

GPU Virtualization on VMware's Hosted I/O Architecture

Micah Dowty
VMware, Inc.
3401 Hillview Ave, Palo Alto, CA 94304
micah@vmware.com

Jeremy Sugerman
VMware, Inc.
3401 Hillview Ave, Palo Alto, CA 94304
yoel@vmware.com

ABSTRACT

Modern graphics co-processors (GPUs) can produce high fidelity images several orders of magnitude faster than general purpose CPUs, and this performance expectation is rapidly becoming ubiquitous in personal computers. Despite this, GPU virtualization is a nascent field of research. This paper introduces a taxonomy of strategies for GPU virtualization and describes in detail the specific GPU virtualization architecture developed for VMware's hosted products (VMware Workstation and VMware Fusion).

We analyze the performance of our GPU virtualization with a combination of applications and microbenchmarks. We also compare against software rendering, the GPU virtualization in Parallels Desktop 3.0, and the native GPU. We find that taking advantage of hardware acceleration significantly closes the gap between pure emulation and native, but that different implementations and host graphics stacks show distinct variation. The microbenchmarks show that our architecture amplifies the overheads in the traditional graphics API bottlenecks: draw calls, downloading buffers, and batch sizes.

Our virtual GPU architecture runs modern graphics-intensive games and applications at interactive frame rates while preserving virtual machine portability. The applications we tested achieve from 86% to 12% of native rates and 43 to 18 frames per second with VMware Fusion 2.0.

Categories and Subject Descriptors

I.3.4 [Graphics Utilities]: Virtual device interfaces

General Terms

Design, Performance

Keywords

I/O Virtualization, Virtual Device, GPU

1. INTRODUCTION

Over the past decade, virtual machines (VMs) have become increasingly popular as a technology for multiplexing both desktop and server commodity x86 computers. Over that time, several critical challenges in CPU virtualization have been overcome, and there are now both software and hardware techniques for virtualizing CPUs with very low overheads [2]. I/O virtualization, however, is still very much an open problem and a wide variety of strategies are used. Graphics co-processors (GPUs) in particular present a challenging mixture of broad complexity, high performance, rapid change, and limited documentation.

Modern high-end GPUs have more transistors, draw more power, and offer at least an order of magnitude more computational performance than CPUs. At the same time, GPU acceleration has extended beyond entertainment (e.g., games and video) into the basic windowing systems of recent operating systems and is starting to be applied to non-graphical high-performance applications including protein folding, financial modeling, and medical image processing. The rise in applications that exploit, or even assume, GPU acceleration makes it increasingly important to expose the physical graphics hardware in virtualized environments. Additionally, virtual desktop infrastructure (VDI) initiatives have led many enterprises to try to simplify their desktop management by delivering VMs to their users. Graphics virtualization is extremely important to a user whose primary desktop runs inside a VM.

GPUs pose a unique challenge in the field of virtualization. Machine virtualization multiplexes physical hardware by presenting each VM with a *virtual device* and combining their respective operations in the hypervisor platform in a way that utilizes native hardware while preserving the illusion that each guest has a complete stand-alone device. Graphics processors are extremely complicated devices. In addition, unlike CPUs, chipsets, and popular storage and network controllers, GPU designers are highly secretive about the specifications for their hardware. Finally, GPU architectures change dramatically across generations and their generational cycle is short compared to CPUs and other devices. Thus, it is nearly intractable to provide a virtual device corresponding to a real modern GPU. Even starting with a complete implementation, updating it for each new GPU generation would be prohibitively laborious. Thus, rather than modeling a complete modern GPU, VMware's

primary approach paravirtualizes: it delivers an idealized software-only GPU and our own custom graphics driver for interfacing with the guest operating system.

The main technical contributions of this paper are (1) a taxonomy of GPU virtualization strategies—both emulated and passthrough-based, (2) an overview of the virtual graphics stack in VMware’s hosted architecture, and (3) an evaluation and comparison of VMware Fusion’s 3D acceleration with other approaches. We find that a hosted model [18] is a good fit for handling complicated, rapidly changing GPUs while the largely asynchronous graphics programming model is still able to efficiently utilize GPU hardware acceleration.

The rest of this paper is organized as follows. Section 2 provides background and some terminology. Section 3 describes a taxonomy of strategies for exposing GPU acceleration to VMs. Section 4 describes the device emulation and rendering thread of the graphics virtualization in VMware products. Section 5 evaluates the 3D acceleration in VMware Fusion. Section 6 summarizes our findings and describes potential future work.

2. BACKGROUND

While CPU virtualization has a rich research and commercial history, graphics hardware virtualization is a relatively new area. VMware’s virtual hardware has always included a display adapter, but it initially included only basic 2D support [21]. Experimental 3D support did not appear until VMware Workstation 5.0 (April 2005). Both Blink [7] and VMGL [10] used a user-level Chromium-like approach [8] to accelerate fixed function OpenGL in Linux and other UNIX-like guests. Parallels Desktop 3.0 [14] accelerates some OpenGL and Direct3D guest applications with a combination of Wine and proprietary code [22, 15], but loses its interposition properties while those applications are running. Finally, at the most recent Intel Developer Forum, Parallels presented a demo that dedicates an entire native GPU to a single virtual machine using Intel’s VT-d [9, 1]. Parallels has only recently started shipping this solution as part of Parallels Workstation 4.0 Extreme. We have not evaluated its behavior.

