# TASO: Optimizing Deep Learning Computation with Automatic Generation of Graph Substitutions
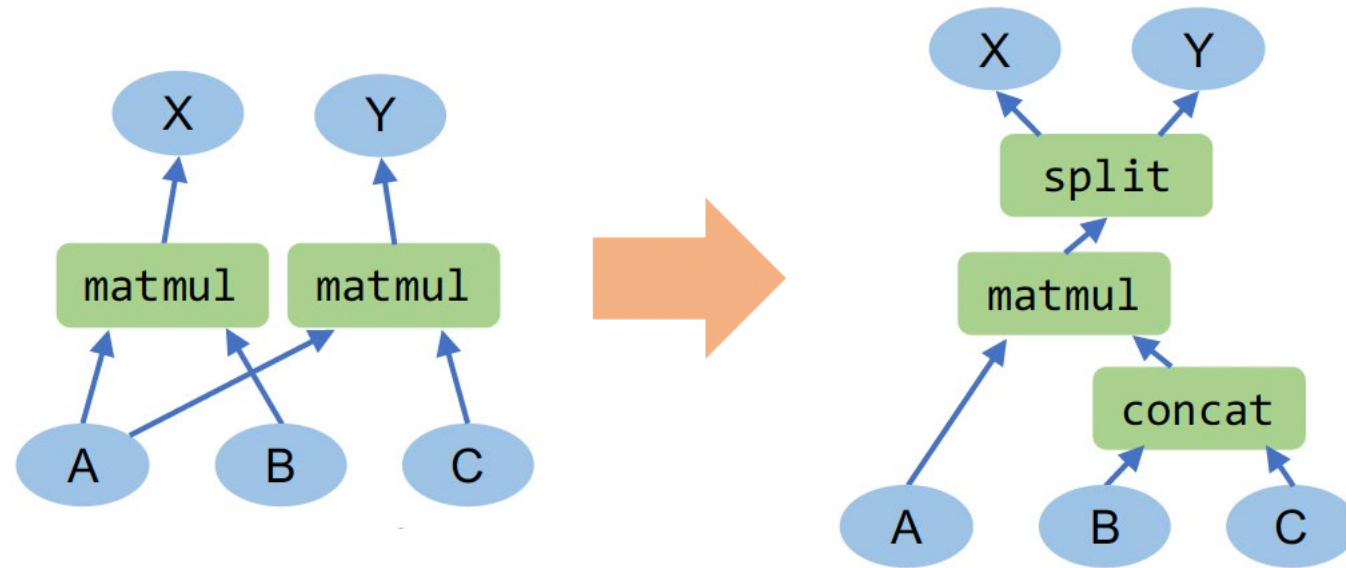
Zhihao Jia, Oded Padon, James Thomas,

