

1 **Avoiding Pitfalls when Using NVIDIA GPUs for** 2 **Real-Time Tasks in Autonomous Systems***

3 **Ming Yang**

4 The University of North Carolina at Chapel Hill, USA
5 yang@cs.unc.edu

6 **Nathan Otterness**

7 The University of North Carolina at Chapel Hill, USA
8 otternes@cs.unc.edu

9 **Tanya Amert**

10 The University of North Carolina at Chapel Hill, USA
11 tamert@cs.unc.edu

12 **Joshua Bakita**

13 The University of North Carolina at Chapel Hill, USA
14 jbakita@cs.unc.edu

15 **James H. Anderson**

16 The University of North Carolina at Chapel Hill, USA
17 anderson@cs.unc.edu

18 **F. Donelson Smith**

19 The University of North Carolina at Chapel Hill, USA
20 smithfd@cs.unc.edu

21 **Abstract**

22 NVIDIA's CUDA API has enabled GPUs to be used as computing accelerators across a wide
23 range of applications. This has resulted in performance gains in many application domains, but
24 the underlying GPU hardware and software are subject to many non-obvious pitfalls. This is
25 particularly problematic for safety-critical systems, where worst-case behaviors must be taken into
26 account. While such behaviors were not a key concern for earlier CUDA users, the usage of GPUs
27 in autonomous vehicles has taken CUDA programs out of the sole domain of computer-vision and
28 machine-learning experts and into safety-critical processing pipelines. Certification is necessary in
29 this new domain, which is problematic because GPU software may have been developed without
30 any regard for worst-case behaviors. Pitfalls when using CUDA in real-time autonomous systems
31 can result from the lack of specifics in official documentation, and developers of GPU software
32 not being aware of the implications of their design choices with regards to real-time requirements.
33 This paper focuses on the particular challenges facing the real-time community when utilizing
34 CUDA-enabled GPUs for autonomous applications, and best practices for applying real-time
35 safety-critical principles.

36 **2012 ACM Subject Classification** Computer systems organization → Heterogeneous (hybrid)
37 systems, Computer systems organization → Embedded software, Computer systems organization
38 → Real-time systems, Computer systems organization → Embedded and cyber-physical systems,
39 Software and its engineering → Scheduling, Software and its engineering → Concurrency control,
40 Software and its engineering → Process synchronization

* Work supported by NSF grants CNS 1409175, CPS 1446631, CNS 1563845, and CNS 1717589, ARO grant W911NF-17-1-0294, and funding from General Motors.



© Ming Yang, Nathan Otterness, Tanya Amert, Joshua Bakita, James H. Anderson, and F. Donelson Smith;

licensed under Creative Commons License CC-BY

30th Euromicro Conference on Real-Time Systems (ECRTS 2018).

Editor: Sebastian Altmeyer; Article No. 20; pp. 20:1–20:21



Leibniz International Proceedings in Informatics

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

41 **Keywords and phrases** real-time systems, graphics processing units, scheduling algorithms, par-
42 allel computing, embedded software

43 **Digital Object Identifier** 10.4230/LIPIcs.ECRTS.2018.20

44 **1** Introduction

45 A fundamental shift is reshaping how real-time analysis is applied in all forms of autonomous
46 systems (*e.g.*, UAVs, robotics, and, especially, self-driving automobiles). These systems are
47 increasingly dependent on escalating computational requirements for various applications
48 based on machine learning (ML). Examples include computer-vision applications that rec-
49 ognize people and objects in high-bit-rate streams from multiple video cameras, and applica-
50 tions that process 3-D models of the surrounding environment from high-volume streams of
51 LIDAR data. These and other ML applications in autonomous vehicles have prompted the
52 adoption of specialized computing accelerators to match computational demands. Graphics
53 processing units (GPUs) are among the most prominent and accessible of these specialized
54 accelerators because of their high-throughput performance. While high throughput is neces-
55 sary for ML applications based on multiple streams of sensor inputs, it alone is not sufficient.
56 *Safe operation of autonomous vehicles also requires temporal correctness from GPU-using*
57 *tasks*—this is where real-time analysis becomes essential for autonomous systems.

58 **Why there is a problem.** Unfortunately, GPUs present many challenges, so model-
59 ing, analyzing, and certifying a safety-critical autonomous system using GPUs is currently
60 beyond the state-of-the-art. One reason is that GPUs are fundamentally different from
61 CPUs. Real-time analysis is based on well-understood scheduling algorithms that allocate
62 CPU capacity. In contrast, GPU hardware and software together implement GPU-specific
63 scheduling algorithms that are proprietary, opaque, and can change without notice. Model-
64 ing and analysis efforts under these conditions are subject to many pitfalls when applied to
65 real-time safety-critical workloads in GPU-using autonomous systems.

66 **Focus of this paper.** Our motivation for this work is to provide guidance, recommenda-
67 tions, and warnings about numerous pitfalls to both research and implementation practi-
68 tioners. We have found that writing programs for real-time tasks that combine CPU and
69 GPU computations is harder than we first thought. Based on several years of study, exper-
70 imentation, and experience with GPU programming, we are presenting here a compendium
71 of specific issues that are essential background for developing task systems where real-time
72 design meets GPUs.

73 **Choice of GPU platforms.** We base our findings on our experiences with NVIDIA GPUs
74 for a number of reasons. The most salient reason is that *NVIDIA GPUs are in cars on the*
75 *road today*. Further, NVIDIA has positioned itself as a market leader in automotive appli-
76 cations. For example, NVIDIA’s “Jetson” line of embedded platforms specifically targets
77 autonomous systems, and is marketed as “the embedded platform for autonomous every-
78 thing” [21]. Three generations of the Jetson series of embedded single-board computers
79 have been produced by NVIDIA; the TK1, TX1, and TX2. NVIDIA also markets a higher-
80 performance line of embedded platforms, the “Drive PX” series, which includes multiple
81 models such as the Drive PX2, Drive PX Xavier, and Drive PX Pegasus.

82 NVIDIA GPUs serve as an exemplar of the push for throughput over predictability in
83 GPUs. Recent developments in the NVIDIA GPU ecosystem are focused on improving ML
84 applications, especially those for autonomous driving. Most of these improvements center
85 around increasing throughput or reducing execution latency, but little, if any, attention

86 has been paid to requirements of the real-time tasks used in autonomous systems. This
87 lack of attention is evident in the sparse efforts by NVIDIA to improve or document GPU
88 scheduling behavior or improve the predictability of GPU execution times.

89 1.1 Contributions

90 The major contribution of this paper lies in discussing pitfalls for real-time GPU usage of
91 relevance to both those conducting research on autonomous systems and those who design
92 and build them. These pitfalls fall within three categories:

93 **Synchronization and blocking.** In any task consisting of a combination of CPU and GPU
94 computations, there are necessary synchronization points (*e.g.*, a CPU program needs to wait
95 until a GPU has produced a result). Synchronization inherently leads to blocking terms in
96 scheduling analysis. Unfortunately, we have learned that why and when synchronization
97 blocking occurs in a GPU-using task is not straightforward to determine. Further, some
98 forms of synchronization can lead to significant capacity loss on both CPUs and GPUs.
99 We have constructed experiments that expose these synchronization effects and carefully
100 describe them along with a list of specific pitfalls the unwary programmer may encounter.
101 This contribution is fully presented in Sec. 3.

102 **GPU concurrency.** We have realized that there is a fundamental trade-off that exists for
103 designing real-time tasks that use a GPU. A conventional choice is to write and execute
104 the task program as an operating system (OS) process in its own non-shared address space.
105 This provides cross-task memory isolation. If this choice is used, however, the NVIDIA
106 GPU programming environment (described in Sec. 2) does not permit any concurrent com-
107 putations on the GPU even if sufficient GPU resources are available. Depending on how
108 GPU programs are organized and written, this can lead to capacity loss on the GPU. The
109 alternate choice is to write and execute a task as a schedulable thread that shares a process
110 address space with other task threads. Cross-task memory isolation is lost, but the GPU
111 programming environment provides mechanisms that allow concurrent computations on the
112 GPU. NVIDIA provides a third option with a middleware environment that is claimed to
113 provide the best of both choices—memory isolation with concurrency enabled. We have
114 performed a case study using algorithms that are exemplars for computer-vision tasks in
115 autonomous vehicles to evaluate these trade-off options. The results and guidelines are fully
116 presented in Sec. 4.

117 **CUDA programming perils.** Our research has necessarily involved constructing many
118 thousands of lines of GPU programming for performing experiments. This experience has
119 been especially enlightening about the perils one can encounter in programming for NVIDIA
120 GPUs. The perils span a spectrum of pain ranging from simple documentation errors to
121 functions that default in strange ways, to programming “gotchas.” We present a list of perils
122 with descriptions and examples of the ones most likely to cause problems in Sec. 5.