The most immediate application for GPU virtualization is desktop virtualization. While server workloads still form the core use case for virtualization, desktop virtualization is now the strongest growth market [12]. Desktop users run a diverse array of applications, including games, video, CAD, and visualization software. Windows Vista, Mac OS X, and recent Linux distributions all include GPU-accelerated windowing systems. Furthermore, an increasing number of ubiquitous applications are adopting GPU acceleration. Adobe Flash Player 10, the latest version of a product which currently reaches 99.0% of Internet viewers [6], includes support for GPU acceleration. There is a user expectation that virtualized applications will “just work”, and this increasingly includes having access to their graphics card.

2.1 GPU Hardware

This section will briefly introduce GPU hardware. A full discussion of GPU architecture and programming is outside

the scope of this paper.

Graphics hardware has experienced a continual evolution from mere CRT controllers to powerful programmable stream processors. Early graphics accelerators could draw rectangles or bitmaps. Later graphics accelerators could rasterize triangles and transform and light them in hardware. With current PC graphics hardware, formerly fixed-function transformation and shading has become generally programmable. Applications use high-level Application Programming Interfaces (APIs) to configure the graphics pipeline, and provide *shader* programs which perform application specific per-vertex and per-pixel processing on the GPU [11, 5].

Future GPUs are expected to continue providing increased programmability. Compute APIs like OpenCL and CUDA already facilitate the use of GPUs for non-graphical computation. Intel recently announced its Larrabee [17] architecture, a potentially disruptive technology which follows this hardware trend to its extreme by implementing a GPU using an array of specialized x86 CPU cores.

With the recent exception of many AMD GPUs, for which open documentation is now available [3], GPU hardware is proprietary. NVIDIA’s hardware documentation, for example, is a closely guarded trade secret. Nearly all graphics applications interact with the GPU via a standardized API such as Microsoft’s DirectX or the vendor-independent OpenGL standard.

3. GPU VIRTUALIZATION TAXONOMY

This section explores the GPU virtualization approaches we have considered at VMware. We use four primary criteria for judging them: performance, fidelity, multiplexing, and interposition.

Performance and fidelity emphasize minimizing the cost of virtualization. We use *performance* to mean the speed of operation, and *fidelity* as the breadth and quality of the available GPU features. Higher performance improves a system’s ability to run interactively. Higher fidelity may increase compatibility with and visual quality of workloads.

Multiplexing and interposition emphasize the added value of virtualization. *Multiplexing* is the ability for multiple virtual machines to share the same physical GPU. To be practical, multiplexing requires some form of secure isolation between VMs. *Interposition* allows virtualization software to mediate access between a virtual machine and the physical hardware. Some degree of interposition is required for basic virtualization features such as execution checkpointing and suspend-to-disk. Higher degrees of interposition can allow for more advanced features: Live migration, fault-tolerant execution, disk image portability, and many other features are enabled by insulating the guest from physical hardware dependencies.

We observe that different use cases weight the criteria differently. For example, a VDI deployment values high VM-to-GPU consolidation ratios (e.g., multiplexing) while a con-

sumer running a VM to access a game or CAD application unavailable on his host values performance and likely fidelity. A tech support person maintaining a library of different configurations and an IT administrator running server VMs are both likely to value portability and secure isolation (interposition).

Since these criteria are often in opposition (e.g., performance at the expense of interposition), we describe several possible designs. Rather than give an exhaustive list, we describe points in the design space which highlight interesting trade-offs and capabilities. Table 1 summarizes these parameters. At a high level, we group GPU virtualization techniques into two categories: front-end (application facing) and back-end (hardware facing).

3.1 Front-end Virtualization

Front-end virtualization introduces a virtualization boundary at a relatively high level in the stack, and runs the graphics driver in the host/hypervisor. This approach does not rely on any GPU vendor- or model-specific details. Access to the GPU is entirely mediated through the vendor provided APIs and drivers on the host while the guest only interacts with software. Current GPUs allow applications many independent “contexts” so multiplexing is easy. Interposition is not a given—unabstracted details of the GPU’s capabilities may be exposed to the virtual machine for fidelity’s sake—but it is straightforward to achieve if desired. However, there is a performance risk if too much abstraction occurs in pursuit of interposition.

Front-end techniques exist on a continuum between two extremes: *API remoting*, in which graphics API calls are blindly forwarded from the guest to the external graphics stack via remote procedure call, and *device emulation*, in which a virtual GPU is emulated and the emulation synthesizes host graphics operations in response to actions by the guest device drivers. These extremes have serious disadvantages that can be overcome by intermediate solutions. Pure API remoting is simple to implement, but completely sacrifices interposition and involves wrapping and forwarding an extremely broad collection of entry points. Pure emulation of a modern GPU delivers excellent interposition and implements a narrower interface, but a highly complicated and under-documented one.

Our hosted GPU acceleration employs front-end virtualization and is described in Section 4. Parallels Desktop 3.0, Blink, and VMGL are other examples of front-end virtualization. Parallels appears to be closest to pure API remoting, as virtual machine execution state cannot be saved to disk while OpenGL or Direct3D applications are running. VMGL uses Chromium to augment its remoting with OpenGL state tracking and Blink implements something similar. This allows them suspend-to-disk functionality and reduces the amount of data which needs to be copied across the virtualization boundary.

