# LithOS: An Operating System for Efficient Machine Learning on GPUs

Patrick H. Coppock, Brian Zhang, Eliot H. Solomon, Vasilis Kypriotis, Leon Yang\*, Bikash Sharma\*, Dan Schatzberg\*, Todd C. Mowry, and Dimitrios Skarlatos

**SOSP 2025** 





## GPU Utilization Is Low Everywhere!



Device: 30%

SM: **15**%

Train: **40**%

(more details in our paper)



(Philly, ATC 2019)

Elababa 10%

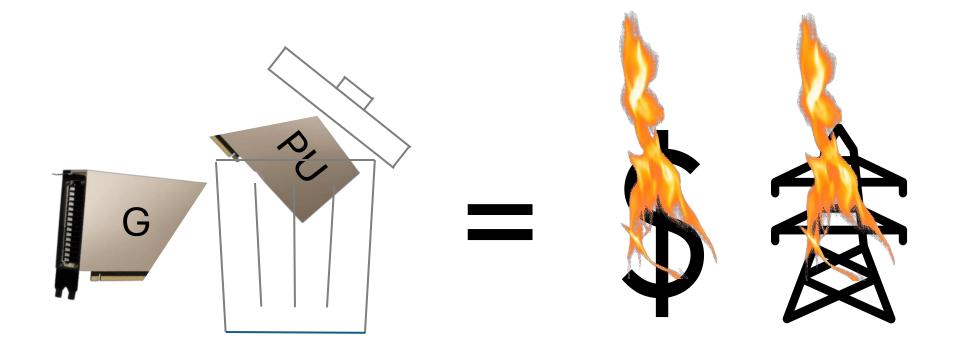
(AntMan, OSDI 2021)



(MuxFlow, arXiv)



#### Low GPU Utilization Is Costly



Throwing GPUs away wastes money and power!



#### What Is the Role of the OS in the GPU Era?

#### Mainframes

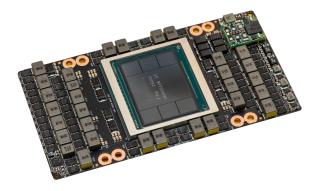


1950s



- Expensive
- Scarce
- Single-Tenant

**GPUs** 



2020s

Operating systems were invented to solve these problems!

## LithOS: An Operating System for GPUs

#### Fully transparent to ML stack

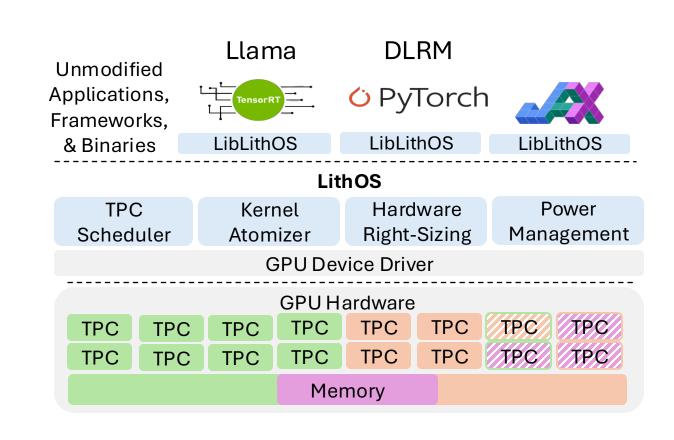
- No PTX, framework, compiler changes
- TPC Scheduling + Stealing
- Kernel Atomizer
- Right-size kernels to TPCs
- Transparent DVFS

#### Better performance

- 13x lower tail latencies vs. NVIDIA
- 3x lower tail latencies vs. SotA<sub>1</sub>
- 1.6x higher throughput vs. SotA<sub>2</sub>

