
4:  // *Step 1: enumerating potential graphs.*
5:  $\mathcal{D} = \{\}$ // $\mathcal{D}$ is a graph hash table indexed by their fingerprints.
6:  BUILD$(1, \emptyset, \mathcal{I})$
7:  **function** BUILD$(n, \mathcal{G}, \mathcal{I})$
8:      **if** $\mathcal{G}$ contains duplicated computation **then**
9:          **return**
10:     $\mathcal{D} = \mathcal{D} + (\text{FINGERPRINT}(\mathcal{G}), \mathcal{G})$
11:     **if** $n < threshold$ **then**
12:         **for** $op \in \mathcal{P}$ **do**
13:             **for** $i \in \mathcal{I}$ and $i$ is a valid input to $op$ **do**
14:                 Add operator $op$ into graph $\mathcal{G}$.
15:                 Add the output tensors of $op$ into $\mathcal{I}$.
16:                 BUILD$(n + 1, \mathcal{G}, \mathcal{I})$
17:                 Remove operator $op$ from $\mathcal{G}$.
18:                 Remove the output tensors of $op$ from $\mathcal{I}$.
19:
20: // *Step 2: testing graphs with identical fingerprint.*
21: $\mathcal{S} = \{\}$
22: **for** $\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{D}$ with the same FINGERPRINT$(\cdot)$ **do**
23:     **if** $\mathcal{G}_1$ and $\mathcal{G}_2$ are equivalent for all test cases **then**
24:         $\mathcal{S} = \mathcal{S} + (\mathcal{G}_1, \mathcal{G}_2)$
25: **return** $\mathcal{S}$

Enumerate potential graphs
→ Fingerprinting + Hash Collisions
(Bansal, 2006)

Collect and test candidate substitutions

**UNIVERSITY OF CAMBRIDGE**

**Table 2.** Operator properties used for verification. The operators are defined in Table 1, and the properties are grouped by the operators they involve. Logical variables $w, x, y,$ and $z$ are of type tensor, and variables $a, c, k, p,$ and $s$ are of type parameter. The variable $a$ is used for the axis of concatenation and split, $c$ for the activation mode of convolution, $k$ for the kernel shape of pooling, $p$ for the padding mode of convolution and pooling, and $s$ for the strides of convolution and pooling.

