

IOS: Inter-Operator Scheduler for CNN Acceleration

Yaoyao Ding¹ Ligeng Zhu² Zhihao Jia³ Gennady Pekhimenko¹ Song Han²

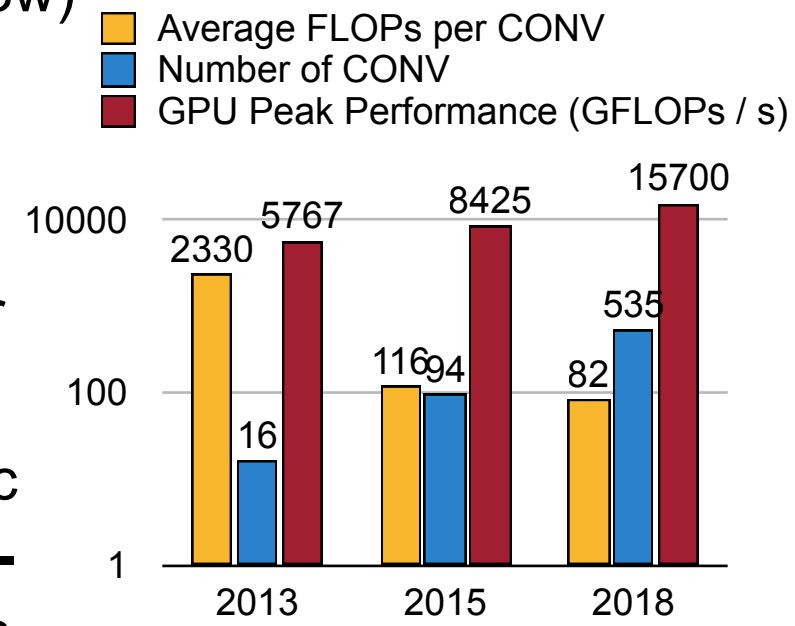
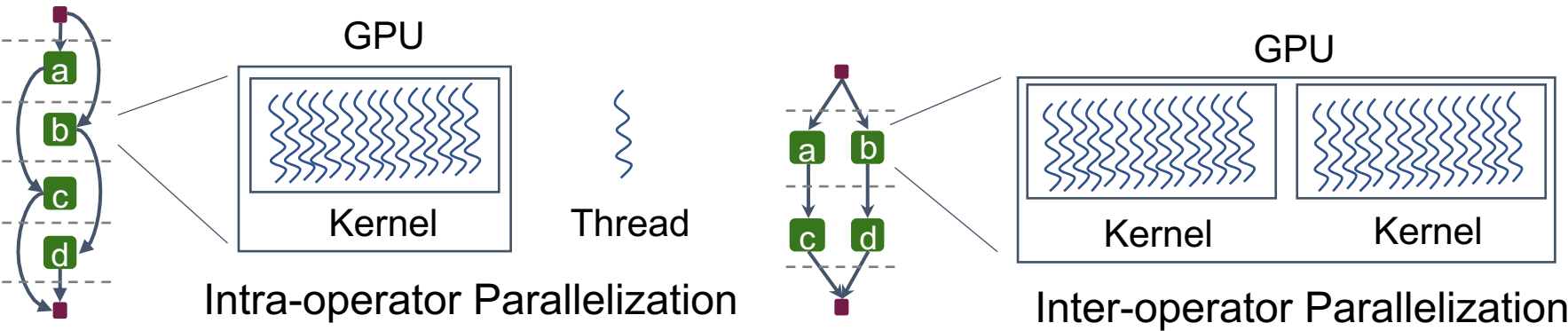
University of Toronto, Massachusetts Institute of Technology, Carnegie Mellon University

We provide scripts to reproduce results in every figure and table!

