

Real-Time GPU Resource Management with Loadable Kernel Modules

Yusuke Fujii*, Yuhei Suzuki*, Takuya Azumi[†], Nobuhiko Nishio*, and Shinpei Kato[‡]

* Graduate School of Information Science and Engineering, Ritsumeikan University

[†] Graduate School of Engineering Science, Osaka University

[‡] Graduate School of Information Science, Nagoya University

Abstract—Graphics processing units (GPUs) are increasingly well designed for high-performance computing. GPU programming environments have also become matured for general-purpose computing on GPUs (GPGPU). Originally, the GPU software stack is tailored to accelerate particular best-effort applications. In recent years, however, GPU technologies have been applied for real-time systems, extending the operating system (OS) modules to support real-time GPU-resource management. Unfortunately, such a system extension makes it difficult to maintain the system with version updates, since the OS kernel and device drivers need to be modified at the source-code level, preventing continuous research and development of GPU technologies for real-time systems. A loadable-kernel-module (LKM) framework, called “Linux-RTXG,” for managing real-time GPU resources on Linux without having to modify the OS kernel and device drivers is proposed and experimentally evaluated. Linux-RTXG provides novel mechanisms for interrupt interception and independent synchronization to achieve real-time scheduling and resource reservation capabilities for GPU applications on top of existing device drivers and runtime libraries. The results of an experimental evaluation of Linux-RTXG demonstrate that the overhead incurred by introducing Linux-RTXG is comparable to that for introducing existing kernel-dependent approaches. Moreover, they demonstrate that multiple GPU applications can be successfully scheduled by Linux-RTXG to meet their priority and quality-of-service (QoS) requirements in real-time.

Keywords—*Graphics Processing Units, Resource Management, Real-Time Scheduling, Operating Systems, Linux*

I. INTRODUCTION

Graphics processing units (GPUs) are now common platforms for various data-parallel and compute-intensive applications. Although GPUs are primarily used to accelerate high-performance-computing (HPC) applications, their performance advantage is increasingly recognized in regard to real-time systems. Examples of such GPU applications include route navigation for autonomous vehicles [1], object detection [2], plasma control for fusion reactors [3], windowing applications [4], and database operators [5]. Benchmarking suites for various workloads [6] are also well deployed.

GPU applications, whose main purpose is to accelerate particular computing blocks, were often best-effort oriented. Conventional GPU technologies are therefore particularly designed for individual data-parallel and compute-intensive workloads. However, due to an emergence of real-time systems using GPUs, supporting system software for real-time management of GPU resources is becoming a more significant problem.

Since the standard system software released by GPU vendors and communities, such as CUDA [7] and OpenCL [8], is not tailored to support real-time systems, extending the system software (including the operating-system (OS) kernel and device drivers) requires a complex undertaking.

In previous work, GPU resource management for real-time systems was developed. In detail, TimeGraph [9] provides GPU scheduling and resource reservation capabilities at the device-driver level, multiplexing GPU commands to support priorities, and quality of service (QoS) for GPU applications. Gdev [10] is a rich set of runtime libraries and device drivers that achieves first-class GPU resource management, where GPU contexts are fully managed in the OS space. The main drawback of these work is that they require detailed information on the implementation of GPU runtime libraries and device drivers, mostly obtained by reverse engineering.

In general, a GPU software stack is encapsulated by an application-programming interface (API) in runtime libraries and an application binary interface (ABI) in device drivers. This means that modifications to system software outside of the API and ABI layers allow the entire software stack to be loadable, enabling more sustainable solutions. CPU scheduling is also important for GPU applications, because CPU time is consumed when GPU functions are launched from the API and ABI layers. To satisfy real-time requirements in regard to GPU computing, therefore, CPU and GPU resources must be managed in coordination.

Presented by Elliott et al., one notable work on coordinated CPU and GPU resource management for real-time systems is GPUSync [11], [12], which was intended to extend CPU scheduling for multiple GPU-aware contexts with budget enforcement. GPUSync was built on top of the API and ABI layers of proprietary software, providing a configurable framework that can verify the combination of task-allocation policies for multicore CPUs and GPU-resource-allocation policies for multiple GPUs.

GPUSync is implemented by using *LITMUS^{RT}* [13], which introduces a significant amount of changes to the OS kernel. TimeGraph and Gdev also make some modifications to the device driver. Those built-in approaches to the OS kernel and device drivers require users to install patches or developers to maintain patches in order to stay up to date with the latest version releases. This so-called “porting” work is a complex undertaking against research and development and

is not sustainable, given that open-source software, such as Linux, especially, is frequently updated with non-trivial code changes.

Linux supports the loadable kernel module (LKM), which can load and unload kernel modules so as to, respectively, add and remove supplementary kernel functions. Existing Linux-based real-time OSs use an LKM to extend real-time capabilities of Linux. In particular, RESCH [14], [15] provides a real-time scheduling framework, called ExSched, which is completely independent of code modifications to the OS kernel and device drivers. Unfortunately, RESCH does not support GPU resource management. It has never been demonstrated that GPU resource management for real-time systems can be fully implemented by using an LKM, without making any modifications to the OS kernel and device drivers.

Contribution: An LKM framework called “Linux-RTXG (Linux Real-Time eXtention with GPUs),” which provides LKM-based real-time GPU resource management in Linux, is proposed. Linux-RTXG allows the system to easily re-configure scheduling algorithms and install their modules at runtime for GPU applications. The most significant contribution of Linux-RTXG is that resource management modules for not only CPUs but also GPUs can be added to Linux without any modification to the OS kernel and device drivers. CPU mechanisms for scheduling and resource reservation are based on RESCH, while GPU mechanisms for scheduling and resource reservation are implemented by using Gdev. In addition to integrating RESCH and Gdev, Linux-RTXG provides a new framework to coordinate CPU and GPU-resource management, freeing them from the built-in OS kernel and device drivers. It is thus a completely “patchless” approach.

Organization The rest of this paper is organized as follows. Section II describes the system model and basic approaches for Linux-RTXG. With a particular emphasis on GPU scheduling, Section III presents the design and implementation of Linux-RTXG. Section IV presents the results of an evaluation of the system overhead and reservation performance of Linux-RTXG. Section V discusses related work, and Section VI concludes the paper.

II. SYSTEM MODEL

A model of GPU programming and GPU scheduling is assumed and described as follows. To highlight an unsolved problem concerning GPU resource management, existing work is introduced first. In addition to that problem, factors that have prevented GPU resource management from patchless implementation are described. As for the model, the system is assumed to be composed of multiple GPUs and multi-core CPUs.

A. GPU Programming Model

General-purpose computing on GPUs (GPGPU) is supported by special programming languages such as CUDA and OpenCL. In this work, it is assumed that GPU application programs are written in CUDA; however, the concept of GPU resource management studied in this work is not limited by programming languages. Conceptually, the contribution of this work is applicable to OpenCL and other languages. A GPU task is defined as a process running on the CPU that launches a GPU kernel to the GPU, whose cyclic execution unit is referred

to as a “GPU job.” The GPU kernel is the process that is executed on the GPU. Multi-tasking environments, in which multiple tasks of GPU applications are allowed to coexist, are also assumed, though the current model of GPU programming does not allow multiple contexts to be executed at the same time on the GPU. In other words, multiple GPU contexts can be created and launched, but they must be executed exclusively on the GPU.