3.2 Back-end Virtualization

Back-end techniques run the graphics driver stack inside the virtual machine with the virtualization boundary between

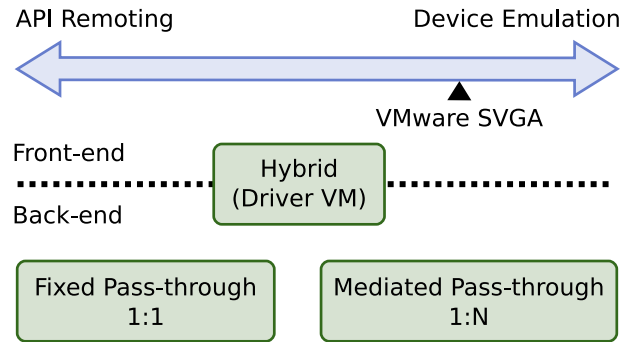


Figure 1: A visual representation of our GPU virtualization taxonomy

the stack and physical GPU hardware. These techniques have the potential for high performance and fidelity, but multiplexing and especially interposition can be serious challenges. Since a VM interacts directly with proprietary hardware resources, its execution state is bound to the specific GPU vendor and possibly the exact GPU model in use. However, exposure to the native GPU is excellent for fidelity: a guest can likely exploit the full range of hardware abilities.

The most obvious back-end virtualization technique is *fixed pass-through*: the permanent association of a virtual machine with full exclusive access to a physical GPU. Recent chipset hardware features, such as Intel’s VT-d, make fixed pass-through practical without requiring any special knowledge of a GPU’s programming interfaces. However, fixed pass-through is not a general solution. It completely forgoes any multiplexing or interposition, and packing machines with one GPU per virtual machine (plus one for the host) is not feasible.

One extension of fixed pass-through is *mediated pass-through*. As mentioned, GPUs already support multiple independent contexts and mediated pass-through proposes dedicating just a context, or set of contexts, to a virtual machine rather than an entire GPU. In this approach, high-bandwidth operations (command buffer submission, vertex and texture DMA) would be performed using memory or MMIO resources which are mapped directly to the physical GPU. Low-bandwidth operations (resource allocation, legacy features) may be implemented using fully virtualized resources. This architecture is illustrated by Figure 2.

This allows multiplexing, but incurs additional costs: the GPU hardware must implement multiple isolated contexts in a way that they can be mapped to different virtual machines efficiently and securely. The host/hypervisor must have enough of a hardware driver to allocate and manage GPU resources such as memory and contexts. Also, the logical GPUs which appear in each VM may or may not have the same hardware interface which would be exposed by an equivalent physical GPU. This means that mediated pass-through may require changes to the guest device drivers.

Mediated pass-through is capable of supporting basic interposition features such as saving execution state to disk, but

Technique	Performance	Fidelity	Multiplexing	Interposition
Software Rendering	very low	high	yes	perfect
Front-end	medium	medium	yes	good
Fixed pass-through	high	high	no	none
Mediated pass-through	high	high	yes	some

Table 1: Design space trade-offs for each virtualized graphics technique

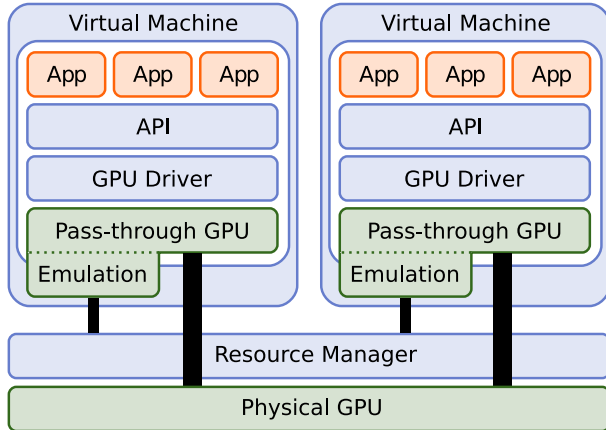


Figure 2: Mediated pass-through architecture

more advanced interposition features such as virtual machine migration present significant challenges, especially if one wants compatibility across different models or families of physical GPU. Features like live migration or deterministic execution replay are particularly difficult. A number of tactics using paravirtualization or standardization of a subset of hardware interfaces can potentially unlock these additional interposition features. Analogous techniques for networking hardware were presented at VMworld 2008 [23].

There may also be long-term security and reliability concerns in a mediated pass-through solution. GPUs often contain hardware errata which must be worked around by drivers. If a virtualizable GPU contained a hardware flaw which allowed one context to access the memory of another context, or which allowed one context to hang the GPU, there might be no efficient and secure way to work around this errata in software. This same concern exists for hardware virtualization of other devices, including CPUs, however the high complexity and low architectural transparency of GPUs could make errata a more significant threat.

3.3 Hybrid Implementations

By combining front-end virtualization and back-end virtualization, it is possible to use front-end virtualization in environments without a native graphics stack, such as a bare metal hypervisor or a security-enhanced operating system. Pass-through becomes a mechanism for securely deploying an off-the-shelf graphics driver in this environment. Front-end virtualization is then used to multiplex the GPU(s) among several virtual machines while maintaining high degrees of VM portability and isolation.

We call this the *Driver VM* approach. The bulk of the GPU driver stack runs in an isolated low-privilege VM. The physical GPU is securely exposed to the driver VM using fixed pass-through, implemented with the use of I/O virtualization support in the chipset. This requires no special knowledge of the GPU’s programming model, and it can be compatible with unmodified GPU drivers. To multiplex the GPU and expose it to other VMs, any form of front-end virtualization may be used. This approach was recently described by HP Laboratories as part of a secure desktop virtualization architecture for Xen [4].

