

Work in Progress: Increasing Schedulability via on-GPU Scheduling

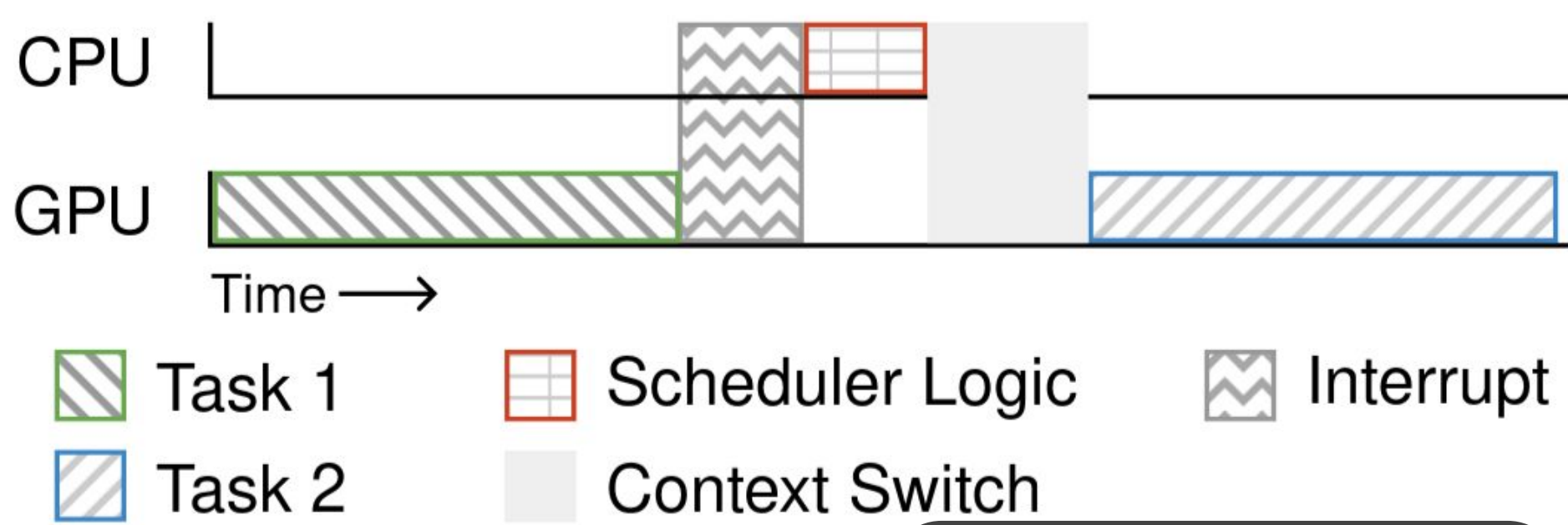
Joshua Bakita (jbakita@cs.unc.edu) and James H. Anderson (anderson@cs.unc.edu)
University of North Carolina at Chapel Hill