GPU programming requires a set of APIs, such as CUDA Driver API and CUDA Runtime API, provided by runtime libraries. A typical approach to GPU programming follows several steps: (i) *cuCtxCreate* creates a GPU context; (ii) *cuMemAlloc* allocates memory space to the device memory; (iii) *cuModuleLoad* and *cuMemcpyHtoD* copy the data and the GPU kernel from the host memory to the memory space of the allocated device; (iv) *cuLaunchGrid* invokes the GPU kernel; (v) *cuCtxSynchronize* synchronizes a GPU task to wait for the completion of GPU kernel; (vi) *cuMemcpyDtoH* transfers the data back to the host memory from the device memory; and (vii) *cuMemFree*, *cuModuleUnload*, and *cuCtxDestroy* release the allocated memory space and the GPU context.

B. GPU Scheduling

Problems concerning resource management in many variants of real-time OSs (RTOSs) [16], [17], [18], [19] have been addressed. Of particular interest includes those in Linux-based RTOSs [13], [20], [21], [22], [14]. GPUs are also supported by Linux. It is assumed in this work that the OS architecture refers to Linux.

To meet real-time requirements in multi-tasking environments, RTOSs should provide *scheduling* and *resource reservation* capabilities. Rate Monotonic (RM) and Earliest Deadline First (EDF) [23] are well-known algorithms for priority-driven scheduling for real-time systems. Resource-reservation algorithms come in many variants, such as Constant Bandwidth Server (CBS) [24] and Total Bandwidth Server (TBS) [25].

GPU computing must deal with the data-transfer bandwidth and compute cores as shared resources. In case of real-time systems, therefore, scheduling and resource reservation for the GPU as well as the CPU must be considered. TimeGraph and Gdev are involved in solving problems concerning GPU resource management, but they are not much aware of the fact that GPU tasks consume CPU time to drive APIs. To support GPUs in real-time systems, the OS scheduler must be able to manage GPU tasks in coordination with CPU tasks on the host side, while monitoring GPU time for completion of GPU kernels. Therefore, the problem of CPU and GPU coordinated scheduling is considered in term of resource-reservation mechanisms.

The recently developed GPUSync framework indeed employs CPU and GPU coordinated scheduling. In the case of fGPUSync, the device driver and runtime library for GPU computing are hidden in black-box modules released by GPU vendors. On the other hand, TimeGraph and Gdev address this problem by using reverse-engineered open-source software. Both approaches are limited to some extent. That is to say, using black-box modules makes it difficult to manage the system in a fine-grained manner. Open-source software tends to lack some functionality due to incomplete reverse engineering. GPUSync manages to incorporate black-box modules in

scheduling and resource reservation by arbitrating interrupt handlers and runtime accesses. This GPUSync approach is preferred in the sense a reliable proprietary driver and library can be utilized, while functions for scheduling and resource reservation can still be provided.

GPU Synchronization: Given that the GPU is a coprocessor connected to the host CPU, synchronization between the GPU and the CPU must be considered to guarantee the correctness of execution logics. Most GPU architectures support two synchronization mechanisms. One is based on memory-mapped registers called “FENCE.” To use FENCE, it is necessary to send special commands to the GPU when a GPU kernel is launched so that the GPU can write the specified value to this memory-mapped space when the GPU kernel is completed. On the host side, a GPU task is polling this value via the mapped space. The other technique is interrupt-based synchronization called “NOTIFY.” In a similar manner to FENCE, to use NOTIFY, it is also necessary to send special commands to the GPU. Instead of writing to memory-mapped registers, NOTIFY raises an interrupt from the GPU to the CPU and, at the same time, writes some associated values to I/O registers of the GPU. On the host side, a GPU task is suspended so that it waits for the occurrence of the interrupt. NOTIFY allows other tasks to use CPU resources when the GPU task is waiting for completion of execution of its GPU kernel, though a scheduling overhead is involved. FENCE is easier to use and is more responsive than NOTIFY, but it is implemented at the expense of CPU utilization. More details about GPU architectures can be found in a previous work [9], [10], [26].

To synchronize the CPU and the GPU, Gdev supports both NOTIFY and FENCE. NOTIFY is primarily used for scheduling of GPU tasks, while FENCE is used in driver-level synchronization. In the case of Gdev, as aforementioned, the implementation of GPU synchronization involves additional commands to the GPU. The device driver must thus be modified to a certain extent. GPUSync indirectly utilizes the NOTIFY technique with tasklet interception [27] on top of the proprietary black-box modules. “Tasklet” refers to Linux’s softirq implementation. GPUSync identifies the interrupt that has invoked a GPU kernel using a callback pointer with a tasklet.

Loadable Kernel Modules: The main concept of our system is to enable both the CPU and GPU to be managed by the OS scheduler without any code changes to the OS kernel and device drivers. CPU scheduling under this constraint has already been demonstrated by RESCH [14], [15]. To schedule GPU tasks, it must be possible to hook the scheduling points where the preceding GPU kernel is completed and the active context is switched to the next GPU kernel. The scheduling points can be hooked by two methods. The API-driven method (presented by RGEM [28]) explicitly awakens the scheduler after GPU synchronization, such as *cuCtxSynchronize()*, invoked by the API. The interrupt-driven method (presented by TimeGraph and Gdev), on the other hand, uses interrupts that can be configured by NOTIFY. GPUSync is also based on this interrupt-driven method. Especially in the case of CUDA, the standard *cuCtxSynchronize()* API synchronizes completion of all GPU kernels. Therefore, the API-driven method can be used if a GPU context issues no more than one GPU kernel. In

other words, if a GPU task invokes multiple GPU kernels, the interrupt-driven method is more appropriate to realize real-time capabilities.

The interrupt-driven method forces TimeGraph, Gdev, and GPUSync to modify the code of either the Linux kernel or device drivers. This modification is necessary because the interrupt service routine (ISR) must be managed to create the scheduling points. Gdev has developed independent synchronization mechanisms on top of the proprietary software; however, for scheduling and resource reservation on Gdev, the Linux kernel still needs some modification. As a result, the available release versions of the Linux kernel and device drivers for TimeGraph, Gdev, and GPUSync are very limited. A core challenge of this work is to develop independent synchronization mechanisms that do not require any modification to the OS kernel and device drivers so that a coordinated CPU and GPU resource-management scheme can be utilized with a wide range of the Linux kernel and device drivers.

Scope and Limitation: GPU resource management has a “non-preemptive” nature. For example, the execution of GPU kernels is non-preemptive. The data transfer between the CPU and the GPU is also non-preemptive. Previous work has addressed the problem of response time by preventing overrun that occurs while the kernel is being divided [29], [30]. The most difficult problem is scheduling of self-suspending tasks because the GPU is treated as an I/O device in the system. For example, GPU tasks are suspended until their GPU kernels are completed. This “self-suspending” problem was presented as an NP-hard problem in previous work [31], [32]. A lot of work on the scheduling of self-suspending tasks, however, is ongoing [33], [34].

Such a “schedulability” problem of GPU scheduling out of the scope of this work. Instead, the design and implementation of GPU scheduling and resource reservation with existing algorithms is focused on. The scope of this work is also restricted to time resources (but not to memory resources). In other words, GPU-scheduling problems, not device-memory-allocation problems, are considered in the following.

The proposed system supports efficient data transfer between the CPU and the GPU by using GPU microcontrollers [35], but that support is not considered as a contribution of this work. Our prototype system is limited to a Linux system and CUDA environment; however, the concept of our method is applicable to other OSs and programming languages, as far as they support the aforementioned FENCE and NOTIFY primitives.

III. DESIGN AND IMPLEMENTATION

The design and implementation of Linux-RTXG, which provides a framework for CPU and GPU coordinated scheduling based on the LKM, is described as follows. In particular, the approach taken for GPU scheduling and its integration to CPU scheduling are explained. Due to a space constraint, the detail of LKM-based CPU scheduling is referred to in the RESCH project [14], [15].

