

KVM/ARM: The Design and Implementation of the Linux ARM Hypervisor

Christoffer Dall

Department of Computer Science
Columbia University
cdall@cs.columbia.edu

Jason Nieh

Department of Computer Science
Columbia University
nieh@cs.columbia.edu

Abstract

As ARM CPUs become increasingly common in mobile devices and servers, there is a growing demand for providing the benefits of virtualization for ARM-based devices. We present our experiences building the Linux ARM hypervisor, KVM/ARM, the first full system ARM virtualization **solution that can run unmodified guest operating systems on ARM multicore hardware. KVM/ARM introduces split-mode virtualization, allowing a hypervisor to split its execution across CPU modes and be integrated into the Linux kernel.** This allows KVM/ARM to leverage existing Linux hardware support and functionality to simplify hypervisor development and maintainability while utilizing recent ARM hardware virtualization extensions to run virtual machines with comparable performance to native execution. KVM/ARM has been successfully merged into the mainline Linux kernel, ensuring that it will gain wide adoption as the virtualization platform of choice for ARM. We provide the first measurements on real hardware of a complete hypervisor using ARM hardware virtualization support. **Our results demonstrate that KVM/ARM has modest virtualization performance and power costs, and can achieve lower performance and power costs compared to x86-based Linux virtualization on multicore hardware.**

Categories and Subject Descriptors C.0 [Computer Systems Organization]: General—Hardware/software interface, System architectures; D.4.7 [Operating Systems]: Organization and Design; D.4.8 [Operating Systems]: Performance

Keywords Virtualization, Hypervisors, Operating Systems, Multicore, ARM, Linux

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ASPLOS '14, March 1–4, 2014, Salt Lake City, Utah, USA.

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-2305-5/14/03...\$15.00.

<http://dx.doi.org/10.1145/2541940.2541946>

1. Introduction

ARM-based devices are seeing tremendous growth across smartphones, netbooks, and embedded computers. While ARM CPUs have benefited from their advantages in power efficiency in these markets, ARM CPUs also continue to increase in performance such that they are now within the range of x86 CPUs for many classes of applications. This is spurring the development of new ARM-based microservers and an upward push of ARM CPUs into traditional server, PC, and network systems.

To help bring the benefits of virtualization to ARM-based devices, ARM CPUs now include hardware support for virtualization. **Although virtualization is commonly used on x86-based systems, there are key differences between ARM and x86 virtualization. First, ARM and x86 virtualization extensions have important differences such that x86 hypervisor designs are not directly amenable to ARM.** These differences impact hypervisor performance and design, especially for multicore systems, but have not been evaluated with real hardware. **Second, unlike x86-based systems, there is no PC-standard hardware equivalent for ARM.** The ARM market is fragmented with many different vertically integrated ARM platforms with non-standard hardware. Virtualizing ARM in a manner that works across the diversity of ARM hardware in the absence of any real hardware standard is a key challenge.

We describe our experiences building KVM/ARM [13, 23], the ARM hypervisor in the mainline Linux kernel. **KVM/ARM is the first hypervisor to leverage ARM hardware virtualization support to run unmodified guest operating systems (OSes) on ARM multicore hardware.** Our work makes four main contributions. First, we **introduce split-mode virtualization**, a new approach to hypervisor design that splits the core hypervisor so that it runs across different privileged CPU modes to take advantage of the specific benefits and functionality offered by each CPU mode. This approach provides key benefits in the context of ARM virtualization. ARM introduced a new CPU mode for running hypervisors called Hyp mode, **but Hyp mode has its own set of features distinct from existing kernel modes. Hyp mode targets running a standalone hypervisor underneath the OS kernel, and was not designed to work well with a hosted hypervisor design**, where the hypervisor is integrated with a host

Why
ARM
and x86
need
separate
hypervis
or logic

One of the
key
contributions,
split mode

kernel. For example, standard OS mechanisms in Linux would have to be significantly redesigned to run in Hyp mode. Split-mode virtualization makes it possible to take advantage of the benefits of a hosted hypervisor design by running the hypervisor in normal privileged CPU modes to leverage existing OS mechanisms without modification while at the same time still using Hyp mode to leverage ARM hardware virtualization features.