4. VMWARE’S VIRTUAL GPU

All of VMware’s products include a virtual display adapter that supports VGA and basic high resolution 2D graphics modes. On VMware’s hosted products, this adapter also provides accelerated GPU virtualization using a front-end virtualization strategy. To satisfy our design goals, we chose a flavor of front-end virtualization which provides good portability and performance, and which integrates well with existing operating system driver models. Our approach is most similar to the *device emulation* approach above, but it includes characteristics similar to those of *API remoting*. The in-guest driver and emulated device communicate asynchronously with VMware’s Mouse-Keyboard-Screen (MKS) abstraction. The MKS runs as a separate thread and owns all of our access to the host GPU (and windowing system in general).

4.1 SVGA Device Emulation

Our virtual GPU takes the form of an emulated PCI device, the *VMware SVGA II* card. No physical instances of this card exist, but our virtual implementation acts like a physical graphics card in most respects. The architecture of our PCI device is outlined by Figure 3. Inside the VM, it interfaces with a device driver we supply for common guest operating systems. Currently only the Windows XP driver has 3D acceleration support. Outside the VM, a user-level device emulation process is responsible for handling accesses to the PCI configuration and I/O space of the SVGA device. The virtual hardware protocol is publicly documented [20].

Our virtual graphics device provides three fundamental kinds of virtual hardware resources: registers, Guest Memory Regions (GMRs), and a FIFO command queue.

PCI BAR0 controls access to a set of I/O registers for infrequent operations that must be emulated synchronously. Other registers may be located in the faster *FIFO Memory* region, which is backed by plain system memory on the host. I/O space registers are used for mode switching, GMR

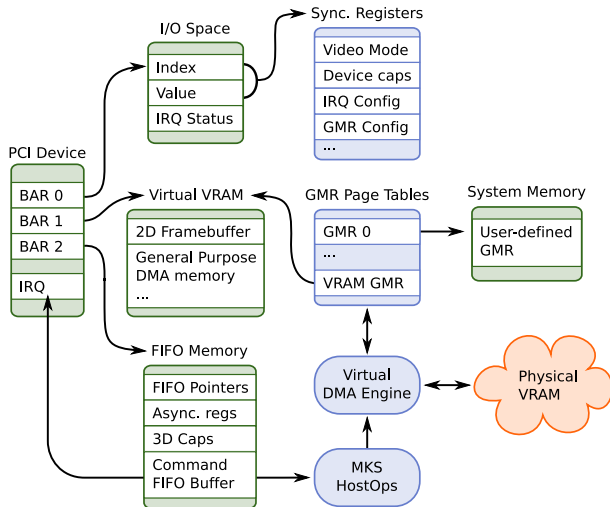


Figure 3: VMware SVGA II device architecture

setup, IRQ acknowledgement, versioning, and for legacy purposes. FIFO registers include large data structures, such as the host’s 3D rendering capabilities, and fast-moving values such as the mouse cursor location—this is effectively a shared memory region between the guest driver and the MKS.

GMRs are an abstraction for guest owned memory which the virtual GPU is allowed to read or write. GMRs can be defined by the guest’s video driver using arbitrary discontinuous regions of guest system memory. Additionally, there always exists one default GMR: the device’s “virtual VRAM.” This VRAM is actually host system memory, up to 128 MB, mapped into PCI memory space via BAR1. The beginning of this region is reserved as a 2D framebuffer.

In our virtual GPU, physical VRAM is not directly visible to the guest. This is important for portability, and it is one of the primary trade-offs made by our front-end virtualization model. To access physical VRAM *surfaces* like textures, vertex buffers, and render targets, the guest driver schedules an asynchronous DMA operation which transfers data between a surface and a GMR. In every surface transfer, this DMA mechanism adds at least one copy beyond the normal overhead that would be experienced in a non-virtualized environment or with back-end virtualization. Often only this single copy is necessary, because the MKS can provide the host’s OpenGL or Direct3D implementation with direct pointers into mapped GMR memory.¹

This virtual DMA model has the potential to far outperform a pure API remoting approach like VMGL or Chromium, not only because so few copies are necessary, but because the guest driver may cache lockable Direct3D buffers directly in GMR memory. Unlike a high-level socket or pipe abstraction, this low-level shared memory system allows the vir-

¹Even this single copy may be avoidable if the virtual device implementation collaborates closely with the physical GPU driver in order to convert virtual DMA operations directly into physical DMA operations that operate on the virtual machine’s memory.

tual machine’s 3D driver stack to implement different GMR memory usage models for different kinds of surfaces. Some surfaces may use the shared memory like a ring buffer for updated regions, some surfaces may store a single cached copy of their data in GMR memory, and some surfaces need not be resident in GMR memory at all.

Like a physical graphics accelerator, the SVGA device processes commands asynchronously via a lockless FIFO queue. This queue, several megabytes in size, occupies the bulk of the FIFO Memory region referenced by BAR2. During un-accelerated 2D rendering, FIFO commands are used to mark changed regions in the framebuffer, informing the MKS to copy from the guest framebuffer to the physical display. During 3D rendering, the FIFO acts as transport layer for our architecture-independent SVGA3D rendering protocol. FIFO commands also initiate all DMA operations, perform hardware accelerated blits, and control accelerated video and mouse cursor overlays.

We deliver host to guest notifications via a virtual interrupt. Our virtual GPU has multiple interrupt sources which may be programmed via FIFO registers. To measure the host’s command execution progress, the guest may insert FIFO *fence* commands, each with a unique 32-bit ID. Upon executing a fence, the host stores its value in a FIFO register and optionally delivers an interrupt. This mechanism allows the guest to very efficiently check whether a specific command has completed yet, and to optionally wait for it by sleeping until a *fence goal* interrupt is received. A fast way to check command execution progress is critical for efficiently managing DMA memory in the guest driver.