123 **Value for autonomous systems.** We believe that this paper will help bridge the gap
124 between research and implementation in autonomous systems. For example, real-time re-
125 searchers may not be familiar with GPU programming for applications of ML and other
126 forms of AI used in real-time tasks. Likewise, programmers responsible for implementations
127 are given little guidance about creating GPU-using task systems amenable to real-time anal-
128 ysis. We provide the necessary understanding required to apply GPUs in real-time tasks
129 while avoiding numerous hidden pitfalls. We also expose GPU-related issues that must be
130 mitigated for real-time guarantees to be possible in autonomous systems. We further believe
131 that the fundamental issues presented herein are relevant to any real-time application using
132 computational accelerators, and likely hold for other manufacturers’ GPUs, digital signal

133 processors (DSPs), or FPGAs.

134 **1.2 Related Work**

135 Treating GPUs as non-shared devices has been a consistent theme in much of the prior
136 research on GPU scheduling for real-time systems. More predictable execution times result
137 from restricting access to the entire GPU (or its independent execution and data movement
138 components) to a single task at a time [11, 15, 16, 27, 28, 29, 31].

139 Other prior research takes a slightly different approach and improves schedulability by
140 simulating preemptive execution [3, 15, 17, 33]. These designs typically split GPU compu-
141 tations into smaller fragments, which can be individually scheduled and preempted. One
142 of these frameworks, called Kernelet [32], even allows GPU sharing as a means to improve
143 utilization, but interference effects caused by sharing are not addressed.

144 The decision to treat GPUs as non-shared devices is largely motivated by a perceived
145 need to work with a greatly simplified *model* of GPU execution (resulting, we believe, pri-
146 marily from a lack of information from GPU manufacturers). If GPU scheduling behavior
147 is an opaque “black box,” it is a rational conclusion that sharing must be avoided because
148 execution ordering and interference effects *cannot be known*. Our research is motivated,
149 however, by an observation that GPU sharing will become essential for effectively utilizing
150 less-capable embedded GPUs. Our research goal is to enable the modeling and analysis of a
151 combined CPU+GPU scheduling framework that allows real-time tasks to share multicore
152 CPUs and one or more GPUs.

153 We began our research by experimentally investigating the impacts of GPU sharing on
154 the NVIDIA Jetson TK1 [24] and TX1 [25]. In these studies, we focused on GPU sharing
155 by CPU processes (tasks) that have separate address spaces. We found that sharing in this
156 context happens only through round-robin time-sliced multiplexing of GPU computations
157 onto the GPU execution hardware. This multiplexing form of scheduling presents many
158 challenges for modeling and analysis. In later work, we experimentally investigated GPU
159 sharing by CPU tasks that share an address space (threads) on both the TX1 [26] and the
160 more-capable TX2 [1]. In these studies, we found that truly concurrent sharing can indeed
161 occur and deduced rules the GPU uses to schedule execution.

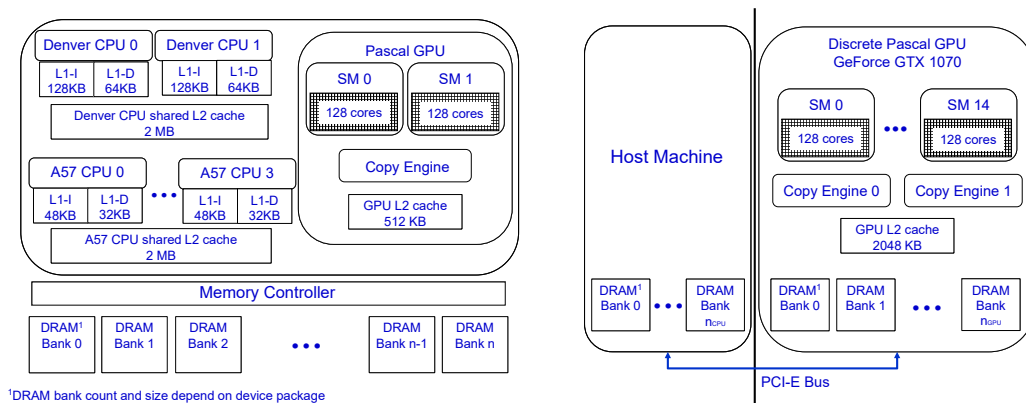
162 The work summarized so far was all directed at scheduling real-time tasks that use a
163 GPU for parts of their executions. Other work has focused on timing analysis for GPU
164 workloads [4, 5, 6, 7, 8], techniques for remedying performance bottlenecks [13], direct I/O
165 communication [2], and techniques for managing or evaluating GPU hardware resources,
166 including the cache and DRAM [9, 10, 12, 14, 18, 19, 30].

167 **2 Background**

168 In this section, we provide background information on the NVIDIA GPUs used in this
169 research. The CUDA programming framework is described, and a simple example of a
170 CUDA program is explained.

171 **2.1 CUDA-Enabled Devices**

172 The work presented here refers to the Kepler, Maxwell, Pascal, and Volta architectures of
173 NVIDIA GPUs. NVIDIA introduced these four different generations of GPU architectures,
174 in that order, within a time span of about five years (2012 - 2017)—a pace of change more



■ **Figure 1** Jetson TX2 Architecture (left) and GeForce GTX 1070 Architecture (right)

175 rapid than normally seen in CPU generations. GPUs are programmed using the CUDA
 176 API, which is an NVIDIA-provided set of libraries and language extensions for C/C++.

177 We consider both *discrete* GPUs and *integrated* GPUs. An integrated GPU, such as
 178 the NVIDIA Jetson TX2 shown on the left in Fig. 1, is part of a System-On-Chip (SoC)
 179 implementation combined with conventional multicore CPUs. The SoC is packaged along
 180 with DRAM and external connectors as a small (approximately 7 inches square) single-
 181 board computer. The integrated GPU shares hardware resources, such as DRAM, with
 182 CPU cores. The TX2 runs the Linux operating system, with additional support from closed-
 183 source binary drivers provided by NVIDIA. The TX2's low size, weight, and power (SWaP)
 184 requirements and low price tag make it a good exemplar of GPU-enabled platforms intended
 185 for embedding in autonomous systems.

186 Fig. 1 (left) shows the high-level architecture of the TX2. The TX2 contains a six-core
 187 heterogeneous ARMv8 CPU, 8 GB of DRAM, and an integrated Pascal GPU. The TX2's
 188 GPU consists of two *streaming multiprocessors* (SMs), each comprised of 128 GPU cores.
 189 The SMs together can be logically viewed as an *execution engine* (EE). Additionally, there
 190 is a hardware *copy engine* (CE) that can copy data between memory regions allocated for
 191 CPU use and those allocated for GPU use. The integrated GPU has fewer GPU cores than
 192 found in typical high-end GPUs used for graphics, gaming, and high-performance computing
 193 applications. We are interested in exploiting any potential for sharing the TX2's GPU by
 194 multiple tasks so that its computing capacity is not unnecessarily wasted.

195 Shown on the right in Fig. 1 is the architecture of the GTX 1070, an example of a discrete
 196 GPU. Discrete GPUs consist only of the SMs and local device memory, typically packaged
 197 on an adapter card for mounting in a PCIe expansion slot on a computer motherboard. Like
 198 all discrete GPUs, the GTX 1070 does not share memory with the host CPU, instead using
 199 the PCIe bus to copy data to and from host memory. This GPU features many more SMs
 200 than the TX2, increasing the potential benefit attainable if shared among multiple tasks. It
 201 also has two CEs, and a larger cache.

Algorithm 1 Vector Addition Pseudocode.

```

1: kernel VECADD(A ptr to int, B: ptr to int, C: ptr to int)
   ▷ Calculate index based on built-in thread and block information
2:   i := blockDim.x * blockIdx.x + threadIdx.x
3:   C[i] := A[i] + B[i]
4: end kernel

5: procedure MAIN
   ▷ (i) Allocate GPU memory for arrays A, B, and C
6:   cudaMalloc(d_A)
7:   ...
   ▷ (ii) Copy data from CPU to GPU memory for arrays A and B
8:   cudaMemcpy(d_A, h_A)
9:   ...
   ▷ (iii) Launch the kernel
10:  vecAdd<<<numBlocks, threadsPerBlock>>>(d_A, d_B, d_C)
   ▷ (iv) Copy results from GPU to CPU array C
11:  cudaMemcpy(h_C, d_C)
   ▷ (v) Free GPU memory for arrays A, B, and C
12:  cudaFree(d_A)
13:  ...