| Operator Property | Comment |
|---|---|
| $\forall x, y, z.\ \mathrm{ewadd}(x, \mathrm{ewadd}(y, z)) = \mathrm{ewadd}(\mathrm{ewadd}(x, y), z)$ | ewadd is associative |
| $\forall x, y.\ \mathrm{ewadd}(x, y) = \mathrm{ewadd}(y, x)$ | ewadd is commutative |
| $\forall x, y, z.\ \mathrm{ewmul}(x, \mathrm{ewmul}(y, z)) = \mathrm{ewmul}(\mathrm{ewmul}(x, y), z)$ | ewmul is associative |
| $\forall x, y.\ \mathrm{ewmul}(x, y) = \mathrm{ewmul}(y, x)$ | ewmul is commutative |
| $\forall x, y, z.\ \mathrm{ewmul}(\mathrm{ewadd}(x, y), z) = \mathrm{ewadd}(\mathrm{ewmul}(x, z), \mathrm{ewmul}(y, z))$ | distributivity |
| $\forall x, y, w.\ \mathrm{smul}(\mathrm{smul}(x, y), w) = \mathrm{smul}(x, \mathrm{smul}(y, w))$ | smul is associative |
| $\forall x, y, w.\ \mathrm{smul}(\mathrm{ewadd}(x, y), w) = \mathrm{ewadd}(\mathrm{smul}(x, w), \mathrm{smul}(y, w))$ | distributivity |
| $\forall x, y, w.\ \mathrm{smul}(\mathrm{ewmul}(x, y), w) = \mathrm{ewmul}(x, \mathrm{smul}(y, w))$ | operator commutativity |
| $\forall x.\ \mathrm{transpose}(\mathrm{transpose}(x)) = x$ | transpose is its own inverse |
| $\forall x, y.\ \mathrm{transpose}(\mathrm{ewadd}(x, y)) = \mathrm{ewadd}(\mathrm{transpose}(x), \mathrm{transpose}(y))$ | operator commutativity |
| $\forall x, y.\ \mathrm{transpose}(\mathrm{ewmul}(x, y)) = \mathrm{ewmul}(\mathrm{transpose}(x), \mathrm{transpose}(y))$ | operator commutativity |
| $\forall x, w.\ \mathrm{smul}(\mathrm{transpose}(x), w) = \mathrm{transpose}(\mathrm{smul}(x, w))$ | operator commutativity |
| $\forall x, y, z.\ \mathrm{matmul}(x, \mathrm{matmul}(y, z)) = \mathrm{matmul}(\mathrm{matmul}(x, y), z)$ | matmul is associative |
| $\forall x, y, w.\ \mathrm{smul}(\mathrm{matmul}(x, y), w) = \mathrm{matmul}(x, \mathrm{smul}(y, w))$ | matmul is linear |
| $\forall x, y, z.\ \mathrm{matmul}(x, \mathrm{ewadd}(y, z)) = \mathrm{ewadd}(\mathrm{matmul}(x, y), \mathrm{matmul}(x, z))$ | matmul is linear |
| $\forall x, y.\ \mathrm{transpose}(\mathrm{matmul}(x, y)) = \mathrm{matmul}(\mathrm{transpose}(y), \mathrm{transpose}(x))$ | matmul and transpose |
| $\forall s, p, c, x, y, w.\ \mathrm{conv}(s, p, c, \mathrm{smul}(x, w), y) = \mathrm{conv}(s, p, c, x, \mathrm{smul}(y, w))$ | conv is bilinear |
| $\forall s, p, x, y, w.\ \mathrm{smul}(\mathrm{conv}(s, p, \mathrm{A_{none}}, x, y), w) = \mathrm{conv}(s, p, \mathrm{A_{none}}, \mathrm{smul}(x, w), y)$ | conv is bilinear |
| $\forall s, p, x, y, z.\ \mathrm{conv}(s, p, \mathrm{A_{none}}, x, \mathrm{ewadd}(y, z)) = \mathrm{ewadd}(\mathrm{conv}(s, p, \mathrm{A_{none}}, x, y), \mathrm{conv}(s, p, \mathrm{A_{none}}, x, z))$ | conv is bilinear |
| $\forall s, p, x, y, z.\ \mathrm{conv}(s, p, \mathrm{A_{none}}, \mathrm{ewadd}(x, y), z) = \mathrm{ewadd}(\mathrm{conv}(s, p, \mathrm{A_{none}}, x, z), \mathrm{conv}(s, p, \mathrm{A_{none}}, y, z))$ | conv is bilinear |
| $\forall s, c, k, x, y.\ \mathrm{conv}(s, \mathrm{P_{same}}, c, x, y) = \mathrm{conv}(s, \mathrm{P_{same}}, c, x, \mathrm{enlarge}(k, y)),$ | enlarge convolution kernel |
| $\forall s, p, x, y.\ \mathrm{conv}(s, p, \mathrm{A_{relu}}, x, y) = \mathrm{relu}(\mathrm{conv}(s, p, \mathrm{A_{none}}, x, y))$ | conv with $\mathrm{A_{relu}}$ applies relu |
| $\forall x.\ \mathrm{relu}(\mathrm{transpose}(x)) = \mathrm{transpose}(\mathrm{relu}(x))$ | operator commutativity |