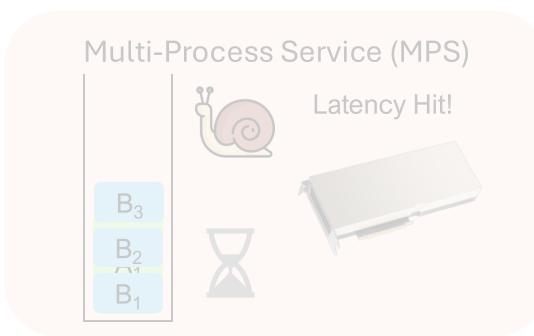
#### Easy resource saving

- Right-Sizing: 25% capacity savings
- DVFS: 25% energy savings

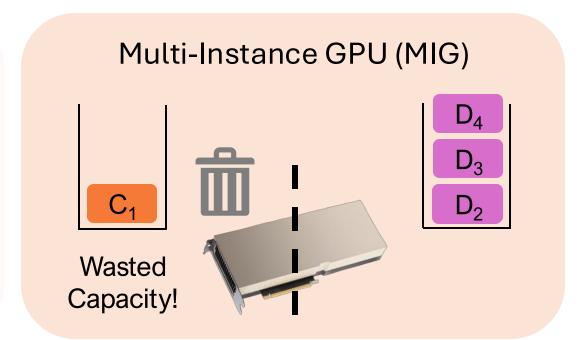


5

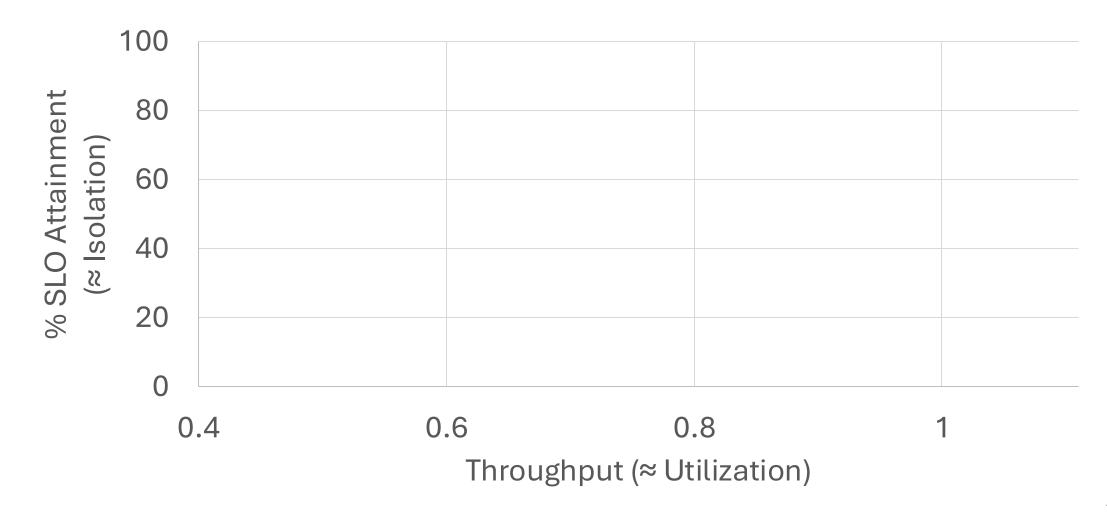
#### SotA GPU Sharing Methods Aren't Effective



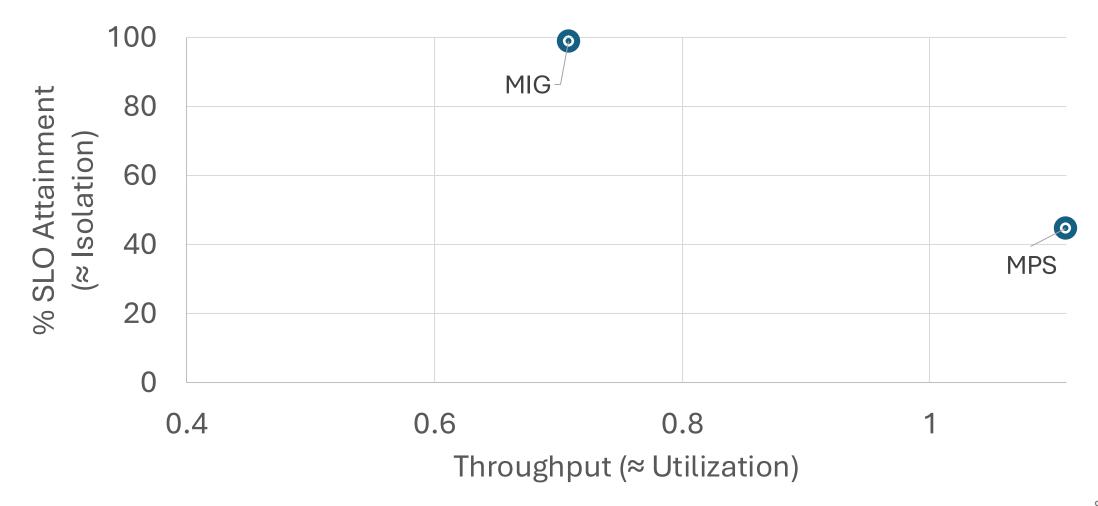
Time-sharing, no preemption → performance interference



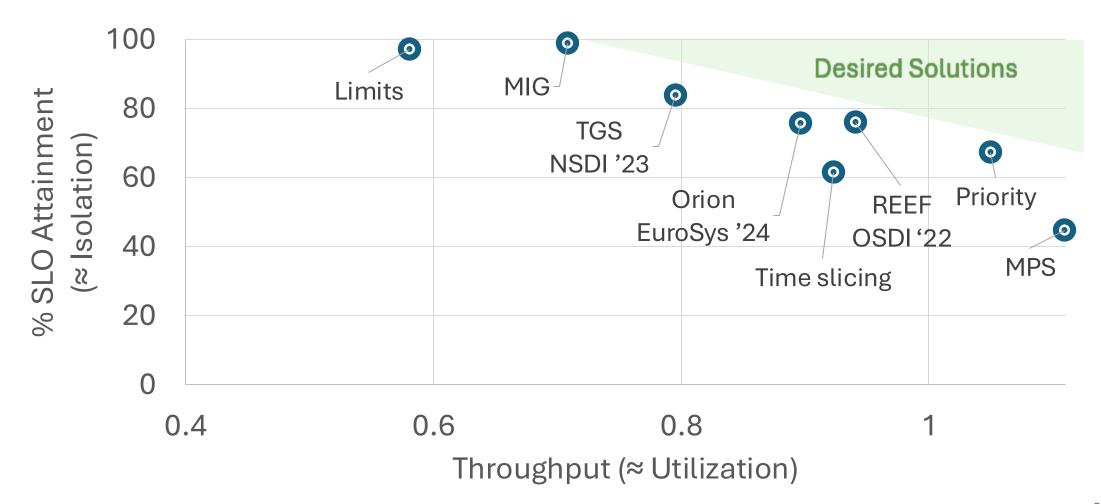
Coarse-grained, fixed partitions  $\rightarrow$  poor utilization