```

202 **2.2 Relevant CUDA Programming Fundamentals**

203 A *CUDA program* runs as a task (process or thread) on a CPU and relies on a GPU for
204 some part of its computational requirements.¹ The general structure of a CUDA program
205 when it needs to interact with the GPU is as follows: **(i)** allocate memory for GPU use;
206 **(ii)** copy input data from CPU memory to GPU memory; **(iii)** launch execution of a GPU
207 program called a *kernel*² to process the data; **(iv)** copy the results from the GPU memory
208 back to the CPU memory; **(v)** free unneeded memory.

209 CUDA kernels are written from the perspective of a single *GPU thread*. As an example,
210 consider the CUDA program expressed in pseudocode in Algorithm 1. It uses the kernel
211 VECADD to add a single pair of elements per GPU thread, storing the sum in a corresponding
212 location in an output array. Line 2 demonstrates the use of special global system-defined
213 variables to determine the array element on which to operate. When the kernel executes,
214 threads will run in lock-step with each thread performing the same operation simultaneously
215 on different data. To avoid confusion with GPU threads, we will henceforth refer to CPU
216 threads as *CPU tasks* (or just *tasks*).

217 A kernel is run on the GPU as a set of thread blocks that can be executed in any order.
218 These thread blocks, or simply *blocks*, are each comprised of a number of threads. As seen
219 in Line 10 of Algorithm 1, the number of blocks and threads per block are programmer-
220 specified and can be set at runtime when a kernel is launched. The GPU scheduler uses
221 these values to assign work to the SMs. *Blocks are the schedulable entities on the GPU*. All
222 threads in a block are always executed on the same SM, and run non-preemptively until
223 completion. A kernel completes when all threads in all blocks have exited.

224 We refer to kernels and memory-copy operations collectively as *GPU operations*. GPU
225 operations are submitted to a GPU in *CUDA streams*. Operations within a stream are exe-
226 cuted in FIFO order. By default, the *NULL stream* is used, but users can submit operations
227 to multiple user-defined streams.³ Kernels from different streams can run concurrently by

¹ Note that both CPU and GPU computations are specified in the same CUDA program.

² Unfortunate terminology, not to be confused with an OS kernel.

³ CUDA documentation only guarantees that operations within a stream are executed in FIFO order,

228 sharing the GPU’s cores if sufficient internal resources are available. Copy operations are
229 handled by the GPU’s CE and can be concurrent with kernel executions on the EE.

230 CUDA API calls can be synchronous or asynchronous; for many calls, a variant of both
231 is available. For example, `cudaMemcpy` and `cudaMemcpyAsync` both copy data between re-
232 gions of CPU memory and GPU memory, or between two regions of GPU memory, but
233 `cudaMemcpyAsync` can return control to the calling CPU task before the copy is completed,
234 whereas `cudaMemcpy` blocks the CPU task until the memory copy completes.

235 Kernel launches are always supposed to be asynchronous. The CUDA documenta-
236 tion⁴ [23], however, uses a narrow definition of “asynchronous” that can be misleading.
237 According to the documentation, “asynchronous library functions that return control to the
238 host thread before the device completes the requested task.” Notably, this definition does
239 not imply that asynchronous API calls are *nonblocking* to the CPU. As noted in Sec. 3, we
240 have found situations in which kernel launches still cause CPU blocking even if the API call
241 returns before the requested kernel completes.

242 **3 Synchronization and Blocking**

243 CPU scheduling has been studied and well-understood for decades; in particular, real-time
244 scheduling analysis of task systems is based on predictable scheduler and task behaviors. A
245 worst-case execution time (WCET) for each task can be determined using clear specifications
246 of the machine’s architecture including the cache, bus, and DRAM operations. Incorporating
247 GPUs into real-time analysis (as with all coprocessors), requires different models with new
248 sets of issues to be considered. In this section, we discuss one set of issues that lead to a
249 surprising number of pitfalls when CUDA GPUs are used: *synchronization*.

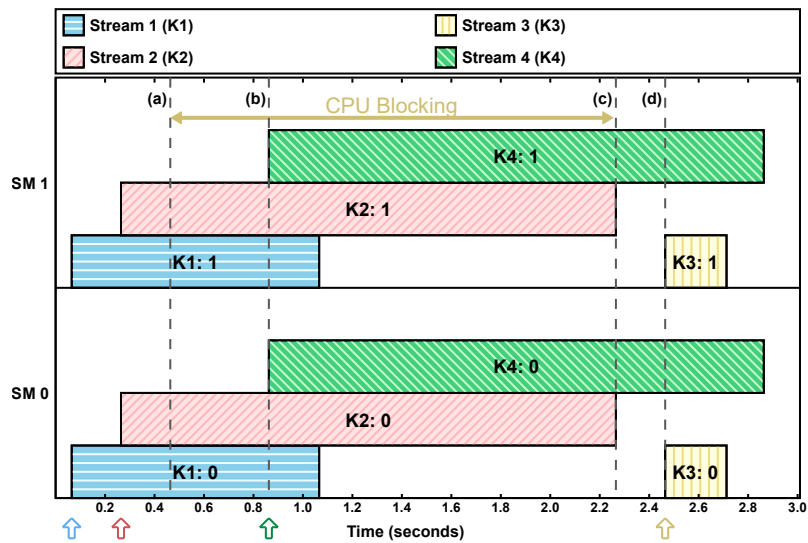
250 In prior work, we investigated the scheduling rules for kernels and copy operations in
251 CUDA programs [1]. However, this investigation focused on a limited context where few
252 CUDA operations beyond kernel launches and memory copies were used. In most real-world
253 CUDA software, programmers will likely encounter (both intentionally and unintentionally)
254 the need for synchronization between CPU and GPU operations. The added complexity
255 of synchronization can result in utilization loss, potentially leading to unbounded response
256 times in task sets with high utilization. In this section, we explore various forms of CPU-GPU
257 synchronization and the resulting implications for real-time systems. We limit attention for
258 now to CPU tasks that share a single Linux address space and create user-defined streams.
259 As covered in detail in Sec. 4, this setup allows potential concurrency among operations on
260 the GPU.

261 **3.1 Overview of GPU Synchronization**

262 Most developers are familiar with the concepts of synchronization in a CPU-only context
263 where two or more tasks must communicate or coordinate their actions. Synchronization
264 becomes more complicated when a CPU task must coordinate with programs executed on
265 the GPU. The common case is that the CPU task must determine when data in GPU mem-
266 ory is safe to access (*e.g.*, copy back to CPU memory). This is accomplished using *GPU*
267 *synchronization*, where the GPU must complete outstanding work and reach a *synchroniza-*
268 *tion point*: a point in time when data access can safely occur. There are also other, less

but does not describe how operations from different streams are ordered.

⁴ Specifically, Section 3.2.5.1 of the Programming Guide for CUDA version 9.1.85.



■ **Figure 2** Explicit synchronization requested before K3, observed on the Jetson TX2.

269 common, cases when GPU synchronization is necessary.

270 In CUDA there are multiple ways to achieve GPU synchronization. They fall into two
 271 broad categories: *explicit synchronization*, which is always programmer-requested, and *im-*
 272 *PLICIT synchronization*, which can occur as a side effect of CUDA API functions intended
 273 for purposes other than synchronization. We have uncovered in our research some unfor-
 274 tunate pitfalls relating to actual GPU synchronization behavior, especially with respect to
 275 *blocking*. So, while these may not be pitfalls for non-safety-critical applications, ignoring the
 276 effects of certain specific mechanisms for achieving synchronization would be perilous in a
 277 safety-critical system where blocking must be anticipated and accounted for in analysis.

278 3.1.1 Explicit Synchronization

279 Explicit synchronization refers to synchronization points that the CUDA programmer ex-
 280 plicitly requests using the CUDA API. Explicit synchronization is typically used after a
 281 program has launched one or more asynchronous CUDA kernels or memory-transfer opera-
 282 tions and must wait for computations to complete. In contrast to implicit synchronization,
 283 the sole purpose of explicit-synchronization functions is to block the calling CPU task
 284 until the GPU reaches a synchronization point.

285 The CUDA documentation⁵ states that explicit synchronization will block the calling
 286 task until “all preceding commands” have completed. For example, if the API function
 287 `cudaDeviceSynchronize` is invoked, “preceding commands” may encompass all commands
 288 issued to the device from all CPU tasks. Other explicit-synchronization options, including
 289 `cudaStreamSynchronize`, will only block until preceding commands from a specified stream
 290 have completed.

291 We carried out experiments using our open-source framework⁶ to investigate the specific
 292 behaviors of GPU synchronization on real GPU hardware. Fig. 2 shows the behavior of
 293 explicit synchronization observed in one such experiment. In Fig. 2 (also in Figs. 3 and 4),

⁵ Section 3.2.5.5.3 of the Programming Guide for CUDA version 9.1.85.

⁶ Available at https://github.com/yalue/cuda_scheduling_examiner_mirror.