The Pitch

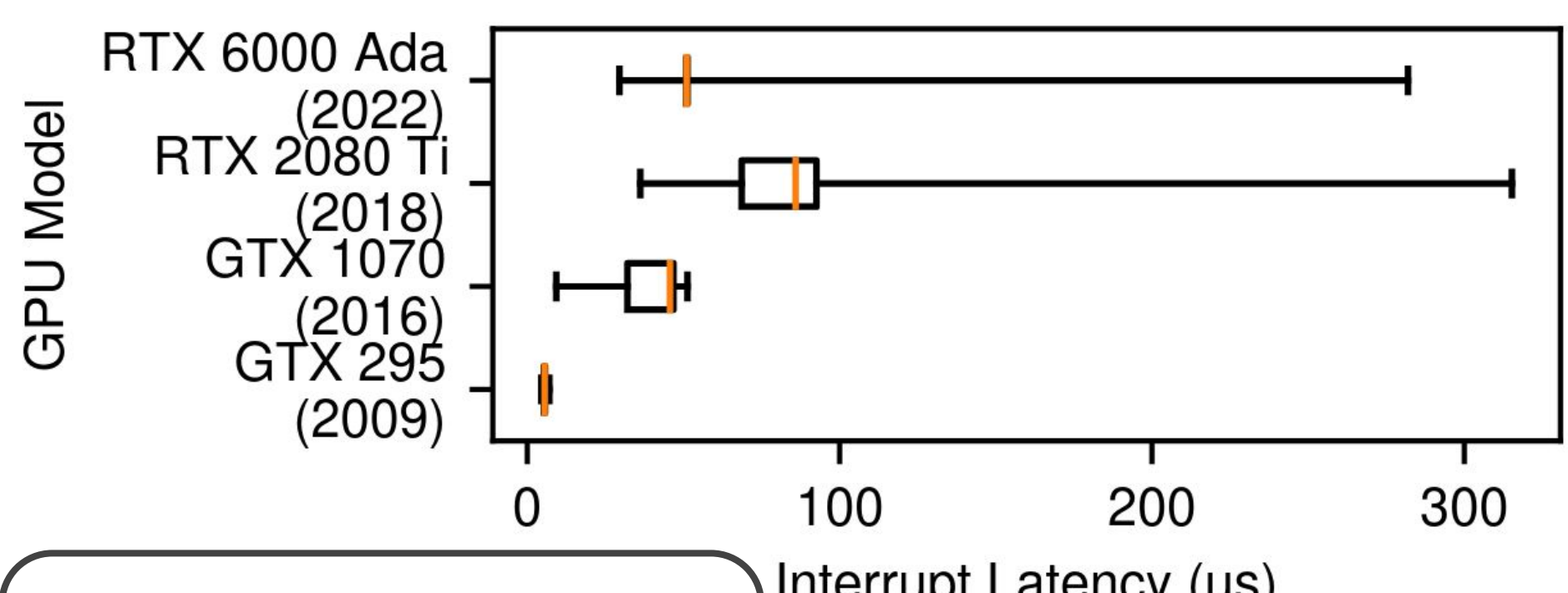
When scheduling multiple tasks on one GPU, by executing the GPU scheduler on the GPU rather than on the CPU, we can reduce absolute overhead, and eliminate capacity loss that occurs on the CPU.

We do this by developing a new scheduling framework for discrete NVIDIA GPUs that runs entirely on the GPU itself.

