

LET: A Way Forward for Safe GPU Co-Scheduling

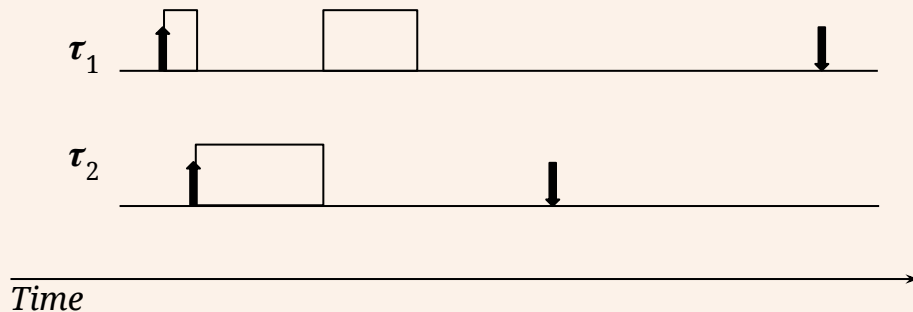


Nathan Otterness, James H. Anderson

Prerequisite: Timing Models

Tasks (implicit deadline):

- $\tau_1: (0.4, 2.0)$
- $\tau_2: (0.2, 1.0)$



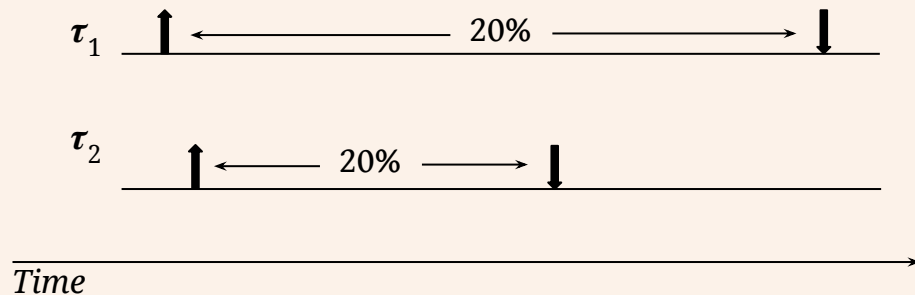
Real-time scheduling often assumes a *Bounded Execution Time* (BET) model:

Tasks must occupy a processor for a specific time interval before their deadline.

Prerequisite: Timing Models

Tasks (implicit deadline):

- $\tau_1: (0.4, 2.0), U = 0.2$
- $\tau_2: (0.2, 1.0), U = 0.2$



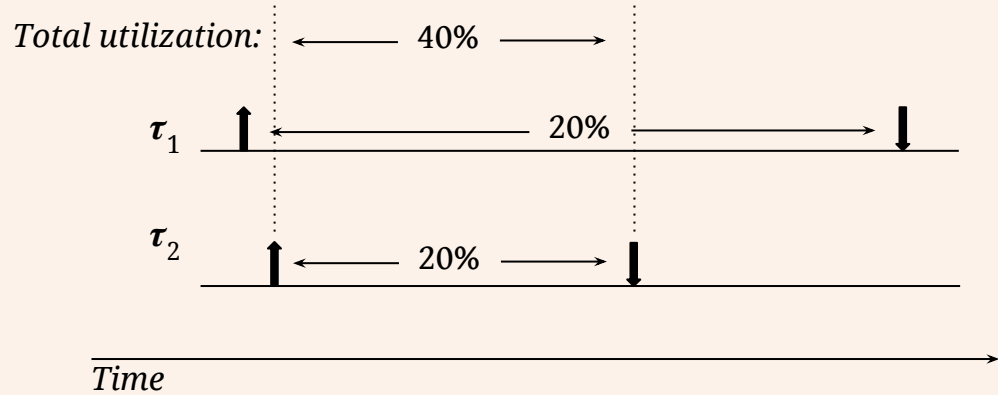
Logical Execution Time (LET):

Tasks occupy a *proportion of resources* for their *entire period*.

Prerequisite: Timing Models

Tasks (implicit deadline):

- $\tau_1: (0.4, 2.0), U = 0.2$
- $\tau_2: (0.2, 1.0), U = 0.2$



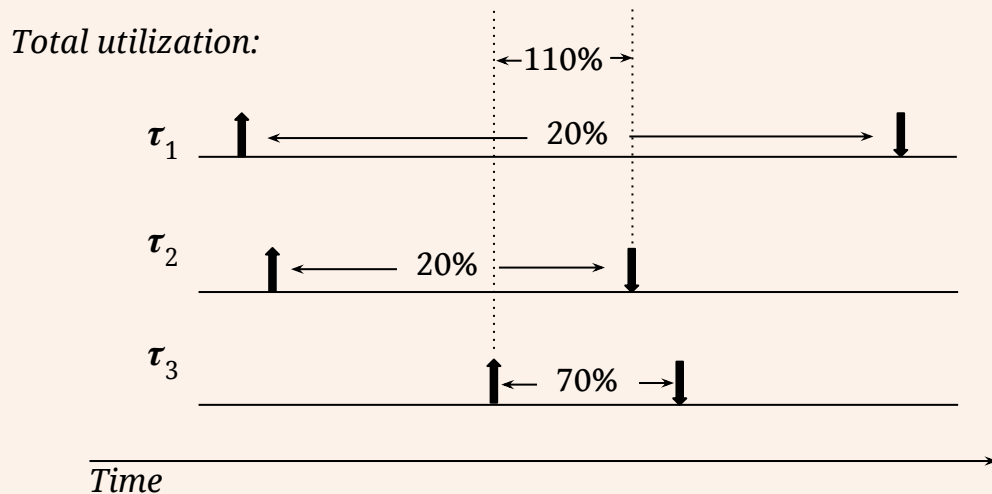
Logical Execution Time (LET):

Tasks occupy a *proportion of resources* for their *entire period*.

Prerequisite: Timing Models

Tasks (implicit deadline):

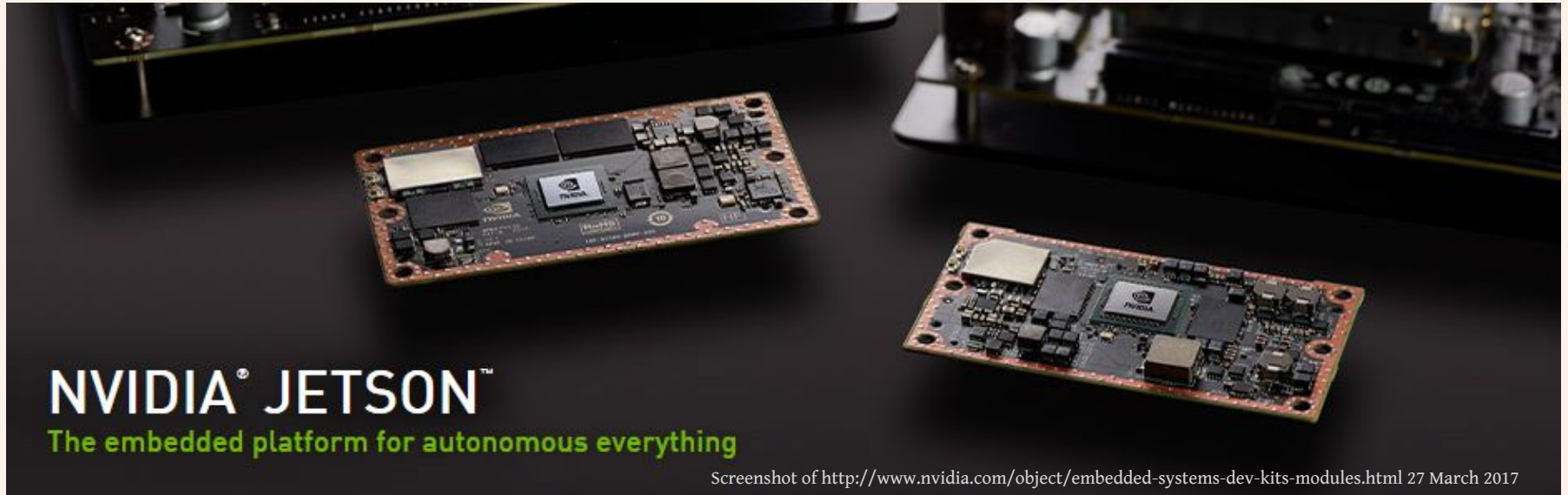
- $\tau_1: (0.4, 2.0), U = 0.2$
- $\tau_2: (0.2, 1.0), U = 0.2$
- $\tau_3: (0.35, 0.5), U = 0.7$



Logical Execution Time (LET):

Tasks occupy a *proportion of resources* for their *entire period*.

Safety-Critical GPUs



NVIDIA® JETSON™
The embedded platform for autonomous everything

Screenshot of <http://www.nvidia.com/object/embedded-systems-dev-kits-modules.html> 27 March 2017

Platforms augmented with graphics processing units (GPUs), such as the NVIDIA Jetson TX1, are increasingly prevalent in embedded systems.

Safety-Critical GPUs

Despite a lack of documentation needed for modeling and *certification*, work is underway to incorporate GPUs into safety-critical systems.

Featured Automotive Partners

NVIDIA partners with some of today's most forward-looking automakers, tier-1 suppliers, research, and startup companies to integrate GPU technology and artificial intelligence to develop self-driving cars, trucks, and shuttles. Their innovations in GPU-based supercomputing enable deep learning, natural language processing, and gesture control that will change how people drive cars—and even enable cars to drive people.



TOYOTA

NVIDIA is collaborating with Toyota to deliver artificial intelligence hardware and software technologies that will enhance the capabilities of autonomous driving systems planned for market introduction within the next few years.

[Read More](#)



AUDI

At CES 2017, Audi and NVIDIA announced an acceleration of a long-running partnership—this new shared goal will put advanced AI cars on the road starting in 2020. Together, Audi and NVIDIA have been delivering automotive breakthroughs for over a decade. Currently, Audi's award-winning Audi connect display systems are powered by NVIDIA and come in every car they make.

[Read More](#)



TESLA

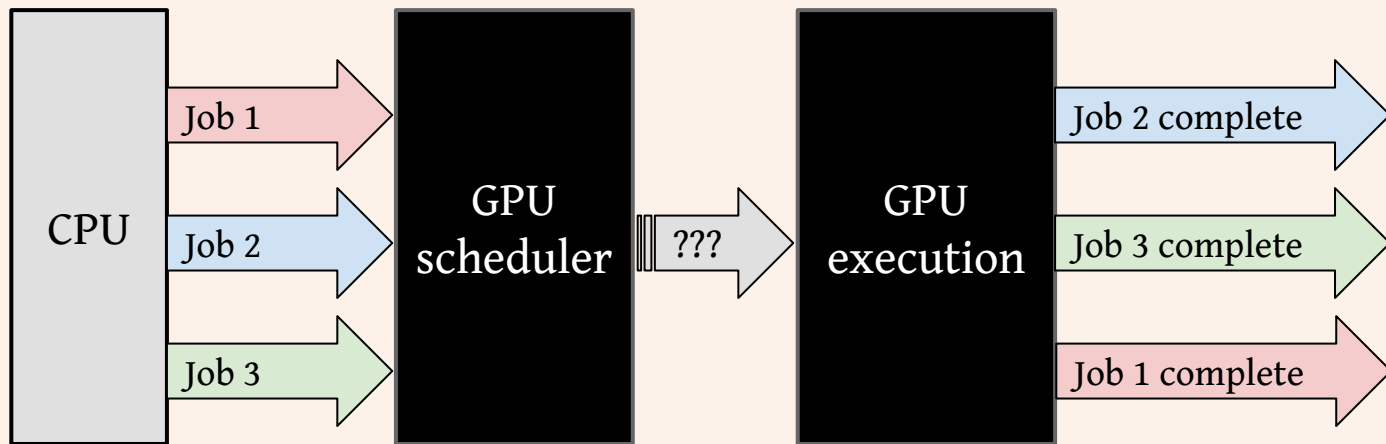
Tesla Motors and NVIDIA have partnered since the early development of the revolutionary Model S.



