

# **Toward Efficient and Realizable Virtualization of Compute Accelerators**

Amogh Akshintala

A dissertation submitted to the faculty of the University of North Carolina at Chapel Hill in partial fulfillment of the requirements for the degree of Doctor of Philosophy in the Department of Computer Science.

Chapel Hill  
2020

Approved by:

Christopher J. Rossbach

Donald E. Porter

Fabian N. Monroe

Kevin Jeffay

Michael Ferdman

©2020  
Amogh Akshintala  
ALL RIGHTS RESERVED

## ABSTRACT

Amogh Akshintala: Toward Efficient and Realizable Virtualization of Compute Accelerators  
(Under the direction of Donald E. Porter, and Christopher J. Rossbach)

The proposed dissertation will focus on software virtualization of specialized compute devices (e.g., GPUs, TPUs) that are programmed through an API. Specialized compute accelerators, such as GPUs and TPUs, are usually controlled through a user-space API. The proposed dissertation will show that unlike with CPUs, where the ISA is the canonical interface exposed to the programmer, ISA virtualization is untenable for specialized compute accelerators. Further, the proposed dissertation will present a novel taxonomy, *IEMTS* for cleanly understanding the design space for virtualizing compute accelerators. Based on insights from this taxonomy, the proposed dissertation will present a novel virtualization technique, *Hypervisor-Interposed Remote Acceleration*, that is at once realizable and performant.

## ACKNOWLEDGEMENTS

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi. Lorem ipsum dolor sit amet,

## TABLE OF CONTENTS

LIST OF TABLES .....	vii
LIST OF FIGURES .....	viii
LIST OF ABBREVIATIONS .....	ix
1 Introduction .....	1
1.1 History of Virtualization .....	2
1.2 Hardware Virtualization .....	3
1.3 Domain Specific Accelerator Virtualization .....	4
2 Background .....	7
2.1 Virtualization Properties .....	7
2.2 Domain Specific Accelerators .....	8
2.2.1 DSA Design .....	8
2.2.2 DSA Software stack .....	8
2.3 DSA Virtualization .....	8
2.3.1 Inefficacy of Traditional GPGPU Virtualization Techniques .....	9
2.3.1.1 Pass-through .....	9
2.3.1.2 Device emulation .....	10
2.3.1.3 Full virtualization .....	10
2.3.1.4 Mediated pass-through .....	10
2.3.1.5 Para-virtualization .....	10
2.3.1.6 API Remoting .....	11
2.3.1.7 Hardware virtualization support .....	12

3	ISA virtualization is untenable for GPUs .....	14
3.1	Introduction.....	14
3.2	Background.....	15
3.2.1	SVGA .....	15
3.2.2	Mesa3D OpenCL Support.....	17
3.2.3	GPU ISAs and IRs.....	17
3.3	Design.....	18
3.4	Implementation .....	20
3.4.1	Trillium .....	20
3.4.2	GPUvm .....	22
3.4.3	User-space API remotng over RPC .....	22
3.4.4	Optimizations .....	23
3.5	Methodology .....	23
3.5.1	Benchmarks.....	23
3.5.2	Control Experiments.....	24
3.6	Evaluation .....	25
3.6.1	The impact of vISA choice.....	25
3.6.2	End-to-End.....	27
3.7	Conclusion .....	27
4	Hypervisor-mediated API-remoting.....	28
5	IEMTS — A new accelerator virtualization taxonomy .....	30
6	Proposed work — vTask .....	33
7	Related Work.....	36
7.1	GPU Virtualization .....	36
7.2	Language-level Virtualization .....	40
	BIBLIOGRAPHY .....	41

## LIST OF TABLES

3.1	EVALUATION BENCHMARKS IN THREE CATEGORIES <sup>a</sup> .....	24
5.1	Comparing virtualization designs using the IEMTS framework.....	31
7.1	Existing GPU virtualization proposals, grouped by approach. Previously published in the Trillium paper [30]. .....	37

## LIST OF FIGURES

2.1	Unfairness in slowdown between needle and hotspot applications in separate VMs running GPU kernels iteratively with BitFusion FlexDirect. When running alone, hotspot has throughput of 126.3 ms/kernel. Fairness is calculated by $ s_1 - s_2  / (s_1 + s_2)$ , where $s_i$ is the slowdown of application $i$ when running concurrently. <a href="#">[REMOVE AVA FROM THIS FIGURE.]</a> .....	11
2.2	Throughput achieved by three instances of QATzip (running in VMs with SR-IOV pass-through) with different block sizes, running separately ( <b>Uncontended</b> ) and concurrently ( <b>Contended</b> ). Slowdown during concurrent execution is dependent on block size, i.e., the QAT HW scheduler cannot guarantee fairness. ....	13
3.1	The design of SVGA. ....	16
3.2	XEN-SVGA and TRILLIUM designs. (a) The TRILLIUM stack. (b) XEN-SVGA approximates the SVGA model extended to support GPU Compute. (c) The design of TRILLIUM with shadow pipe. ....	19
3.3	Xen-based virtualization designs. (a) Trap-based virtualization: GPUvm. (b) User-space API remoting over RPC—dashed arrows indicate API-REMOTE-CPU, while solid ones indicate API-REMOTE-GPU.....	20
3.4	End-to-end execution times of benchmarks on virtualization prototypes, relative to end-to-end execution time on the NVIDIA CUDA runtime in a native setting. The gRPC transport overhead is removed from the reported measurements, which is up to 10% of the total execution time for API remoting, and 40% for TRILLIUM. ....	25
3.5	Kernel execution slowdown due to virtual ISAs. TGSI: the LLVM TGSI back-end compiler used in XEN-SVGA. LLVM: LLVM NVIDIA PTX (NVPTX) back-end used in TRILLIUM. No IR: native NVIDIA compiler. ....	26
4.1	An accelerator silo. The public API and the interfaces with striped backgrounds are interposition candidates. All interfaces with backgrounds are proprietary and subject to change. ....	28
6.1	Data processed by two API stacks must pass through the guest application .....	33



## **LIST OF ABBREVIATIONS**

DSA	Domain Specific Architectures
I/O	Input/Output
IPC	Inter-Process Communication
IPI	Inter-Processor Interrupt
AvA	Accelerated Virtualization of Accelerators
AYO	Add Your Own in alphabetic order. . .

## CHAPTER 1: INTRODUCTION

This dissertation is concerned with fair, efficient and safe sharing of domain specific accelerators between mutually distrustful users.

Domain Specific Accelerators (DSAs) are programmable compute units that are specialized to a particular class of computation in order to improve performance, or to optimize energy usage (or both) for that class of computation. Domain Specific Accelerators are seeing wide adoption in data centers as they afford the data center provider a very low cost/performance on the specific computation they are specialized for. For example, a 2013 projection by Google showed that they would need to double the installed CPU compute capacity in order to support just three minutes of Google Voice Search per user per day, using speech recognition DNNs [69]. This realization led Google to design and adopt the Tensorflow Processing Unit (TPU), a DSA specialized for the Tensor-based computation popular in Neural Networks. The very first generation of TPUs, which were deployed to Google data centers in 2015, were empirically found to have  $200\times$  and  $79\times$  higher Performance/Watt respectively over the CPUs and Nvidia k80 GPUs that were prevalent in their data centers at the time [70].

Data centers themselves are proliferating rapidly: a recent report by Synergy Research [111] found that there were more than 500 hyperscale data centers operational across the world at the end of 2019. This proliferation of data centers is primarily motivated by the increased need for computation resulting from the unquestionable permeation of digital services into every facet of our lives (e.g., Google Maps, video streaming, digital communication). Data centers leverage the cost efficiency of centralized computing, while using techniques such as virtualization to provide the flexibility of dedicated/distributed computing.

Virtualization is a cornerstone technology in our current computing landscape. Unlike physical compute, virtualized compute resources can be securely multiplexed among many users, which allows for high utilization of the installed compute capacity. Further, virtualization provides customers with the ability to scale elastically: when the demand for a digital service is high, the customer may

request and utilize a large number of machines; when the spike in demand subsides, the customer can just release the excess compute resources, thereby only paying for what they are actively using.

Virtualization of domain specific accelerators remains an open problem: DSAs are dedicated to single users instead of being multiplexed in today’s data centers. Cloud providers such as Amazon expose Graphical Processing Units (GPU) and Field Programmable Gate Arrays (FPGA) to individual VMs via PCIe-passsthrough, thereby bypassing the hypervisor, and giving up the consolidation and fault tolerance benefits of virtualization. This dissertation explores the trade-offs in the design space for software-based virtualization schemes for DSAs and proposes a low-overhead novel virtualization scheme that interposes the user-facing Application Programmer’s Interface (API) while using automation to compensate for the resulting development complexity.

## 1.1 History of Virtualization

The story of virtualization is long and tumultuous, stretching from the very early days of computing, right to today. Memory was the first resource to be virtualized: German physicist Fritz-Rudolf Güntsch described the basic idea of virtual memory in his doctoral dissertation in 1956 [62]. The Cambridge University/Ferranti Inc. Atlas [74] computer was the first to openly commercialize virtual memory. Virtually all computers since then have supported virtual memory, with most providing hardware units— *Memory Management Unit (MMU)*—to accelerate the virtualization of memory.

The next era dealt with the virtualization of the individual machine with its processor and peripherals. *Hardware virtualization* [27]—the idea of virtualizing the entire computer to enable the simultaneous execution of multiple Operating Systems (OS)—was invented in 1962, and commercialized in the IBM VM-370 [46] hypervisor for the IBM 370 computer. Virtualization was briefly forgotten through the 1980s and 1990s, as the mainframe computer became all but obsolete during the Personal Computer (PC) revolution. Intel’s x86 Instruction Set Architecture (ISA), which came to dominate the PCs that transplanted the mainframes, was not designed to be traditionally virtualizable [92], and was widely considered unvirtualizable [55, 41]. Multiple vendors introduced *software emulation* based solutions to enable the execution of one OS on top of another (e.g., Insignia SoftPC, Connectix VirtualPC, VMware Workstation). Over time, better techniques were devised to virtualize the x86 ISA, including ISA extensions (e.g., AMD-V and Intel VT-x) to enable native

execution of virtualized applications, only trapping to the hypervisor when the application attempts to perform sensitive operations.

Hardware interfaces weren't the only ones virtualized, and in fact, most computing environments today rely on hardware virtualization and one or more of the following *software virtualization* schemes working in tandem. Sun Microsystems popularized *application virtualization* in the 1990s with the Java programming language: applications are written to an abstract machine—the *Java Virtual Machine (JVM)*—which is backed by a runtime system that ensures the program can execute on any platform. This scheme eschews *compatibility*—the ability to execute unmodified legacy applications—for *portability*. *Operating system-level virtualization* (e.g., Library OSes, Containers) virtualizes yet another layer in the software stack: the operating system's interfaces (e.g., system calls, kernel name-spaces). This style of virtualization preserves compatibility by transparently modifying the interfaces the application uses to access system resources, and results in low overhead execution.

## 1.2 Hardware Virtualization

This dissertation is primarily concerned with hardware virtualization of domain specific accelerators. Hardware virtualization is vital to high utilization of available physical resources in large computing installations, e.g., hardware virtualization is foundational to *cloud computing*. There have been many attempts to define the fundamental characteristics of hardware virtualization: Popek and Goldberg came up with three properties that a virtualization scheme must exhibit—*equivalence, performance, and safety*, while Bugnion, Nieh and Tsafirir [42] provided a more succinct definition—“*the application of the layering principle with enforced modularity such that the exposed resource is identical to the underlying resource*”. For the purposes of this dissertation, we concern ourselves with *realizable, fair, isolated, and efficient sharing of domain specific accelerators among mutually distrustful entities*.

Hardware virtualization typically involves mediating access to the shared resource either by exposing an interface that is identical to that of the physical resource (*full-virtualization*), or by exposing an alternative interface, operations on which are in-turn synthesized to the native interface (*para-virtualization*). The exposed interface is considered *virtual*, as it is not controlled by the physical underlying hardware, and instead is entirely under the control of supervisory virtualization

software, the *hypervisor* (also known as the *Virtual Machine Monitor*). While operations in the resulting *virtual machine* may be directly executed on the physical hardware for improved performance, as in the case of hardware-assisted virtualization schemes like AMD-V and Intel VT-x, all privileged operations must still trap to the hypervisor.

Four decades of attention from both the academic community and industry has given rise to a large body of techniques that enable efficient virtualization of CPUs. Software techniques, such as binary translation and device emulation, are well established, and are still used in practice due to lower overhead for sequences of sensitive instructions that need to be emulated [29]. Dominant ISAs (e.g., x86 and ARM) provide extensions to enable low-overhead virtualization (e.g., Intel VT-x, AMD-V). Both of these are foundational building blocks for cloud computing [1].

### 1.3 Domain Specific Accelerator Virtualization

Virtualizing Domain Specific Accelerators is a delicate act of balancing the essential characteristics of a virtualization scheme—compatibility, interposition, sharing, isolation—with the need to preserve the raw performance these processors provide. Virtualization techniques developed for CPUs (ISA virtualization) can not be applied to these accelerators: their control interfaces are closer to those of I/O devices than the ISAs of CPUs. Techniques developed for I/O devices, such as NICs, are also untenable for DSAs as they sacrifice one or more essential virtualization characteristics. Full-virtualization based schemes (e.g., GPUvm [108]) suffer from massive overheads that essentially negate the speedup that makes the domain specific accelerator attractive in the first place. Para-virtual systems that interpose on low-level interfaces such as the kernel driver (e.g., SVGA [51]), introduce much lower overhead than full-virtualization based schemes but have poor compatibility. The introduction of an artificial abstract interface constructed expressly for the purpose of interposition necessitates massive engineering effort to support new hardware in the host and new software frameworks in the guest. User-space API-remoting solutions [119, 53, 95] interpose on the user-space API in the guest and forward the interposed operation to the host as an RPC. This approach introduces very low overhead and can evolve with the hardware easily, but has traditionally eschewed hypervisor interposition, thereby making it difficult to enforce safety and isolation among guests.

Virtualizing a Graphics Processing Unit (GPU) for the purposes of graphics rendering is a well studied problem, with existing commercial solutions (e.g., VMware’s SVGA [51]). Over the last decade, GPUs have also been re-purposed for general purpose compute (commonly known as GPGPU). Chapter 2 of this dissertation presents background on domain specific accelerators, how their software stacks are different from those of CPUs and I/O devices, and how and why previous software virtualization solutions do not fare well for DSAs. In order to understand the trade-offs with each of the canonical virtualization techniques when applied to GPGPU virtualization, we present an end-to-end evaluation of representative systems. Chapter 2 also considers the notion of the modern DSA stack being a proprietary silo, i.e., that it’s only stable and publicly available interfaces are at the top and the bottom. Later chapters build on this notion of silo-ed DSA software stacks to design an effective virtualization scheme.

Chapter 3 of this dissertation presents our findings from attempting to extend the SVGA model of GPU virtualization to virtualize GPGPUs. We find that the tight coupling between ISA virtualization and device virtualization in SVGA leads to poor performance for GPGPU compute. Eliding ISA virtualization enables the resulting prototype, Trillium [30], to outperform all other traditional virtualization schemes that retaining hypervisor interposition. However, Trillium is still 2—3× slower than RPC-style API forwarding.