A. Linux-RTXG

An architectural overview of Linux-RTXG is shown in Figure 1. The system architecture of Linux-RTXG can

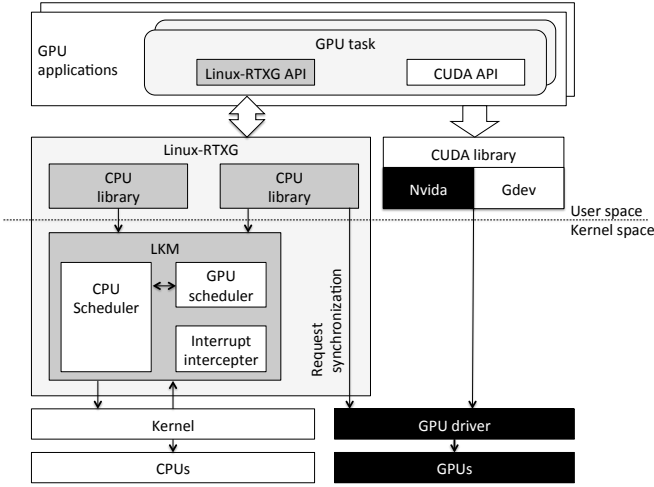


Fig. 1. Architectural overview of Linux-RTXG.

be divided into two parts. First, the Linux-RTXG core contains a CPU scheduler and a GPU scheduler with a resource-reservation mechanism. The implementation of the Linux-RTXG core is provided in the kernel space by an LKM. It can thus use exported Linux kernel functions, such as *schedule()*, *mod_timer()*, *wake_up_process()*, and *set_cpus_allowed_ptr()*. These functions can be called from the user-space interface by using the input/output control (*ioctl*) system call, which is a standard system call for device drivers. Secondly, the Linux-RTXG library provides an independent synchronization method for coordinated management of CPU and GPU resources. The independent synchronization method can be used on top of a proprietary driver [7] as well as an open-source driver [36]. Note that this method is required to manage interrupts for GPU scheduling without modifying any codes of the OS kernel and device drivers.

B. GPU Scheduling

As well as being mostly based on the interrupt-driven method for GPU synchronization, Linux-RTXG is partly based on the API-driven method. The scheduler is invoked only when computation requests are submitted. The basic APIs supported by Linux-RTXG are listed in Table I. Note that some APIs have arguments and others do not. Being independent of GPU runtimes, Linux-RTXG does not modify the existing CUDA API to cope with proprietary software. However, CUDA application programs must add the Linux-RTXG APIs to use the functionality of Linux-RTXG.

The sample code including the Linux-RTXG APIs is shown in Figure 2. GPU tasks are provided with function calls to Linux-RTXG at strategic points.

The execution flow of the GPU task managed by the Linux-RTXG APIs is shown in Figure 3. Note that this example is restricted to a single GPU kernel. The GPU task can control the timing that the GPU kernel is invoked by calling *rtx_gpu_launch()*. After this function call, the task is suspended until it receives an interrupt so that other GPU kernels can be launched. When execution of the GPU kernel is completed, an interrupt is raised by the GPU, and the corresponding interrupt handler is executed by the Linux kernel. The interrupt interceptor awakens some suspended

```
void gpu_task(){
    /* variable initialization */
    /* calling RESCH API */
    dev_id = rtx_gpu_device_advice(dev_id);
    cuDeviceGet(&dev, dev_id);
    cuCtxCreate(&ctx, SYNC_FLAG, dev);
    rtx_gpu_open(&handle, vdev_id);
    /* Module load and set kernel function */
    /* Device memory allocation */
    /* Memory copy to device from host */
    rtx_gpu_launch(&handle);
    cuLaunchGrid(function, grid_x, grid_y);
    rtx_gpu_notify(&handle);
    rtx_gpu_sync(&handle);
    /* Memory copy to host from device */
    /* Release allocated memory */
}
```

Fig. 2. Sample code with the Linux-RTXG APIs.

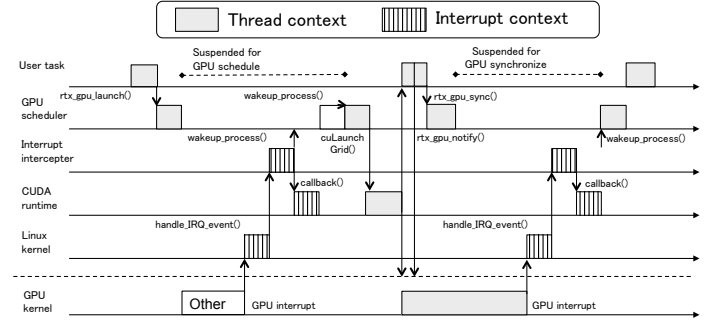


Fig. 3. Execution flow of a GPU task.

tasks according to the priority. The awakened task proceeds to launch the GPU kernel by using a CUDA API, such as *cuLaunchGrid()*. After the GPU kernel is launched, the task is going to register NOTIFY to set up an interrupt by *rtx_gpu_notify()*, and it is put into sleep mode again until it receives the interrupt by *rtx_gpu_sync()*. Dispatching the subsequent task is performed by the GPU scheduler, which is called upon by the interrupt from the GPU. Linux-RTXG manages the order of task execution according to this flow.

A hierarchical scheduling mechanism that applies the concept of virtual GPUs to combine specified GPU tasks by a group is presented in the following. The virtual GPUs are created by a resource-reservation mechanism, while GPU scheduling uses a priority mechanism. Specifically, each invocation of the GPU kernel is associated with a scheduling entity, and Linux-RTXG allocates the scheduling entities to virtual GPUs. The virtual GPUs can belong to any physical GPUs. In the case of Linux-RTXG, computing resources are distributed to virtual GPUs.

The pseudo-code of the Linux-RTXG scheduler is shown in Figure 4. This code works under the assumption that *on_arrival* is called when a GPU task requests to launch the GPU kernel. When *on_arrival* is called, the GPU task checks whether the given execution permission is held by the allocated virtual GPU as well as the allocated scheduling entity. If the virtual GPU to which the GPU task belongs does

TABLE I. A BASIC SET OF APIS FOR LINUX-RTXG.

<code>rtx_gpu_open()</code>	Registers itself to Linux-RTXG and creates a scheduling entity. It must be called first.
<code>rtx_gpu_device_advice()</code>	Obtains recommendations for which GPU devices to use.
<code>rtx_gpu_launch()</code>	Controls timing of launch of the GPU kernel, (i.e., a scheduling entry point). It must be called before the CUDA launch API.
<code>rtx_gpu_sync()</code>	Waits for the completion of execution of the GPU kernel by sleeping with TASK UNINTERRUPTIBLE status.
<code>rtx_gpu_notify()</code>	Sends a NOTIFY/FENCE command to the GPU. Either FENCE or NOTIFY is selected by a flag that is set by an argument.
<code>rtx_gpu_close()</code>	Releases the scheduling entity.

```

se: Scheduling entity
se->vgpu: Group that belongs to se
se->task: Task that is associated with se
vgpu->parent: Physical GPU identifier

void on_arrival(se) {
    check_permit_vgpu(se->vgpu)
    while (!check_permit_se(se)) {
        enqueue(se->vgpu, se);
        sleep_task(se->task);
    }
}

void on_completion(se) {
    reset_the_permit(se->vgpu, se)
    n_vgpu = pick_up_the_next_vgpu(se->vgpu->parent)
    se = pick_up_the_next_se(n_vgpu)
    if (se) {
        dequeue(se->vgpu, se);
        wakeup_task(se->task);
    }
    set_the_permit(se->vgpu, se)
}

```

Fig. 4. High-level pseudo-code of the Linux-RTXG scheduler.

not hold the execution permission, the GPU task is enqueued to *wait_queue* and suspended. Otherwise, the GPU task can launch the GPU kernel. After a while, *on_completion* is called by the scheduler thread when the launched GPU kernel is completed, and the next set of a virtual GPU and a GPU task can be selected. At the end of *on_completion*, the selected GPU task is woken up.