The SVGA3D protocol is a simplified and idealized adaptation of the Direct3D API. It has a minimal number of distinct commands. Drawing operations are expressed using a single flexible vertex/index array notation. All host VRAM resources, including 2D textures, 3D textures, cube environment maps, render targets, and vertex/index buffers are represented using a homogeneous surface abstraction. Shaders are written in a variant of Direct3D’s bytecode format, and most fixed-function render states are based on Direct3D render state.

This protocol acts as a common interchange format for GPU commands and state. The guest contains API implementations which produce SVGA3D commands rather than commands for a specific GPU. This provides an opportunity to actively trade capability for portability. The host can control which of the physical GPU’s features are exposed to the guest. As a result, VMs using SVGA3D are widely portable between different physical GPUs. It is possible to suspend a live application to disk, move it to a different host with a different GPU or MKS backend, and resume it. Even if the destination GPU exposes fewer capabilities via SVGA3D, in some cases our architecture can use its layer of interposition as an opportunity to emulate missing features in software. This portability assurance is critical for preventing GPU virtualization from compromising the core value propositions of machine virtualization.

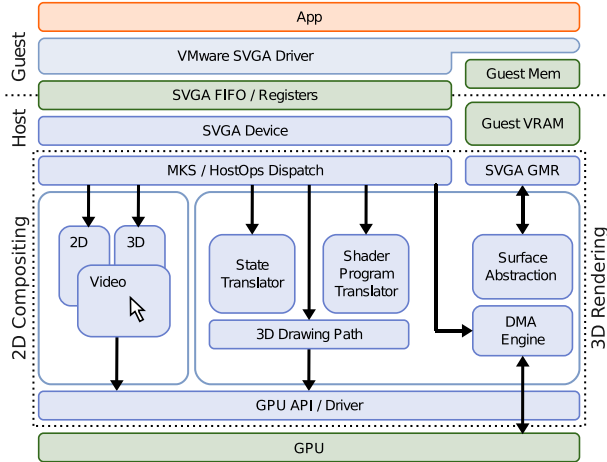


Figure 4: The virtual graphics stack. The MKS/HostOps Dispatch and rendering occur asynchronously in their own thread.

4.2 Rendering

This FIFO design is inherently asynchronous. All host-side rendering happens in the MKS thread, while the guest’s virtual CPUs execute concurrently. As illustrated in Figure 4, access to the physical GPU is mediated first through the GPU vendor’s driver running in the host OS, and secondly via the Host Operations (HostOps) backends in the MKS. The MKS has multiple HostOps backend implementations including GDI and X11 backends to support basic 2D graphics on all Windows and Linux hosts, a VNC server for remote display, and 3D accelerated backends written for both Direct3D and OpenGL. In theory we need only an OpenGL backend to support Windows, Linux, and Mac OS hosts; however we have found Direct3D drivers to be of generally better quality, so we use them when possible. Additional backends could be written to access GPU hardware directly.

The guest video driver writes commands into FIFO memory, and the MKS processes them continuously on a dedicated rendering thread. This design choice is critical for performance, however it introduces several new challenges in synchronization. In part, this is a classic producer-consumer problem. The FIFO requires no host-guest synchronization as long as it is never empty nor full, but the host must sleep any time the FIFO is empty, and the guest must sleep when it is full. The guest may also need to sleep for other reasons. The guest video driver must implement some form of flow control, so that video latency is not unbounded if the guest submits FIFO commands faster than the host completes them. The driver may also need to wait for DMA completion, either to recycle DMA memory or to read back results from the GPU. To implement this synchronization efficiently, the FIFO requires both guest to host and host to guest notifications.

The MKS may periodically poll the command FIFO when a console is attached, in order to provide a “virtual vertical refresh rate” of 25 to 100 Hz for unaccelerated 2D graphics. During synchronization-intensive 3D rendering, we need a

lower latency guest to host notification. The guest can write to the *doorbell*, a register in I/O space, to explicitly ask the host to poll the command FIFO immediately.

5. EVALUATION

We conducted two categories of tests: application benchmarks, and microbenchmarks. All tests were conducted on the same physical machine: a 2nd generation Apple Mac Pro, with a total of eight 2.8 GHz Intel Xeon cores and an ATI Radeon HD2600 graphics card. All VMs used a single virtual CPU. With one exception, we found that the non-virtualized tests were unaffected by the number of CPU cores enabled.

5.1 Application Benchmarks

Application	Resolution	FPS
RTHDRIBL	1280 × 1024	22
RTHDRIBL	640 × 480	27.5
Half Life 2: Episode 2	1600 × 1200	22.2
Half Life 2: Episode 2	1024 × 768	32.2
Civilization 4	1600 × 1200	18
Max Payne 2	1600 × 1200	42

Table 2: Absolute frame rates with VMware Fusion 2.0. All applications run at interactive speeds (18–42 FPS).

The purpose of graphics acceleration hardware is to provide higher performance than would be possible using software alone. Therefore, in this section we will measure both the performance impact of virtualized graphics relative to non-virtualized GPU hardware, and the amount of performance improvement relative to TransGaming’s SwiftShader [19] software renderer, running in a VMware Fusion virtual machine.