Domain Specific Accelerators (e.g., Google TPU, Intel QAT, etc.) are typically exposed to developers via a user-space API. The API is typically implemented by proprietary software that interacts with the hardware through opaque interfaces. Chapters 4, 5 and 6 generalize the lessons learned for GPGPUs to all API-controlled domain specific accelerators by interposing their user-space APIs. Chapter 4 presents an overview of AvA, a framework that enables automated virtualization of accelerator APIs. AvA combines on a novel virtualization scheme called Hypervisor Interposed Remote Acceleration (HIRA), with automation based on a Domain Specific Language, LAPIS, which is used to capture semantic information of the interposed APIs. This dissertation is primarily concerned with the HIRA virtualization scheme. Chapter 5 focuses on the performance implications of API-remoting based virtualization of a single specialized accelerator. Chapters 4 and 5 draw on material that appeared in a HotOS workshop paper [128] and an ASPLOS’20 paper [129]. Chapter 6 explores performance issues that arise when an application uses multiple API-remoted virtual accelerators in a pipelined fashion.

Virtualization schemes are traditionally taxonomized according to the core techniques employed (e.g. emulation, full- or para-virtualization, API remotng, etc.), and evaluated in a property trade-off space comprising performance, compatibility, interposition, and isolation. We argue that both the de facto taxonomy and the property trade-off space are illustrative but not informative for GPGPU virtualization: there is a large body of research that has had little influence on practice. We suggest an alternative framework called **IEMTS** that teases apart design axes that are implicitly and unnecessarily intertwined in much of the literature. By focusing on the **I**nterface interposed, the interposition **E**ndpoints, the **M**echanism of interposition, the **T**ransport used to move the interposed operations between the guest and the host, and the mechanism used to **S**ynthesize the interposed interface, **IEMTS** enables a clearer understanding of trade-offs in prior designs and provides a model for comparison of alternative designs. **IEMTS** will be presented in Chapter 7 of the proposed dissertation, along with analysis of traditional virtualization techniques in the context of GPGPUs.

## CHAPTER 2: BACKGROUND

Bugnion, Nieh and Tsafirir [42] define virtualization as “*the application of the layering principle with enforced modularity such that the exposed resource is identical to the underlying resource*”. To put it in simpler words, virtualization is about controlling the interface to a hardware resource (*interposition*) in order to multiplex it among multiple users in a safe manner (*isolation*), without any of them being the wiser (*compatibility*). Virtualization is a huge area with a long and storied history, as alluded to in the introduction; see *Bugnion, Nieh and Tsafirir’s Hardware and software support for virtualization* [42] for comprehensive treatment.

### 2.1 Virtualization Properties

Given our focus accelerator virtualization, let us consider the following key properties: *interposition*, *compatibility*, and *isolation*.

**Interposition.** Virtualization decouples a logical resource from a physical one through an indirection layer, intercepting guest interactions with a virtual resource and providing the requested functionality using a combination of software and the underlying physical resource. Thus, virtualization works by *interposing* an interface, and changing or adding to its behavior. Interposition is fundamental to virtualization and provides well-known benefits [121]. The choice of interface and the mechanism for interposing it profoundly impacts the resulting system’s practicality. *Inefficient* interposition of an interface (e.g. trapping frequent MMIO access) undermines performance [108, 130]; *incomplete* interposition compromises the hypervisor’s ability to enforce isolation.

**Compatibility** as applied to virtualization captures multiple related dimensions, from robustness to evolution of interposed interfaces and adjacent stack layers, to applicability across multiple platforms or related devices. For example, full virtualization of a accelerators’s hardware interface has *poor* compatibility in that it works only with that device. However, it has *good* compatibility with guest software, which will work without modification, assuming the operating system has appropriate



drivers for the device. Current accelerator virtualization techniques reflect a compromise between these two forms of compatibility.

**Isolation.** Cross-VM isolation is a critical requirement for multi-tenancy: when a resource is multiplexed among mutually distrustful tenants, tenants must not be able to see/alter each other’s data (*data isolation*), or adversely affect each other’s performance (*performance isolation*). A poor choice of interposition mechanism and/or interface limits the system’s ability to provide these guarantees: e.g., API remotng [3, 53, 13] has poor isolation in the common case, as the hypervisor is bypassed. Using separate servers for each protection provides isolation.

## 2.2 Domain Specific Accelerators

Domain Specific Accelerators (DSAs) are programmable compute units that are specialized to a particular class of computation in order to improve performance, and to optimize energy usage for that class of computation. The slowing down of Moore’s law coupled with the rise of Dark Silicon [54] has made Domain Specific Accelerators extremely attractive as they exhibit high computation/Watt efficiency in the computation domain they are specialized to. For example, consider Google’s Tensorflow Processing Unit (TPU) [69], a DSA for the Tensor-based computation popular in Neural Networks. The first generation of TPUs were empirically found to have  $200\times$  and  $79\times$  higher Performance/Watt respectively over the CPUs and Nvidia k80 GPUs that were prevalent in Google’s data centers at the time [70].

### 2.2.1 DSA Design

...

### 2.2.2 DSA Software stack

...

## 2.3 DSA Virtualization

Despite mounting evidence of accelerator under-utilization [17, 122, 88, 89, 16], and abundant prior research into multiplexing Domain Specific Accelerators (DSAs) [89, 16, 133, 122, 88, 126], DSAs

remain dedicated exclusively to a single guest in shared computing environments. This section provides background on the well studied, but mostly unresolved problem of GPGPU virtualization to explain this trend.

Existing GPU virtualization solutions [51, 80] support graphics frameworks like Direct3D [39], OpenGL [99]. In principle, there should be no fundamental difference between GPU virtualization for graphics versus *compute* workloads, as “compute shaders” are implemented by the hardware as an additional stage in the graphics pipeline [5]. In practice, they have significantly different goals: for graphics, virtualization designs target an interactive frame rate (18-30 fps [21]); for GPGPU compute, virtualization designs must preserve the raw speedup achieved by the hand-optimized GPGPU application, which is a considerably harder target to hit. As a result, GPGPU virtualization remains an open problem. While graphics devices have long enjoyed well-defined OS abstractions and interfaces [86], research attention to OS abstractions for GPGPUs [97, 98, 104, 71, 72, 76] has yielded little consensus. Persistent vendor-specificity of programming frameworks further impedes both interposition and compatibility.

### 2.3.1 Inefficacy of Traditional GPGPU Virtualization Techniques

An ideal GPGPU virtualization design would require no modification of guest applications, libraries and OSes (compatibility), arbitrate fair and isolated sharing of GPU resources between mutually distrustful VMs (sharing and isolation) at the native performance of the hardware (performance), while allowing virtualized software and physical hardware to evolve independently (encapsulation). We briefly describe each technique, and look at their strengths and shortcomings in this section. Refer to Related Work 7 for details on individual prior work under each technique.

#### 2.3.1.1 Pass-through

PCIe pass-through, the current *de facto* standard technique for GPGPU virtualization, provides a VM with full exclusive access to a physical GPU. The GPU’s hardware interface is directly exposed to the guest OS, and therefore can’t be multiplexed as the hypervisor does not interpose *any* interface. Virtualization hardware extensions (e.g., Intel VT-d [26]) are device-agnostic, making PCIe pass-through easily adaptable to any DSA. Pass-through provides native performance at the cost of *sharing, interposition, compatibility and isolation*.

### 2.3.1.2 Device emulation

Device Emulation [37] provides a full-fidelity software-backed virtual device which yields excellent compatibility, interposition, and isolation. However, device emulation can't support hardware acceleration making it untenable for virtualizing GPGPUs.

### 2.3.1.3 Full virtualization

The hypervisor interposes GPU's hardware interface to provides a virtual environment in which unmodified GPGPU programs run on unmodified guest software stacks. For DSAs, this interface tends to be memory mapped I/O (MMIO), necessitating trap-based interposition (e.g. using memory protection or de-privileging), leading to devastating performance slowdowns (e.g.,  $100\times$  slowdown with GPUvm [108, 130]). DSA hardware interfaces tend to be proprietary and device-specific, so full virtualization based solutions have poor compatibility, even across different devices of the same type (e.g. AMD vs NVIDIA GPUs). Full virtualization solutions also typically rely on reverse engineering of proprietary control interfaces, rendering them extremely tedious to build, maintain and evolve.

### 2.3.1.4 Mediated pass-through

A hybridization of pass-through and full virtualization, Mediated pass-through [113, 90, 125] uses pass-through for data plane operations, and provides a privileged control plane interface for sensitive operations. Mediated pass-through can preserve some of the raw speedup of acceleration and allows guests to use native drivers and libraries. However, limited interposition limits a hypervisor's ability to effectively manage resource sharing. More importantly, hardware support is required. To our knowledge, Intel integrated GPUs are the only accelerators with such support.

### 2.3.1.5 Para-virtualization

Rather than interposing an existing interface in the stack, para-virtualization [108, 51, 115, 81, 63, 58, 85, 93, 107, 123, 33] *creates* an efficiently interposable interface in software and adjusts adjacent stack layers to use it. The driver and runtime libraries in every supported OS must be modified to work in concert with the virtualization layer. Para-virtualization enables encapsulation of diverse hardware behind a single interposable interface, but compromises compatibility. Guest software

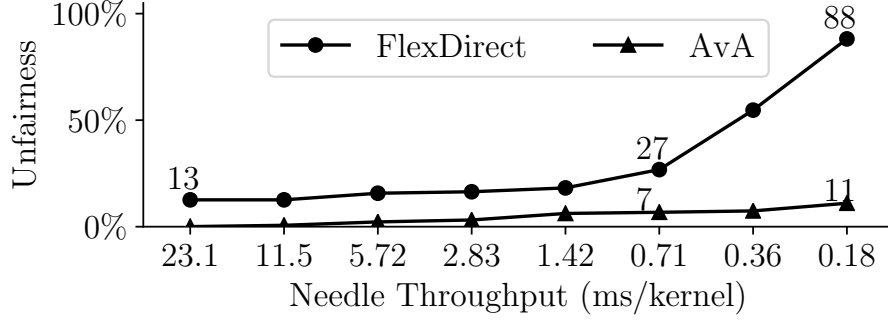


Figure 2.1: Unfairness in slowdown between needle and hotspot applications in separate VMs running GPU kernels iteratively with BitFusion FlexDirect. When running alone, hotspot has throughput of 126.3 ms/kernel. Fairness is calculated by  $|s_1 - s_2| / (s_1 + s_2)$ , where  $s_i$  is the slowdown of application  $i$  when running concurrently. [REMOVE AVA FROM THIS FIGURE.]

must be modified, and the para-virtual device interface must be maintained as interfaces evolve. For example, VMware’s SVGA II [51] encapsulates multiple GPU programming frameworks, but keeping up with the evolution of those frameworks has proved untenable: SVGA remains multiple versions behind current frameworks [20, 18].

### 2.3.1.6 API Remoting

User-space API Remoting based virtualization designs interpose application-level APIs (e.g. CUDA, OpenCL) by shimming a dynamic library and remote them to the corresponding framework in the host [101, 61, 56], on a dedicated appliance VM [119], or on a remote server [53, 95, 84, 75, 35, 52, 83]. API Remoting is similar to RPC [?] or system call interposition [33, 110, 34, 79]. Limited interposition frequency, batching opportunities [53] and high-speed networks [6, 3] reduce overheads, making this class of designs appealing to industry. Dell XaaS [67], BitFusion FlexDirect [3], and Google Cloud TPUs [9] currently use it to support GPUs, FPGAs, and TPUs. However, API remoting compromises compatibility if multiple APIs or API versions must be supported. Moreover the technique bypasses the hypervisor, giving up the interposition required for hypervisor-enforced resource management. Our experiments with commercial systems like BitFusion FlexDirect [3] show vulnerability to massive unfairness pathologies. Figure 2.1 shows the problem on an NVIDIA GTX 1080: FlexDirect is unfair (up to 88.1%) when running two applications with different kernel run-lengths (126.3 ms/kernel vs 0.18 ms/ kernel in the worst case).

Deferring enforcement to a trusted surrogate in the host or remote machine is a tenable alternative. However, the co-ordination required to integrate with hypervisor-level resource management means that current solutions do not support it, and the engineering effort required would be substantial. Existing accelerator API remoting systems are by themselves massive undertakings without any hypervisor integration: systems like Bitfusion FlexDirect [3], and rCUDA [53] reflect multi-year system-building efforts.

### 2.3.1.7 Hardware virtualization support

Hardware support for virtualization (e.g., Single Root I/O Virtualization (SR-IOV)) enables a single physical device to present itself as multiple virtual devices. A hypervisor can manage and distribute these virtual devices to guests, effectively deferring virtualization, scheduling, and resource management to the hardware. NVIDIA and AMD both ship GPU cards targeted at the VDI market that use SR-IOV to export multiple virtual GPUs from the hardware.

SR-IOV exhibits close to native performance [49], but this is achieved at the cost of interposition — the hypervisor can’t interpose on any interactions with the hardware. SR-IOV also suffers from the multiple administrator problem: the hardware controller and the hypervisor/OS may make mutually inconsistent decisions leading to unpredictable behavior. SR-IOV provides an interface and protocol for managing VFs, but the device vendor must *implement* any cross-VF sharing support *in silicon*. The technique can provide strong virtualization guarantees [48, 50], but hardware-level resource management is inflexible and slow to evolve: current implementations are trivially vulnerable to fragmentation and unfairness pathologies that cannot be changed.

Hardware designers tend to favor simple resource management policy implementations, easily leading to pathologies. To illustrate the problem, we measured compression throughput when three VMs contend on an Intel QuickAssist [?] (with SR-IOV). The three VMs configured the accelerator to compress at different chunk sizes. Figure 2.2 shows the results of this experiment, with and without contention. Each VM was assigned a PCIe Virtual Function (VF) exposed by the same Physical Function (PF), causing the hardware to schedule requests round-robin. When there is no contention, each application achieves a similar throughput. However, when the 3 applications were executed concurrently, the throughput achieved was a function of offload chunk size used, *yielding unfairness that cannot be fixed without changing the hardware*.

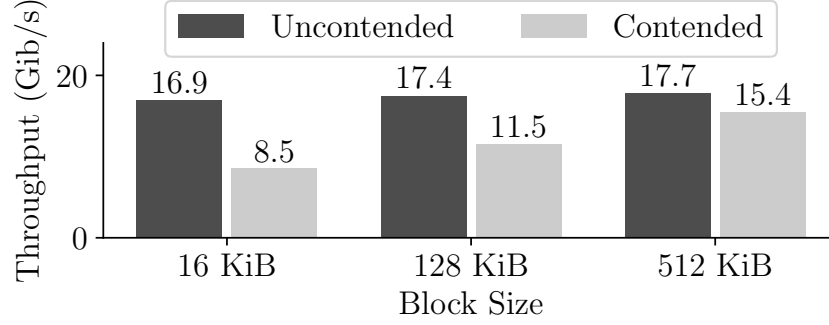


Figure 2.2: Throughput achieved by three instances of QATzip (running in VMs with SR-IOV pass-through) with different block sizes, running separately (**Uncontended**) and concurrently (**Contended**). Slowdown during concurrent execution is dependent on block size, i.e., the QAT HW scheduler cannot guarantee fairness.