C. GPU Synchronization

The independent-synchronization and interrupt-interception mechanisms are described next. The independent-synchronization mechanism invokes NOTIFY and FENCE without using the GPU runtime API. The interrupt-interception mechanism enables interrupt-driven invocation of the scheduler without modifying either the OS kernel or device drivers. By this means, Linux-RTXG is not required to modify the OS kernel and device drivers, while being able to create the scheduling points for GPU tasks.

Independent-synchronization Mechanism: The independent-synchronization mechanism using NOTIFY and FENCE is described as follows. This mechanism invokes an interrupt using NOTIFY, and writes the fence value using the GPU microcontrollers and FENCE. NVIDIA's proprietary software uses the ioctl interface to communicate with the kernel space and the user space. These ioctl interfaces provide driver functions, such as device-memory allocation, obtaining GPU states and memory mapping. Gdev provides a runtime library that can control the GPU on top of NVIDIA's proprietary driver by using these ioctl interfaces. The mechanism also uses an ioctl interface similar to Gdev in order to send commands to the GPU. Specifically, the

mechanism is divided into two processes: (i) initialization and (ii) notification.

The initialization process generates a dedicated GPU context. This process creates a virtual address space, allocates an indirect buffer object for commands, and creates a context object that is required to employ the FIFO engine, followed by allocation of a kernel-memory object and mapping of the FIFO-engine registers to host memory space through a memory-mapped I/O (MMIO). The FIFO engine is a GPU microcontroller that decodes and dispatches the commands sent from the host CPU side.

The notification process sends commands to a GPU compute engine or a GPU data-copy engine by the *iowrite* function associated with the mapped FIFO-engine registers so that an interrupt will be raised from the GPU to the CPU. The compute engine and the data-copy engine are such GPU microcontrollers that control the states of GPU computation and data transfer. They are also used to switch GPU contexts on the GPU computation and data transfer. Note that this independent-synchronization mechanism requires the information of ioctl interfaces. Therefore, it depends on the GPU architecture and implementation of device drivers.

Interrupt Interception: Interrupts are handled by the ISR that is registered to the Linux kernel by the device driver. The scheduler function is required to receive the interrupts and identify them by reading the GPU status register. The GPU status register must be read by the OS scheduler before it is reset by the ISR.

The Linux kernel has a structure that holds interrupt parameters, called *irq_desc*, for each interrupt number. This structure has an internal sub-structure called *irq_action*, which includes the ISR callback pointer. The *irq_desc* structure is allocated to the global kernel memory space, and is freely accessible from the kernel space. Therefore, not only the Linux kernel but also external LKMs can obtain the information of *irq_desc* and the ISR callback pointer. The ISR callback pointer associated with the GPU device driver is obtained, and a new interrupt interception ISR is registered in the Linux kernel. Finally, interrupts from the GPU through the ISR can be intercepted, and the callback pointer can be retained. In addition, I/O registers are mapped from the PCIe base address registers (BAR) to the kernel memory space by the device driver [35], [37]. Therefore, Linux-RTXG remaps the BAR0 to the allocated space by using *ioremap()* when the ISR is initialized. The interrupt-interception mechanism can identify the source of every interrupt by reading this remapped space.

D. Scheduler Integration

The mainline Linux scheduler implements the following three real-time scheduling policies:

- *SCHED_DEADLINE*

- *SCHED_FIFO*
- *SCHED_RR*

SCHED_DEADLINE is the implementation of CBS and EDF, which is the latest real-time scheduler for Linux introduced in version 3.14.0, while *SCHED_FIFO* and *SCHED_RR* represent fixed-priority scheduling. Unfortunately, synchronization does not work with the *SCHED_DEADLINE* scheduling policy for GPU tasks. Two problems concerning implementation must therefore be solved. The first problem is attributed to the implementation of *sched_yield()*. Note that *sched_yield()* uses *yield()* in the kernel space. Releasing the CPU by *sched_yield()* while waiting for I/O in polling makes it possible to utilize CPU time more efficiently. However, *sched_yield()* will set the remaining execution time of the polling task to zero by treating it as a parameter of *SCHED_DEADLINE*. As a result, the task cannot be executed until the runtime is replenished in the next period. This means that *sched_yield()* should not be called while polling in *SCHED_DEADLINE*. However, *sched_yield()* is frequently used by device drivers and libraries. Real-time performance of GPU computing, even in CUDA, could be affected by this problem. This problem is addressed by limiting the GPU synchronization method to NOTIFY, namely, eliminating the potential of FENCE, in the *SCHED_DEADLINE* scheduling policy.

The second problem is subject to the implementation of wake-up and sleep functions, particularly the check equation (1) when restoring a task from the sleep state. If equation (1) holds, the runtime is replenished, and the absolute deadline is set to the next cycle deadline.

$$\frac{Absolute_Deadline - Current_Time}{Remaining_Runtime} > \frac{Relative_Deadline}{Period} \quad (1)$$

This check condition is revised so that the GPU execution time is subtracted from *RemainingRuntime* when a task is restored by the GPU kernel execution, except when a task is restored by the period.

IV. EVALUATION

The runtime overhead and real-time performance of Linux-RTXG was evaluated. The runtime overhead is classified into interrupt interception, independent synchronization, and priority-driven scheduling in Linux. It is demonstrated in the following that the overhead of the LKM-based real-time scheduler for GPU applications is acceptable. The real-time performance is verified in terms of QoS management and prioritization using both synthetic workloads and real-world applications on top of three different device drivers, namely, NVIDIA, Nouveau, and Gdev. It is also demonstrated that multiple GPU applications co-scheduled by the LKM-based real-time scheduler are successfully prioritized and maintained at the desired frame rate even under high CPU load.

Given that CPU scheduling performance has already been demonstrated in previous work [14], only GPU scheduling performance was experimentally evaluated in this work. Considering the results of this work and those from the previous work, it is clarified that both emerging GPU applications

and traditional CPU applications can be scheduled according to priorities and resource reserves by using LKMs, without modifying the Linux kernel and device drivers.

A. Experimental Setup

Experiments were conducted using the Linux kernel 3.16.0, an NVIDIA GeForce GTX680 GPU, a 3.40-GHz Intel Core i7 2600 (eight cores including two hyper-threading cores), and 8-GB main memory. GPU application programs were written in CUDA and compiled by NVCC v6.0.1. A NVIDIA 331.62 driver and a Nouveau Linux-3.16.0 driver with NVIDIA CUDA 6.0 and Gdev were used.

B. Interrupt-interception Overhead

The overhead incurred by interrupt interception was measured by using the Nouveau GPU driver to quantify the overhead of interception in the case of varied interrupt types. As performance metrics, elapsed time from the beginning to the end of the ISR was adopted.

The overhead for interrupt interception is shown in Figure 5. “Raw ISR” represents the original implementation of ISR. “ISR Intercept” represents the modified ISR with the interrupt interception mechanism, while “ISR Intercept w/Func” includes the overhead for identifying the ISR and calling the scheduler thread in addition to “ISR Intercept”. The average execution time for 1,000 executions is presented with error bars.