MERCEDES-BENZ

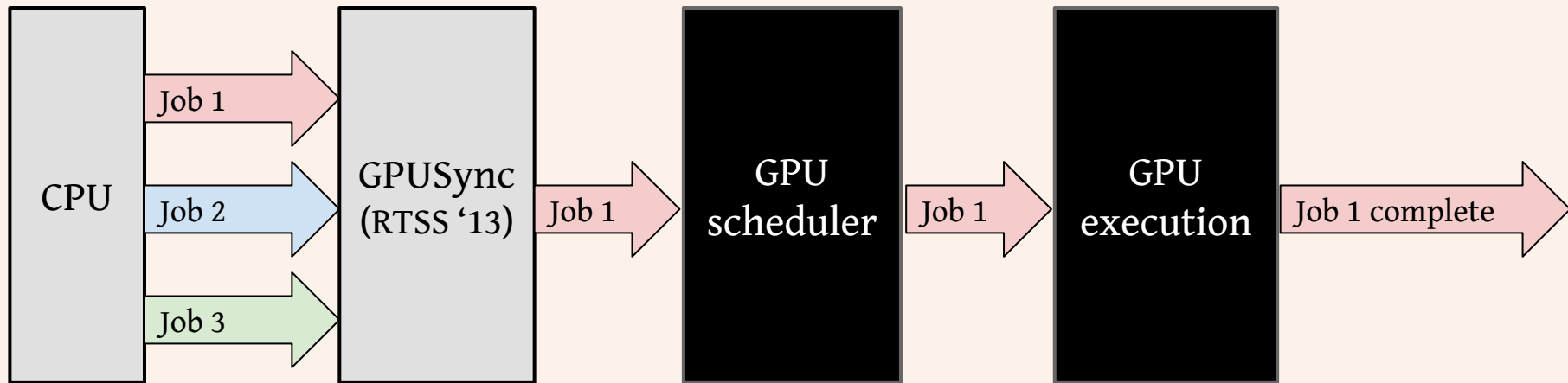
Mercedes-Benz and NVIDIA have announced a partnership to bring an NVIDIA AI-powered car to

GPU Co-Scheduling



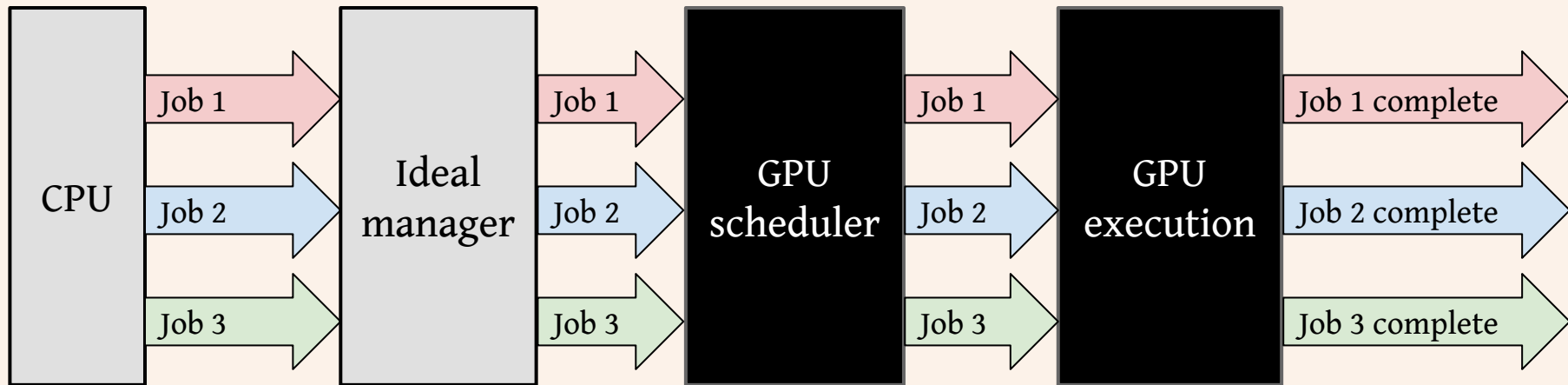
The problem with GPU co-scheduling is that a lack of information leads to a lack of predictability.

GPU Co-Scheduling



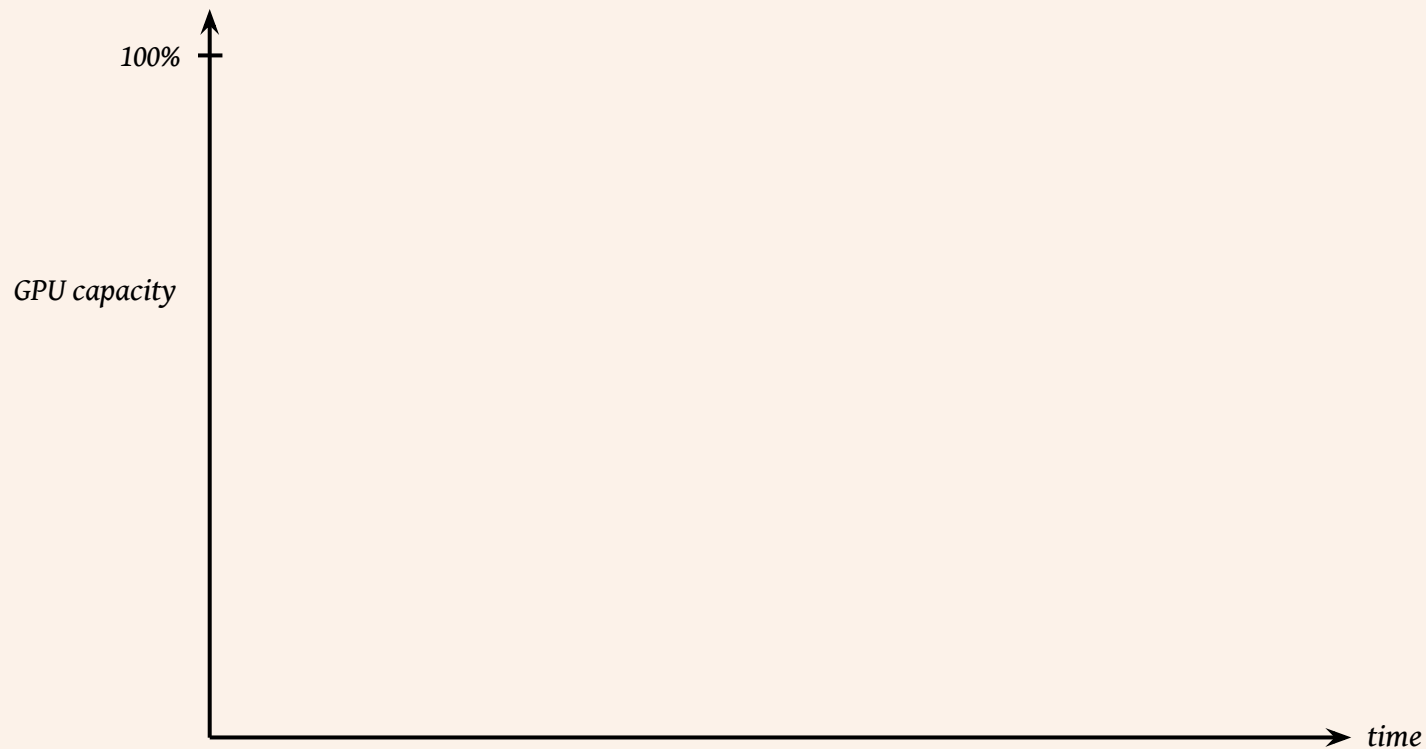
Earlier systems work around this problem by enforcing exclusive access to GPUs.

GPU Co-Scheduling

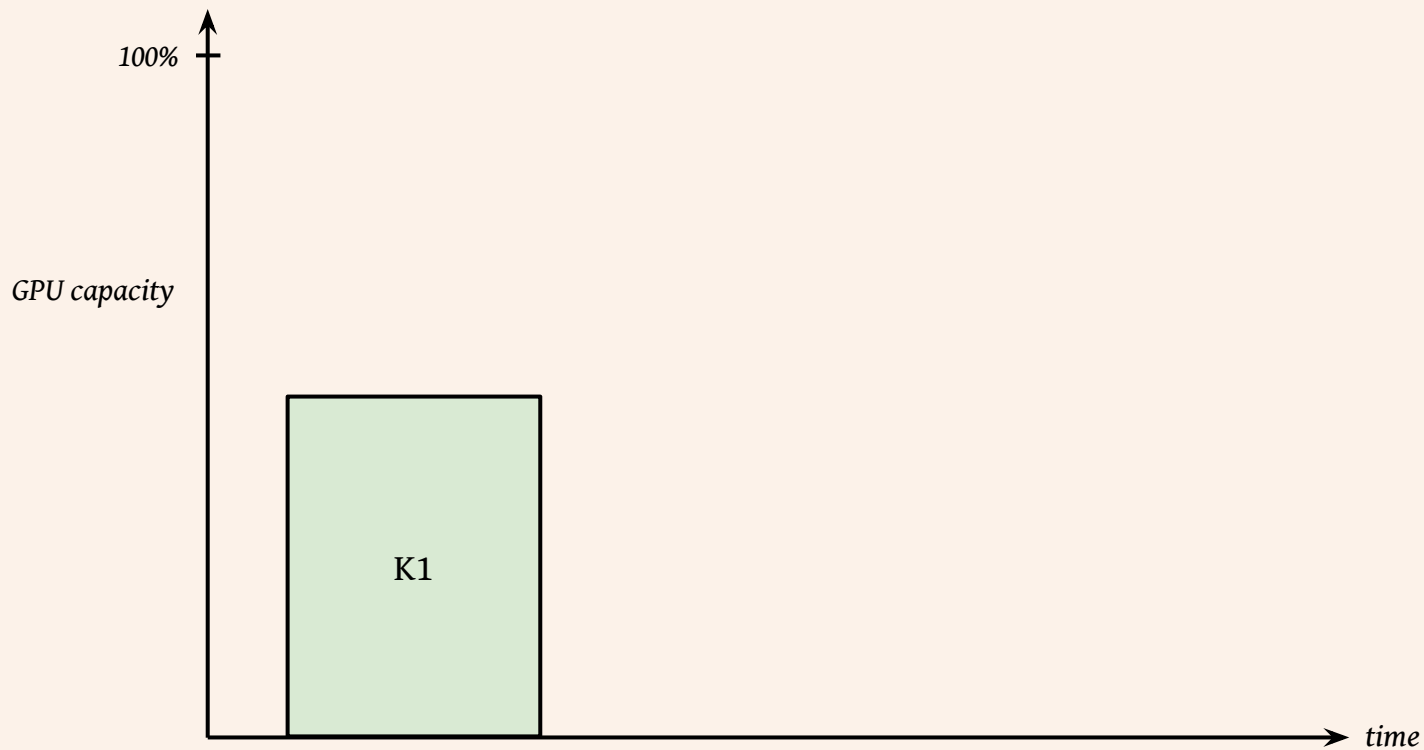


An ideal management system will enable both predictability and concurrency.