Further, evidence is scant that broad SR-IOV support will emerge for accelerators: only two current GPUs support it [2, 64], none of the TPUs we evaluate support it; and SR-IOV *interface* IP blocks from FPGA vendors (used by [118, 131, 91, 65]) do not implement resource management.

While hardware support for virtualization is the desired end game, we do not expect a better standard for such support to emerge soon: incentives for investing in the significant engineering effort required for such a standard are scarce. This dissertation focuses on virtualization schemes that can be purely implemented in software so as to sidestep this challenge.

## CHAPTER 3: ISA VIRTUALIZATION IS UNTENABLE FOR GPUS

### 3.1 Introduction

In many parallel computing domains, compute density and programmability [14, 106, 60] have made GPUs the clear choice for efficiency and performance [10]. Popular machine learning frameworks such as Caffe [68], Tensorflow [25], Microsoft CNTK [127], and Torch7 [45] rely on GPU acceleration heavily. GPUs have made significant inroads in HPC as well: five of the top seven supercomputers in the world are powered by GPUs [24].

Despite much prior research [120, 66, 28, 116] on GPGPU virtualization, practical options currently available to providers of virtual infrastructure all involve bypassing the hypervisor. The most commonly adopted technique is to dedicate GPUs to single VM instances via PCIe pass-through [32, 113], thereby giving up the consolidation and fault tolerance benefits of virtualization. More recently, industry players such as VMware, Dell and BitFusion have introduced user-space API-remoting [38, 75, 95, 119, 53] based solutions as an alternative to pass-through. API-remoting recovers the consolidation and encapsulation benefits of virtualization but bypasses hypervisor interposition. The absence of hypervisor interposition results in multiple disjoint resource managers (the remote user-space API executor and the hypervisor) with no insight into each others' decisions, thereby leading to poor decision making, and priority-inversion problems [97].

To recover hypervisor interposition while maintaining low-overhead, we retrofit GPGPU support into a virtual GPU device: We added support for OpenCL to an implementation of the SVGA [?] (see § 3.2.1) design in Xen, by implementing the key missing component—a compiler for SVGA's TGSi virtual ISA. This effort helped us realize that because GPUs already support vendor-specific virtual ISAs (vISAs), the additional vISA provides little benefit. In fact, we found that it harms performance by necessitating a translation layer that obscures the program's semantic information from the final vendor-provided compiler. Drawing on this lesson, we adapted TRILLIUM to take a more flexible approach to ISA virtualization: eliding it entirely when the host GPU stack bundles a

compiler (most do), and using LLVM IR, when necessary, to provide a common target for GPGPU drivers.

TRILLIUM represents an unexplored point in the GPGPU virtualization design space: hypervisor-mediated API-remoting. TRILLIUM is an existence proof of a viable alternative design that preserves desirable virtualization properties such as consolidation, hypervisor interposition, isolation, encapsulation, etc., without requiring full hardware virtualization. TRILLIUM outperforms a full virtualization system from the literature by up to  $14\times$  ( $5.5\times$  on average) and outperforms the para-virtual SVGA-like design by as much as  $7.3\times$  ( $5.4\times$  on average).

## 3.2 Background

Existing GPU virtualization solutions [51, 80] support graphics frameworks like Direct3D [39], OpenGL [99]. In principle, there should be no fundamental difference between GPU virtualization for graphics versus *compute* workloads. In practice, they have significantly different goals: For graphics, virtualization designs target an interactive frame rate (18-30 fps [21]). For GPGPU compute, virtualization designs must preserve the raw speedup achieved by the hand-optimized GPGPU application, which is a considerably harder target to hit. As a result, GPGPU virtualization remains an open problem. While graphics devices have long enjoyed well-defined OS abstractions and interfaces [86], research attention to OS abstractions for GPGPUs [97, 98, 104, 71, 72, 76] has yielded little consensus.

### 3.2.1 SVGA

SVGA [51] remotes DirectX and OpenGL over an emulated (software) PCIe device. The SVGA virtual device behaves like a physical GPU, by exporting virtual resources in the form of registers, extents of guest memory accessible to the virtual device, and a command queue. I/O registers (used for mode switching, IRQs, memory allocation) are mapped in an interposed PCIe Base Address Register (BAR) to enable synchronous emulation. Access to GPU memory is supported through asynchronous DMA. Figure 3.1 presents an overview of SVGA.

SVGA combines many aspects of full-, para-virtual and API remoting designs. Unmodified guests can transparently use SVGA as a VGA device, making full virtualization possible where



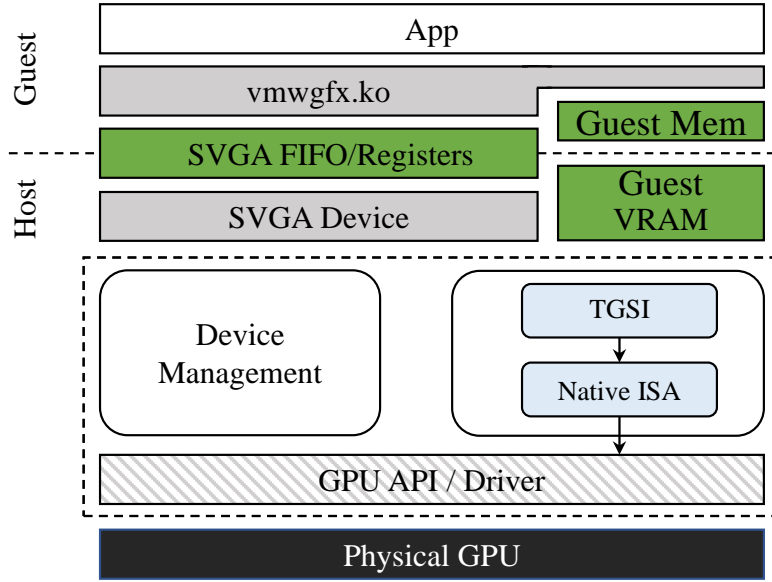


Figure 3.1: The design of SVGA.

necessary. However, access to GPU acceleration requires para-virtualization through VMware’s guest driver. SVGA processes commands from a memory mapped command queue; the command queue functions as a transport layer for protocols between the guest graphics stack and the hypervisor.

SVGA uses the DirectX [39] API as its internal protocol, thereby realizing an API-remoting design. The transport layer and protocol are completely under the control of the hypervisor, enabling many of the benefits of API-remoting while ameliorating its downsides. However, using the DirectX API as a transport protocol requires that the driver and hypervisor translate guest interactions into DirectX whether they are natively expressed in DirectX or not. Coupling the transport layer with a particular version of the DirectX protocol has led to serious complexity and compatibility *challenges*: supporting each new version of the API takes many person-years (VMware introduced support for DirectX 10 (introduced in 2006) in 2015!). SVGA also supports a virtual GPU ISA called TGSi [117]. TGSi maps naturally to the graphics features of the ISAs it was originally designed to encapsulate, but has failed to keep up with GPU ISAs that have evolved to support general purpose computation primitives.

### 3.2.2 Mesa3D OpenCL Support

The Mesa3D Graphics Library [23] is an open-source graphics framework that implements graphics runtime libraries (e.g., OpenGL [99], Vulkan [73], Direct3D [39], and OpenCL [106]) on most GNU/Linux installations. It also includes official device drivers, written in a common framework, Gallium3D [22], for Intel and AMD GPUs. Support for NVIDIA GPUs is provided via reverse-engineered open-source Nouveau driver. Gallium3D imposes TGSI as the common virtual ISA for compute shaders, and decomposes drivers into two components: *state trackers*, which keep track of the device state, and *pipe drivers*, which provide an interface for controlling the GPU’s graphics pipeline.

OpenCL support was first introduced in Mesa3D 9.0 with the release of the Clover state tracker. It was envisioned that Clover would leverage the LLVM [82] compiler to lower the OpenCL source to TGSI. Despite much effort by the open-source community [7, 11], an LLVM TGSI back-end has remained incomplete. Clover currently supports an incomplete set of OpenCL 1.1 APIs on AMD GPUs and fails to operate correctly on NVIDIA GPUs.

### 3.2.3 GPU ISAs and IRs

GPU front-end compilers produce code in virtual ISAs (NVIDIA PTX and LLVM IR for AMD) which are subsequently finalized using JIT compilers in the GPU driver to the native ISA (SASS and GCN). The vISA remains stable across generations to preserve compatibility, while the physical ISA is free to evolve. TGSI, the virtual ISA used in both the Mesa stack and SVGA, plays a similar role—enabling interoperability between graphics frameworks and GPUs from different vendors. An improved virtual ISA, SPIR-V, has been proposed as a new standard [73] and an effort is under way to replace TGSI with SPIR-V in the Mesa3D stack.

LLVM has become the de-facto standard for building compilers: both NVIDIA and AMD use it to implement their virtual ISA compilers, as do all the compilers in the Mesa stack including the TGSI compiler we implemented. LLVM IR is in a unique position to become a standard IR.

### 3.3 Design

TRILLIUM exports an abstract virtual device and a para-virtual guest driver, which we use to interpose and forward the OpenCL and CUDA APIs to the host. Unlike SVGA, which requires translation layers to ensure that all graphics frameworks APIs can be mapped to the SVGA protocol, Trillium forwards the lowest layer in the GNU/Linux Graphics stack: the pipe-driver, effectively remotng OpenCL/CUDA API calls in the guest to the OpenCL/CUDA library in the host.

Our experience implementing the required TGSI vISA support in the Mesa graphics stack led us to believe that the TGSI layer is unnecessary. Not only does this translation introduce additional complexity in the guest stack, it also hurts performance, as we demonstrate in Section 3.6. The guest OpenCL compiler cannot target the native GPU architecture, and semantic information is lost to the host compiler. Further, while incorporating a TGSI compiler is possible in open frameworks like OpenCL, the task is significantly more daunting for closed frameworks like CUDA. Attempts to translate between TGSI and NVIDIA SASS in the reverse-engineered Nouveau driver understandably results in code that is significantly less performant than that produced by the proprietary stack.

TRILLIUM takes a different approach: TRILLIUM forwards API calls for compiling OpenCL code to the hypervisor. The OpenCL compiler in the host OpenCL framework (optimized for the physical hardware by the hardware vendor) is invoked on the forwarded OpenCL code to lower it directly to the physical device ISA.

Figure 3.2a shows the Trillium design layers in a generic hypervisor stack. The OpenCL API is forwarded from the driver similar to the SVGA model. The OpenCL compute kernel (to be run on the GPU), can be passed through to the host via hypercalls in the driver, without being translated to any vISA, where it will be translated and optimized for the physical GPU in a virtual appliance (Dom 2 in Figure 3.2a).

TRILLIUM does not currently guarantee performance isolation and relies on the hardware scheduler. Performance isolation can easily be implemented via a rate-limiting API scheduler in the hypervisor, such as in GPUvm [109].

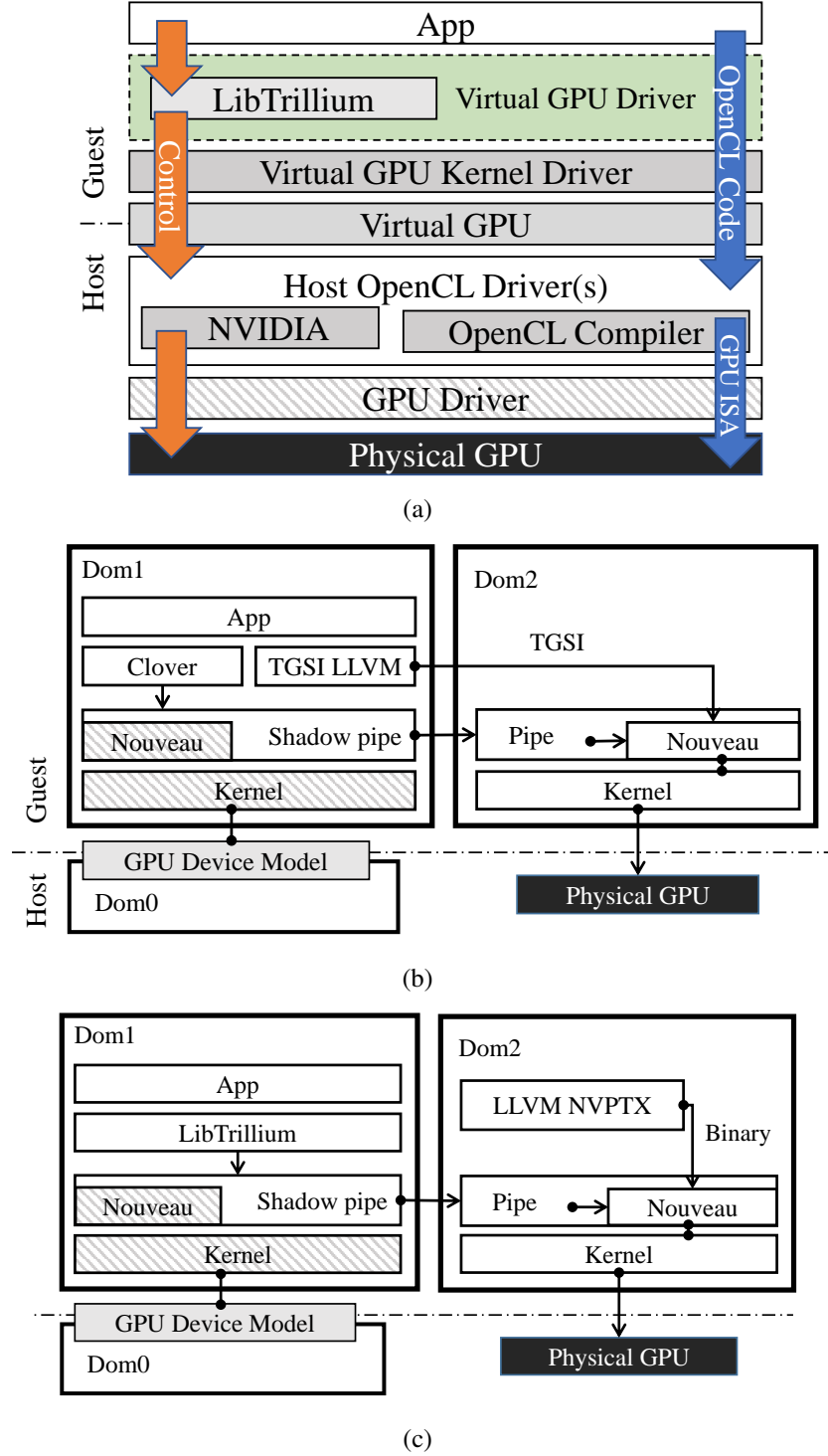


Figure 3.2: XEN-SVGA and TRILLIUM designs. (a) The TRILLIUM stack. (b) XEN-SVGA approximates the SVGA model extended to support GPU Compute. (c) The design of TRILLIUM with shadow pipe.

