Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems

Joshua Bakita 🖂 🕼

University of North Carolina at Chapel Hill, NC, USA

James H. Anderson 🖂 🗈 5

University of North Carolina at Chapel Hill, NC, USA 6

- Abstract

As GPU-using tasks become more common in embedded, safety-critical systems, efficiency demands 8 necessitate sharing a single GPU among multiple tasks. Unfortunately, existing ways to schedule multiple tasks onto a GPU often either result in a loss of ability to meet deadlines, or a loss of 10 11 efficiency. In this work, we develop a system-level spatial compute partitioning mechanism for NVIDIA GPUs and demonstrate that it can be used to execute tasks efficiently without compromising 12 timing predictability. Our tool, called nvtaskset, supports composable systems by not requiring 13 task, driver, or hardware modifications. In our evaluation, we demonstrate sub-1- μs overheads, 14 stronger partition enforcement, and finer-granularity partitioning when using our mechanism instead 15 of NVIDIA's Multi-Process Service (MPS) or Multi-instance GPU (MiG) features. 16

2012 ACM Subject Classification Computer systems organization \rightarrow Heterogeneous (hybrid) systems; 17 Computer systems organization \rightarrow Real-time systems; Software and its engineering \rightarrow Schedul-18

ing; Software and its engineering \rightarrow Concurrency control; Computing methodologies \rightarrow Graphics 19

processors; Computing methodologies \rightarrow Concurrent computing methodologies 20

Keywords and phrases Real-time systems, composable systems, graphics processing units, CUDA 21

Digital Object Identifier 10.4230/LIPIcs.ECRTS.2025.18 22

Funding Work supported by NSF grants CPS 2038960, CPS 2038855, CNS 2151829, CPS 2333120, 23 and ONR contract N0001424C1127. 24

1 Introduction 25

Rapid developments in artificial intelligence (AI)—especially deep neural networks (DNNs) 26 running on GPUs [18]—have led to new cyber-physical systems, from intelligent assistants 27 to self-driving cars. Real-world safety or usability concerns impose practical response-time 28 deadlines on these systems, which may also need to run multiple AI tasks—such as one 29 DNN for a conversational interface alongside others for object detection or planning in a 30 self-driving car. However, this raises a problem-how to schedule GPU-using tasks onto a 31 GPU efficiently while reliably meeting deadlines? When scheduling a GPU, generally either 32 competitive sharing [42, 27, 46]—tasks run concurrently and fight for resources—or mutual 33 exclusion [12, 11, 3]—one task runs at a time—are recommended. 34

Unfortunately, competitive sharing increases efficiency at the cost of timing predictability 35 [11, 35, 2, 41, 7, 40], whereas mutual exclusion gives up efficiency for predictability. Without 36 timing predictability, one cannot guarantee met deadlines. This tension puts embedded 37 system designers in a difficult position. The easy option—trading off predictability for 38 efficiency—is dangerous for safety-critical systems like self-driving cars. 39

This problem is exacerbated by the issue of composability. Each GPU-using task may 40 be developed by different groups, may not be modifiable by the scheduler, and may change 41 out-of-step with other tasks during a device's lifetime. This puts further burden on the 42 scheduling system, as it must guarantee efficient and predictable execution for each task, 43 even as tasks change opaquely. 44



licensed under Creative Commons License CC-BY 4.0 37th Euromicro Conference on Real-Time Systems (ECRTS 2025). Editor: Renato Mancuso; Article No. 18; pp. 18:1-18:24

© Joshua Bakita and James H. Anderson:

Leibniz International Proceedings in Informatics

LIPICS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

18:2 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems

Mechanism	D OTZAL,	Lostica, 1	Treaties	Low. Cont. C	Harder dead	Dyname Ent.		(B)
Software [13, 39, 15, 37, 43]	1	×	Х	Х	×	\checkmark	1	
libsmctrl [4]	1	×	×	\checkmark	1	\checkmark	\checkmark	
NVIDIA MiG [29]	X	\checkmark	1	1	1	×	×	
NVIDIA MPS [27]	1	\sim	\sim	1	\sim	×	1	
nvtaskset (ours)	1	\sim	1	1	1	1	1	

Table 1 Comparison of spatial GPU compute partitioning mechanisms.

In this work, we demonstrate that spatial partitioning of GPU compute cores is an 45 effective path to resolving this problem. To show this, we uncover and repurpose hardware 46 capabilities in all NVIDIA GPUs to build a new system-level spatial partitioning mechanism. 47 Spatial partitioning—a way to run tasks concurrently on mutually exclusive sets of cores-48 allows concurrent task execution for efficiency, while minimizing shared-resource interference 49 between tasks to protect timing predictability. Our work builds on two key insights: GPUs are 50 architecturally well-suited to spatial partitioning, and all NVIDIA GPUs contain hardware 51 capabilities that can be leveraged to implement spatial partitioning. 52

While our work is motivated by the computational needs of DNNs, it is *not* constrained to DNNs—we consider arbitrary *unmodified* CUDA-using GPU tasks. We focus on NVIDIA GPUs for their market-leading technology and adoption, and on CUDA-using tasks because of an intermediate-software limitation—we believe that the hardware capabilities we unveil and leverage could be applied to spatially partition any NVIDIA GPU workload in the future.

Prior work. Prior work on the spatial partitioning of NVIDIA GPU cores is limited. Table 1 58 (returned to with rigorous definitions in Sec. 2.5 and Sec. 3) classifies key works. Most prior 59 work ("Software" in Table 1) modifies tasks to cooperatively yield unallocated compute 60 cores [39, 15, 37, 43, 13, 47]—these approaches suffer the inability to enforce partitions on 61 misbehaving tasks. Our prior work, libsmctrl [4] addresses this enforcement problem, but 62 requires task modification and a shared address space among all tasks. Multi-instance GPU 63 (MiG) from NVIDIA [29] can hardware-enforce partitions without these compromises, but 64 its partitions are static, less fine-grained, and only available on data-center GPUs. The only 65 NVIDIA-provided option for all their GPUs, the Execution Resource Provisioning feature of 66 the Multi-Process Service (MPS), does not enforce robust partition boundaries. We revisit 67 these mechanisms at length in Sec. 3. 68

⁶⁹ Contributions. In this work, we:

- Reveal how NVIDIA MPS (current state-of-the-art) is implemented, unveiling previouslyunknown hardware capabilities and quantifying its real-time safety.
- Discover and mitigate efficiency and predictability pitfalls of NVIDIA MPS, including
 how it may assign two partitions of 50% to the same 50% of the GPU.
- 74 3. Develop a new system-level spatial-partitioning tool for NVIDIA GPUs—nvtaskset—
 75 built on principles from libsmctrl and NVIDIA MPS.
- 4. Show that nvtaskset supports unmodified tasks without measurable overheads or port ability limitations.

78 5. Demonstrate that nvtaskset has more efficient and granular partition enforcement than

⁷⁹ NVIDIA MiG and MPS, and equivalent or better timing predictability.

 $_{\rm 80}$ 6. Uncover that NVIDIA MiG underperforms due to inherent overheads of 6–15%.

Organization. We introduce our system model, overview GPU architecture, and discuss spatial partitioning in Sec. 2. We review prior work in Sec. 3. In Sec. 4, we elucidate the implementation and pitfalls of NVIDIA MPS, and then in Sec. 5 apply the hardware capabilities used by MPS to develop nvtaskset. We evaluate the overheads, partition enforcement, and granularity of nvtaskset in Sec. 6, and conclude in Sec. 7.

86 2 Background

In this section, we summarize necessary background on the GPU (from our prior works [4, 5]),
and discuss why spatial partitioning enables efficiency and timing predictability.

89 2.1 System Model

We assume an x86_64 or aarch64 platform containing at least one embedded or discrete
 NVIDIA GPU of the Volta (2018) generation or newer.

Tasks are assumed to be closed-source, unmodifiable, CUDA-using Linux processes. Non-CUDA GPU-using tasks may coexist, but may not use spatial partitioning.

We focus on embedded, real-time systems, where resources are limited but execution-time deadlines must be met.

96 2.2 GPU Usage Model

⁹⁷ Each task executing on a GPU has an associated GPU *context*. This context includes
⁹⁸ per-GPU-task state, such as an on-GPU virtual address space. Unless explicitly stated
⁹⁹ otherwise, we assume a one-to-one mapping of GPU-context to CPU-task.

Work is dispatched into a context via one or more compute kernels. A kernel is launched 100 by passing a GPU-executable binary, and a number of GPU threads, to a GPU-usage library 101 such as CUDA or OpenCL. The library will then compose and pass a Task Metadata 102 Descriptor (TMD) to the GPU for execution. Each GPU thread concurrently executes the 103 same instructions, but operates at a different data offset. For example, a kernel for an 104 element-wise array addition could replace a for loop over every index value with a GPU 105 thread per value. The threads would then execute the loop body over all index values in 106 107 parallel, rather than having to iterate.

