

# IOS: Inter-Operator Scheduler for CNN Acceleration

Yaoyao Ding<sup>1 2</sup>, Ligeng Zhu<sup>3</sup>, Zhihao Jia<sup>4</sup>,  
Gennady Pekhimenko<sup>1 2</sup>, Song Han<sup>3</sup>

1



2



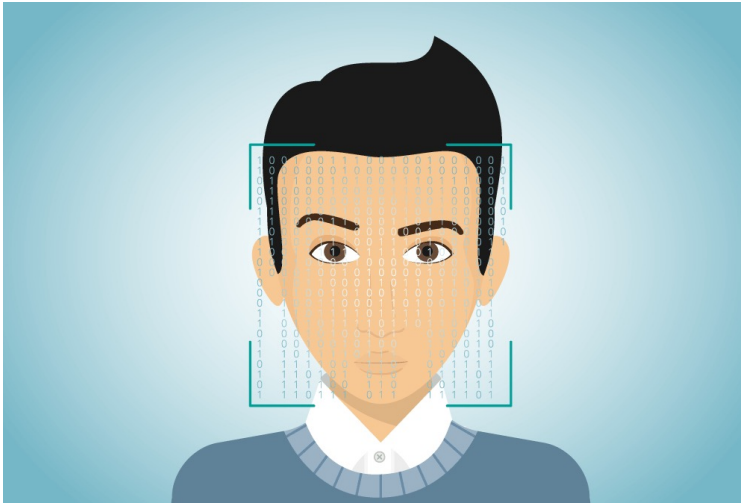
3



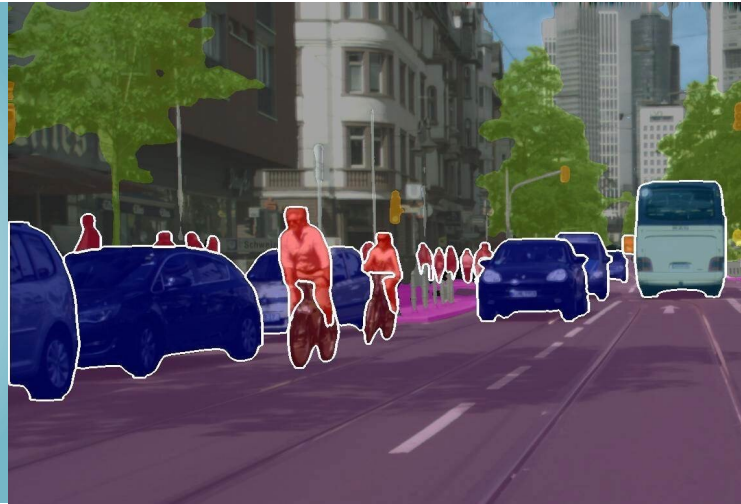
4



# Efficient Deployment of CNNs is Important



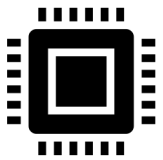
Face Recognition



Self Driving



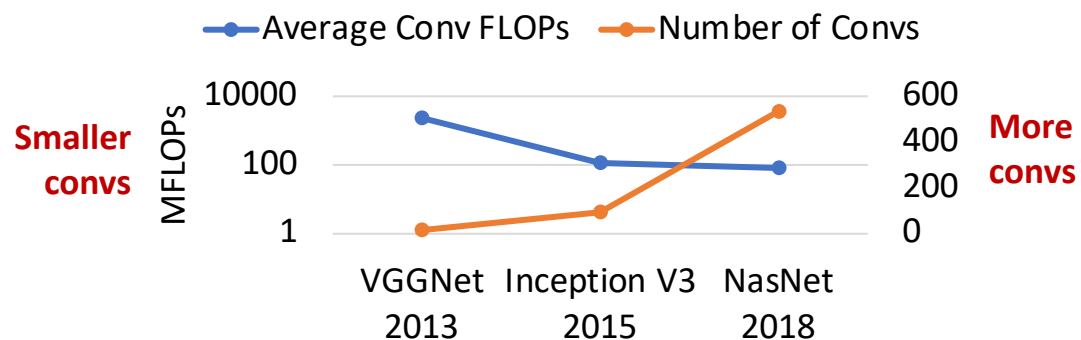
Language Translation



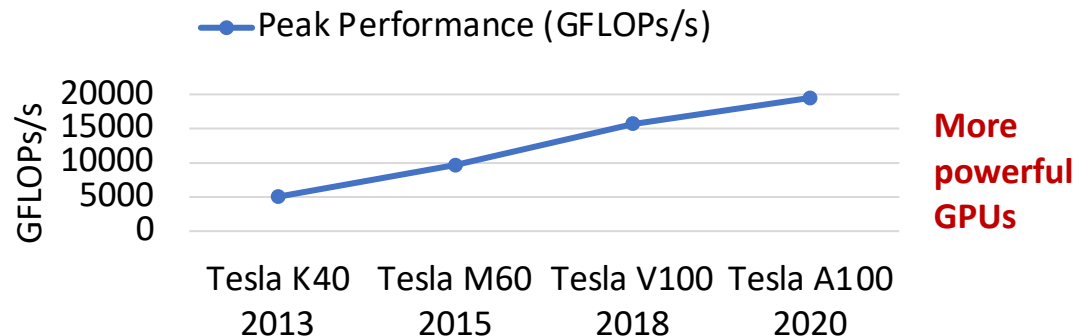
Is CNN inference in current DL libraries well utilizing underlying hardware?

