

Toward Efficient and Realizable Hardware Virtualization

AMOGH AKSHINTALA

Abstract

The proposed dissertation will focus on developer effort and compatibility in software virtualization of CPU ISAs, and software virtualization of specialized compute devices (e.g., GPUs, TPUs) that are programmed through an API.

Although binary translation is a well-established software ISA virtualization technique, given the size and complexity of today’s dominant ISAs, developers are routinely forced to adopt ad-hoc techniques to prioritize development effort. The proposed dissertation will present a principled approach to determine priority among different parts of the ISA. We believe this data will be useful to designers of virtual ISAs and [AA: what was the other thing?](#) as well.

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 provided 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-mediated API-remoting*, that is at once realizable and performant.

1 Introduction

Virtualization has a long and tumultuous history. Virtual memory was first described by German physicist Fritz-Rudolf Güntsch in his doctoral dissertation in 1956 [31] and commercialized [36] in the Cambridge University/Ferranti Inc. Atlas computer. Virtually all computers since then support virtual memory, with most providing hardware units—*Memory Management Unit (MMU)*—to accelerate the virtualization of memory.

Hardware virtualization [7]—the idea of virtualizing the entire computer to enable the simultaneous execution of multiple Operating Systems (OS)—was invented in 1962, and commercialized as the IBM VM-370 [23] 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 [44], and was widely considered unvirtualizable [AA: track down Pat Gelsinger note.](#) 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 have been devised to the problem of virtualizing the x86 ISA, including ISA extensions (e.g., AMD-V and Intel VT-x) to enable the execution of virtualized applications at native speed, only trapping to the hypervisor when the application attempts to perform sensitive operations.

Several other forms of virtualization have also been considered. 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.

The proposed dissertation is primarily concerned with hardware virtualization. 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 hardware virtualization, from Popek and Goldberg’s classical virtualization properties—*equivalence*, *performance*, and *safety* to Bugnion, Nieh and Tsafir’s [21] definition of virtualization—“*the application of the layering principle with enforced modularity such that the exposed resource is identical to the underlying resource*”. While technically correct, all of these definitions are too contrite to be useful. Instead, for the purposes of this dissertation, we concern ourselves with the following overarching goal—*realizable, fair, isolated, and efficient sharing of hardware resources 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 *virtual*, in that it is not directly exposed 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 performance reasons, as in the case of hardware-assisted virtualization schemes like AMD-V and Intel VT-x, all privileged operations still trap to the hypervisor. The interface interposed may be a hardware interface (ISA, Memory, I/O Protocols, etc.) or a software interface (Syscalls, APIs, etc.).

1.1 CPU virtualization

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. While dominant ISAs, such as x86 and ARM, even provide extensions to enable low-overhead virtualization, binary translation results in lower overhead for sequences of sensitive instructions that need to be emulated [10].

When implementing a new binary translator [27] or developing a secure virtual instruction set, **AA: what was 3rd?**, the developer is left to their own devices to answer questions regarding *realizability*, e.g., to prioritize different parts of the ISA during development or to understand the relative value of different parts of the ISA to backwards compatibility. Predictably, developers have typically adopt ad-hoc methodologies to overcome this challenge [20]. We hypothesize that a principled approach to answering these questions lies in understanding the distribution of importance, to users, in the ISA being virtualized. Chapter 1 will present a methodology for determining user preference, and the resulting dataset. Briefly, we estimate the importance of an instruction in the ISA by measuring its frequency of occurrence in applications, and then weighting the frequency data with the likelihood of users installing those applications. This is completed work [11].