294 each shaded rectangle corresponds to a separate thread block. The left and right endpoints of
295 each rectangle correspond to the times at which the block started and completed execution,
296 as measured on the GPU. Each rectangle's height represents its size in CUDA threads.
297 Additionally, the vertical axis is subdivided by SM. The particular experiments presented
298 in Figs. 2-4 were performed using the Jetson TX2, which features two SMs. Up to 2,048
299 CUDA threads can be assigned to a single SM at once.

300 The CUDA program executed to produce Fig. 2 consists of four CPU tasks all sharing a
301 single address space. Each CPU task launched one kernel in a separate user-defined stream.
302 Kernel launches were separated by a small amount of time. Each kernel consisted of two
303 blocks of 512 threads, and the figure shows that one block from each kernel was scheduled
304 on each SM. Each thread performed a busy-loop for a set amount of time.

305 An explicit-synchronization command, `cudaDeviceSynchronize`, was issued at time **(a)**
306 by the CPU task responsible for launching kernel K3. This caused K3's CPU task to be
307 blocked until the prior commands, the execution of kernels K1 and K2, had both completed
308 at time **(c)**. This behavior is exactly what one would expect, given the description of ex-
309 plicit synchronization from official documentation. However, our experiments also uncovered
310 Pitfall 1 for the unwary:

311 ► **Pitfall 1.** *Explicit synchronization does not block future commands issued by other tasks.*

312
313 The fact that the launch of K4 by its CPU task was not blocked at time **(b)** is an example
314 of this pitfall. Implicit synchronization, which we cover next, presents even more serious
315 pitfalls.

316 3.1.2 Implicit Synchronization

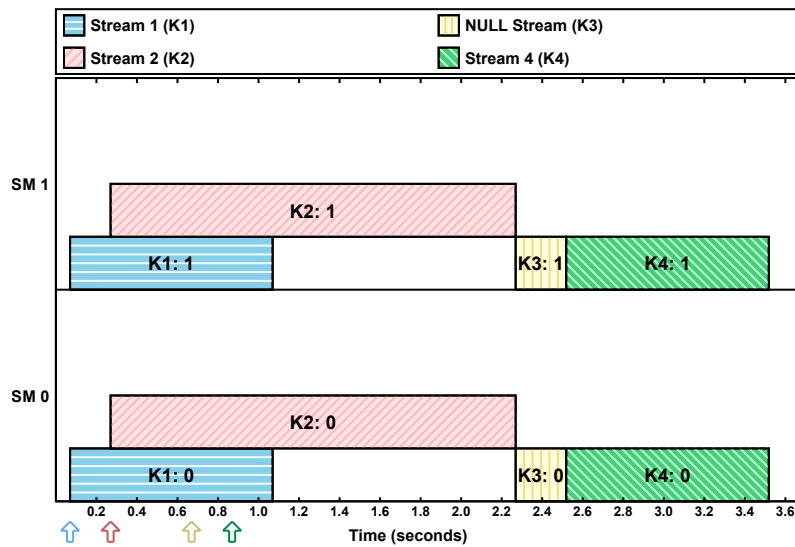
317 Implicit synchronization occurs as a side effect of CUDA API calls that are otherwise un-
318 related to synchronization. For example, implicit GPU synchronization may occur due to
319 freeing GPU memory or launching a kernel to the default stream. Presumably, this is because
320 some modifications to GPU device state can only occur while no kernels are executing. The
321 CUDA documentation about implicit synchronization⁷ states that “two commands from
322 different streams cannot run concurrently if any one of the following operations is issued
323 in-between them by the host thread:

- 324 1. A page-locked host memory allocation
- 325 2. A device memory allocation
- 326 3. A device memory set
- 327 4. A memory copy between two addresses to the same device memory
- 328 5. Any CUDA command to the NULL stream”

329 Unlike the relatively straightforward documentation about explicit synchronization, our ex-
330 periments revealed that this list includes several operations that do not necessarily cause
331 implicit synchronization, and fails to include some functions that do. We consider this par-
332 ticularly problematic for real-time systems, where the ability to accurately model blocking
333 is critical.

334 ► **Pitfall 2.** *Documented sources of implicit synchronization may not occur.*

⁷ Section 3.2.5.5.4 of the Programming Guide for CUDA version 9.1.85.



■ **Figure 3** Implicit synchronization caused by launching kernel K3 in the NULL stream.

335 Pitfall 2 became apparent to us when, in all of our experiments, we never observed
 336 implicit synchronization as a result of a device-memory operation (allocation, set, or copy) or
 337 a page-locked host memory allocation. Our experiments covered the two most recent CUDA
 338 versions, 8.0 and 9.0, and the three most recent NVIDIA GPU architectures, Maxwell,
 339 Pascal, and Volta. This, of course, does not prove that implicit synchronization can *never*
 340 happen under such circumstances, but it does indicate that the documentation’s statement
 341 that “two commands cannot run concurrently” is not a reliable rule. The only case (from
 342 this list) in which we did observe implicit synchronization was launching GPU operations in
 343 the NULL stream.

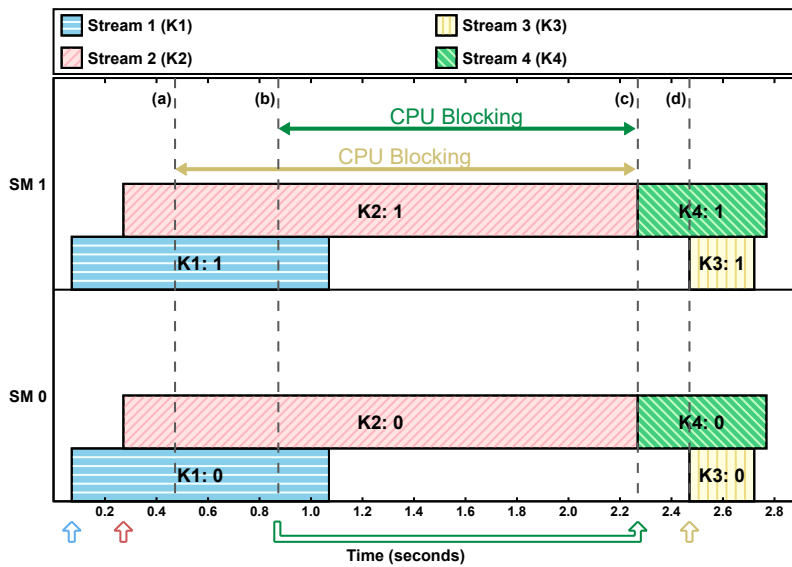
344 Fig. 3 shows a similar scenario to the one in Fig. 2, with one key difference: the CPU task
 345 for K3 did not call `cudaDeviceSynchronize` before K3 was launched, but instead launched
 346 K3 in the NULL stream. The implicit synchronization, and resulting loss of concurrency, is
 347 clearly visible in the figure. Execution of K3 must wait for the first two kernels to complete,
 348 and, in contrast to explicit synchronization, K4 is also prevented from running concurrently.
 349 Even though this loss of concurrency may be striking, it is notably explicitly documented,
 350 and can be used (or avoided) in a careful design for a real-time task.

351 We found, however, a different source of implicit synchronization that is a far more
 352 problematic pitfall, and is not even listed in the documentation on synchronization: *freeing*
 353 *device memory*.

354 ► **Pitfall 3.** *The CUDA documentation neglects to list some functions that cause implicit*
 355 *synchronization.*

356 ► **Pitfall 4.** *Some CUDA API functions will block future, unrelated, CUDA tasks on the*
 357 *CPU.*

358 Fig. 4 shows the results of an experiment identical to the one in Fig. 2, but this time
 359 the call to `cudaDeviceSynchronize` at time (a) was replaced with a call to `cudaFree`,
 360 which was used to de-allocate memory on the GPU. Pitfalls 3 and 4 can be observed in
 361 this plot. The fact that this blocked the calling CPU thread until all prior GPU work had
 362 completed at time (c) indicates that `cudaFree` created implicit synchronization. Similar



■ **Figure 4** Implicit synchronization causing additional CPU blocking due to `cudaFree`.

363 to the NULL-stream behavior, implicit synchronization also prevented subsequent kernels
 364 from starting to execute until `cudaFree` completed at time (c). We speculate that this
 365 behavior by `cudaFree` is necessary because alterations to memory-mapping state requires a
 366 quiescent execution environment. However, the most surprising effect was not that K4 was
 367 blocked, but that K4’s task was blocked *on the CPU* until time (c), even though it issued an
 368 “asynchronous” kernel launch. This reveals a pitfall that can harm real-time analysis that
 369 does not consider the fact that CPU tasks can experience blocking from GPU operations
 370 that are launched from unrelated tasks.