# Motivation for Inter-Operator Parallelization

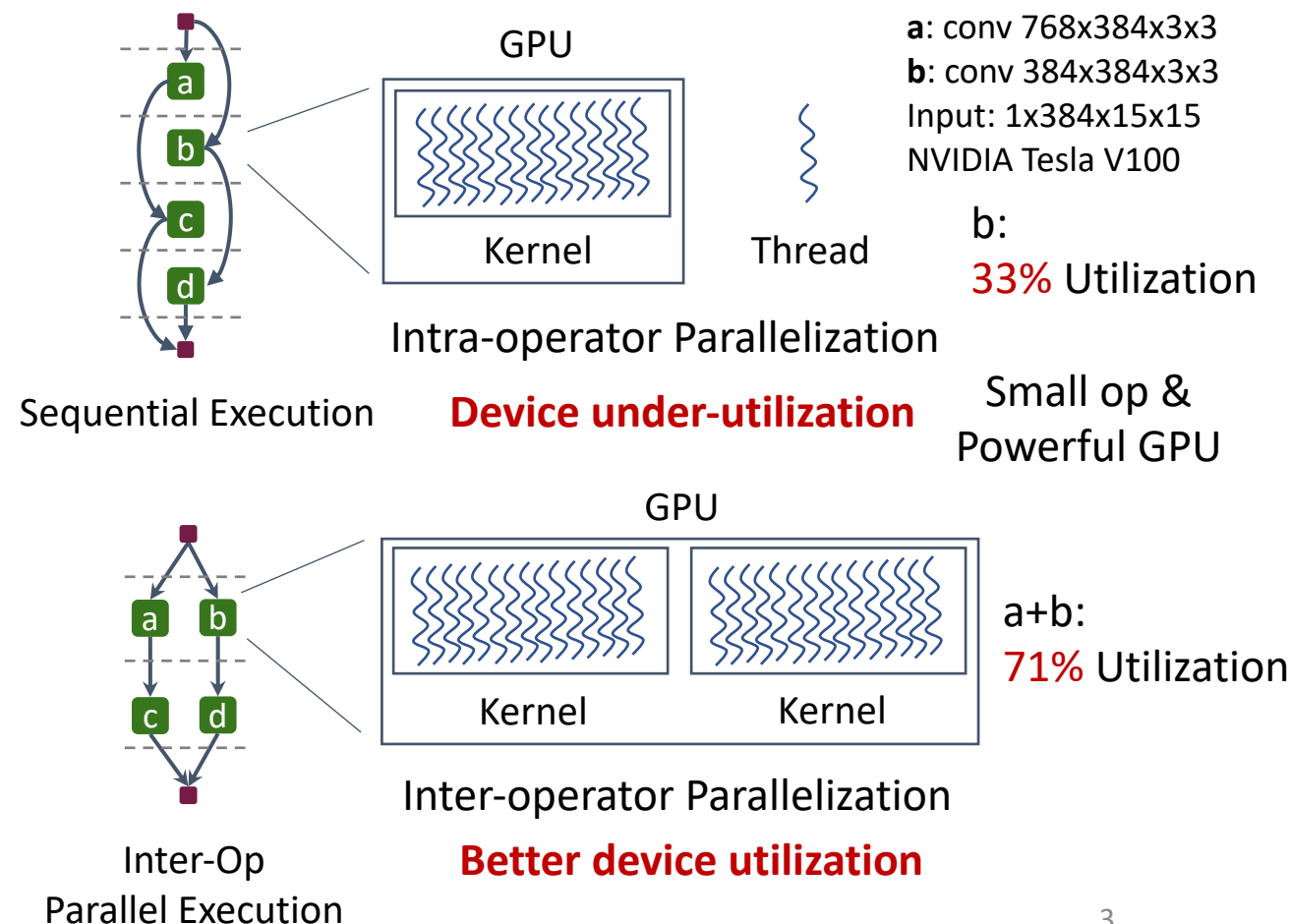
## 1. More small convs in CNN design.