Threads are organized into groups known as *blocks*, and all threads within a single block execute on the same SM (Streaming Multiprocessor; a group of compute cores). A series of kernel executions may be serialized via first-in-first-out command queues known as *streams* in CUDA. Only kernel executions in the same stream are serialized; kernels in separate streams are permitted to execute concurrently.

113 2.3 GPU Architecture

In Fig. 1 we illustrate the typical architecture of an NVIDIA GPU chip, using the Adageneration AD102 as an exemplar (used in the RTX 4090 and RTX 6000 Ada GPU models).
NVIDIA GPUs are subdivided into several independent *engines*, with the most significant
being the Compute/Graphics Engine. Smaller engines handle special functions such as bulk

18:4 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems



Figure 2 Number of CUDA cores in NVIDIA's top-end GPUs; 2011 to 2023. (Select points annotated with chip IDs.)

data movement (the Copy Engines) and video processing. Each engine may have a different context active on it at a time [5].

The Compute/Graphics Engine is subdivided into General Processing Clusters (GPCs). Each GPC is independently connected to each DRAM controller and the co-located last-level, level-two (L2) cache slices. The GPCs subdivide into Thread Processing Clusters (TPCs) of two Streaming Multiprocessors (SMs) each. Each SM contains dozens of compute cores (128 each on the AD102) and a level-one (L1) cache. The AD102 contains 18,432 compute cores total; Fig. 2 illustrates how core counts have increased in recent years.

¹²⁶ We next discuss how GPU contexts are scheduled onto hardware engines.

127 2.4 GPU Scheduling

We present the scheduling pipeline from CPU task to GPU compute cores in Fig. 3. CPU tasks insert kernel launch commands (represented as TMDs) into CUDA streams, and those CUDA streams are mapped onto a smaller or equivalent number of GPU *channels*.¹ Which GPU channel(s) may access the Compute/Graphics Engine at a time is dictated by the *runlist*, making the runlist the central arbitrator.

GPU runlists (typically one for each GPU engine [5]) are composed of *time-slice groups*

¹ Strictly, channels contain a queue called a *pushbuffer*. The channel count is important—using more streams than channels causes blocking [5].

On-CPU	Streams	Channels	Runlist	On-GPU
GPU-Using Task 1	Ū		J.	
GPU-Using Task 2	Ū			

Figure 3 The scheduling pipeline for two tasks using the GPU Compute/Graphics Engine—the runlist arbitrates which task uses the engine at a time. (Dashed boxes at right represent intra-engine scheduling stages we skip—see [4].)

	Task 1's Runlist Entries; Time-Slice Group 0			Task 2's Runlist Entries; Time-Slice Group 1			Task 3's Runlist Entries; Time-Slice Group 2			
	Header	Channel 0	Channel 1	Header	Channel 2	Channel 3	Header	Channel 4	Channel 5	
	Timeslice: 2ms	[queue state	[queue state	Timeslice: 2ms	[queue state	[queue state	Timeslice: 2ms	[queue state	[queue state	
	Channels: 2	pointer]	pointer]	Channels: 2	pointer]	pointer]	Channels: 2	pointer]	pointer]	
Ċ		3	2	6	4 Offset (bytes)	9	6	1:	28	

Figure 4 Example Compute/Graphics Engine runlist of three tasks.

(TSGs) of channels. TSGs are executed in a work-conserving, preemptive round-robin order 134 by default. While a TSG is active, commands from all its contained channels are pulled 135 and executed on the runlist-associated engine. Each TSG has an associated time-slice; upon 136 expiration of this timeslice the runlist scheduler preempts any in-progress commands and 137 switches to the next TSG. Such a switch is also triggered by the exhaustion of commands 138 from all the TSG's channels. Note that all a TSG's channels must be from the same context.² 139 We show a runlist for three tasks in Fig. 4, following the in-memory layout. In this 140 example, TSG 0 would be executed for 2 ms, and during that time, commands from channels 141 0 and 1 would be received and executed by the Compute/Graphics Engine. After 2 ms, any 142 active commands from channels 0 and 1 would be preempted, and the GPU would move on 143 to commands from the channels of TSG 1 for 2 ms. This process repeats for TSG 2, then 144 loops back to TSG 0. The time-slice can be set differently for each TSG. 145



Figure 5 GPU utilization over one inference of the YoloV2 DNN in Darknet on the 4,352-core NVIDIA RTX 2080 Ti.

This GPU scheduling algorithm ensures that the Compute/Graphics Engine has only one context active at a time. This can lead to significant under-utilization, as many GPU-using tasks are unable to saturate all GPU compute cores on their own [42, 27, 46]. We demonstrate this for one inference of the YOLOv2 image-detection network in Fig. 5. At no point is the network able to utilize the entire GPU—it utilizes only 40% of the GPU's SMs on average.

² See line 293 in manuals/ampere/ga100/dev_ram.ref.txt [30].

18:6 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems

¹⁵¹ Increasing GPU core counts (Fig. 2) have only worsened this problem.

¹⁵² To efficiently use the GPU, idle compute capacity must be reclaimed for other tasks.

¹⁵³ Unfortunately, concurrently running multiple tasks on the GPU (e.g., by sharing a single ¹⁵⁴ context) leads to a new set of problems.

154 context) leads to a new set of problems.

155 2.5 Interference and Spatial Partitioning

When multiple tasks execute concurrently on a GPU, shared-resource interference can occur.
This is where contention for shared hardware resources such as caches or compute cores
creates slowdowns between tasks. Such interference creates unpredictability, as knowing
when and how tasks will interfere on what shared resource is an unresolved area of study.
Thus, to ensure timing predictability, few hardware resources should be shared.

Prior work has avoided such interference on the GPU by only allowing one task to access
the GPU at a time [12, 11, 3]. This can lead to underutilization—concurrently running
multiple tasks would allow for more efficient use of GPU hardware.

Spatial partitioning allows for concurrency without interference. It prevents hardware
 resources from being shared between tasks by partitioning them into mutually exclusive
 subsets, and assigning each subset to a concurrently running task. Without shared resources,
 interference is prevented, and execution times remain predictable.

In designing a spatial partitioning mechanism, it should satisfy the following properties:
 Portable: The mechanism should work on a wide set of GPU models.

Logically Isolated: The mechanism should preserve logical isolation between tasks, *e.g.*,
 virtual address space isolation and independent exception handling.

172 Transparent: The mechanism should not require changes to tasks, e.g., recompilation.

Low-overhead: The mechanism should add negligible overhead to critical-path operations,
 e.g., a kernel launch.

Hardware-enforced: Partitions should be enforced by hardware to protect from malicious or
 misbehaving tasks.

177 Dynamic: Partitions should be reconfigurable without restarting a task.

Granular: Partitions should be defined in granular units, *e.g.*, a TPC for compute partitioning.
The first four properties are desirable for any software system, but the last three are

¹⁸⁰ partitioning-specific. We now discuss prior work in light of these requirements.

181 3 Related Work

Efficiently using a GPU in a real-time system requires concurrently running tasks on the GPU Compute/Graphics Engine, with partitioning of shared hardware resources. Prior work has identified the GPU DRAMs (with co-located L2 cache slices) and SMs (with co-located L1 caches) as needing to be partitioned [40, 4, 15]. This section covers prior work on and towards such partitioning.

187 3.1 DRAM Partitioning

The history of DRAM partitioning is extensive ([19, 20, 44] are exemplars), and it has been extended onto the GPU by Fractional GPUs [15] and SGDRC [45]. These works use a memory-organization approach known as *page-coloring* [19] to force each task onto prescribed GPU DRAM and L2 units. Such memory reorganization requires difficult-to-obtain modelspecific hashing functions, but the principle is generally applicable to any GPU. NVIDIA

has since developed a proprietary alternative, but this is not available on most GPUs and
 cannot be enabled independently of MiG [29, 9].

The DRAM partitioning technique of Fractional GPUs and SGDRC satisfies all desired spatial isolation properties stated in Sec. 2.5, so we assume the application of such an approach and focus on the remaining problem: compute partitioning.

3.2 Compute Partitioning

Prior work on spatial partitioning for GPU compute cores can be divided into academic provided and NVIDIA-provided solutions. We summarize these works in Table 1.

Academic-provided solutions. Third-party solutions for spatial partitioning on NVIDIA 201 GPUs have been limited to cooperation-based software mechanisms until recently. One 202 system [39], which has been recently improved [15, 37], is commonly used in papers that 203 claim to partition NVIDIA GPUs. This approach depends on kernels launching blocks on all 204 SMs, and on each block aborting if it finds itself executing on an SM outside of its partition. 205 This approach is vulnerable to a full loss of partitioning if any block misbehaves and does not 206 cooperatively yield an unassigned SM. Another cooperation-based variant known as persistent 207 threads [43, 13], still used in recent work [47], is similarly vulnerable to misbehaving tasks. 208 Given the inability of these works to *enforce* partitioning, they are vulnerable to a loss of 209 isolation. We group these works under the "Software" heading in Table 1, as they implement 210 partitioning via task modification rather than via hardware features. 211

Our prior work, libsmctrl [4], addresses the enforcement problem by leveraging undocumented hardware capabilities to enforce partitions of TPCs. Unfortunately, like earlier mechanisms, libsmctrl requires merging tasks into the same context to concurrently execute them, compromising logical isolation and transparency.