Todd Warszawski, Matei Zaharia, Alex Aiken

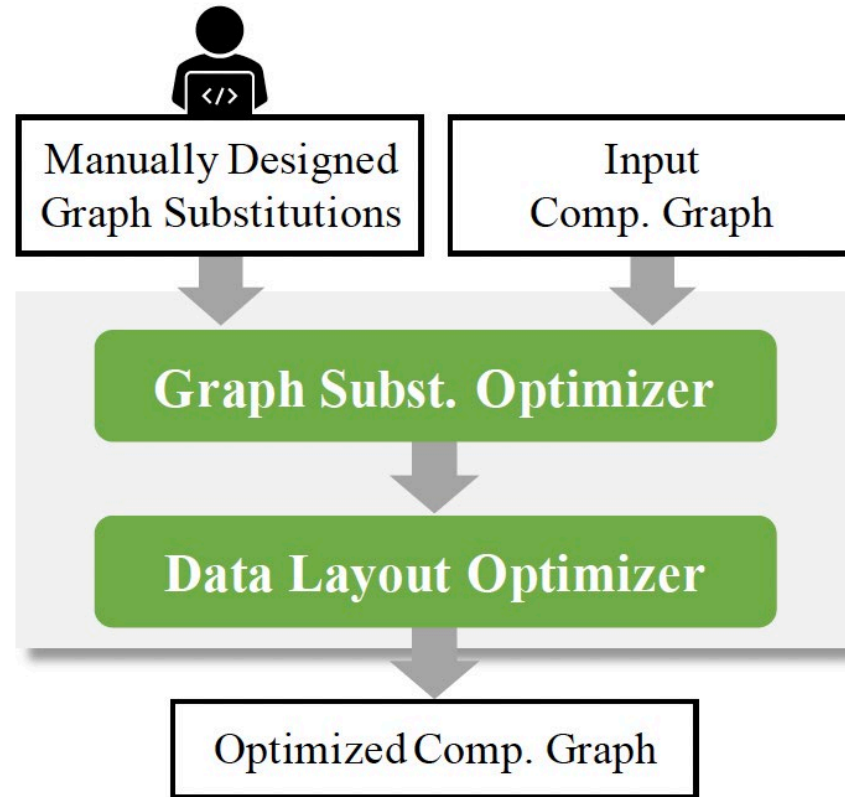Presented By

Pranav Talluri

# Background

- DNNs are expressed as computation graphs

- Multiple formulations can achieve the same goal, with differing costs

- Introduces the desire to optimise DNN computation graphs

- Before TASO, the specific optimisations were manually designed by human experts

- TASO automates the generation of graph substitutions in order to programmatically optimise DNN graphs
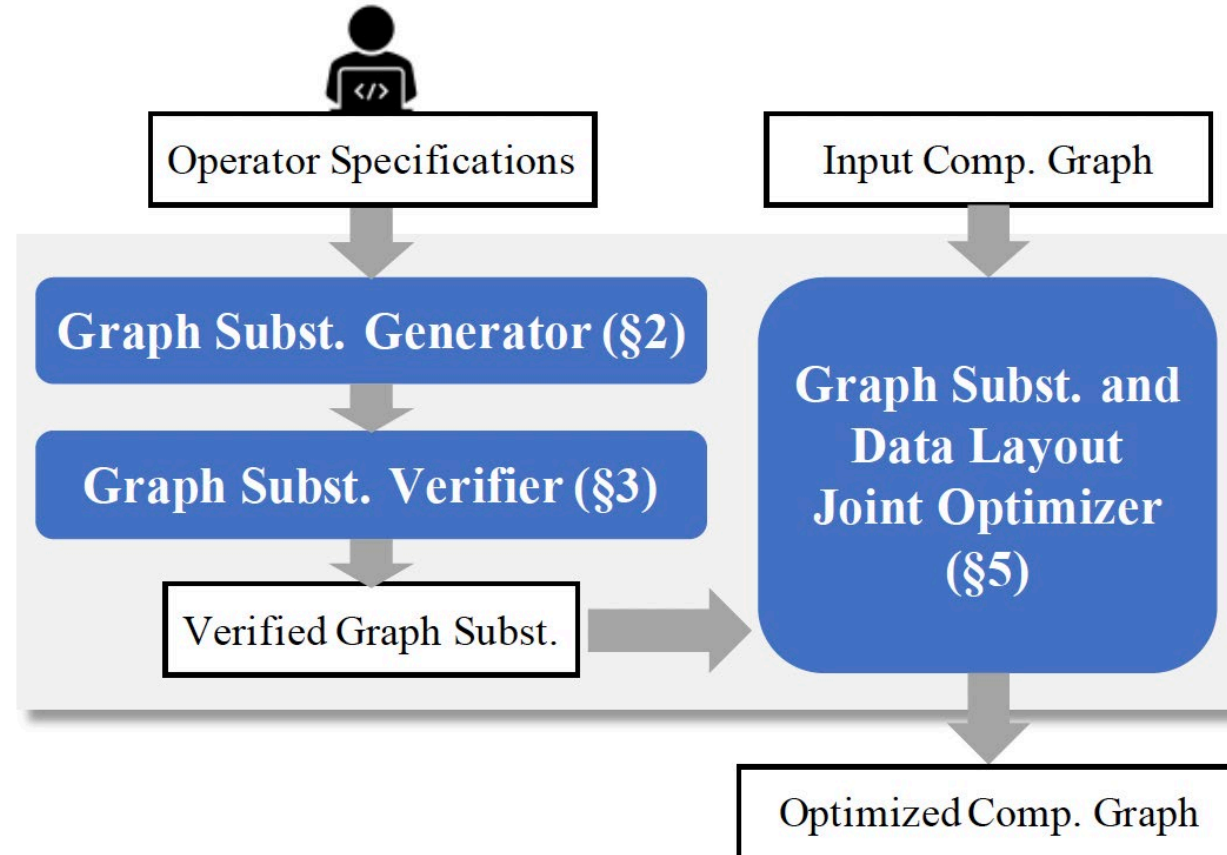
# Background

# Background

# Overview

- TASO automates generation of graph substitutions

- Framework agnostic (cuDNN + TVM)

- Takes operator specifications as an input

- Does so in a few stages:
  - Programmatically generate candidate graph substitutions
    - Generate
    - Quick test to prune impossible substitutions
  - Formally verify validity
  - Cost-based backtracking search to find an optimised graph
    - Includes co-optimisation of data locality