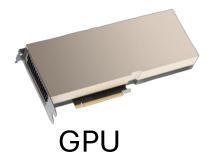


# CPUs Can't Be Shared Efficiently



Sharing Happens Transparently Threads
Scheduled to
Specific Cores

Preemption



Invasive Modifications Required No Control of Where Kernels Execute

No Preemption



Interposition +
Reverse
Engineering

TPC Scheduling

Kernel Atomization

## LithOS Interposition + Reverse Engineering

ML Application

ML Framework

CUDA Runtime

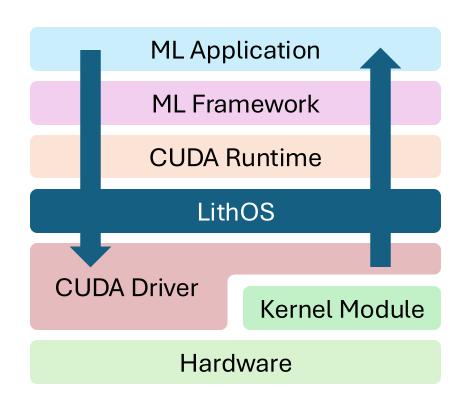
CUDA Driver

Kernel Module

Hardware



# LithOS Interposition + Reverse Engineering



TPC Masking

Prelude Kernels

#### Advantages

- No need to modify application software (CUDA C++, PTX, etc.)
- Pulls kernel scheduling logic out of closed-source software and hardware
- Easy to port across hardware and driver versions

#### Vision

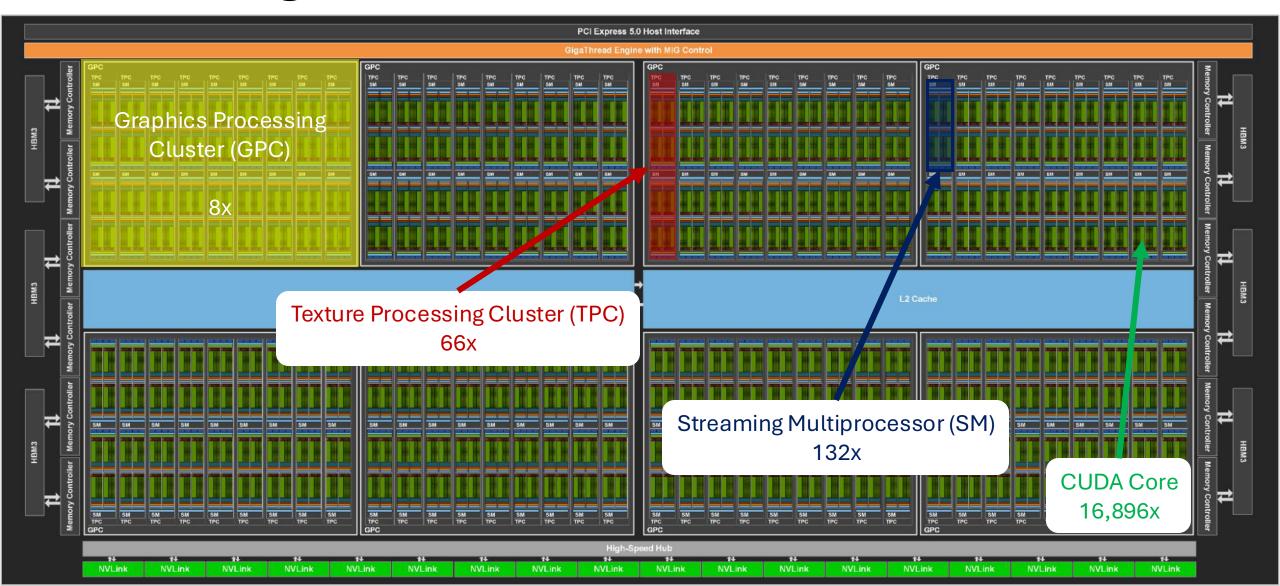
 Software from LithOS down is part of an integrated, open OS



MIG partitions on GPC boundaries

## **GPU Organization 101**

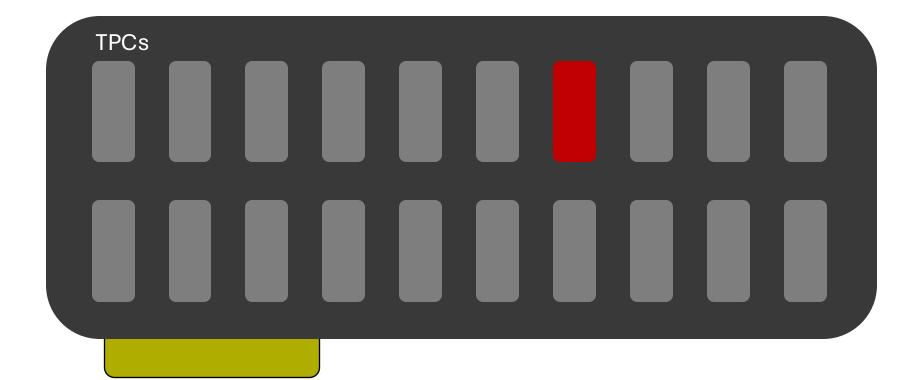
LithOS partitions on TPC boundaries









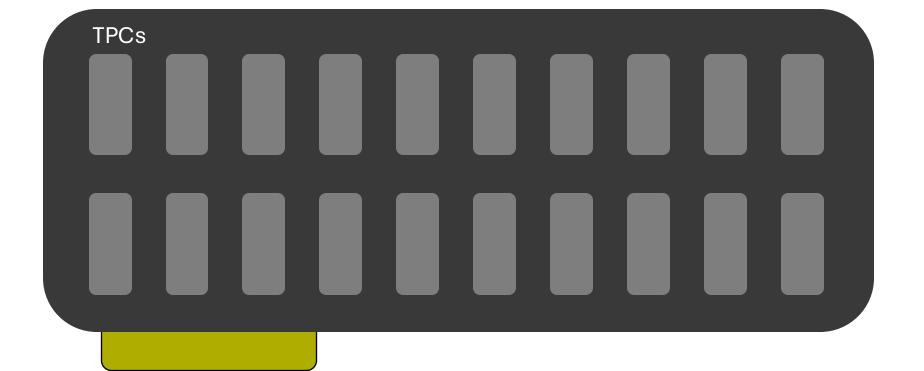


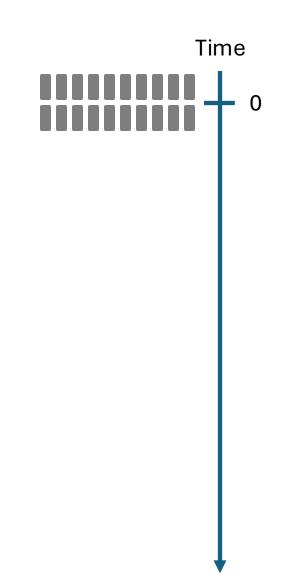


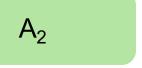
A<sub>1</sub>

B<sub>1</sub>

 $C_1$ 

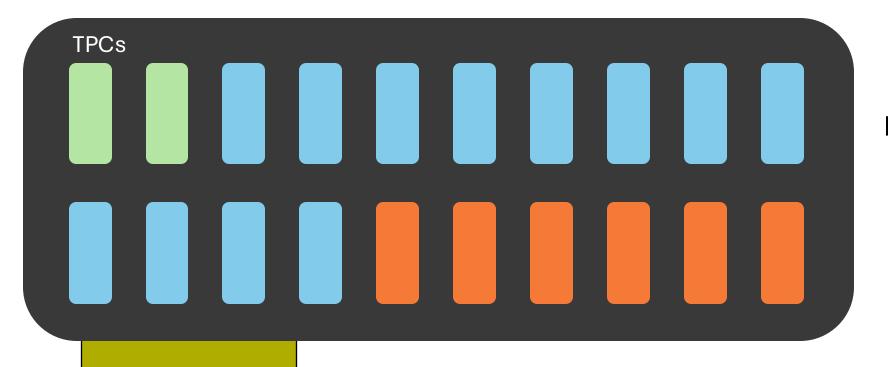


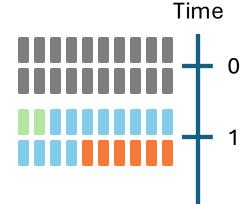








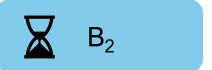




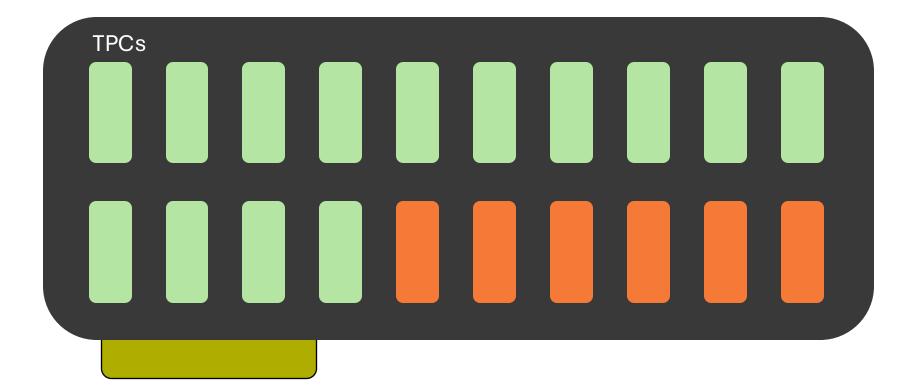
TPC Resources Idle!

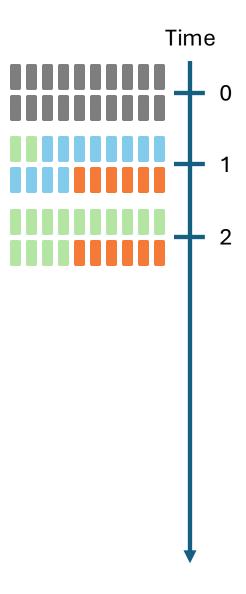
TPC Stealing!



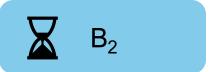




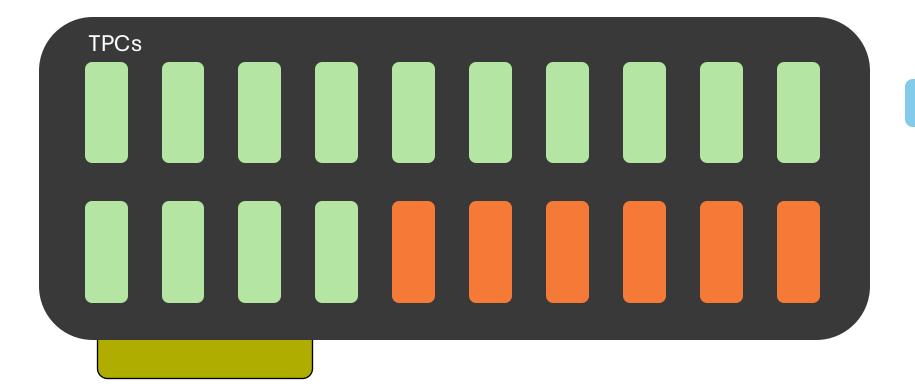


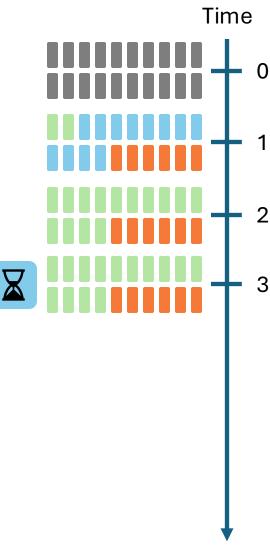




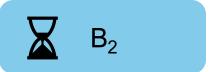




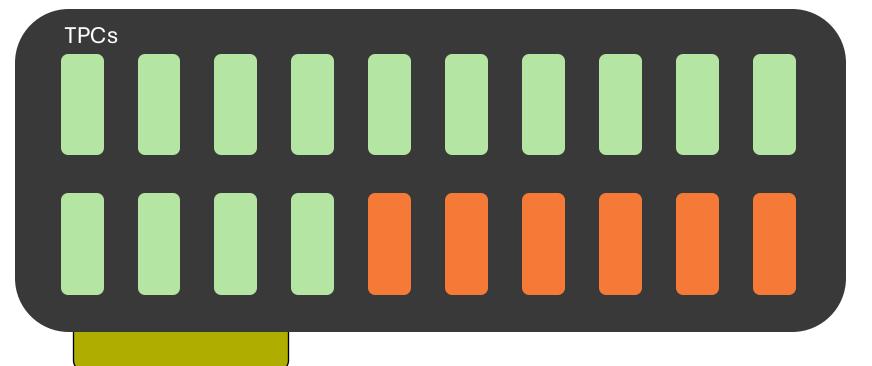


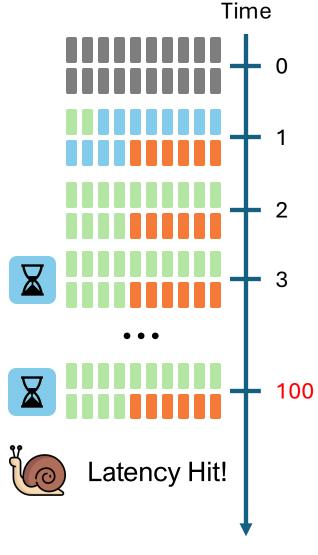






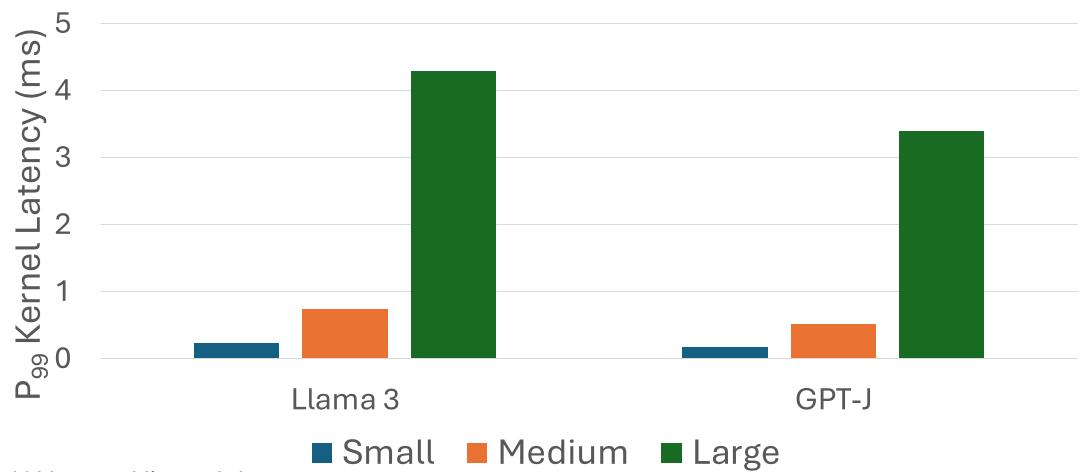


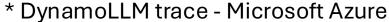




#### Kernel Runtimes Can Be Long

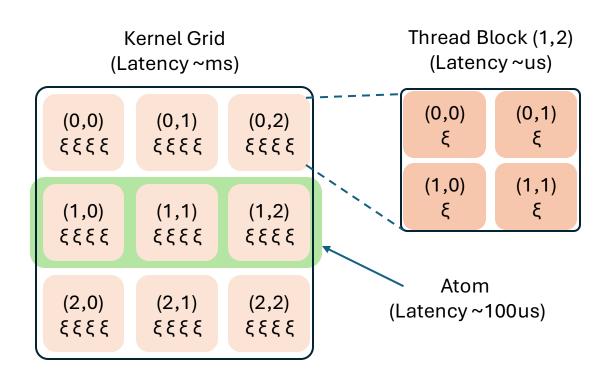
More data in the paper!

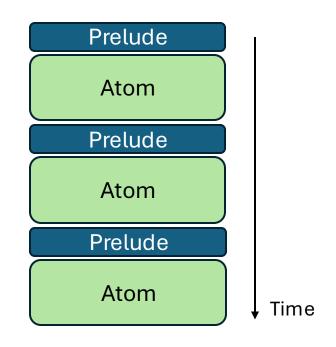






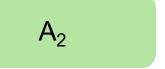
#### Transparent Kernel Atomization For Preemption





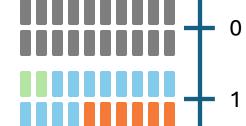
Thread-block-level scheduling in software!



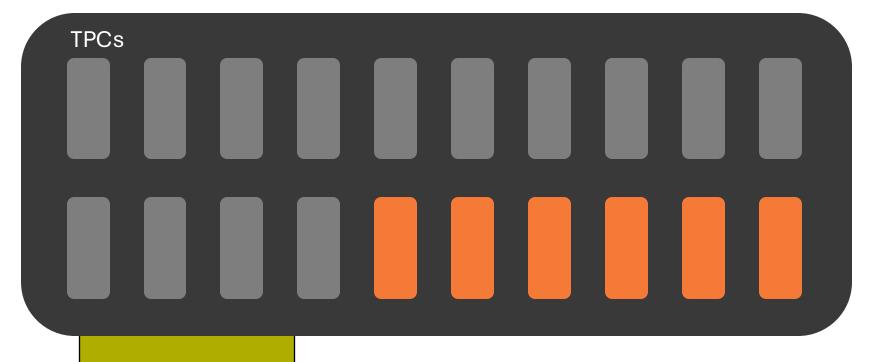






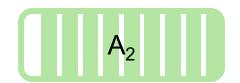


Time



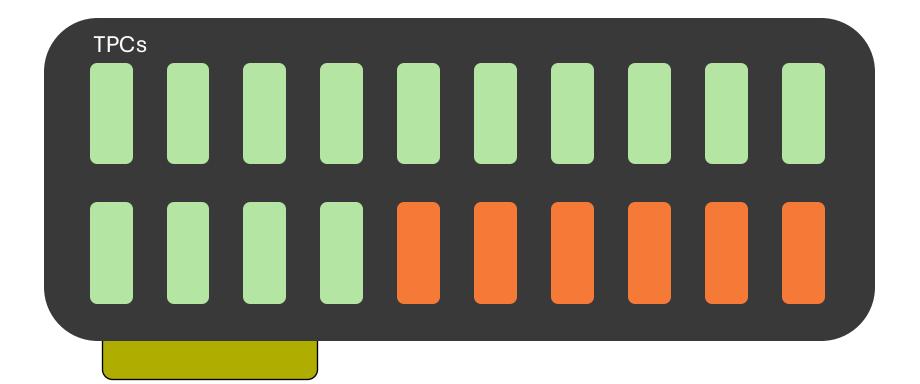
TPC Resources Idle!