In addition to VMware Fusion 2.0, which uses the architecture described above, we measured Parallels Desktop 3.0 where possible (three of our configurations do not run). To distinguish the effects that can be caused by API translation and by the host graphics stacks, we also ran our applications on VMware Workstation 6.5. These used our Direct3D rendering backend on the same hardware, but running Windows XP using Boot Camp.

It is quite challenging to measure the performance of graphics virtualization implementations accurately and fairly. The system under test has many variables, and they are often difficult or impossible to isolate. The virtualized operating system, host operating system, CPU virtualization overhead, GPU hardware, GPU drivers, and application under test may each have a profound effect on the overall system performance. Any one of these components may have opaque fast and slow paths—small differences in the application under test may cause wide gaps in performance, due to subtle and often hidden details of each component’s implementation. For example, each physical GPU driver may have different undocumented criteria for transferring vertex data at maximum speed.

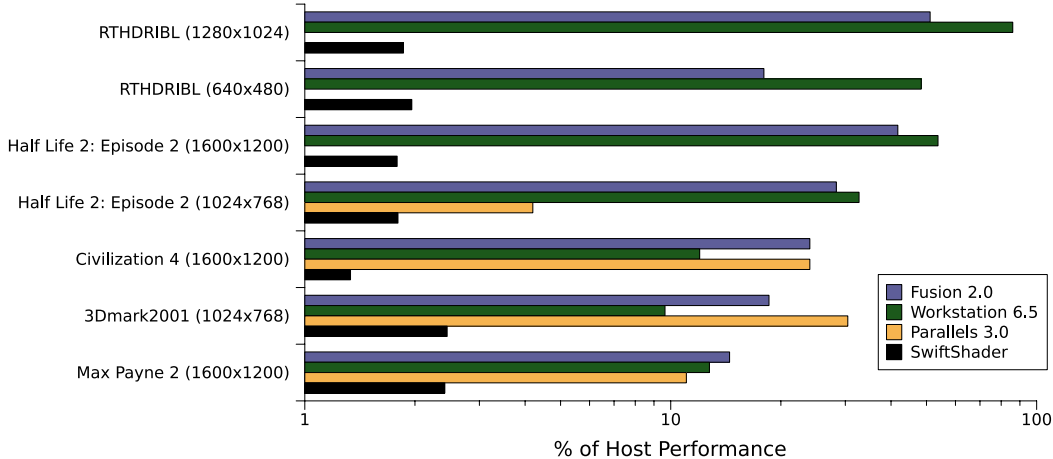


Figure 5: Relative performance of software rendering (SwiftShader) and three hardware accelerated virtualization techniques. The log scale highlights the huge gap between software and hardware acceleration versus the gap between virtualized and native hardware.

Additionally, the matrix of possible tests is limited by incompatible graphics APIs. Most applications and benchmarks are written for a single API, either OpenGL or Direct3D. Each available GPU virtualization implementation has a different level of API support. Parallels Desktop supports both OpenGL and Direct3D, VMware Fusion supports only Direct3D applications. VMGL and Chromium were not tested, as they only support OpenGL.

Figure 5 summarizes the application benchmark results. All three virtualization products performed substantially better than the fastest available software renderer, which obtained less than 3% of native performance in all tests. Applications which are mostly GPU limited, RTHDRIBL² [16] and *Half Life 2: Episode 2*, ran at closer to native speeds. Max Payne exhibits low performance relative to native, but that reflects the low ratio of GPU load to API calls. As an older game, it stresses the GPU less per drawing call, so the virtualization overhead occupies a higher proportion of the total execution time. In absolute terms, though, Max Payne has the highest frame rate of the applications we tested.

Table 2 reports the actual frame rates exhibited with these applications under VMware Fusion. While our virtualized 3D acceleration still lags native performance, we make two observations: it still achieves interactive frame rates, and it closes most of the gap between software rendering and native performance. For example, at 1600×1200 , VMware Fusion renders *Half-Life 2* at 22 frames per second, which is 23.35x faster than software rendering and only 2.4x slower than native.

5.2 Microbenchmarks

To better understand the nature of front-end virtualization’s performance impact, we performed a suite of microbench-

marks based on triangle rendering speed under various conditions. For all microbenchmarks, we rendered unlit untextured triangles using Vertex Shader 1.1 and the fixed-function pixel pipeline. This minimizes our dependency on shader translation and GPU driver implementation. The microbenchmark source code and results are available at [13].

Each test renders a fixed number of frames, each containing a variable number of draw calls with a variable length vertex buffer. For robustness against jitter or drift caused by timer virtualization, all tests measured elapsed time via a TCP/IP server running on an idle physical machine. Parameters for each test were chosen to optimize frame duration, so as to minimize the effects of noise from time quantization, network latency, and vertical sync latency.

Figures 6 through 9 show the results from our tests. All graphs show elapsed time on the Y axis, so lower is better.

The static vertex test, Figure 6, tests performance scalability when rendering vertex buffers which do not change contents. In Direct3D terms, this tests the managed buffer pool. Very little data must be exchanged between host and guest in this test, so an ideal front-end virtualization implementation would do quite well.

VMware Workstation manages to get just over 80% of the host’s performance in this test. In contrast, Parallels Desktop and VMware Fusion each get around 30%. In our experience, the wide gap between the OpenGL test results (both Parallels Desktop and VMware Fusion) as compared with the Direct3D results (VMware Workstation and host) are best explained by inefficiency in the Vertex Buffer Object implementation within Mac OS’s OpenGL stack. The difference in slope between Parallels Desktop and VMware Fusion could be explained by the two products producing GPU configurations (shaders, states, memory layouts) which have different performance characteristics.