1.2 Accelerator virtualization

Compute heavy and data parallel workloads such as graph processing and machine learning have precipitated a Cambrian explosion of specialized processors. These emerging compute devices (e.g., GPGPUs, TPUs, IPUs, IO accelerators), however, pose a challenge to virtualization developers, who once again find themselves 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) are not applicable to these specialized 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 specialized compute devices as they result in the sacrifice of one or more of the essential characteristics listed above. Full-virtualization based schemes, such as GPUvm [49], suffer from massive overheads that essentially negate the speedup that makes the specialized compute unit attractive in the first place. Para-virtual systems, such as SVGA [24] that interpose on low-level interfaces, such as the kernel driver, introduce much lower overhead than full-virtualization based schemes but have poor compatibility, i.e., 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 [55, 26, 45] 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 [25]. Over the last decade, GPUs have been re-purposed for parallel general purpose compute (commonly known as GPGPU). Chapter 2 will present our findings from attempting to extend the SVGA model of GPU virtualization to cover GPGPU virtualization as well. We find that the tight coupling between ISA virtualization and device virtualization in SVGA leads to poor performance for GPGPU compute. We propose a new virtualization scheme, Trillium, that doesn’t rely on ISA virtualization and show that Trillium outperforms all other traditional virtualization schemes while retaining hypervisor interposition. Material presented in Chapter 2 will be drawn from a published paper [12].

Specialized compute units (e.g., Google TPU, Intel QAT, etc.) are typically exposed to developers via a user-space API. The API is typically implemented by a combination of proprietary software that interacts with the hardware through opaque interfaces. Chapters 3, 4 and 5 will explore the performance implications of virtualizing the user-space API for specialized compute accelerators. Chapter 3 will present an overview of AvA, a framework that enables automated virtualization of accelerator APIs. Chapter 4 will focus on the performance implications of API-remoting based virtualization of a single specialized accelerator. Chapters 3 and 4 will draw on material that appeared in a HotOS workshop paper [59] and a full paper that is currently under submission. Chapter 5 (proposed work) will explore 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 remoting, 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 **E**ndpoints interposed, 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 6, along with analysis of traditional virtualization techniques in the context of GPGPUs.

Concretely, the proposed dissertation will evaluate the following hypotheses:

- H 1:** Priority among instructions in an ISA, in the context of binary translation, can be automatically inferred from user preferences. (Chapter 1)
- H 2:** ISA virtualization is untenable for performant virtualization of compute accelerators. (Chapter 2)
- H 3:** Hypervisor-mediated API-remoting is a low-overhead virtualization scheme for API-controlled compute accelerators. (Chapters 3, 4, and 5)
- H 4:** The characteristics of a virtualization technique can be succinctly described by a scheme that explicitly captures the *Interface* interposed, the *Endpoints* interposed on, the *Mechanism* of interposition, the *Transport* used to connect the interposed endpoints, and the mechanism used to *Synthesize* the interposed operation on the host. (Chapter 6)

2 CPU virtualization

In order to understand CPU Virtualization as it is today, it is illuminating to consider the history of the technique, even though a complete treatment of this subject is out of the scope of this proposal (interested readers are instead referred to Bugnion, Nieh, and Tsafirir’s book on this topic—*Hardware and software support for virtualization* [21].)

CPU virtualization as a technique was first considered for the IBM 360 in 1970 [39]. The idea then, as it is today is, was to provide each user with the illusion of having a dedicated machine to themselves, by simultaneously running multiple operating systems on the same machine. The IBM 370 was specifically designed to be *virtualizable* [23], while a concurrent machine, the PDP-10, wasn't. The inability to virtualize the PDP-10 led Popek and Goldberg [44] to formalize the requirements of a Virtual Machine Monitor—*equivalence, safety, and performance*—and three theorems about the virtualizability of an Instruction Set Architecture (ISA). To summarize briefly, a VMM or a hypervisor must meet the following criteria: the virtual machine constructed by the VMM must be indistinguishable from the native machine (equivalence), must not be able to access resources not allocated to it or influence the way these resources are used by other VMs or the hypervisor (safety), and the performance of software executing in the VM must be comparable to native (performance). The first theorem defined what it means for an ISA to be virtualizable. The second theorem was concerned with recursive or nested virtualization. The third theorem presents the necessary conditions for a *Hybrid Virtual Machine*—one that combines direct execution of non-sensitive instructions with emulation for all sensitive instructions—to be designed for ISAs that violate the first theorem.

