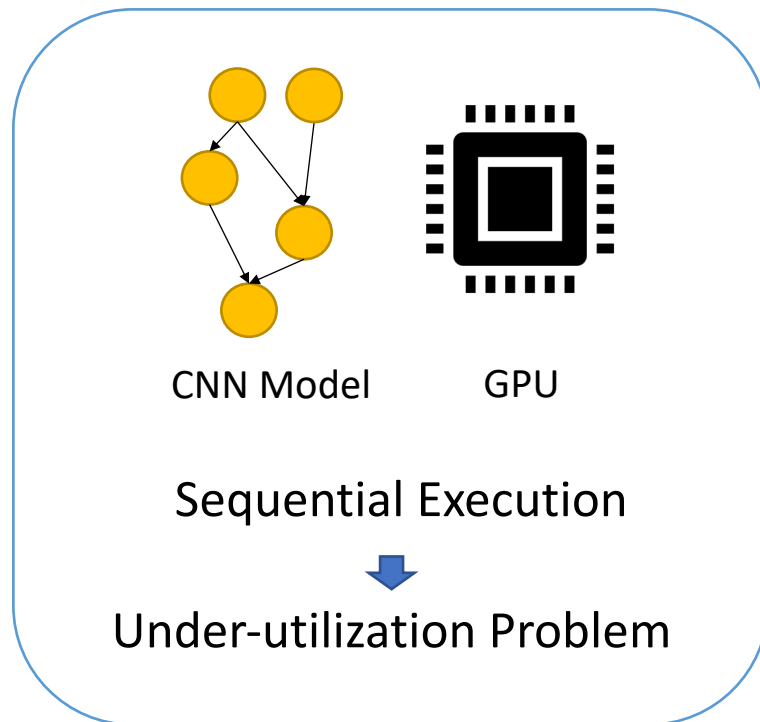# IOS: Inter-Operator Scheduler for CNN Acceleration

**Yaoyao Ding**[1,2], Ligeng Zhu[3], Zhihao Jia[4],

Gennady Pekhimenko[1,2], Song Han[3]
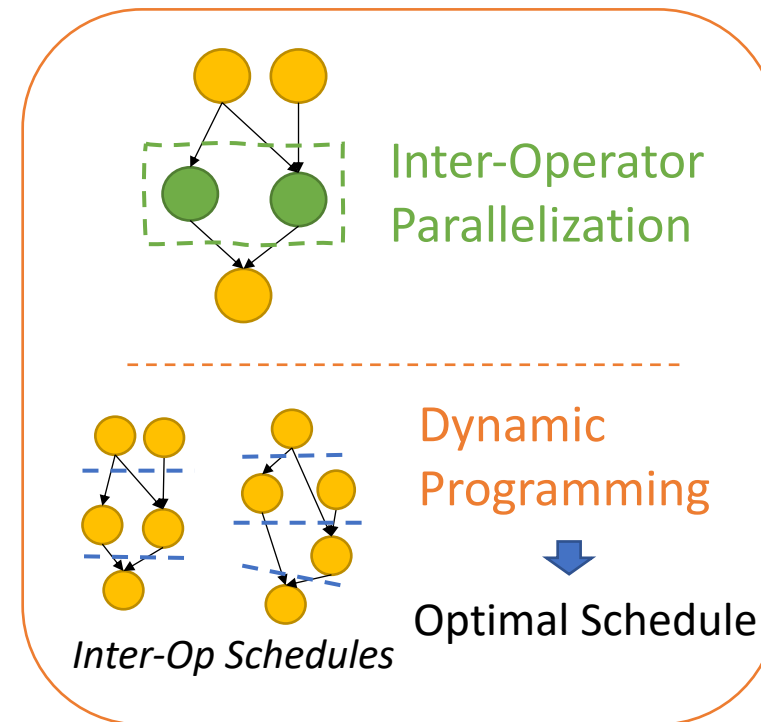
[1] UNIVERSITY OF TORONTO

[2] VECTOR INSTITUTE

[3] MIT

[4] Carnegie Mellon University

# Executive Summary

**Motivation**



CNN Model          GPU

Sequential Execution

Under-utilization Problem

**Inter-Operator Scheduler**



Inter-Operator Parallelization

Dynamic Programming

Inter-Op Schedules          Optimal Schedule
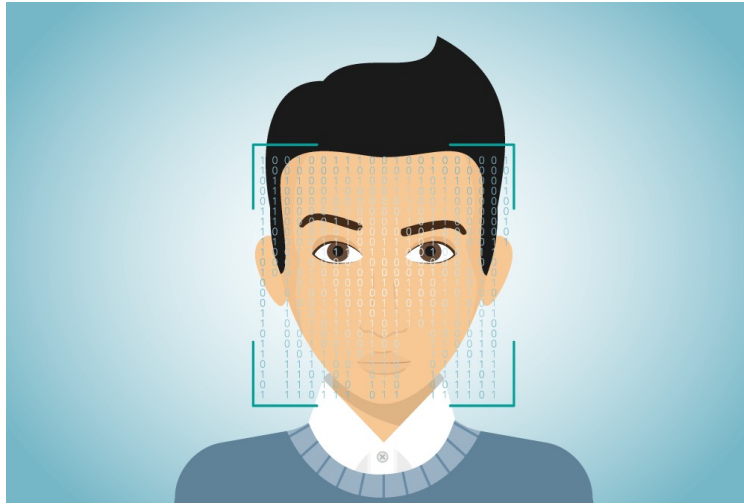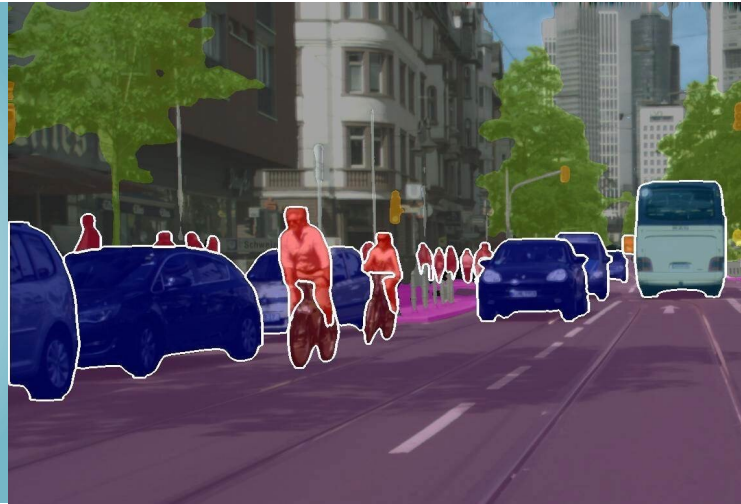
**1.1-1.5x speedup**

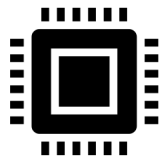# Efficient Deployment of CNNs is Important