Comparing Raw ISR and ISR Intercept indicates that the overhead for introducing the interrupt-interception mechanism was 247ns, which is only 0.8% of Raw ISR. Similarly, the overhead for introducing ISR Intercept w/Func was 790ns, which is only 2.6% of Raw ISR. This result indicates that the interrupt-interception overhead is insignificant in regard to total performance.

As shown in Figure 6, response times of ISR with and without interrupt interception are compared next. The response time is defined as the elapsed time from the beginning of interrupt processing to the point where the interrupt type (e.g., timer, compute, FIFO, and GPIO) is identified. According to the result, interrupt interception makes response times 1.4-times longer.

Response times when interrupt interception is self-contained in the ISR (top-half) and when it is expanded to the tasklet (bottom-half) were also compared. This comparison differentiates Linux-RTXG from GPUSync in terms of performance, since GPUSync intercepts the *tasklet_schedule()* on top of the proprietary driver, while Linux-RTXG realized ISR-based interception. In this case, the start time of measuring the response time is set to that of the function call to *do_IRQ*. The result of this comparison is shown in Figure 7. As can be seen from this figure, the tasklet-based approach requires about five-times longer response times than the IRQ-based approach. This is because the tasklet is typically called after significant ISR processing steps.

C. Independent-synchronization Overhead

The overhead of the independent synchronization mechanism was quantified next. This mechanism must call *rtx_gpu_notify()* at the time of requested synchronization

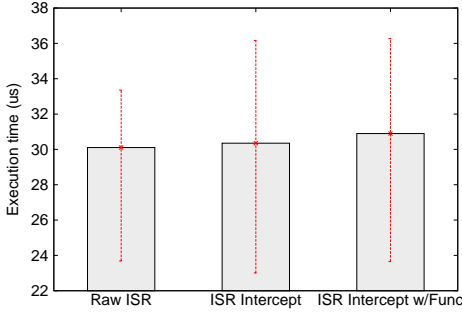


Fig. 5. Interrupt interception overhead.

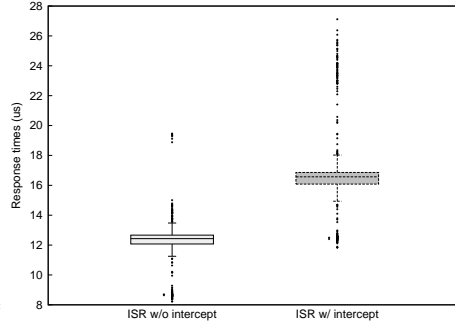


Fig. 6. Impact of interrupt interception.

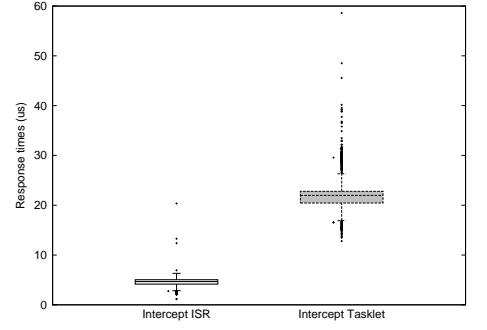


Fig. 7. Comparison of ISR and tasklet.

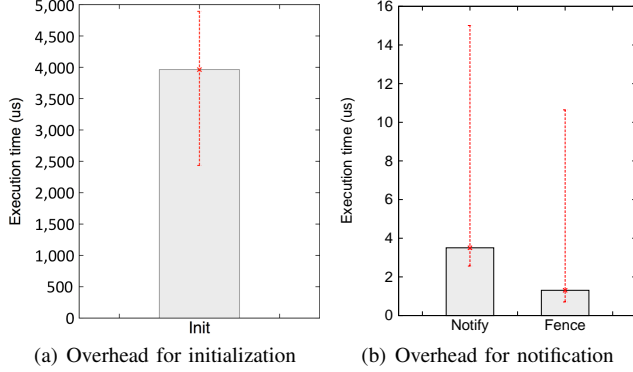


Fig. 8. Overhead for independent synchronization.

(e. g. , after launching a GPU kernel). To use this mechanism, an initialization procedure is required. In the measurement of the initialization overhead, the overhead is defined by the execution time of corresponding APIs.

As shown in Figure 8, the initialization overhead could reach $5,000us$, whereas the notification overhead is no more than a few microseconds. The initialization procedure calls a Linux process to allocate indirect buffers and control registers for several GPU engines. Albeit a significant overhead, the application program is not much affected by this procedure, since it is called only once at the beginning. On the other hand, the notification overhead is not a major consideration for the application program, though it is a little scattered due to an `ioctl` system call. When `NOTIFY` is used to set the notification, the overhead was $3.5us$, while it was reduced to $2us$ when `FENCE` is used.

D. Scheduling Overhead

The scheduling overhead incurred by the presented Linux-RTXG scheduler was evaluated next. Three synthetic tasks – (i) vanilla, (ii) mutex, and (iii) rtx – to measure the overhead were executed. These tasks are based on the microbenchmark program provided by the Gdev project, which tests a GPU loop function. The program was modified to generate multiple GPU tasks by using the `fork()` system call. Each task releases a job ten times periodically, each of which includes data transfer between the CPU and the GPU, followed by the execution of the GPU kernel. The rtx task was scheduled by Linux-RTXG. The mutex task was limited to a single GPU kernel by using explicit mutual-exclusion control similar to the rtx task. The vanilla task was not changed from the original microbenchmark program. The CPU scheduling policy was set to `SCHED_FIFO`, while the GPU scheduling policy

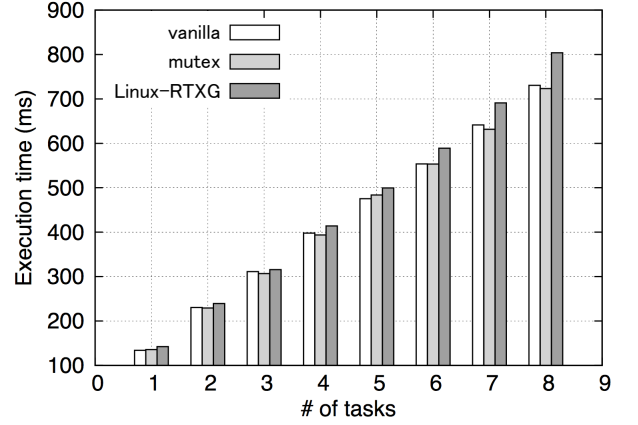


Fig. 9. Scheduling overhead.

was based on fixed priorities with GPU-resource reservation, called the BAND scheduling policy [10]. The independent synchronization mechanism is employed with `NOTIFY`.

We measured average time in 100 times of GPU task execution (1,000 jobs). The result is provided in Figure 9. The scheduling overhead increases in proportion as the number of tasks, because the time consumed in queueing GPU tasks is increased. The maximum overhead corresponds to 10% of the vanilla task for eight tasks.

E. Prioritization and QoS Performance

Experiments were also performed to evaluate prioritization and QoS performance for GPU applications provided by Linux-RTXG. These performances of Linux-RTXG were evaluated on both the proprietary driver (NVIDIA) and the open-source driver (Nouveau), as compared to the built-in-kernel approach (Gdev).

QoS management indicates if GPU time of the corresponding task is guaranteed. QoS performance was evaluated by observing how well the tasks are isolated. First, GPU utilization when running two GPU tasks was measured. Each GPU task has a unique workload with a different size of GPU reserve. One task was allocated to VGPU0, and given 40% of GPU reserve. This task (with a high priority) has 1.2 times the higher workload than the other, which was allocated to VGPU1 with 60% of GPU reserve. The VGPU1 task was scheduled to start approximately 5s after the VGPU0 task.