NVIDIA-provided solutions. Partitioning solutions from NVIDIA are able to enforce
partitioning at a hardware level, but were not designed with the needs of an embedded realtime system in mind. The two principal options provided by NVIDIA are the Multi-instance
GPU (MiG) feature, and the Multi-Process Service (MPS).

NVIDIA MiG [29] allows multiple contexts to concurrently run on the GPU by splitting 220 the GPU into static, fixed-size partitions. Each partition may then concurrently run a 221 different task. Partition options are highly limited, with at best four possible partition sizes, 222 and the smallest partition size is 14 SMs. MiG is implemented by duplicating every part 223 of the hardware scheduling pipeline for every GPC [9], and is only available on NVIDIA's 224 highest-end datacenter GPUs. This approach cannot provide fine-granularity partitioning, 225 and requires hardware modifications that NVIDIA has shown no intent to make widely 226 available. This leads us to classify MiG as non-granular and non-portable in Table 1. 227

Conversely, NVIDIA MPS [27] is available on any recent discrete NVIDIA GPU. MPS (since the Volta architecture),³ enables the Compute/Graphics Engine to concurrently execute multiple tasks, but does not control which SMs each task is assigned to. This may result in two tasks sharing an SM [27], which can cause as much as a 100x performance degradation [40]. Due to this limitation, we classify MPS in Table 1 as only partially providing granular,

³ Pre-Volta-generation MPS works differently, providing little-to-no isolation between co-running tasks. When we refer to MPS in this paper, we refer to the Volta-generation-and-newer version.

18:8 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems

hardware-enforced partitioning. Other properties of MPS are undocumented, and we derive those later in this work.

235 3.3 Barriers to Better Compute Partitioning

²³⁶ Unfortunately, devising better solutions for spatial partitioning of GPU compute has been
²³⁷ stymied by the scarce details released about NVIDIA GPU architecture and capabilities—
²³⁸ even instruction encodings are secret [14, 17, 16]. Many years of investigation have begun to
²³⁹ contravene this limitation.

Otterness et al. [34] and Amert et al. [2] elucidated the scheduling of CUDA kernels via 240 black-box experiments and introduced the cuda_scheduling_examiner tool, which we use. 241 Olmedo et al. [32] and our prior work [4] used these results and other sources to construct 242 a model of how the underlying GPU kernel dispatch hardware works. Parallel work by 243 Capodieci et al. [6], Spliet and Mullins [38], and us [5] clarified the preemption and high-level 244 scheduling capabilities of NVIDIA GPUs. This last work [5] also introduced two tools we 245 use: the nvdebug tool for directly extracting GPU state, and the gpu-microbench suite for 246 examining scheduling behavior. Outside of the academic community, the Nouveau [24] and 247 Mesa [23] reverse-engineered NVIDIA GPU driver projects have documented GPU hardware 248 capabilities and CPU-to-GPU interfaces. We lean heavily on all these prior works throughout 249 our paper. 250

²⁵¹ **4** How Does MPS Work?

In search for a better spatial partitioning mechanism for real-time embedded systems, we begin by investigating the only portable mechanism for co-running unmodified tasks on NVIDIA GPUs—MPS. What hardware capabilities does MPS leverage, and how MPS fall short of the properties we desire in a spatial partitioning mechanism? Given the absence of prior work, we investigate these questions experimentally. This is relevant both to our work, and to the safety of other works which propose using MPS in embedded systems.

258 4.1 Methodology

To determine how MPS interacts with the GPU scheduling pipeline, we applied the nvdebug tool [5], gpu-microbench suite [5], and cuda_scheduling_examiner toolkit [34] to test and observe GPU state and behavior with and without MPS. Adding MPS changed more aspects of GPU state than nvdebug was able to display, so we extended it on Ada (2022) and older GPUs, drawing layout information from open-source NVIDIA code [30, 26, 28] and the nouveau driver [24]. Our improved version of nvdebug is available.⁴

To understand the semantic meaning of this newly accessible GPU state, we leveraged context and definitions from NVIDIA patents touching on the runlist scheduler [10], kernel scheduling pipeline [36, 1], MPS [8] and MiG [9]. Unfortunately, patents may describe infeasible or impossible devices, and so we only tenuously relied on them, verifying described behavior with experiments.

⁴ Available at http://rtsrv.cs.unc.edu/cgit/cgit.cgi/nvdebug.git and within our artifact.



Figure 6 Compute/Graphics Engine runlist of three tasks; reillustrated from Fig. 4 with detail.



Figure 7 Compute/Graphics Engine runlist of three tasks, with Tasks 2 and 3 run as MPS clients.

4.2 MPS Terminology

MPS uses a client-server paradigm, where each CUDA-using task is a client of at most one MPS server task. The MPS server acts as an intermediary to the GPU for its clients, and allows clients to run concurrently with one another. Relative to the GPU, a system of two MPS clients and one MPS server would appear as a single task, since clients only access the GPU through the MPS server. Multiple MPS servers may exist, each with different clients.⁵ In this case, each MPS server would appear to be a separate task to the GPU.

4.3 How MPS Modifies Runlist Scheduling

We now investigate how adding MPS modifies arbitration between GPU-using tasks, *i.e.*,
how it modifies the runlist and associated data structures. To enable our subsequent analysis
of MPS's pitfalls, this section is particularly detailed.

We begin by reillustrating Fig. 4, with detail from our extensions to nvdebug, in Fig. 6. This figure retains the same runlist as Fig. 4, but expands on the configuration of each channel. The channel-configuration data structure is known as the channel's *instance block*, and besides describing the command queue (not shown for space)—specifies the virtual address space to be used for commands from the channel. Virtual address spaces are configured in a peculiar way; each channel includes an indexed list of page tables, and the page table to use is selected by specifying an index into that list. The list of page tables is called the "Subcontext

⁵ Support for multiple MPS servers is only implicitly documented. The environment variable CUDA_MPS_PIPE_DIRECTORY can be set to control where an MPS server advertises itself, and where CUDA searches for the MPS server. Two MPS servers cannot advertise to the same path, and so this variable must be set uniquely for each to allow multiple servers to exist.

18:10 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems

Table," and the channel's index into it is called the channel's "Subcontext ID."⁶ Adding to the confusion, this mechanism only selects the page table for the Compute/Graphics Engine; the instance block includes a separate page table configuration field for other engines, such as Copy.⁷ In our experiments, we always observed the same page table configuration for all engines; we mirror that in Fig. 6. Note that all channels in a TSG have an identical subcontext table—this is a requirement for channels sharing a context.⁸ All of these details are important in order to discuss how MPS effects address-space isolation.

In Fig. 7, we illustrate how runlist and channel configurations change when MPS is enabled. This is still a system of the same three GPU-using tasks, but an MPS server has been started with Tasks 2 and 3 as clients (Task 1 remains independent of MPS). Both the runlist location and virtual address space configuration have changed for the tasks now running as MPS clients. We discuss each change in turn.

With MPS enabled, the two tasks now running as MPS clients no longer have independent 300 TSGs in the runlist. The runlist only contains two TSGs: one for Task 1 (TSG 0), and one 301 for the MPS server (TSG 1). While Tasks 2 and 3 do not retain independent TSGs, they do 302 retain independent channels within the MPS server's TSG (Channels 3–4 for Task 2 and 5–6 303 for Task 3). (Note that the MPS server has taken Channel 2 for its own use, and so Tasks 2 304 and 3 are forced to use channels with higher IDs.) When this runlist is scheduled, Task 1 305 will, as before, execute for 2 ms before its budget expires and it is preempted. The GPU will 306 then switch to the next TSG in the runlist, the one for the MPS server. This will likewise be 307 run for 2 ms before the GPU loops back to Task 1. While the MPS server's TSG is active, 308 commands from all its channels are sent to the Compute/Graphics Engine concurrently. 309 The effect is such that all the MPS client tasks's CUDA streams concurrently execute as 310 though they were in the same program. Note that the two MPS clients share a single 2 ms 311 time-slice with a 4 ms period, whereas before they each got a single 2 ms time-slice at a 6 ms 312 period. This is a reduction in GPU time available jointly to the two tasks, from two-thirds 313 to one-half. This is a side-effect of enabling MPS to be aware of: if any other tasks in the 314 system continue to run without MPS, the total GPU time available to all MPS clients will 315 be less than the total time those tasks would have collectively received if run apart.⁹ 316

We now consider how virtual address space configurations change—or stay the same—with the addition of MPS. Visible in Fig. 7, at the bottom of the Channel 2–6 Instance Blocks, the Subcontext Table for every channel of the MPS server TSG is triple the size of the tables in Fig. 6. This allows the table to include separate page tables for MPS, Task 2, and Task 3; the appropriate one is selected by the Subcontext ID on each channel. Each channel also has a non-compute page table identical to the one selected by its Subcontext ID. In this manner, Tasks 2 and 3 retain distinct virtual address spaces, just as without MPS.¹⁰

324 4.4 Evaluating Spatial Partitioning in MPS

What does this mean for MPS's suitability to real-time embedded systems? We answer by considering MPS under each of our desired properties (Sec. 2.5). Without loss of generality,

⁶ Subcontext ID is also "Virtual Engine ID" (VEID) in some sources.