While virtualization was well understood and in production in the 1970s, all the lessons learned during that time were seemingly lost or ignored by later computer architects. The new ISAs that established their dominance in computing as the mainframe computers of the 1970s were rendered mostly obsolete by the Personal Computer (Intel x86, DEC ALPHA, SUN SPARC, IBM POWER, MIPS, ARM, etc.) were all unvirtualizable. Even though some of them initially supported virtualizable modes, (e.g., Intel x86's 16-bit virtual mode—v8086—allowed the execution of 16-bit software in 32-bit mode.), this support was left on the wayside as the ISA evolved. Virtualization was considered a quirky bad idea from the 1970s.

Hardware virtualization made a come back in the late 1990s and early 2000s [19, 57, 15] as researchers noticed that virtualization was a promising alternative approach to the problem of multi-core scalability. Despite the non-virtualizable nature of the dominant Intel x86 ISA, researchers realized that they could leverage Popek and Goldberg's third theorem to build a Hybrid Virtual Machine using techniques like *binary translation*—the hypervisor inspects the executing binary and replaces any sensitive instructions with one or more non-sensitive instructions that provide equivalent behavior—and shadow paging.

2.1 Determining priority among instructions for Binary Translation

CPU vendors added support for trap-and-emulate based virtualization as extensions to their CPUs in the mid 2000s (e.g., Intel VT-x, AMD-V). These extensions introduced a new execution mode with support for nested paging, and support for automatically trapping to the hypervisor when the processor attempts to execute a sensitive instruction. This newly introduced hardware support for virtualization, was a mixed blessing for the virtualization software.

Researchers at VMware found that a combination of hardware virtualization and binary translation was essential to achieve optimal performance for the application [10]. On the other hand,

naive application of the hardware support of virtualization led to worse performance than when executing the application a software-emulated virtual machine [8].

Binary translation of a gargantuan ISA like Intel x86-64, which has ~3800 instructions, requires prioritizing development effort on the most “important” instructions—trading completeness for simplicity and a quicker development cycle. For instance, the authors of VMware workstation describe an “on-demand implementation” process, where the x86 binary translator focused on just the instructions needed for a target OS; the entire ISA was never supported, and guest OSes such as OS/2 did not work [20]. Similarly, Amit et al. [14] showed that KVM cannot correctly implement certain obscure x86 behaviors in a guest OS. Prioritizing instruction support is a natural and ubiquitous engineering trade-off. Some instructions appear in program binaries more frequently than others, e.g., the MOV instruction (used to move data) is the most common x86-64 instruction. On the contrary, the VFMADDSD instruction, used to express a fused multiplication and addition operation, is relatively rare. Further, many instructions perform similar operations, albeit with subtle distinctions.

What, then, is the basis for assigning priority to instructions? Common approaches include analyzing benchmark suites [18, 32, 16], or execution traces collected in target environments [34]. The ad-hoc nature of this approach leaves many useful questions unanswered: Is the chosen test suite actually representative? What is the path of least effort to support a new ISA in a software tool? What minimum set of instructions must be implemented to run at least one application? What instruction sub-set is sufficient to run the majority of deployed applications?

To paraphrase Hennessy and Patterson [42], the best thing to measure is what actually runs on the user’s system. This chapter will present and analyze a dataset collected from static analysis of all x86-64 ELF binaries in the Ubuntu 16.04 GNU/Linux distribution. We leverage package installation frequency, an approximation of a package’s importance to users, from Ubuntu and Debian popularity contest data [50, 43], to infer the relative importance of an instruction from the percentage of binaries on a given system that contain that instruction. We adapt metrics from a prior study of OS API compatibility [53], specifically, *instruction importance* — the relative importance of a given instruction, and *weighted completeness* — the completeness of a system that implements a subset of the ISA.

This chapter will present:

- An instruction occurrence dataset gathered using static analysis of 9,337 open-source applications in the Ubuntu 16.04 repositories.
- Evaluation of conventional wisdom about ISA usage.
- An iterative plan for developing new tools that use the x86-64 ISA.
- Empirical validation of standard benchmarks.
- An instruction occurrence data visualization tool, and the analysis framework used in this study are available at <http://x86instructionpop.com/>.