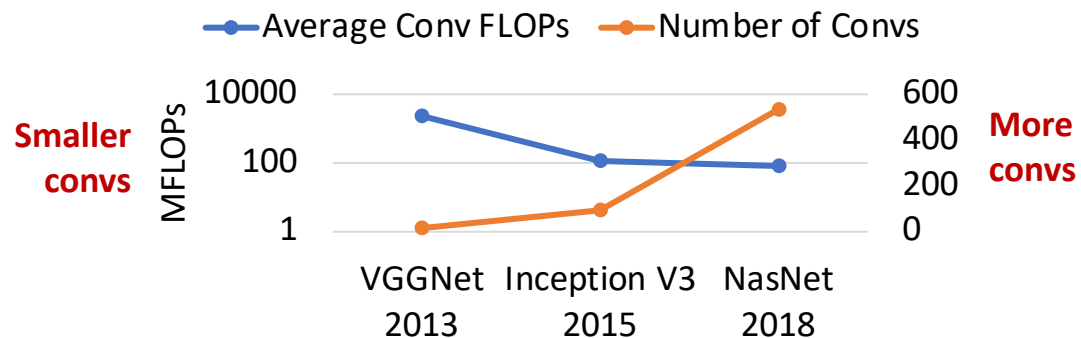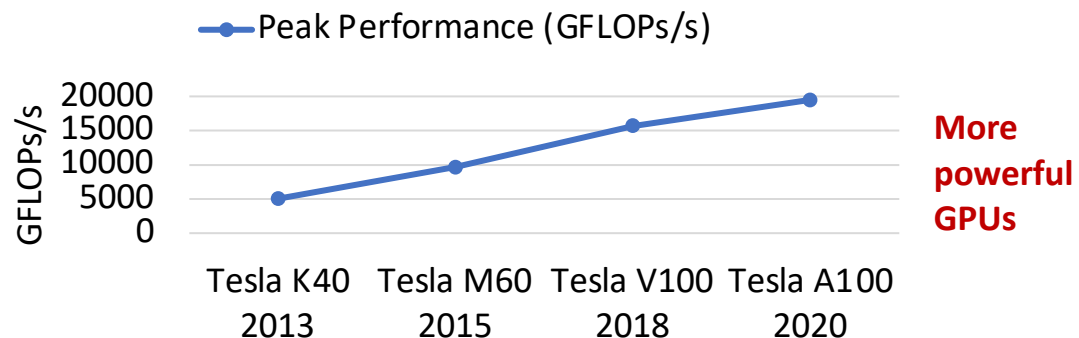Face Recognition

Self Driving

Language Translation

Is CNN inference in current DL libraries well utilizing underlying hardware?
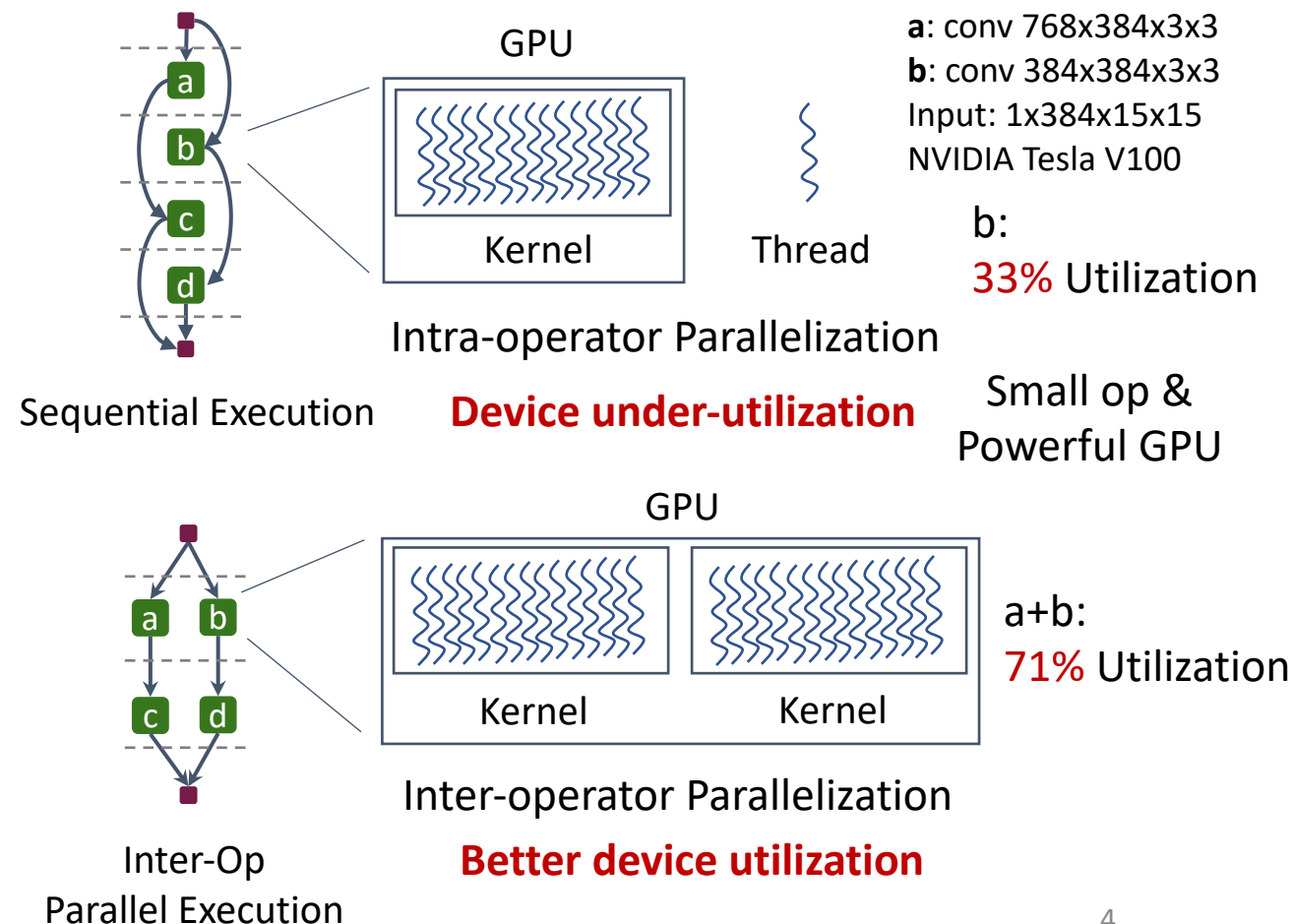
# Motivation for Inter-Operator Parallelization

## 1. More small convs in CNN design.



Average Conv FLOPs — Number of Convs

**Smaller convs**

**More convs**

VGGNet 2013    Inception V3 2015    NasNet 2018

## 2. GPU peak performance increased



Peak Performance (GFLOPs/s)