### 3.4 Implementation

We evaluated the TRILLIUM design against a representative of each traditional virtualization technique: full-virtualization, user-space API-remoting and SVGA. Due to the difficulty of implementing a trap-based virtualization scheme, we chose to evaluate against GPUvm [109], the only existing open-source implementation. GPUvm is tightly coupled with the Xen hypervisor [36]. As a result, all the other prototypes were built on the Xen hypervisor to keep the platform common for fair comparison.

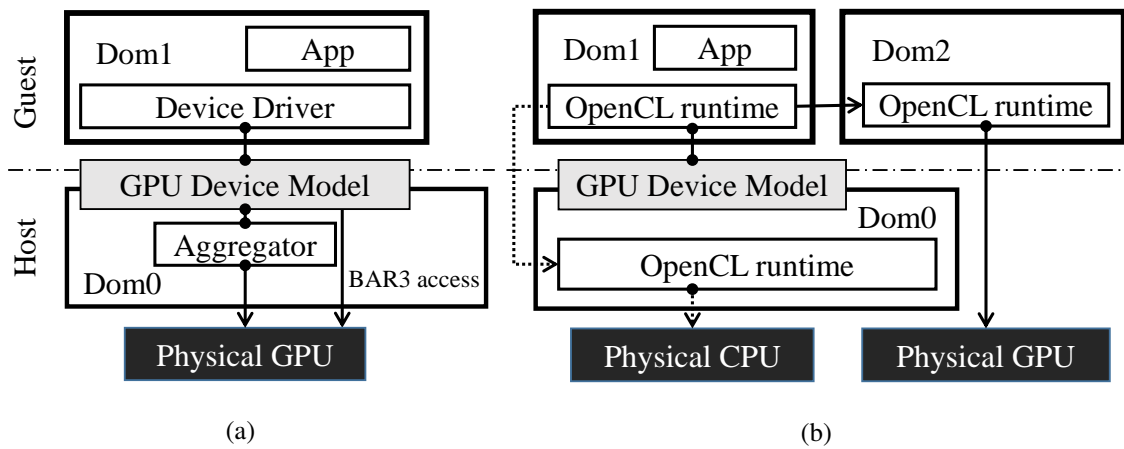


Figure 3.3: Xen-based virtualization designs. (a) Trap-based virtualization: GPUvm. (b) User-space API remoting over RPC—dashed arrows indicate API-REMOTE-CPU, while solid ones indicate API-REMOTE-GPU.

#### 3.4.1 Trillium

We initially implemented TRILLIUM on Xen following the SVGA design, by implementing OpenCL support in a virtual device and extending the Mesa stack with TGSI support (see Section 3.4.1 for details). The generated TGSI is sent to the host via RPC, and then finalized to a binary that can be run on the physical NVIDIA GPU using the open source Nouveau driver. Upon empirically finding that TGSI is a performance bottleneck, we revisited the basic design. We preserve the original prototype, hereafter called XEN-SVGA, as a baseline representative of the original SVGA design: this design is shown in Figure 3.2b. The current TRILLIUM design is shown in Figure 3.2c.

XEN-SVGA and TRILLIUM, implement API-forwarding in a custom pipe-driver in Gallium3D, that we call `shadow-pipe`. We chose to forward the pipe-driver as it presents a narrow interpo-

sition interface in the graphics driver. However, given that each OpenCL API call is decomposed into many different pipe-driver calls, other APIs higher up in the graphics stack may be better suited for interposition. The `shadow-pipe` is in the *application domain*'s graphics stack, and shims the pipe-driver interface as RPC calls to the actual Nouveau pipe-driver in the *privileged domain*.

XEN-SVGA manages user-level contexts, command queues and memory objects; and translates the input OpenCL GPGPU kernel to TGSI in the application domain. TRILLIUM skips the compilation phase in the application domain. The OpenCL kernel is forwarded to the privileged domain via RPC, where it is parsed and compiled by the LLVM NVPTX back-end in parallel. This binary is then loaded onto the GPU when the pipe-driver hits the binary loading phase. TRILLIUM can also emit LLVM IR if an OpenCL compiler is not available in the host.

Our implementation relies on gRPC as a transport mechanism between the guest and the host, as an implementation convenience. As zero-copy transfer [44, 112] and hypercall [94] mechanisms are well-studied, and a production-ready version of TRILLIUM would rely on these mechanisms, we measure and remove transport overhead from our reported measurements in Section 3.6. The overheads stem from remoting calls to the privileged domain over the network, which is especially significant since a single OpenCL API call may be decomposed into many pipe-driver APIs, and from the large amount of kernel input data that must be copied between VMs.

**LLVM TGSI Back-end** The Mesa3D stack implements OpenCL support via a state-tracker called Clover. Clover provides the library for the OpenCL application to link against, while most of the compilation is handled by invoking the OpenCL and C++ front-ends of the LLVM [82] compiler framework. Clover provides much of the front-end infrastructure required to support GPGPU computing in XEN-SVGA and TRILLIUM.

Historically, lack of a working TGSI back-end in LLVM, despite several attempts at building one in the past 5 years [7, 11], has left OpenCL support for NVIDIA GPUs and SVGA in Mesa3D incomplete. In order to support OpenCL in XEN-SVGA, we implemented an LLVM TGSI back-end. While the TGSI back-end is not yet mature, we have added support for a majority of the 32-bit integer and floating point operations, intrinsics, memory barriers, and control flow. Using this backend we are able to compile and run 10 out of the 12 Rodinia benchmarks [43] used to benchmark GPUvm.

Because the compiler is built using the LLVM framework, it enjoys all of the IR-level optimizations in LLVM.

LLVM IR handles control flow by using conditional and unconditional branches to and from Basic Blocks. A majority of the usual optimizations (constant propagation, loop unrolling, etc) are applied on the IR. On the other hand, TGSi assumes a linear control flow through the program, using higher level constructs such as IF-THEN-ELSE, FOR and WHILE loops. To accommodate this difference in control flow techniques, we leveraged a similar implementation in the AMDGPU back-end which calculates a Strongly-Connected-Components (SCC) graph from the Basic Block-based control flow in the LLVM IR, and then duplicates Basic Blocks as necessary. It is a testament to the maturity and flexibility of LLVM that the infrastructure to produce an SCC, and an example of how to use it to raise the control flow abstraction level were readily available.

### 3.4.2 GPUvm

GPUvm [109] is an open-source trap-based interposition design (a simplified block-diagram representation is shown in Figure 3.3a). The application domain (Dom 1) is presented with a GPU Device Model, which is emulated in the privileged domain (Dom 0). The emulation layer in Dom 0 interposes, validates, and fulfills all attempts to access the GPU. GPUvm has not been maintained: The last release, in 2012, is based on Xen 4.2.0 and runs on Fedora 16 [130]. In order to compare all prototypes on the same modern platform, we ported GPUvm to Ubuntu 16.04 with Xen 4.8.2.

### 3.4.3 User-space API remoting over RPC

In order to faithfully mimic user-level API-remoting-over-RPC systems [53, 75, 38], OpenCL API calls are trapped by a user-space shim library and forwarded via RPC from one appliance VM, which is the OpenCL “client”, to another appliance VM, which acts as the OpenCL “server”. Figure 3.3b shows the setup of the two API-remoting schemes we considered: API-REMOTE-GPU and API-REMOTE-CPU. The black arrows indicate the workflow of API-REMOTE-GPU, where the OpenCL server runs the OpenCL commands on a physical GPU using the NVIDIA OpenCL framework. The grey arrows show the API-REMOTE-CPU setup, where the OpenCL commands are executed on a multi-core CPU (Intel CPU Xeon E5-2643) using the Intel OpenCL SRB 5.0 framework. RPC is implemented using gRPC 1.6 (based on Google ProtocolBuffers 3.4.0) and inter-service

communications are implemented over XML-RPC 1.39. Lower-overhead data-movement techniques, such as zero-copy, can be applied when both the client and the server are on a local machine.

### 3.4.4 Optimizations

TRILLIUM interposes at the pipe-driver API yielding fine-grained interposition, and therefore finer-grained multiplexing of the GPGPU. However, interposing at this layer also results in significant transport overhead. Many pipe-driver functions are responsible for context management and information retrieval—operations that do not result in interaction with the GPU. We reduce communication overhead by batching these types of API-calls, taking care to fall back to synchronous API-forwarding when any pipe-driver API calls that interact with the physical GPU are invoked.

We optimize the API-REMOTE-GPU and API-REMOTE-CPU systems by preinitializing the device and preallocating contexts and command queues on the privileged domain. These contexts are assigned to applications as they execute context creation APIs and are reclaimed asynchronously.

## 3.5 Methodology

All experiments were run on a Dell Precision 3620 workstation with NVIDIA Quadro 6000 GPU and Intel Xeon CPU E5-2643 (3.40GHz) CPU. We implemented or ported all prototypes and benchmarks on Ubuntu 16.04 with Xen 4.8.2. VMs were hardware-accelerated via Xen Hardware Virtual Machines (HVM) with 2 virtual CPUs (pinned) and 4 GB memory.

Of the GPU hardware available to us, the NVIDIA Quadro 6000 GPU was the only one that GPUvm, the full-virtualization baseline ran on. GPUvm depends on GDev [72] an open source CUDA runtime (released in 2012) implemented using Nouveau [8] GPU drivers, and the CUDA 4.2 compiler on Linux Kernel 3.6.5. GDev has not been maintained since 2014, and the effort to update it was too onerous. Experiments to control for hardware versions are reported in 3.5.2.

### 3.5.1 Benchmarks

XEN-SVGA depends on the TGSi back-end compiler that we implemented to leverage the Clover OpenCL runtime in Mesa3D. API-REMOTE-GPU and API-REMOTE-CPU leverage the NVIDIA and Intel OpenCL library respectively and support all of the Rodinia benchmarks. GPUvm is built on



top of the GDev CUDA runtime. Care was taken to ensure that the CUDA and OpenCL versions of the benchmarks use the same parameters, datasets, memory barriers, sync points, etc. Experiments to control for the programming framework are reported in 3.5.2.

Table 3.1: EVALUATION BENCHMARKS IN THREE CATEGORIES<sup>a</sup>

Benchmark	Description	Type
backprop	Back propagation (pattern recognition)	R
gaussian	256x256 matrix Gaussian elimination	D
lud	256x256 matrix LU decomposition	M
nn	$k$ -nearest neighbors classification	D
nw	Needleman-Wunsh (DNA-seq alignment)	M
pathfinder	Search shortest paths through 2-D maps	R

a. Interposition-dominant, interposition-rare, and moderate workloads.

The 10 Rodinia benchmarks that our TGSI compiler could compile were categorized based on frequency of interposition: **I**nterposition-**D**ominant workloads run kernels hundreds or thousands of times requiring frequent interposition to set arguments, etc. **I**nterposition-**R**are workloads run a small number of long-running kernels, requiring very little interposition. **M**oderate-interposition workloads lie somewhere in between the other two. Two benchmarks were selected from each category to be used in the evaluation (the optimizations described in 3.4.4 take significant manual effort).

### 3.5.2 Control Experiments

Software and platform version dependencies necessitated that our experimental environments vary slightly for the systems under evaluation — different front-end programming languages (CUDA vs. OpenCL), different runtime implementations (GDev CUDA vs. NVIDIA CUDA), or different drivers (Nouveau vs. NVIDIA). Resolving all of these differences would have taken monumental effort, but control experiments showed that these variables had negligible impact on our measurements.

**OpenCL vs. CUDA** GPUvm relies on the GDev implementation of the CUDA framework, while all the other designs rely on OpenCL. To assess the impact of different front-end languages on performance, we measured execution times for all benchmarks in both CUDA and OpenCL (Rodinia includes both implementations) holding all other variables constant, and found that the front-end language has near negligible impact, and the harmonic mean of differences in kernel execution time across all benchmarks is less than 1%; the worst (maximal) case is 15%. We also found negligible

difference in performance between kernels compiled using CUDA 8.0 and the CUDA 4.2 required by GDev.

**Hardware Generations.** The performance improvements over the span of generations between the Quadro 6000 and modern cards is substantial. To estimate the effect of this variable we ran all benchmarks on both Quadro 6000 and a more recent GPU, Quadro P6000. While overall execution times are improved substantially, and the ratio of time spent on the host to time spent on the GPU changes as a result, the relative speedups are uniform across all benchmarks. This suggests that the trends that we observe on the Quadro 6000 still hold on newer hardware. We re-iterate that software dependencies of the GPUvm baseline prevent us from using more recent hardware. Our evaluation is performed on the newest (several generations older) GPU hardware that all our systems can run on.

## 3.6 Evaluation

We are interested in understanding the impact of a vISA on end-to-end performance, the effect of interposition frequency on performance, and the effectiveness of our proposed design, TRILLIUM.

### 3.6.1 The impact of vISA choice

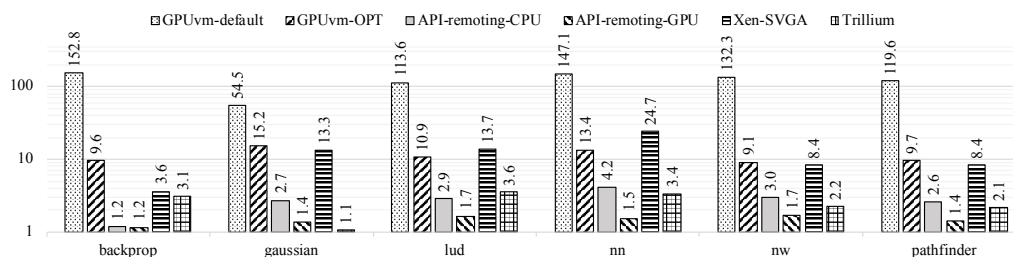


Figure 3.4: End-to-end execution times of benchmarks on virtualization prototypes, relative to end-to-end execution time on the NVIDIA CUDA runtime in a native setting. The gRPC transport overhead is removed from the reported measurements, which is up to 10% of the total execution time for API remoting, and 40% for TRILLIUM.