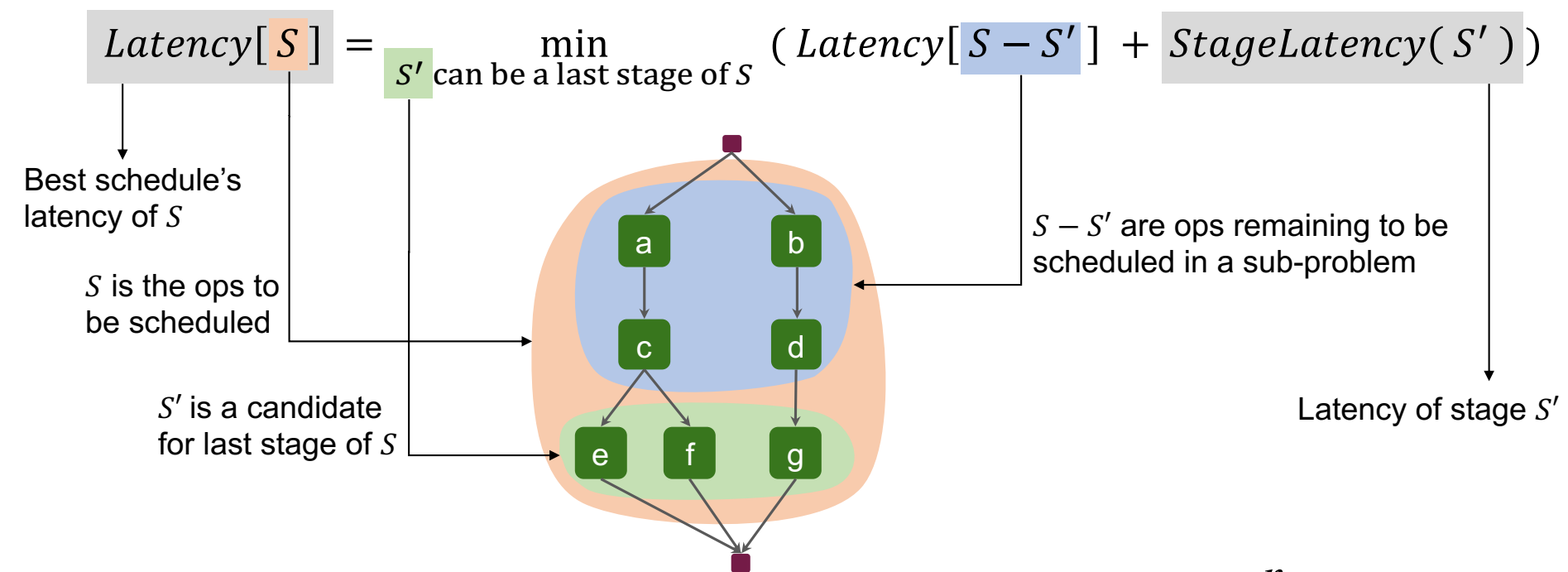


Overview

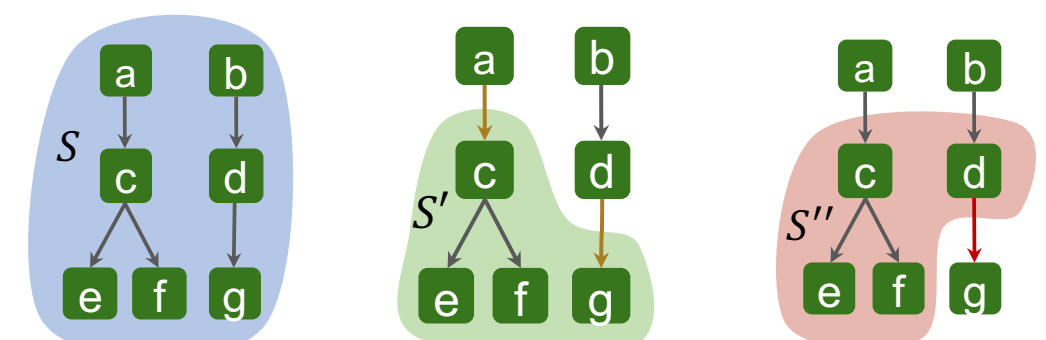
- Existing frameworks (PyTorch, TensorFlow) focus on **intra-operator parallelization**.
- Only utilizing intra-operator parallelism suffers from device under-utilization problem, especially for small op & power GPU.
- Therefore, we propose IOS — a dynamic programming algorithm scheduling **inter-operator parallelization** of CNN models.