**More powerful GPUs**

Tesla K40 2013    Tesla M60 2015    Tesla V100 2018    Tesla A100 2020

## 3. Intra- and Inter-operator Parallelization



**a**: conv 768x384x3x3
**b**: conv 384x384x3x3
Input: 1x384x15x15
NVIDIA Tesla V100

GPU

Kernel

Thread

b:
**33% Utilization**

Sequential Execution

Intra-operator Parallelization

**Device under-utilization**

Small op & Powerful GPU

GPU

Kernel    Kernel

a+b:
**71% Utilization**

Inter-Op Parallel Execution

Inter-operator Parallelization

**Better device utilization**
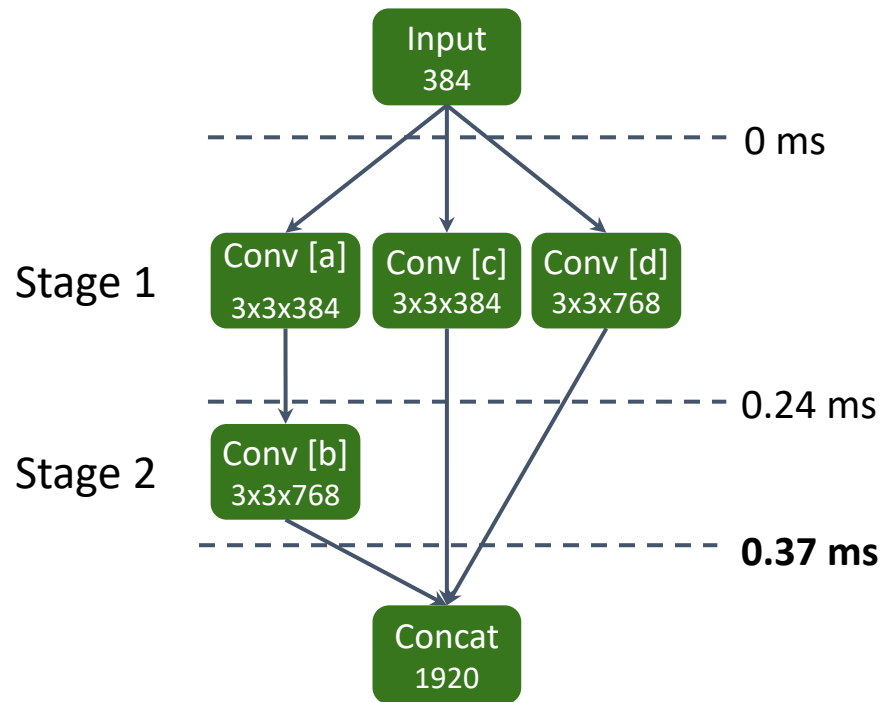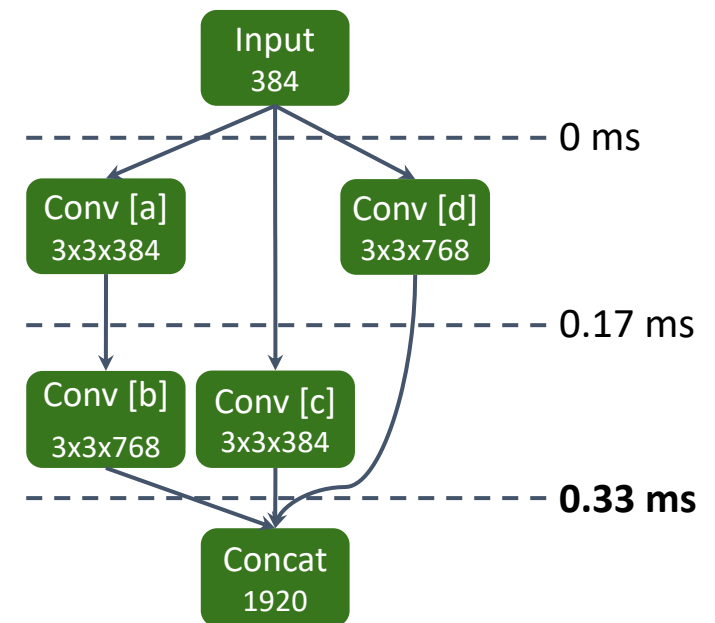
# Background: Wavefront Schedule Policy

Wavefront Schedule Policy: Execute **all available** operators stage by stage
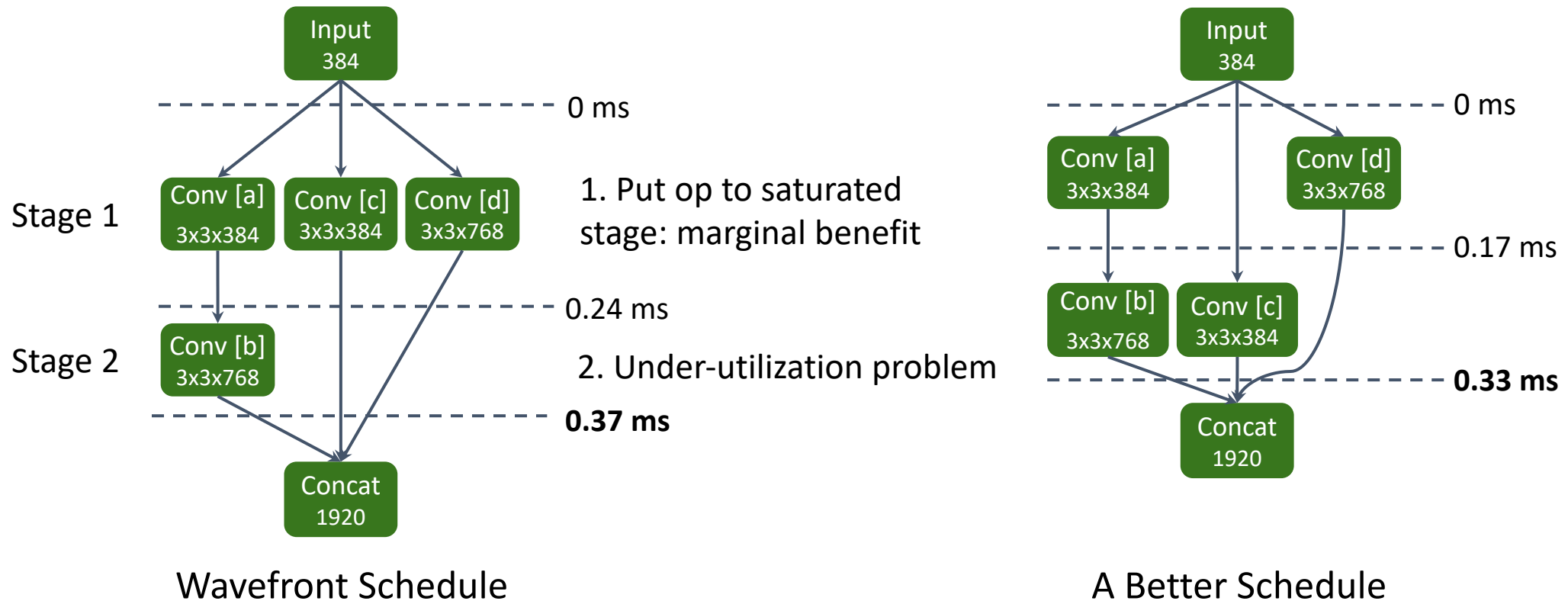


Wavefront Schedule

Move Conv [c] from Stage 1 to Stage 2

A Better Schedule

# Background: Wavefront Schedule Policy
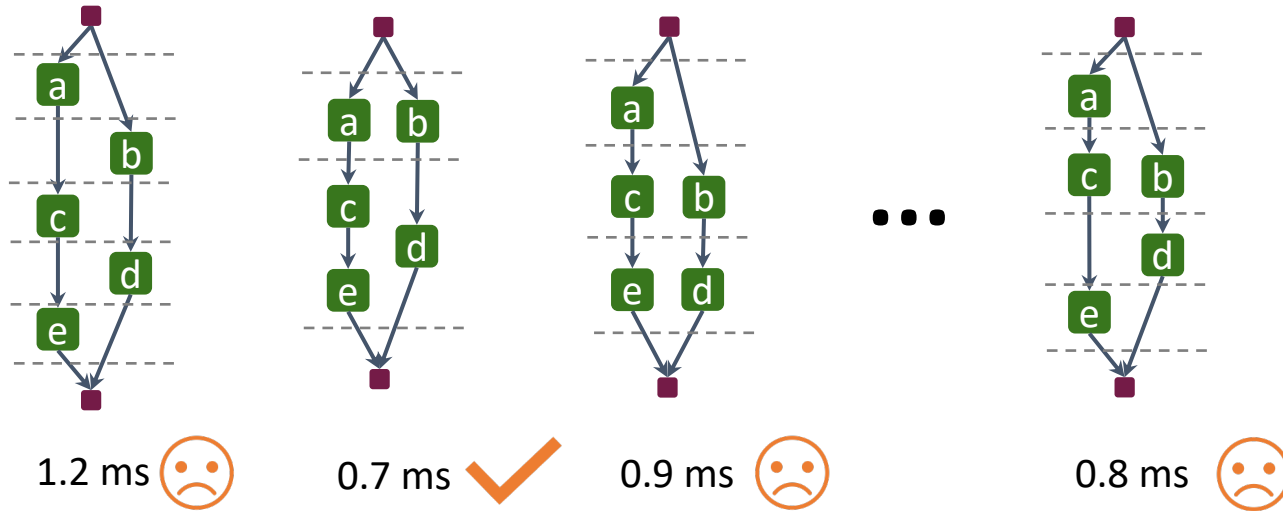
Wavefront Schedule Policy: Execute **all available** operators stage by stage