| $\forall s, p, x, k.\ \mathrm{conv}(s, p, \mathrm{A_{none}}, x, \mathrm{C_{pool}}(k)) = \mathrm{pool_{avg}}(k, s, p, x)$ | pooling by conv. with $\mathrm{C_{pool}}$ |
| $\forall k, x.\ \mathrm{conv}(1, \mathrm{P_{same}}, \mathrm{A_{none}}, x, \mathrm{I_{conv}}(k)) = x$ | identity kernel |
| $\forall x.\ \mathrm{matmul}(x, \mathrm{I_{matmul}}) = x$ | identity matrix |
| $\forall x.\ \mathrm{ewmul}(x, \mathrm{I_{ewmul}}) = x$ | ewmul identity |
| $\forall a, x, y.\ \mathrm{split_0}(a, \mathrm{concat}(a, x, y)) = x$ | split definition |
| $\forall a, x, y.\ \mathrm{split_1}(a, \mathrm{concat}(a, x, y)) = y$ | split definition |
| $\forall x, y, z, w.\ \mathrm{concat}(0, \mathrm{concat}(1, x, y), \mathrm{concat}(1, z, w)) = \mathrm{concat}(1, \mathrm{concat}(0, x, z), \mathrm{concat}(0, y, w))$ | geometry of concatenation |
| $\forall a, x, y, w.\ \mathrm{concat}(a, \mathrm{smul}(x, w), \mathrm{smul}(y, w)) = \mathrm{smul}(\mathrm{concat}(a, x, y), w)$ | operator commutativity |
| $\forall a, x, y, z, w.\ \mathrm{concat}(a, \mathrm{ewadd}(x, y), \mathrm{ewadd}(z, w)) = \mathrm{ewadd}(\mathrm{concat}(a, x, z), \mathrm{concat}(a, y, w))$ | operator commutativity |
| $\forall a, x, y, z, w.\ \mathrm{concat}(a, \mathrm{ewmul}(x, y), \mathrm{ewmul}(z, w)) = \mathrm{ewmul}(\mathrm{concat}(a, x, z), \mathrm{concat}(a, y, w))$ | operator commutativity |
| $\forall a, x, y.\ \mathrm{concat}(a, \mathrm{relu}(x), \mathrm{relu}(y)) = \mathrm{relu}(\mathrm{concat}(a, x, y))$ | operator commutativity |
| $\forall x, y.\ \mathrm{concat}(1, \mathrm{transpose}(x), \mathrm{transpose}(y)) = \mathrm{transpose}(\mathrm{concat}(0, x, y))$ | concatenation and transpose |
| $\forall x, y, z.\ \mathrm{concat}(1, \mathrm{matmul}(x, y), \mathrm{matmul}(x, z)) = \mathrm{matmul}(x, \mathrm{concat}(1, y, z))$ | concatenation and matrix mul. |
| $\forall x, y, z, w.\ \mathrm{matmul}(\mathrm{concat}(1, x, z), \mathrm{concat}(0, y, w)) = \mathrm{ewadd}(\mathrm{matmul}(x, y), \mathrm{matmul}(z, w))$ | concatenation and matrix mul. |
| $\forall s, p, c, x, y, z.\ \mathrm{concat}(0, \mathrm{conv}(s, p, c, x, z), \mathrm{conv}(s, p, c, y, z)) = \mathrm{conv}(s, p, c, \mathrm{concat}(0, x, y), z)$ | concatenation and conv. |
| $\forall s, p, c, x, y, z.\ \mathrm{concat}(1, \mathrm{conv}(s, p, c, x, y), \mathrm{conv}(s, p, c, x, z)) = \mathrm{conv}(s, p, c, x, \mathrm{concat}(0, y, z))$ | concatenation and conv. |
| $\forall s, p, x, y, z, w.\ \mathrm{conv}(s, p, \mathrm{A_{none}}, \mathrm{concat}(1, x, z), \mathrm{concat}(1, y, w)) =$ $\mathrm{ewadd}(\mathrm{conv}(s, p, \mathrm{A_{none}}, x, y), \mathrm{conv}(s, p, \mathrm{A_{none}}, z, w))$ | concatenation and conv. |
| $\forall k, s, p, x, y.\ \mathrm{concat}(1, \mathrm{pool_{avg}}(k, s, p, x), \mathrm{pool_{avg}}(k, s, p, y)) = \mathrm{pool_{avg}}(k, s, p, \mathrm{concat}(1, x, y))$ | concatenation and pooling |
| $\forall k, s, p, x, y.\ \mathrm{concat}(0, \mathrm{pool_{max}}(k, s, p, x), \mathrm{pool_{max}}(k, s, p, y)) = \mathrm{pool_{max}}(k, s, p, \mathrm{concat}(0, x, y))$ | concatenation and pooling |
| $\forall k, s, p, x, y.\ \mathrm{concat}(1, \mathrm{pool_{max}}(k, s, p, x), \mathrm{pool_{max}}(k, s, p, y)) = \mathrm{pool_{max}}(k, s, p, \mathrm{concat}(1, x, y))$ | concatenation and pooling |