Inter-Operator Scheduler



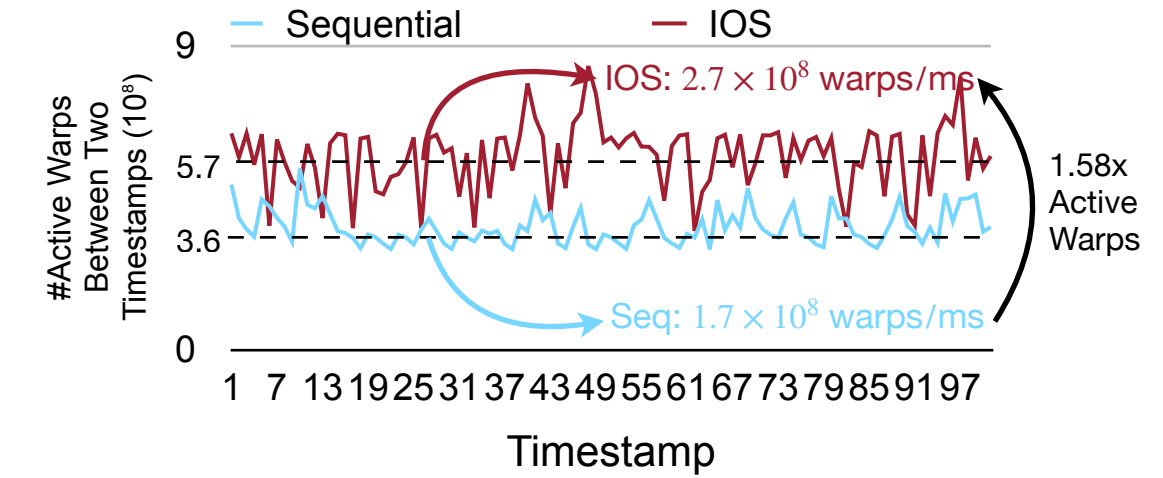
The time complexity of the dynamic programming is: $O((\frac{n}{d} + 1)^{2d})$



S : ops to be scheduled S' : Last stage candidate S'' : NOT Last stage candidate

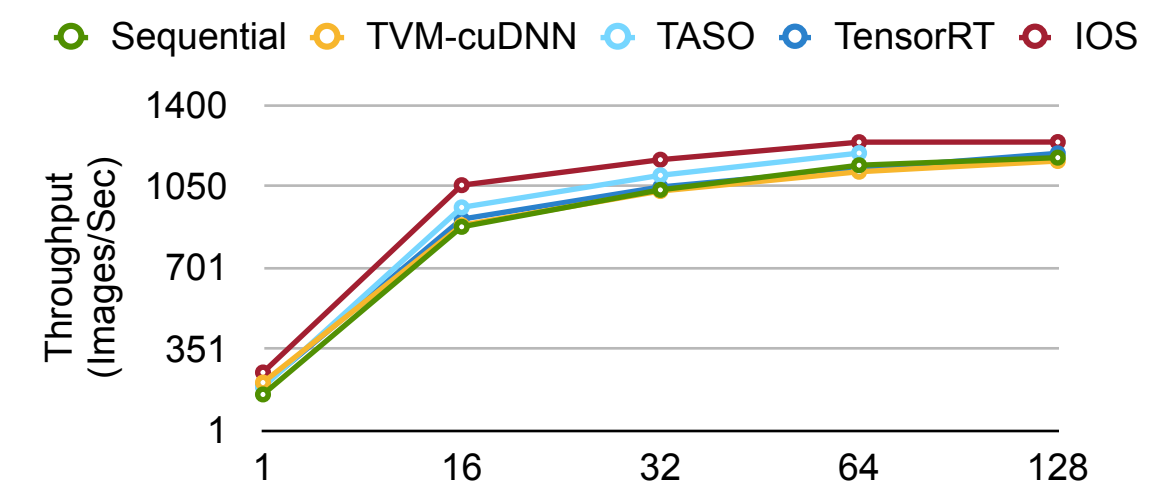
S' can be a last stage of $S \Leftrightarrow$ There is no edge from S' to $S - S'$