1. Put op to saturated stage: marginal benefit

2. Under-utilization problem

Wavefront Schedule

A Better Schedule

**Wavefront schedule policy is sub-optimal**

# Inter-Operator Scheduler (IOS)

**General Idea:** Explore the schedule space **exhaustively**.



1.2 ms      0.7 ms ✓     0.9 ms     ...     0.8 ms

# Inter-Operator Scheduler (IOS)

**General Idea:** Explore the schedule space **exhaustively**.

**Challenge:** The number of schedules is **exponential** in the number of operators.

e.g., NASNet has more than $10^{12}$ schedules

**Prohibitive to enumerate**

# Inter-Operator Scheduler (IOS)

**General Idea:** Explore the schedule space **exhaustively**.

**Challenge:** The number of schedules is **exponential** in the number of operators.

**Observation 1**: Optimal schedule for a subgraph can be reused



Model

Incomplete Schedule

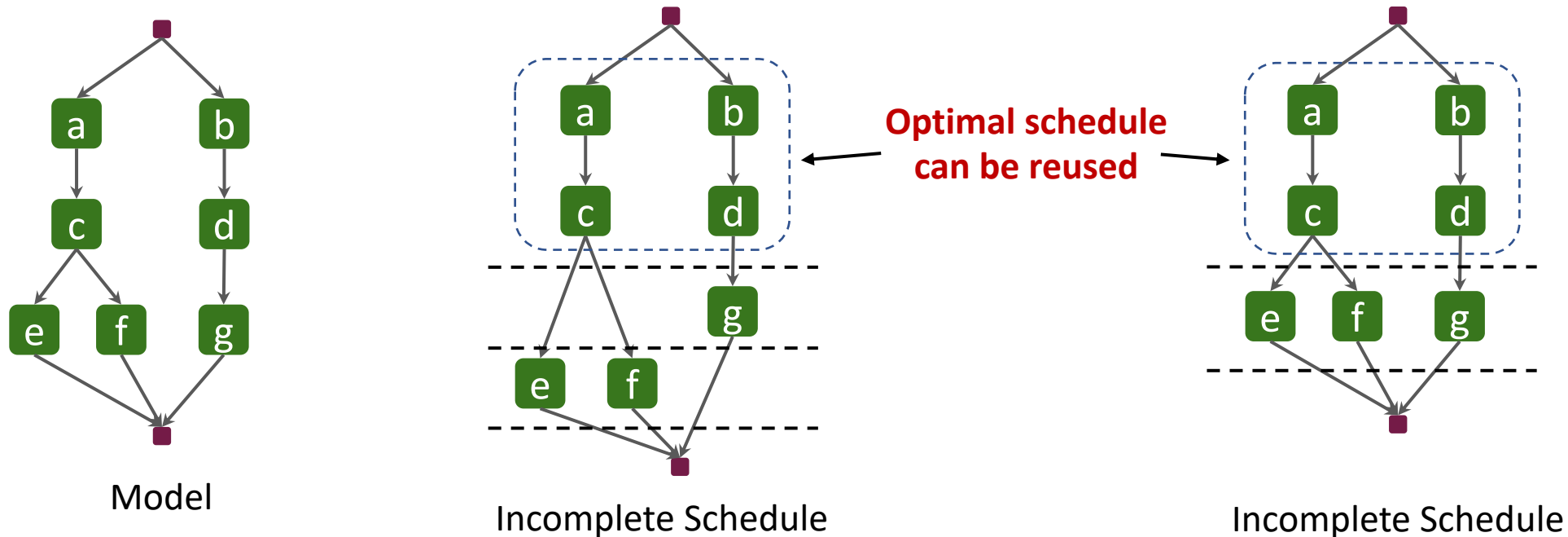**Optimal schedule can be reused**

Incomplete Schedule

# Inter-Operator Scheduler (IOS)

**General Idea:** Explore the schedule space **exhaustively**.

**Challenge:** The number of schedules is **exponential** in the number of operators.

**Observation 1**: Optimal schedule for a subgraph can be reused

    **Key Idea**: Dynamic Programming
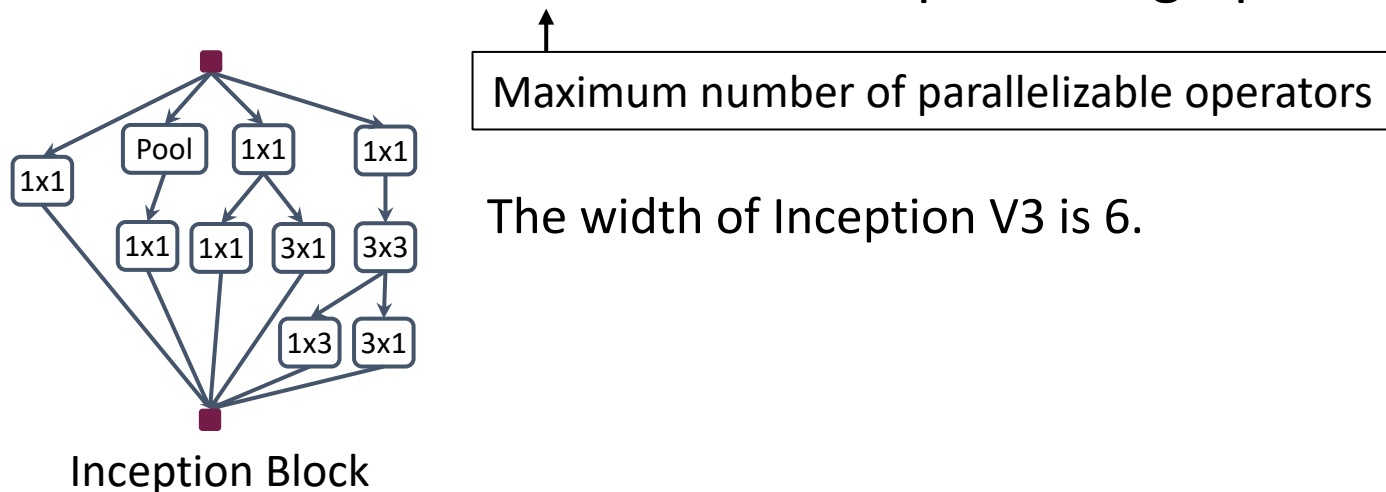
# Inter-Operator Scheduler (IOS)

**General Idea:** Explore the schedule space **exhaustively**.

**Challenge:** The number of schedules is **exponential** in the number of operators.

**Observation 1**: Optimal schedule for a subgraph can be reused

      **Key Idea**: Dynamic Programming

**Observation 2**: The width of the computation graph is usually small



Maximum number of parallelizable operators

The width of Inception V3 is 6.