The result obtained under the FIFO scheduling policy on the NVIDIA Driver is shown in Figure 10 (a), while that under

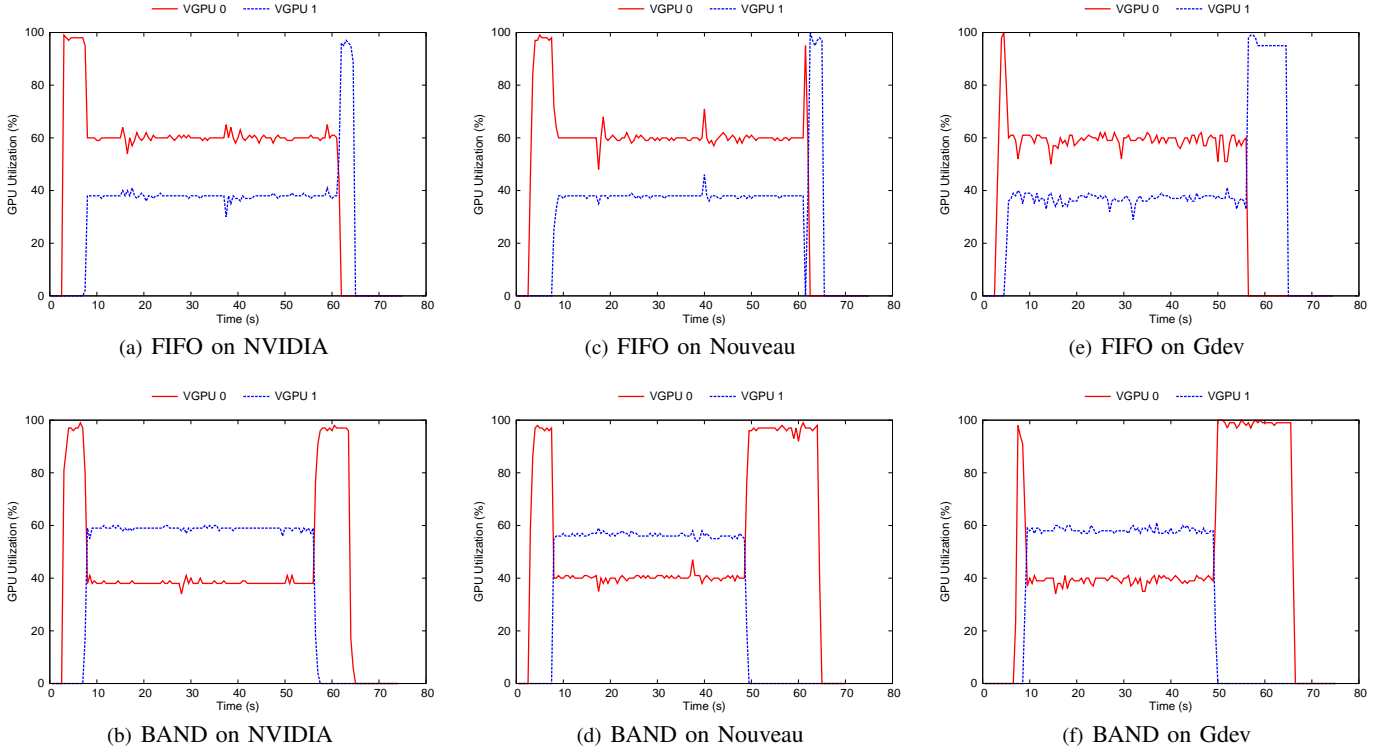


Fig. 10. GPU utilization of two tasks. Each task executes different types of GPU-intensive workloads with specified GPU reserves (VGPU0 = 40%; VGPU1 = 60%).

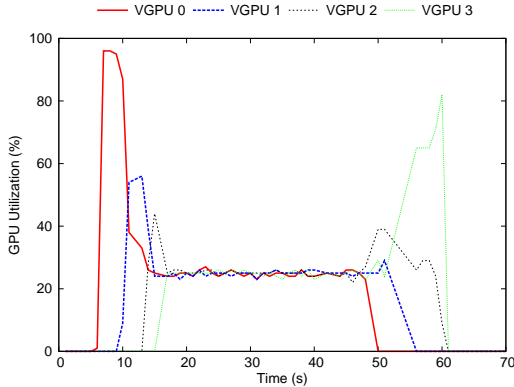


Fig. 11. Utilization of four tasks with the Linux-RTXG's BAND VGPU scheduling. Each task had an equal workload and an equal GPU-resource allocation.

the BAND scheduling policy is shown in Figure 10 (b). The corresponding results on the Nouveau GPU driver are shown in Figures 10 (c) and (d).

Since VGPU0 has a higher workload than VGPU1, the GPU tasks are performed in accordance with their workload, as shown in Figures 10 (a) and (c). On the other hand, the GPU resource-reservation mechanism provided by Linux-RTXG can successfully perform the GPU tasks in accordance with their reserves as indicated in Figures 10 (b) and (d). It is important to note that these prioritization and resource reservations for GPU applications can be achieved without modifying the OS kernel and device drivers in Linux-RTXG.

The maximum error of the VGPU1 task was approximately 3% under the BAND scheduling policy on the NVIDIA driver, while that of the VGPU0 task was approximately 5%. When

the Nouveau driver was used, those values were 2% and 6%, respectively. These large increases occurred due to GPU kernel overruns.

In addition to the performance of the proposed loadable-kernel approach, the performance was compared with that of a prior work, Gdev. The Gdev scheduling results are shown in Figures 10 (e) and (f). Almost no performance loss is imposed by Linux-RTXG as compared to Gdev. In detail, when the Gdev scheduler was used, the maximum error of the VGPU1 task was approximately 3% under the BAND scheduling policy, while that of the VGPU0 task was 5%. There is also a large variation on a time basis when the Gdev scheduler was used. This is because the runtime functions of Gdev must be called in the kernel space, where other system calls can easily block their operation.

GPU utilization running four GPU tasks was measured next. Each GPU task has the same workload and the same GPU reserve on a different VGPU. The isolation result for this scenario is shown in Figure 11. The maximum error of VGPU's was at most 9%, which is incurred due to the timing of budget replenishment and synchronization latency. In fact, this result almost matches that reported in previous work on Gdev [10] using the built-in kernel approach. As a result, the independent-synchronization mechanism employed in Linux-RTXG does not sacrifice scheduling performance.