⁷ See line 297 in manuals/ampere/ga100/dev_ram.ref.txt [30].

⁸ See line 293 in manuals/ampere/ga100/dev_ram.ref.txt [30].

⁹ The time-slice length of the MPS server could be extended to ensure tasks retain access to an equivalent fraction of GPU time, but this is unsupported by NVIDIA's software.

¹⁰Why the convoluted TSG-wide Subcontext Table, rather than a per-channel page table? This may speed up context switches by allowing a TSG's page tables to be read all at once, rather than requiring a scan of all a TSG's channels.

³²⁷ we assume that all MPS clients are associated with a single server in this subsection.

328 4.4.1 Portability

The version of MPS we study is supported on all of NVIDIA's GPUs since Volta (2018), including their embedded "Tegra" GPUs (as of CUDA 12.5).

331 4.4.2 Logical Isolation

Without MPS, contexts are isolated from one another in their GPU addresses spaces, hardware 332 scheduling decisions, and exception handling. MPS preserves isolation in only the first area. 333 The MPS documentation states that MPS clients have fully isolated virtual address 334 spaces [27, Sec. 1.1.2]. Our findings support this—each MPS client exclusively uses its 335 own page table. This follows from the per-subcontext page tables discussed in Sec. 4.3. 336 Specifically, the unique-per-MPS-client subcontext ID is passed along with commands to the 337 Compute/Graphics Engine,¹¹ and this ID is used to access and maintain per-subcontext page 338 table state throughout the scheduling and execution pipeline [8]. This isolation also covers 339 non-compute engines, as their page $table^6$ is always configured to match the per-subcontext 340 page table. Only the kernel-level driver can change a channel's page table or subcontext 341 ID,¹² ensuring that no client may reconfigure itself to access another client's pages. 342

³⁴³ Other areas lack isolation, bringing us to our first pitfall:

▶ Pitfall 1. MPS clients share a per-server limit on the number of concurrently executing
 kernels.

The use of subcontexts does not isolate MPS clients from one another in the GPU's 346 hardware scheduling pipeline. Our prior work on this pipeline [4] found that two tasks co-347 running in a single context may conflict due to a hardware limit on the number of concurrent 348 kernels.¹³ While an NVIDIA patent [8] suggests this limit is maintained per subcontext, 349 we did not observe this, even on the most-recent Hopper- and Ada-generation GPUs. We 350 tested by launching several hundred kernels, and found that the number of kernels a task can 351 concurrently execute is reduced by the number of kernels concurrently executing in other 352 MPS clients. 353

▶ Pitfall 2. A crash in any MPS client may crash all MPS clients.

The MPS documentation warns that MPS clients are not isolated from "fatal GPU faults" in other MPS clients [27, Sec. 2.2.3]. Such faults include errors such as out-of-bounds memory accesses in kernels. We tested and found that this lack of isolation persists on NVIDIA's latest Hopper- and Ada-generation GPUs (at least for out-of-bounds memory accesses).

NVIDIA's drivers report exceptions on a per-subcontext basis,¹⁴ and other documentation hints that exceptions do not necessarily halt the entire TSG.¹⁵ We urge NVIDIA to leverage these capabilities to enhance MPS's fault isolation.

In summary, MPS preserves address space isolation, but does not fully isolate the scheduling pipeline nor prevent on-GPU exceptions from crashing other MPS clients. This partial isolation is better than libsmctrl and software partitioning, but worse than MiG.

 $^{^{11}\,\}mathrm{See}$ line 2525 in manuals/ampere/ga100/dev_pbdma.ref.txt [30].

¹²See line 2545 in manuals/ampere/ga100/dev_pbdma.ref.txt and line 525 in dev_ram.ref.txt [30].

¹³Specifically, "task slot exhaustion" in our prior work [4].

¹⁴See line 1115 in src/nvidia/src/kernel/gpu/mmu/arch/volta/ kern_gmmu_gv100.c [26].

 $^{^{15}}$ See line 666 in manuals/ampere/ga100/dev_runlist.ref.txt

18:12 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems



Figure 8 Four streams, each with two one-block kernels launched in them on an otherwise-idle NVIDIA Titan V GPU. Left is without MPS, right is with MPS defaults.

365 4.4.3 Transparency

We define a transparent partitioning mechanism as one that does not require any changes to tasks (Sec. 2.5). This property encompasses more than no binary modification; tasks should not need to be modified to account for a different set of available features. While MPS does not require modifying task binaries, it changes the set of supported features and the scheduling behavior in a semi-transparency-compromising way.

Pitfall 3. MPS clients cannot launch kernels from on the GPU.

CUDA Dynamic Parallelism (CDP) [25, Sec. 9]¹⁶ allows for one CUDA kernel to launch another without involving the CPU. MPS omits support for this feature [27, Sec. 2.3.2]. We verified that this restriction persists even on NVIDIA's latest Ada-generation GPUs; any MPS client attempting to use this feature will get an error during initialization. We could not identify a conclusive reason why this feature is unsupported with MPS.

Pitfall 4. MPS clients support fewer concurrent CUDA streams.

The MPS server also changes at least one scheduling-related property for its clients: the 378 number of channels available to a task. Each MPS client only has access to two channels by 379 default, whereas each non-MPS task has access to eight by default. An insufficient number 380 of channels can result in implicit synchronization [5], so this changed default can significantly 381 impact scheduling behavior, as shown in Fig. 8. Each subfigure shows how eight single-block 382 kernels in four streams execute over time (x-axis) on the GPU's SMs (y-axis). Kernel launch 383 times are indicated with arrows at bottom, and the stream of each kernel is color-and-pattern 384 coded. On the left, we show the system running without MPS; on the right, we show it 385 running with MPS. No other work is running in this system. With MPS (right), we see that 386 kernel launches are blocked due to channel exhaustion (as in [5]). 387

Pitfall 5. MPS clients receive fake SM IDs.

¹⁶ CDP is also called CUDA Native Parallelism (CNP) or "GWC" [22, 26].

18:13

In GPU kernels, the special register "smid "returns the SM identifier on which a particular 389 thread is executing" [31, Sec. 10.8]. Unfortunately, we found this register returns inconsistent 390 values across MPS clients. For example, each MPS client's kernel's blocks start on %smid 391 zero, and subsequent blocks start on sequentially-increasing SM IDs—even if another client 392 claimed to be executing on those SMs. This indicates that the "smid register is emulated for 393 each MPS client, as suggested in an NVIDIA patent [8]. We tested and found this behavior 394 on Volta-, Turing-, Ampere-, and Ada-generation GPUs. This pitfall primarily hinders GPU 395 study by obfuscating the SM assignment algorithm. 396

In summary, MPS affects available CUDA features, scheduling concurrency, and hardware
 behavior, making it only partially transparent.

399 4.4.4 Overheads

GPU commands, such as kernel launches, are sent to the GPU via the queue encapsulated 400 within a channel. MPS clients have direct access to their channel queues, and so no overheads 401 are added to the kernel-launch critical path. We verified this on Volta-, Turing-, and 402 Ada-generation GPUs via the measure_launch_oh benchmark we add to gpu-microbench. 403 Task startup overheads, such as the time for library loading, are affected by MPS. The 404 MPS server is lazily initialized, meaning the first MPS client pays an extra startup overhead. 405 As this can be avoided by using a dummy task to pull forward server initialization, we do 406 not consider it a pitfall. We show other startup overheads to be negligible in Sec. 6. 407

408 4.4.5 Partitioning Capability

MPS only supports a static, per-task limit on what fraction of a GPU's TPCs a task may use.
This limit, set via an environment variable,¹⁷ is called "Execution Resource Provisioning."

▶ **Pitfall 6.** MPS partitions are not bound to specific SMs.

No API is provided to specify a set of SMs, TPCs, or GPCs onto which an MPS client should be partitioned, and we find that MPS also does not select a set internally.

The MPS documentation does not clarify how partitioning is enforced. However, we find that when MPS is running, a GPU register¹⁸ is set to enable "dynamic partitioning" in the Work Distribution Unit (WDU). The WDU is the GPU hardware unit responsible for dispatching blocks of pending kernels to TPCs [4]. A patent filed by NVIDIA when execution resource provisioning was added to MPS appears to describe this feature [8].