Deferring the compilation of front-end code to the host not only eliminates redundant translations, and the need to have a compiler in the guest driver, but also ensures that the compiler has a high-fidelity view of the physical hardware. Typically, the execution/compilation framework is extremely tightly coupled with the vISA used, making the choice of vISA even more tenuous as it leads to the

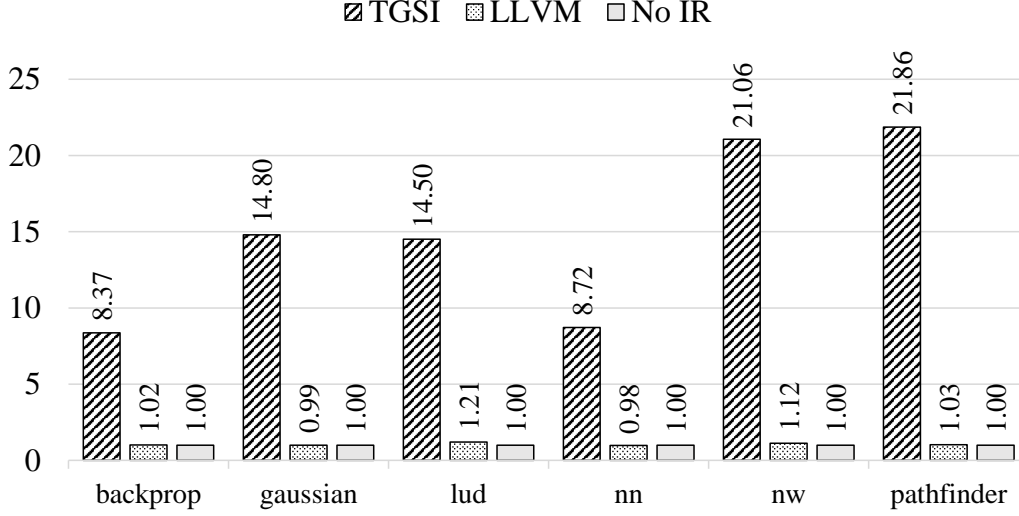


Figure 3.5: Kernel execution slowdown due to virtual ISAs. TGSI: the LLVM TGSI back-end compiler used in XEN-SVGA. LLVM: LLVM NVIDIA PTX (NVPTX) back-end used in TRILLIUM. No IR: native NVIDIA compiler.

second order effect of having to rely on a particular implementation of the compute framework (e.g., Mesa3D OpenCL vs NVIDIA OpenCL).

To understand the impact of the virtual ISA on the quality of the generated GPU code we measured GPU execution time for NVIDIA SASS kernels generated in 3 ways: a) using the Mesa3d OpenCL stack (OpenCL→ TGSI→ SASS), b) using the LLVM OpenCL stack (OpenCL→ LLVM IR→ SASS), and c) using the native NVIDIA OpenCL compiler (OpenCL→ SASS). These measurements are reported in Figure 3.5 relative to kernel execution time in a native setting.

Code generated from TGSI IR is dramatically slower in all cases than code generated by the NVIDIA OpenCL framework. We observe slowdowns of up to  $22\times$ , with a harmonic mean of  $13\times$  across the 6 benchmarks that were optimized for evaluation. While we predicted the basic trend these experiments show, we were surprised by the magnitude of the difference. We found quality of the kernel generated by the LLVM NVPTX compiler to be comparable to native, at least in terms of execution time. This is unsurprising given recent efforts [124] to optimize the LLVM tool-chain for NVIDIA GPUs.

The TRILLIUM design uses LLVM IR as the common virtual ISA for GPGPU applications, where necessary: OpenCL code is compiled to PTX using the LLVM NVPTX back-end in the guest, and then finalized and executed in the host using the NVIDIA CUDA framework.

### 3.6.2 End-to-End

We compare TRILLIUM against full-virtual (GPUVM-DEFAULT and GPUVM-OPT), API remoting (API-REMOTE-GPU and API-REMOTE-CPU) and para-virtual (XEN-SVGA) systems. XEN-SVGA approximates an SVGA-like design in Xen (Mesa3D with TGSI). TRILLIUM bypasses translation from OpenCL to TGSI. To characterize the behavior a full-virtualization design, we measure GPUvm [109] in its default configuration (worst case) shown as GPUVM-DEFAULT, and in its fully optimized configuration, labeled GPUVM-OPT.

Figure 3.4 shows the end-to-end execution time (relative to native GPU execution) for the six chosen benchmarks for all the systems evaluated. As expected, traditional API remoting designs incur the lowest overhead, which is achieved by giving up hypervisor interposition. TRILLIUM fares well with best case performance of just  $1.1\times$  over native, and within  $3.6\times$  at worst. XEN-SVGA is sensitive to the performance lost in GPU kernel code resulting from redundant compilation through TGSI (which adds significant overheads as previously shown in Figure 3.5). GPUVM-OPT exhibits about  $9.1\times$  slowdown for applications with short-lived kernels (e.g. Needleman-Wunsh algorithm); the overhead can be as high as  $15.2\times$  when the workload has long-running kernels (e.g. Gaussian Elimination).

We find that remoting calls intended to a CPU is uniformly more performant than full-virtualization of the GPU, and sometimes performs just as well as (backprop) or better than remoting to the GPU ( $1.6\times$  faster for the bfs benchmark. The performance gain from accelerating the bfs kernel on the GPU is severely dwarfed by the cost of initialization on the GPU). GPGPU compute is only economical when it provides acceleration over the CPU; if overheads make the CPU competitive, the profitability threshold has been crossed. Further, the competitiveness of API-REMOTE-CPU suggests opportunity: systems could back a virtual GPU with CPU if they can detect when it is profitable to do so.

## 3.7 Conclusion

TRILLIUM represents a local optima in the GPGPU virtualization space—by decoupling device virtualization from GPU ISA virtualization, it maintains the virtualization benefits of a para-virtual system, while exhibiting the performance of a user-space remoting system.

## CHAPTER 4: HYPERVISOR-MEDIATED API-REMOTING

Practical virtualization must support sharing and isolation under flexible policy with minimal overhead. The structure of current accelerator stacks makes this extremely difficult to achieve. Accelerator stacks are *silos* (Figure 4.1) comprising proprietary layers communicating through memory mapped interfaces. This opaque organization makes it *impossible* to interpose intermediate layers cleanly to form a virtualization boundary. Practically interposable alternatives leave designers with a Hobson’s choice between critical virtualization properties such as interposition and compatibility.

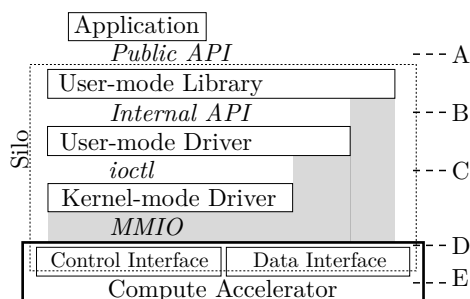


Figure 4.1: An accelerator silo. The public API and the interfaces with striped backgrounds are interposition candidates. All interfaces with backgrounds are proprietary and subject to change.

We present AVA, a system that addresses the fundamental limitations of existing accelerator virtualization techniques. AVA combines API-agnostic para-virtual I/O stack components with a Domain-Specific Language (DSL) and toolchain to automate construction and deployment of guest libraries and API servers. AVA uses an abstract para-virtual device to serve as a transport endpoint for forwarding the public APIs of vendor-provided frameworks (e.g. CUDA or TensorFlow). Unlike currently popular user-space API remoting solutions [3, 67, 119, 53, 96], AVA preserves hypervisor-level resource management and strong isolation using a novel technique called *Hypervisor Interposed Remote Acceleration*. AVA forwards API calls over hypervisor-managed communication channels, inserting automatically-generated resource management components between traditional front- and back-ends to enforce policies described in the DSL specification. Critically, *automation*

from AVA enables hypervisors to keep up with fast accelerator evolution: automatic generation of components minimizes engineering effort.

AVA supports a broad range of currently-shipping compute accelerators: We virtualized ten accelerators including NVIDIA and AMD GPUs, Google TPUs, and Intel QuickAssist. Virtualizing an API framework using AVA requires modest developer effort: a single developer virtualized OpenCL in a handful of days, a stark contrast to the person-years of developer effort for VMware’s SVGA II or Bitfusion’s FlexDirect [3]. Experiments show that AVA provides near-native performance (e.g., 2.4% slowdown for TensorFlow and 5.6% for CUDA), enforces isolation and fair sharing across guests, and supports live migration.

The proposed chapter will make the following arguments:

- The chapter demonstrates feasibility of automatically constructed virtual accelerator support, showing that a single technique can deal with many architectures, APIs, versions, and policies.
- We introduce Hypervisor Interposed Remote Acceleration (HIRA) to enable hypervisor-enforced isolation and sharing policies unachievable with current SR-IOV and API remoting systems.
- We utilize a novel DSL, LAPIS, for describing API functions, resources, and policies to enable automatic construction of virtual stacks from native header files.
- Our evaluation shows low developer effort, strong isolation, and good performance.

This chapter will draw from joint work with Hangchen Yu, Arthur Peters and Christopher J Rossbach. Part of this work was published as a workshop paper [128], and a longer paper that will appear at ASPLOS’20. As with any big system building effort, it was a team effort. This material will be a part of Hangchen’s and Arthur’s dissertations as well, and we have agreed to split the intellectual contributions from the project thus: While I will focus on the core virtualization technique, i.e., hypervisor-mediated API-remoting in my dissertation, Hangchen’s dissertation will deal with the challenges of designing API-remoting systems for both virtualization as well as use from within an OS kernel; Arthur’s dissertation will focus on the specification and design of LAPIS, a DSL for specifying and generating remoting code for arbitrary C APIs.

## CHAPTER 5: IEMTS — A NEW ACCELERATOR VIRTUALIZATION TAXONOMY

Traditionally, virtualization designs have been taxonomized according to the core techniques employed (e.g. emulation, full- or para-virtualization, API remoting, etc.), and evaluated in a property trade-off space comprising performance, compatibility, interposition, and isolation. *Isolation* ensures that mutually distrustful guests cannot access each other’s data or harm each other’s performance. *Compatibility*, characterizes how well a design preserves the freedom of hardware and software components to evolve independently: e.g. changes in the hypervisor should not force changes to guest software. Virtualization provides an indirection layer between logical and physical resources by *interposing* a well-defined interface. The quality of interposition determines the nature of benefits (e.g. extent of consolidation) afforded by a virtualized system [121].

Virtualization techniques are well explored, yielding conventional wisdom about their fundamental trade-offs. For example, *full virtualization* interposes the software-hardware interface to provide a virtual view of the underlying hardware. This enables guests to run unmodified OS and application binaries, yielding high compatibility. However, hardware interfaces for GPUs rely heavily on MMIO and communication through memory, which necessitates page-fault-based interposition [114, 12, 77, 87] techniques that cripple performance. *Para-virtual* designs export an abstract device to the guest, but require hypervisor-specific drivers and runtime libraries in the guest, trading compatibility for improved performance. *API remoting (or forwarding)* [61, 51, 57, 101] aggregates high-level API calls issued in VMs, running them on the host or in a dedicated appliance VM. This technique can provide near native performance because API calls are infrequent, but has poor compatibility because it requires changes in guest applications or libraries.

We argue that the current *de facto* taxonomy and property trade-off space are illustrative but not informative for GPUs: there is a large body of research that has had little influence on practice. First, Classifying virtualization designs as API-remoting vs. full vs. para-virtual captures important concepts, and emergent properties compactly, but doesn’t explain their correlation to properties like performance. Second, virtualization properties such as compatibility, isolation, and interposition

	GPUvm	VMware SVGA		rCUDA
		Control Interface	GPU ABI	
Interposed Interface	MMIO/BAR	DirectX APIs	Device ISA	Userspace API
Interposition Source	Trap handler	Guest driver/libs	Guest Driver	Guest Library
Interposition Destination	Host driver	Host framework	Host Driver	Host/Server Daemon
Interposition Mechanism	Trap	Guest library	Compilation to vISA	Guest Library Shim
Transport	Fault	Hypervisor FIFOs	Hypervisor FIFOs	RPC
Synthesis	Emulation	Call host API	Binary translation	Call Server API

Table 5.1: Comparing virtualization designs using the IEMTS framework.

have highly context-dependent meaning and their relative value to system designers can be hard to quantify. Consider compatibility: there are many dimensions to compatibility (library, hardware, OS, etc.), and each of those are commonly achieved by separate technical, and non-technical means (e.g., TGSi is the common vISA for both the VMware and GNU/Linux graphics stacks; this is *not a lucky coincidence*).

We argue that practical design goals, such as providing a virtualization layer with specific characteristics, get obscured when these properties are considered as a set of constraints that must be preserved, without first refining for context. Further, production systems, such as VMware SVGA [51], compose multiple virtualization techniques in order to leverage the best properties of each technique, especially in the presence of multiple interfaces.

To enable a cleaner separation of concerns, we draw on the observation that *all* virtualization relies on encapsulation and interposition, and note that a design can be clearly understood by identifying:

- the **I**nterface that is interposed,
- the **E**nd-points (source and destination) the interposed event is transported between,
- the **M**echanism used to interpose,
- the **T**ransport mechanism used to communicate between endpoints,
- the mechanisms used to **S**ynthesize or implement the desired functionality at the destination. We call this the **IEMTS** framework.

Table 5.1 presents analysis of three prior GPU virtualization systems under the IEMTS framework. A quick glance at the table tells us that GPUvm’s [108] performance woes arise from its realiance on trapping and emulating the guest’s MMIO accesses. VMware’s SVGA has two entries in the table because there are two interfaces being virtualized: the control interface (the Direc-



tX/OpenGL API) and the shader ISA. Explicitly separating the two interfaces helped us realize that ISA-virtualization is not necessary for accelerator virtualization in the Trillium project. Further, it became obvious to us that the control interface virtualization in SVGA and the API-remoting in rCUDA look almost the same under IEMTS with the exception of the transport. This observation led us to design AvA's hypervisor-mediated API-remoting scheme, and shift our attention to solving the compatibility effort via automation.

## CHAPTER 6: PROPOSED WORK — VTASK

The previous chapters in the proposed dissertation showed that hypervisor-mediation API-remoting is an effective mechanism for sharing API-controlled compute devices among mutually distrustful tenants, e.g., in a cloud computing environment. Virtualization vendors, such as VMware, have begun adopting API-remoting based solutions for accelerator virtualization [19].

API-remoting works by interposing on API calls invoked by the application in the guest OS, and executing them in a surrogate, the API-server, in the host. Typically, API-servers are associated with a single API framework (for modularity and failure isolation between APIs/accelerators, and in order to be able to use remote resources) and each API-server is a surrogate for a single guest (to preserve isolation between guests). Applications that use multiple accelerator API frameworks will, therefore, be associated with multiple API-servers, one per framework.

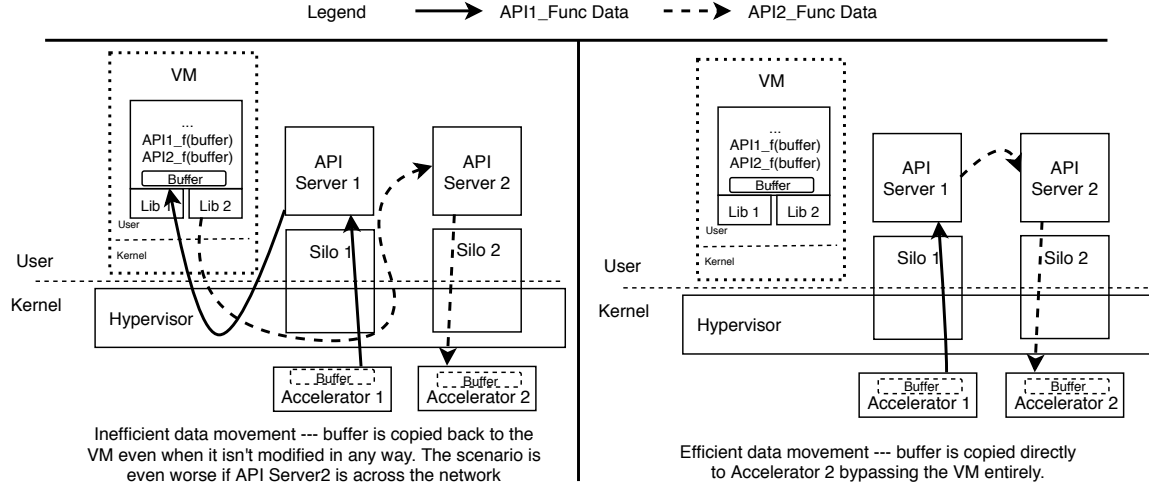


Figure 6.1: Data processed by two API stacks must pass through the guest application

Under a typical API-remoting system, applications that pipeline disparate accelerator frameworks are burdened with redundant data movement. All inter-accelerator data movement must take place in the guest application as that is where the accelerators are in the same logical address space. Figure 6.1 illustrates this scenario: when an *API-1* function is invoked, associated data is copied from the *guest*

*application* to *API-server-1*, and then to *Device-1*'s memory. Once the function finishes executing on *Device-1*, the result is copied back to the *guest application*. When a function from *API-2* is invoked, the same data (i.e., the output of the *API-1 function*) is copied from the *guest application* to *API-server-2* and then to *Device-2* to be processed.