**Table 3.** The number of remaining graph substitutions after applying the pruning techniques in order.

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

**Algorithm 2** Cost-Based Backtracking Search

1: **Input**: an input graph $\mathcal{G}_{in}$, verified substitutions $\mathcal{S}$, 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**
6:      $\mathcal{G} = \mathcal{P}.\text{dequeue}()$
7:      **for** substitution $s \in \mathcal{S}$ **do**
8:          // $\text{LAYOUT}(\mathcal{G}, s)$ returns possible layouts applying $s$ on $\mathcal{G}$.
9:          **for** layout $l \in \text{LAYOUT}(\mathcal{G}, s)$ **do**
10:             // $\text{APPLY}(\mathcal{G}, s, l)$ applies $s$ on $\mathcal{G}$ with layout $l$.
11:             $\mathcal{G}' = \text{APPLY}(\mathcal{G}, s, l)$
12:             **if** $\mathcal{G}'$ is valid **then**
13:                 **if** $Cost(\mathcal{G}') < Cost(\mathcal{G}_{opt})$ **then**
14:                    $\mathcal{G}_{opt} = \mathcal{G}'$
15:                 **if** $Cost(\mathcal{G}') < \alpha \times Cost(\mathcal{G}_{opt})$ **then**
16:                    $\mathcal{P}.\text{enqueue}(\mathcal{G}')$
17: **return** $\mathcal{G}_{opt}$

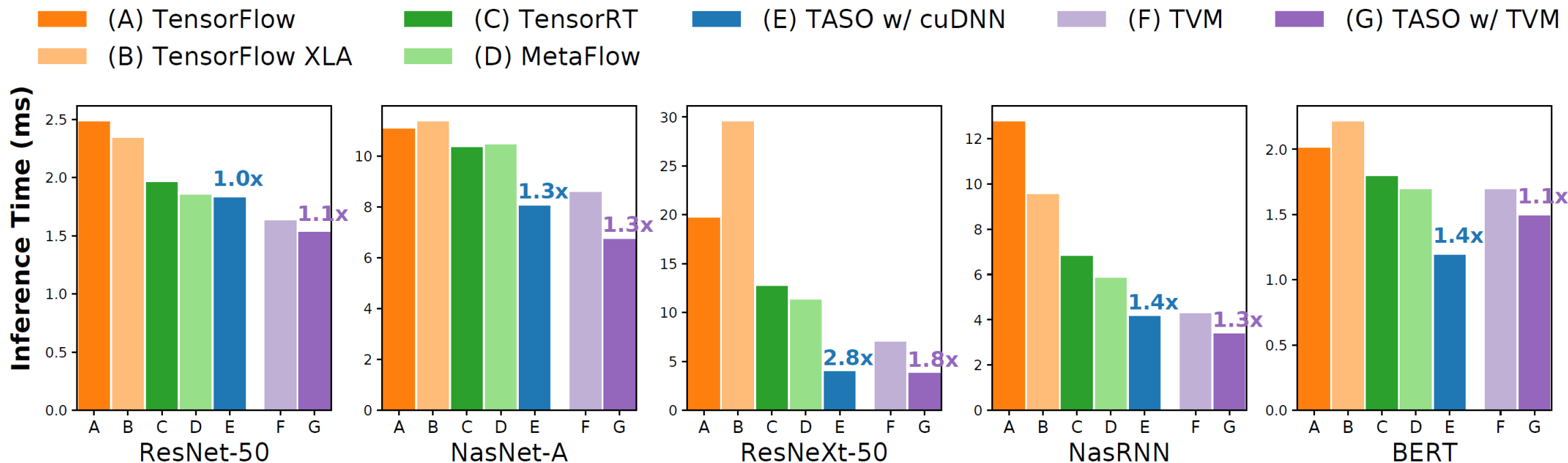Cost-based backtracking algorithm from MetaFlow
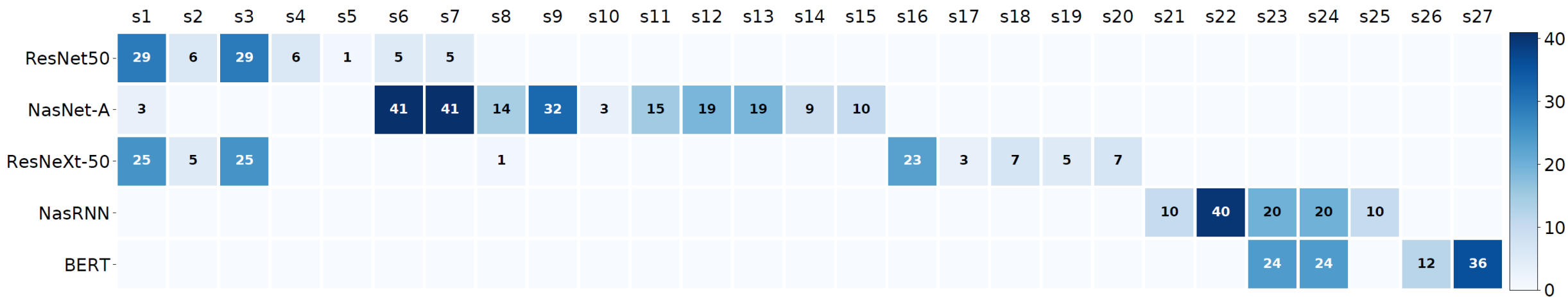(Jia, 2019)