This chapter will be drawn from joint work [11] with Bhushan Jain, Chia-che Tsai, Michael Ferdman and Donald E. Porter. **AA: Add some results here from this work.**

3 ISA virtualization is untenable for GPUs

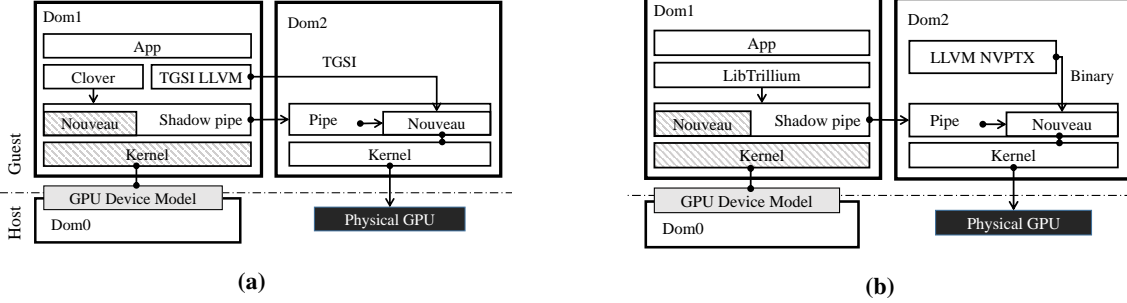


Figure 1: *Xen-SVGA and Trillium designs. (a) The Trillium stack. (b) Xen-SVGA approximates the SVGA model extended to support GPU Compute. (c) The design of Trillium with shadow pipe.*

In many parallel computing domains, compute density and programmability [4, 48, 29] have made GPUs the clear choice for efficiency and performance [2]. Popular machine learning frameworks such as Caffe [35], Tensorflow [6], CNTK [58], and Torch7 [22] 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 [5].

Despite much prior research [56, 33, 9, 54] 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 [13, 51], 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 [17, 37, 45, 55, 26] 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 [46].

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 [25] design in Xen (shown in Figure 1a), 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. Instead, we found that it harms performance, as shown in Figure 2, 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. Figure 1b visually presents the Trillium design.

Trillium is an existence proof of a viable alternative design—hypervisor-mediated API-remoting—

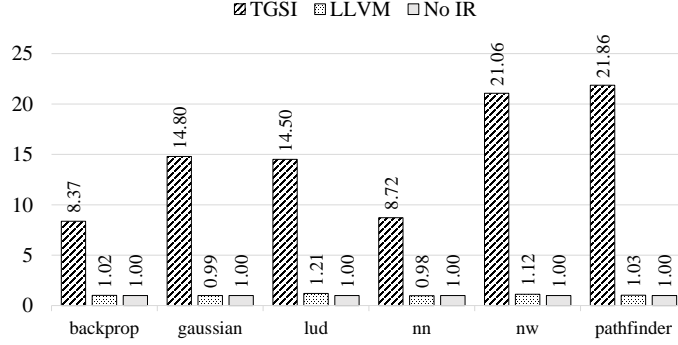


Figure 2: Kernel execution slowdown due to virtual ISAs. *TGSI*: the LLVM TGSI back-end compiler used in Xen-SVGA. *LLVM*: LLVM NVPTX back-end used in Trillium. *No IR*: native NVIDIA compiler.

that preserves desirable virtualization properties such as consolidation, hypervisor interposition, isolation, encapsulation, etc., without requiring full hardware virtualization. While Trillium outperforms GPUvm [49], a full virtualization system, by up to $14\times$ ($5.5\times$ on average) and the para-virtual SVGA-like design by as much as $7.3\times$ ($5.4\times$ on average), it performs worse than a userspace API-remoting framework on average. We believe this is because of a poor choice of API to forward: Trillium forwards the nouveau kernel graphics driver API, which is low enough in the stack that each userspace API function is broken into multiple RPC calls. A better approach would be to forward the userspace API itself (presented in the next chapter).