²A *Real-time High Dynamic Range Image-Based Lighting* demo. This program uses complex shaders, floating point render targets, and several rendering passes.

The dynamic vertex test, Figure 7, switches from the managed buffer pool back to the default Direct3D buffer pool, and uploads new vertex data prior to each of the 100 draws per frame. It tests the driver stack’s ability to stream data to the GPU, and manage the re-use of buffer memory. Of the tests presented, this one is the most sensitive to differences in the resource management architecture and implementation of a virtual GPU.

The next test, Figure 8, is intended to test virtualization overhead while performing a GPU-intensive operation. While triangles in previous tests had zero pixel coverage, this tests renders triangles covering half the viewport. Ideally, this test would show nearly identical results for any front-end virtualization implementation. The actual results are relatively close, but on VMware’s platform there is a substantial amount of noise in the results. This appears to be due to the irregular completion of asynchronous commands when the physical GPU is under heavy load. Also worth noting is the fact that VMware Fusion, on average, performed better than the host machine. It’s possible that this test is exercising a particular drawing state which is more optimized in ATI’s Mac OS OpenGL driver than in their Windows Direct3D driver.

The final test, Figure 9, measures the overhead added to every separate draw. This was the only test where we saw variation in host performance based on the number of enabled CPU cores. This microbenchmark illustrates why the number of draw calls per frame is, in our experience, a relatively good predictor of overall application performance with front-end GPU virtualization.

VMware Workstation scored with performance very close to native, since its Direct3D backend performs fairly little work on each drawing operation. VMware Fusion’s translation from SVGA3D state to OpenGL state requires more time on each drawing operation, so this explains some of the performance gap. Parallels must perform a similar amount of work on each draw call, due to the Wine-based Direct3D to OpenGL API translation that it performs in the guest. It is also likely that some portion of the performance gap is due to the difference between ATI’s Mac OS OpenGL driver and their Windows Direct3D driver.

6. CONCLUSION

In VMware’s hosted architecture, we have implemented front-end GPU virtualization using a virtual device model with a high level rendering protocol. We have shown it to run modern graphics-intensive games and applications at interactive frame rates while preserving virtual machine interposition.

There is much future work in developing reliable benchmarks which specifically stress the performance weaknesses of a virtualization layer. Our tests show API overheads of about 2 to 120 times that of a native GPU. As a result, the performance of a virtualized GPU can be highly dependent on subtle implementation details of the application under test.

Back-end virtualization holds much promise for perfor-

mance, breadth of GPU feature support, and ease of driver maintenance. The current *fixed pass-through* implementations suffer from significant limitations, but *mediated pass-through* architectures can provide GPU multiplexing which preserves performance, fidelity, and some degree of interposition.

Front-end virtualization currently shows a substantial degradation in performance and GPU feature set relative to native hardware. Nevertheless, it is already enabling virtualized applications to run interactively that could never have been virtualized before, and is a foundation for virtualization of tomorrow’s GPU-intensive software. Even as back-end virtualization gains popularity, front-end virtualization can fill an important role for VMs which must be portable among diverse GPUs.

7. ACKNOWLEDGMENTS

Many people have contributed to the VMware SVGA and 3D code over the years. We would specifically like to thank Tony Cannon and Ramesh Dharan for their work on the foundations of our display emulation. Matt Ginzton has provided invaluable help in diverse areas of our graphics architecture over the years. Aaditya Chandrasekhar pioneered our shader translation architecture and continues to advance our Direct3D virtualization. Shelley Gong, Alex Corscadden, Mark Sheldon, and Stephen Johnson all actively contribute to our 3D emulation.

We would also like to thank Muli Ben-Yehuda, Alan Cox, and Scott Rixner for organizing the first USENIX Workshop on I/O Virtualization (WIOV 2008) where this paper was originally published.

8. REFERENCES

- [1] D. Abramson, J. Jackson, S. Muthrasanallur, G. Neiger, G. Regnier, R. Sankaran, I. Schoinas, R. Uhlig, B. Vembu, and J. Wiegert. Intel virtualization technology for directed I/O. *Intel Technology Journal*, 10, August 2006.
- [2] K. Adams and O. Agesen. A comparison of software and hardware techniques for x86 virtualization. *Operating Systems Review*, 40(5):2–13, Dec. 2006.
- [3] AMD developer guides and manuals. <http://developer.amd.com/documentation/guides/Pages/default.aspx>.
- [4] C. I. Dalton, D. Plaquin, W. Weidner, D. Kuhlmann, B. Balacheff, and R. Brown. Trusted virtual platforms: a key enabler for converged client devices. *SIGOPS Operating Systems Review*, 43(1):36–43, 2009.
- [5] K. Fatahalian and M. Houston. A closer look at GPUs. *Communications of the ACM*, 51(10):50–57, 2008.
- [6] Flash player penetration. http://www.adobe.com/products/player_census/flashplayer/.
- [7] J. G. Hansen. Blink: Advanced display multiplexing for virtualized applications. In *Proceedings of NOSSDAV 2007*, June 2007.
- [8] G. Humphreys, M. Houston, R. Ng, R. Frank, S. Ahern, P. Kirchner, and J. Klosowski. Chromium: A stream-processing framework for interactive rendering on clusters. In *SIGGRAPH 2002 Conference*