In order to eliminate redundant data movement when an application uses multiple accelerators via API-remoting, the hypervisor must track the data passed to these API calls. The hypervisor must keep track of where the data flowed from and to, the validity of different copies of the data (e.g., if the data is modified on the accelerator, but hasn't been copied back to the guest application), and eliminate redundant data movement. As an example, if a guest application were to invoke the `cudaMemcpyDtoH()` function to copy data back from an Nvidia GPU, and then invoke the Intel QAT compression function `cpaDCCompressData2()` on the same data without modifying it in any way, the hypervisor should be able to detect this and elide the copying of data to and from the guest application. Further optimization may also be possible: peer-to-peer data copy between the devices if they are on the same machine, or by directly copying the data from the first API-server on one remote machine to the second API-server on another remote machine.

We propose to build vTask, an application-transparent data orchestration system that optimizes data movement among accelerators virtualized via API-remoting. vTask will leverage information from API annotations [128] to track data buffers across the guest application, the API-servers servicing API calls made by the guest, and the accelerator hardware. vTask will optimize data movement across these components while ensuring that a coherent view of the data buffer is presented to anyone attempting to read the data. Ideally, vTask will require no changes to the guest application or extra annotations of any kind from the application programmer. We hypothesize that annotations provided to virtualize the API (by the device or virtualization vendor) will be sufficient to infer the semantics of the data buffers managed.

We will prototype vTask in AvA, a state-of-the-art para-virtual API-remoting system for KVM. vTask will rely on device-side buffer allocation and deallocation API calls, and special annotations provided by LAPIS, AvA's API description language, to determine buffer lifetime. Further, vTask will implement a simple MESI-style coherence protocol to track spatial validity of data (i.e., to track where the latest data is present). vTask will leverage optimizations such as shared memory,

Unified Virtual Memory, and PCIe Peer-to-Peer (P2P) data transfer where available, but does not make assumptions about their universal availability.

With vTask, AvA will be able to handle data movement between both local and remote devices. When API-remoting to a remote system, the devices used by the guest application may be present on separate machines. We hypothesize that vTask will be able to eliminate costly data transfers over the network by adhering to the principle of lazy loading wherever possible, i.e., data is not moved until a demand fault occurs.

## CHAPTER 7: RELATED WORK

Virtualization has such a long and storied history that Attempting to capture the entire story is an exercise in futility. The introduction 1 captures the history of CPU virtualization in broad strokes. This section then focuses on a major theme of the proposed dissertation: accelerator virtualization.

Accelerating specific computation is not a new idea—support for specialized computation is extremely commonplace in CPUs (e.g., Floating Point Units (FPU), Vector Processing Units). These specialized compute units are typically exposed to the programmer as extensions to the Instruction Set Architecture (ISA). Virtualizing these specialized compute units, therefore, is no different from virtualizing the rest of the CPU and ISA virtualization is well explored [27, 46, 92, 40, 42, 41].

Processors specialized for complex computational tasks, such as graphical rendering, largely evolved as discrete devices separate from the CPU (although some CPUs do integrate GPUs). These devices are not typically integrated into the CPU ISA; instead, they appear to system software as I/O devices with memory-mapped command-queues and I/O registers. I/O virtualization is well understood [121, 33, 78, 103, 132, 26], but these techniques aren't enough to virtualize programmable accelerators. Although programmable accelerators look like I/O devices, they are also general computing platforms, i.e., they load binaries, have their own memory, and are typically exposed to the application programmer via an API.

### 7.1 GPU Virtualization

GPU virtualization has received a lot of attention since the late 2000s. This section presents an overview of all prior work. Table 7.1 presents a comprehensive overview of prior accelerator virtualization techniques in terms of traditional virtualization properties. The **lib unmod** and **OS unmod** columns indicate ability to support unmodified guest libraries and OS/driver. The **lib-compat** and **hw-compat** indicate the ability (compatibility) to support a GPU device abstraction that is independent of *framework* or *hardware* actually present on the host. **sharing**, **isolation** and **sched. policy** indicate cross-domain sharing, isolation and some attempt to support fairness or

Technique	System	lib unmod	OS unmod	lib-compat	hw-compat	sharing	isolation	migration	sched. policy	graphics	GPGPU	I/D	benchmark	slowdown	native speedup	virtual speedup
Full-virtual	GPUvm [109]	✓		✓		✓	✓		XC, BAND		✓	<i>D</i>	Rodinia	141×	11.4×	0.08×
	gVirt [113]	✓		✓		✓	✓	✓	QoS	✓		<i>I</i>	2D [15], 3D [4]	1.6×	N/A	N/A
PCIe Pass-thru	AWS GPU [31]	✓	✓							✓	✓	<i>D</i>	Any	1×		
API remoting	GViM [61]				✓	✓	✓		RR, XC		✓	<i>D</i>	CUDA 1.1 SDK	1.16×	22×	19×
	gVirtuS [56]				✓	✓	✓		RR		✓	<i>D</i>	CUDA 2.3 MM	3.1×	11.1×	3.6×
	vCUDA [101]		✓		✓			✓	HW		✓	<i>D</i>	CUDA 4.0 SDK	1.91×	6×	3.1×
	vmCUDA [119]		✓		✓	✓			HW		✓	<i>D</i>	CUDA 5.0 SDK	1.04×	33×	31.7×
Distributed API remoting	rCUDA [53, 95]		✓		✓	✓	✓		RR		✓	<i>D</i>	CUDA 3.1 SDK	1.83×	49.8×	27.2×
	GridCuda [84]		✓		✓	✓	✓		FIFO		✓	<i>D</i>	CUDA MM, SOR	1.23×		
	SnuCL [75]		✓		✓	✓	✓				✓	<i>D</i>	SNU NPB [100]			
	VCL [35]		✓		✓	✓	✓				✓	<i>D</i>	Stencil2D [47]			
Para-virtual	GPUvm [109]					✓	✓		XC, BAND		✓	<i>D</i>	Rodinia	5.9×	11.4×	1.9×
	HSA-KVM [66]	✓				✓	✓		HW		✓	<i>I</i>	AMD OCL SDK	1.1×		
	LoGV [59]	✓		✓		✓	✓	✓	RR		✓	<i>D</i>	Rodinia	1.01×	11.4×	11.3×
	SVGA2 [51]	✓				✓	✓	✓		✓		<i>D</i>	2D, gaming	3.9×		
	Paradice [33]	✓		✓		✓	✓		HW, QoS	✓	✓	<i>D</i>	OpenGL, OpenCL	1.1×		
	VGVM [115]				✓	✓	✓		HW		✓	<i>D</i>	CUDA 5.0 SDK	1.02×	33×	32.3×

Table 7.1: Existing GPU virtualization proposals, grouped by approach. Previously published in the Trillium paper [30].

performance isolation (policies such as RR Round-Robin, XC XenoCredit, HW hardware-managed, etc.). The **migration** shows support for VM migration. **I/D** indicates it supports either integrated or discrete GPU. The table also includes performance entries for each system including the geometric-mean slowdown (execution time relative to native execution) across all reported benchmarks. We additionally include the benchmarks used, and where possible, a report (or estimate) of the geometric-mean speedup one should *expect* for using GPUs over CPUs using hardware similar to that used in this paper. The final column is the expected geometric-mean speedup for the given benchmarks running in the virtual GPGPU system over running on native CPUs. Values in this column were computed by dividing the expected speedup from using a GPU by the slowdown introduced by virtualization. Entries where overheads eclipse GPU-based performance gains are marked in **red**. Performance profitable entries are **blue**. Greyed out cells indicate the metric is meaningless for that design. Light grey cells indicate that the data was not available.

## Full Virtualization

**GPUvm** [108] virtualizes CUDA on Kepler and Fermi (NVIDIA) GPUs for Xen. GPUvm presents a GPU device model to each VM. Attempts to access the GPU from all VMs are routed through a GPU Aggregator. The aggregator maintains shadow page tables, shadow channels, implements a “fair share scheduler”, and modifies requests to enforce isolation. GPUvm interposes on communication between guest device driver and the GPU device model, by trapping and forwarding MMIO writes.

The authors also explore a number of optimizations: lazy shadowing, bar remap, para-virtualization, and multi-call batching. Despite these optimizations, GPUvm remains non-viable due to its high overhead—the most optimal configuration of GPUvm induces a  $6 \times$  slowdown on average.

**gVirt** [114] is a *graphics*-only virtualization technique for Intel GPUs. The GPU is multiplexed among multiple VMs via pass-through for access to performance critical resources (command buffer and frame buffer), and trap-emulate for resources generally accessed off the critical path (PTEs, I/O Registers). Initialization and power management are done by the native driver in DOM0. Memory is multiplexed with a combination of partitioning and “ballooning”. gVirtus is geared toward graphics and can’t be easily extended to support GPGPU computing.

### API Remoting

**GVim** [61] supports a straightforward split-driver API remoting approach to virtualization of CUDA in Xen. CUDA API calls made by applications in the Guest VM are interposed through a front-end driver (using Xen event channels) and forwarded to a back-end driver in DOM0, which exercises the CUDA driver and runtime. While GVim’s split-driver design is very similar to AvA’s HIRA, AvA presents an accelerator-agnostic framework that can be used to implement hypervisor-mediated API-remoting for arbitrary devices, can enforce flexible policies via callbacks and tackles the compatibility issues inherent in API-remoting.

**vCUDA** [102] is another CUDA API-remoting system. CUDA API calls are redirected through an interposer library (“vCUDA”) to a stub in the host OS, which interacts with the device using pass through. RPC turns out to be the primary performance term. The authors explores RPC batching as an optimization. The system has no support for interposition.

**vmCUDA** [119] observes that while pass-through is performant, it precludes sharing, and VM migration. vmCUDA employs a split-driver model with a front-end driver in the guest that provides an interposition point, and a back-end driver in the control domain which interacts with the CUDA runtime and driver. CUDA applications in the guest are linked against an interposer library which forward calls and data to the appliance VM. As with vCUDA, the system guarantees no isolation among VMs.

*gVirtuS* [56] is an API remoting framework that claims to provide transparent support for CUDA, OpenCL, and OpenGL on Xen, KVM, and VMware ESXi, using a split-driver design to provide a formal abstraction layer for GPUs that is independent of VMM.

*rCUDA* [53, 95], *GridCuda* [84], *SnuCL* [75] and *VCL* [35] are all user-mode middle-ware systems for multiplexing GPUs and CUDA/OpenCL across a cluster. A client library encapsulates access to a (potentially) remote GPU. While the basic design is isomorphic to the API remoting design, virtual machines need not be present.

### **Para-virtualization**

*LoGV* [59] describes an approach to GPGPU virtualization that uses GPU protection hardware in the hypervisor to enforce cross-VM isolation. This strategy has two important consequences. First, cross-VM isolation is easy to enforce, and the overheads of virtualization induced by the hypervisor component is minimal. Second, guests are left with no hardware mechanism to enforce cross-process isolation, which pushes the responsibility on the guest driver, which may be forced to use high overhead techniques to provide such guarantees. We classify LoGV as a para-virtualization technique because it ultimately forces change on the guest OS in the form of changes to support process isolation. The virtualization overheads reported in the paper are exceptionally rosy (indeed, the virtualized solution is faster than native in 3 of 4 reported benchmarks). However, the evaluation prototype makes no effort to enforce isolation in guests, which ultimately hides a significant cost and makes the reported numbers meaningless.

*HSA-KVM* [66] is a para-virtual system for Heterogeneous System Architecture (HSA) compliant systems. HSA has the CPU and GPU integrated into the same physical address space. The GPU exposes multiple Architected Queues (Command Queues) that can be allocated to different guests. HSA-KVM comes closest to the flexibility, compatibility and performance of CPU virtualization. However, the design espoused requires high levels of co-operation from the accelerator hardware, and as alluded earlier, this level of hardware support is still missing in most accelerators.

*SVGA2* [51] is VMware’s *graphics-only* GPU virtualization scheme. SVGA2 employs a para-virtual split-driver design with a custom GPU ISA for shader programs (TGSI). SVGA2 supports multiple front-end libraries (OpenGL, DirectX, etc.) via a common driver that is shipped with most mainstream operating systems. SVGA2 uses DirectX as an internal transport mechanism, effectively



realizing API-remoting through the split-driver design. Trillium attempts to extend the SVGA2 model to GPGPU computing. We find that the SVGA2 model is a poor fit for general purpose accelerators due to drastically different constraints. As an example, consider performance. SVGA2 has a much lower target to hit: frames per second (fps), while GPGPU virtualization must preserve the raw speedup over computing with a CPU. This makes the multiple translations needed to implement the many-to-many multiplexing viable for graphics rendering but not for general purpose computation.

## 7.2 Language-level Virtualization

*Dandelion* [98] abstracts accelerators at the language runtime by compiling sequential .Net code to the accelerator (GPU or FPGA). vTask draws inspiration from the buffer management in Dandelion and PTask [97], the underlying accelerator abstraction layer. *HPVM* [105] explores the design of a virtual ISA (vISA) for abstracting heterogeneous compute devices. The HPVM vISA serves as a portable compilation target for managed language run-times built on top of the LLVM compiler infrastructure. HPVM can serve as a replacement for LLVM in Trillium. HPVM, however, doesn't absolve the language run-time implementation from interacting with the accelerator silo. *TornadoVM* is an example of such a managed language runtime implementation (Graal VM).