Concretely, this chapter will show that ISA virtualization is harmful for GPU virtualization, and will lay the groundwork for a new hypothesis—Hypervisor-mediated API-remoting (of the userspace programming framework API) is a realizable, performant, safe and composable virtualization scheme for API-controlled accelerators.

The proposed chapter will draw material from joint work [12] with Hangchen Yu, Arthur M. Peters, and Christopher J. Rossbach.

4 Hypervisor-mediated API-remoting

Hypervisors have not kept up with the pace of accelerator innovation. Specialized hardware and frameworks emerge far faster than hypervisors support them. Many factors contribute to the growing gap, but lack of demand is *not* among them, evinced by the wide variety of accelerators currently available from cloud providers [1, ?, ?, ?, ?, 40, ?]. The challenge is technical: virtualizing accelerators is hard.

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 ??) 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.

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 [?, ?, 55, 26, ?], AVA preserves hypervisor-level resource management and strong isolation using a novel technique called *Hypervisor Interposed Remote Acceleration (HIRA)*. 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. As Figure ?? shows, a solution that tracks API framework evolution can track hardware evolution as well.

AVA supports a broad range of currently-shipping compute offload 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 [?]. 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. We make the following contributions:

- We demonstrate 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, LAPI, 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.

5 IEMTS

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 phys-

ical 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 [?].

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 [52, 3, 38, 41] 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)* [30, 25, 28, 47] 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. We suggest an alternative framework called IEMTS that teases apart design axes that are implicitly and unnecessarily intertwined in much of the literature. IEMTS enables clearer understanding of trade-offs in prior designs and provides a model for comparison of alternative designs. While we hypothesize that our findings generalize to a broad array of accelerators we focus on GPGPUs because they are easily available and have rich research literature.

AA: add the table and some discussion of what IEMTS actually is.

6 vTask

Specialized compute accelerators (e.g., TPUs, GPUs, IPUs) are rapidly proliferating in the data center. This trend is fueled by the emergence of workloads with an appetite for massive amounts of both data and compute, at a time when transistor shrinkage (Moore’s law) has slowed down and the Dark Silicon problem looms large. Hyperscalers are even deploying custom silicon in their data centers, e.g., Google Cloud TPUs, Amazon Nitro.

API-remoting based virtualization has been shown [12, 59] to be the only effective mechanism for sharing these specialized 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 [cite BitFusion acquisition].

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

per framework.

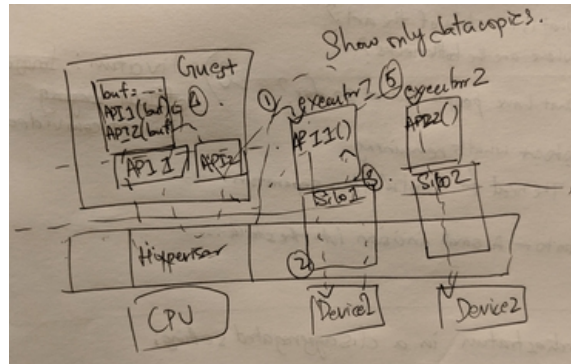


Figure 3: Data processed by two API stacks must pass through the guest application

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 3 illustrates this scenario: when an API-1 function is invoked, associated data is copied from the guest application to executor-1, and then to device-1’s memory. Once the function finishes executing on device-1, the result is copied back to the guest VM. When a function from API-2 is invoked, the same data (i.e., the output of the API-1 function) is copied to executor-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 accelerator silo, or 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 `cpaDCCCompressData2()` 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.

vTask is an application-transparent data orchestration system that optimizes data movement among accelerators virtualized via API-remoting. vTask leverages information from API annotations (cite AvA) to track data buffers across the guest application, the executors servicing API calls made by the guest, and the accelerator hardware. vTask optimizes data movement across these components while ensuring that a coherent view of the data buffer is presented to anyone attempting to read the data. vTask requires no changes to the guest application or annotations of any kind from the application programmer; annotations provided to virtualize the API (by the device or virtualization vendor) are leveraged to infer semantics of the data buffers managed.