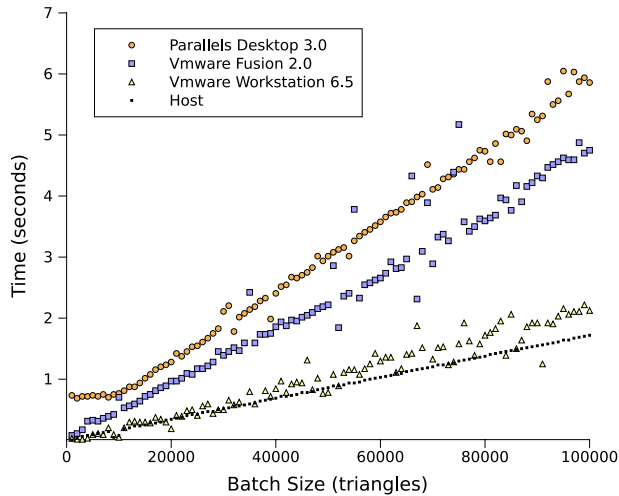


Figure 6: Static vertex rendering performance



Figure 7: Dynamic vertex rendering performance

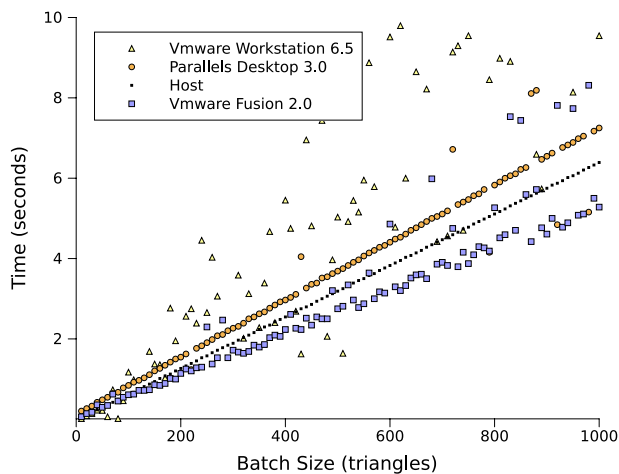


Figure 8: Filled triangle rendering performance

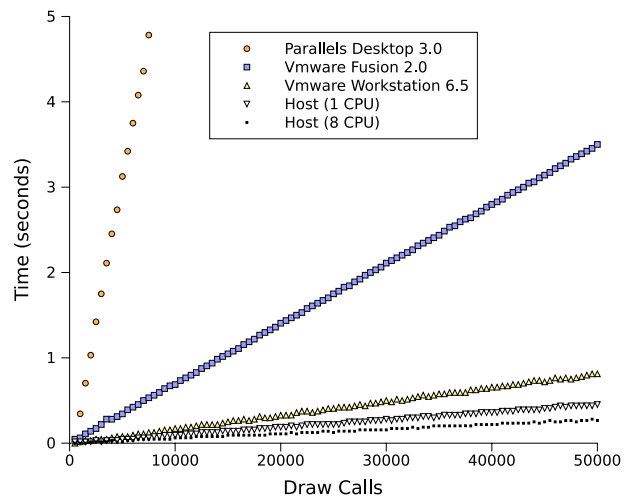


Figure 9: Draw call overhead

Proceedings, Annual Conference Series, pages 693–702. ACM Press/ACM SIGGRAPH, 2002.

- [9] IDF SF08: Parallels and Intel virtualization for directed I/O. http://www.youtube.com/watch?v=EiqMR5Wx_r4.
- [10] H. A. Lagar-Cavilla, N. Tolia, M. Satyanarayanan, and E. de Lara. VMM-independent graphics acceleration. In *Proceedings of the 3rd International Conference on Virtual Execution Environments, VEE 2007*, pages 33–43. ACM, June 2007.
- [11] A. Lefohn, M. Houston, C. Boyd, K. Fatahalian, T. Forsyth, D. Luebke, and J. Owens. Beyond programmable shading. In *SIGGRAPH '08: ACM SIGGRAPH 2008 classes*. <http://s08.idav.ucdavis.edu/>.
- [12] A. Mann. *Virtualization and Management: Trends, Forecasts, and Recommendations*. Enterprise

Management Associates, 2008.

- [13] Microbenchmark source code and results. <http://vmware-svga.sourceforge.net/wiov/>.
- [14] Parallels Desktop. <http://www.parallels.com/en/desktop/>.
- [15] Parallels on the Wine project wiki. <http://wiki.winehq.org/Parallels>.
- [16] Real-time high dynamic range image-based lighting demo. <http://www.daionet.gr.jp/~masa/rthdribl/>.
- [17] L. Seiler, D. Carmean, E. Sprangle, T. Forsyth, M. Abrash, P. Dubey, S. Junkins, A. Lake, J. Sugerman, R. Cavin, R. Espasa, E. Grochowski, T. Juan, and P. Hanrahan. Larrabee: a many-core x86 architecture for visual computing. *ACM Transactions on Graphics*, 27(3), 2008.
- [18] J. Sugerman, G. Venkitachalam, and B.-H. Lim. Virtualizing I/O devices on VMware workstation's

hosted virtual machine monitor. In *Proceedings of the USENIX Annual Technical Conference, General Track*, pages 1–14. USENIX, 2001.

- [19] SwiftShader. <http://www.transgaming.com/products/swiftshader/>.
- [20] VMware SVGA device developer kit. <http://vmware-svga.sourceforge.net/>.
- [21] VMware SVGA device interface and programming model. In X.org source repository, `xf86-video-vmware` driver README.
- [22] Wine project web site. <http://winehq.org/>.
- [23] H. Xu, S. Varley, and P. Thakkar. TA2644: Networking I/O virtualization. In *VMworld 2008*.