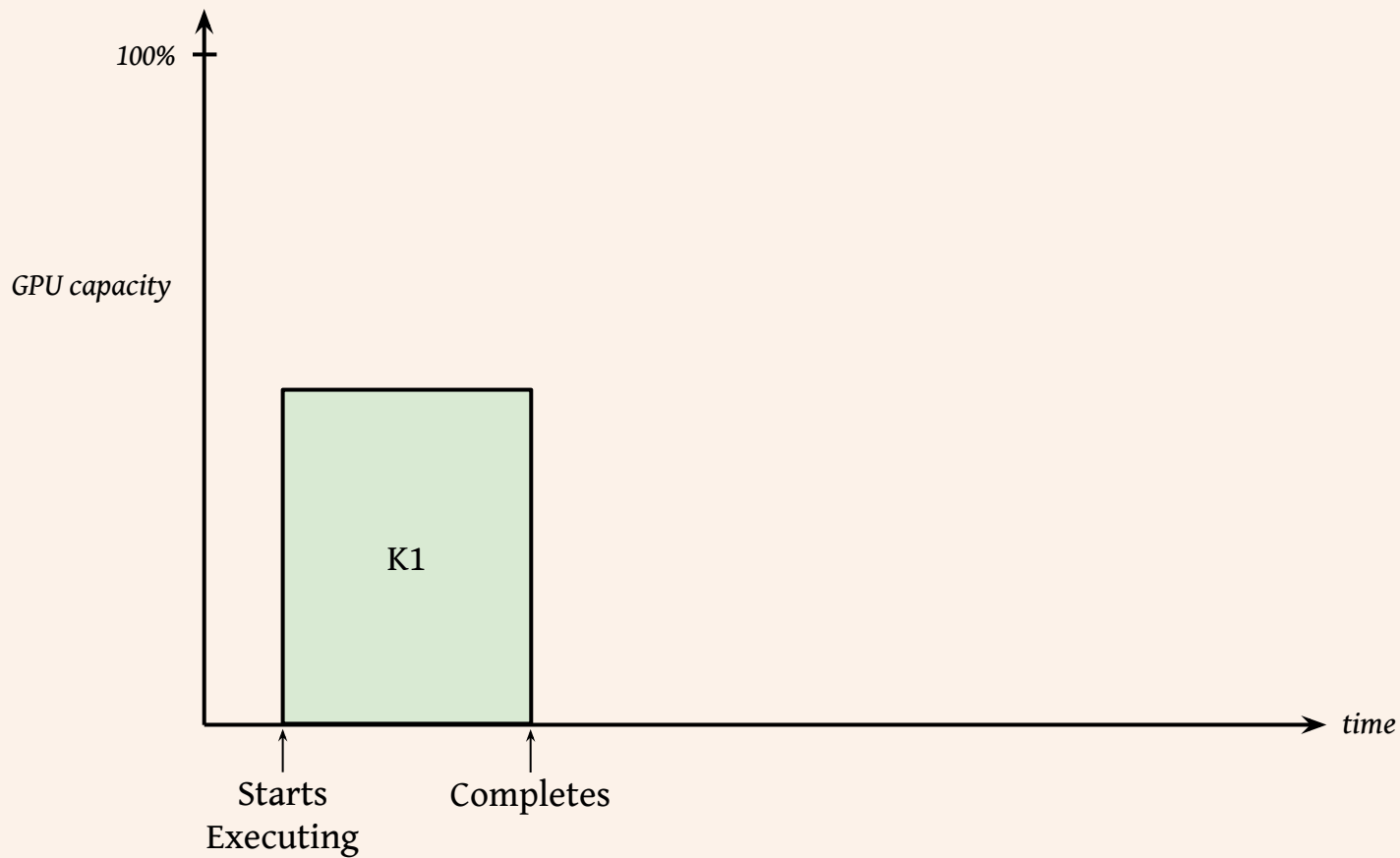
A Simplified GPU Model



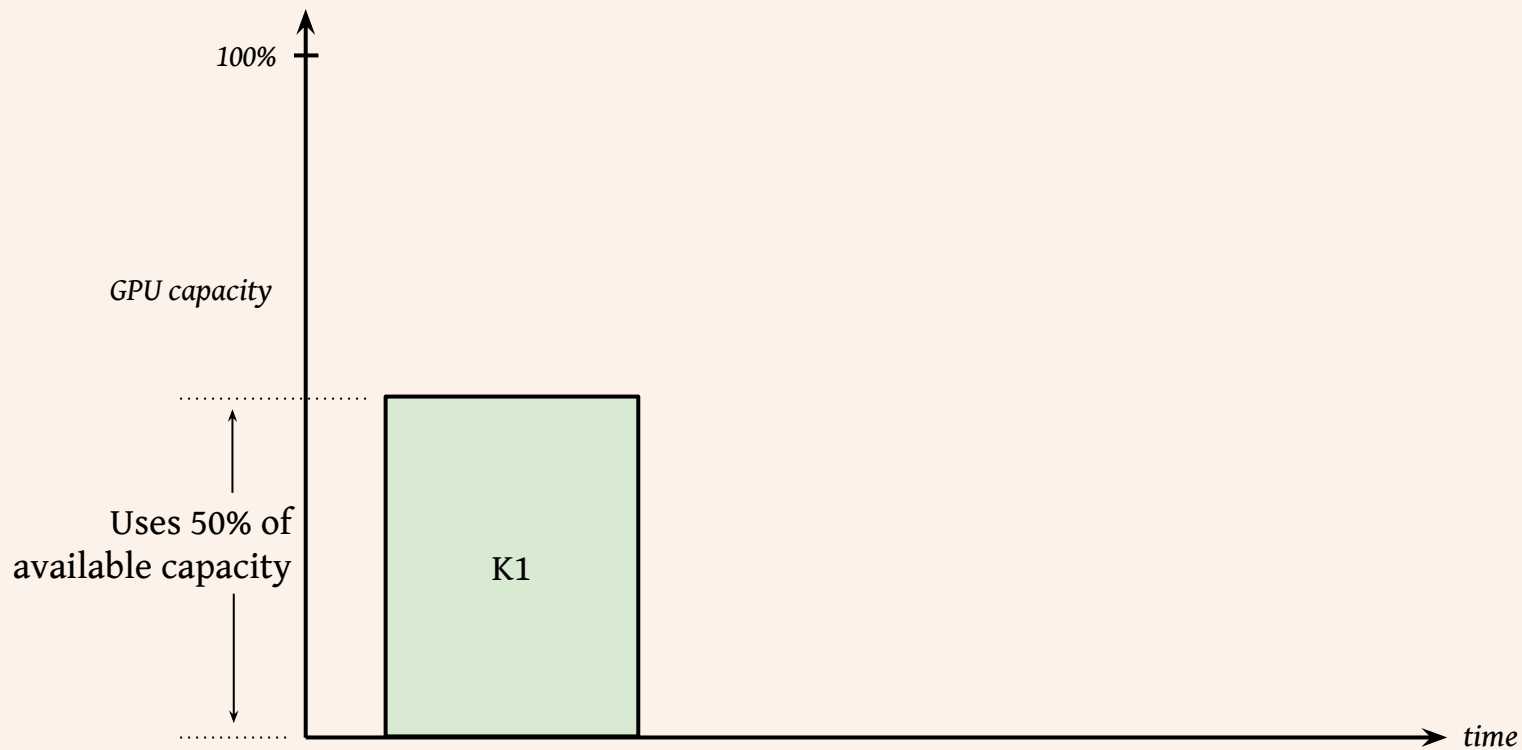
A Simplified GPU Model



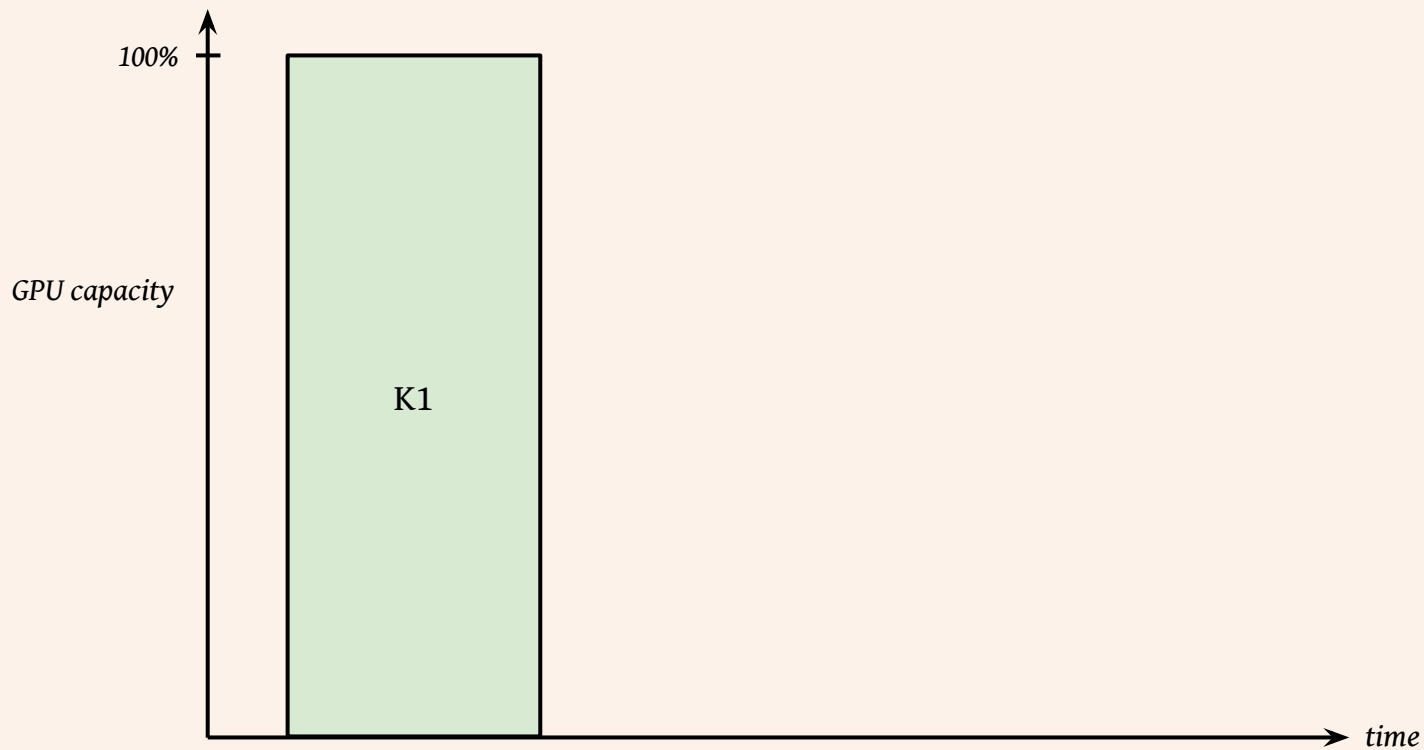
A Simplified GPU Model



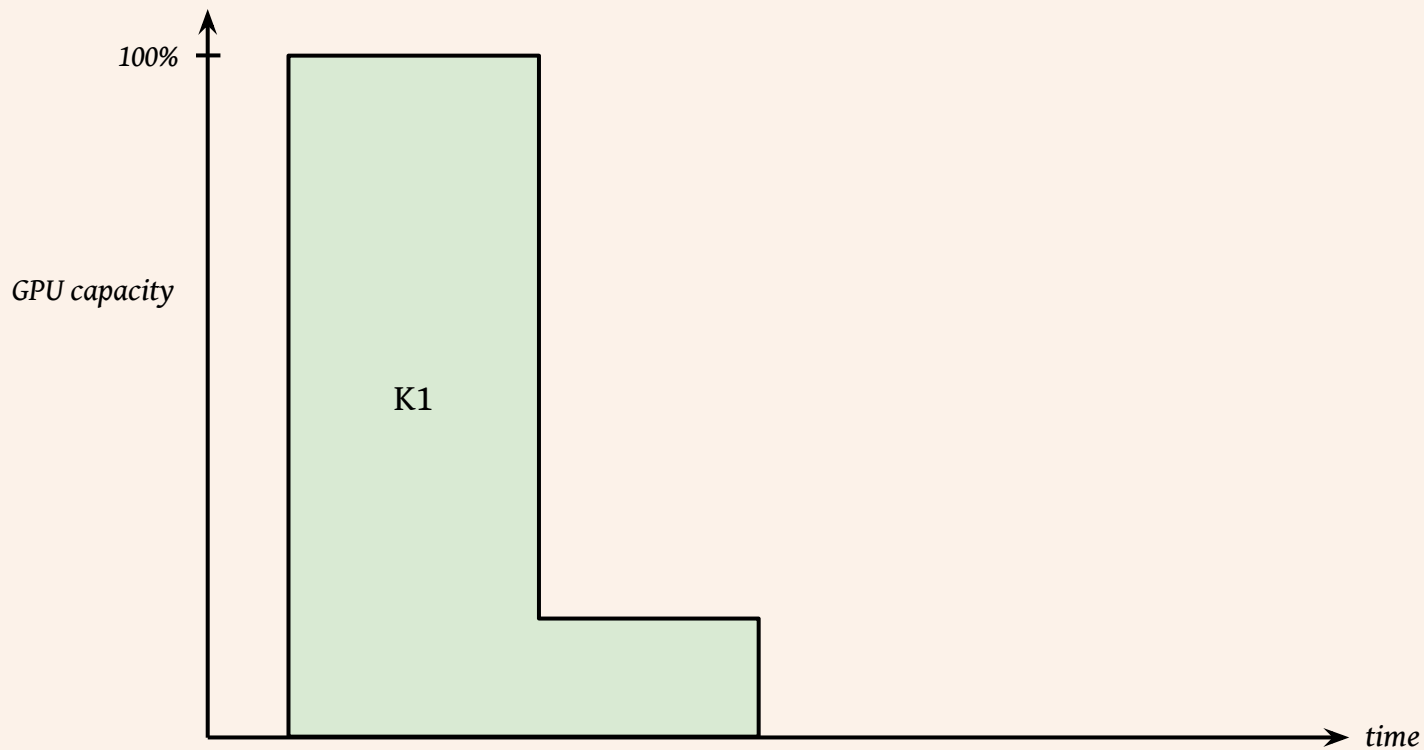
A Simplified GPU Model



A Simplified GPU Model



A Simplified GPU Model



Approaches to Co-Scheduling

Different ways to Co-Schedule GPU Tasks:

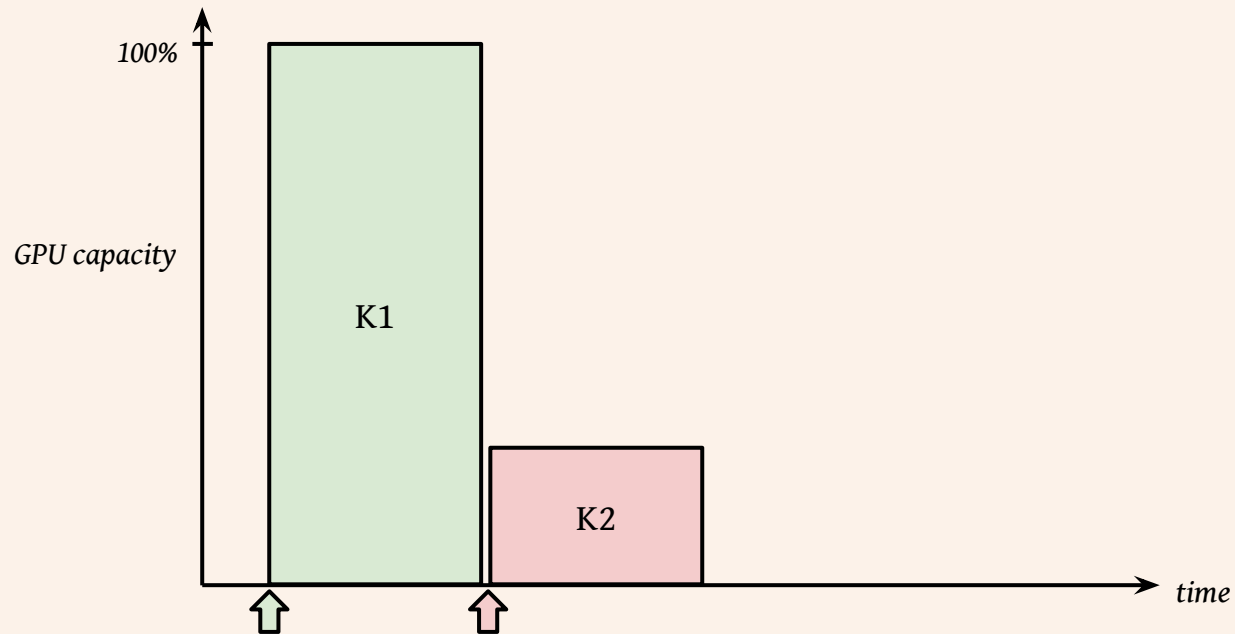
- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, using NVIDIA's MPS middleware.

Approaches to Co-Scheduling

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, using NVIDIA's MPS middleware.