We prototyped vTask in AvA, a state-of-the-art para-virtual API-remoting system for KVM. vTask relies on device-side buffer allocation and deallocation API calls to determine buffer lifetime. vTask also leverages special annotations to explicitly specify buffer lifetime provided by

LAPIS, AvA’s API description language. vTask implements a simple MESI-style coherence protocol to track spatial validity of data (i.e., to track where the latest data is present). vTask relies on 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.

vTask can 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. vTask takes care 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.

Evaluating vTask on Y applications that different combinations of the ten accelerators supported by AvA shows that vTask improves performance by X% on average when all devices are on the same machine, and Z% when the devices are on remote hosts.

This work makes the following contributions:

- We identify and measure performance problems with using multiple accelerators via API-remoting
- We propose an application-transparent data orchestration framework, vTask, which transparently manages data buffers in the hypervisor and elides redundant data movement.
- We prototyped vTask in AvA with support for both local and remote accelerators, and used it to evaluate the performance of applications that use different combinations of the ten accelerators supported by AvA. vTask improves application performance by XX

7 Plan of Work

Task	Deadline
Dissertation proposal and oral exam	Nov. 2019
Implement vTask in AvA (submit to ATC’20)	15 Jan. 2020
Dissertation draft to committee	1 Mar. 2020
Dissertation defense	late Mar. 2020
Submit dissertation to graduate school	12 Apr. 2020

Table 1: *Proposed timeline.*

References

- [1] Amazon EC2 instance types. <https://aws.amazon.com/ec2/instance-types/>. Accessed: 2017-04.
- [2] GPU Applications Catalog. <https://www.nvidia.com/en-us/data-center/gpu-accelerated-applications/catalog/>. Jan. 2018.

- [3] KVMGT - the implementation of intel gvt-g(full gpu virtualization) for KVM. <https://lwn.net/Articles/624516/>. 2014.
- [4] NVIDIA CUDA 4.0. <http://developer.nvidia.com/cuda-toolkit-40>. 2011.
- [5] TOP500 Supercomputer Sites. <https://www.top500.org/lists/2018/11/>, 2019.
- [6] 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.
- [7] R.J. Adair. *A Virtual Machine System for the 360/40*. IBM Cambridge Scientific Center report. International Business Machines Corporation, Cambridge Scientific Center, 1966.
- [8] Keith Adams and Ole Agesen. A Comparison of Software and Hardware Techniques for x86 Virtualization. *SIGARCH Comput. Archit. News*, 34(5):2–13, October 2006.
- [9] 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.
- [10] 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.
- [11] Amogh Akshintala, Bhushan Jain, Chia-Che Tsai, Michael Ferdman, and Donald E Porter. x86-64 instruction usage among c/c++ applications. In *Proceedings of the 12th ACM International Conference on Systems and Storage*, pages 68–79. ACM, 2019.
- [12] 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.
- [13] Inc or Its Affiliates Amazon Web Services. Amazon EC2 P3 Instances. <https://aws.amazon.com/ec2/instance-types/p3/>. Accessed: 2018-2-6.
- [14] Nadav Amit, Dan Tsafir, Assaf Schuster, Ahmad Ayoub, and Eran Shlomo. Virtual cpu validation. In *ACM Symposium on Operating Systems Principles (SOSP)*, October 2015.
- [15] 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 *Proceedings of the Nineteenth ACM Symposium on Operating Systems Principles, SOSP ’03*, pages 164–177, New York, NY, USA, 2003. ACM.
- [16] Christian Bienia. *Benchmarking Modern Multiprocessors*. PhD thesis, Princeton University, January 2011.
- [17] 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.