# Overview

# Approach: Generate Substitutions

- Substitution = source, target, mapping
- Configuration parameter dependent operators
- Generation algorithm
  - Enumerate potential graphs
  - Create graphs iteratively
  - Collect fingerprints
  - Test graphs with identical fingerprints
- Special Cases

---

**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:     $\mathcal{D} = \mathcal{D} + ($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** $\mathcal{S}$

# Approach: Formal Verification

- Verify generated substitutions

- Operator properties expressed in FOL
  - Manually written and reviewed
  - Further validated using symbolic execution
  - Properties are added when required
  - Checked for consistency and redundancies are removed

- Uses Z3 (SMT Solver)

- Shapes of tensors are not modelled

- Data layout not included

# Approach: Formal Verification

| Name | Description | Parameters |
|---|---|---|
| | Tensor Operators | |
| ewadd | Element-wise addition | |
| ewmul | Element-wise multiplication | |
| smul | Scalar multiplication | |
| transpose | Transpose | |
| matmul | Batch matrix multiplication[#] | |
| conv | Grouped convolution[%] | stride, padding, activation |
| enlarge | Pad conv. kernel with zeros[†] | kernel size |
| relu | Relu operator | |
| $pool_{avg}$ | Average pooling | kernel size, stride, padding |
| $pool_{max}$ | Max pooling | kernel size, stride, padding |
| concat | Concatenation of two tensors | concatenation axis |
| $split_{\{0,1\}}$ | Split into two tensors | split axis |
| | Constant Tensors | |
| $C_{pool}$ | Average pooling constant | kernel size |
| $I_{conv}$ | Convolution id. kernel | kernel size |
| $I_{matmul}$ | Matrix multiplication id. | |
| $I_{ewmul}$ | Tensor with 1 entries | |

# Approach: Formal Verification

| 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 |

...

# Approach: Pruning Redundant Substitutions

- Redundant substitutions are subsumed by more general, valid substitutions
- Input tensor renaming
- Common subgraph



source graph: A x (B x A)

target graph: (A x B) x A

source graph: A + (B x C)