Multi-process Co-Scheduling

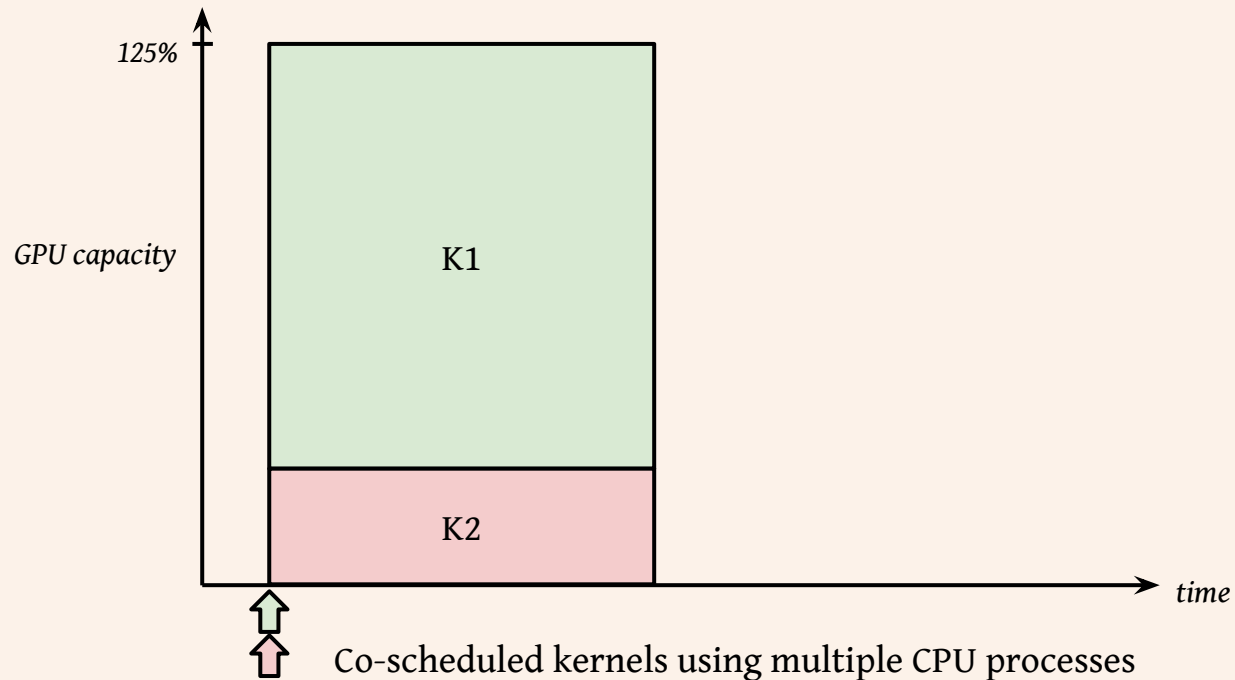


Starting Example: Sequential Execution

Reported in:

- *An Evaluation of the NVIDIA TX1 for Supporting Real-time Computer-Vision Workloads (RTAS'17)*
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed (RTSS'17)*

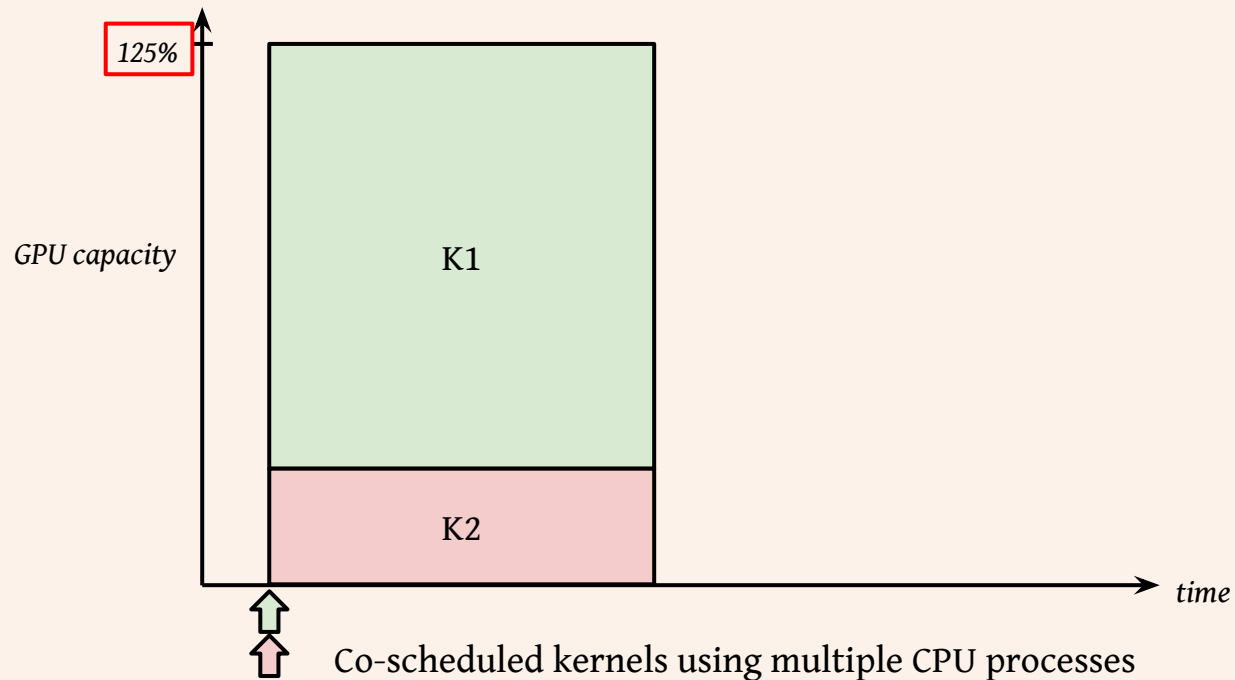
Multi-process Co-Scheduling



Reported in:

- *An Evaluation of the NVIDIA TX1 for Supporting Real-time Computer-Vision Workloads (RTAS'17)*
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed (RTSS'17)*

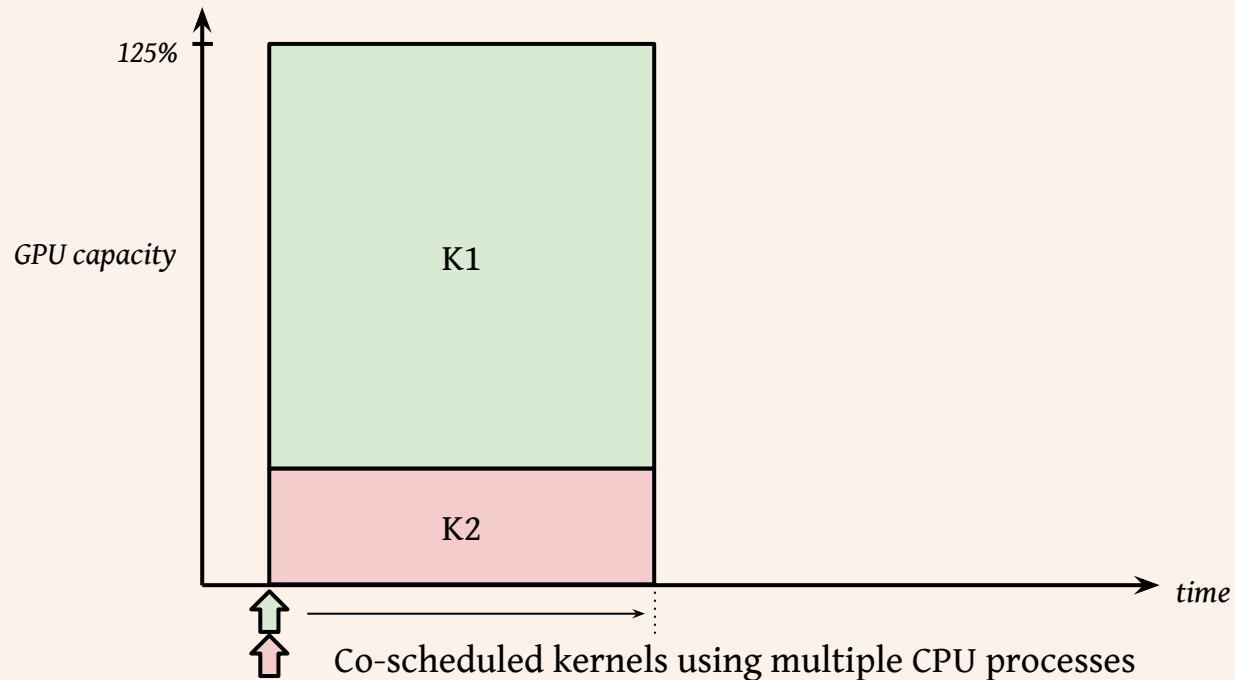
Multi-process Co-Scheduling



Reported in:

- *An Evaluation of the NVIDIA TX1 for Supporting Real-time Computer-Vision Workloads (RTAS'17)*
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed (RTSS'17)*

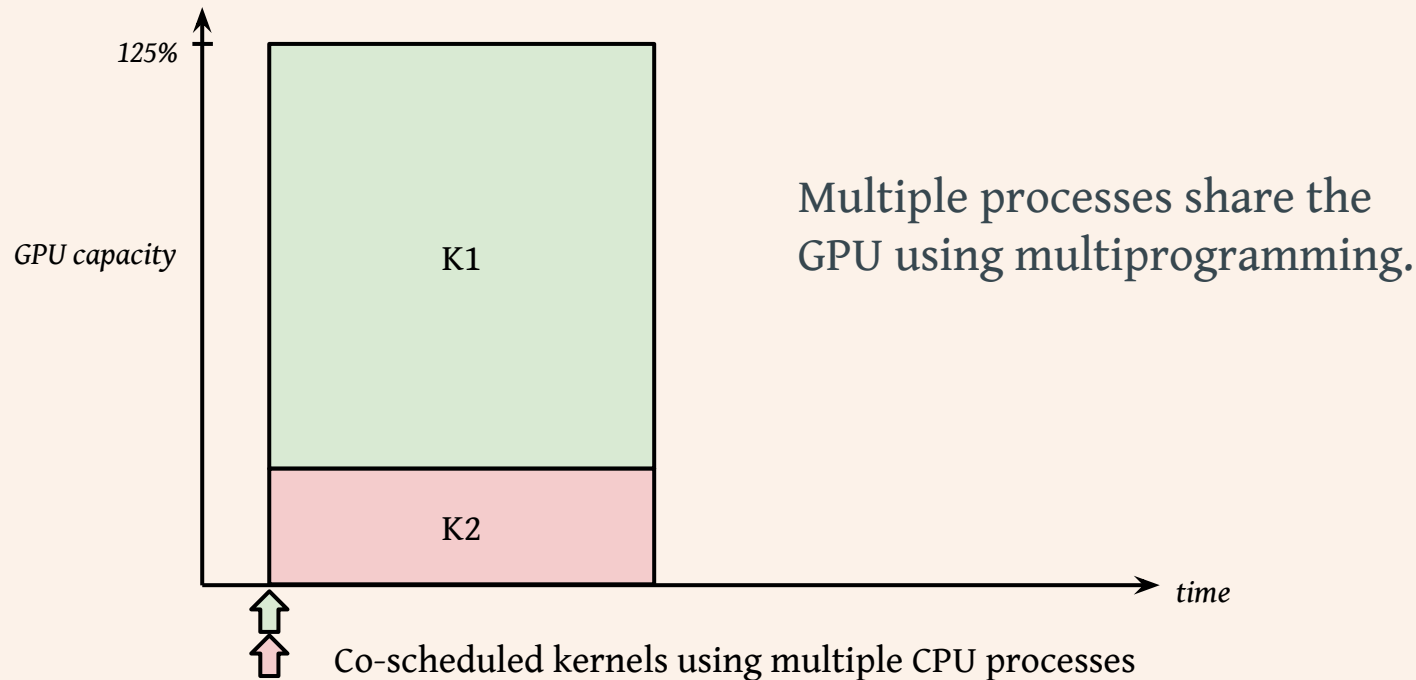
Multi-process Co-Scheduling



Reported in:

- *An Evaluation of the NVIDIA TX1 for Supporting Real-time Computer-Vision Workloads (RTAS'17)*
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed (RTSS'17)*

Multi-process Co-Scheduling



Reported in:

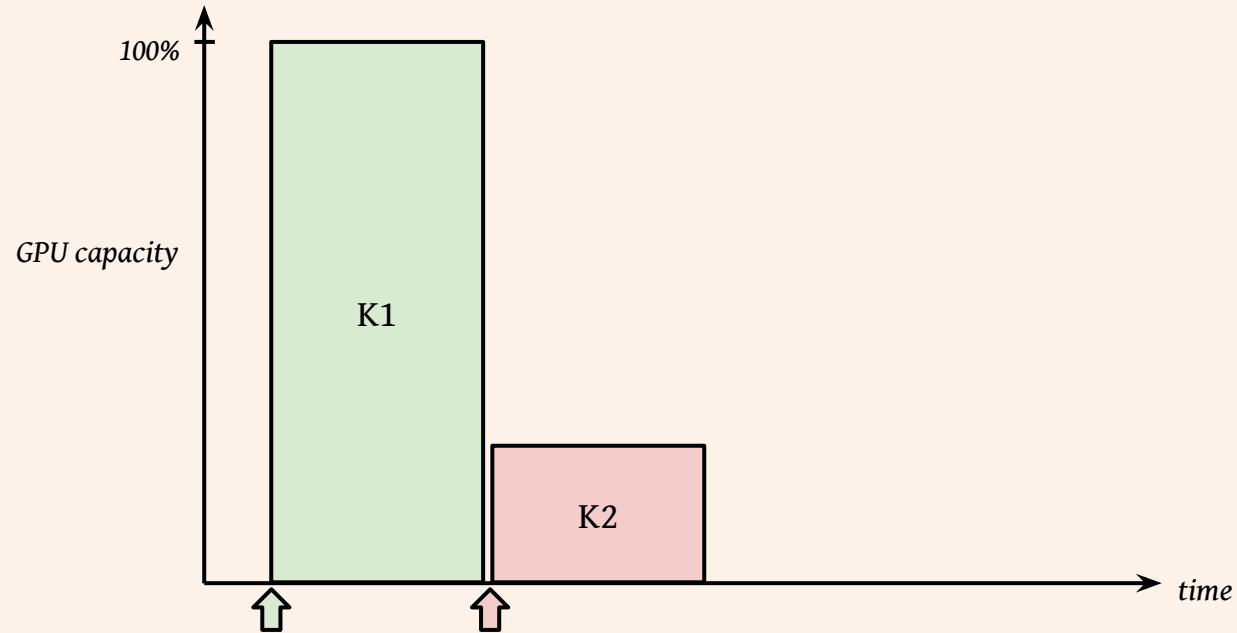
- *An Evaluation of the NVIDIA TX1 for Supporting Real-time Computer-Vision Workloads (RTAS'17)*
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed (RTSS'17)*

Approaches to Co-Scheduling

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, using NVIDIA's MPS middleware.