In the patent, NVIDIA describes their dynamic execution resource partitioning system as 419 associating each subcontext (MPS client) with a configurable number of credits. Every time 420 a block of a kernel is dispatched to a TPC not previously occupied by its subcontext, the 421 number of credits is decremented. Once a subcontext reaches zero credits, the WDU will 422 only dispatch blocks from that subcontext to TPCs already in-use by that subcontext. When 423 a TPC is vacated by all work from the subcontext, the number of credits is incremented. In 424 effect, this caps the number of TPCs that any subcontext may concurrently have kernels 425 executing on, but makes no guarantees about which TPCs are assigned to each subcontext. 426 We find this behavior holds when partitioning with MPS. This is a problem, as GPU 427 partitions mis-aligned with hardware can lead to cache, bus, TLB, and other interference [4, 428 40]. We give an experimental example in the following pitfall. 429

 $^{^{17}{\}tt CUDA_MPS_ACTIVE_THREAD_PERCENTAGE}$

¹⁸ NV_PGRAPH_PRI_CWD_PARTITION_CTL [30]; CWD is a synonym for the WDU [4].

18:14 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems





(b) Percentage of TPCs partially or fully occupied.

Figure 9 Three tasks run with MPS's Execution Resource Provisioning feature on an NVIDIA Titan V GPU. Arrows indicate release times. MPS quirks leave half the TPCs idle after time t_2 , despite both remaining tasks (Tasks 1 and 3) having outstanding work, and both being allowed up to 50% of the GPU.

▶ Pitfall 7. The hardware implementation of MPS's partitioning feature is prone to assigning
 two tasks to the same set of SMs, leaving other SMs idle.

As MPS makes no guarantees about which TPCs are assigned to which MPS client, two 432 clients with partitions of 50% may each be assigned to the same set of TPCs. We demonstrate 433 this with an experimental result in Fig. 9 for a system of three tasks running as MPS clients, 434 each with a 50% limit on the amount of the GPU they may use. Our tasks' executions 435 briefly overlap, meaning that there will be a transient period where the sum of the allocated 436 GPU capacity is 150%, forcing tasks to share portions of the GPU. However, such sharing 437 unexpectedly persists, even after the sum of allocated GPU capacity returns to at most 100%. 438 This is due to a flaw in the design of MPS's execution resource provisioning feature. 439

Since we are using MPS, Pitfall 5 applies, and we cannot generate a plot that shows 440 which specific TPCs each task is applied to. Instead, we plot the total number of GPU cores 441 in use in Fig. 9a, and illustrate the total number of TPCs with any cores in use in Fig. 9b. 442 Both plots show the same time interval for the same experiment, and the release times for 443 each task are indicated with arrows. For this example, it is important to know that the 444 WDU tries to spread work out to as many TPCs as possible, and assigns blocks of large 445 kernels to less-occupied TPCs first [32]. This means that all TPCs can be occupied, even if 446 only a fraction of a GPU's cores are busy. 447

Task 1 starts first. Task 1 is a light workload and only requires about a quarter of the 448 GPU's cores, but the WDU spreads that work across as many TPCs as possible, causing 449 Task 1 to partially occupy 50% of the TPCs. Task 2 is a heavy workload, and upon joining, 450 it fully occupies the 50% of the TPCs unoccupied by Task 1. At this point, t_0 , all TPCs 451 are occupied, but not all cores. When Task 3 joins at t_1 it is placed on the least-occupied 452 TPCs. Since Task 2 has fully utilized its TPCs, Task 3 is placed on the partially occupied 453 TPCs of Task 1. This poses a problem when Task 2 terminates at t_2 : Task 3 is unable to 454 migrate onto the newly freed TPCs. Because both Task 1 and Task 3 have already maxed 455 out their active-on-50%-of-TPCs limit, both remain stuck on the same set of TPCs. This 456 leaves 50% of the GPU unused, despite pending work from two running tasks each with a 457 50% GPU allocation. Ideally, after Task 2 terminated, Task 3 should have been migrated to 458 the newly-freed TPCs, keeping 100% of the TPCs occupied. 459

The partition settings we use, where the sum of allocated compute exceeds 100%, is suggested by NVIDIA as a "more optimal strategy" [27, Sec. 2.3.5]. We suggest that NVIDIA

462 remove said recommendation.

This pitfall could be worse. The WDU will normally use an alternate assignment algorithm for small kernels, packing kernels unto SMs before spreading them onto idle SMs [32]. If this behavior persisted between kernels of different MPS clients, a system of only two MPS clients could be bound to to the same set of TPCs. Fortunately, after repeating a variant of the experiments from prior work [32], we did not find a similar behavior between MPS clients.¹⁹ This unexpected behavior appears triggered by the aforementioned GPU register¹³ for dynamic partitioning—zeroing this register restores the normal algorithm.

470 ▶ **Pitfall 8.** The partition size for an MPS client is static.

The partition size is specified at context creation, and NVIDIA provides no API to change the partition size for an already created context.

In summary, MPS's partition sizes cannot be changed dynamically, its partition boundaries
are weak, and its partitions may unexpectedly overlap. This completes our classification of
MPS in Table 1, bringing us to the conclusion that no prior work is adequate for spatial
partitioning GPU compute cores in an embedded, real-time system.

477 **5** Compute Partitioning with nvtaskset

We build a new tool called nvtaskset that allows for transparent GPU partitioning between
unmodified CUDA-using tasks. We do this by combining the logical isolation and transparency
of MPS with the partitioning capability of libsmctrl.

481 5.1 nvtaskset

491

nvtaskset works like Linux's taskset tool. While taskset sets the CPU cores that any
task may use, nvtaskset sets the GPU cores that a task may use. To co-run two unmodified
DNN tasks on the GPU with nvtaskset, one would run the shell commands:

```
485
486 # Launch and co-run dnn-one and dnn-two with 16 TPCs each
487 $ nvtaskset -t 0-15 ./dnn-one
488 $ nvtaskset -t 16-31 ./dnn-two
```

⁴⁹⁰ The partitions could later be resized:

```
492 # Assuming dnn-one is PID 100, change its allocation to 8 TPCs
433 $ nvtaskset -p -t 0-7 100
```

⁴⁹⁵ And tasks started without nvtaskset could later be moved to a partition:

```
496
497 # Start dnn-three (PID 102) unpartitioned, then move it to TPCs 8-15
498 $ ./dnn-three &
499 [1] 102
500 $ nvtaskset -p -t 8-15 102
```

⁵⁰² Partitions can also be specified in terms of GPCs instead of TPCs (if nvdebug is loaded):

¹⁹ To workaround Pitfall 5, we observed the change in block distribution by limiting each task to one TPC, and chaining a large kernel after the small one that would be assigned to the same TPC without MPS. The large kernels do not limit the utilization of each other, indicating that the preceding small kernel of each client ran on a different TPC.



Figure 10 Same experiment as in Fig. 9, but with libsmctrl used to partition Task 1 onto different TPCs than Tasks 2 and 3.

504# Start dnn-one on GPC 0 and 2\$05\$ nvtaskset -g 0,2 ./dnn-one

503

nvtaskset is distributed as an extension to libsmctrl, has no dependencies beyond CUDA,
 and can be built and installed via a simple make install invocation.

509 5.2 Implementation of nvtaskset

nvtaskset uses MPS to allow tasks to co-run, uses parts of libsmctrl to enforce partitioning,
 uses shared-library interception to apply to every CUDA-using task, uses shared memory
 to communicate dynamic partition updates, and uses GPU topology registers to associate
 GPCs with TPCs.

Co-running built on MPS. nvtaskset associates all tasks launch via it with an automatically started MPS server. This allows tasks to co-run. Note that subsequently launched
 CUDA tasks will automatically connect to this MPS server, allowing them to co-run and be partitioned, even if not launched via nvtaskset.²⁰ In order to mitigate Pitfall 4 from MPS, nvtaskset configures the number of channels for each CUDA-using task to eight.²¹

⁵¹⁹ **Partitioning built on libsmctrl.** MPS provides no rigorous or dynamic means for parti-⁵²⁰ tioning compute cores, but libsmctrl [4] is able to provide both of these properties.

libsmctrl implements partitioning by modifying the TMD for each kernel immediately
 before it is uploaded to the GPU. The TMD contains a field that specifies which TPCs the
 hardware may run the kernel on, and libsmctrl modifies this field to affect partitioning.

Tasks dispatch kernels the same with and without MPS—by inserting launch commands into their channel queues. This allows libsmctrl to be used with MPS. Fig. 10 shows how using libsmctrl for partitioning, instead of MPS's Execution Resource Provisioning feature, addresses Pitfalls 6 and 7 of MPS. Specifically, we adjusted the experiment of Fig. 9 to use libsmctrl, assigned Task 1 the lower 50% of TPCs, and assigned Task 2 and Task 3 the

 $^{^{20}}$ To force a task to run independently of MPS, set the CUDA_MPS_PIPE_DIRECTORY environment variable for that task to /dev/null.

²¹ MPS clients mirror the behavior of non-MPS-tasks once the number of channels per task is configured to its non-MPS default (8 channels). We do this by setting the environment variable CUDA_DEVICE_MAX_CONNECTIONS to 8 before launching an MPS client.

⁵²⁹ upper 50%. The result is that Task 3 completes a quarter-second faster, since it does not get ⁵³⁰ stuck on the same TPCs as Task 1.