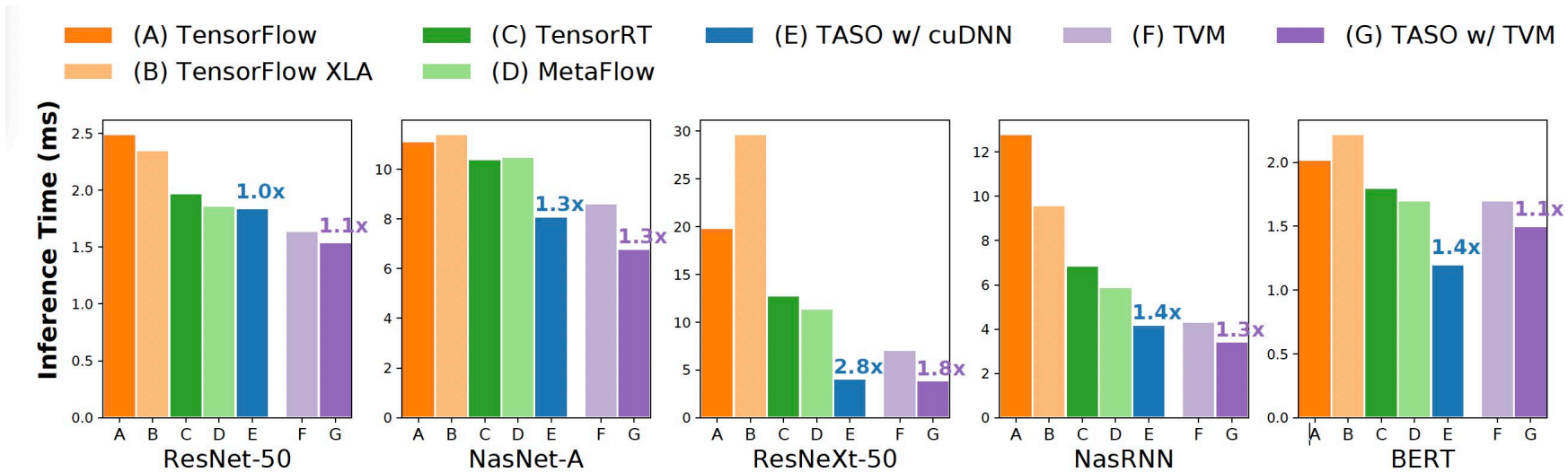target graph: (B x C) + A

# Approach: Joint Optimisation

- Utilises MetaFlow cost-based backtracking search algorithm

- Considers data layout optimisation opportunities

- Joint optimisation uncovers otherwise impossible optimisations

- Costs are given by execution times of specific operators

- Cycle removal

- Alpha parameter prunes search space

---

**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:         // LAYOUT$(\mathcal{G}, s)$ returns possible layouts applying $s$ on $\mathcal{G}$.
9:         **for** layout $l \in$ LAYOUT$(\mathcal{G}, s)$ **do**
10:             // APPLY$(\mathcal{G}, s, l)$ applies $s$ on $\mathcal{G}$ with layout $l$.
11:             $\mathcal{G}' = $ 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:                     $\mathcal{P}.\text{enqueue}(\mathcal{G}')$
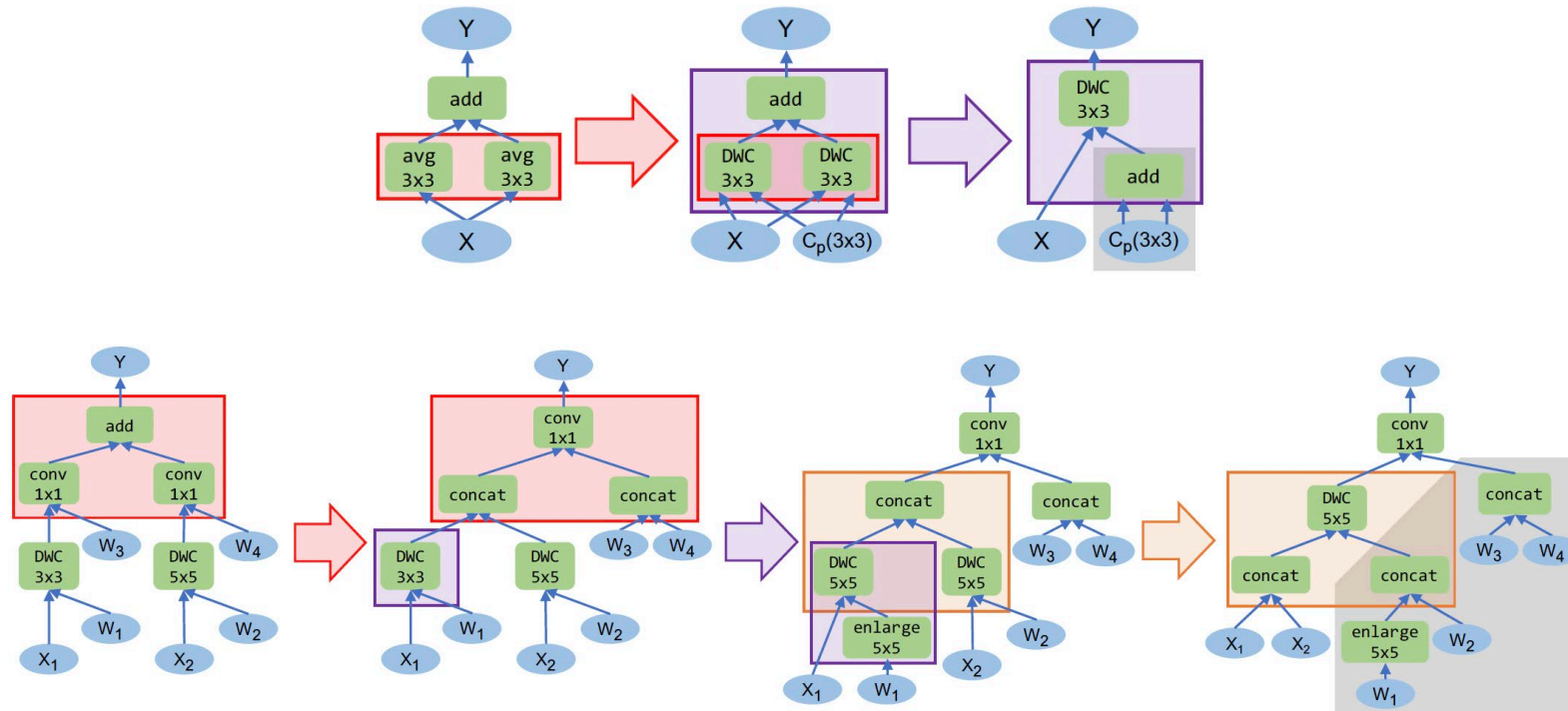17: **return** $\mathcal{G}_{opt}$