Multi-thread Co-Scheduling

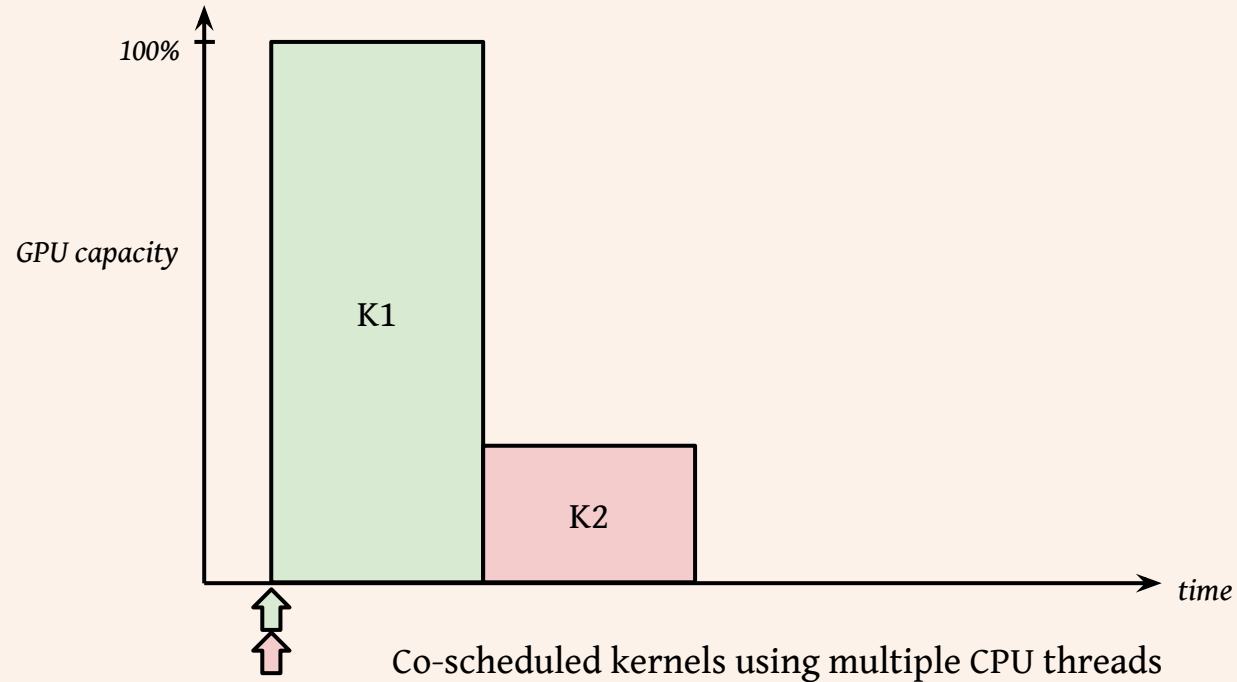


Starting Example: Sequential Execution

Reported in:

- *Inferring the Scheduling Policies of an Embedded CUDA GPU* (OSPRT '17)
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed* (RTSS'17)

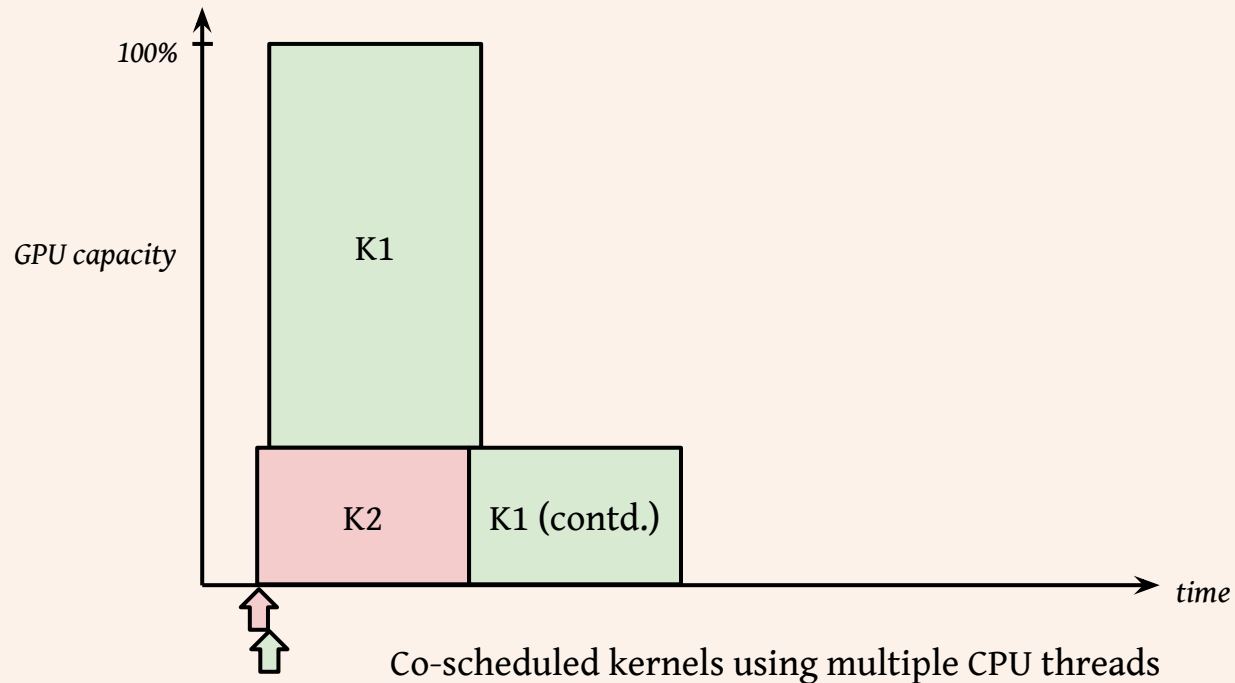
Multi-thread Co-Scheduling



Reported in:

- *Inferring the Scheduling Policies of an Embedded CUDA GPU* (OSPRT '17)
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed* (RTSS'17)

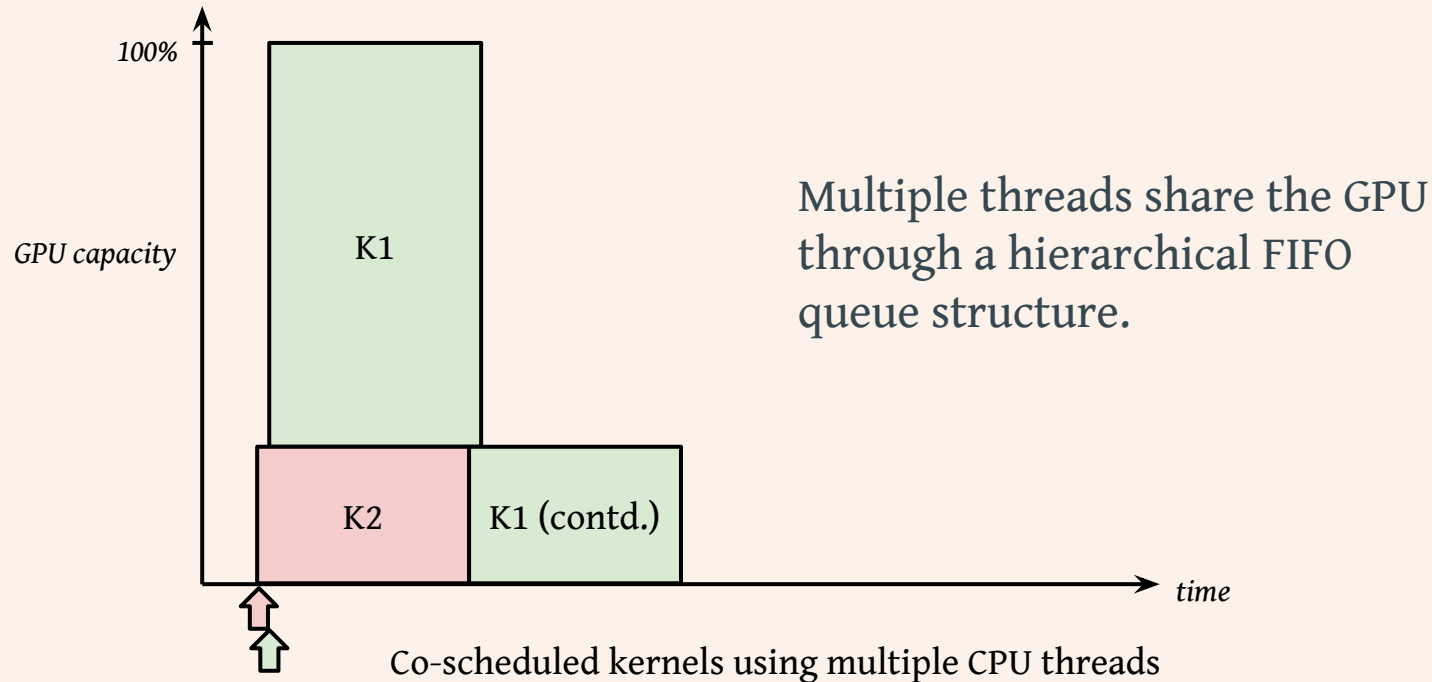
Multi-thread Co-Scheduling



Reported in:

- *Inferring the Scheduling Policies of an Embedded CUDA GPU* (OSPRT '17)
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed* (RTSS'17)

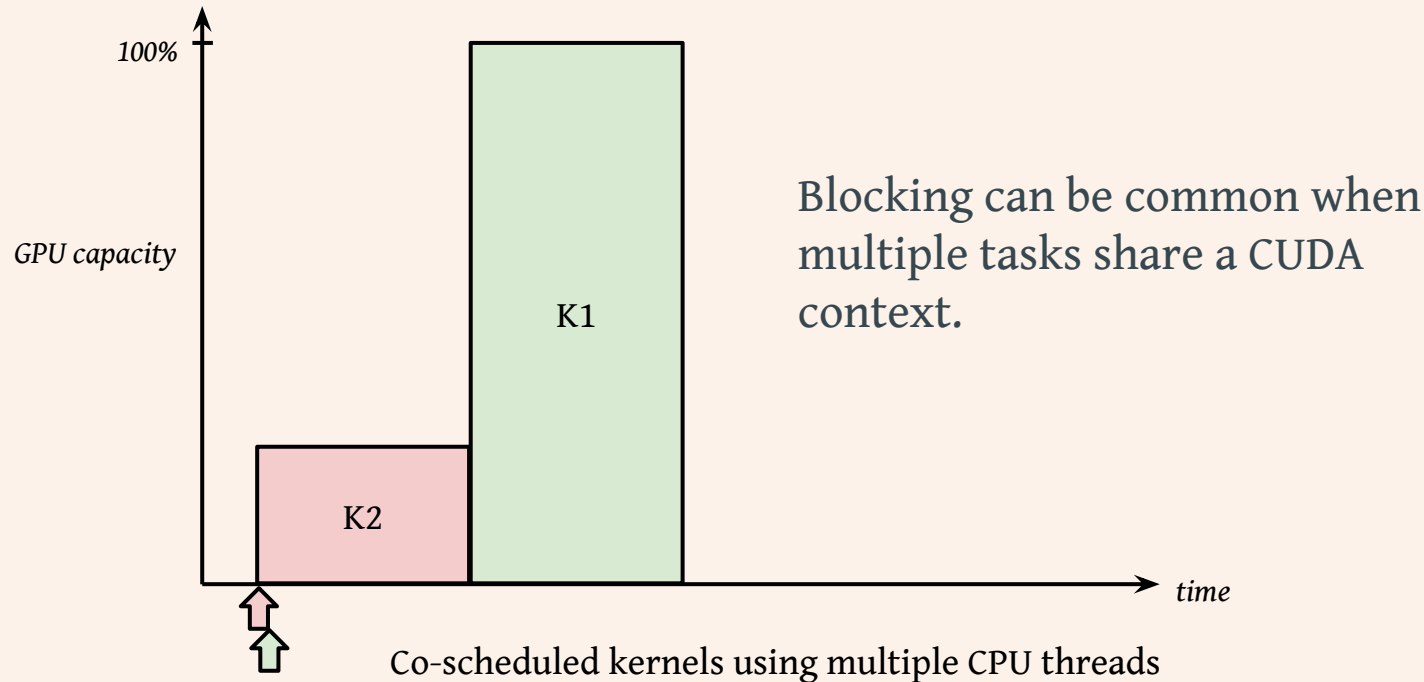
Multi-thread Co-Scheduling



Reported in:

- *Inferring the Scheduling Policies of an Embedded CUDA GPU* (OSPRT '17)
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed* (RTSS'17)