More Active Warps



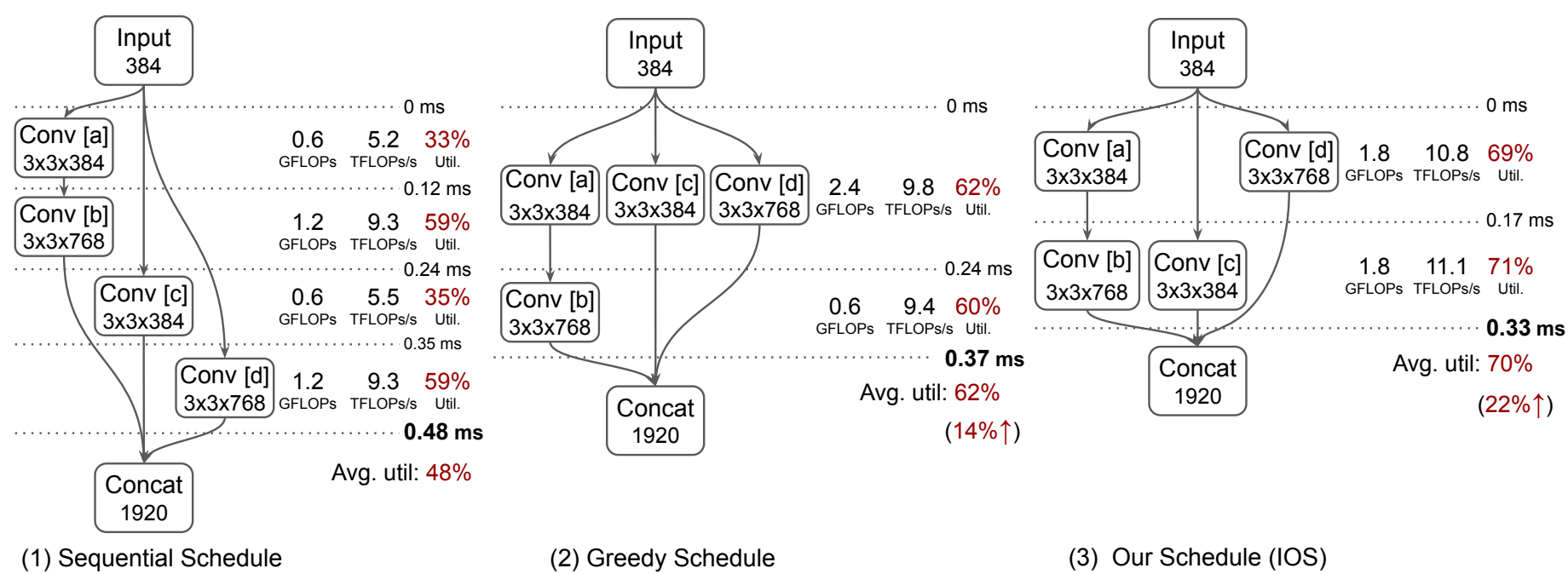
IOS Schedule has More Active Warps per ms