# Evaluation: Optimisation

- Setup – tested on 5 DNNs

- Successful automatic optimisation – inference time reduction
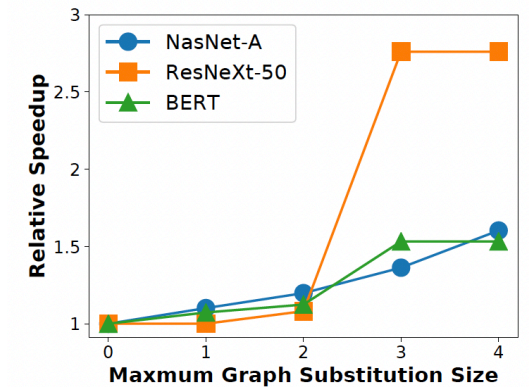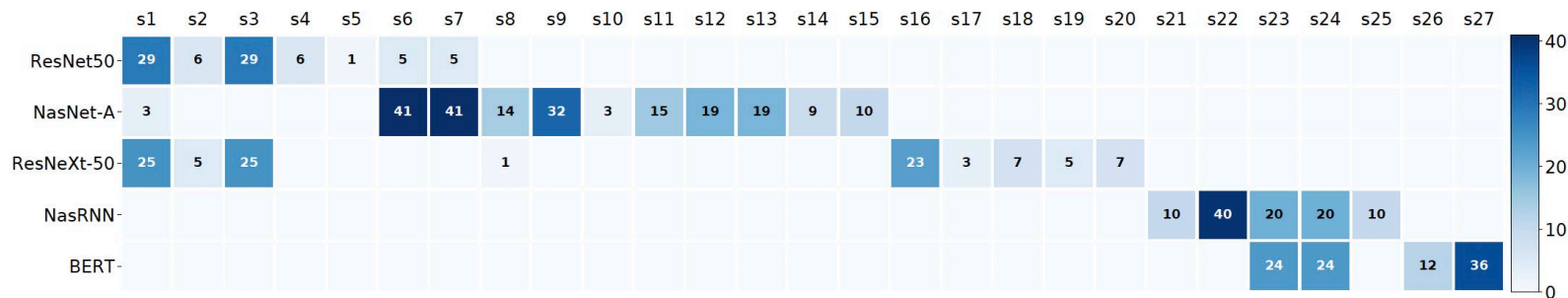  - cuDNN: 1.3x to 2.8x
  - TVM:    1.1x to 1.8x

# Evaluation: Substitutions

- NasNet was produced using neural architecture search
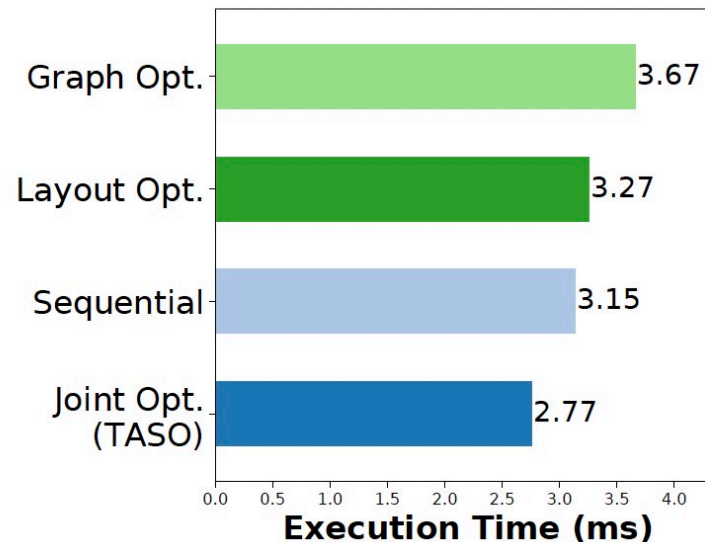- Unconventional optimisations were discovered

# Evaluation: Substitutions

- Different DNNs used different optimisations, showing usefulness of TASO

- Scalability
  - Larger operator substitutions could be useful

# Evaluation: Substitutions

- Joint optimisation
  - Better than individual or sequential
  - $(A \times B) \to ((B^T \times A^T)^T)$ with $B^T$ in row-major and $A^T$ in column-major
  - Phase ordering?
- Relatively quick - <10 minutes for each DNN

# Review

## Positives

- Novel idea

- Successful execution
  - Improves DNN performance
  - Reduces human effort
  - Extensible framework

- Seminal work in an exciting research area
  - Graph transformation backend still in use

## Negatives

- Reliant on user provided operator properties

- Scalability of generator

- Phase ordering problem + search procedure

- Cost model has issues

# Future Works

- Future works have built on this approach
- PET
  - Partially equivalent optimisations
- TENSAT
  - Equality saturation
- X-RLflow
  - RL approach to searching optimisation space
- REGAL
  - Transfer knowledge