=

Priority Queue

+

Nested for loop over all substitutions/data layouts
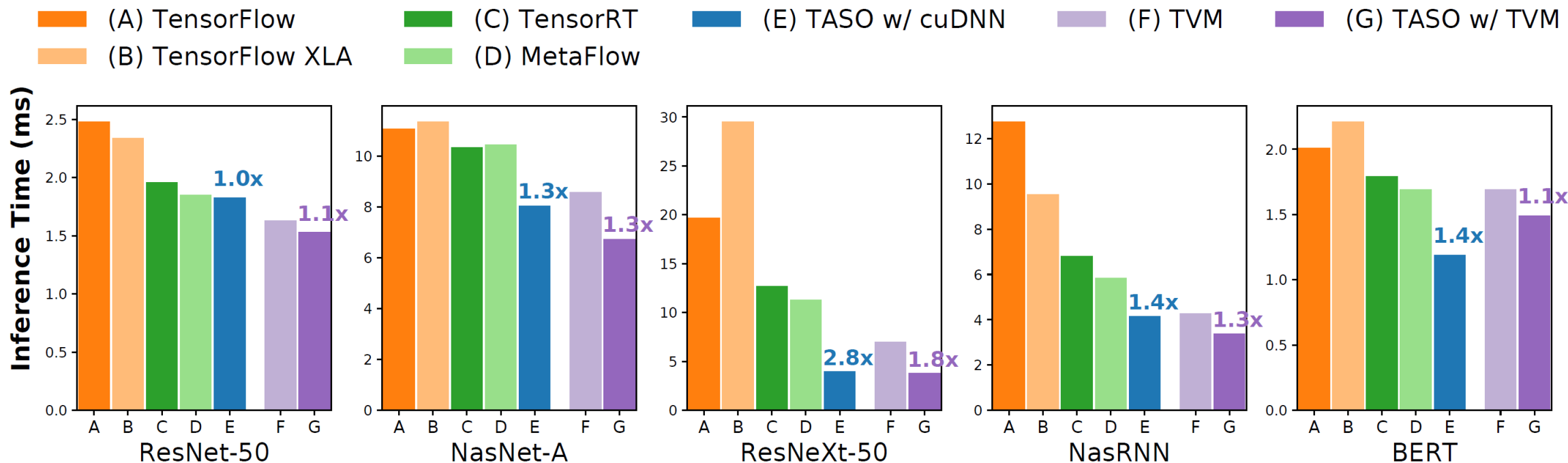
UNIVERSITY OF
CAMBRIDGE