## BIBLIOGRAPHY

- [1] Amazon EC2 instance types. <https://aws.amazon.com/ec2/instance-types/>. Accessed: 2017-04.
- [2] AMD Multi-user GPU. <http://www.amd.com/Documents/Multiuser-GPU-White-Paper.pdf>. Accessed: 2018-07.
- [3] Bitfusion: The elastic AI infrastructure for multi-cloud. <https://bitfusion.io>. Accessed: 2019-04.
- [4] Cairo-perf-trace. <http://www.cairographics.org>. Jan. 2018.
- [5] Compute shader. [https://www.khronos.org/opengl/wiki/Compute\\_Shader](https://www.khronos.org/opengl/wiki/Compute_Shader). Jan. 2017.
- [6] Deploying rCUDA in cloud computing environments. [http://www.rcuda.net/pub/white\\_paper\\_cloud\\_v2.pdf](http://www.rcuda.net/pub/white_paper_cloud_v2.pdf). Published: 2016-12.
- [7] Francisco Jerez's TGSi back-end. <https://github.com/curro/llvm>. Jan. 2018.
- [8] Freedesktop nouveau open-source driver. <http://nouveau.freedesktop.org>. Accessed: 2017-04.
- [9] Google Cloud TPU. <https://cloud.google.com/tpu>. Accessed: 2019-04.
- [10] GPU Applications Catalog. <https://www.nvidia.com/en-us/data-center/gpu-accelerated-applications/catalog/>. Jan. 2018.
- [11] Hans de Goede's TGSi back-end. <https://cgit.freedesktop.org/~jwrdegoede/llvm>. Jan. 2018.
- [12] KVMGT - the implementation of intel gvt-g(full gpu virtualization) for KVM. <https://lwn.net/Articles/624516/>. 2014.
- [13] Multi-Process Service. [https://docs.nvidia.com/deploy/pdf/CUDA\\_Multi\\_Process\\_Service\\_Overview.pdf](https://docs.nvidia.com/deploy/pdf/CUDA_Multi_Process_Service_Overview.pdf). Accessed: 2019-08.
- [14] NVIDIA CUDA 4.0. <http://developer.nvidia.com/cuda-toolkit-40>. 2011.
- [15] Phoronix test suites. <http://phoronix-test-suite.com>. Jan. 2018.
- [16] Project Fiddle: Fast and Efficient Infrastructure for Distributed Deep Learning. <https://www.microsoft.com/en-us/research/project/fiddle>. Accessed: 2019-04.
- [17] Underutilizing Cloud Computing Resources. <https://www.gigenet.com/blog/underutilizing-cloud-computing-resources/>. Published: 2017-11.
- [18] VMware SVGA3D Guest Driver. <https://www.mesa3d.org/vmware-guest.html>. Accessed: 2019-08.

- [19] Vmware to acquire bitfusion. <https://blogs.vmware.com/vsphere/2019/07/vmware-to-acquire-bitfusion.html>. Accessed: 2019-10.
- [20] VMware Workstation Version History. [https://en.wikipedia.org/wiki/VMware\\_Workstation#Version\\_history](https://en.wikipedia.org/wiki/VMware_Workstation#Version_history). Accessed: 2019-08.
- [21] Why frame rate and resolution matter. <https://www.polygon.com/2014/6/5/5761780/frame-rate-resolution-graphics-primer-ps4-xbox-one>. 2011.
- [22] Gallium3d technical overview, 2017.
- [23] The mesa 3d graphics library, 2017.
- [24] TOP500 Supercomputer Sites. <https://www.top500.org/lists/2018/11/>, 2019.
- [25] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- [26] Darren Abramson, Jeff Jackson, Sridhar Muthrasanallur, Gil Neiger, Greg Regnier, Rajesh Sankaran, Ioannis Schoinas, Rich Uhlig, Balaji Vembu, and John Wiegert. Intel virtualization technology for directed i/o. *Intel technology journal*, 10(3), 2006.
- [27] R.J. Adair. *A Virtual Machine System for the 360/40*. IBM Cambridge Scientific Center report. International Business Machines Corporation, Cambridge Scientific Center, 1966.
- [28] Neha Agarwal, David Nellans, Mike O’Connor, Stephen W Keckler, and Thomas F Wenisch. Unlocking bandwidth for GPUs in CC-NUMA systems. In *HPCA*, 2015.
- [29] Ole Agesen, Jim Mattson, Radu Rugina, and Jeffrey Sheldon. Software techniques for avoiding hardware virtualization exits. In *Presented as part of the 2012 USENIX Annual Technical Conference (USENIX ATC 12)*, pages 373–385, Boston, MA, 2012. USENIX.
- [30] Amogh Akshintala, Hangchen Yu, Arthur Peters, and Christopher J Rossbach. Trillium: The code is the ir. In *The Second Special Session on Virtualization in High Performance Computing and Simulation (VIRT 2019)*, Dublin, Ireland, 2019.
- [31] Amazon. *Amazon Elastic Compute Cloud*, 2015.
- [32] Inc or Its Affiliates Amazon Web Services. Amazon EC2 P3 Instances. <https://aws.amazon.com/ec2/instance-types/p3/>. Accessed: 2018-2-6.
- [33] Ardalan Amiri Sani, Kevin Boos, Shaopu Qin, and Lin Zhong. I/o paravirtualization at the device file boundary. In *Proceedings of the 19th International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS ’14*, pages 319–332, New York, NY, USA, 2014. ACM.
- [34] Ardalan Amiri Sani, Kevin Boos, Min Hong Yun, and Lin Zhong. Rio: A system solution for sharing i/o between mobile systems. In *Proceedings of the 12th annual international conference on Mobile systems, applications, and services*, pages 259–272. ACM, 2014.

- [35] A. Barak, T. Ben-Nun, E. Levy, and A. Shiloh. A package for opencl based heterogeneous computing on clusters with many gpu devices. In *Cluster Computing Workshops and Posters (CLUSTER WORKSHOPS), 2010 IEEE International Conference on*, pages 1–7, Sept 2010.
- [36] Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, and Andrew Warfield. Xen and the art of virtualization. In *ACM SIGOPS operating systems review*, volume 37, pages 164–177. ACM, 2003.
- [37] Fabrice Bellard. Qemu, a fast and portable dynamic translator. In *USENIX Annual Technical Conference, FREENIX Track*, pages 41–46, 2005.
- [38] BitFusion Inc. Bitfusion FlexDirect Virtualization Technology White Paper. <http://bitfusion.io/wp-content/uploads/2017/11/bitfusion-flexdirect-virtualization.pdf>, 2019. Accessed: 2019-2-28.
- [39] David Blythe. The direct3d 10 system. *ACM Trans. Graph.*, 25(3):724–734, 2006.
- [40] Edouard Bugnion, Scott Devine, Kinshuk Govil, and Mendel Rosenblum. Disco: Running commodity operating systems on scalable multiprocessors. *ACM Transactions on Computer Systems (TOCS)*, 15(4):412–447, 1997.
- [41] Edouard Bugnion, Scott Devine, Mendel Rosenblum, Jeremy Sugerman, and Edward Y Wang. Bringing virtualization to the x86 architecture with the original vmware workstation. *ACM Transactions on Computer Systems (TOCS)*, 30(4):12, 2012.
- [42] Edouard Bugnion, Jason Nieh, and Dan Tsafirir. Hardware and software support for virtualization. *Synthesis Lectures on Computer Architecture*, 12(1):1–206, 2017.
- [43] Shuai Che, Michael Boyer, Jiayuan Meng, David Tarjan, Jeremy W Sheaffer, Sang-Ha Lee, and Kevin Skadron. Rodinia: A benchmark suite for heterogeneous computing. In *Workload Characterization, 2009. IISWC 2009. IEEE International Symposium on*, pages 44–54. Ieee, 2009.
- [44] Hsiao-keng Jerry Chu. Zero-copy tcp in solaris. In *Proceedings of the 1996 annual conference on USENIX Annual Technical Conference*, pages 21–21. Usenix Association, 1996.
- [45] Ronan Collobert, Koray Kavukcuoglu, and Clément Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, number EPFL-CONF-192376, 2011.
- [46] R. J. Creasy. The origin of the vm/370 time-sharing system. *IBM J. Res. Dev.*, 25(5):483–490, September 1981.
- [47] Anthony Danalis, Gabriel Marin, Collin McCurdy, Jeremy S Meredith, Philip C Roth, Kyle Spafford, Vinod Tipparaju, and Jeffrey S Vetter. The scalable heterogeneous computing (shoc) benchmark suite. In *Proceedings of the 3rd Workshop on General-Purpose Computation on Graphics Processing Units*, pages 63–74. ACM, 2010.
- [48] Yaozu Dong, Xiaowei Yang, Jianhui Li, Guangdeng Liao, Kun Tian, and Haibing Guan. High performance network virtualization with SR-IOV. *Journal of Parallel and Distributed Computing*, 72(11):1471–1480, 2012.

- [49] Yaozu Dong, Zhao Yu, and Greg Rose. Sr-iov networking in xen: Architecture, design and implementation. In *Proceedings of the First Conference on I/O Virtualization*, WIOV'08, pages 10–10, Berkeley, CA, USA, 2008. USENIX Association.
- [50] Yaozu Dong, Zhao Yu, and Greg Rose. Sr-iov networking in xen: Architecture, design and implementation. In *Workshop on I/O Virtualization*, 2008.
- [51] Micah Dowty and Jeremy Sugerman. Gpu virtualization on vmware's hosted i/o architecture. *ACM SIGOPS Operating Systems Review*, 43(3):73–82, 2009.
- [52] José Duato, Francisco D. Igual, Rafael Mayo, Antonio J. Peña, Enrique S. Quintana-Ortí, and Federico Silla. An efficient implementation of gpu virtualization in high performance clusters. In *Proceedings of the 2009 International Conference on Parallel Processing*, Euro-Par'09, pages 385–394, Berlin, Heidelberg, 2010. Springer-Verlag.
- [53] Jose Duato, Antonio J. Pena, Federico Silla, Juan C. Fernandez, Rafael Mayo, and Enrique S. Quintana-Orti. Enabling CUDA acceleration within virtual machines using rCUDA. In *Proceedings of the 2011 18th International Conference on High Performance Computing*, HIPC '11, pages 1–10, Washington, DC, USA, 2011. IEEE Computer Society.
- [54] H Esmailzadeh, E Blem, R St.Amant, K Sankaralingam, and D Burger. Dark silicon and the end of multicore scaling. In *Computer Architecture (ISCA), 2011 38th Annual International Symposium on*, pages 365–376, June 2011.
- [55] Patrick Paul "Pat" Gelsinger. Private Communication, 1998.
- [56] G. Giunta, R. Montella, G. Agrillo, and G. Coviello. A gpgpu transparent virtualization component for high performance computing clouds. *Euro-Par 2010-Parallel Processing*, page 379–391, 2010.
- [57] Giulio Giunta, Raffaele Montella, Giuseppe Agrillo, and Giuseppe Coviello. A gpgpu transparent virtualization component for high performance computing clouds. In *European Conference on Parallel Processing*, pages 379–391. Springer, 2010.
- [58] Abel Gordon, Nadav Har'El, Alex Landau, Muli Ben-Yehuda, and Avishay Traeger. Towards exitless and efficient paravirtual i/o. In *Proceedings of the 5th Annual International Systems and Storage Conference*, SYSTOR '12, pages 10:1–10:6, New York, NY, USA, 2012. ACM.
- [59] M. Gottschlag, M. Hillenbrand, J. Kehne, J. Stoess, and F. Bellosa. Logv: Low-overhead gpgpu virtualization. In *High Performance Computing and Communications 2013 IEEE International Conference on Embedded and Ubiquitous Computing (HPCCEUC), 2013 IEEE 10th International Conference on*, pages 1721–1726, Nov 2013.
- [60] Kate Gregory and Ade Miller. C++ amp: accelerated massive parallelism with microsoft visual c++. 2014.
- [61] Vishakha Gupta, Ada Gavrilovska, Karsten Schwan, Harshvardhan Kharche, Niraj Tolia, Vanish Talwar, and Parthasarathy Ranganathan. Gvim: Gpu-accelerated virtual machines. In *Proceedings of the 3rd ACM Workshop on System-level Virtualization for High Performance Computing*, pages 17–24. ACM, 2009.

- [62] Fritz-Rudolf Güntsch. *Logical Design of a Digital Computer with Multiple Asynchronous Rotating Drums and Automatic High Speed Memory Operation*. Doctoral dissertation, Technische Universität Berlin, 1956.
- [63] Nadav Har’El, Abel Gordon, Alex Landau, Muli Ben-Yehuda, Avishay Traeger, and Razya Ladelsky. Efficient and scalable paravirtual i/o system. In *Proceedings of the 2013 USENIX Conference on Annual Technical Conference*, USENIX ATC’13, pages 231–242, Berkeley, CA, USA, 2013. USENIX Association.
- [64] Alex Herrera. Nvidia grid: Graphics accelerated vdi with the visual performance of a workstation. *Nvidia Corp*, 2014.
- [65] Chun-Hsian Huang and Pao-Ann Hsiung. Hardware resource virtualization for dynamically partially reconfigurable systems. *IEEE Embedded Systems Letters*, 1(1):19–23, 2009.
- [66] Yu-Ju Huang, Hsuan-Heng Wu, Yeh-Ching Chung, and Wei-Chung Hsu. Building a kvm-based hypervisor for a heterogeneous system architecture compliant system. In *Proceedings of the 12th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments*, VEE ’16, pages 3–15, New York, NY, USA, 2016. ACM.
- [67] JAIN Jayant, Anirban Sengupta, Rick Lund, Raju Koganty, Xinhua Hong, and Mohan Parthasarathy. Configuring and operating a XaaS model in a datacenter, November 13 2018. US Patent App. 10/129,077.
- [68] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages 675–678. ACM, 2014.
- [69] Norman P. Jouppi, Cliff Young, Nishant Patil, and David Patterson. A domain-specific architecture for deep neural networks. *Commun. ACM*, 61(9):50–59, August 2018.
- [70] Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, and et al. In-datacenter performance analysis of a tensor processing unit. *SIGARCH Comput. Archit. News*, 45(2):1–12, June 2017.
- [71] Shinpei Kato, Karthik Lakshmanan, Ragunathan Rajkumar, and Yutaka Ishikawa. Timegraph: Gpu scheduling for real-time multi-tasking environments. In *Proceedings of the 2011 USENIX Conference on USENIX Annual Technical Conference*, USENIXATC’11, pages 2–2, Berkeley, CA, USA, 2011. USENIX Association.
- [72] Shinpei Kato, Michael McThrow, Carlos Maltzahn, and Scott Brandt. Gdev: First-class gpu resource management in the operating system. In *Proceedings of the 2012 USENIX Conference on Annual Technical Conference*, USENIX ATC’12, pages 37–37, Berkeley, CA, USA, 2012. USENIX Association.
- [73] Khronos Group. *Vulkan 1.0.64 - A Specification*, 2017.
- [74] Tom Kilburn, David BG Edwards, Michael J Lanigan, and Frank H Sumner. One-level storage system. *IRE Transactions on Electronic Computers*, (2):223–235, 1962.