371 3.2 Overcoming Synchronization-Related Pitfalls

372 GPU synchronization has two problematic effects—introducing indeterminate amounts of
 373 blocking and reducing GPU concurrency. This means that programmers who develop real-
 374 time systems must understand the pitfalls inherent in explicit and implicit synchronization.
 375 This is especially true if the schedulability of a real-time task system relies on minimizing
 376 blocking or high GPU utilization. Avoiding pitfalls can be accomplished through careful
 377 construction of CUDA programs to, for example, avoid using the NULL stream or free-
 378 ing memory outside of certain time intervals. A more robust method would be to adopt
 379 middleware that handles such problems transparently.

380 Our experiments indicate that GPU synchronization does not extend across GPU-using
 381 tasks that are isolated in separate address spaces. If synchronization is the dominant limiting
 382 factor on schedulability, it may be desirable to place each task in a separate address space
 383 (OS process). As explained in the next section, this organization means that that CUDA
 384 kernels from different tasks can no longer execute concurrently, but it may still be beneficial
 385 overall if synchronization-related blocking is a greater limiting factor.

386 It turns out that NVIDIA may be aware of this issue. Even though it is not currently
 387 available for embedded platforms such as the TX2, NVIDIA does provide useful middleware
 388 for discrete GPUs: the *CUDA Multi-Process Service* (MPS). MPS allows kernels from mul-
 389 tiple processes to execute concurrently on a single GPU, while maintaining the desirable

390 property that GPU synchronization from one process will not affect other processes. We
 391 explore the benefits of MPS further in Sec. 4.

392 **4 Concurrency and Performance**

393 In prior work, we investigated different GPU scheduling behavior when running GPU-using
 394 real-time task systems in two contexts: **(i)** each task has its own distinct address space,
 395 *i.e.*, it runs as an OS process, and **(ii)** all tasks belong to the same address space, *i.e.*,
 396 each task runs as a schedulable thread within a process. We refer to these two contexts as
 397 *process-based* and *thread-based* tasks, respectively.

398 While process-based tasks have the advantage of memory protection, they do not actually
 399 execute on the GPU concurrently; instead, GPU operations are multiprogrammed in a way
 400 that makes predictable scheduling of GPU-related resources difficult if not impossible to
 401 achieve [1]. When operations are multiprogrammed on a GPU, their execution times depend
 402 on contention for shared GPU resources, making it hard to bound a task’s overall execution
 403 time. Additionally, concurrency among GPU operations may be important in order to
 404 avoid wasting GPU processing cycles, especially when a single kernel cannot fully utilize the
 405 GPU’s resources. Although this may be avoided by running tasks with user-defined streams
 406 in a shared address space, a shared address space may actually *reduce* concurrency in task
 407 systems where tasks regularly interfere with each other via implicit synchronization (Sec. 3).
 408 Fortunately, NVIDIA provides a third option: middleware called the *Multi-Process Service*
 409 (*MPS*) [20].

410 **4.1 Multi-Process Service (MPS)**

411 MPS enables concurrent execution of GPU operations launched by independent CPU address
 412 spaces. It has the potential to combine the advantages of both thread- and process-based
 413 tasks. Programs written using the CUDA API require no changes to use MPS—if MPS
 414 is running, CUDA programs transparently issue requests to MPS rather than directly to a
 415 GPU. Official documentation reports that MPS operates as a server process with its own
 416 CUDA context, and that CUDA API requests are redirected from client processes to the
 417 MPS server. Because the server’s CUDA context is effectively shared, GPU operations
 418 launched by separate processes can execute concurrently on a shared GPU, providing the
 419 benefits of thread-based tasks. However, MPS also continues to preserve the advantage of
 420 process-based tasks: separate processes will not block each other with implicit or explicit
 421 synchronization.

422 It is not clear from available documentation how MPS actually schedules GPU operations
 423 and whether the GPU scheduling rules revealed in prior work [1] are followed under MPS.
 424 For example, the documentation for MPS only mentions possible overlap between kernels
 425 and copy operations.⁸ Given the documentation flaws discussed in Secs. 3 and 5.2, one could
 426 be skeptical of the veracity of this claim, so we verified experimentally that those scheduling
 427 rules are also followed under MPS. We omit from this paper the experimental methods used
 428 for verifying the scheduling rules; readers can refer to [1].

429 Maximizing the utilization of GPU resources using streams in thread-based tasks is sug-
 430 gested by NVIDIA’s “Best Practices Guide” [22]. However, it would be unwise to simply take

⁸ “MPS allows kernel and memory copy operations from different processes to overlap on the GPU, achieving higher utilization and shorter running times” [20].

	Multiple Process-based Tasks	Multiple Thread-based Tasks
Without MPS	MP	MT
With MPS	MP(MPS)	MT(MPS)

■ **Table 1** Abbreviations used for our four experimental scenarios.

431 this recommendation at face value when choosing between MPS or a process- or thread-based
 432 task organization in a safety-critical system. Additionally, *MPS is not yet supported on*
 433 *embedded ARM platforms* like the Jetson TX2, so the other management systems are
 434 still necessary on some systems. Therefore, we conducted a case study on computer-vision
 435 software, demonstrating the performance differences among the available configurations.

436 4.2 Case Study of Computer-Vision Tasks

437 Our motivation primarily remains autonomous driving, so we chose to study algorithms for
 438 computer-vision tasks that provide functions commonly used for autonomous driving. In
 439 evaluating the results from this case study, we consider that the real-time tasks that use
 440 GPUs for autonomous driving may have multiple levels of criticality. Some may be safety-
 441 critical with hard deadlines and be provisioned for worst-case execution plus a margin for
 442 safety. Others may have only bounded tardiness requirements, or even be background work
 443 that can be provisioned for average-case execution.

444 We focus here on five programs from NVIDIA’s provided sample code for VisionWorks:

- 445 ■ **Video Stabilization.** Smooths shaky video content. This is often a preprocessing step
 446 for a computer-vision pipeline.
- 447 ■ **Feature Tracking.** Tracks features between consecutive frames. This algorithm is used
 448 to track the positions of objects in a scene.
- 449 ■ **Motion Estimation.** Estimates the direction of moving pixels, which is fundamental
 450 to calculating trajectories of moving objects, *e.g.*, pedestrians and other vehicles.
- 451 ■ **Hough Transform.** (*Hough*) A feature-extraction algorithm; the provided sample de-
 452 tects circles and lines in images.
- 453 ■ **Stereo Matching.** Uses input from two cameras to generate depth information by
 454 matching features in both frames.

455 **Methodology.** We adapted NVIDIA’s VisionWorks samples to be compatible with our
 456 open-source experimental framework.⁹ These samples generally only use a single CUDA
 457 stream. We ran four instances of the same sample program in each experiment. We config-
 458 ured each instance to process 1,000 frames from a video sequence while recording per-frame
 459 response times. Our framework allows running each program instance in a shared address
 460 space (multiple thread-based tasks, MT) or in independent address spaces (multiple process-
 461 based tasks, MP), both with and without the MPS server active. This produces experiments
 462 for each algorithm in four different scenarios as summarized in Tbl. 1. Experimental results
 463 under MT(MPS) were always similar to MT with slight overheads caused by MPS, so we
 464 omit it in all of our results for clarity. We conducted these experiments on a Maxwell-
 465 architecture discrete GPU with CUDA 9.0. We briefly summarize results on other devices
 466 and different CUDA versions later.

⁹ Again, https://github.com/yalue/cuda_scheduling_examiner_mirror.

VisionWorks Samples	Scenarios	Max	99 th %	90 th %	Mean	Median
Video Stabilization	MP	17.55	12.88	5.43	3.31	2.69
	MP (MPS)	36.73	11.12	5.37	2.81	2.06
	MT	17.0	13.87	8.94	4.72	3.63
Feature Tracking	MP	5.64	3.87	1.45	1.08	0.96
	MP (MPS)	14.73	6.04	1.51	1.31	1.09
	MT	31.11	20.86	11.51	4.68	2.68
Motion Estimation	MP	28.64	21.25	17.33	16.75	17.24
	MP (MPS)	33.05	22.66	15.75	14.3	14.89
	MT	42.86	26.12	16.53	15.07	15.14
Hough Transform	MP	13.56	11.61	7.28	5.68	5.7
	MP (MPS)	18.35	11.66	6.44	3.74	3.18
	MT	58.65	22.64	15.82	9.12	8.94
Stereo Matching	MP	75.13	50.54	30.42	24.14	24.77
	MP (MPS)	59.73	45.05	26.87	22.59	24.41
	MT	125.96	58.82	34.36	20.75	18.95

■ **Table 2** Per-frame response time data (in milliseconds) of VisionWorks samples. The fastest scenario for each time metric is indicated by bold text.