Inception Block

# Inter-Operator Scheduler (IOS)

**General Idea:** Explore the schedule space **exhaustively**.

**Challenge:** The number of schedules is **exponential** in the number of operators.
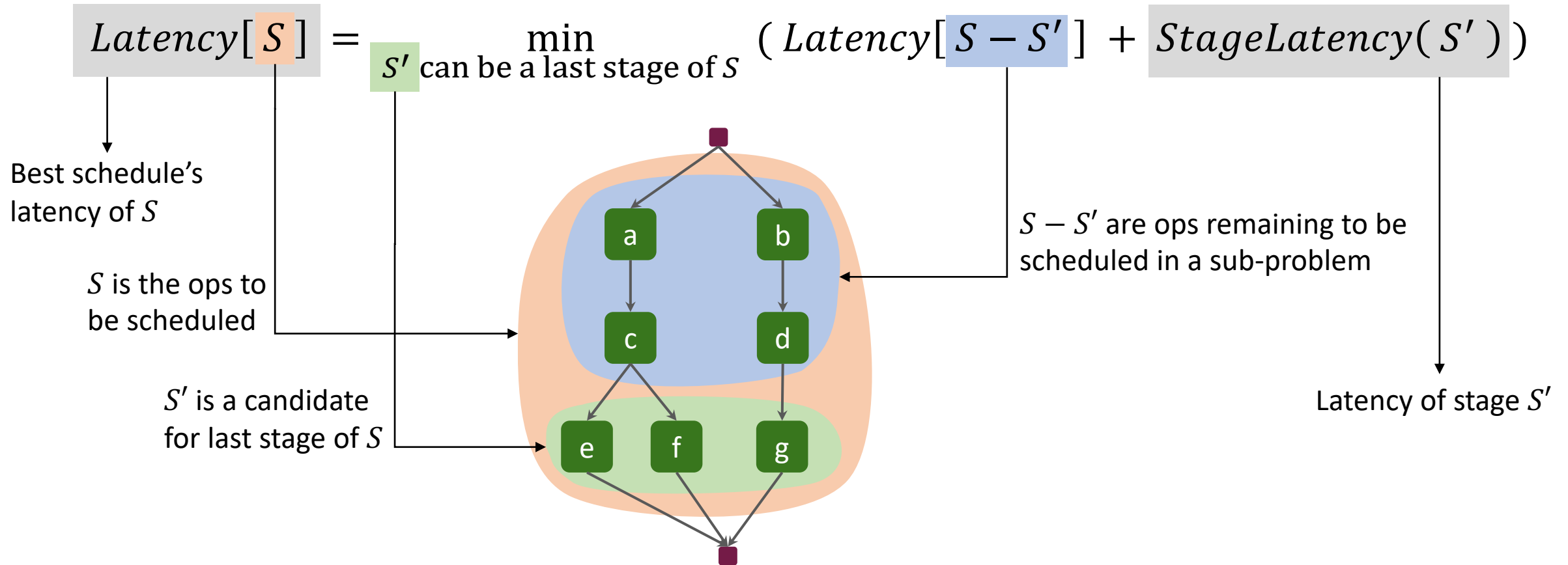
**Observation 1**: Optimal schedule for a subgraph can be reused

   **Key Idea**: Dynamic Programming

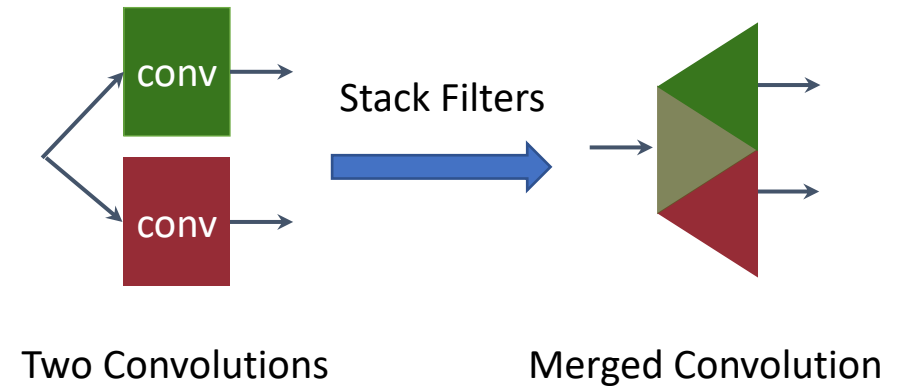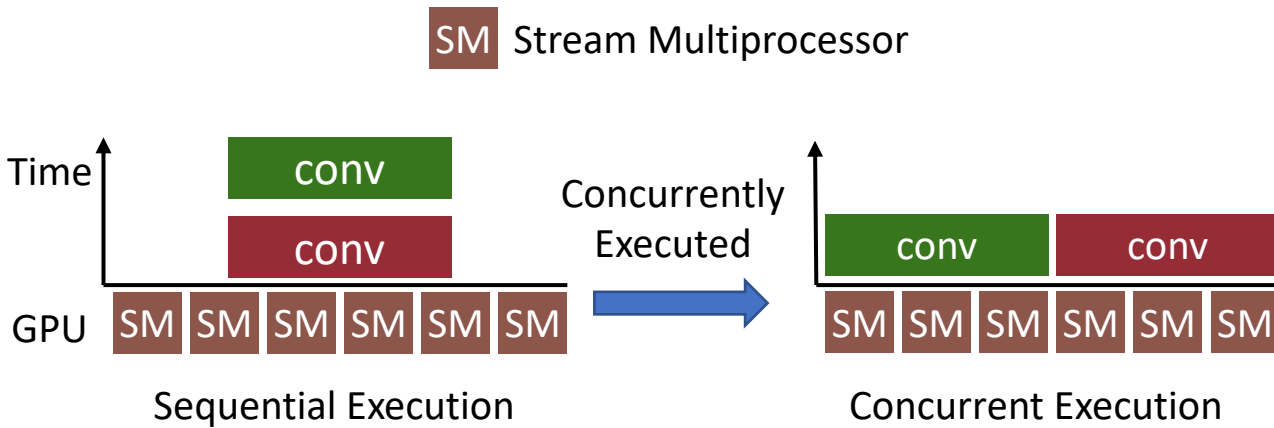**Observation 2**: The width of the computation graph is usually small

   **Key Result**: Time complexity is only exponential in the width

# Inter-Operator Scheduler (IOS)

$$Latency[\,S\,] = \min_{S' \text{ can be a last stage of } S} (\,Latency[\,S - S'\,] + StageLatency(\,S'\,)\,)$$

Best schedule's latency of $S$

$S$ is the ops to be scheduled

$S'$ is a candidate for last stage of $S$

$S - S'$ are ops remaining to be scheduled in a sub-problem

Latency of stage $S'$

# Parallelization Strategy Selection

$$Latency[\,S\,] = \min_{S' \text{ can be a last stage of } S} (\,Latency[\,S - S'\,] + StageLatency(\,S'\,))$$



SM  Stream Multiprocessor

Time

GPU

Concurrently Executed

Sequential Execution

Concurrent Execution

Stack Filters

Two Convolutions

Merged Convolution

Concurrent Execution

**Profile & Select**
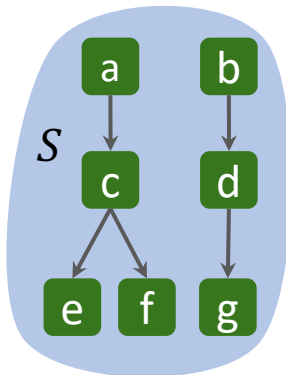
General
Sub-optimal performance

Operator Merge

Specialized
Usually better performance

# Last Stage Candidates

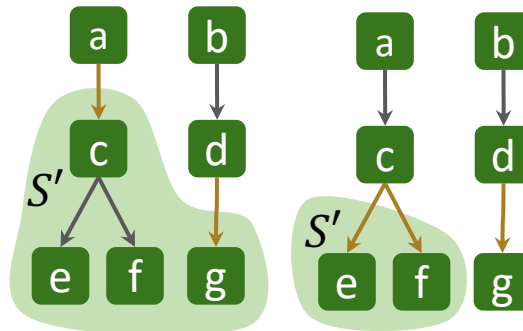$$Latency[\,S\,] = \min_{\substack{S' \text{ can be a last stage of } S}} (\,Latency[\,S - S'\,] + StageLatency(S')\,)$$