- [18] James Bucek, Klaus-Dieter Lange, and J3akim v. Kistowski. Spec cpu2017: Next-generation compute benchmark. In *Companion of the 2018 ACM/SPEC International Conference on Performance Engineering, ICPE '18*, pages 41–42, New York, NY, USA, 2018. ACM.
- [19] 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.
- [20] 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.
- [21] Edouard Bugnion, Jason Nieh, and Dan Tsafir. Hardware and software support for virtualization. *Synthesis Lectures on Computer Architecture*, 12(1):1–206, 2017.
- [22] Ronan Collobert, Koray Kavukcuoglu, and Cl3ment Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, number EPFL-CONF-192376, 2011.
- [23] R. J. Creasy. The origin of the vm/370 time-sharing system. *IBM J. Res. Dev.*, 25(5):483–490, September 1981.
- [24] Micah Dowty and Jeremy Sugerman. Gpu virtualization on vmware’s hosted i/o architecture. *SIGOPS Oper. Syst. Rev.*, 43(3):73–82, July 2009.
- [25] Micah Dowty and Jeremy Sugerman. Gpu virtualization on vmware’s hosted i/o architecture. *ACM SIGOPS Operating Systems Review*, 43(3):73–82, 2009.
- [26] 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.
- [27] Alex Fishman, Mike Rapoport, Evgeny Budilovsky, and Izik Eidus. HVX: Virtualizing the cloud. In *Presented as part of the 5th USENIX Workshop on Hot Topics in Cloud Computing*, San Jose, CA, 2013. USENIX.
- [28] 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.
- [29] Kate Gregory and Ade Miller. C++ amp: accelerated massive parallelism with microsoft visual c++. 2014.
- [30] 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.

- [31] 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.
- [32] John L Henning. SPEC CPU2006 benchmark descriptions. *SIGARCH Comput. Archit. News*, 34(4):1–17, 2006.
- [33] 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.
- [34] R. Jagtap, S. Diestelhorst, A. Hansson, M. Jung, and N. When. Exploring system performance using elastic traces: Fast, accurate and portable. In *2016 International Conference on Embedded Computer Systems: Architectures, Modeling and Simulation (SAMOS)*, pages 96–105, July 2016.
- [35] 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.
- [36] 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.
- [37] 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.
- [38] 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.
- [39] R. A. Meyer and L. H. Seawright. A virtual machine time-sharing system. *IBM Systems Journal*, 9(3):199–218, Sep 1970.
- [40] Microsoft. Microsoft azure goes back to rack servers with project olympus, 2017.
- [41] 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.
- [42] David A. Patterson and John L. Hennessy. *Computer Organization and Design, Fifth Edition: The Hardware/Software Interface*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 5th edition, 2013.
- [43] Avery Pennarun, Bill Allombert, and Petter Reinholdtsen. Debian Popularity Contest. <http://popcon.debian.org>, 2018.

- [44] Gerald J. Popek and Robert P. Goldberg. Formal requirements for virtualizable third generation architectures. *Commun. ACM*, 17(7):412–421, July 1974.
- [45] 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.
- [46] 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.
- [47] 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.
- [48] 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.
- [49] Yusuke Suzuki, Shinpei Kato, Hiroshi Yamada, and Kenji Kono. Gpvm: Why not virtualizing gpus at the hypervisor? In *USENIX Annual Technical Conference*, pages 109–120, 2014.
- [50] The Ubuntu Web Team, Avery Pennarun, Bill Allombert, and Petter Reinholdtsen. Ubuntu Popularity Contest. <http://popcon.ubuntu.com>, 2018.
- [51] 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.
- [52] 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.
- [53] Chia-Che Tsai, Bhushan Jain, Nafees Ahmed Abdul, and Donald E Porter. A Study of Modern Linux API Usage and Compatibility: What to Support when You’re Supporting. In *Proceedings of the ACM European Conference on Computer Systems (EuroSys)*, London, United Kingdom, 2016.
- [54] 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.
- [55] 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.

- [56] 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.
- [57] Andrew Whitaker, Marianne Shaw, and Steven D. Gribble. Scale and performance in the denali isolation kernel. *SIGOPS Oper. Syst. Rev.*, 36(SI):195–209, December 2002.
- [58] 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.
- [59] 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.