Large Batch Size



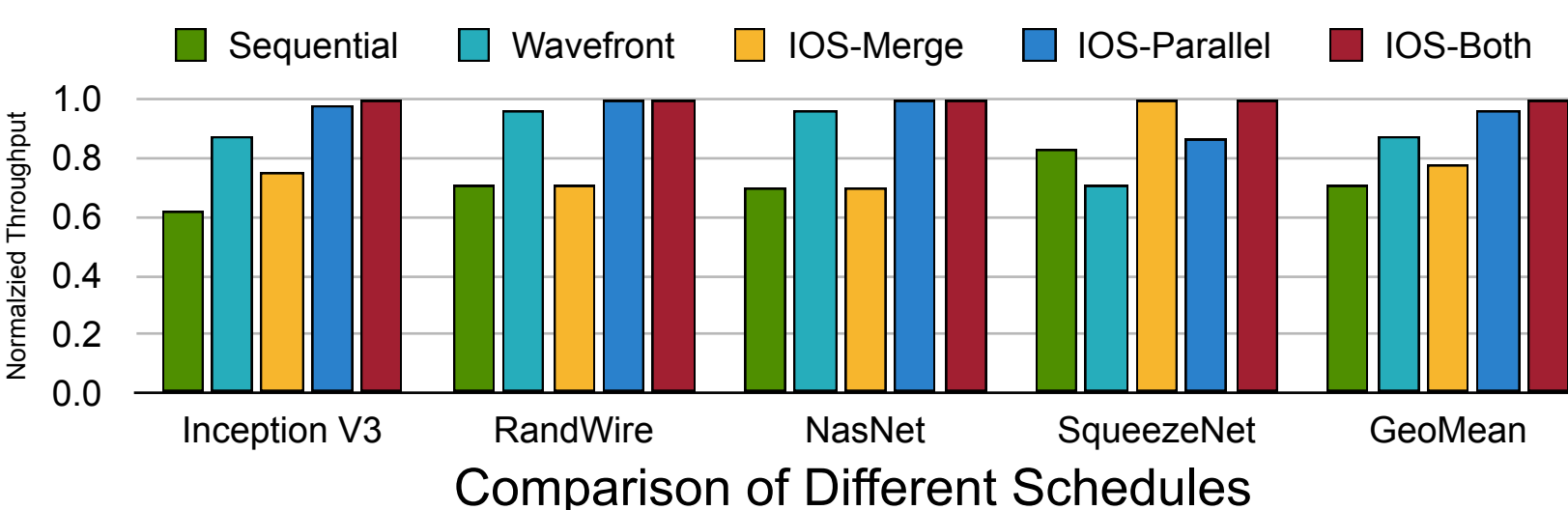
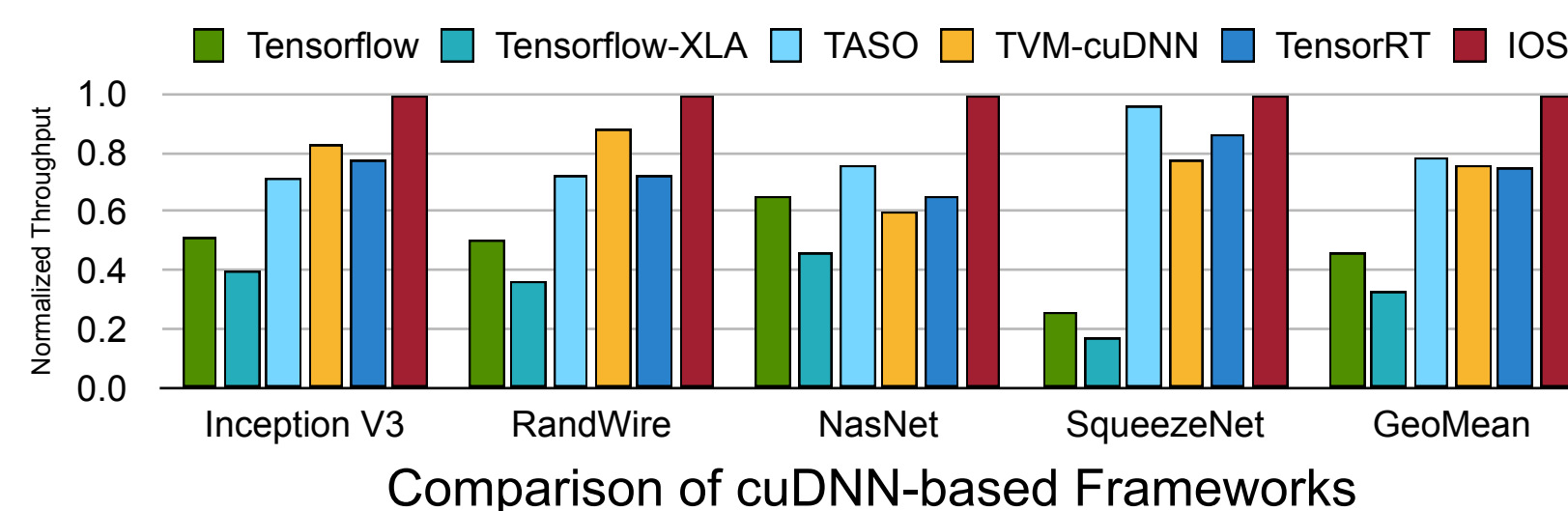
Consistent Improvement for Larger Batch Sizes (Inception V3 is used as benchmark)