$S'$ can be a last stage of $S \iff$ There is no edge from $S'$ to $S - S'$
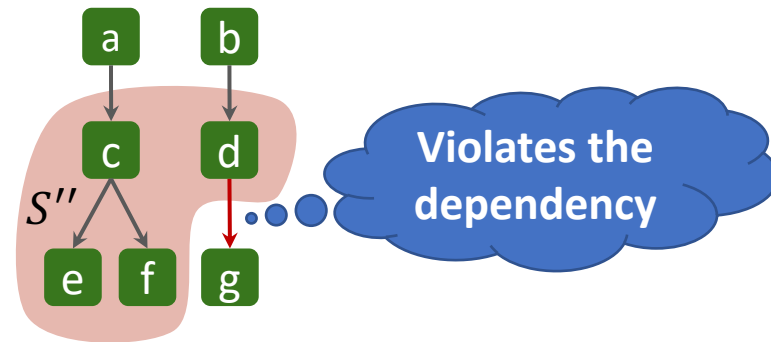


Operators $S$ to be scheduled

Edge from $S - S'$ to $S'$

$S'$ can be a last stage of $S$

Edge from $S'$ to $S - S'$

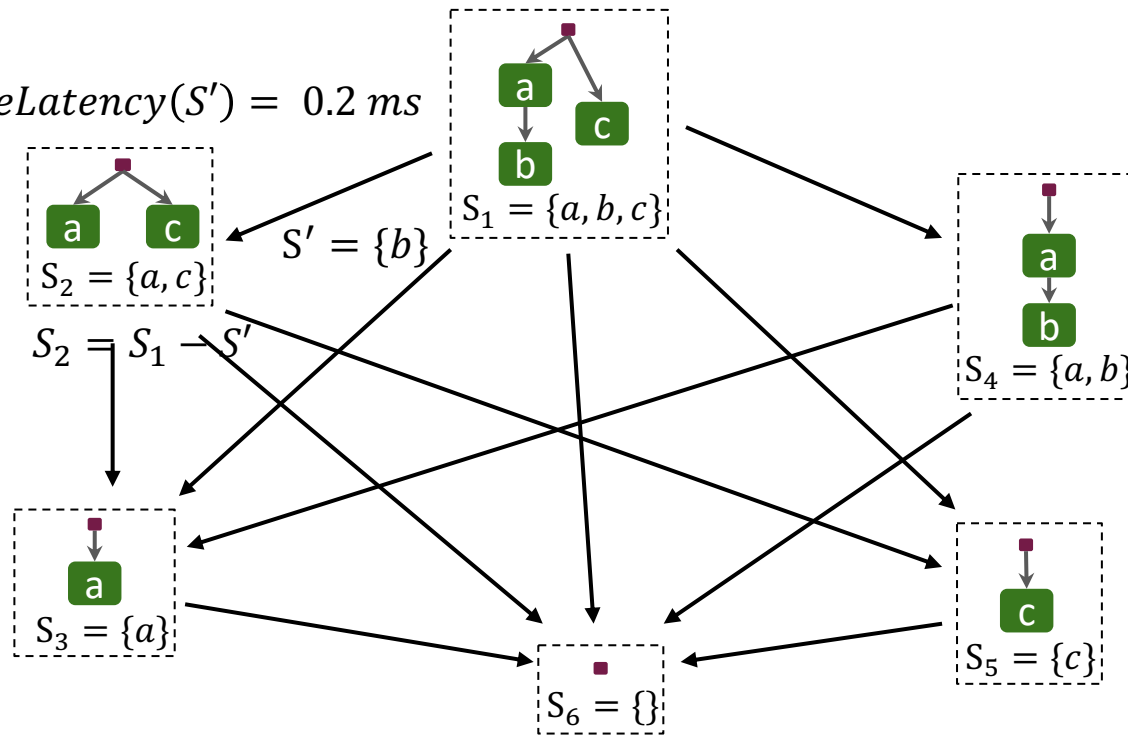Violates the dependency

$S''$ can **NOT** be a last stage of $S$

# Transition Graph and Time Complexity

$$Latency[\,S\,] = \min_{S' \text{ can be a last stage of } S} (\,Latency[\,S - S'\,] + StageLatency(\,S'\,)\,)$$

**Vertices:** all valid *state S*

**Edges:** $S \rightarrow (S - S')$

$StageLatency(S') = 0.2\ ms$



A Simple Model

$S_2 = \{a, c\}$

$S_1 = \{a, b, c\}$

$S' = \{b\}$

$S_2 = S_1 - S'$

$S_4 = \{a, b\}$
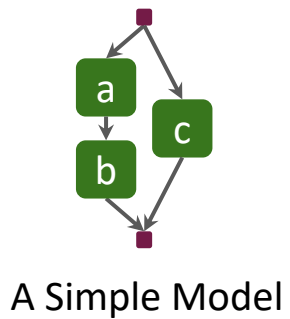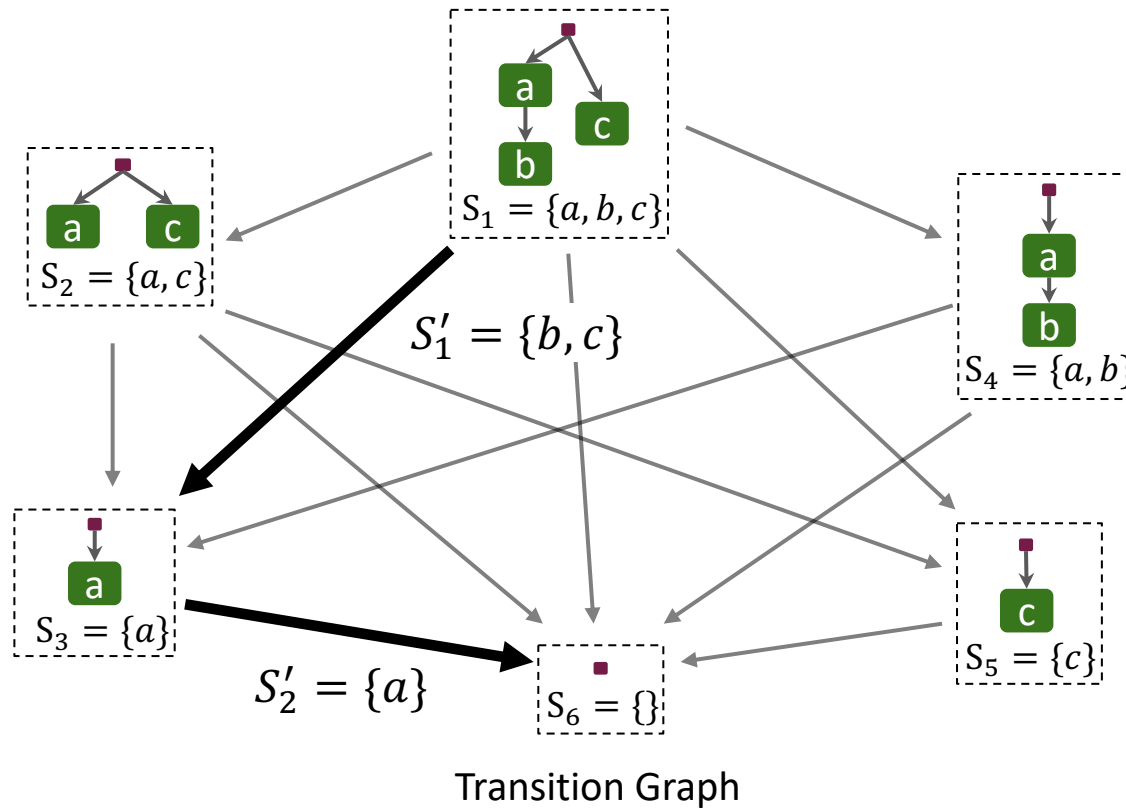
$S_3 = \{a\}$

$S_6 = \{\}$

$S_5 = \{c\}$

Transition Graph

# Transition Graph and Time Complexity

$$Latency[\,S\,] = \min_{S' \text{ can be a last stage of } S} (\,Latency[\,S - S'\,] + StageLatency(\,S'\,)\,)$$