F. Real-world Application

The performance of Linux-RTXG on a real-world application was demonstrated as follows. The demonstration was done with a GPU-accelerated object-detection program [2] based on the well-known DPM and HOG algorithms. It was assumed

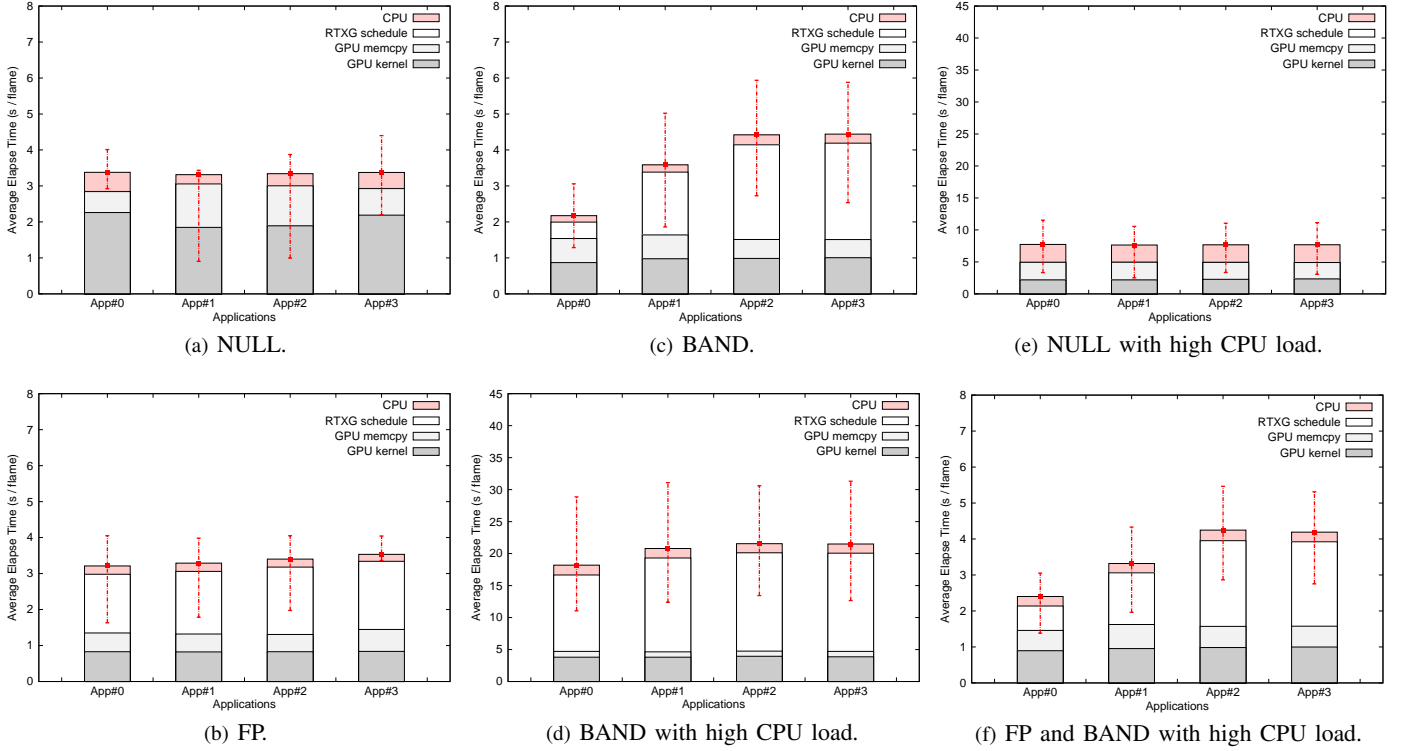


Fig. 12. Execution time of the GPU-accelerated object-detection program.

that the monitoring system covered the four directions of the compass (east, west, south, and north).

The execution time per frame was measured in six different setups, as shown in Figure 12, where (a) no scheduling was applied, (b) fixed priorities were applied with GPU scheduling, (c) GPU resource reservation was further applied with the BAND scheduling policy, and (d), (e), and (f) high CPU load was contending with the *SCHED_OTHER* hackbench task. The GPU tasks were allocated with 60%, 20%, 10%, and 10% of GPU reserves, respectively.

The results of this experiment indicate that if a high-CPU load task exists, the performance of GPU tasks could be seriously affected, since GPU tasks cannot acquire CPU time needed to execute the GPU API. Linux-RTXG can provide both CPU scheduling and GPU scheduling with resource-reservation capabilities to solve this problem, as shown in Figure 12 (f). As seen in the figure, the target programs are not affected by the high-CPU-load task thanks to the *SCHED_FIFO* scheduling, while GPU execution can be maintained at the desired frame rate thanks to the GPU-scheduling and resource-reservation capabilities. These results demonstrate that Linux-RTXG can supply prioritization and QoS performance to real-world applications.

V. RELATED WORK

RGEM [28] and GPU-Sparc [30] provide real-time GPU resource management without the OS kernel and device driver modifications. However, their synchronization mechanism depends on proprietary software. TimeGraph [9], Gdev [10], Ptask [38], and GPUSync [11] realize independent synchronization mechanisms but require modifying the OS kernel and device drivers. To the best of our knowledge, Linux-RTXG

is the only solution that can provide real-time GPU resource management with a synchronization mechanism, without modifying the OS kernel and device drivers.

Table II shows a comparison of Linux-RTXG and previous work. Note that Data/Computing Overlap indicates whether perform the data transfer and computing simultaneously, and Closed Source Compatible indicates whether use an NVIDIA CUDA. GPUSync supports the fixed-priority and the EDF scheduling policies for CPU tasks, while GPU-Sparc employs the *SCHED_FIFO* scheduling policy. As shown from this table, Linux-RTXG has demonstrated all features shown in Table II. In particular, the resource management modules of Linux-RTXG are all loadable and are freed from the detailed implementation of runtime libraries, device drivers, and the OS kernel.

More in-depth resource management would require detailed information about the execution mechanisms in the black-box GPU stack. Menychtas et al. presented interception-based OS-level GPU scheduling, called NEON, without modifying software stack by using state machine inferencing [39], [40]. GPU resource management using GPU micro-controllers [26] and in-kernel runtime functions [10] has also been demonstrated to manage the GPU. For these pieces of open-source work, the Nouveau project has been used as a baseline driver [36].

VI. CONCLUSION

Linux-RTXG, a Linux extension with the LKM-based real-time GPU-resource-management module, was proposed and evaluated. Linux-RTXG can manage GPU resources without having to modify the OS kernel and device drivers. In addition, it is compared to CPU-resource management proposed in

TABLE II. LINUX-RTXG VS PRIOR WORK

	CPU Fixed-Priority	CPU EDF	GPU Priority Scheduling	Budget Enforcement	Data/Computing Overlap	Closed Source Compatible	Patchless	GPU Runtime independent
RGEM			x			x	x	
Gdev			x	x	x			
PTask			x	x	x	x		x
NEON			x		x	x		x
GPUSync	x	x	x	x	x	x		x
GPU-Sparc	x		x		x	x	x	
Linux-RTXG	x	x	x	x	x	x	x	x

previous work [14], [15]. A new scheme for GPU-resource management was also developed. By intercepting the interrupts from the GPU in the top-half ISRs, the developed synchronization mechanism can wait for completion of execution of specific GPU kernels without the need of modifying the OS kernel and device drivers. Experimental evaluation of Linux-RTXG indicated the chosen performance unit (task overhead) varied within about 10% for eight tasks and within about 4% for four tasks. Furthermore, superior prioritization and QoS performance of Linux-RTXG were demonstrated with a real-world object-detection application, where GPU programs can be protected from high CPU load while the desired frame rate on the GPU is successfully maintained. To the best of our knowledge, this is the first complete work Linux extension that can manage GPU resources with loadable kernel modules.