Explore More Schedules is Important

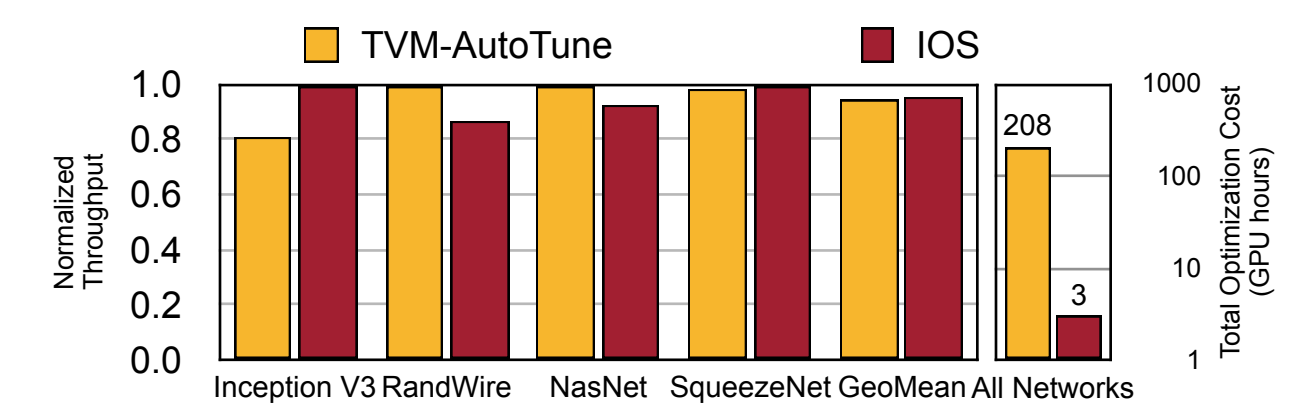


- Sequential Schedule**: the default choice for most frameworks, but leads to **insufficient utilization** as only one operator at a stage.
- Wavefront Schedule**: a greedy method that execute all available operators stage by stage. It is sub-optimal due to **unbalanced** schedule.
- IOS Schedule (ours)**: explores schedule space **exhaustively**, balances the computation in each stage, and best utilizes the hardware.

IOS Accelerates Inference



IOS v.s. AutoTVM



AutoTVM and IOS are **orthogonal** and can be combined to further boost the performance

Schedule Specialization

Specialization for Different Batch Sizes	Optimized for		
	1	32	128
Execute on 1	4.03	4.50	4.63
Execute on 32	29.21	27.44	27.93
Execute on 128	105.98	103.74	103.29

Specialization for Different Devices	Optimized for	
	K80	V100
Execute on K80	13.87	14.65
Execute on V100	4.49	4.03

Specialization for Batch Sizes Specialization for Devices
Specialized Schedules achieves the best performance