**Vertices:** all valid state $S$

**Edges:** $S \rightarrow (S - S')$



A Simple Model

$S_1 = \{a, b, c\}$

$S_2 = \{a, c\}$

$S_4 = \{a, b\}$

$S_1' = \{b, c\}$

$S_3 = \{a\}$

$S_5 = \{c\}$

$S_2' = \{a\}$

$S_6 = \{\}$

Transition Graph

IOS: Find the **shortest** path

$S_2' = \{a\}$
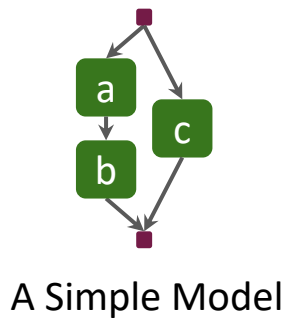
$S_1' = \{b, c\}$

Any path from $S_1$ to $S_6$ is a schedule
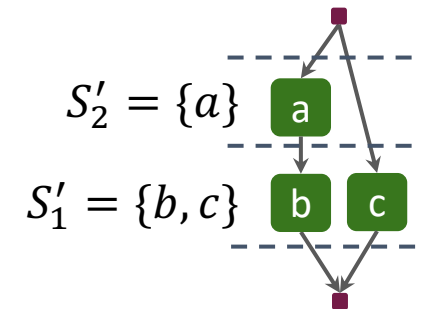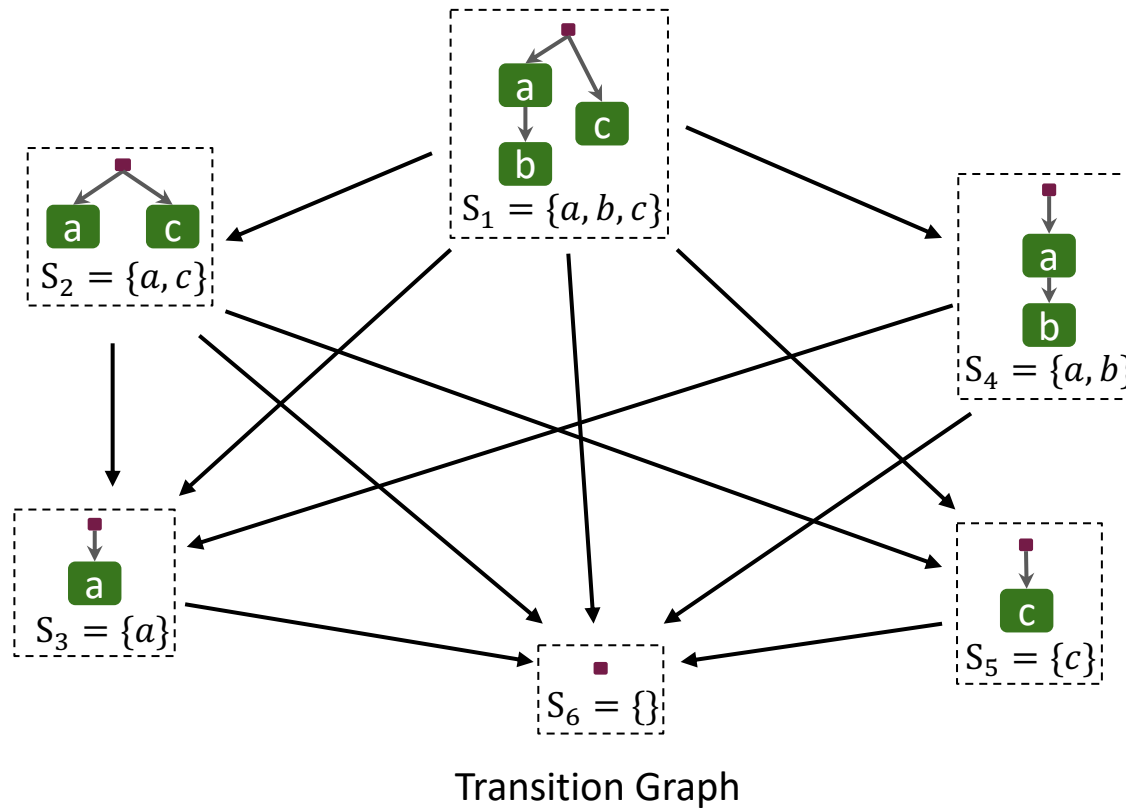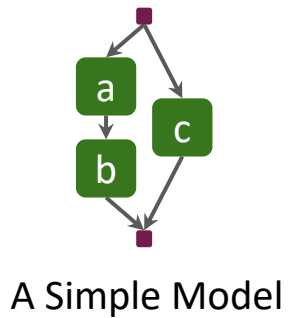
# Transition Graph and Time Complexity

$$Latency[\,S\,] = \min_{S' \text{ can be a last stage of } S} (\,Latency[\,S - S'\,] + StageLatency(\,S'\,)\,)$$

**Vertices:** all valid state $S$

**Edges:** $S \to (S - S')$



A Simple Model

$S_1 = \{a, b, c\}$

$S_2 = \{a, c\}$

$S_4 = \{a, b\}$

$S_3 = \{a\}$

$S_5 = \{c\}$

$S_6 = \{\}$

Transition Graph

IOS: Find the **shortest** path

**Time Complexity** of IOS:

$$\mathcal{O}\left((n/d + 1)^{2d}\right)$$

$\boldsymbol{n}$: number of operators

$\boldsymbol{d}$: max number of parallelizable ops

18

# Methodology

- **Benchmarks**
  - Inception V3 ⎤
  - SqueezeNet ⎦ Expert Designed
  - Randwire ⎤
  - NasNet ⎦ Neural Architecture Search

- **Environment**

   NVIDIA Tesla V100

   10.2   7.6.5

- **Baselines**
  - State-of-the-art Frameworks (cuDNN-based)

     TensorFlow   **TASO**

  - Different schedules on IOS Runtime

- **IOS Implementation**
  - cuDNN kernels
  - CUDA Stream

# Comparison of cuDNN-based Frameworks

**Tensorflow:** A popular machine learning framework.

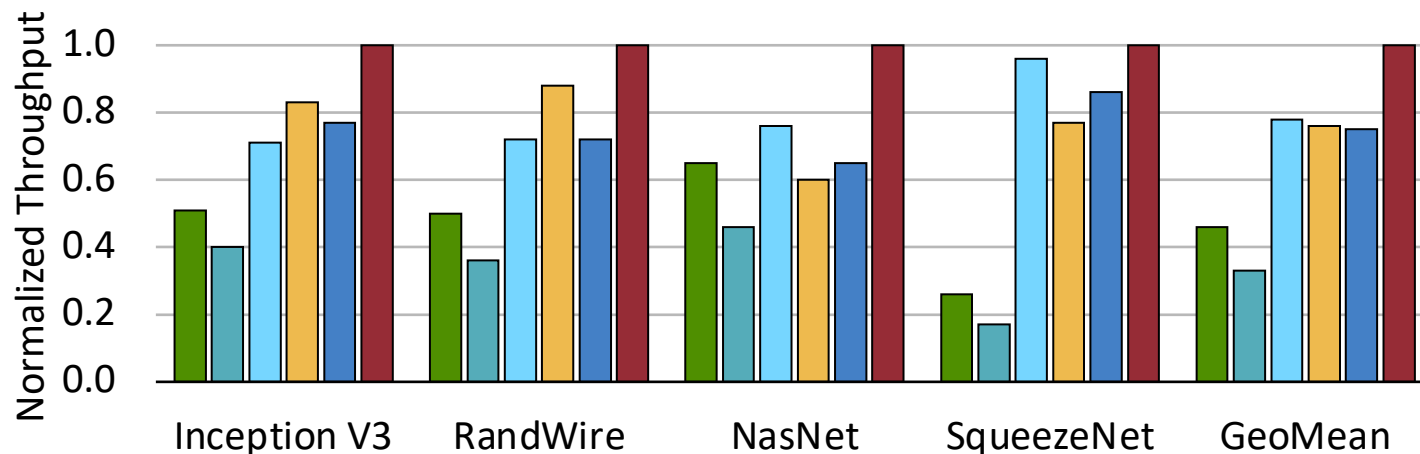**Tensorflow-XLA:** TensorFlow with compilation optimization.