467 **Results.** We show cumulative distribution function (CDF) and kernel density estimation
468 (KDE)¹⁰ plots of Hough and feature tracker as representatives in Figs. 5–8. The KDE curve
469 was produced using the Python package `scipy.stats.gaussian_kde`. In both the CDF
470 and KDE plots, each curve represents the recorded response-time data in an experimental
471 scenario. For example, the curves labeled “x4 MP” in Figs. 5 and 6 represent the per-frame
472 response time distributions where each of four Hough instances is run in a separate process.
473 Result data for all five algorithms is summarized in Tbl. 2, which lists the maximum, 99th-
474 percentile, 90th-percentile, mean, and median frame times for each scenario and algorithm.

475 ► **Observation 1.** *MP(MPS) exhibits good average-case performance.*

476 Obs. 1 is supported by the data in Tbl. 2. 90th-percentile, mean, and median performance
477 under configuration MP(MPS) were consistently good with the top performance for three
478 of the five algorithms. For Feature Tracking, MP was best in all metrics, and for Stereo
479 Matching, MT had better mean and median performance. The results for average-case
480 performance indicate that using MP(MPS) would likely be an attractive option for soft-real-
481 time systems, *e.g.*, systems that can occasionally drop a video frame without compromising
482 safety. We conjecture that the average-case performance advantage of MP(MPS) over MP
483 in most cases is due to improved concurrency and lower GPU context-switching overheads.

484 Feature Tracking was the most notable exception to Obs. 1. In this case, MP was
485 only slightly better than MP(MPS) when comparing the 90th-percentile, mean, and median
486 performance. We conducted additional experiments using NVIDIA’s CUDA-profiling tool,
487 `nvprof`, to gain some insight into this behavior. We found that Feature Tracking’s overall
488 execution time is heavily influenced by a large number of memory transfers, rather than

¹⁰ KDE is a statistical method for estimating a continuous probability density function (PDF) from a set of discrete sample values.

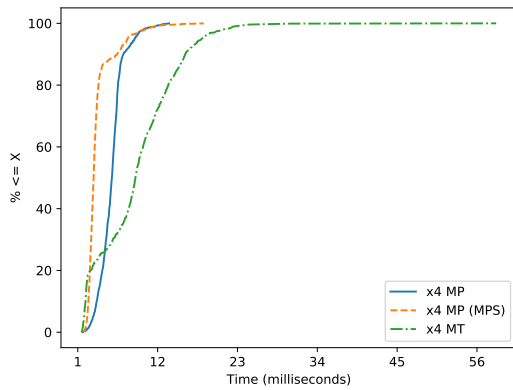


Figure 5 Per-frame response time CDFs for Hough.

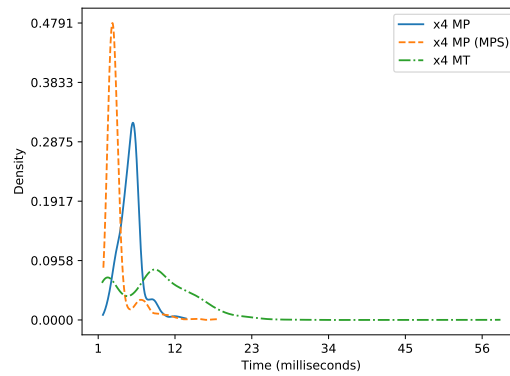


Figure 6 Per-frame response time KDEs for Hough.

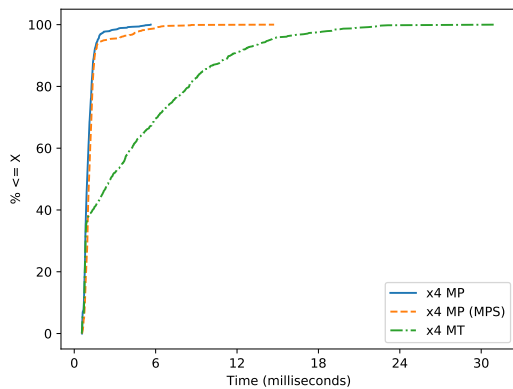


Figure 7 Per-frame response time CDFs for Feature Tracking.

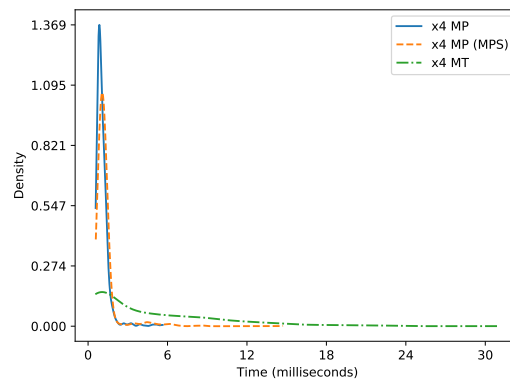


Figure 8 Per-frame response time KDEs for Feature Tracking.

489 CUDA kernel executions. This likely means that MPS only provides limited GPU concur-
 490 rency benefits to Feature Tracking, which failed to outweigh other MPS-related overheads.

491 ► **Observation 2.** *Worst-case and 99th-percentile runtimes were typically better under MP.*
 492

493 While MP(MPS) largely resulted in average-case improvements, Tbl. 2 shows three of
 494 our five applications (Feature Tracking, Motion Estimation, and Hough Transform) showed
 495 the smallest worst-case and 99th-percentile execution times under MP. This indicates that
 496 MP may be a better option for certain task systems where worst-case performance is more
 497 important than average-case. Our results illustrate why the trade-offs between process-based
 498 and thread-based designs for tasks must be evaluated for individual algorithms.

499 ► **Observation 3.** *MP and MP(MPS) exhibit more predictable execution times than MT.*

500 Obs. 3 is supported by Figs. 6 and 8, where the KDE shows a tight unimodal distribution
 501 for MP and MP(MPS) but not MT. A unimodal distribution function with little dispersion
 502 indicates that the response times exhibit low variance. MT, in contrast, shows both bimodal
 503 (in Fig. 6) and unimodal (in Fig. 8) distributions with significant dispersions (indicating high
 504 variance). Even if specific “spikes” are more difficult to observe in the corresponding CDF
 505 plots, the difference in response-time ranges are also apparent from the endpoints of the
 506 CDF curves in Figs. 5 and 7.

507 ► **Observation 4.** *The MT configuration generally performed poorly.*

508 Obs. 4 is supported by Tbl. 2 and the plots. The only metrics where MT outperformed
 509 the other scenarios were the mean and median times for Stereo Matching, and worst-case
 510 response time for Video Stabilization (where MT was only slightly better than MP).

511 **Other Results.** In addition to the results presented above, we also conducted this case
 512 study using CUDA 8.0 on a Maxwell discrete GPU (GTX 860M) and CUDA 9.0 on Pascal
 513 discrete GPUs (GTX 1050 and GTX 1070). Even though we chose to omit tables of results
 514 from the other GPUs and CUDA versions in this paper, we made similar observations ex-
 515 cepting that the performance of all configurations was better on a Pascal GPU. Additionally,
 516 the experimental results with CUDA 8.0 on the same Maxwell GPU stayed nearly identical
 517 to those using CUDA 9.0.

518 **Summary.** Our case study compared the impact of different GPU-sharing approaches on
 519 the performance of computer-vision algorithms. The results we obtained for these algorithms
 520 ran contrary to some of our observations regarding GPU concurrency from prior work [1, 26].

521 ► **Pitfall 5.** *The suggestion from NVIDIA’s documentation to exploit concurrency through
 522 user-defined streams may be of limited use for improving performance in thread-based tasks.*

524 We assumed that enabling concurrent GPU execution was of significant importance for
 525 limiting capacity loss in real-time workloads on embedded systems, and therefore fell victim
 526 to Pitfall 5. Instead, our results show that MT rarely outperforms tasks running as multiple
 527 processes, even without MPS. Additionally, any performance improvement via fine-tuned
 528 stream organization for MT can also be achieved with MP(MPS). That being said, even
 529 though enabling concurrency using MP(MPS) is generally beneficial, it unfortunately is not
 530 an option on ARM-based embedded platforms like the Jetson TX2. We would encourage
 531 NVIDIA to consider this shortcoming in hope that one day it may be addressed.

532 **5 Perils of CUDA Programming for Real-Time Tasks**

533 In the previous sections we presented several specific pitfalls in correctly designing and
 534 running CUDA programs for real-time tasks. Elements of both CUDA’s design and docu-
 535 mentation contribute to this ensemble of perils to avoid. In this section, we discuss some of
 536 the broader categories of pitfalls.

537 **5.1 Synchronous Defaults**

538 As hinted in Sec. 3, one of the primary pitfalls when designing a real-time task system
 539 that uses a GPU is that *all possible* blocking must be accounted for in analysis. Therefore,
 540 reducing the amount of blocking on both the CPU and GPU is essential. On the GPU, this
 541 requires issuing all CUDA operations to user-defined (non-NULL) streams, and carefully
 542 controlling the use of other API functions, like `cudaFree`, that cause blocking via implicit
 543 synchronization.