**Figure 7.** End-to-end inference performance comparison among existing DNN frameworks and TASO. The experiments were performed using a single inference sample, and all numbers were measured by averaging 1,000 runs on a NVIDIA V100 GPU. We evaluated the TASO's performance with both the cuDNN and TVM backends. For each DNN architecture, the numbers above the TASO bars show the speedup over the best existing approach with the same backend.

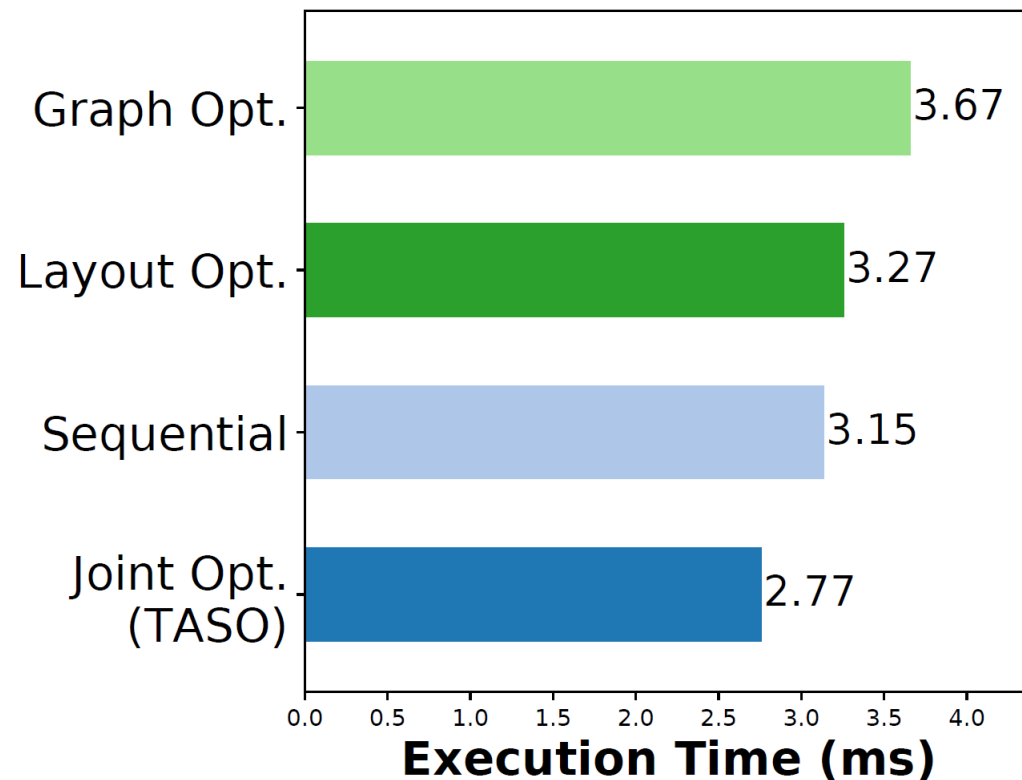|        | s1 | s2 | s3 | s4 | s5 | s6 | s7 | s8 | s9 | s10 | s11 | s12 | s13 | s14 | s15 | s16 | s17 | s18 | s19 | s20 | s21 | s22 | s23 | s24 | s25 | s26 | s27 |
|--------|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| ResNet50   | 29 | 6  | 29 | 6  | 1  | 5  | 5  |    |    |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| NasNet-A   | 3  |    |    |    |    | 41 | 41 | 14 | 32 | 3   | 15  | 19  | 19  | 9   | 10  |     |     |     |     |     |     |     |     |     |     |     |     |
| ResNeXt-50 | 25 | 5  | 25 |    |    |    |    | 1  |    |     |     |     |     |     |     | 23  | 3   | 7   | 5   | 7   |     |     |     |     |     |     |     |
| NasRNN     |    |    |    |    |    |    |    |    |    |     |     |     |     |     |     |     |     |     |     |     | 10  | 40  | 20  | 20  | 10  |     |     |
| BERT       |    |    | 12 |    |    |    |    |    |    |     |     |     |     |     |     |     |     |     |     |     |     |     | 24  | 24  |     | 12  | 36  |