Multi-thread Co-Scheduling



Reported in:

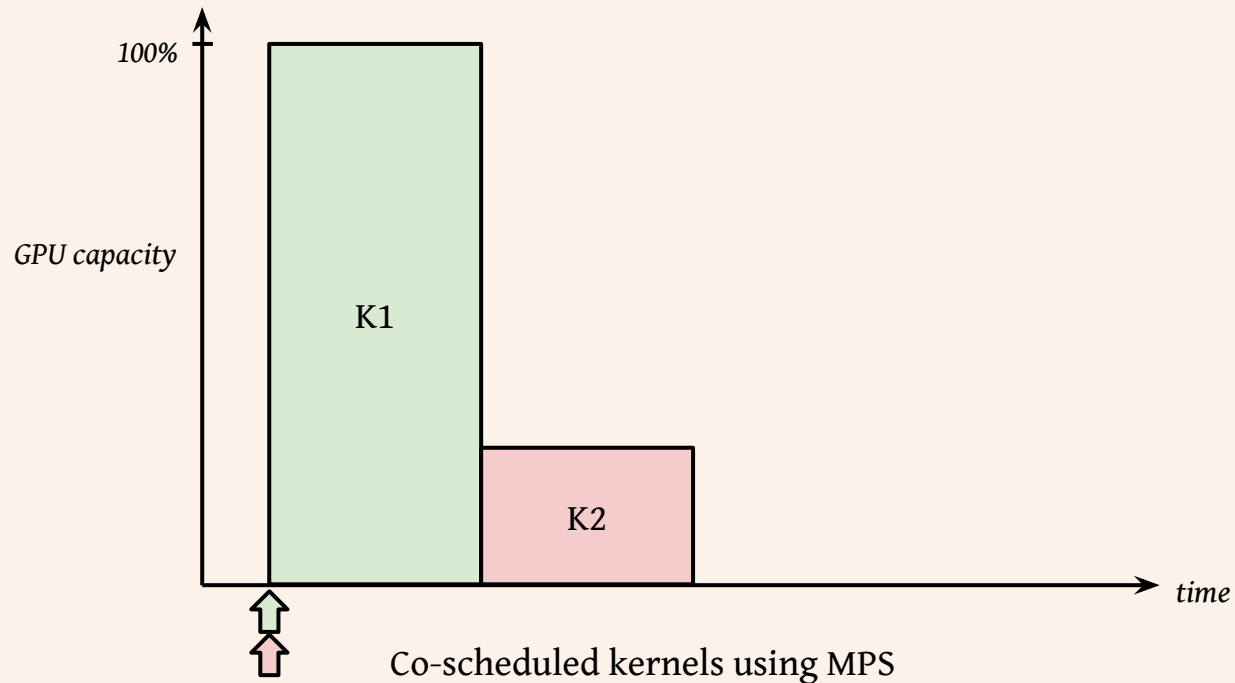
- *Inferring the Scheduling Policies of an Embedded CUDA GPU* (OSPRT '17)
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed* (RTSS'17)

Approaches to Co-Scheduling

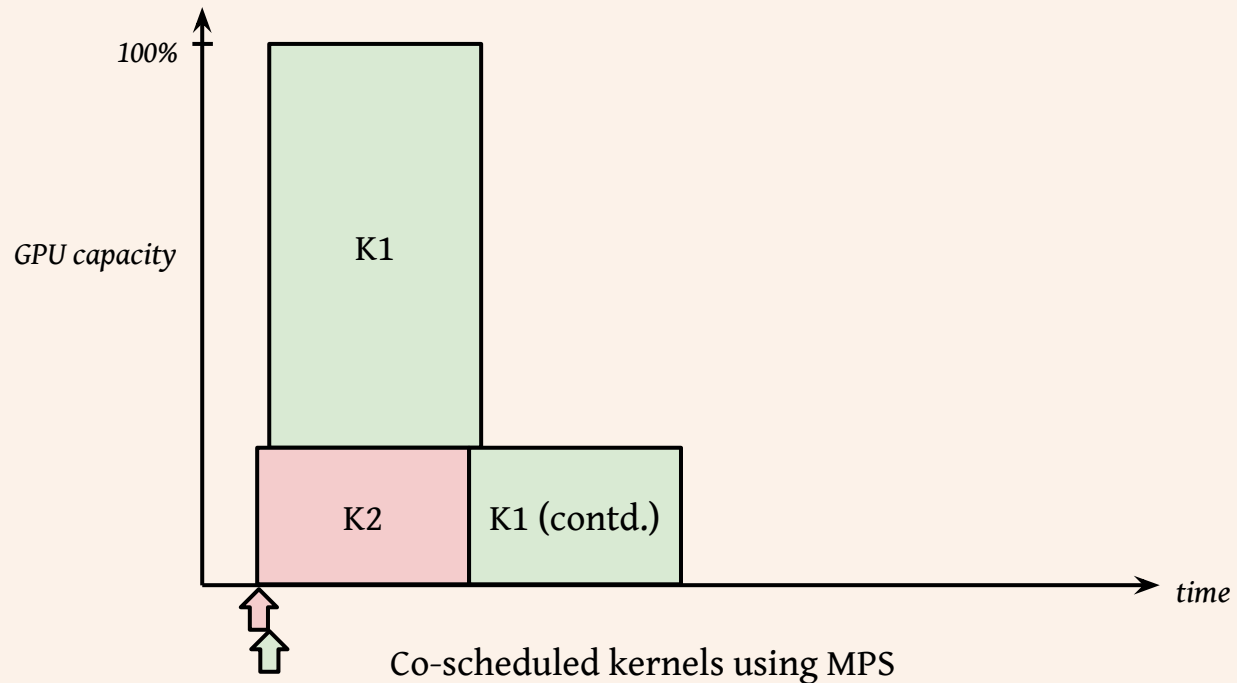
Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, using NVIDIA's MPS middleware.

Multi-process Co-Scheduling with MPS



Multi-process Co-Scheduling with MPS



Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware.

Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware.

The *best* system:

- Is easy to schedule (e.g. supports preemption).
- Does not require modifying existing CUDA programs.
- Fully utilizes the GPU when possible.

Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware.

The *best* system:

- Is easy to schedule (e.g. supports preemption). ✘
- Does not require modifying existing CUDA programs. ✘
- Fully utilizes the GPU when possible. ✘

Which one is best for safety-critical systems?

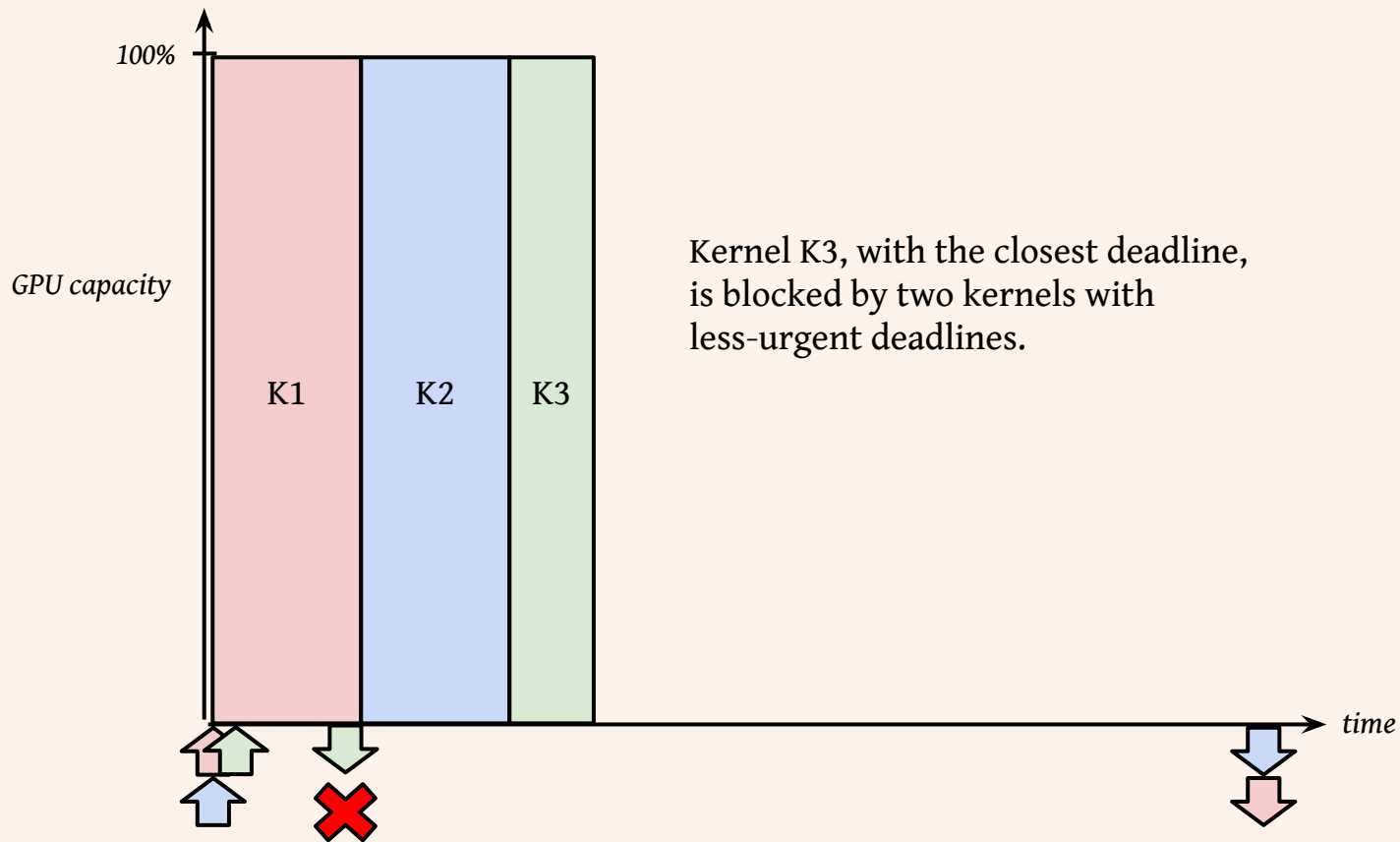
Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware.

The *best* system:

- Is easy to schedule (e.g. supports preemption). ❌
- Does not require modifying existing CUDA programs. ✔️
- Fully utilizes the GPU when possible. ✔️

Sporadic kernels under MPS (or threads)



Which one is best for safety-critical systems?

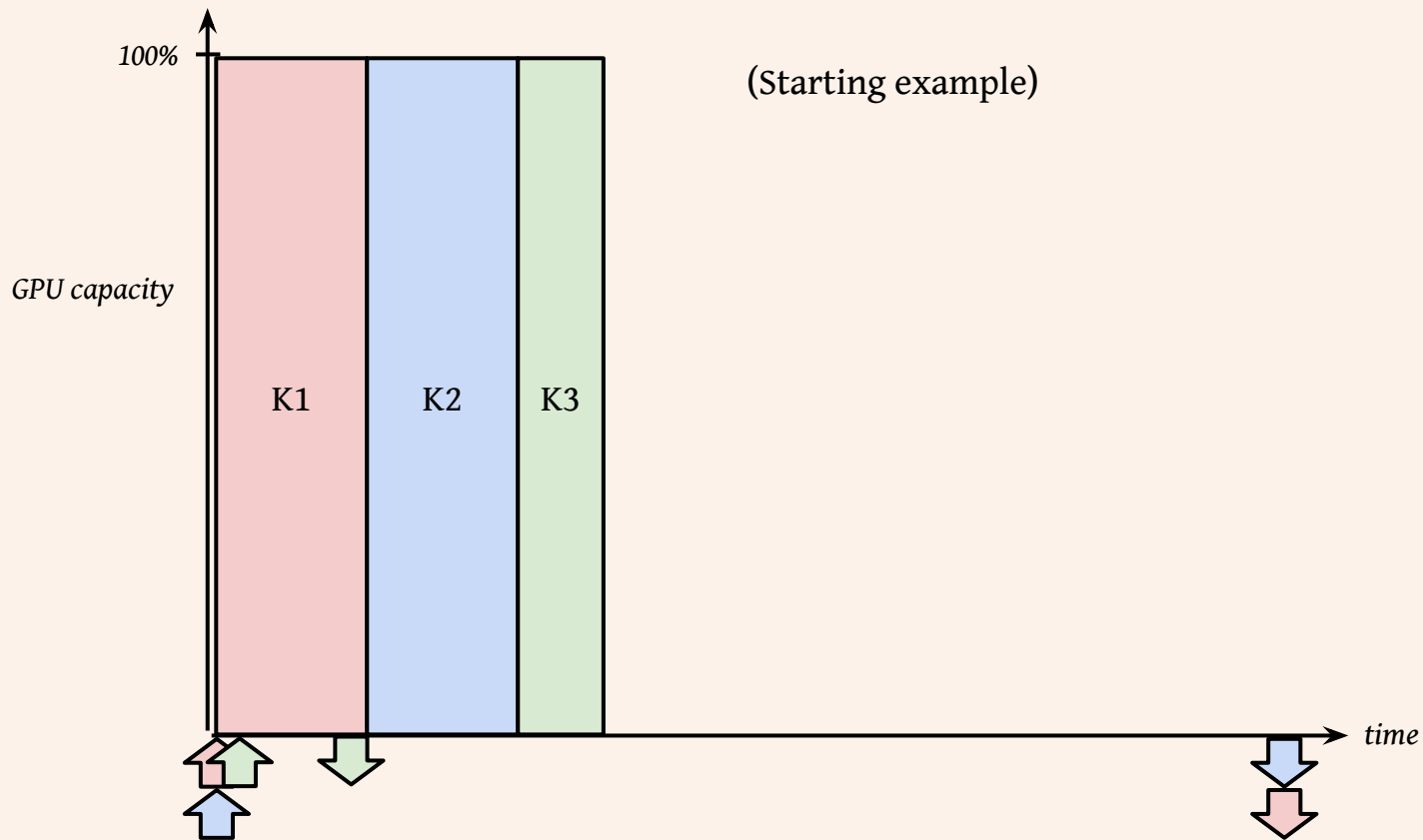
Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware.