544 Even though it may seem like an easy task for a programmer to just specify a user-
 545 defined stream as opposed to the NULL stream, we note that simple mistakes in doing so
 546 may be easy to miss. This is particularly true when using the `Async` versions of CUDA API
 547 functions, such as `cudaMemsetAsync`. For example, consider the code snippets in Listings 1
 548 and 2, which present a particular example of Pitfall 6 below.

■ **Listing 1** Causes implicit synchronization.

```

549 if (!CheckCUDAError(cudaMemsetAsync(
    state->device_block_smids, 0,
    data_size))) {
    return 0;
}

```

■ **Listing 2** Correctly asynchronous.

```

if (!CheckCUDAError(cudaMemsetAsync(
    state->device_block_smids, 0,
    data_size, state->stream))) {
    return 0;
}

```

550 ► **Pitfall 6.** *Async CUDA functions use the GPU-synchronous NULL stream by default.*

551 Listing 1’s call to `cudaMemsetAsync` is missing a final argument specifying a user-defined
552 stream, which causes the NULL stream to be used by default. As pointed out in Sec. 3.1.2,
553 NULL-stream usage causes implicit synchronization and hence blocking. This mistake is
554 corrected in Listing 2. This specific mistake actually led to *months* of mystifyingly incon-
555 sistent results in our own experiments—despite our relatively deep experience examining
556 the subtleties of CUDA behavior (note that these code snippets are parts of much larger
557 listings). *Would an ML application developer catch such a mistake or appreciate its impact?*
558 Note that NVIDIA’s CUDA compiler does not catch this mistake because the compiler is
559 based on the C++ programming language, which allows default arguments to functions.

560 Even though the examples in Listings 1 and 2 only use `cudaMemsetAsync`, Pitfall 6
561 applies to other CUDA API functions as well, such as `cudaMemcpyAsync`. The fact that
562 the CUDA documentation indicates that these functions cause implicit synchronization, as
563 discussed in Sec. 3 and Sec. 5.2, makes potential programmer errors even harder to notice in
564 cases where synchronization is due to NULL-stream usage rather than memory operations.

565 To summarize this discussion, CUDA provides a brittle programming environment:
566 difficult-to-spot mistakes can have profound consequences for real-time tasks.

567 5.2 Flawed Documentation

568 Another substantial danger stems from the inaccurate official documentation provided by
569 NVIDIA. While function signatures and data structures seem to receive accurate (but often
570 sparse) official documentation, scheduling and synchronization remain under-discussed. Our
571 group’s past work includes demystifying some scheduling rules [1]. In our work to demystify
572 implicit synchronization (see definition in Sec. 3.1.2), however, we came across not only
573 missing documentation, but incorrect documentation.

574 ► **Pitfall 7.** *Observed CUDA behavior often diverges from what the documentation states or*
575 *implies.*

576 Consider Tbl. 3. In all but one of the cases we investigated, the documentation claims
577 implicit synchronization will occur when it does not. While this absence of synchronization
578 may positively benefit performance, it also may cause incorrect timing analysis. Further-
579 more, program logic may be broken in the (albeit unlikely) case that the program relies on
580 a function like `cudaMemsetAsync` to trigger GPU synchronization.

581 Unfortunately, the documentation also contains less-benign flaws. Take `cudaFree` and
582 `cudaFreeHost` as an example. Our experiments in Sec. 3 found these functions to not only
583 cause implicit synchronization, but block other CPU tasks from proceeding while `cudaFree`
584 waits on the GPU. Much to our surprise, the documentation mentions neither of these side
585 effects, leaving the reader to assume that these functions behave similarly to other CUDA
586 functions and have no side effects.

587 Our experiments also revealed that `cudaMalloc` and `cudaMallocHost` may also cause
588 cross-task CPU blocking in a similar manner to `cudaFree` in certain situations, even though

Source	Observed Behavior			Documented Behavior	
	Blocks Other CPU Tasks	Implicit Sync. (Sec. 3.1.2)	Caller Must Wait for GPU	Implicit Sync. (Sec. 3.1.2)	Caller Must Wait for GPU
<code>cudaDeviceSynchronize</code>	No	No	Yes	No	Yes
<code>cudaFree</code>	Yes	Yes	Yes	No (undoc.)	No (impl.)
<code>cudaFreeHost</code>	Yes	Yes	Yes	No (undoc.)	No (impl.)
<code>cudaMalloc</code>	?	No	No	Yes	No (impl.)
<code>cudaMallocHost</code>	?	No	No	Yes	No (impl.)
<code>cudaMemcpyAsync</code> D-D	No	No	No	Yes	No
<code>cudaMemcpyAsync</code> D-H	No	No	No	Yes *	No
<code>cudaMemcpyAsync</code> H-D	No	No	No	Yes *	No
<code>cudaMemset</code> (sync.)	No	Yes	No	Yes	No
<code>cudaMemsetAsync</code>	No	No	No	Yes	No
<code>cudaStreamSynchronize</code>	No	No	Yes	No	Yes

■ **Table 3** Observed vs. documented synchronization sources in CUDA. For `cudaMemcpyAsync` we distinguish the direction of copy between device and host: (D-D) internal to GPU memory; (D-H) GPU memory to CPU memory; (H-D) CPU memory to GPU memory. *The documentation is contradictory for these instances, but the more detailed option indicates that these functions only cause synchronization if host memory is not page-locked. We were unable to observe this regardless of whether host memory was page-locked or not.

589 these functions do not trigger implicit synchronization. As we have not yet determined the
590 specific causes for this behavior, this property is indicated by an entry of ‘?’ in certain cells
591 in Tbl. 3. In any case, we failed to find any mention of this variant of CPU blocking in
592 the CUDA documentation, and investigating these functions remains an open topic that we
593 plan to explore in future work.

594 An especially worrying pitfall is the following:

595 ► **Pitfall 8.** *CUDA documentation can be contradictory.*

596 In one case, namely `cudaMemcpyAsync`, we discovered that the CUDA documentation
597 actively contradicts itself. Section 3.2.5.1 of the CUDA Programming Guide states “The
598 following device operations are asynchronous with respect to the host: . . . Memory copies
599 performed by functions that are suffixed with **Async**,” but Section 2 of the CUDA Runtime
600 API documentation states “For transfers from device memory to pageable host memory,
601 [`cudaMemcpyAsync`] will return only once the copy has completed.” This raises further doubts
602 about the correctness of other parts of the CUDA documentation.

603 We note that the CUDA API contains 146 non-deprecated or compatibility-related func-
604 tions, and we have only tested a small fraction of these in depth. Therefore, it is likely that
605 our findings with Pitfalls 7 and 8 apply to other portions of the documentation that we have
606 yet to observe.

607 5.3 Unknown Future

608 All of the pitfalls discussed in this paper, as well as the need to compare the alternatives
609 considered in Sec. 4 empirically, can be attributed to a single overarching problem: the
610 black-box nature of current GPU-enabled platforms means that *developers do not have a*
611 *reliable model of GPU behavior*. Much of our group’s prior work has focused on developing
612 such a model. However, this highlights what is perhaps the most important pitfall:

613 ► **Pitfall 9.** *What we learn about current black-box GPUs may not apply in the future.*

614 Despite the fact that we validated our experimental results on several of the most recent
615 CUDA versions and GPU architectures, there is no guarantee that our results will hold after
616 future GPU-architecture or CUDA-version updates. This applies not only to rules about
617 scheduling or blocking, but also may apply to performance characteristics like memory-access
618 times, as we found in prior work [26].

619 Even though other safety-critical hardware inevitably undergoes changes and updates,
620 future-proof programs can still be developed against a stable specification. Likewise, the
621 only way to truly mitigate Pitfall 9 is for GPU manufacturers to release stable, accurate
622 documentation about their GPU platforms, along, preferably, with giving developers greater
623 control over GPU scheduling and synchronization. Only then we can have a reliable GPU
624 model upon which to base real-time analysis and certification. We hope that work such as
625 ours signals to manufacturers like NVIDIA that greater openness is a desirable feature when
626 marketing in safety-critical domains.

627 Unfortunately, there is little indication that NVIDIA plans to move towards open hard-
628 ware or software in the immediate future. In the meantime, one of our continuing objectives
629 is to produce tools, such as our experimental framework, that can be quickly adapted to
630 new GPU hardware. So far, our tools have allowed us to quickly re-validate our prior results
631 every time NVIDIA updates its black-box hardware or software.