**TASO:** Transformation-based optimizer.

**TVM-cuDNN:** TVM backed with cuDNN convolution kernel.

**TensorRT:** NVIDIA high-performance inference engine.

**IOS:** Our method

Under-utilization due to sequential execution



Performance is normalized to the best framework

**IOS** outperforms all frameworks and achieves **1.1-1.5x** speedup.

# Comparison of Different Schedules

**Under-utilize Device**

**Sequential:** Run each op sequentially.

**Unbalanced Schedule**

**Wavefront:** Run all available ops stage by stage.

**IOS Runtime**

**IOS-Merge:** IOS with only "operator merge" policy.

**No trade off between parallelization strategies**

**IOS-Parallel:** IOS with only "parallel execution" policy.

**IOS-Both:** IOS with both policies
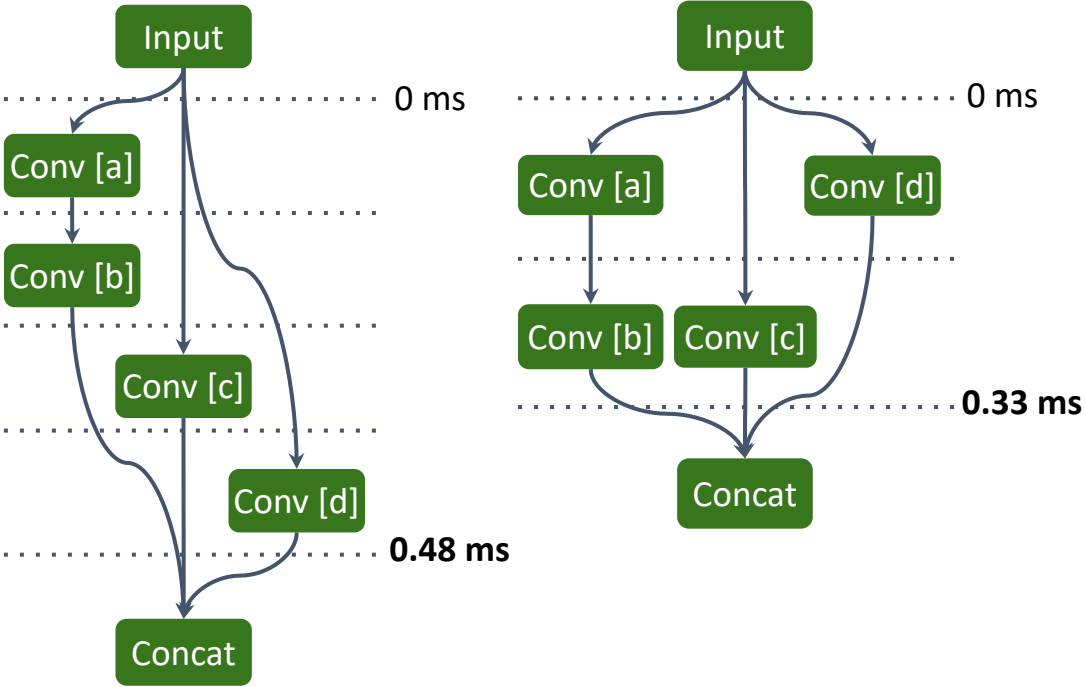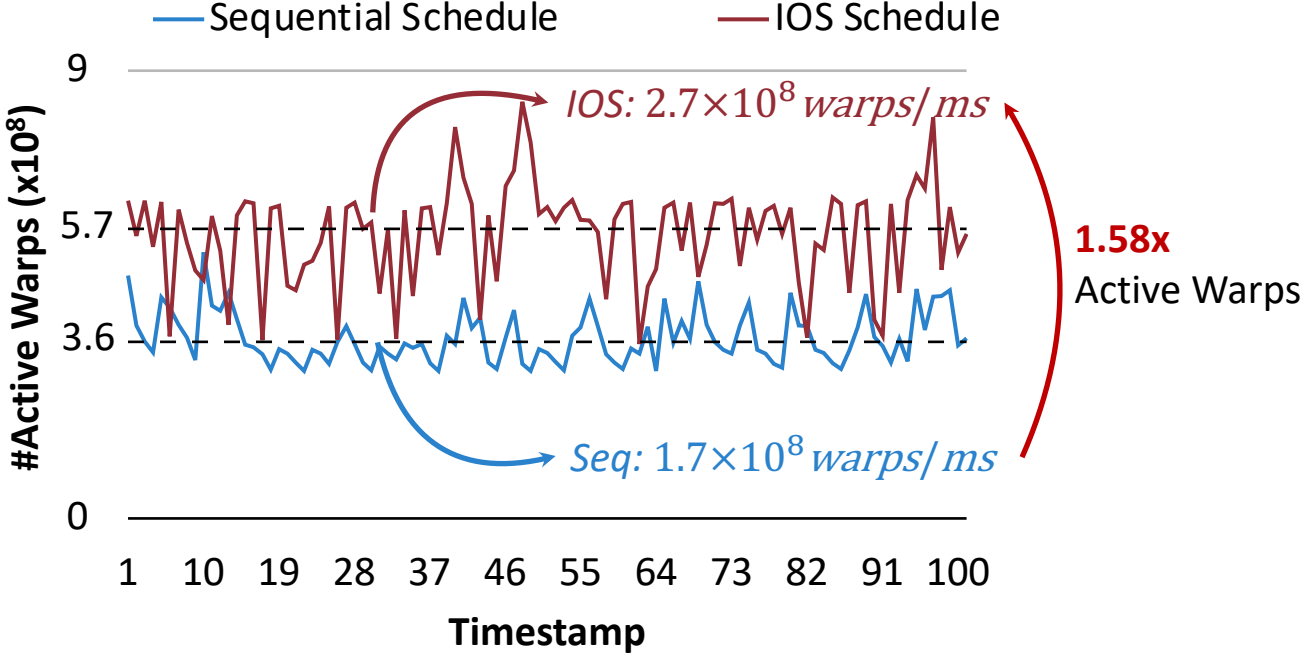


Performance is normalized to the best schedule

**IOS-Both** achieves the best performance

# More Active Warps Improve Utilization



Sequential Schedule

IOS Schedule

$IOS: 2.7 \times 10^8 \ warps/ms$

$Seq: 1.7 \times 10^8 \ warps/ms$

**1.58x** Active Warps

#Active Warps (x10⁸)

Timestamp

*NVIDIA CUPTI profile frequency is every 2.1 ms.*

More active warps ➡ More eligible warps to execute at each cycle ➡ Higher Device-Utilization

# Conclusion

- Sequential execution suffers from **under utilization** problem.

- **I**nter-**O**perator **S**cheduler (**IOS**):
  - Utilize both intra- and **inter-operator parallelism** in CNNs.
  - **Dynamic-programming** explores the schedule space **exhaustively**.
  - Time Complexity: $\mathcal{O}\left((n/d + 1)^{2d}\right)$, $d$ is usually small.

- Key Results: **1.1-1.5x** speedup on diverse CNNs.

https://github.com/mit-han-lab/inter-operator-scheduler
We provide scripts to reproduce results in every figure and table!