Unmodified task support. Unfortunately, libsmctrl requires minor task modification
 and recompilation to use, compromising its transparency. We address this limitation in
 nvtaskset. libsmctrl affects partitioning through interactions with the CUDA library;
 tasks only need a single API call to enable a partition. nvtaskset eliminates the need for
 modification by effectively making the loader perform this API call.

Specifically, we make nvtaskset fully transparent to a task by tricking the loader into 536 loading the libsmctrl library before loading CUDA. (We do this by naming the compiled 537 form of our library libcuda.so.1, and then linking our library against the real libcuda.so. 538 libcuda.so.1 is always dynamically loaded, even in staticly linked binaries, so our library is 539 always loaded. Our library does not intercept CUDA library functions; we rely on the recursive 540 search behavior of the loader to resolve all the CUDA symbols to the real libcuda.so.) 541 Once leaded, our library uses a load-time constructor function to load CUDA, set up the 542 task for partitioning changes via **nvtaskset**, and register the TMD interception callback 543 with CUDA to enable partitioning. The loader than resumes as normal, and whenever the 544 task executes a kernel launch, our pre-registered callback will be triggered inside CUDA and 545 the TPC mask will be applied—no task modification required. 546

⁵⁴⁷ Dynamically changeable partitions. nvtaskset supports dynamic partition changes by
exposing a shared-memory region from each CUDA-using task that contains the current
partition setting. Changes to this setting (*e.g.*, via nvtaskset -p) are automatically detected
and applied to subsequent kernel launches by our callback, avoiding Pitfall 8 of MPS.

Specifying partitions by GPC. nvtaskset supports specifying partitions as a set of GPCs, 551 and internally translates that into a set of TPCs by correcting and utilizing libsmctrl's API 552 for obtaining GPC-to-TPC mappings. libsmctrl assumed that SM IDs are assigned linearly 553 to on-chip GPCs, such that the first n SM IDs would be in GPC 0, the next n in GPC 1, 554 and so on. We uncover that SM ID to GPC mappings are arbitrary, and configured by the 555 NVIDIA driver in a striping-like configuration by default. Knowing these mappings is critical 556 to align partitions on GPC boundaries, and so we fixed the API for determining SM-to-GPC 557 mappings in libsmctrl and used this in nvtaskset. We repeated the "Partitioning Strategy" 558 experiments from the libsmctrl paper [4] on nvtaskset, but found no significant differences. 559 We recommend partitioning on GPC boundaries when possible, as NVIDIA's Thread 560 Block Groups feature will only work for a task when it has access to at least one full GPC. 561

562 5.3 Limitations of nvtaskset

⁵⁶³ nvtaskset has several limitations inherited from MPS and libsmctrl. nvtaskset:

- ⁵⁶⁴ Is only compatible with CUDA-using tasks.
- 565 Does not support tasks that use multiple GPUs.
- 566 Cannot affect partition changes on already-launched kernels.
- ⁵⁶⁷ Is still subject to Pitfalls 1–3 of MPS.
- ⁵⁶⁸ Only supports up to 15 co-running tasks.²²

²²TSGs are limited to a maximum of 128 channels [30], so an MPS server can only service 15 clients of 8 channels each (after subtracting 8 MPS-server-internal channels). We verified this 15-client limit on

18:18 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems



Figure 11 Overheads of each partitioning mechanism.

nvtaskset is also subject to Pitfall 5 by default. While the %smid register can be restored 569 to its consistent, non-MPS behavior by toggling off the WDU's dynamic partitioning register, 570 we chose not to do this. To better support tasks which ignore the CUDA programming model 571 and attempt to execute work on self-selected SMs (such as persistent threads [43, 13]), we 572 can "hide" the existence of unallocated SMs with the WDU's dynamic partitioning capability. 573 To do this, leave on fake SM IDs, and set MPS's partition limit to match the number of 574 SMs allocated by nvtaskset. This can ensure that tasks only see allocated SMs, with 575 seemingly-contiguous SM IDs. 576

⁵⁷⁷ However, **nvtaskset** is portable to any recent discrete NVIDIA GPU, does not require task ⁵⁷⁸ modifications, preserves address space isolation, and is available as open-source software.²³

579 **6** Evaluation

We evaluate nvtaskset against MiG, MPS, and (where applicable) against no partitioning. We compare overhead impact, strength of partition enforcement, and granularity of partitioning. Our other target properties—portability, logical isolation, transparency, and dynamic reconfigurability—are largely binary properties, and have already been discussed.

Experiments in this section were run on an NVIDIA A100 40GB GPU running CUDA 12.4 on a 6-core, Linux 5.4 machine with the NVIDIA 550.78 GPU driver. At time of writing, the A100 is the only GPU that supports all of MiG, MPS, and nvtaskset. We use the SE-packed partitioning strategy [4, 33] to assign TPCs to nvtaskset partitions. As MiG cannot perform a 50/50 split of the A100, we use a 57/43 split for all mechanisms.²⁴

6.1 Evaluating Partitioning Overheads

We measure startup overhead and launch overheads for all mechanisms via benchmarks run in the 57% partition. Startup overhead is the time from exec() to first kernel running on the GPU of a minimal program, and launch overhead is the time from cuLaunchKernel() to the

both Volta- and Ada-generation GPUs.

²³ Available at http://rtsrv.cs.unc.edu/cgit/cgit.cgi/libsmctrl.git and within our artifact.

 $^{^{24}}$ MiG can only partition on GPC boundaries, and there are seven GPCs in the A100, meaning that a 57/43 split is as close as MiG can get to 50/50.

kernel running on the GPU. We added these benchmarks to the gpu-microbench suite [5],²⁵ and ran each experiment once to prime caches before gathering the data displayed in Fig. 11.

• Observation 1. No partitioning mechanism adds startup or launch overheads.

As shown in Fig. 11a and Fig. 11b, no partitioning mechanism worsens observed worst- or average-case startup or launch times.²⁶

598 • Observation 2. Only MiG reduces launch overheads.

⁵⁹⁹ Uniquely, MiG statically binds a hardware scheduling pipeline to each partition's TPCs [9] ⁶⁰⁰ such that only a subset of TPCs have to be set up, and considered as part of a kernel launch. ⁶⁰¹ We suspect this is what lowers launch and startup overheads with MiG.

602 • Observation 3. MPS and **nvtaskset** have the lowest startup overheads.

With MPS, newly launched tasks (MPS clients) only initialize channels and subcontexts, with the MPS server providing the parent context and TSG. This appears to significantly reduce startup overheads for MPS and the MPS-based nvtaskset, more than compensating for the added client-server communication cost of MPS.

607 6.2 Evaluating Partition Enforcement

To evaluate how strongly partitions are enforced, we measured the execution time of a $6144 \times$ 608 6144 matrix multiply (yielding 36×2^{10} blocks of 1024 threads) ("MM6144") in a 57% partition 609 while interfering tasks executed on the remainder of the GPU. We used the mandelbrot and 610 random walk tasks from the cuda scheduling examiner toolkit as compute- and memory-611 heavy interfering tasks respectively. Interfering tasks executed continuously, out of sync with 612 each other and the matrix multiply. We carefully configured this experiment to sidestep the 613 pitfalls of MPS discussed in Sec. 4.4, focusing exclusively on how well a partition is enforced 614 in an otherwise-ideal scenario. Good partition enforcement means that the execution time 615 distribution of our MM6144 task does not shift in the presence of interfering tasks. 616

The results against memory- and compute-heavy interference are shown in Fig. 12a and Fig. 12b respectively. Each plot includes a baseline "Scaled" time for comparison—this is *not* a measured value, but is a scaled-up value derived from the execution time of the benchmark without any partitioning or interference. For example, if MM6144 took 244 ms when run alone on the GPU, the "Scaled" value for a 57% partition would be 144/0.57 = 253 ms. This value is the minimum execution time possible for the MM6144 task in a 57% partition before accounting for sympathetic caching effects.

We specifically choose this matrix multiply task as it was most sensitive to interference among several other synthetic tasks we tested, and it could be more-precisely timed than a full neural network. Since we use the same partitioning mechanism as libsmctrl, the enforcement results previously shown for real tasks [4] continue to apply.

▶ Observation 4. MPS's partitioning mechanism has strictly worse predictability and efficiency, even when used correctly.

²⁵ Available at http://rtsrv.cs.unc.edu/cgit/cgit.cgi/gpu-microbench.git and within our artifact.

²⁶ Fig. 11b omits 100th percentile (max) times, as use of Linux's isolcpus, sched_yield() and nohz=full

options were insufficient to eliminate all noise, potentially due to SMIs [21] or interrupts.



(a) Partition enforcement versus memory-heavy competition (random walks of GPU DRAM).

(b) Partition enforcement versus compute-heavy competition (Mandelbrot set generation). Axis cropped.

Figure 12 Partition enforcement vs. a memory-heavy or compute-heavy competitor. Specifically, 0, 25, 50, 75, and 100th percentile time to execute many large matrix multiplies in a 57% partition for each partitioning mechanism while competing tasks run in the remaining 43% partition (1,000 samples each).