Second, we designed and implemented KVM/ARM from the ground up as an open source project that would be easy to maintain and integrate into the Linux kernel. This is especially important in the context of ARM systems which lack standard ways to integrate hardware components, features for hardware discovery such as a standard BIOS or PCI bus, and standard mechanisms for installing low-level software. A standalone bare metal hypervisor would need to be ported to each and every supported hardware platform, a huge maintenance and development burden. Linux, however, is supported across almost all ARM platforms and by integrating KVM/ARM with Linux, KVM/ARM is automatically available on any device running a recent version of the Linux kernel. By using split-mode virtualization, we can leverage the existing Kernel-based Virtual Machine (KVM) [22] hypervisor interface in Linux and can reuse substantial pieces of existing kernel code and interfaces to reduce implementation complexity. KVM/ARM requires adding less than 6,000 lines of ARM code to Linux, a much smaller code base to maintain than standalone hypervisors. KVM/ARM was accepted as the ARM hypervisor of the mainline Linux kernel as of the Linux 3.9 kernel, ensuring its wide adoption and use given the dominance of Linux on ARM platforms. Based on our open source experiences, we offer some useful hints on transferring research ideas into implementations likely to be adopted by the open source community.

Third, we demonstrate the effectiveness of KVM/ARM on real multicore ARM hardware. Our results are the first measurements of a hypervisor using ARM virtualization support on real hardware. We compare against the standard widely-used Linux KVM x86 hypervisor and evaluate its performance overhead for running application workloads in virtual machines (VMs) versus native non-virtualized execution. Our results show that KVM/ARM achieves comparable performance overhead in most cases, and significantly lower performance overhead for two important applications, Apache and MySQL, on multicore platforms. These results provide the first comparison of ARM and x86 virtualization extensions on real hardware to quantitatively demonstrate how the different design choices affect virtualization performance. We show that KVM/ARM also provides power efficiency benefits over Linux KVM x86.

Finally, we make several recommendations regarding future hardware support for virtualization based on our experiences building and evaluating a complete ARM hypervisor. We identify features that are important and helpful to reduce the software complexity of hypervisor implementations, and discuss mechanisms useful to maximize hypervisor performance, especially in the context of multicore systems.

This paper describes the design and implementation of KVM/ARM. Section 2 presents an overview of the ARM virtualization extensions and a comparison with x86. Section 3 describes the design of the KVM/ARM hypervisor. Section 4 discusses the implementation of KVM/ARM and our experiences releasing it to the Linux community and having it adopted into the mainline Linux kernel. Section 5 presents experimental results quantifying the performance and energy efficiency of KVM/ARM, as well as a quantitative comparison of real ARM and x86 virtualization hardware. Section 6 makes recommendations for designing future hardware virtualization support. Section 7 discusses related work. Finally, we present some concluding remarks.

what kinda support is required
for virtualization?

2. ARM Virtualization Extensions

Because the ARM architecture is not classically virtualizable [27], ARM introduced hardware virtualization support as an optional extension in the latest ARMv7 [6] and ARMv8 [7] architectures. For example, the Cortex-A15 [4] is a current ARMv7 CPU with hardware virtualization support. We present a brief overview of the ARM virtualization extensions.

CPU Virtualization Figure 1 shows the CPU modes on the ARMv7 architecture, including TrustZone (Security Extensions) and a new CPU mode called Hyp mode. TrustZone splits the modes into two worlds, secure and non-secure, which are orthogonal to the CPU modes. A special mode, monitor mode, is provided to switch between the secure and non-secure worlds. Although ARM CPUs always power up starting in the secure world, ARM bootloaders typically transition to the non-secure world at an early stage. The secure world is only used for specialized use cases such as digital rights management. TrustZone may appear useful for virtualization by using the secure world for hypervisor execution, but this does not work because trap-and-emulate is not supported. There is no means to trap operations executed in the non-secure world to the secure world. Non-secure software can therefore freely configure, for example, virtual memory. Any software running at the highest non-secure privilege level therefore has access to all non-secure physical memory, making it impossible to isolate multiple VMs running in the non-secure world.

TrustZone??

Hyp mode was introduced as a trap-and-emulate mechanism to support virtualization in the non-secure world. Hyp mode is a

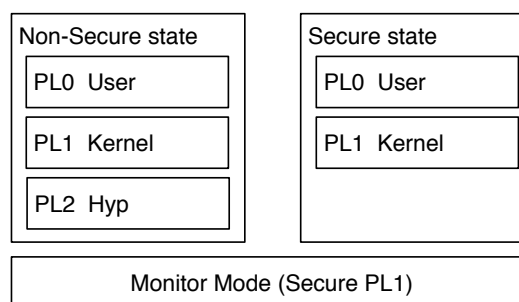


Figure 1: ARMv7 Processor Modes

CPU mode that is strictly more privileged than other CPU modes, user and kernel modes. Software running in Hyp mode can configure the hardware to trap from kernel mode into Hyp mode on various sensitive instructions and hardware interrupts. To run VMs, the hypervisor must at least partially reside in Hyp mode. The VM will execute normally in user and kernel mode until some condition is reached that requires intervention of the hypervisor. At this point, the hardware traps into Hyp mode giving control to the hypervisor, which can then manage the hardware and provide the required isolation across VMs. Once the condition is processed by the hypervisor, the CPU can be switched back into user or kernel mode and the VM can continue executing.