**Figure 10.** A heat map of how often the verified substitutions are used to optimize the five DNN architectures. Only substitutions used in at least one DNN are listed. For each architecture, the number indicates how many times a substitution is used by TASO to obtain the optimized graph.

UNIVERSITY OF
CAMBRIDGE

**Figure 7.** End-to-end inference performance comparison among existing DNN frameworks and TASO. The experiments were performed using a single inference sample, and all numbers were measured by averaging 1,000 runs on a NVIDIA V100 GPU. We evaluated the TASO's performance with both the cuDNN and TVM backends. For each DNN architecture, the numbers above the TASO bars show the speedup over the best existing approach with the same backend.

**Figure 12.** End-to-end inference performance comparison on BERT using different strategies to optimize graph substitution and data layout.

# Evaluation

Con:
- 743 graph substitutions > 43 operator properties

# Evaluation: 743 substitutions > 43 operator properties

All substitutions found by TASO are valid given the user provided properties.

# Evaluation: 743 substitutions > 43 operator properties

All substitutions found by TASO are valid given the user provided properties.

→ There is a sequence of applications of user provided properties proving the correctness of the substitution.

UNIVERSITY OF
CAMBRIDGE

# Evaluation: 743 substitutions > 43 operator properties

All substitutions found by TASO are valid given the user provided properties.

→ There is a sequence of applications of user provided properties proving the correctness of the substitution.

Each user provided property is equivalent to a substitution.

UNIVERSITY OF
CAMBRIDGE

# Evaluation: 743 substitutions > 43 operator properties

All substitutions found by TASO are valid given the user provided properties.

→ There is a sequence of applications of user provided properties proving the correctness of the substitution.

Each user provided property is equivalent to a substitution.

→ **All substitutions are sequences of substitutions provided by the user.**

UNIVERSITY OF
CAMBRIDGE

# Evaluation

Con:

- 743 graph substitutions > 43 operator properties
- Why prune all redundant substitutions?

# Evaluation

Con:

- 743 graph substitutions > 43 operator properties
- Why prune all redundant substitutions?
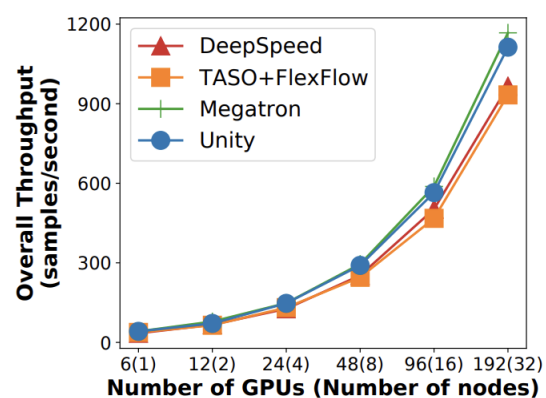- Cost measure of joint optimizer seems oversimplified

# Evaluation

Con:

- 743 graph substitutions > 43 operator properties

- Why prune all redundant substitutions?

- Cost measure of joint optimizer seems oversimplified

Pro:

- Lots of real-world considerations

UNIVERSITY OF
CAMBRIDGE

# Evaluation

Con:

- 743 graph substitutions > 43 operator properties

- Why prune all redundant substitutions?