- [75] J. Kim, S. Seo, J. Lee, J. Nah, G. Jo, and J. Lee. Snuc1: an opencl framework for heterogeneous cpu/gpu clusters. In *Proceedings of the 26th ACM international conference on Supercomputing*, page 341–352. ACM, 2012.
- [76] Sangman Kim, Seonggu Huh, Yige Hu, Xinya Zhang, Emmett Witchel, Amir Wated, and Mark Silberstein. Gpunet: Networking abstractions for gpu programs. In *Proceedings of the 11th USENIX Conference on Operating Systems Design and Implementation, OSDI’14*, pages 201–216, Berkeley, CA, USA, 2014. USENIX Association.
- [77] Volodymyr V Kindratenko, Jeremy J Enos, Guochun Shi, Michael T Showerman, Galen W Arnold, John E Stone, James C Phillips, and Wen-mei Hwu. Gpu clusters for high-performance computing. In *Cluster Computing and Workshops, 2009. CLUSTER’09. IEEE International Conference on*, pages 1–8. IEEE, 2009.
- [78] Yossi Kuperman, Eyal Moscovici, and Joel Nider. Paravirtual Remote I/O.
- [79] Yossi Kuperman, Eyal Moscovici, Joel Nider, Razya Ladelsky, Abel Gordon, and Dan Tsafir. Paravirtual remote i/o. In *ACM SIGARCH Computer Architecture News*, volume 44, pages 49–65. ACM, 2016.
- [80] H. Andres Lagar-Cavilla, Niraj Tolia, M. Satyanarayanan, and Eyal de Lara. Vmm-independent graphics acceleration. In *Proceedings of the 3rd International Conference on Virtual Execution Environments, VEE ’07*, pages 33–43, New York, NY, USA, 2007. ACM.
- [81] H. Andres Lagar-Cavilla, Niraj Tolia, M. Satyanarayanan, and Eyal de Lara. Vmm-independent graphics acceleration. In *Proceedings of the 3rd International Conference on Virtual Execution Environments, VEE ’07*, pages 33–43, New York, NY, USA, 2007. ACM.
- [82] Chris Lattner and Vikram Adve. Llvm: A compilation framework for lifelong program analysis & transformation. In *Proceedings of the international symposium on Code generation and optimization: feedback-directed and runtime optimization*, page 75. IEEE Computer Society, 2004.
- [83] Teng Li, Vikram K. Narayana, Esam El-Araby, and Tarek El-Ghazawi. Gpu resource sharing and virtualization on high performance computing systems. In *Proceedings of the 2011 International Conference on Parallel Processing, ICPP ’11*, pages 733–742, Washington, DC, USA, 2011. IEEE Computer Society.
- [84] Tyng-Yeu Liang and Yu-Wei Chang. Gridcuda: A grid-enabled cuda programming toolkit. In *Advanced Information Networking and Applications (WAINA), 2011 IEEE Workshops of International Conference on*, pages 141–146, March 2011.
- [85] Jiuxing Liu, Wei Huang, Bulent Abali, and Dhabaleswar K. Panda. High performance vmm-bypass i/o in virtual machines. In *Proceedings of the Annual Conference on USENIX ’06 Annual Technical Conference, ATEC ’06*, pages 3–3, Berkeley, CA, USA, 2006. USENIX Association.
- [86] Microsoft Inc. *Windows GDI*, 2017.
- [87] Raffaele Montella, Giuseppe Coviello, Giulio Giunta, Giuliano Laccetti, Florin Isaila, and Javier Blas. A general-purpose virtualization service for hpc on cloud computing: an application to gpus. *Parallel Processing and Applied Mathematics*, pages 740–749, 2012.

- [88] Veynu Narasiman, Michael Shebanow, Chang Joo Lee, Rustam Miftakhutdinov, Onur Mutlu, and Yale N Patt. Improving gpu performance via large warps and two-level warp scheduling. In *Proceedings of the 44th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 308–317. ACM, 2011.
- [89] Johns Paul, Jiong He, and Bingsheng He. Gpl: A gpu-based pipelined query processing engine. In *Proceedings of the 2016 International Conference on Management of Data*, pages 1935–1950. ACM, 2016.
- [90] Bo Peng, Haozhong Zhang, Jianguo Yao, Yaozu Dong, Yu Xu, and Haibing Guan. Mdev-nvme: A nvme storage virtualization solution with mediated pass-through. In *2018 USENIX Annual Technical Conference (USENIX ATC’18)*, pages 665–676, 2018.
- [91] Sébastien Pinneterre, Spyros Chiotakis, Michele Paolino, and Daniel Raho. vfpgamanager: A virtualization framework for orchestrated fgpa accelerator sharing in 5g cloud environments. In *2018 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*, pages 1–5. IEEE, 2018.
- [92] Gerald J. Popek and Robert P. Goldberg. Formal requirements for virtualizable third generation architectures. *Commun. ACM*, 17(7):412–421, July 1974.
- [93] Himanshu Raj and Karsten Schwan. High performance and scalable i/o virtualization via self-virtualized devices. In *Proceedings of the 16th International Symposium on High Performance Distributed Computing, HPDC ’07*, pages 179–188, New York, NY, USA, 2007. ACM.
- [94] Kaushik Kumar Ram, Jose Renato Santos, and Yoshio Turner. Redesigning xen’s memory sharing mechanism for safe and efficient i/o virtualization. In *Proceedings of the 2nd conference on I/O virtualization*, pages 1–1. USENIX Association, 2010.
- [95] C. Reano, A. J. Pena, F. Silla, J. Duato, R. Mayo, and E. S. Quintana-Orti. Cu2rcu: Towards the complete rcuda remote gpu virtualization and sharing solution. *20th Annual International Conference on High Performance Computing*, 0:1–10, 2012.
- [96] Carlos Reaño, Antonio J Peña, Federico Silla, José Duato, Rafael Mayo, and Enrique S Quintana-Ortí. Cu2rcu: Towards the complete rcuda remote gpu virtualization and sharing solution. In *2012 19th International Conference on High Performance Computing*, pages 1–10. IEEE, 2012.
- [97] Christopher J Rossbach, Jon Currey, Mark Silberstein, Baishakhi Ray, and Emmett Witchel. PTask: operating system abstractions to manage GPUs as compute devices. In *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles*, pages 233–248. ACM, 2011.
- [98] Christopher J. Rossbach, Yuan Yu, Jon Currey, Jean-Philippe Martin, and Dennis Fetterly. Dandelion: a compiler and runtime for heterogeneous systems. *SOSP’13: The 24th ACM Symposium on Operating Systems Principles*, November 2013.
- [99] Mark Segal and Kurt Akeley. The opengl graphics system: A specification. Technical report, Silicon Graphics Inc., December 2006.
- [100] Sangmin Seo, Gangwon Jo, and Jaejin Lee. Performance characterization of the nas parallel benchmarks in opencl. In *Workload Characterization (IISWC), 2011 IEEE International Symposium on*, pages 137–148. IEEE, 2011.



- [101] Lin Shi, Hao Chen, Jianhua Sun, and Kenli Li. *vcuda: Gpu-accelerated high-performance computing in virtual machines*. *IEEE Transactions on Computers*, 61(6):804–816, 2012.
- [102] Lin Shi, Hao Chen, Jianhua Sun, and Kenli Li. *vcuda: Gpu-accelerated high-performance computing in virtual machines*. *IEEE Trans. Comput.*, 61(6):804–816, June 2012.
- [103] Pci Sig. Single Root I/O Virtualization and Sharing Specification Revision 1.1. Technical report, January 2010.
- [104] M. Silberstein, B. Ford, I. Keidar, and E. Witchel. *Gpufs: Integrating a file system with gpus*. 2013.
- [105] Prakalp Srivastava, Maria Kotsifakou, and Vikram S. Adve. *HPVM: A portable virtual instruction set for heterogeneous parallel systems*. *CoRR*, abs/1611.00860, 2016.
- [106] John E Stone, David Gohara, and Guochun Shi. *Opencl: A parallel programming standard for heterogeneous computing systems*. *Computing in science & engineering*, 12(3):66–73, 2010.
- [107] Jeremy Sugerman, Ganesh Venkitachalam, and Beng-Hong Lim. *Virtualizing i/o devices on vmware workstation’s hosted virtual machine monitor*. In *Proceedings of the General Track: 2001 USENIX Annual Technical Conference*, pages 1–14, Berkeley, CA, USA, 2001. USENIX Association.
- [108] Yusuke Suzuki, Shinpei Kato, Hiroshi Yamada, and Kenji Kono. *Gpuvmm: Why not virtualizing gpus at the hypervisor?* In *USENIX Annual Technical Conference*, pages 109–120, 2014.
- [109] Yusuke Suzuki, Shinpei Kato, Hiroshi Yamada, and Kenji Kono. *Gpuvmm: Why not virtualizing gpus at the hypervisor?* In *Proceedings of the 2014 USENIX Conference on USENIX Annual Technical Conference*, USENIX ATC’14, pages 109–120, Berkeley, CA, USA, 2014. USENIX Association.
- [110] Michael M Swift, Brian N Bershad, and Henry M Levy. *Improving the reliability of commodity operating systems*. In *ACM SIGOPS operating systems review*, volume 37, pages 207–222. ACM, 2003.
- [111] Synergy Research Group, Reno, and NV. *Hyperscale Data Center Count Passed the 500 Milestone in Q3*. <https://www.srgresearch.com/articles/hyperscale-data-center-count-passed-500-milestone-q3>, October 2019. Accessed: 2020-6-9.
- [112] Hiroshi Tezuka, Francis O’Carroll, Atsushi Hori, and Yutaka Ishikawa. *Pin-down cache: A virtual memory management technique for zero-copy communication*. In *Parallel Processing Symposium, 1998. IPPS/SPDP 1998. Proceedings of the First Merged International... and Symposium on Parallel and Distributed Processing 1998*, pages 308–314. IEEE, 1998.
- [113] Kun Tian, Yaozu Dong, and David Cowperthwaite. *A full gpu virtualization solution with mediated pass-through*. In *Proceedings of the 2014 USENIX Conference on USENIX Annual Technical Conference*, USENIX ATC’14, pages 121–132, Berkeley, CA, USA, 2014. USENIX Association.
- [114] Kun Tian, Yaozu Dong, and David Cowperthwaite. *A full gpu virtualization solution with mediated pass-through*. In *USENIX Annual Technical Conference*, pages 121–132, 2014.

- [115] Dimitrios Vasilas, Stefanos Gerangelos, and Nectarios Koziris. VGVM: efficient GPU capabilities in virtual machines. In *International Conference on High Performance Computing & Simulation, HPCS 2016, Innsbruck, Austria, July 18-22, 2016*, pages 637–644, 2016.
- [116] Jan Vesely, Arkaprava Basu, Mark Oskin, Gabriel H. Loh, and Abhishek Bhattacharjee. Observations and Opportunities in Architecting Shared Virtual Memory for Heterogeneous Systems. In *ISPASS*, 2016.
- [117] VMware, X.org, Nouveau. *Tungsten Graphics Shader Infrastructure*, 2012.
- [118] Duy Viet Vu, Oliver Sander, Timo Sandmann, Steffen Baehr, Jan Heidelberger, and Juergen Becker. Enabling Partial Reconfiguration for Coprocessors in Mixed Criticality Multicore Systems using PCI Express Single-Root I/O Virtualization. In *ReConFigurable Computing and FPGAs (ReConFig), 2014 International Conference on*, pages 1–6. IEEE, 2014.
- [119] Lan Vu, Hari Sivaraman, and Rishi Bidarkar. Gpu virtualization for high performance general purpose computing on the esx hypervisor. In *Proceedings of the High Performance Computing Symposium, HPC '14*, pages 2:1–2:8, San Diego, CA, USA, 2014. Society for Computer Simulation International.
- [120] Carl Waldspurger, Emery Berger, Abhishek Bhattacharjee, Kevin Pedretti, Simon Peter, and Chris Rossbach. Sweet spots and limits for virtualization. In *Proceedings of the 12th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments, VEE '16*, pages 177–177, New York, NY, USA, 2016. ACM.
- [121] Carl Waldspurger and Mendel Rosenblum. I/o virtualization. *Commun. ACM*, 55(1):66–73, January 2012.
- [122] Zhenning Wang, Jun Yang, Rami Melhem, Bruce Childers, Youtao Zhang, and Minyi Guo. Simultaneous multikernel gpu: Multi-tasking throughput processors via fine-grained sharing. In *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 358–369. IEEE, 2016.
- [123] Paul Willmann, Scott Rixner, and Alan L. Cox. Protection strategies for direct access to virtualized i/o devices. In *USENIX 2008 Annual Technical Conference, ATC'08*, pages 15–28, Berkeley, CA, USA, 2008. USENIX Association.
- [124] Jingyue Wu, Artem Belevich, Eli Bendersky, Mark Heffernan, Chris Leary, Jacques Pienaar, Bjarke Roune, Rob Springer, Xuetian Weng, and Robert Hundt. Gpucc: An open-source gpgpu compiler. In *Proceedings of the 2016 International Symposium on Code Generation and Optimization, CGO '16*, pages 105–116, New York, NY, USA, 2016. ACM.
- [125] Lei Xia, Jack Lange, Peter Dinda, and Chang Bae. Investigating virtual passthrough i/o on commodity devices. *ACM SIGOPS Operating Systems Review*, 43(3):83–94, 2009.
- [126] Tsung Tai Yeh, Amit Sabne, Putt Sakdhnagool, Rudolf Eigenmann, and Timothy G Rogers. Pagoda: Fine-grained gpu resource virtualization for narrow tasks. In *ACM SIGPLAN Notices*, volume 52, pages 221–234. ACM, 2017.
- [127] Dong Yu, Adam Eversole, Mike Seltzer, Kaisheng Yao, Zhiheng Huang, Brian Guenter, Oleksii Kuchaiev, Yu Zhang, Frank Seide, Huaming Wang, et al. An introduction to computational networks and the computational network toolkit. *Microsoft Technical Report MSR-TR-2014-112*, 2014.

- [128] Hangchen Yu, Arthur M Peters, Amogh Akshintala, and Christopher J Rossbach. Automatic virtualization of accelerators. In *Proceedings of the Workshop on Hot Topics in Operating Systems*, pages 58–65. ACM, 2019.
- [129] Hangchen Yu, Arthur Michener Peters, Amogh Akshintala, and Christopher J. Rossbach. Ava: Accelerated virtualization of accelerators. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '20, page 807–825, New York, NY, USA, 2020. Association for Computing Machinery.
- [130] Hangchen Yu and Christopher J Rossbach. Full virtualization for gpus reconsidered. In *14th Workshop on Duplicating, Deconstructing, and Debunking (WDDD), ISCA*, 2017.
- [131] Jose Fernando Zazo, Sergio Lopez-Buedo, Yury Audzevich, and Andrew W Moore. A PCIe DMA Engine to Support the Virtualization of 40 Gbps FPGA-accelerated Network Appliances. In *ReConFigurable Computing and FPGAs (ReConFig), 2015 International Conference on*, pages 1–6. IEEE, 2015.
- [132] Lingfang Zeng, Yang Wang, Wei Shi, and Dan Feng. An improved xen credit scheduler for i/o latency-sensitive applications on multicores. In *Cloud Computing and Big Data (CloudCom-Asia), 2013 International Conference on*, pages 267–274. IEEE, 2013.
- [133] Kai Zhang, Bingsheng He, Jiayu Hu, Zeke Wang, Bei Hua, Jiayi Meng, and Lishan Yang. G-net: Effective gpu sharing in nfv systems. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*. USENIX Association, 2018.