The ARM architecture allows each trap to be configured to trap directly into a VM's kernel mode instead of going through Hyp mode. For example, traps caused by system calls or page faults from user mode can be configured to trap to a VM's kernel mode directly so that they are handled by the guest OS without intervention of the hypervisor. This avoids going to Hyp mode on each system call or page fault, reducing virtualization overhead. Additionally, all traps into Hyp mode can be disabled and a single non-virtualized kernel can run in kernel mode and have complete control of the system.

ARM designed the virtualization support around a separate CPU mode distinct from existing kernel mode, because they envisioned a standalone hypervisor underneath a more complex rich OS kernel [14]. They wanted to make it simpler for hypervisor developers to implement the hypervisors, and therefore reduced the number of control registers available in Hyp mode compared to kernel mode. Similarly, they mandated certain bits to be set in the page table entries, because they did not envision a hypervisor sharing page tables with software running in user space, which is for example what the Linux kernel does with kernel mode.

Memory Virtualization ARM also provides hardware support to virtualize physical memory. When running a VM, the physical addresses managed by the VM are actually *Intermediate Physical Addresses* (IPAs), also known as *guest physical addresses*, and need to be translated into *physical addresses* (PAs), also known as *host physical addresses*. Similarly to nested page tables on x86, ARM provides a second set of page tables, Stage-2 page tables, which translate from IPAs to PAs corresponding to guest and host physical addresses, respectively. Stage-2 translation can be completely disabled and enabled from Hyp mode. Stage-2 page tables use ARM's new LPAE page table format, with subtle differences from the page tables used by kernel mode.

Interrupt Virtualization ARM defines the Generic Interrupt Controller (GIC) architecture [5]. The GIC routes interrupts from devices to CPUs and CPUs query the GIC to discover the source of an interrupt. The GIC is especially important in multicore configurations, because it is used to generate Inter-Processor Interrupts (IPIs) from one CPU core to another. The GIC is split in two parts, the distributor and the CPU interfaces. There is only one distributor in a system, but each CPU core has a GIC CPU interface. Both the CPU interfaces and the distributor are accessed over a Memory-Mapped interface (MMIO). The dis-

tributor is used to configure the GIC, for example, to configure the CPU core affinity of an interrupt, to completely enable or disable interrupts on a system, or to send an IPI to another CPU core. The CPU interface is used to acknowledge (ACK) and to signal End-Of-Interrupt (EOI). For example, when a CPU core receives an interrupt, it will read a special register on the GIC CPU interface, which ACKs the interrupt and returns the number of the interrupt. The interrupt will not be raised to the CPU again before the CPU writes to the EOI register of the CPU interface with the value retrieved from the ACK register.

Interrupts can be configured to trap to either Hyp or kernel mode. Trapping all interrupts to kernel mode and letting OS software running in kernel mode handle them directly is efficient, but does not work in the context of VMs, because the hypervisor loses control over the hardware. Trapping all interrupts to Hyp mode ensures that the hypervisor retains control, but requires emulating virtual interrupts in software to signal events to VMs. This is cumbersome to manage and expensive because each step of interrupt and virtual interrupt processing, such as ACKing and EOling, must go through the hypervisor.

The GIC v2.0 includes hardware virtualization support in the form of a virtual GIC (VGIC) so that receiving virtual interrupts does not need to be emulated in software by the hypervisor. The VGIC introduces a VGIC CPU interface for each CPU and a corresponding hypervisor control interface for each CPU. VMs are configured to see the VGIC CPU interface instead of the GIC CPU interface. Virtual interrupts are generated by writing to special registers, the *list registers*, in the VGIC hypervisor control interface, and the VGIC CPU interface raises the virtual interrupts directly to a VM's kernel mode. Because the VGIC CPU interface includes support for ACK and EOI, these operations no longer need to trap to the hypervisor to be emulated in software, reducing overhead for receiving interrupts on a CPU. For example, emulated virtual devices typically raise virtual interrupts through a software API to the hypervisor, which can leverage the VGIC by writing the virtual interrupt number for the emulated device into the list registers. This causes the VGIC to interrupt the VM directly to kernel mode and lets the guest OS ACK and EOI the virtual interrupt without trapping to the hypervisor. Note that the distributor must still be emulated in software and all accesses to the distributor by a VM must still trap to the hypervisor. For example, when a virtual CPU sends a virtual IPI to another virtual CPU, this will cause a trap to the hypervisor, which emulates the distributor access in software and programs the list registers on the receiving CPU's GIC hypervisor control interface.