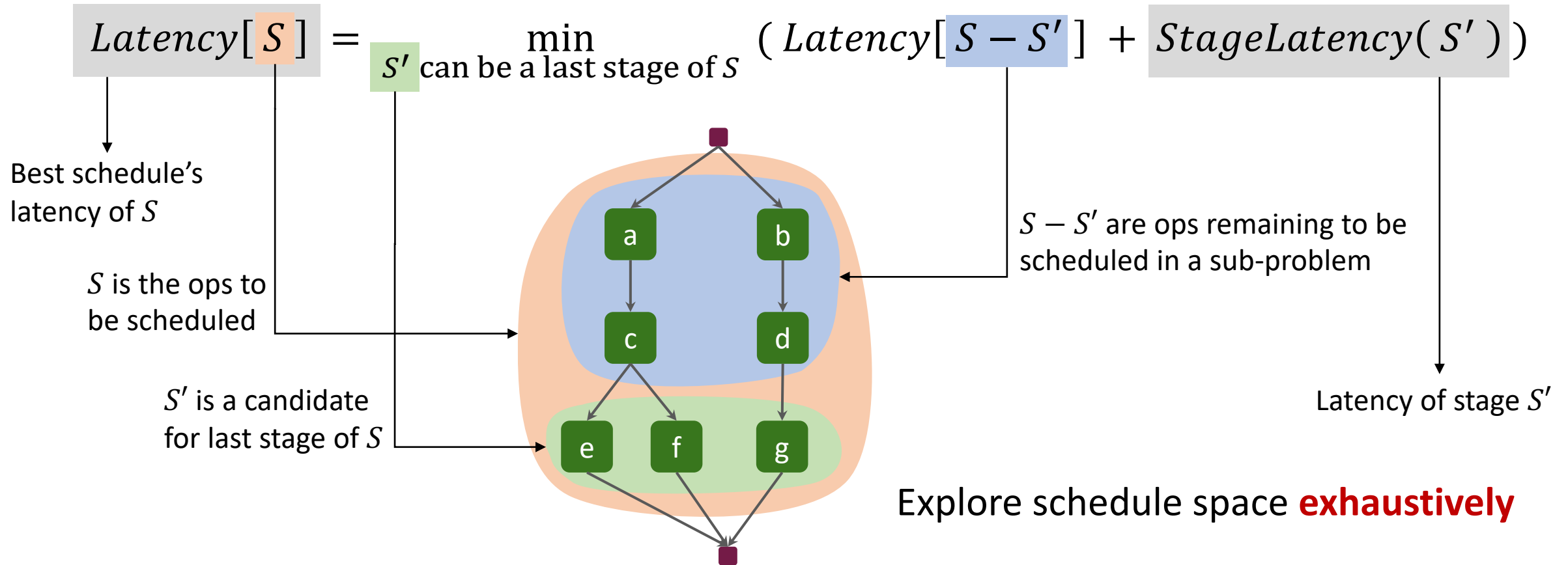
## 2. GPU peak performance increased



## 3. Intra- and Inter-operator Parallelization



# Inter-Operator Scheduler (IOS)



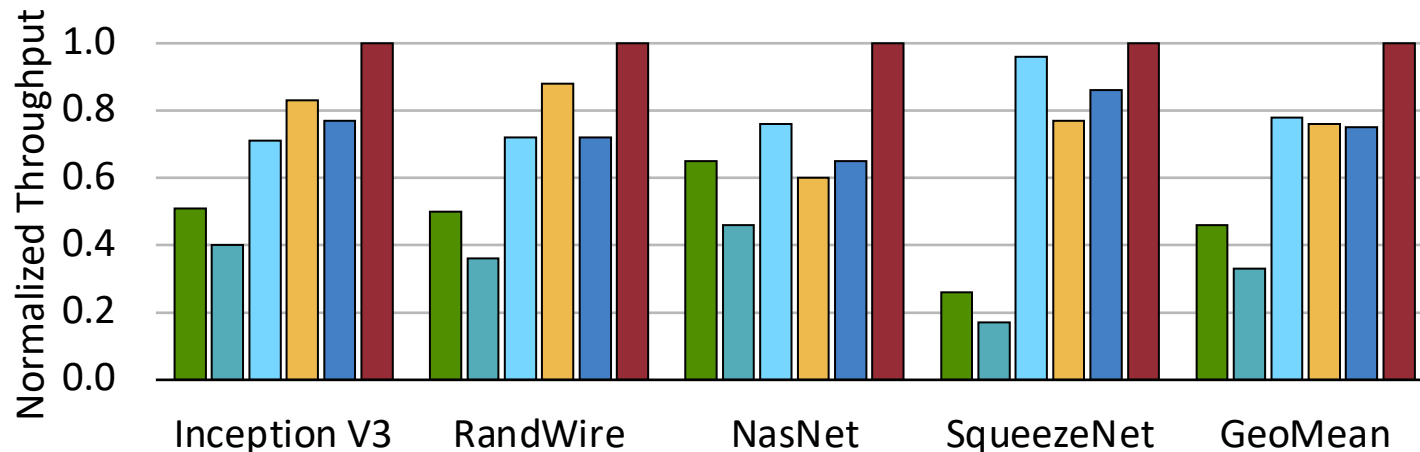
**Time Complexity** of IOS:  $\mathcal{O}((n/d + 1)^{2d})$

$n$ : number of operators  
 $d$ : width of computation graph  
 (max number of parallelizable ops)

# Comparison of cuDNN-based Frameworks

- TensorFlow:** A widely-used 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



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

Performance is normalized to the best framework

# Conclusion

- Sequential execution suffers from **under utilization** problem.
- Inter-Operator Scheduler (**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!