When competing work is memory-bound, as in Fig. 12a, MPS does nothing to prevent memory contention, only limiting the compute available for MM6144. This results in average (line in box) and worst-case (top of whisker) execution times for our MM6144 task only slightly better than with no partitioning at all.

▶ Observation 5. *nvtaskset* can provide partition enforcement approaching MIG without requiring hardware modifications.

For memory-bound competing work (Fig. 12a), MiG beats nvtaskset, likely because MiG also partitions DRAM. Interestingly, even though nvtaskset includes *no explicit memory partitioning*, its implicit partitioning of the L0, L1, and TLB caches by aligning partitions to GPCs appears to be enough to beat MPS and approach MiG's performance.

Surprisingly, for compute-bound competing work (Fig. 12b), nvtaskset *beats MiG* (and every other approach) on average- and worst-case execution times—without the hardware modifications of MiG. Upon further investigation, we found that MiG compromises compute speed to serve other design goals. This unexpected behavior is *not documented*, and causes MiG to fail in other ways, as we will explore under Obs. 7.

In all, nvtaskset enforces partitioning much better than MPS, and even better than MiG in some cases—all without requiring task, driver, or hardware modification.

647 6.3 Evaluating Partition Granularity

To evaluate partitioning granularity, we recorded the total execution time of a 8192×8192 matrix multiply (yielding 64×2^{10} blocks of 1024 threads) ("MM8192") at every possible partition size for each partitioning method and plot the results as points in Fig. 13. The lines in Fig. 13 represent the closest available configuration which allocates at most the specified number of TPCs. For example, there is no MiG configuration for 10 TPCs, so we plot time for the closest available allocation of no more than 10 TPCs—7 TPCs.

654 • Observation 6. *nvtaskset is the most granular partitioning mechanism.*



Figure 13 Partitioning granularity comparison. Specifically, mean time to execute a matrix multiply at each partition size, for each partitioning mechanism (10 samples).

Both MPS and nvtaskset can specify partition sizes down to the per-TPC level, visible as the different MM8196 execution times for each setting in Fig. 13. However, nvtaskset can assign specific TPCs to a partition, in contrast to MPS's generic percentage value. For this 54-TPC GPU, that means nvtaskset supports a total of 2⁵⁴ different partition settings per-task, MPS supports 54 per-task, and MiG only supports 5 per-task.

▶ Observation 7. MiG cannot access 9% of the A100 GPU cores (5 TPCs).

Visible in Fig. 13, the largest-available MiG partition contains only 49 TPCs, whereas MPS 661 or nvtaskset are able to access up to 54 TPCs. Upon further investigation, we find that no 662 configuration of MiG partitions on the A100 can access more than 49 TPCs total. Every 663 partition size is an even multiple of 7 TPCs, and 7 does not divide 54 evenly, wasting the 664 remainder—5 TPCs. This surprising issue is not documented, and means that enabling MiG 665 immediately and inherently disables 9% of the A100 GPU. Concerningly, we found that 666 this issue is even worse on NVIDIA's newer GPUs. On the H100 GPU (SXM-80GB version 667 tested), we found a loss of 6-15% (depending on the MiG partition size chosen).²⁷ 668

669 7 Conclusion

In this work, we developed a system-level spatial partitioning mechanism for NVIDIA GPU
 compute cores, nvtaskset. Our mechanism allows for GPU-using tasks to run *both* efficiently
 and time-predictably by running concurrently on disjoint sets of GPU cores.

We demonstrated that our mechanism is portable, transparent, and low-overhead, and has the ability to provide granular, dynamic, logically-isolated, and hardware-enforced partitions. As part of this work, we exposed critical pitfalls of NVIDIA's MPS-based partitioning mechanism, and revealed previously-undocumented capacity loss issues inherent to NVIDIA MiG. In future work, we aim to extend nvtaskset to support multiple GPUs and non-CUDA workloads, and to build GPU schedulers on nvtaskset that efficiently and predictably schedule tasks across both time and space.

²⁷We suspect the capacity loss stems from a decision to make unit-size MiG slices appear identical, despite differences in the underlying GPCs. Due to floorsweeping, some GPCs will have more working TPCs than others—those TPCs must be disabled to emulate identical GPCs when using MiG.

18:22 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems

680		References
681	1	Karim M Abdalla, Lacky V Shah, Jerome F Duluk Jr. Timothy John Purcell, Tanmov Mandal,
682		and Gentaro Hirota. Scheduling and execution of compute tasks. Jun 2015. U.S. Patent
683		9,069,609.
684	2	Tanva Amert, Nathan Otterness, Ming Yang, James H. Anderson, and F. Donelson Smith.
685		GPU scheduling on the NVIDIA TX2: Hidden details revealed. In <i>Proceedings of the 38th IEEE</i>
686		Real-Time Systems Symposium, pages 104–115, Dec 2017. doi:10.1109/RTSS.2017.00017.
687	3	Tanya Amert, Zelin Tong, Sergey Voronov, Joshua Bakita, F Donelson Smith, and James H
688		Anderson. TimeWall: Enabling time partitioning for real-time multicore+accelerator platforms.
689		In Proceedings of the 42nd IEEE Real-Time Systems Symposium, pages 455-468, Dec 2021.
690		doi:10.1109/RTSS52674.2021.00048.
691	4	Joshua Bakita and James H Anderson. Hardware compute partitioning on NVIDIA GPUs.
692		In Proceedings of the 29th IEEE Real-Time and Embedded Technology and Applications
693		Symposium, pages 54-66, May 2023. doi:10.1109/RTAS58335.2023.00012.
694	5	Joshua Bakita and James H Anderson. Demystifying NVIDIA GPU internals to enable reliable
695		GPU management. In Proceedings of the 30th IEEE Real-Time and Embedded Technology and
696		Applications Symposium, pages 294–305, May 2024. doi:10.1109/RTAS61025.2024.00031.
697	6	Nicola Capodieci, Roberto Cavicchioli, Marko Bertogna, and Aingara Paramakuru. Deadline-
698		based scheduling for GPU with preemption support. In Proceedings of the 39th IEEE Real-Time
699		Systems Symposium, pages 119–130, Dec 2018. doi:10.1109/RTSS.2018.00021.
700	7	Nicola Capodieci, Roberto Cavicchioli, Ignacio Sañudo Olmedo, Marco Solieri, and Marko
701		Bertogna. Contending memory in heterogeneous SoCs: Evolution in NVIDIA Tegra embedded
702		platforms. In Proceedings of the 26th IEEE International Conference on Embedded and Real-
703		Time Computing Systems and Applications, pages 1–10, Aug 2020. doi:10.1109/RTCSA50079.
704		2020.9203722.
705	8	Jerome F Duluk Jr, Luke Durant, Ramon Matas Navarro, Alan Menezes, Jeffrey Tuckey,
706		Gentaro Hirota, and Brian Pharris. Dynamic partitioning of execution resources, Apr 2022.
707	_	U.S. Patent 11,307,903.
708	9	Jerome F Duluk Jr, Gregory Scott Palmer, Jonathon Stuart Ramsey Evans, Shailendra Singh,
709		Samuel H Duncan, Wishwesh Anil Gandhi, Lacky V Shah, Eric Rock, Feiqi Su, James Leroy
710		Deming, et al. Techniques for configuring a processor to function as multiple, separate
711	10	processors, Feb 2022. U.S. Patent 11,249,905.
712	10	Samuel H Duncan, Lacky V Shah, Sean J Treichler, Daniel Elliot Wexler, Jerome F Duluk Jr,
713		independent streams in multi channel time slice groups. Sep 2016, U.S. Detent 0.442,750
714	11	Clear A Ellistt – Deal time scheduling for CDU with andiesting in changed activities
715	11	Glenn A Elliott. Real-time scheduling for GPUs with applications in davanced automotive
716		systems. Find thesis, The University of North Catolina at Chapel IIII, 2015. doi:10.17015/
/1/	10	Clann A Elliott Drugen C Word and James H Anderson CDUSung, A framework for real time
718	12	CPU management. In Proceedings of the 2/th Real Time Systems Symposium pages 33-44
719		Dec 2013 doi:10 1109/RTSS 2013 12
720	13	Kshitii Gunta Leff A Stuart and John D Owens A study of persistent threads style GPU
721	15	programming for GPGPII workloads. In Proceedings of the 2012 IEEE Innovative Parallel
723		Computing Conference, pages 1–14. May 2012. doi:10.1109/InPar.2012.6339596.
724	14	Ari B Haves Fei Hua Jin Huang Vanhao Chen and Eddy Z Zhang Decoding CUDA binary
724	14	In Proceedings of the 2019 IEEE/ACM International Sumposium on Code Generation and
726		<i>Optimization</i> , pages 229–241, Feb 2019, doi:10.1109/CG0.2019.8661186.
727	15	Saksham Jain, Iljoo Baek, Shige Wang, and Ragunathan Raikumar. Fractional GPUs: Software-
728	-	based compute and memory bandwidth reservation for GPUs. In <i>Proceedings of the 25th IEEE</i>
729		Real-Time and Embedded Technology and Applications Symposium, pages 29–41, Apr 2019.
730		doi:10.1109/RTAS.2019.00011.