632 **6 Conclusion**

633 Vehicles on the road today are already running highly complex GPU-accelerated applica-
634 tions. We anticipate a future where safety-critical autonomous vehicles must be certified,
635 but this will require a change in the GPU-programming paradigm. Currently, computer-
636 vision applications are developed with little guidance about how to achieve temporal safety.
637 Even if a single programmer or application avoids some mistakes, it is increasingly diffi-
638 cult to avoid all of them, especially as applications and task systems grow in complexity.
639 This necessitates work such as ours, which seeks to reduce the gap between computer-vision
640 application developers and those responsible for certifying new systems' real-time safety.

641 With little openness in NVIDIA's hardware and software ecosystem, this paper con-
642 tributes a list of potential pitfalls when developing CUDA applications for real-time sys-
643 tems. Reasons for these pitfalls include GPU synchronization, application performance,
644 and problems with documentation. We uncovered these pitfalls via microbenchmark ex-
645 periments, examining the performance of real-world computer-vision applications, and a
646 careful reading of official GPU documentation. While there is no guarantee of stability in
647 our observations as NVIDIA's hardware and software continues to evolve, we hope that our
648 open-source experimental system will at least make it apparent when changes do occur.

649 This paper is part of an ongoing project with the aim of developing an abstract model of
650 GPU execution. In the future, we plan to continue this investigation and eventually develop
651 middleware capable of intercepting and reordering or delaying GPU operations. Our hope is
652 that the control afforded by such middleware will enable us to produce reasonable analytical
653 bounds on blocking and response times, while maintaining high GPU utilization wherever
654 possible. However, even with better management, certifiable safety in the face of GPU
655 sharing requires a *guarantee* that pitfalls including blocking due to GPU synchronization
656 are controlled, which is only possible if developers of GPU-using software are aware of the
657 consequences and how to avoid them. Fortunately, the best practices we have laid out
658 herein should alleviate much of the strain on application developers on their first foray into
659 real-time systems.

660 In addition to NVIDIA’s GPU, we will also investigate other GPU implementations,
 661 *e.g.*, AMD’s open-source GPU runtime and driver stack. Given the chances of modifying
 662 AMD’s open-source implementation, we are interested in improving the real-time guarantees
 663 of AMD’s GPUs and comparing them with NVIDIA’s GPUs.

664 **References**

- 665 **1** T. Amert, N. Otterness, M. Yang, J. Anderson, and F. D. Smith. GPU scheduling on the
 666 NVIDIA TX2: Hidden details revealed. In *RTSS 2017*.
- 667 **2** J. Aumiller, S. Brandt, S. Kato, and N. Rath. Supporting low-latency CPS using GPUs
 668 and direct I/O schemes. In *RTCSA ’12*.
- 669 **3** C. Basaran and K. Kang. Supporting preemptive task executions and memory copies in
 670 GPGPUs. In *ECRTS ’12*.
- 671 **4** K. Berezovskyi, K. Bletsas, and B. Andersson. Makespan computation for GPU threads
 672 running on a single streaming multiprocessor. In *ECRTS ’12*.
- 673 **5** K. Berezovskyi, K. Bletsas, and S. Petters. Faster makespan estimation for GPU threads
 674 on a single streaming multiprocessor. In *ETFA ’13*.
- 675 **6** K. Berezovskyi, F. Guet, L. Santinelli, K. Bletsas, and E. Tovar. Measurement-based
 676 probabilistic timing analysis for graphics processor units. In *ARCS ’16*.
- 677 **7** K. Berezovskyi, L. Santinelli, K. Bletsas, and E. Tovar. WCET measurement-based and
 678 extreme value theory characterisation of CUDA kernels. In *RTNS ’14*.
- 679 **8** A. Betts and A. Donaldson. Estimating the WCET of GPU-accelerated applications using
 680 hybrid analysis. In *ECRTS ’13*.
- 681 **9** N. Capodieci, R. Cavicchioli, P. Valente, and M. Bertogna. SiGAMMA: Server based
 682 integrated GPU arbitration mechanism for memory accesses. In *RTNS 2017*.
- 683 **10** R. Cavicchioli, N. Capodieci, and M. Bertogna. Memory interference characterization
 684 between CPU cores and integrated GPUs in mixed-criticality platforms. In *RTNS 2017*.
- 685 **11** G. Elliott, B. Ward, and J. Anderson. GPUSync: A framework for real-time GPU man-
 686 agement. In *RTSS ’13*.
- 687 **12** B. Forsberg, A. Marongiu, and L. Benini. Gpuguard: Towards supporting a predictable
 688 execution model for heterogeneous SoC. In *DATE ’17*.
- 689 **13** A. Horga, S. Chattopadhyayb, P. Eles, and Z. Peng. Systematic detection of memory
 690 related performance bottlenecks in GPGPU programs. In *JSA ’16*.
- 691 **14** P. Houdek, M. Sojka, and Z. Hanzálek. Towards predictable execution model on ARM-
 692 based heterogeneous platforms. In *ISIE ’17*.
- 693 **15** S. Kato, K. Lakshmanan, A. Kumar, M. Kelkar, Y. Ishikawa, and R. Rajkumar. RGEM:
 694 A responsive GPGPU execution model for runtime engines. In *RTSS ’11*.
- 695 **16** S. Kato, K. Lakshmanan, R. Rajkumar, and Y. Ishikawa. TimeGraph: GPU scheduling
 696 for real-time multi-tasking environments. In *USENIX ATC ’11*.
- 697 **17** H. Lee and M. Abdullah Al Faruque. Run-time scheduling framework for event-driven
 698 applications on a GPU-based embedded system. In *TCAD ’16*.
- 699 **18** A. Li, G. van den Braak, A. Kumar, and H. Corporaal. Adaptive and transparent cache
 700 bypassing for GPUs. In *SIGHPC ’15*.
- 701 **19** X. Mei and X. Chu. Dissecting GPU memory hierarchy through microbenchmarking. In
 702 *TPDS ’16*.
- 703 **20** Multi-process service. Online at [https://docs.nvidia.com/deploy/pdf/CUDA_Multi-
 704 Process_Service_Overview.pdf](https://docs.nvidia.com/deploy/pdf/CUDA_Multi-Process_Service_Overview.pdf).
- 705 **21** NVIDIA. Embedded systems developer kits and modules. Online at [http://www.nvidia.
 706 com/object/embedded-systemsdev-kits-modules.html](http://www.nvidia.com/object/embedded-systemsdev-kits-modules.html).

- 707 **22** NVIDIA. Best practices guide. Online at [http://docs.nvidia.com/cuda/](http://docs.nvidia.com/cuda/cuda-c-best-practices-guide/index.html)
708 [cuda-c-best-practices-guide/index.html](http://docs.nvidia.com/cuda/cuda-c-best-practices-guide/index.html), 2017.
- 709 **23** NVIDIA. CUDA toolkit documentation v9.1.85. Online at [http://docs.nvidia.com/](http://docs.nvidia.com/cuda/)
710 [cuda/](http://docs.nvidia.com/cuda/), 2018.
- 711 **24** N. Otterness, V. Miller, M. Yang, J. Anderson, F.D. Smith, and S. Wang. GPU sharing
712 for image processing in embedded real-time systems. In *OSPERT '16*.
- 713 **25** N. Otterness, M. Yang, T. Amert, J. Anderson, and F.D. Smith. Inferring the scheduling
714 policies of an embedded CUDA GPU. In *OSPERT '17*.
- 715 **26** N. Otterness, M. Yang, S. Rust, E. Park, J. Anderson, F.D. Smith, A. Berg, and S. Wang.
716 An evaluation of the NVIDIA TX1 for supporting real-time computer-vision workloads. In
717 *RTAS '17*.
- 718 **27** U. Verner, A. Mendelson, and A. Schuster. Batch method for efficient resource sharing in
719 real-time multi-GPU systems. In *ICDCN '14*.
- 720 **28** U. Verner, A. Mendelson, and A. Schuster. Scheduling periodic real-time communication
721 in multi-GPU systems. In *ICCCN '14*.
- 722 **29** U. Verner, A. Mendelson, and A. Schuster. Scheduling processing of real-time data streams
723 on heterogeneous multi-GPU systems. In *SYSTOR '12*.
- 724 **30** H. Wong, M. Papadopoulou, M. Sadooghi-Alvandi, and A. Moshovos. Demystifying GPU
725 microarchitecture through microbenchmarking. In *ISPASS '10*.
- 726 **31** Y. Xu, R. Wang, T. Li, M. Song, L. Gao, Z. Luan, and D. Qian. Scheduling tasks with
727 mixed timing constraints in GPU-powered real-time systems. In *ICS '16*.
- 728 **32** J. Zhong and B. He. Kernelet: High-throughput GPU kernel executions with dynamic slic-
729 ing and scheduling. *IEEE Transactions on Parallel and Distributed Systems*, 25:1522–1532,
730 2014.
- 731 **33** H. Zhou, G. Tong, and C. Liu. GPES: A preemptive execution system for GPGPU com-
732 puting. In *RTAS '15*.