Timer Virtualization ARM defines the Generic Timer Architecture which includes support for timer virtualization. Generic timers provide a counter that measures passing of time in real-time, and a timer for each CPU, which is programmed to raise an interrupt to the CPU after a certain amount of time has passed. Timers are likely to be used by both hypervisors and guest OSes, but to provide isolation and retain control, the timers used by the hypervisor cannot be directly configured and manipulated by guest OSes. Such timer accesses from a guest OS would need to

Important point on why trapping to Hyp mode is not efficient

How does the emulation of interrupts look like?

Does VGIC also have distributor and CPU interface...if so how do these interrupts get handled by the hypervisor?

trap to Hyp mode, incurring additional overhead for a relatively frequent operation for some workloads. Hypervisors may also wish to virtualize VM time, which is problematic if VMs have direct access to counter hardware.

ARM provides virtualization support for the timers by introducing a new counter, the *virtual counter* and a new timer, the *virtual timer*. A hypervisor can be configured to use physical timers while VMs are configured to use virtual timers. VMs can then access, program, and cancel virtual timers without causing traps to Hyp mode. Access to the physical timer and counter from kernel mode is controlled from Hyp mode, but software running in kernel mode always has access to the virtual timers and counters.

Comparison with x86 There are a number of similarities and differences between the ARM virtualization extensions and hardware virtualization support for x86 from Intel and AMD. Intel and AMD extensions are very similar, so we limit our comparison to ARM and Intel. ARM supports virtualization through a separate CPU mode, Hyp mode, which is a separate and strictly more privileged CPU mode than previous user and kernel modes. In contrast, Intel has root and non-root mode [20], which are orthogonal to the CPU protection modes. While sensitive operations on ARM trap to Hyp mode, sensitive operations can trap from non-root mode to root mode while staying in the same protection level on Intel. A crucial difference between the two hardware designs is that Intel's root mode supports the same full range of user and kernel mode functionality as its non-root mode, whereas ARM's Hyp mode is a strictly different CPU mode with its own set of features. A hypervisor using ARM's Hyp mode has an arguably simpler set of features to use than the more complex options available with Intel's root mode.

Both ARM and Intel trap into their respective Hyp and root modes, but Intel provides specific hardware support for a VM control block which is automatically saved and restored when switching to and from root mode using only a single instruction. This is used to automatically save and restore guest state when switching between guest and hypervisor execution. In contrast, ARM provides no such hardware support and any state that needs to be saved and restored must be done explicitly in software. This provides some flexibility in what is saved and restored in switching to and from Hyp mode. For example, trapping to ARM's Hyp mode is potentially faster than trapping to Intel's root mode if there is no additional state to save.

ARM and Intel are quite similar in their support for virtualizing physical memory. Both introduce an additional set of page tables to translate guest to host physical addresses. ARM benefited from hindsight in including Stage-2 translation whereas Intel did not include its equivalent Extended Page Table (EPT) support until its second generation virtualization hardware.

ARM's support for virtual timers have no real x86 counterpart, and until the recent introduction of Intel's virtual APIC support [20], ARM's support for virtual interrupts also had no x86 counterpart. Without virtual APIC support, EOLing interrupts in an x86 VM requires traps to root mode, whereas ARM's virtual GIC avoids the cost of trapping to Hyp mode for those interrupt

handling mechanisms. Executing similar timer functionality by a guest OS on x86 will incur additional traps to root mode compared to the number of traps to Hyp mode required for ARM. Reading a counter, however, is not a privileged operation on x86 and does not trap, even without virtualization support in the counter hardware.

3. Hypervisor Architecture

Instead of reinventing and reimplementing complex core functionality in the hypervisor, and potentially introducing tricky and fatal bugs along the way, KVM/ARM builds on KVM and leverages existing infrastructure in the Linux kernel. While a standalone bare metal hypervisor design approach has the potential for better performance and a smaller Trusted Computing Base (TCB), this approach is less practical on ARM. ARM hardware is in many ways much more diverse than x86. Hardware components are often tightly integrated in ARM devices in non-standard ways by different device manufacturers. ARM hardware lacks features for hardware discovery such as a standard BIOS or a PCI bus, and there is no established mechanism for installing low-level software on a wide variety of ARM platforms. