

Xsched: Preemptive Scheduling for Diverse XPUs

Weihang Shen Mingcong Han Jialong Liu Rong Chen Haibo Chen

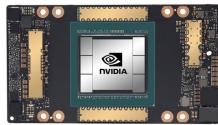
IPADS, Shanghai Jiao Tong University



XPUs are Widely Adopted



GPUs



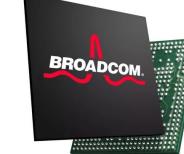
NPUs



TPUs



ASICs

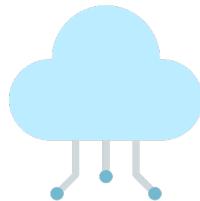


FPGAs



...

Typical scenarios:



Cloud Services

Autonomous Vehicles



Copilot + PC



Edge Devices

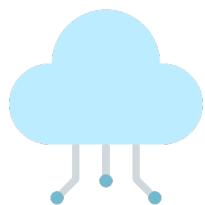
XPU Multitasking is Common

Cloud: improve utilization, save cost



Edge: limited XPU resources

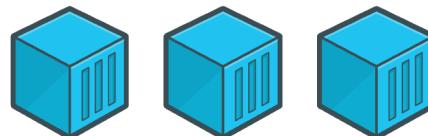
One XPU can serve *tens* of job instances



Serverless Functions

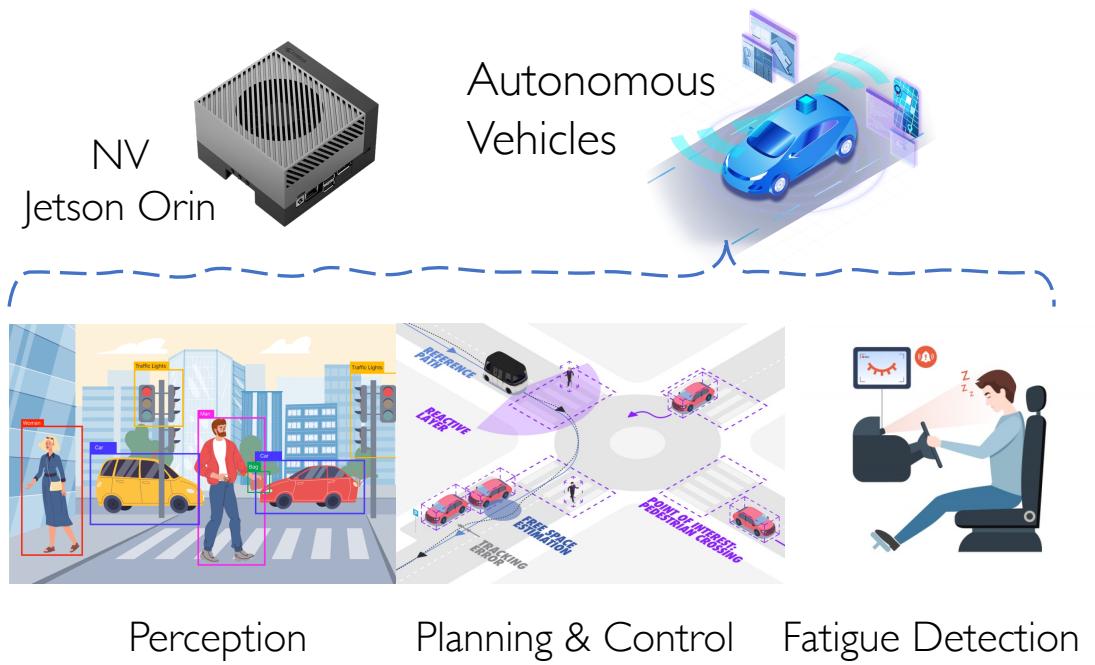


GPU/NPU Cloud

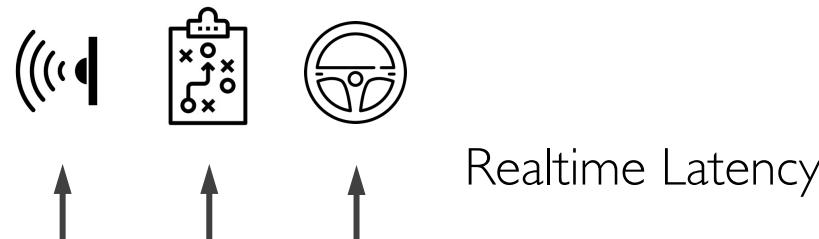


Tenants / Container Jobs

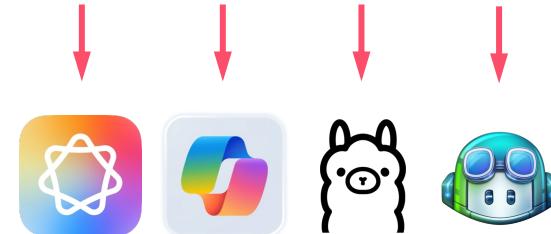
30+ tasks are deployed on 1~2 Orin GPU



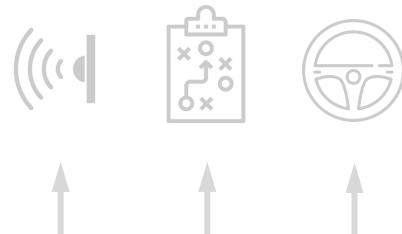
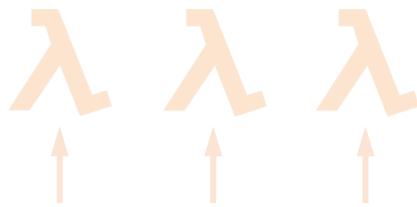
Diverse Scheduling Requirements of Apps



Responsiveness



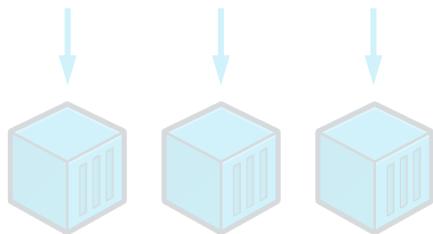
Diverse Scheduling Requirements of Apps



Realtime Latency

Task preemption is a crucial mechanism
to meet *diverse* application needs

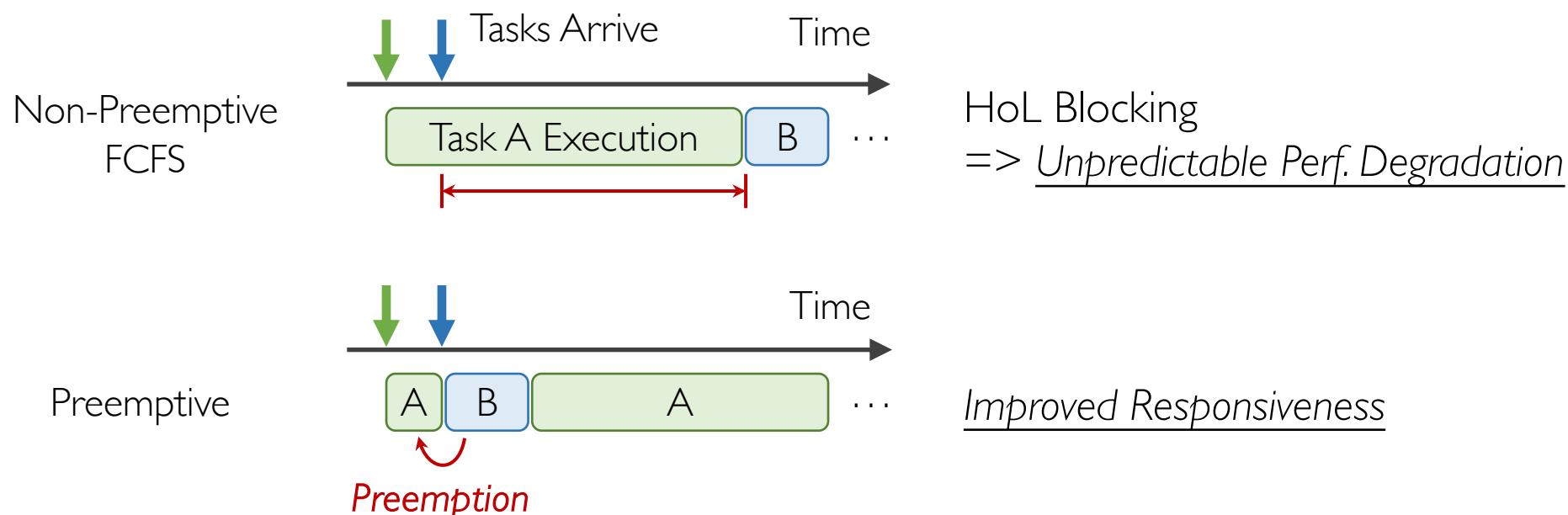
Fairness / Throughput



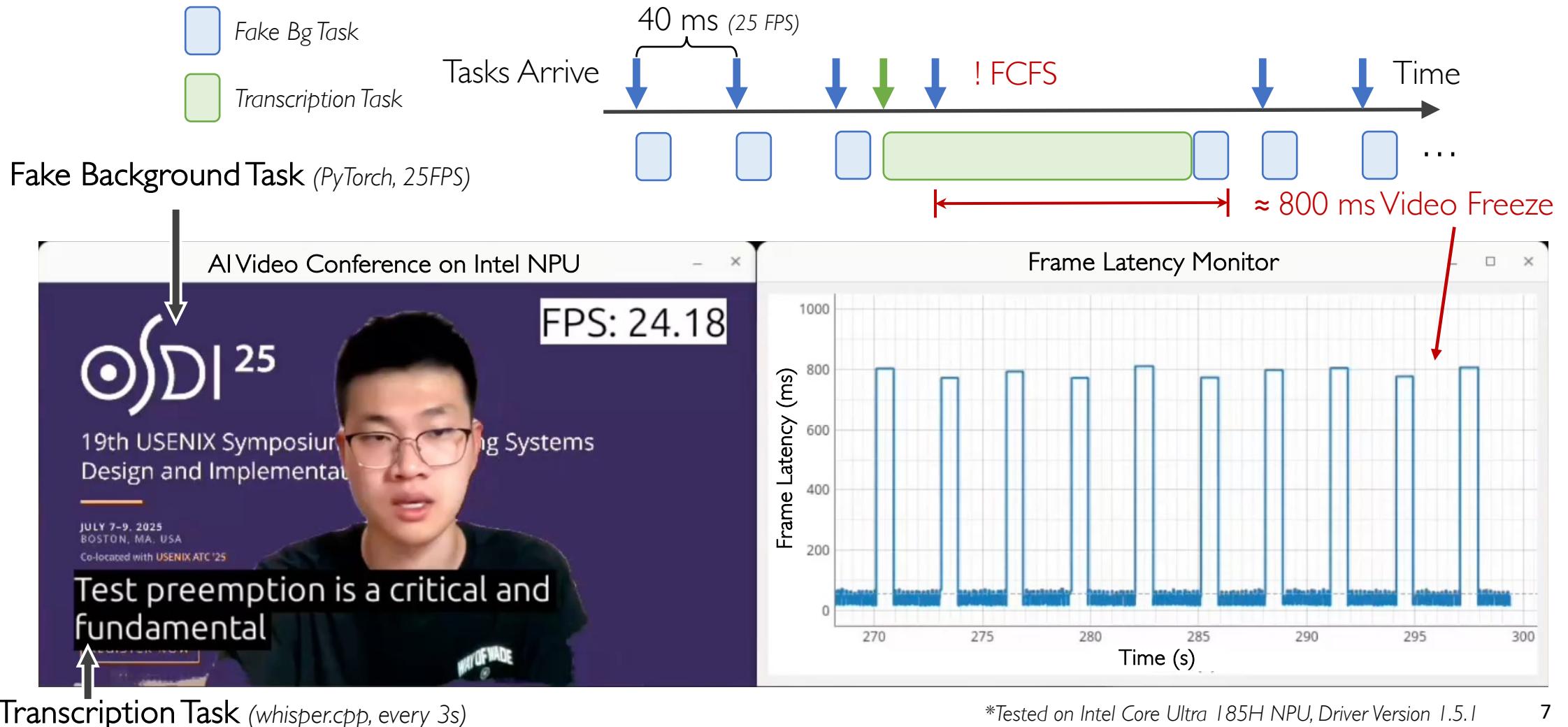
Issues of Native Hardware Scheduling

#1 Most XPU^s *lack preemption* support

Non-preemptive FCFS policy is naturally embedded in the HW, FW, and drivers.



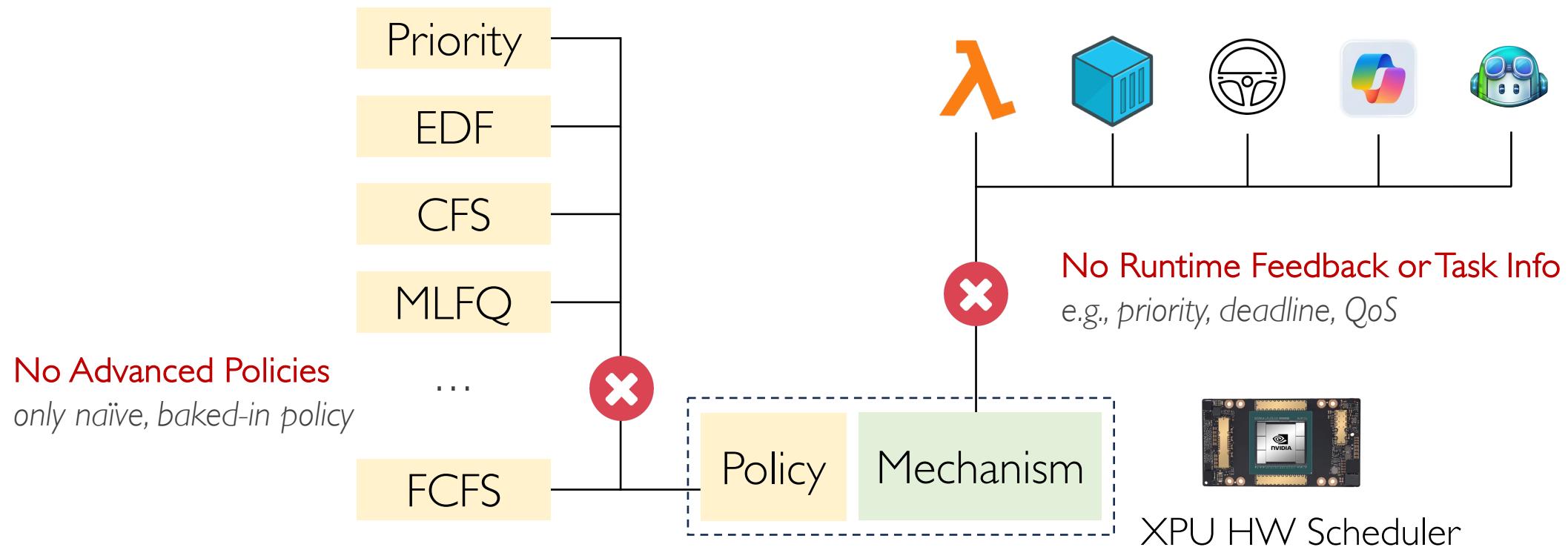
Example: AI Conference on Intel NPU



*Tested on Intel Core Ultra 185H NPU, Driver Version 1.5.1

Issues of Native Hardware Scheduling

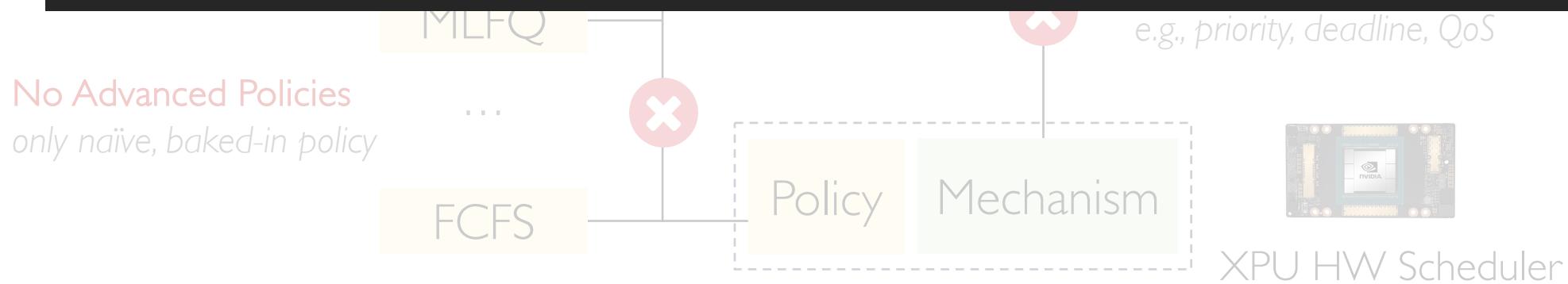
#2 Hardware schedulers are *not flexible*



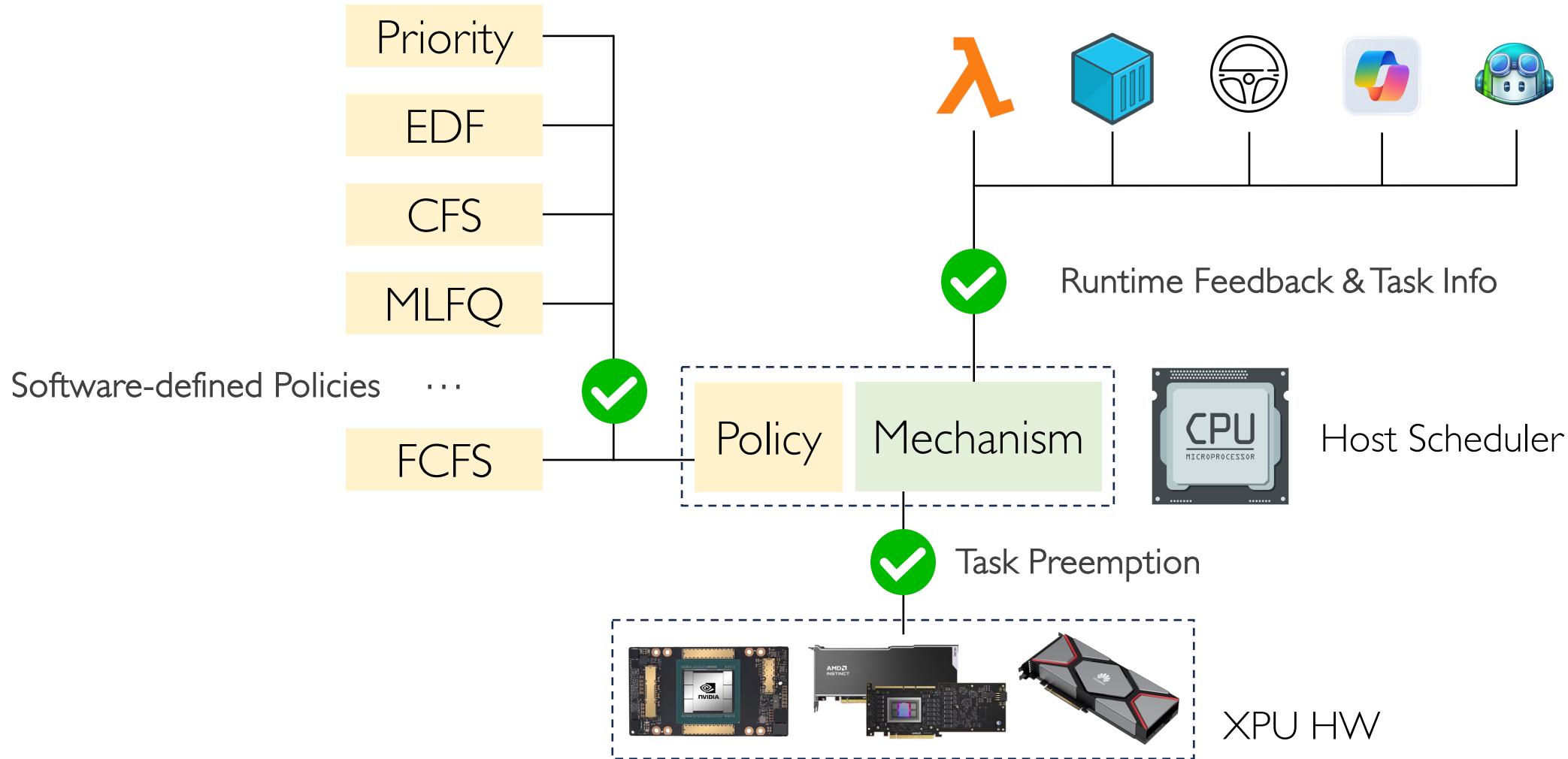
Issues of Native Hardware Scheduling

#2 Hardware schedulers are *not flexible*

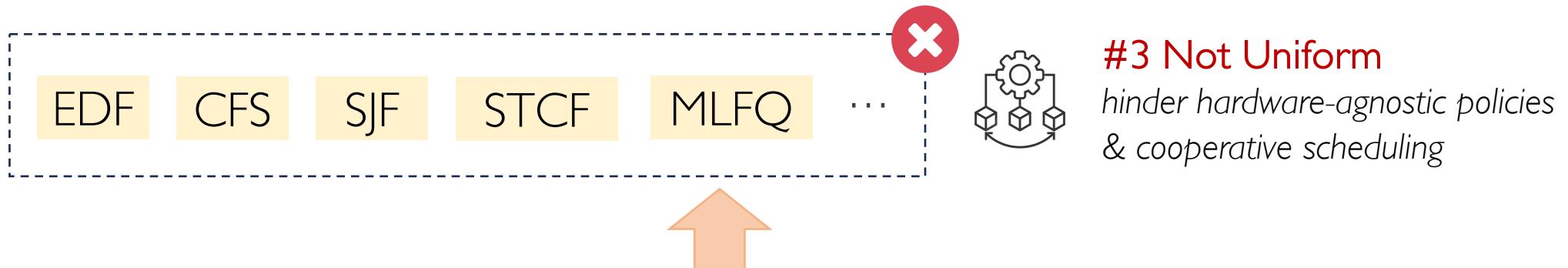
XPU hardware scheduling
mismatches application requirements



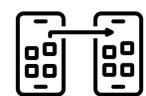
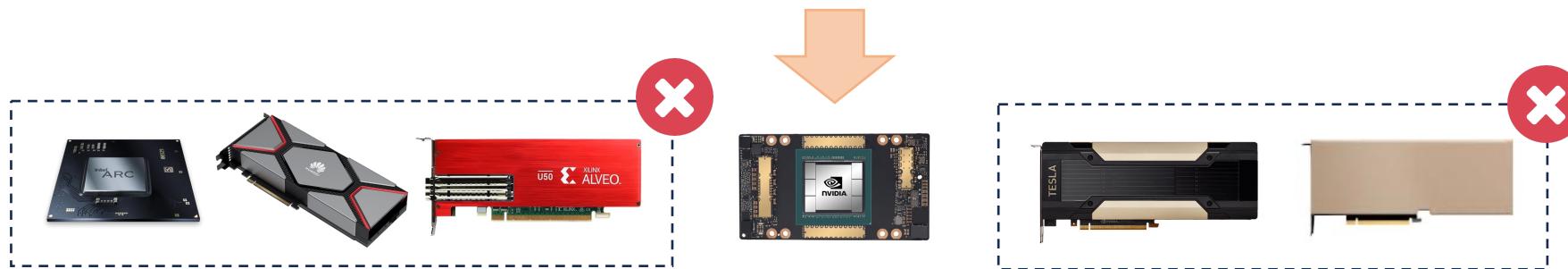
Approach: Host-Managed Scheduling



Existing Solutions are not GENERAL



Existing Solutions, e.g., *Effisha* [PPoPP'17], *FLEP* [ASPLOS'17], *REEF* [OSDI'22], ...



#1 Not Portable

designed for certain GPU,
while new XPU emerging



#2 Not Evolvable

for new / deprecated HW features,
while XPU keep advancing

Existing Solutions are not GENERAL



Our Goal: *GENERAL* & *PREEEMPTIVE*
Host-managed XPU Scheduling



#1 Not Portable

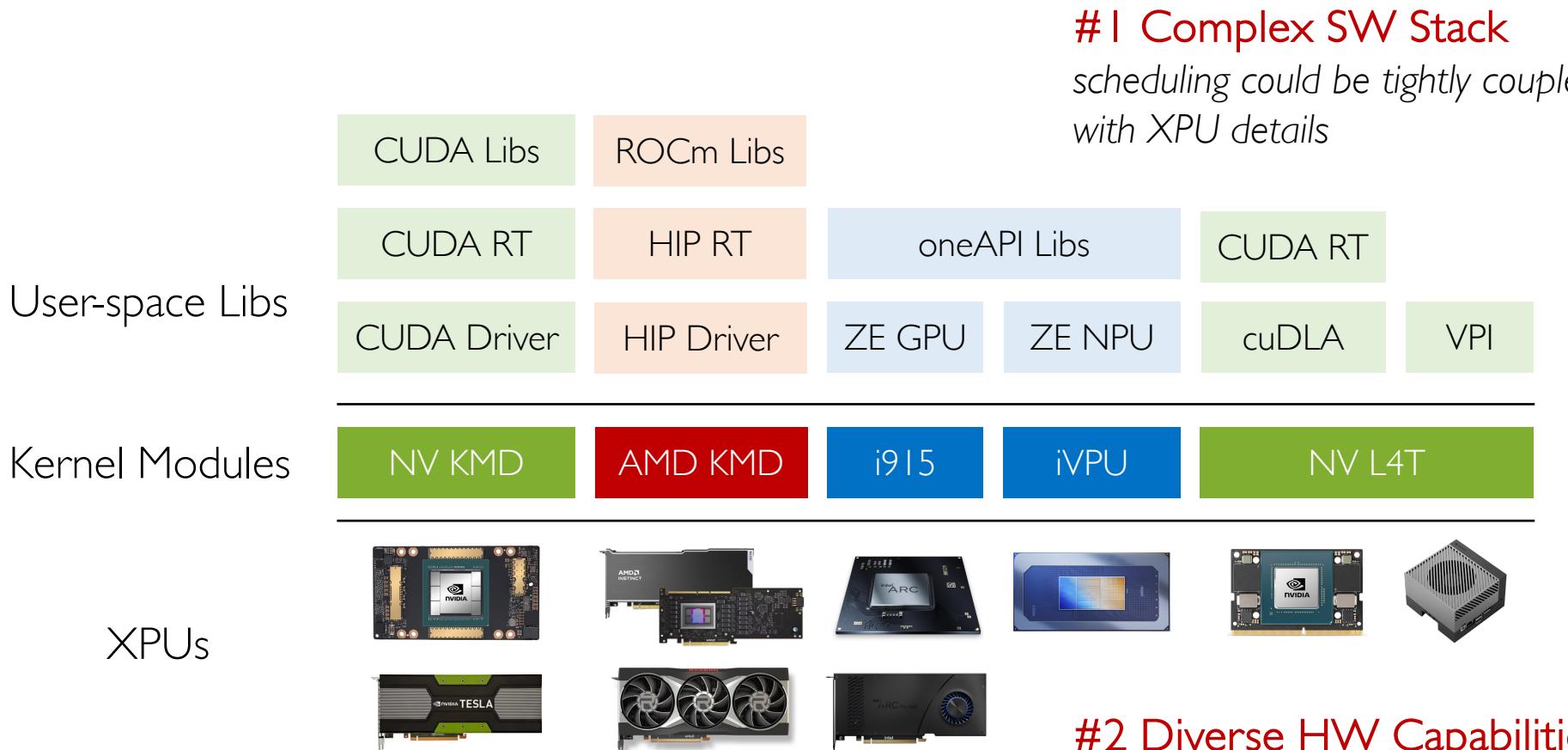
designed for certain GPU,
while new XPUs emerging



#2 Not Evolvable

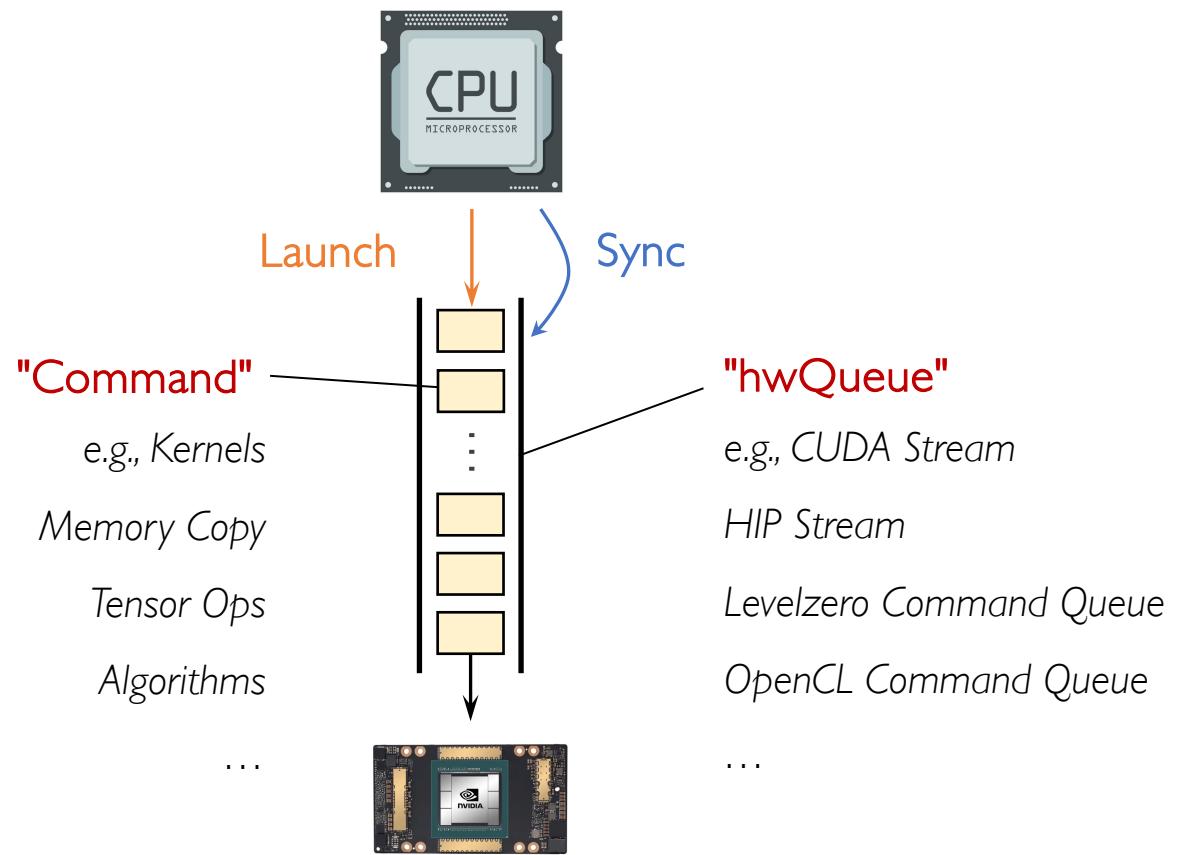
for new / deprecated HW features,
while XPUs keep advancing

Challenges for General XPU Scheduling

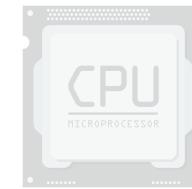


Opportunity: XPU Programming Paradigm

```
def infer_task():
    memcpyH2D(...) # memcpy cmd
    conv(...)      # kernel cmds
    relu(...)
    ...
    softmax(...)
    memcpyD2H(...) # memcpy cmd
    sync(...)      # wait for completion
```



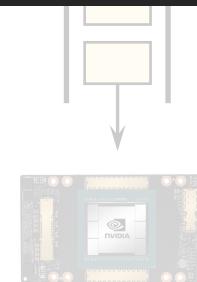
Opportunity: XPU Programming Paradigm



An XPU task can be abstracted as a sequence of *commands* executed on a *command queue*

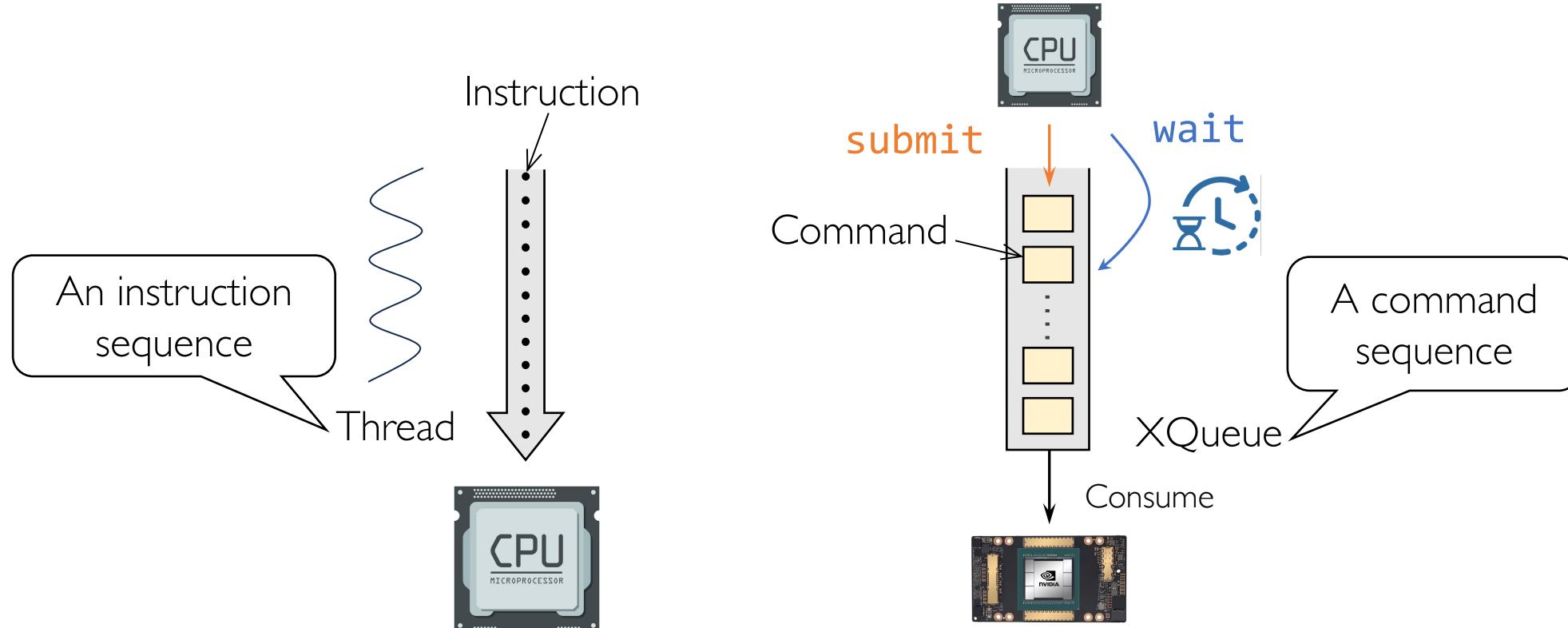
```
memcpyD2H(...) # memcpy cmd  
sync(...) # wait for completion
```

Tensor Ops
Algorithms
...

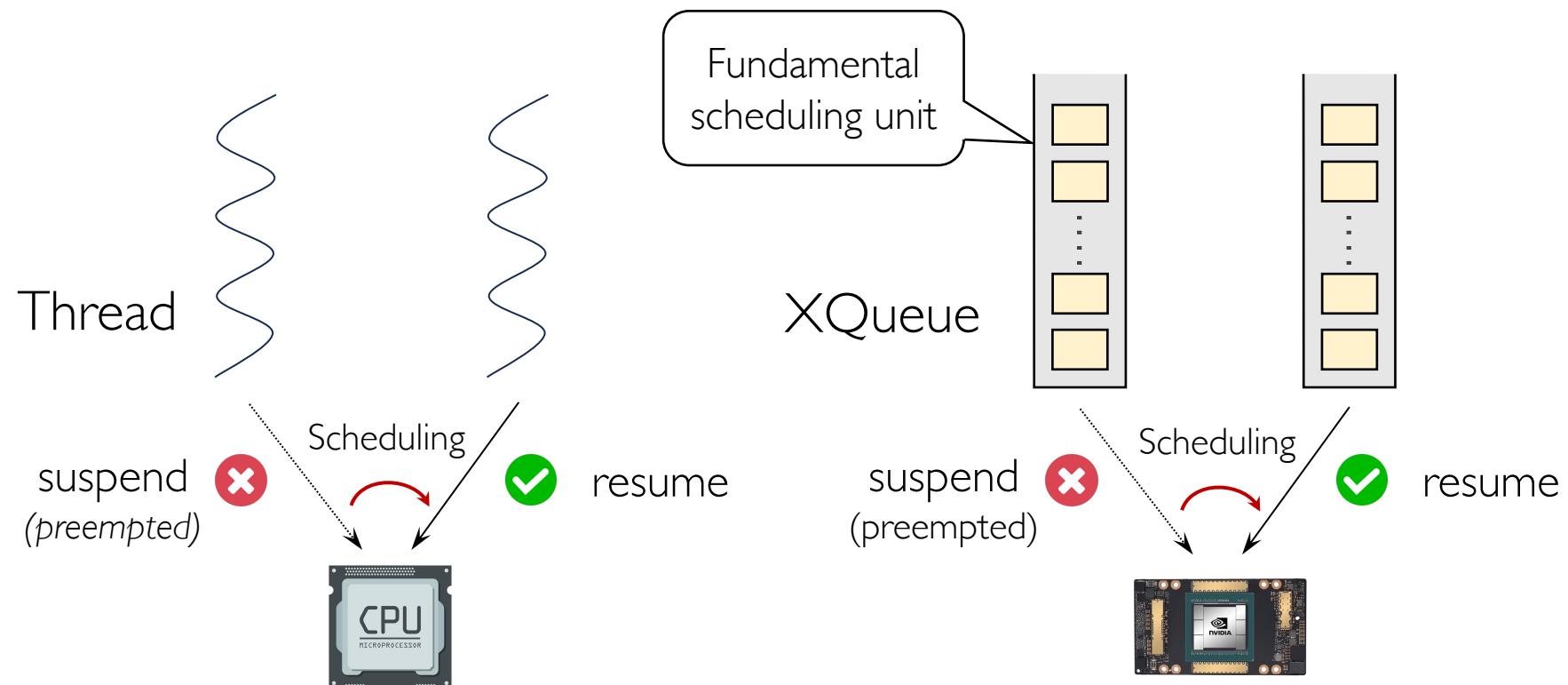


Levelzero Command Queue
OpenCL Command Queue
...

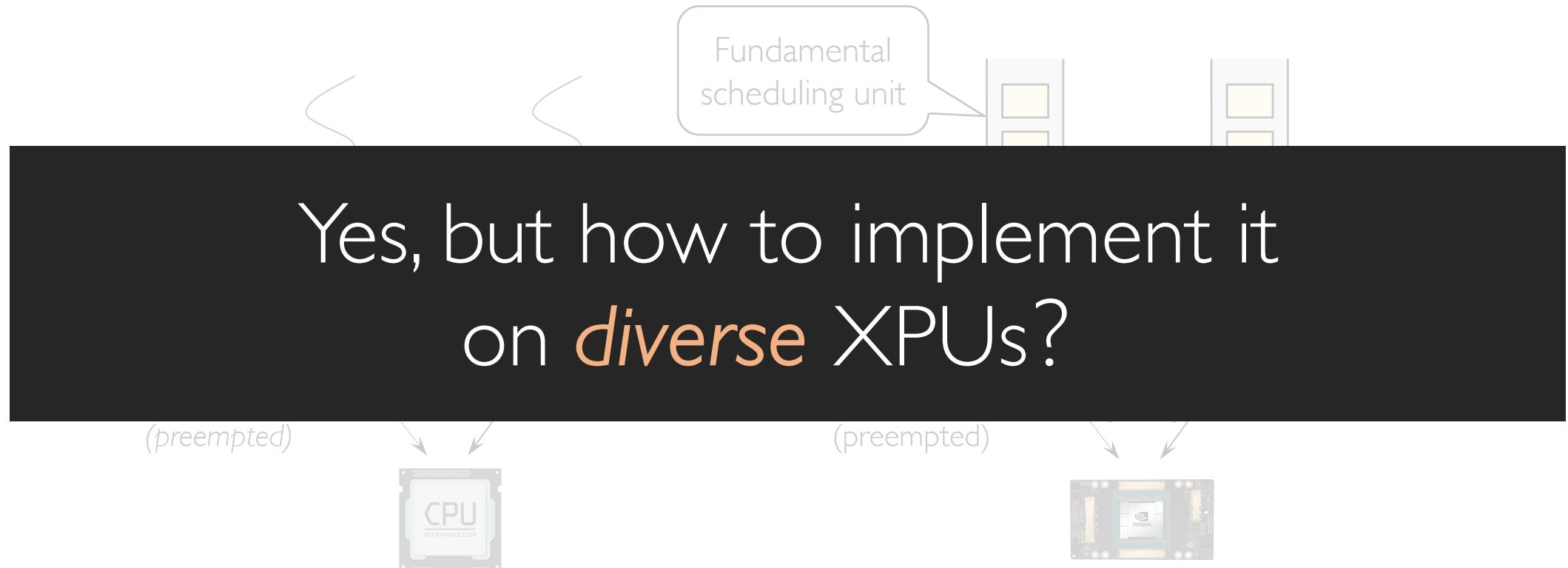
XQueue Abstraction: XPU Command Queue



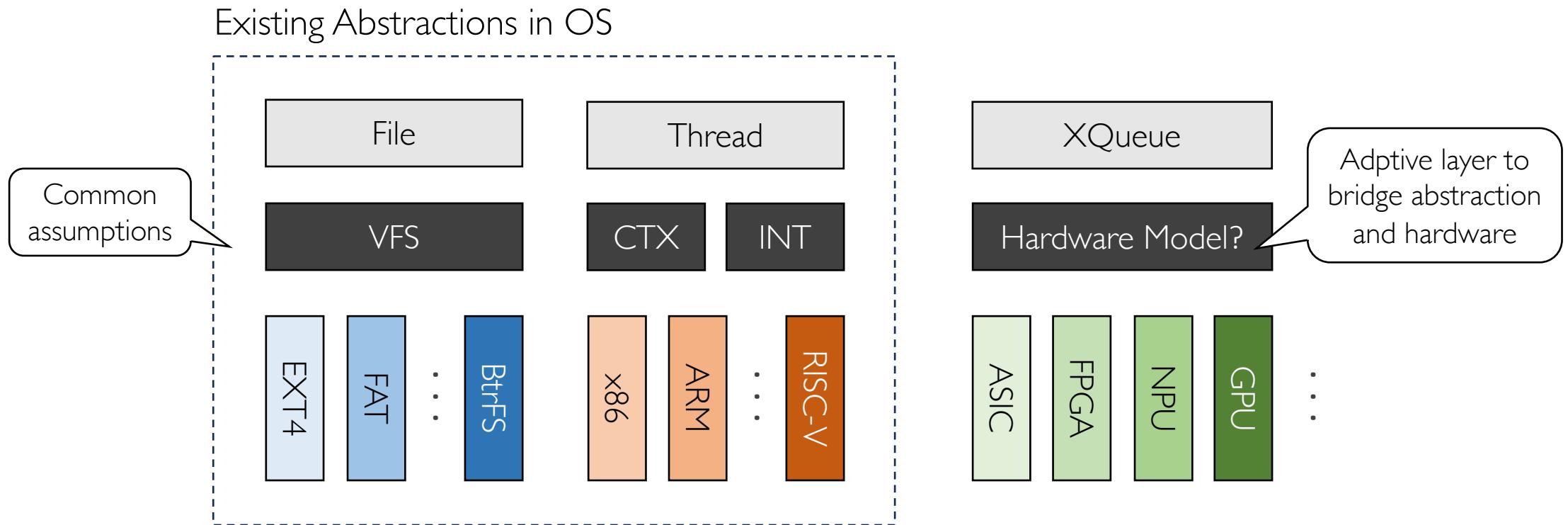
XQueue: *Preemptible* XPU Command Queue



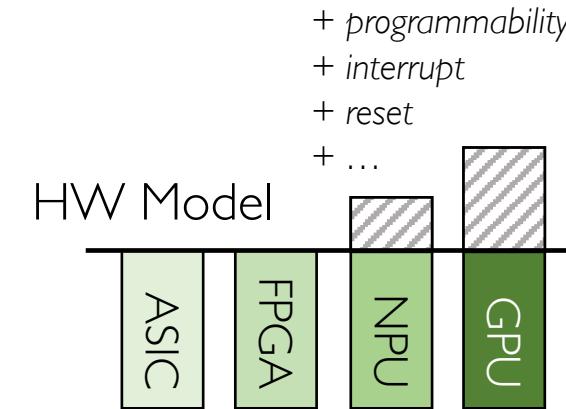
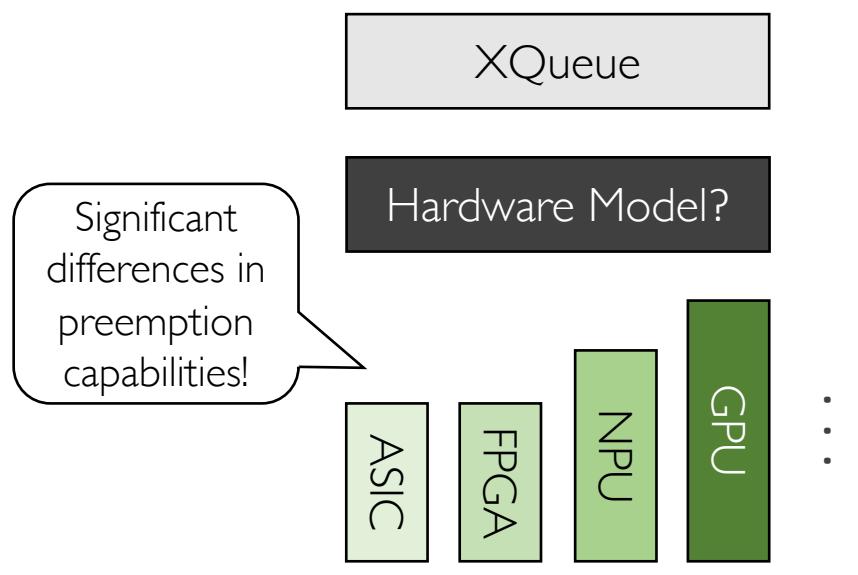
XQueue: *Preemptible* XPU Command Queue



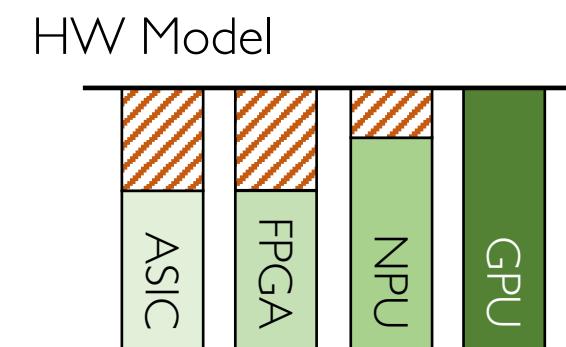
Implementing XQueue on Diverse XPUs



Dilemma: Compatibility or Performance?

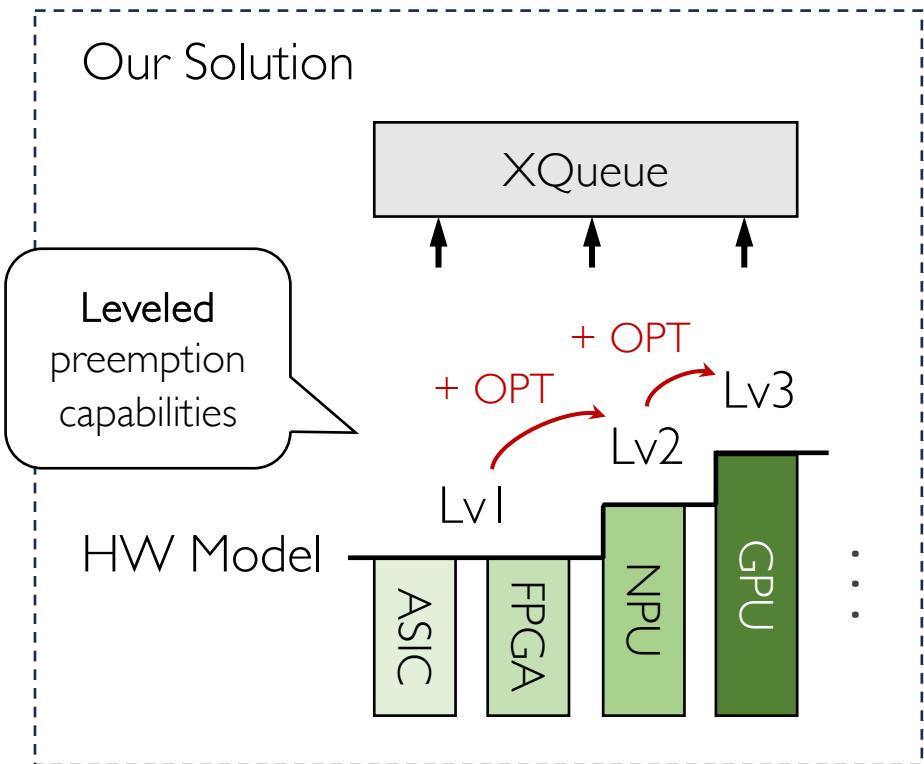


- **Good** compatibility
- **Sacrifice** preemption performance on advanced XPU



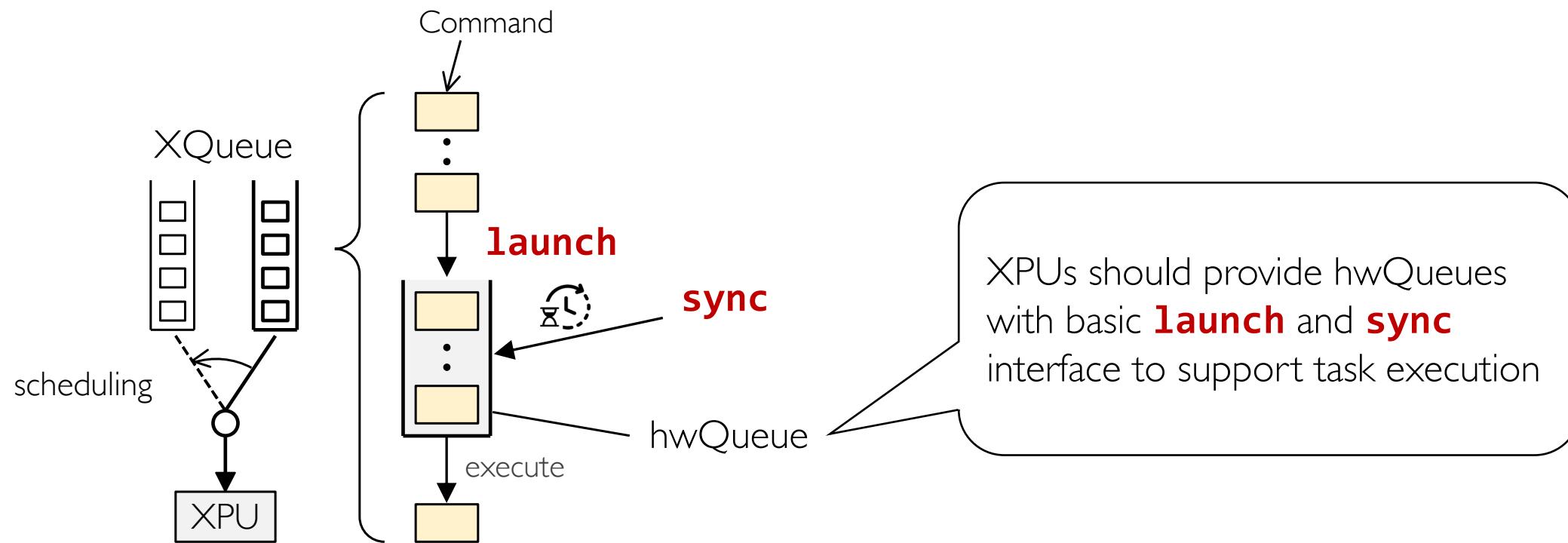
- **Bad** compatibility
- **Full potential** for advanced XPU

Our Solution: Multi-level Hardware Model

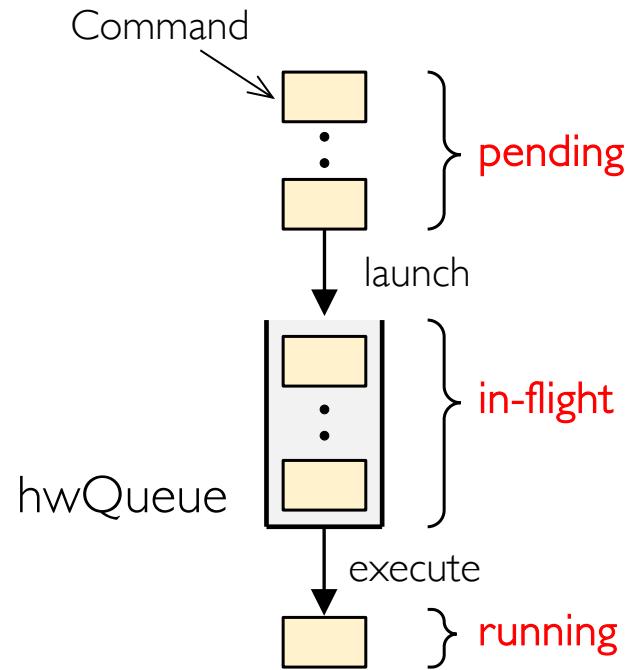


- Categorize preemption-related capabilities into **three incremental levels** (Lv1—Lv3)
- Higher levels posses all capabilities of lower ones and introduce advanced capabilities
- The lowest level (Lv1) is supported by **ALL** XPU
- The highest level (Lv3) unlocks **full potentials** of advanced XPU

Basic Assumption: hwQueue



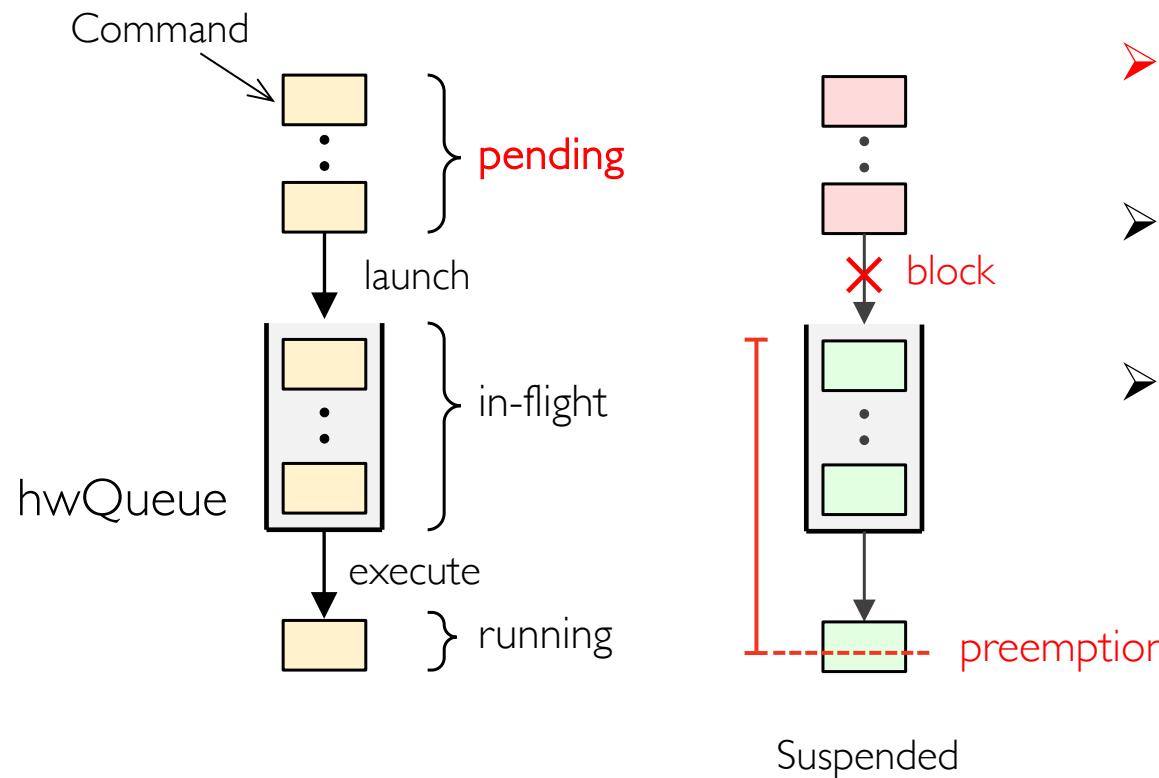
From the Perspective of Command States



- Pending commands
Created and waiting on the host side
- In-flight commands
Launched to a hwQueue, escapes host control
- Running commands
Being executed on the XPU

Level-I: Preempt Pending Command

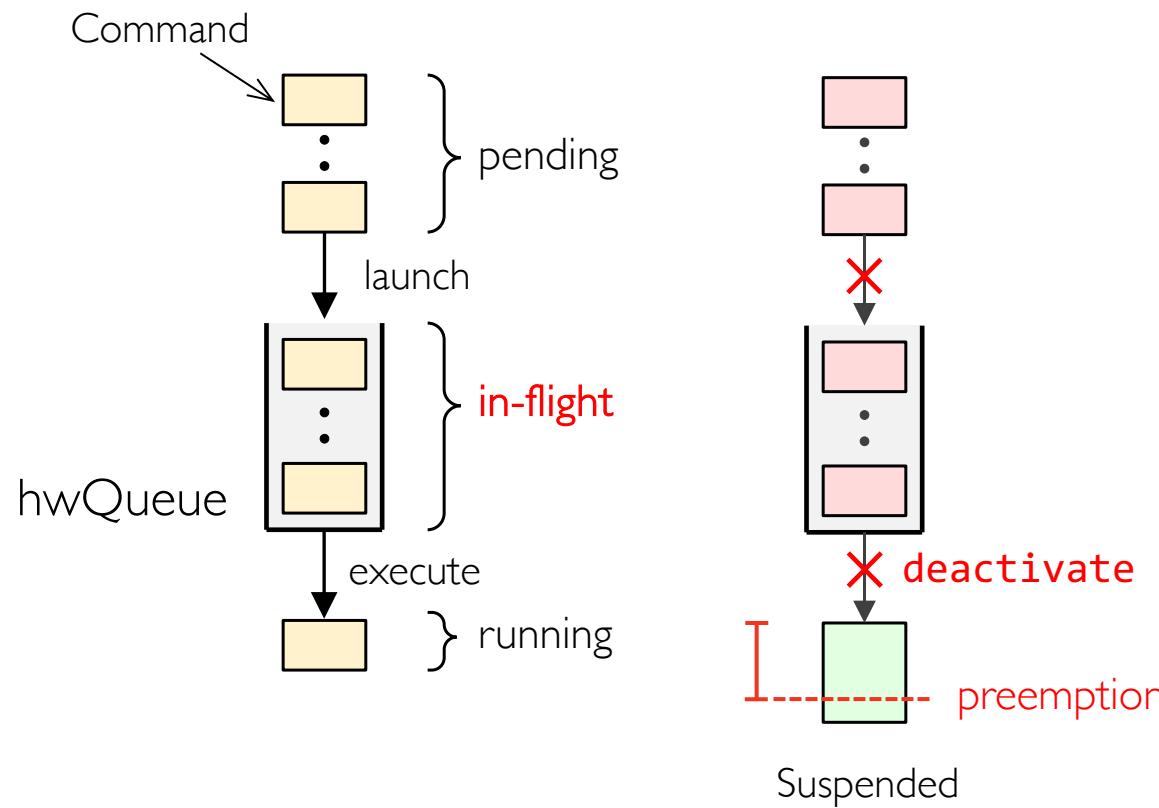
— Preemption Latency Preempted Run-to-complete



- **Block** host CPU from launching new commands, no hardware requirements
- Preemption done until **in-flight** and **running** commands are finished
- OPT: **Progressively** launches the pending commands and **limits the number of in-flight commands** to reduce preemption latency

Level-2: Preempt In-flight Command

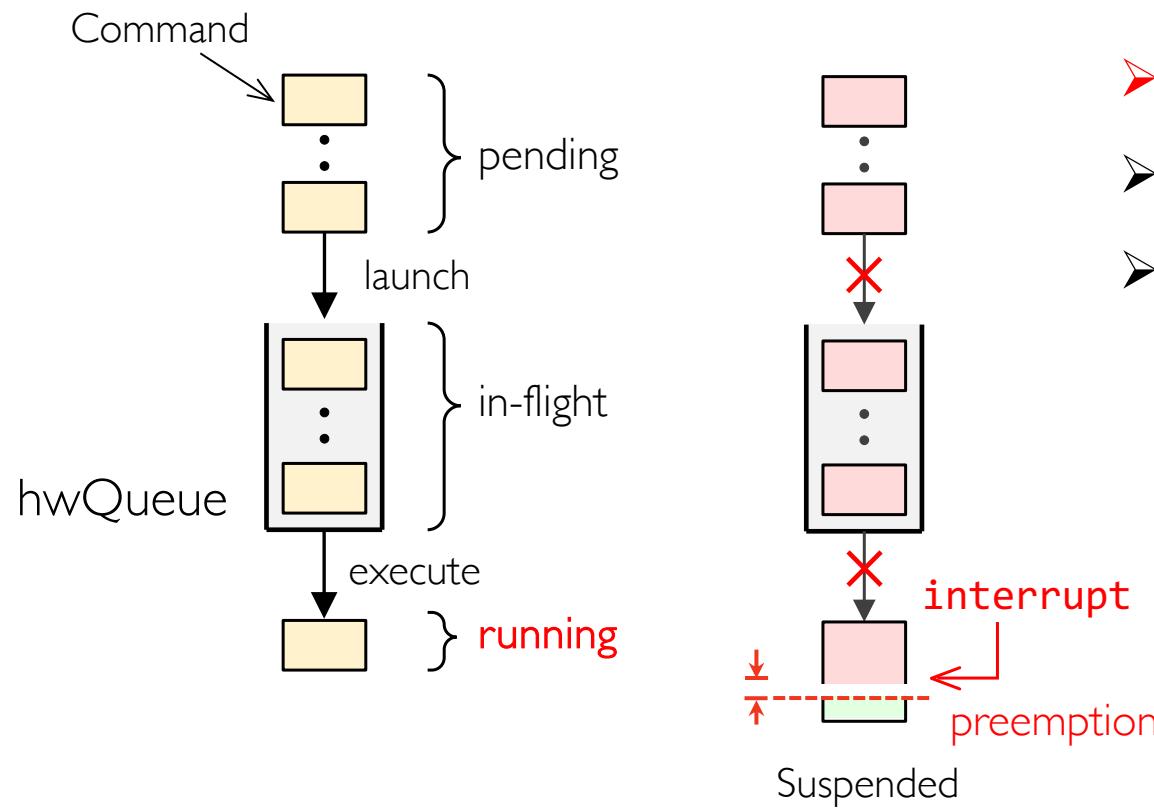
— Preemption Latency Preempted Run-to-complete



- **Deactivate** an **hwQueue** so that its commands will not be executed further
- Significantly reduced preemption latency
- Approach 1: instruct the *μ -controllers* in modern XPU, e.g., Intel NPU, NV & AMD GPU to stall command dispatching
- Approach 2: leverage command **programmability** to abort in-flight commands

Level-3: Preempt Running Command

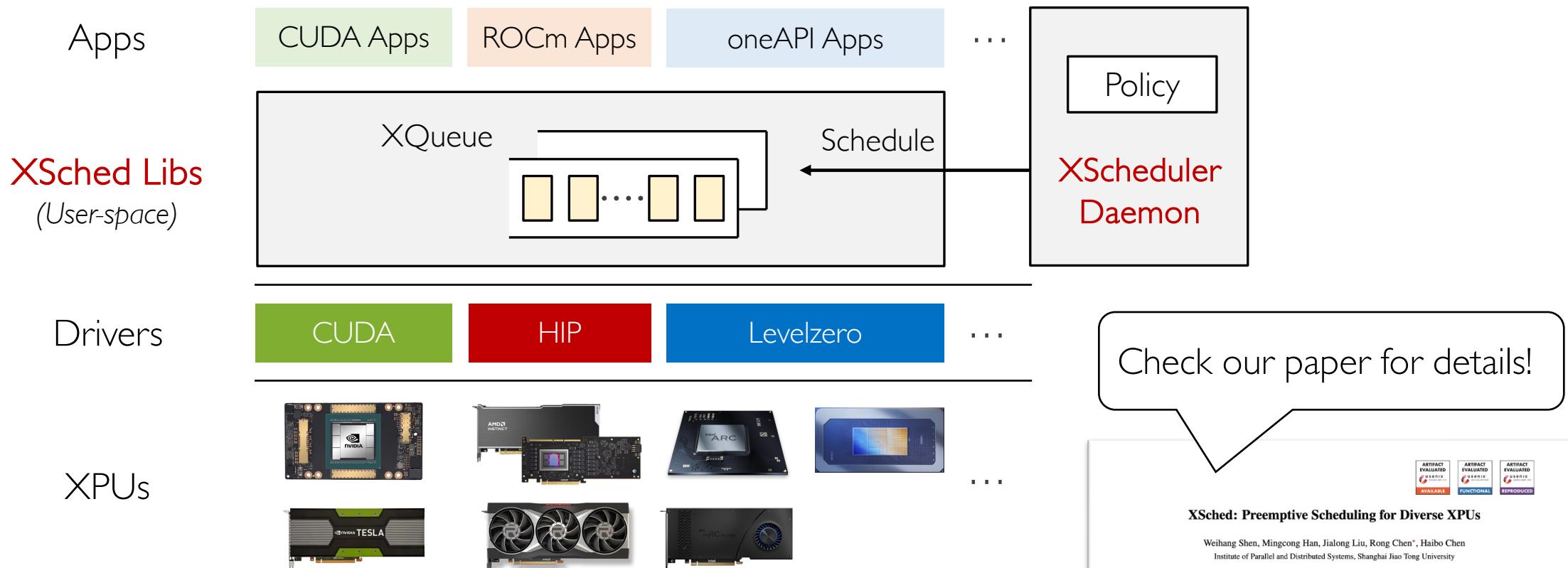
— Preemption Latency Preempted Run-to-complete



- **Interrupt** the running command
- Ultra-low & stable preemption latency
- Supported in modern GPUs

Xsched: A Preemptive XPU Scheduling Framework

based-on the XQueue Abstraction and the Multi-level Hardware Model



Supported & Evaluated XPU



NV GV100 GPU



AMD MI50 GPU



Ascend 910b NPU

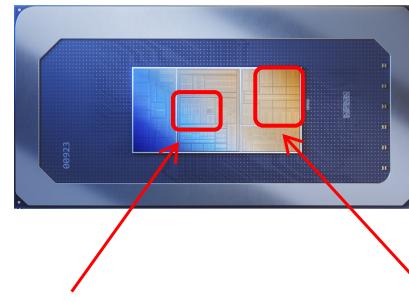


Xilinx VU9P FPGA

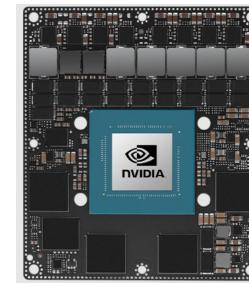
*integrated with AWS-FI



NV K40m GPU

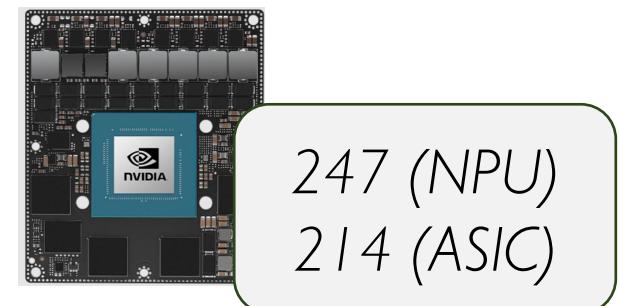
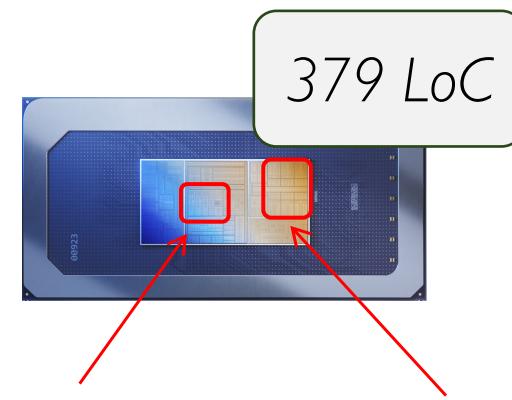


Intel NPU3720 + Arc iGPU
*integrated with Core Ultra 185H



NV DLA (NPU) + OFA (ASIC) + PVA (ASIC)
*integrated with Jetson AGX Orin

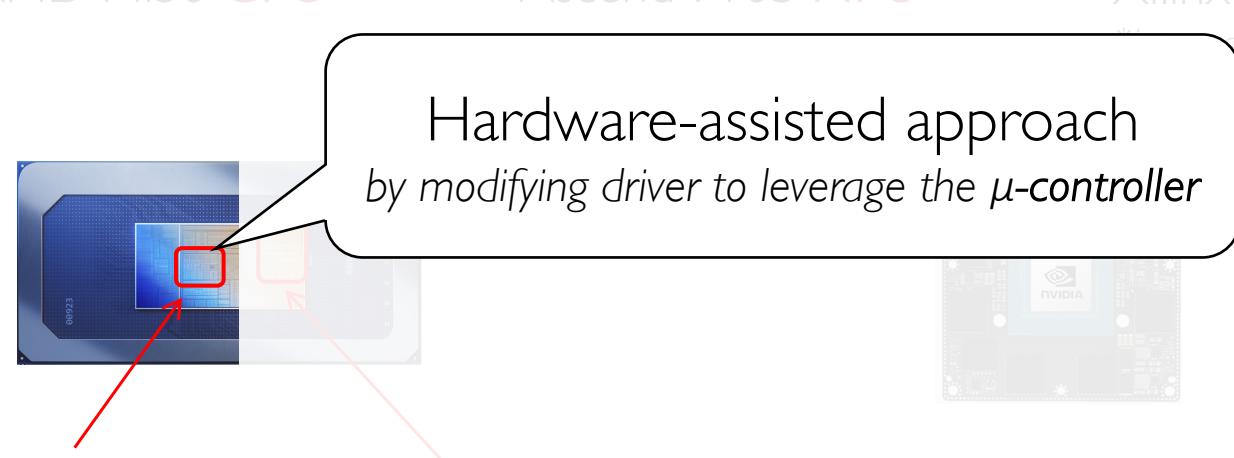
Level-I: *General* Design for *Effortless Integration*



Level-2 Preemption Support



Pure software approach
by modifying GPU kernel binary



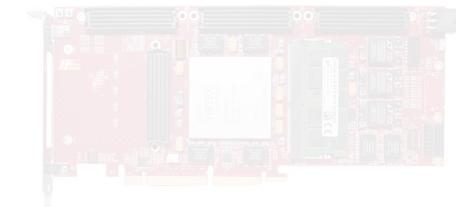
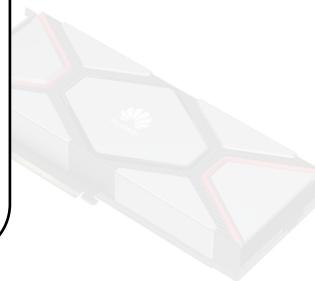
Intel NPU3720 + Arc iGPU
*integrated with Core Ultra 185H

NV DLA (NPU) + OFA (ASIC) + PVA (ASIC)
*integrated with Jetson AGX Orin

Level-3 Preemption Support



Modify GPU *trap handler*
and proactively trigger
GPU *interrupts*



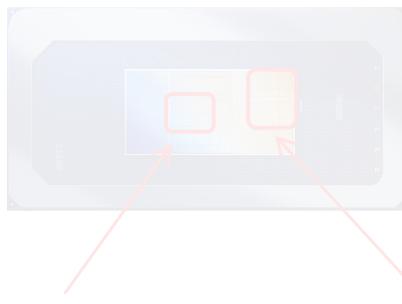
NV GV100 GPU

AMD MI50 GPU

Ascend 910b NPU

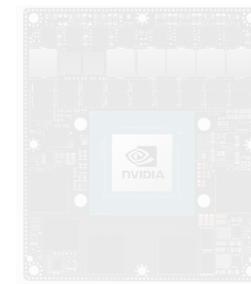
Xilinx VU9P FPGA

*integrated with AWS-FI



NV K40m GPU

Intel NPU3720 + Arc iGPU
*integrated with Core Ultra 185H

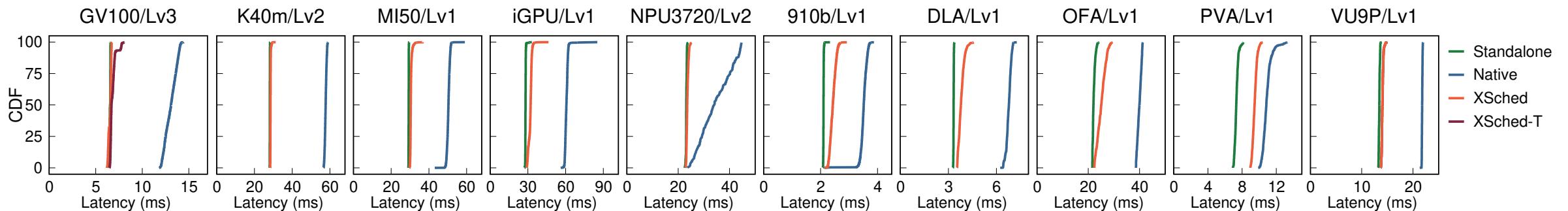


NV DLA (NPU) + OFA (ASIC) + PVA (ASIC)
*integrated with Jetson AGX Orin

Is XSched Effective?

- Policy: Fixed-priority policy
- Result: With **XSched**, latency close to **standalone**, reduce the P99 latency by up to **2.11x** compared with **native**

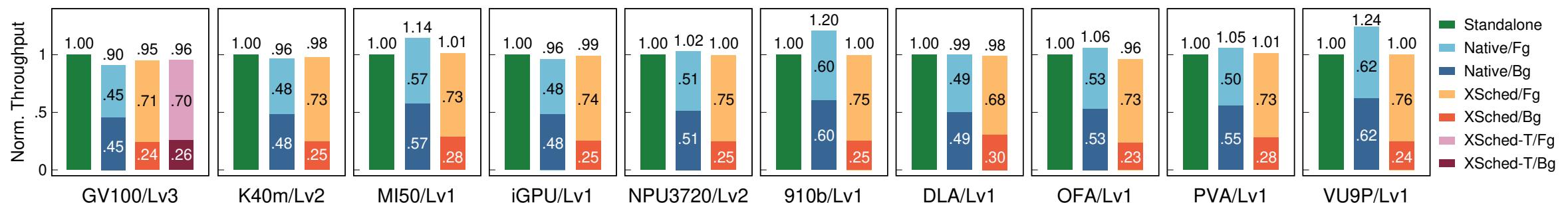
Latency CDF of the high-priority task



Is XSched Effective?

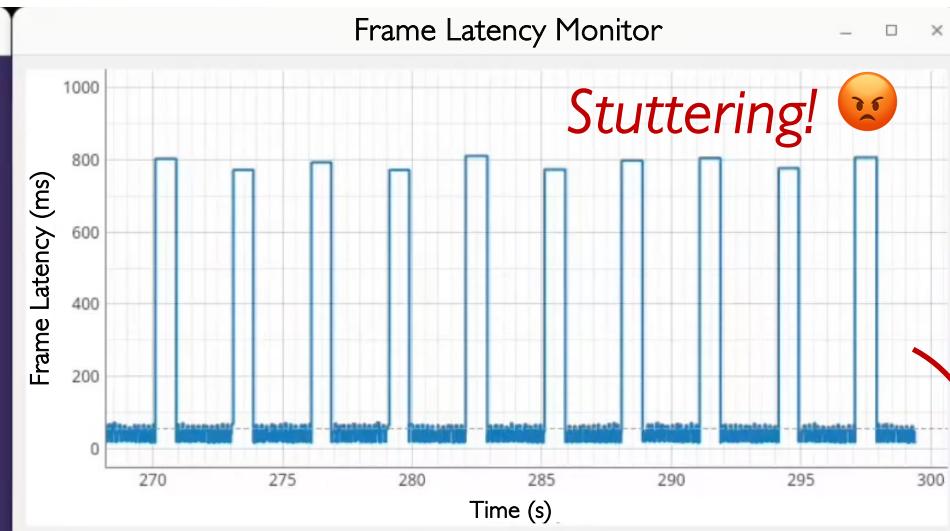
- Policy: Bandwidth partition policy, partition ratio = 75% : 25%
- Result: Guarantee desired throughput ratio with **1.5% overhead**

Normalized throughput of the two process

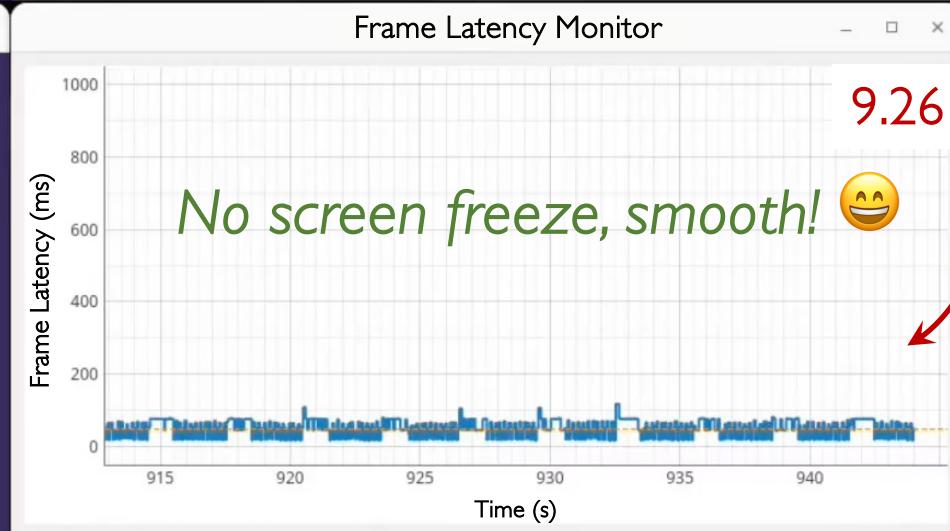


Application: AI Video Conference with XSched

Native

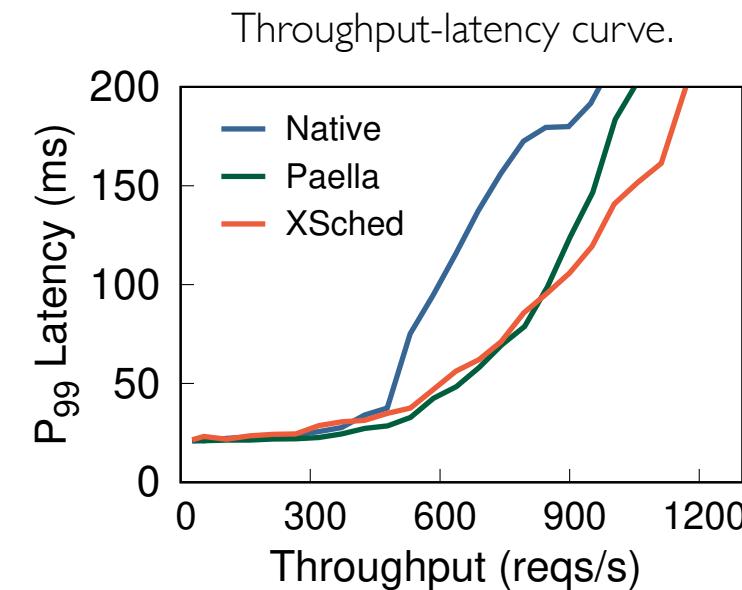
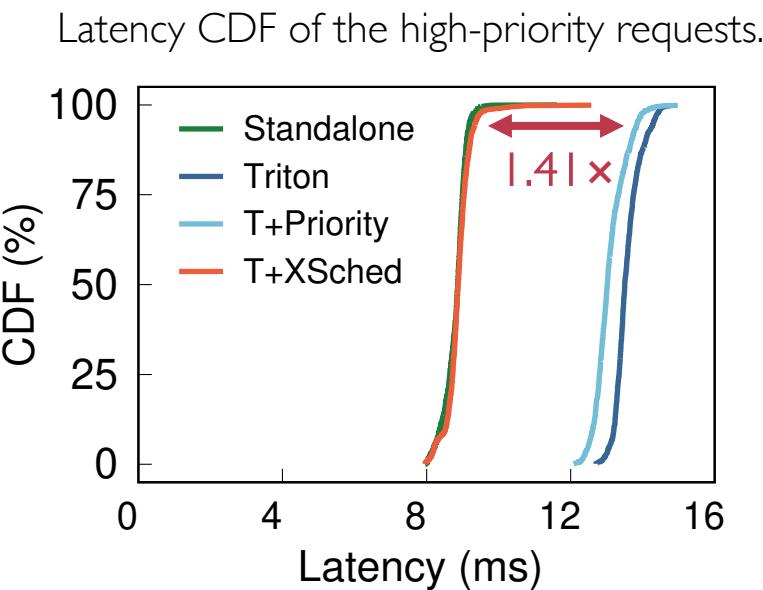


XSched
Laxity-based policy



Application: Inference Serving with XSched

- Integrated with Triton
- $1.41 \times$ reduction on P99
- Comparable with ad-hoc solutions, Paella [SOSP'23]



Conclusion



XSched — A preemptive scheduling framework that supports diverse XPU

Challenges

- Complex XPU software stacks
- Diverse and evolving hardware capabilities



Code

Key Idea

- XQueue — A preemptible command queue abstraction
- Multi-level hardware model (optimal scheduling performance and compatibility)



Publicly available at <https://github.com/XpuOS/xsched>