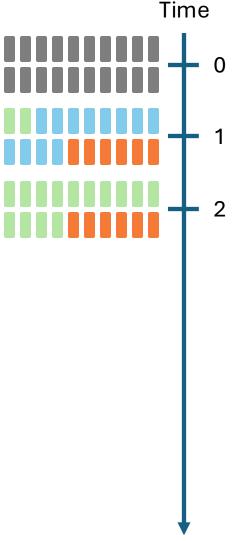
TPC Stealing!



B<sub>2</sub>



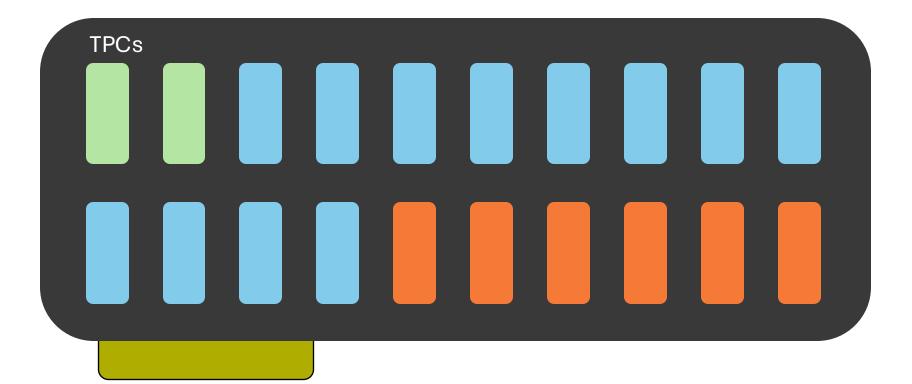


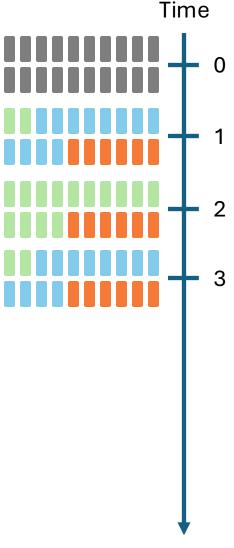


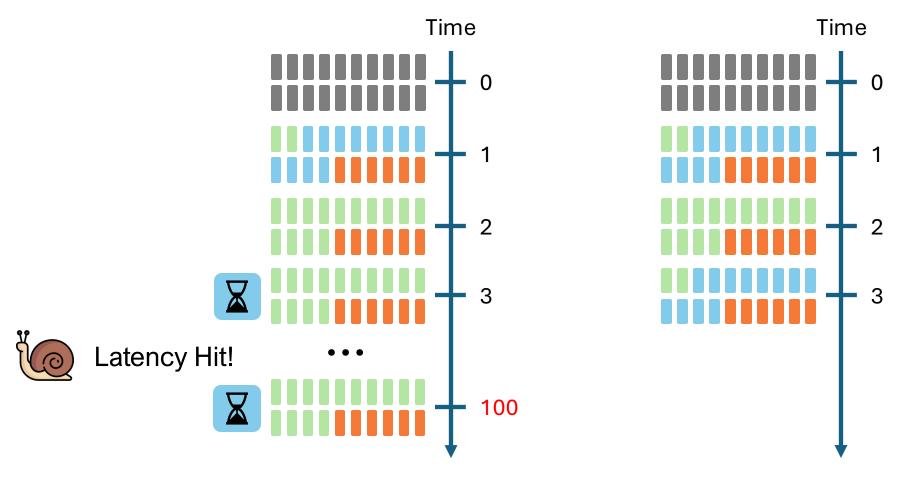








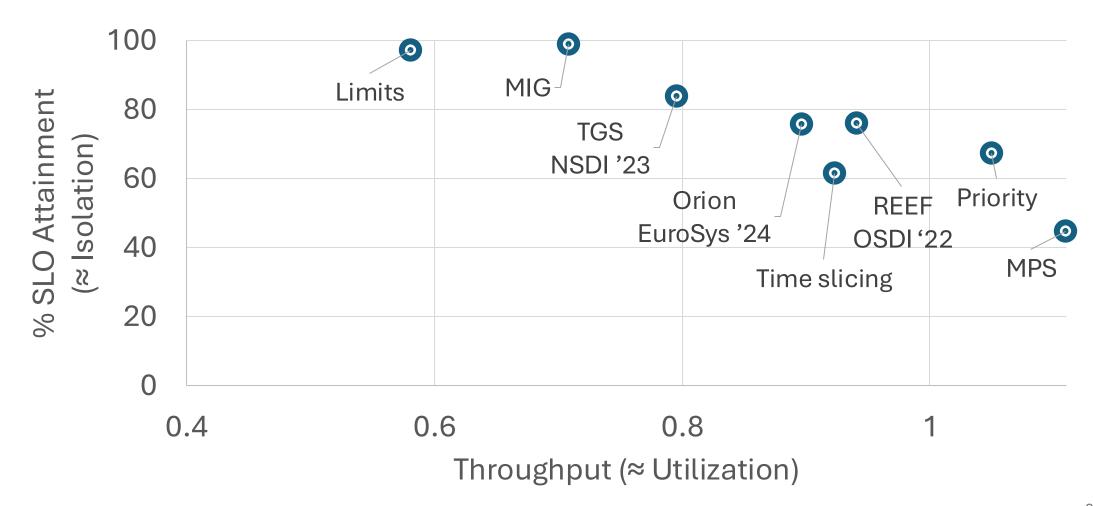




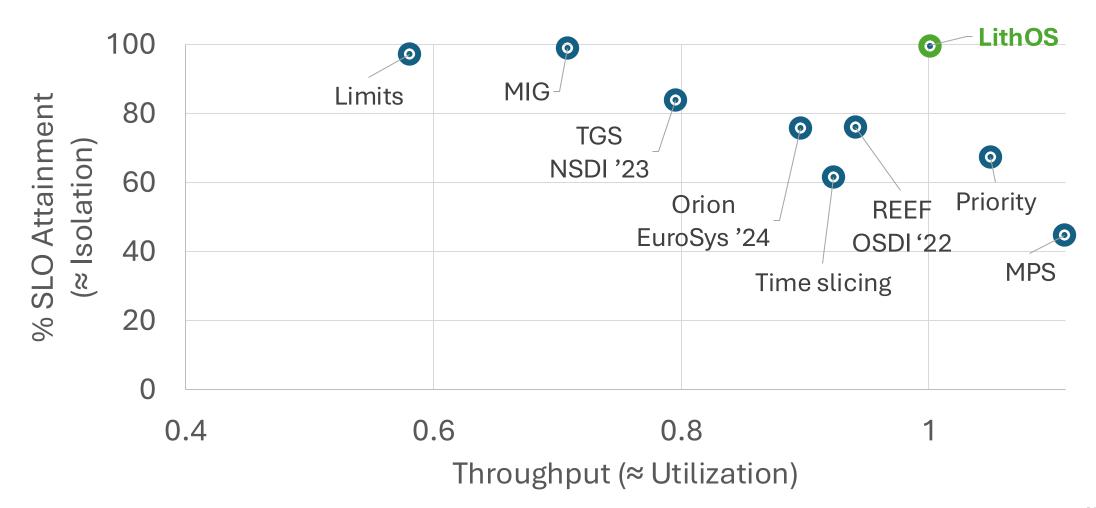










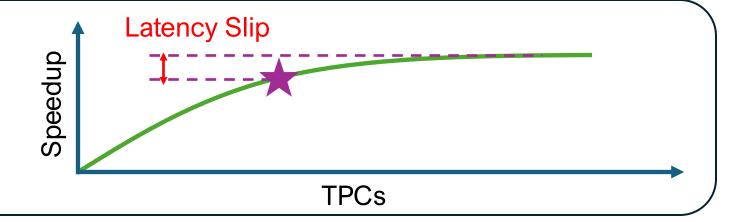




## TPC Right-Sizing and Fine-Grained DVFS

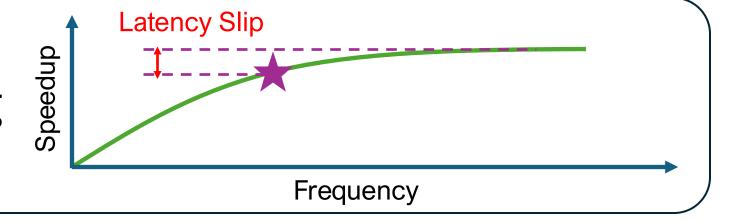
Right-Sizing

25% capacity savings



Dynamic Voltage Frequency Scaling

25% energy savings





#### Check Out the Paper for More Details!

#### Discussion

Abstractions for GPU OSes

#### Mechanisms

- Kernel/TPC Right-Sizing
- DVFS
- Online latency predictor
- Reverse engineering techniques

#### Results

- Meta GPU fleet profiling
- Kernel scaling vs. compute/freq.
- Inference-Inference stacking
- Inference-Training stacking
- Ablation studies

# LithOS: An Operating System for Efficient Machine Learning on GPUs



## LithOS: An Operating System for GPUs

#### Fully transparent to ML stack

- No PTX, framework, compiler changes
- TPC Scheduling + Stealing
- Kernel Atomizer
- Right-size kernels to TPCs
- Transparent DVFS

#### Better performance

- 13x lower tail latencies vs. NVIDIA
- 3x lower tail latencies vs. SotA<sub>1</sub>
- 1.6x higher throughput vs. SotA<sub>2</sub>

#### Easy resource saving

- Right-Sizing: 25% capacity savings
- DVFS: 25% energy savings