The *best* system:

- Is easy to schedule (e.g. supports preemption). ?
- Does not require modifying existing CUDA programs. ✓
- Fully utilizes the GPU when possible. ✗

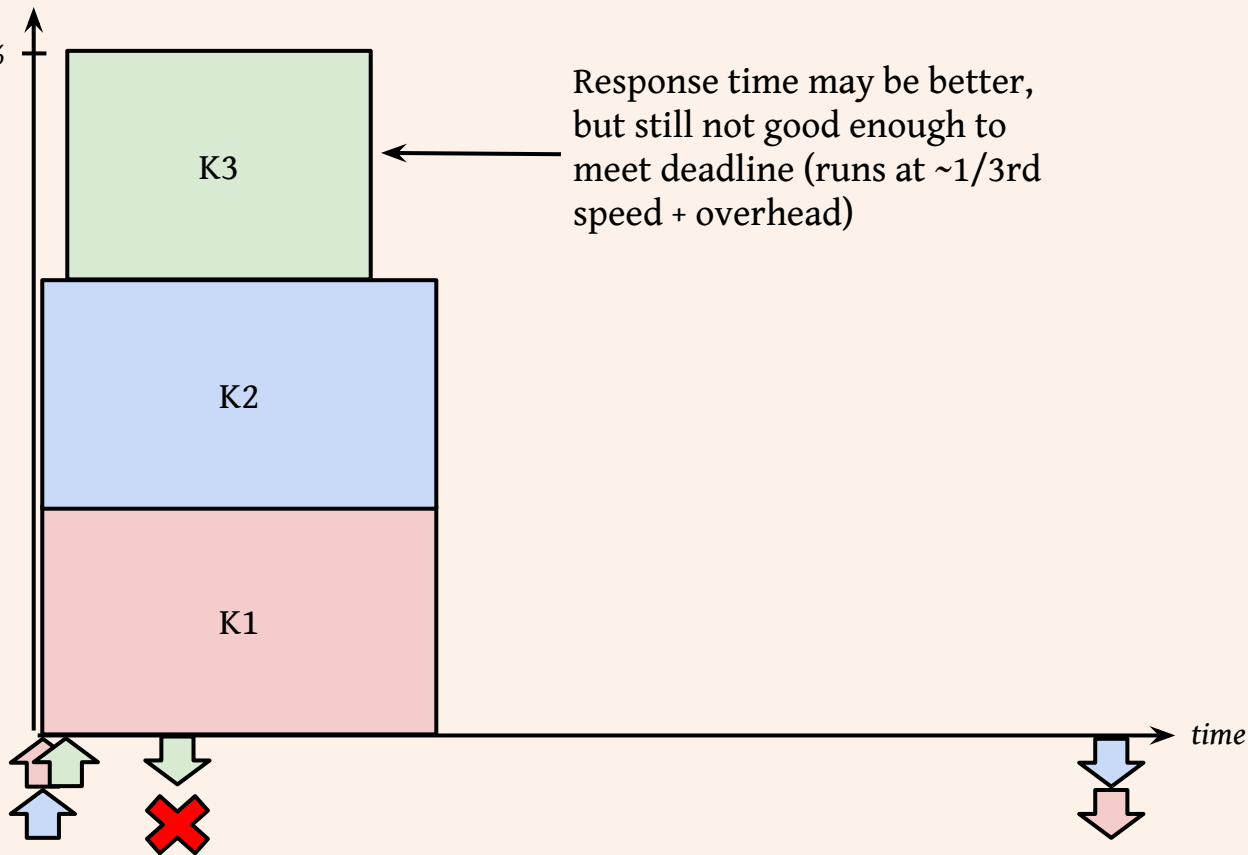
Sporadic kernels using Multiprogramming



Multiprogrammed kernels

Adjusted
"capacity": 300%

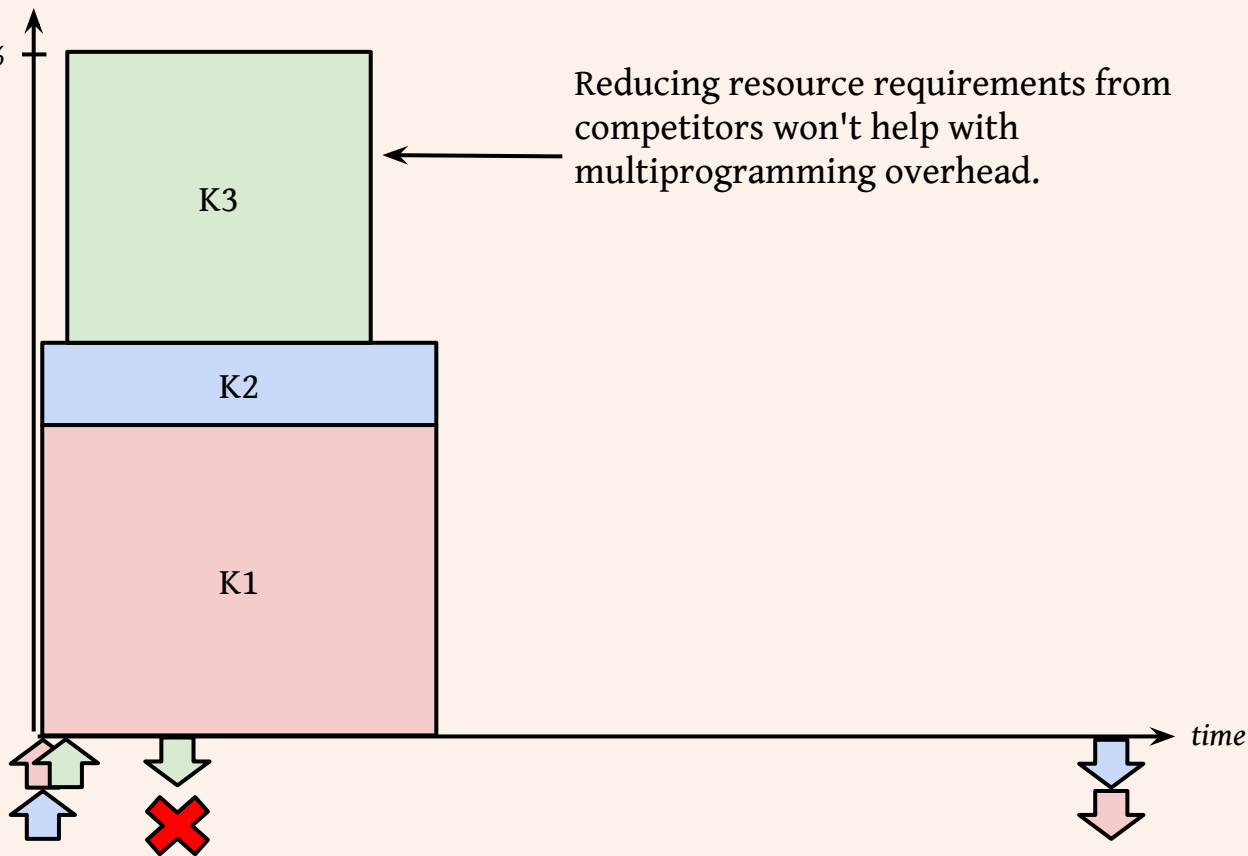
GPU capacity



Multiprogrammed kernels

Adjusted
"capacity": 225%

GPU capacity



Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware.

The *best* system:

- Is easy to schedule (e.g. supports preemption).
- Does not require modifying existing CUDA programs.
- Fully utilizes the GPU when possible.

Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware *on a Volta-architecture GPU to limit per-task computing resources.*

The *best* system:

- Is easy to schedule (e.g. supports preemption).
- Does not require modifying existing CUDA programs.
- Fully utilizes the GPU when possible.

Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware *on a Volta-architecture GPU to limit per-task computing resources.*

The *best* system:

- Is easy to schedule (e.g. supports preemption).
- Does not require modifying existing CUDA programs.
- Fully utilizes the GPU when possible.

Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware *on a Volta-architecture GPU to limit per-task computing resources.*

The *best* system:

- Is easy to schedule (e.g. supports preemption).
- Does not require modifying existing CUDA programs.
- Fully utilizes the GPU when possible.

Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware *on a Volta-architecture GPU to limit per-task computing resources.*

The *best* system:

- Is easy to schedule (e.g. supports preemption).
- Does not require modifying existing CUDA programs. ✓
- Fully utilizes the GPU when possible.

Which one is best for safety-critical systems?

Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware *on a Volta-architecture GPU to limit per-task computing resources.*

The *best* system:

- Is easy to schedule (*approximates fluid scheduling*). ✓
- Does not require modifying existing CUDA programs. ✓
- Fully utilizes the GPU when possible.

Which one is best for safety-critical systems?

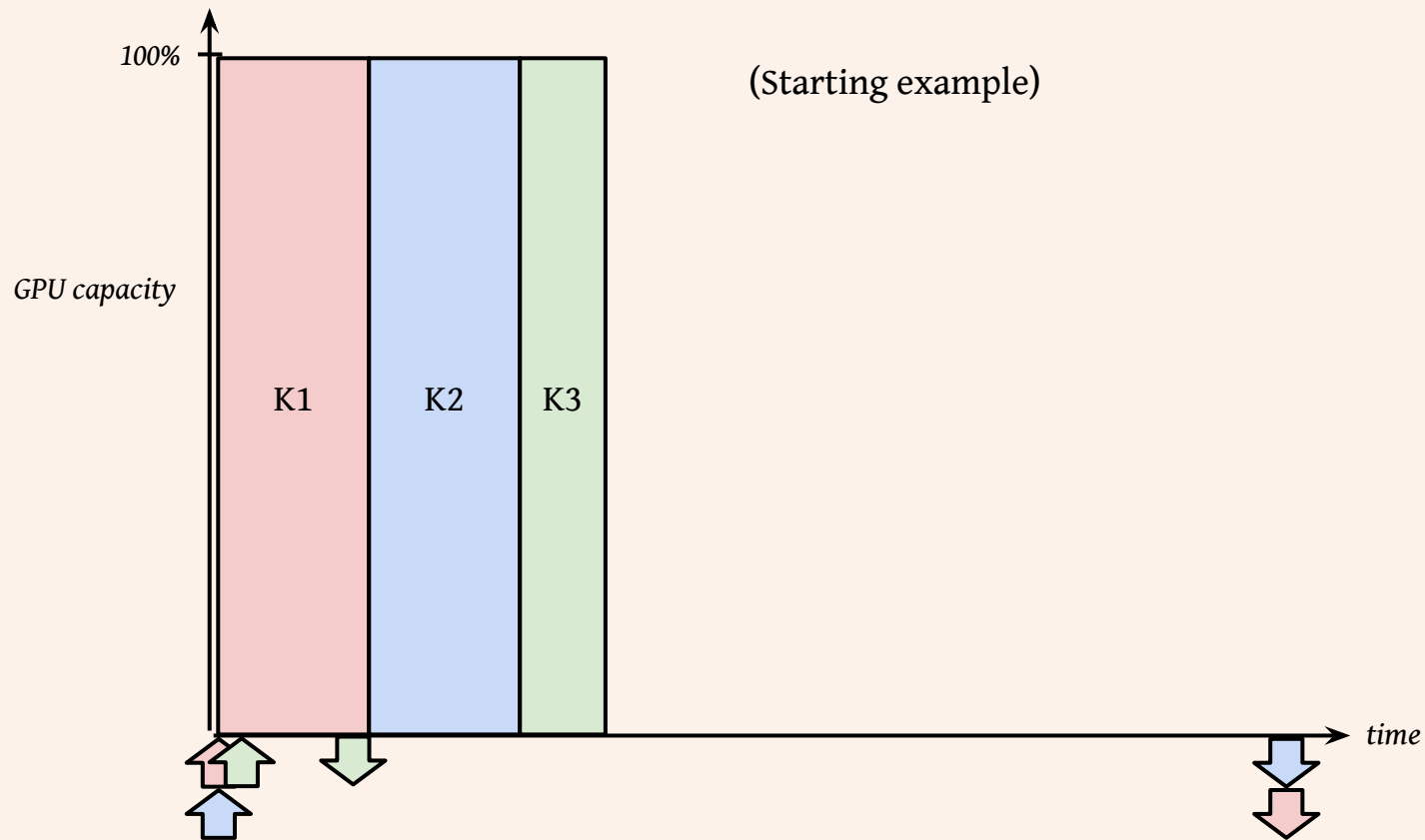
Different ways to Co-Schedule GPU Tasks:

- Every task is a separate CPU process (no GPU middleware).
- Every task is a separate CPU thread in a single process.
- Every task is a separate CPU process, but using NVIDIA's MPS middleware *on a Volta-architecture GPU to limit per-task computing resources.*