Problem

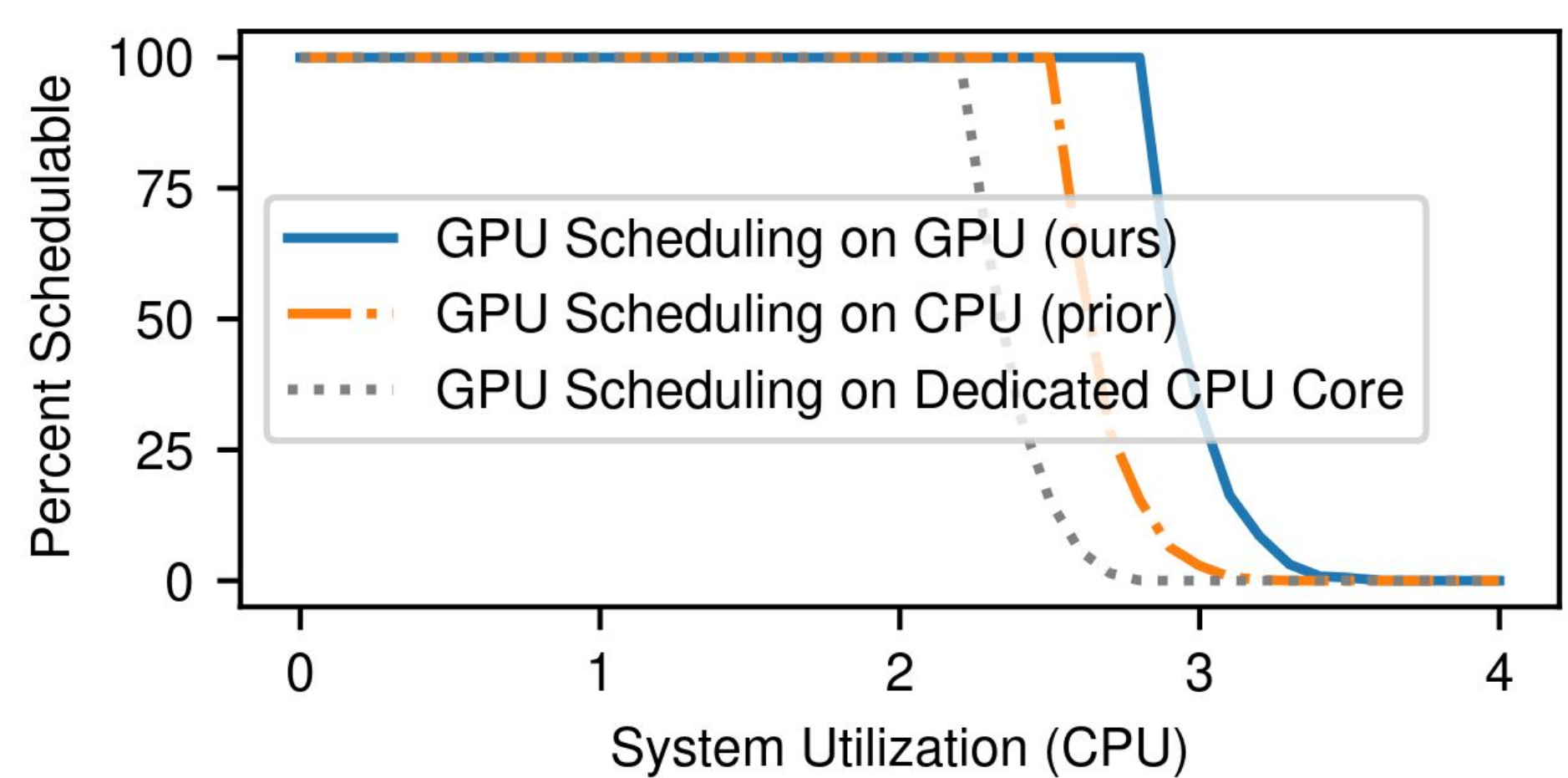


GPU interrupts to the CPU to run scheduler ①



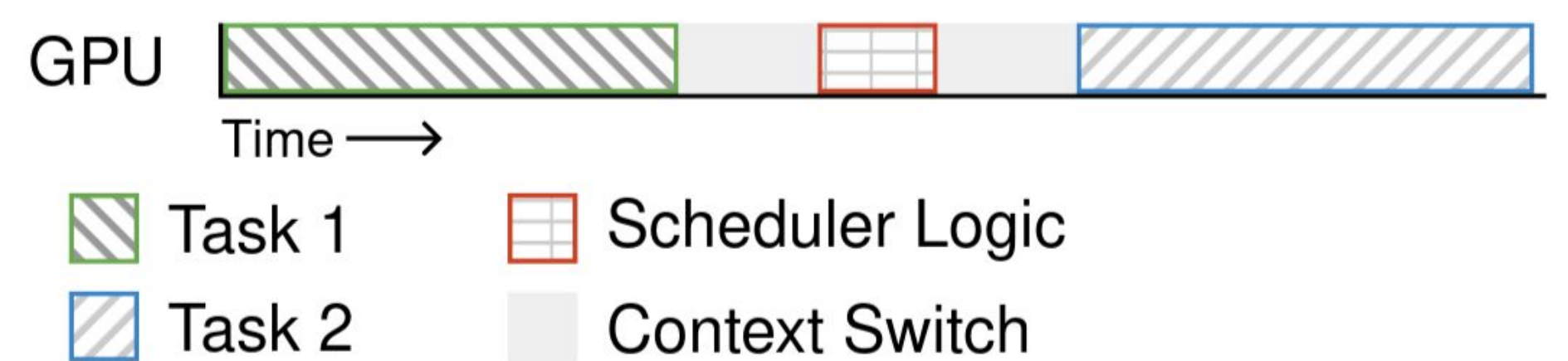
But interrupt latencies can be in the hundreds of microseconds ②

And dedicating a CPU core for spin-based synchronization is expensive ③



Leading to reduced schedulability, esp. when the scheduler must be run frequently ④

Our Solution



We run GPU scheduling synchronously on the GPU ⑤

OVERHEAD OF ON-GPU SCHEDULING ON RTX 6000 ADA

# Tasks	Time w/ Built-in (baseline) (ms)	Time w/ on-GPU (ours) (ms)	Overhead (per invocation)
2	94,218	98,578	95 μ s
3	141,334	147,922	95 μ s
4	188,450	198,021	104 μ s
5	235,540	248,463	112 μ s