16 Zhe Jia, Marco Maggioni, Jeffrey Smith, and Daniele Paolo Scarpazza. Dissecting the NVIDIA 731 Turing T4 GPU via microbenchmarking, Mar 2019. doi:10.48550/arXiv.1903.07486. 732 Zhe Jia, Marco Maggioni, Benjamin Staiger, and Daniele P. Scarpazza. Dissecting the NVIDIA 17 733 Volta GPU architecture via microbenchmarking, Apr 2018. doi:10.48550/arXiv.1804.06826. 734 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep 735 18 convolutional neural networks. Advances in Neural Information Processing Systems, 25:1097-736 1105, Dec 2012. 737 Jochen Liedtke, Hermann Hartig, and Michael Hohmuth. Os-controlled cache predictability 19 738 for real-time systems. In Proceedings of the 3rd IEEE Real-Time Technology and Applications 739 Symposium, pages 213-224, Jun 1997. doi:10.1109/RTTAS.1997.601360. 740 Lei Liu, Zehan Cui, Mingjie Xing, Yungang Bao, Mingyu Chen, and Chengyong Wu. A software 20 741 memory partition approach for eliminating bank-level interference in multicore systems. In 742 Proceedings of the 21st International Conference on Parallel Architecture and Compilation 743 Techniques, pages 367-376, Sept 2012. doi:10.1145/2370816.2370869. 744 Sizhe Liu, Rohan Wagle, James H Anderson, Ming Yang, Chi Zhang, and Yunhua Li. Autonomy 21 745 today: Many delay-prone black boxes. In Proceedings of the 36th Euromicro Conference on 746 Real-Time Systems, pages 12:1-12:27, July 2024. doi:10.4230/LIPIcs.ECRTS.2024.12. 747 22 Albert Meixner. System and method for launching data parallel and task parallel application 748 threads and graphics processing unit incorporating the same, Mar 2016. U.S. Patent 9,286,114 749 B2. 750 23 Mesa Project Authors. The Mesa 3D graphics library, 2022. URL: https://www.mesa3d.org/. 751 24 Nouveau Project Authors. Nouveau: Accelerated open source driver for nVidia cards, 2022. 752 URL: https://nouveau.freedesktop.org/. 753 NVIDIA. CUDA C++ programming guide, 2022. Version PG-02829-001 v11.8. 25 754 NVIDIA. Linux open GPU kernel module source, 2024. URL: https://github.com/NVIDIA/ 26 755 open-gpu-kernel-modules. 756 27 NVIDIA. Multi-process service, 2024. Version R555. 757 28 NVIDIA. nvgpu git repository, 2024. URL: git://nv-tegra.nvidia.com/linux-nvgpu.git. 758 29 NVIDIA. NVIDIA multi-instance GPU user guide, 2024. Version RN-08625-v2.0. 759 Open GPU documentation, 2024. 30 NVIDIA. URL: https://github.com/NVIDIA/ 760 open-gpu-doc. 761 31 NVIDIA. Parallel thread execution ISA, 2024. Version 8.5. 762 32 Ignacio Sañudo Olmedo, Nicola Capodieci, Jorge Luis Martinez, Andrea Marongiu, and Marko 763 Bertogna. Dissecting the CUDA scheduling hierarchy: a performance and predictability 764 perspective. In Proceedings of the 26th IEEE Real-Time and Embedded Technology and 765 Applications Symposium, pages 213–225, Apr 2020. doi:10.1109/RTAS48715.2020.000-5. 766 Nathan Otterness and James H Anderson. Exploring AMD GPU scheduling details by 33 767 experimenting with "worst practices". In Proceedings of the 29th International Conference on 768 Real-Time Networks and Systems, pages 24-34, Apr 2021. doi:10.1145/3453417.3453432. 769 Nathan Otterness, Ming Yang, Tanya Amert, James Anderson, and F Donelson Smith. 34 770 Inferring the scheduling policies of an embedded CUDA GPU. In Proceedings of the 13th 771 Annual Workshop on Operating Systems Platforms for Embedded Real Time Applications, 772 pages 47–52, Jul 2017. 773 Nathan Otterness, Ming Yang, Sarah Rust, Eunbyung Park, James H. Anderson, F. Donelson 35 774 Smith, Alex Berg, and Shige Wang. An evaluation of the NVIDIA TX1 for supporting real-775 time computer-vision workloads. In Proceedings of the 23rd IEEE Real-Time and Embedded 776 Technology and Applications Symposium, pages 353-364, Apr 2017. doi:10.1109/RTAS.2017.3. 777 Timothy John Purcell, Lacky V Shah, and Jerome F Duluk Jr. Scheduling and management 778 36 of compute tasks with different execution priority levels. U.S. Patent Application 13/236,473. Sujan Kumar Saha, Yecheng Xiang, and Hyoseung Kim. STGM: Spatio-temporal GPU 37 780 management for real-time tasks. In Proceedings of the 25th IEEE International Conference 781

18:24 Hardware Compute Partitioning on NVIDIA GPUs for Composable Systems

- on Embedded and Real-Time Computing Systems and Applications, pages 1–6, Aug 2019.
 doi:10.1109/RTCSA.2019.8864564.
- Roy Spliet and Robert Mullins. The case for limited-preemptive scheduling in GPUs for
 real-time systems. In *Proceedings of 14th Annual Workshop on Operating Systems Platforms* for Embedded Real-Time Applications, pages 43–48, Jul 2018.
- 39 Bo Wu, Guoyang Chen, Dong Li, Xipeng Shen, and Jeffrey Vetter. Enabling and exploiting flexible task assignment on GPU through SM-centric program transformations. In *Proceedings* of the 29th ACM on International Conference on Supercomputing, ICS, pages 119–130, Jun 2015. doi:10.1145/2751205.2751213.
- Tyler Yandrofski, Leo Chen, Nathan Otterness, James H Anderson, and F Donelson Smith.
 Making powerful enemies on NVIDIA GPUs. In *Proceedings of the 43rd IEEE Real-Time* Systems Symposium, pages 383–395, dec 2022. doi:10.1109/RTSS55097.2022.00040.
- Ming Yang, Nathan Otterness, Tanya Amert, Joshua Bakita, James H Anderson, and F Donelson Smith. Avoiding pitfalls when using NVIDIA GPUs for real-time tasks in autonomous systems. In *Proceedings of the 30th Euromicro Conference on Real-Time Systems*, pages 20:1–20:21, Jul 2018. doi:10.4230/LIPIcs.ECRTS.2018.20.
- Ming Yang, Shige Wang, Joshua Bakita, Thanh Vu, F Donelson Smith, James H Anderson,
 and Jan-Michael Frahm. Re-thinking CNN frameworks for time-sensitive autonomous-driving
 applications: Addressing an industrial challenge. In *Proceedings of the 25th IEEE Real- Time and Embedded Technology and Applications Symposium*, pages 305–317, Apr 2019.
 doi:10.1109/RTAS.2019.00033.
- 43 Chao Yu, Yuebin Bai, Hailong Yang, Kun Cheng, Yuhao Gu, Zhongzhi Luan, and Depei
 Qian. SMGuard: A flexible and fine-grained resource management framework for GPUs.
 IEEE Transactions on Parallel and Distributed Systems, 29(12):2849–2862, Jun 2018. doi:
 10.1109/TPDS.2018.2848621.
- Heechul Yun, Renato Mancuso, Zheng-Pei Wu, and Rodolfo Pellizzoni. PALLOC: DRAM
 bank-aware memory allocator for performance isolation on multicore platforms. In *Proceedings* of the 20th IEEE Real-Time and Embedded Technology and Applications Symposium, pages
 155–166, Apr 2014. doi:10.1109/RTAS.2014.6925999.
- 45 Yongkang Zhang, Haoxuan Yu, Chenxia Han, Cheng Wang, Baotong Lu, Yunzhe Li, Zhifeng
 Jiang, Yang Li, Xiaowen Chu, and Huaicheng Li. SGDRC: Software-defined dynamic resource
 control for concurrent dnn inference on nvidia gpus. In *Proceedings of the 30th ACM SIGPLAN*Annual Symposium on Principles and Practice of Parallel Programming, pages 267–281, Mar
 2025. doi:10.1145/3710848.3710863.
- 46 Husheng Zhou, Soroush Bateni, and Cong Liu. S³DNN: Supervised streaming and scheduling
 for GPU-accelerated real-time DNN workloads. In *Proceedings of the 24th IEEE Real-Time* and Embedded Technology and Applications Symposium, pages 190–201, Apr 2018. doi:
 10.1109/RTAS.2018.00028.
- An Zou, Jing Li, Christopher D Gill, and Xuan Zhang. RTGPU: Real-time GPU scheduling
 of hard deadline parallel tasks with fine-grain utilization. *IEEE Transactions on Parallel and Distributed Systems*, 34(5):1450–1465, May 2023. doi:10.1109/TPDS.2023.3235439.