REFERENCES

- [1] M. McNaughton, C. Urmson, J. Dolan, and J.-W. Lee, "Motion planning for autonomous driving with a conformal spatiotemporal lattice," in *Proc. ICRA*. IEEE, 2011, pp. 4889–4895.
- [2] M. Hirabayashi, S. Kato, M. Edahiro, K. Takeda, T. Kawano, and S. Mita, "GPU implementations of object detection using HOG features and deformable models," in *Proc. Int. Conf. CPSNA*. IEEE, 2013, pp. 106–111.
- [3] N. Rath, S. Kato, J. Levesque, M. Mauel, G. Navratil, and Q. Peng, "Fast, multi-channel real-time processing of signals with microsecond latency using graphics processing units," *Review of Scientific Instruments*, vol. 85, no. 4, p. 045114, 2014.
- [4] S. Kato, K. Lakshmanan, Y. Ishikawa, and R. Rajkumar, "Resource sharing in GPU-accelerated windowing systems," in *Proc. RTAS*. IEEE, 2011, pp. 191–200.
- [5] P. Bakkum and K. Skadron, "Accelerating SQL database operations on a GPU with CUDA," in *Proc. Workshop GPGPU*. ACM, 2010, pp. 94–103.
- [6] S. Che, M. Boyer, J. Meng, D. Tarjan, J. W. Sheaffer, S.-H. Lee, and K. Skadron, "Rodinia: A benchmark suite for heterogeneous computing," in *Proc. IISWC*. IEEE, 2009, pp. 44–54.
- [7] "CUDA Zone," <https://developer.nvidia.com/category/zone/cuda-zone>, accessed January 12, 2015.
- [8] J. E. Stone, D. Gohara, and G. Shi, "Opencl: A parallel programming standard for heterogeneous computing systems," *IEEE Des. Test. Comput.*, vol. 12, no. 3, pp. 66–73, 2010.
- [9] S. Kato, K. Lakshmanan, R. R. Rajkumar, and Y. Ishikawa, "Time-Graph: GPU scheduling for real-time multi-tasking environments," in *Proc. USENIC ATC*, 2011, p. 17.
- [10] S. Kato, M. McThrow, C. Maltzahn, and S. A. Brandt, "Gdev: First-class GPU resource management in the operating system," in *Proc. USENIC ATC*, 2012, pp. 401–412.
- [11] G. A. Elliott, B. C. Ward, and J. H. Anderson, "GPUSync: A framework for real-time GPU management," in *Proc. RTSS*. IEEE, 2013, pp. 33–44.
- [12] G. A. Elliott and J. H. Anderson, "Exploring the multitude of real-time multi-GPU configurations," in *Proc. RTSS*. IEEE, 2014.
- [13] J. Calandrino, H. Leontyev, A. Block, U. Devi, and J. Anderson, "*LITMUS^{RT}*: A testbed for empirically comparing real-time multiprocessor schedulers," in *Proc. RTSS*, 2006, pp. 111–126.
- [14] S. Kato, R. Rajkumar, and Y. Ishikawa, "A loadable real-time scheduler suite for multicore platforms," Tech. Rep. CMU-ECE-TR09-12, Tech. Rep., 2009.
- [15] M. Asberg, T. Nolte, S. Kato, R. Rajkumar, and Y. Ishikawa, "Exsched: An external CPU scheduler framework for real-time systems," Tech. Rep. CMU-ECE-TR09-12, Tech. Rep., 2009.
- [16] J. A. Stankovic and K. Ramamritham, "The spring kernel: A new paradigm for real-time systems," *IEEE Softw.*, vol. 8, no. 3, pp. 62–72, 1991.
- [17] T. Yang, T. Liu, E. D. Berger, S. F. Kaplan, and J. E. B. Moss, "Redline: First class support for interactivity in commodity operating systems," in *Proc. OSDI*, vol. 8, 2008, pp. 73–86.
- [18] H. Takada and K. Sakamura, "μITRON for small-scale embedded systems," *IEEE Micro*, vol. 15, no. 6, pp. 46–54, 1995.
- [19] R. Rajkumar, K. Juvva, A. Molano, and S. Oikawa, "Resource kernels: A resource-centric approach to real-time and multimedia systems," in *Photonics West'98 Electronic Imaging*. International Society for Optics and Photonics, 1997, pp. 150–164.
- [20] S. Oikawa and R. Rajkumar, "Portable RK: A portable resource kernel for guaranteed and enforced timing behavior," in *Proc. RTAS*. IEEE, 1999, pp. 111–120.
- [21] P. Mantegazza, E. Dozio, and S. Papacharalambous, "RTAI: Real time application interface," *Linux J.*, vol. 2000, no. 72es, p. 10, 2000.
- [22] V. Yodaiken, "The RTLinux manifesto," in *Proc. Fifth Linux Expo*, 1999.
- [23] C. L. Liu and J. W. Layland, "Scheduling algorithms for multiprogramming in a hard-real-time environment," *J. ACM*, vol. 20, no. 1, pp. 46–61, 1973.
- [24] L. Abeni and G. Buttazzo, "Integrating multimedia applications in hard real-time systems," in *Proc. RTSS*. IEEE, 1998, pp. 4–13.
- [25] M. Spuri and G. Buttazzo, "Scheduling aperiodic tasks in dynamic priority systems," *Real-Time Systems*, vol. 10, no. 2, pp. 179–210, 1996.
- [26] Y. Fujii, T. Azumi, N. Nishio, and S. Kato, "Exploring microcontrollers in GPUs," in *Proc. APSys*. ACM, 2013, p. 2.
- [27] G. A. Elliott and J. H. Anderson, "Robust real-time multiprocessor interrupt handling motivated by GPUs," in *Proc. ECRTS*, 2012, pp. 267–276.
- [28] S. Kato, K. Lakshmanan, A. Kumar, M. Kelkar, Y. Ishikawa, and R. Rajkumar, "RGEM: A responsive GPGPU execution model for runtime engines," in *Proc. RTSS*. IEEE, 2011, pp. 57–66.
- [29] C. Basaran and K.-D. Kang, "Supporting preemptive task executions and memory copies in GPGPUs," in *Proc. ECRTS*, 2012, pp. 287–296.
- [30] W. Han, H. Bae, H. Kim, J. Lee, and I. Shin, "Gpu-sparc: Accelerating parallelism in multi-GPU real-time systems," Tech. Rep., 2014.
- [31] F. Ridouard, P. Richard, and F. Cottet, "Negative results for scheduling independent hard real-time tasks with self-suspensions," in *Proc. RTSS*. IEEE, 2004, pp. 47–56.
- [32] K. Lakshmanan, S. Kato, and R. R. Rajkumar, "Open problems in scheduling self-suspending tasks," in *Proc. Real-Time Scheduling Open Problems Seminar*, 2010.
- [33] B. Chattopadhyay and S. Baruah, "Limited-preemption scheduling on multiprocessors," in *Proc. Int. Conf. RTNS*. ACM, 2014, p. 225.
- [34] J. Kim, B. Andersson, D. d. Niz, and R. R. Rajkumar, "Segment-fixed priority scheduling for self-suspending real-time tasks," in *Proc. RTSS*. IEEE, 2013, pp. 246–257.
- [35] Y. Fujii, T. Azumi, N. Nishio, S. Kato, and M. Edahiro, "Data transfer matters for GPU computing," in *Proc. ICPADS*. IEEE, 2013, pp. 275–282.

- [36] “Nouveau,” <http://nouveau.freedesktop.org/wiki/>, accessed January 12, 2015.
- [37] S. Kato, J. Aumiller, and S. Brandt, “Zero-copy I/O processing for low-latency GPU computing,” in *Proc. Int. Conf. ICCPS*. ACM, 2013, pp. 170–178.
- [38] C. J. Rossbach, J. Currey, M. Silberstein, B. Ray, and E. Witchel, “PTask: Operating system abstractions to manage GPUs as compute devices,” in *Proc. SOSP*. ACM, 2011, pp. 233–248.
- [39] K. Menychtas, K. Shen, and M. L. Scott, “Enabling OS research by inferring interactions in the black-box GPU stack,” in *Proc. USENIX ATC*, 2013, pp. 291–296.
- [40] —, “Disengaged scheduling for fair, protected access to fast computational accelerators,” in *Proc. ASPLOS*, 2014, pp. 301–316.