Tbl. 1. Absolute overhead of our approach (vs. 421–783 μ s typical with GCAPS)

Yielding lower absolute overheads ⑥

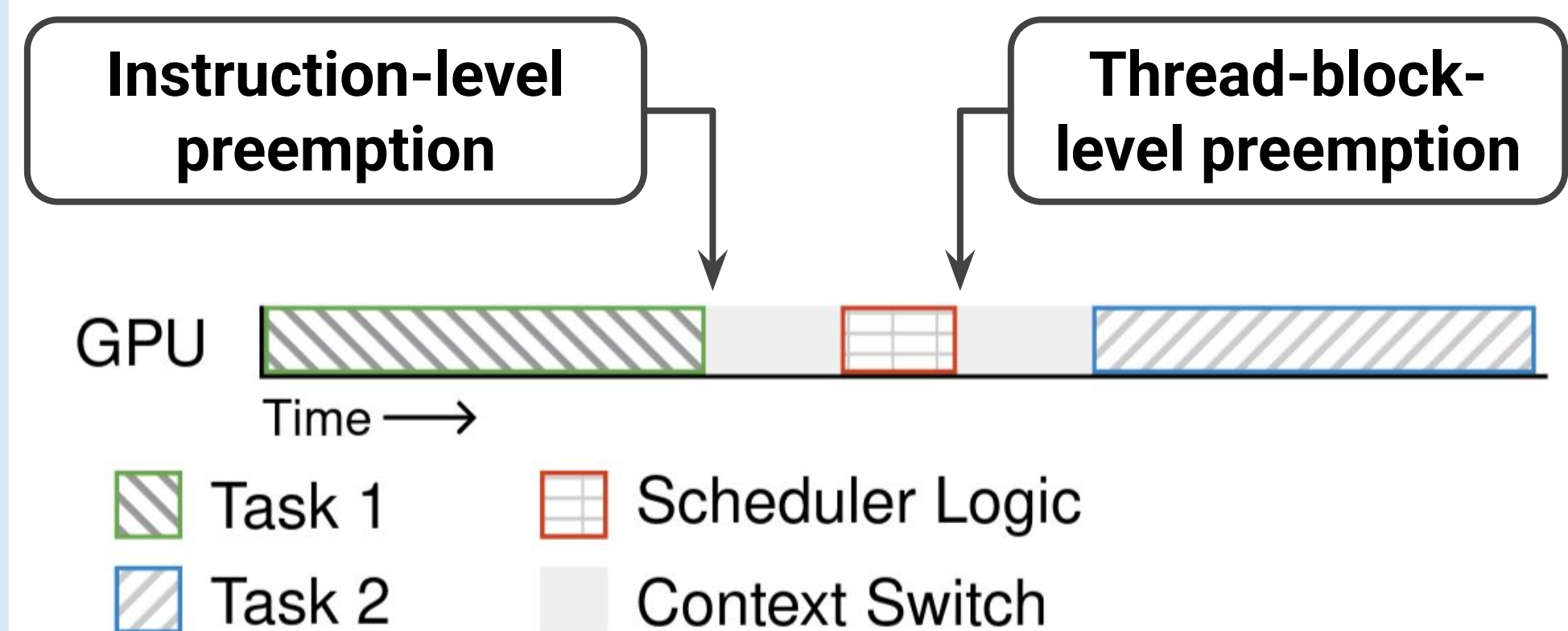
And eliminating any CPU time cost ⑦

Framework Features

Our GPU scheduling framework:

- Works **on many NVIDIA GPUs** (2018+)
- Supports **arbitrary schedulers** via a plugin-based interface
- **No task modifications** required
- **No hardware modifications** required
- Works for **any task type** (CUDA, OpenGL, Vulkan, etc.)
- **Preserves logical isolation** between tasks
- Uses **hardware-enforced** time budgeting

Limiting Preemption Overhead



RTAS
2025

UNC
REAL-TIME SYSTEMS GROUP



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

Full implementation of our work is under submission to appear at an operating systems conference later in 2025.