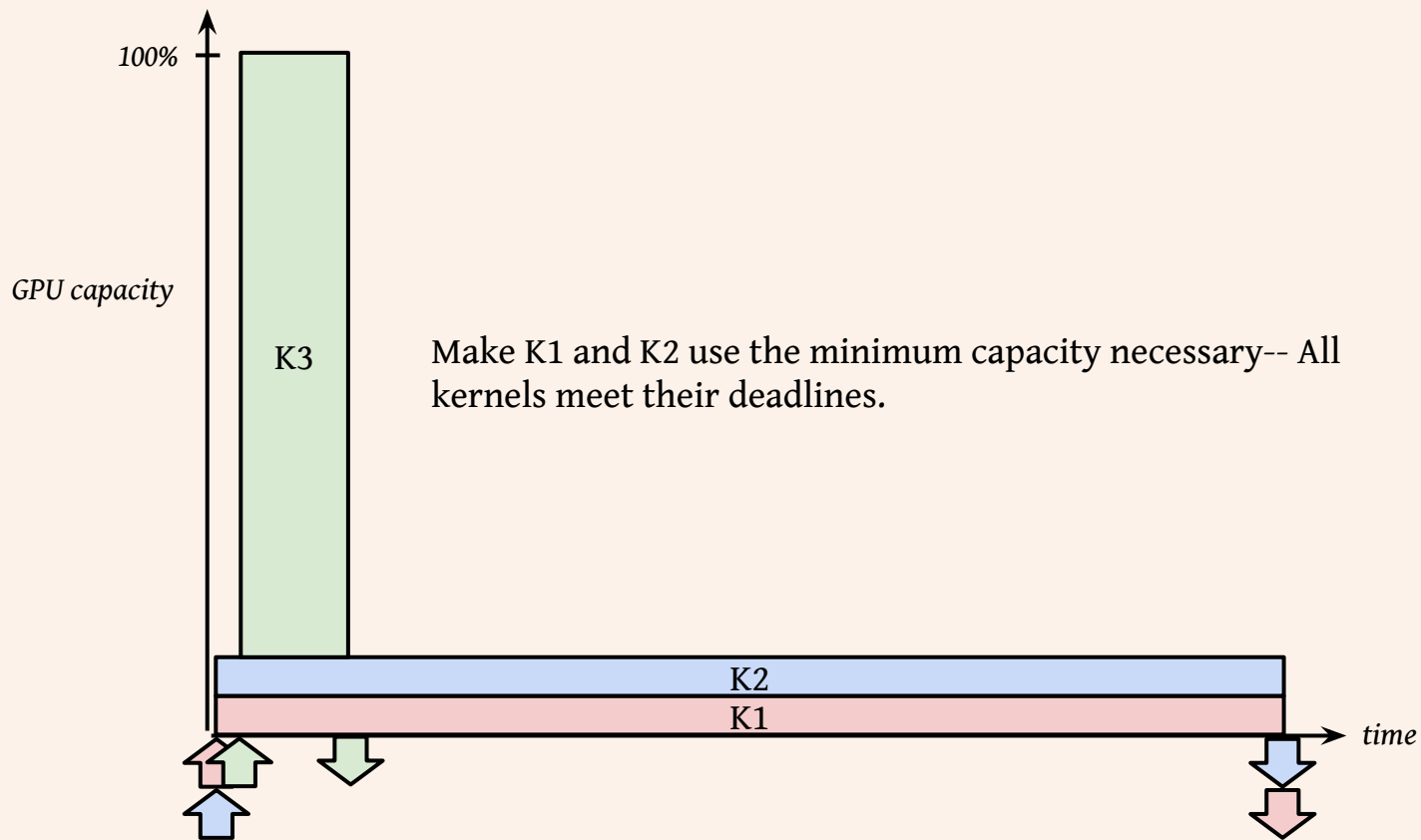
The *best* system:

- Is easy to schedule (*approximates fluid scheduling*). ✓
- Does not require modifying existing CUDA programs. ✓
- Fully utilizes the GPU when *necessary*. ✓

Applying LET to Volta MPS



Applying LET to Volta MPS



What needs to be done?

This is still ongoing work. The next steps include:

1. Determining formulas relating GPU utilization to execution time.
(This can actually be measured *per-task* rather than *per-kernel*.)
2. Write a management system that dynamically sets utilization limits based on the formulas and tasks' deadlines.

Potential Problems

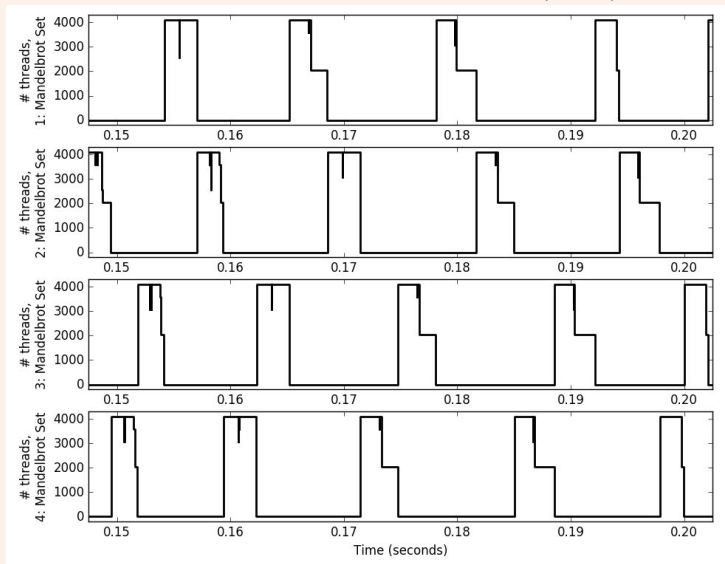
- Volta GPUs are currently expensive and in short supply.
- Embedded GPUs (so far) do not support MPS, regardless of GPU architecture.
- There's no guarantee that future GPU architectures will support setting resource limits.

Conclusion

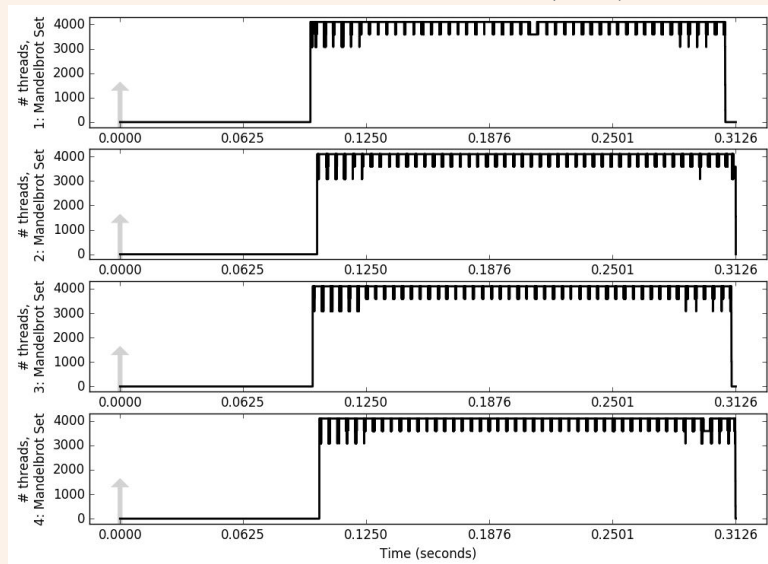
Safe, predictable real-time scheduling seems possible, when applying the principles of LET to GPU resource partitioning on Volta-architecture GPUs.

GPU Co-Scheduling with Processes

Maxwell GPU Architecture (TX1)



Pascal GPU Architecture (TX2)



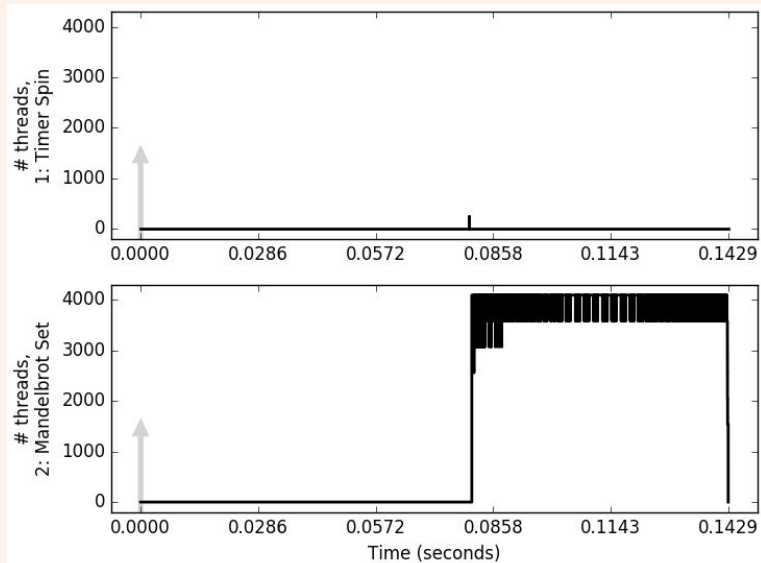
When GPU tasks are launched from separate CPU processes (CUDA contexts), co-scheduling is achieved via multiprogramming.

Reported in:

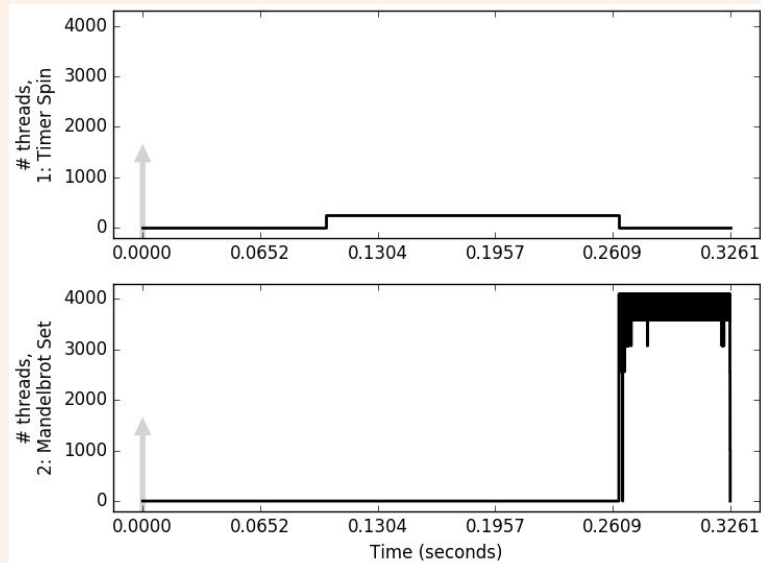
- *An Evaluation of the NVIDIA TX1 for Supporting Real-time Computer-Vision Workloads (RTAS'17)*
- *GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed (RTSS'17)*

GPU Co-Scheduling with Processes

Small, short-lived competing workload



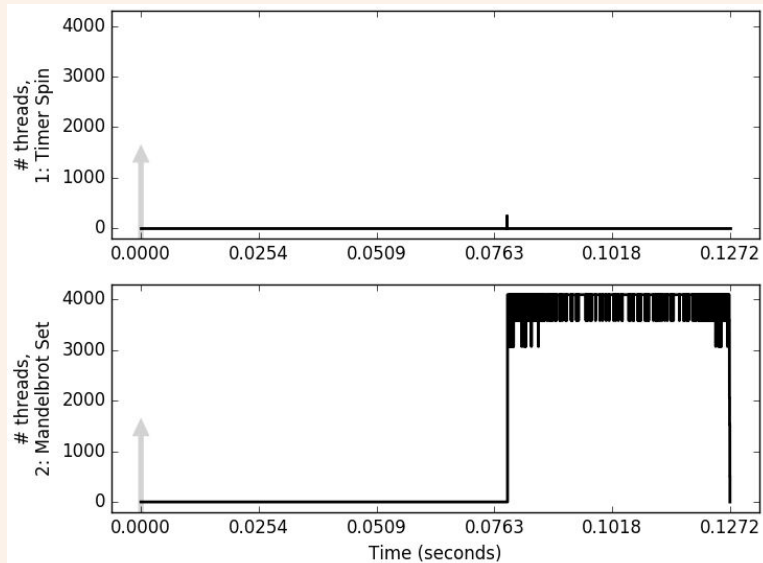
Small, long-lived competing workload



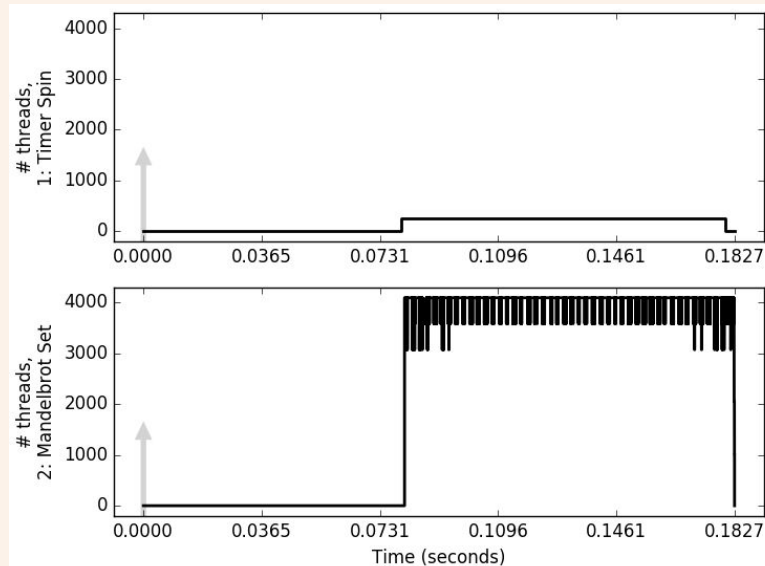
Multiprogramming on Maxwell GPUs leads to blocking.

GPU Co-Scheduling with Processes

Small, short-lived competing workload



Small, long-lived competing workload



Multiprogramming on Pascal GPUs leads to disproportionate performance loss.