- Cost measure of joint optimizer seems oversimplified

Pro:

- Lots of real-world considerations

- Many individual influential contributions

UNIVERSITY OF
CAMBRIDGE

# Impact and Follow Up Work

- PET: non-fully equivalent graph substitutions
  *(Wang, 2021)*

- Equality Saturation: improves joint optimization algorithm
  *(Yang, 2021)*

- Unity: jointly optimizes graph substitutions, data layout, and parallelism
  *(Unger, 2022)*
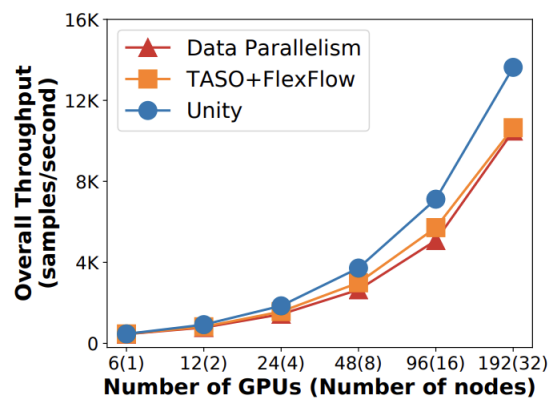
UNIVERSITY OF
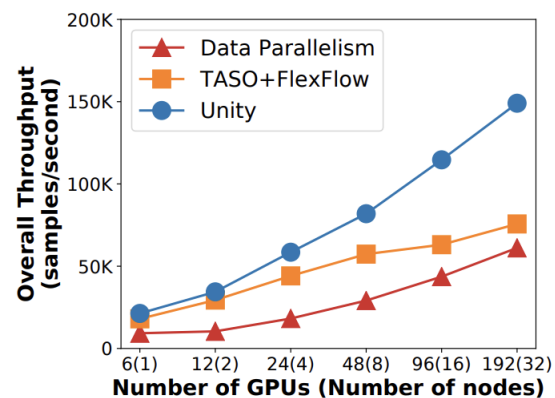CAMBRIDGE
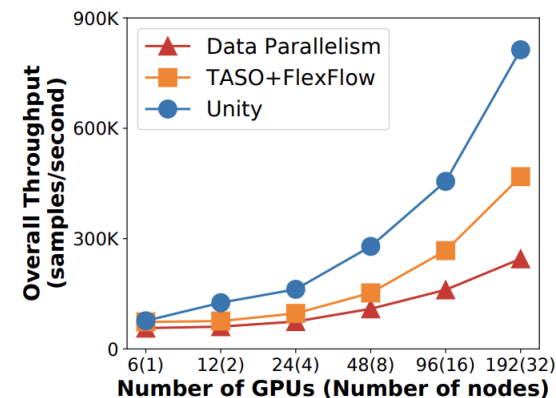
(a) ResNeXt-50.

(b) BERT-Large.

(c) DLRM.

(d) CANDLE-Uno.

(e) Inception-v3.

(f) MLP.

(g) XDL.

(Unger, 2022)

# References

- Bansal, Sorav, and Alex Aiken. "Automatic generation of peephole superoptimizers." ACM SIGARCH Computer Architecture News 34.5 (2006): 394-403.

- Jia, Zhihao, et al. "Optimizing DNN computation with relaxed graph substitutions." Proceedings of Machine Learning and Systems 1 (2019): 27-39.

UNIVERSITY OF
CAMBRIDGE

# References

- Jia, Zhihao, et al. "TASO: optimizing deep learning computation with automatic generation of graph substitutions." Proceedings of the 27th ACM Symposium on Operating Systems Principles. 2019.

- Unger, Colin, et al. "Unity: Accelerating {DNN} Training Through Joint Optimization of Algebraic Transformations and Parallelization." 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22). 2022.

UNIVERSITY OF
CAMBRIDGE

# References

- Wang, Haojie, et al. "{PET}: Optimizing Tensor Programs with Partially Equivalent Transformations and Automated Corrections." 15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21). 2021.

- Yang, Yichen, et al. "Equality saturation for tensor graph superoptimization." Proceedings of Machine Learning and Systems 3 (2021): 255-268.

UNIVERSITY OF
CAMBRIDGE

# UNIVERSITY OF CAMBRIDGE

## TASO: Optimizing Deep Learning Computation with Automatic Generation of Graph Substitutions

*R244 Large-scale data processing and optimisation*
*Presentation by Martin Graf on 16/11/2022*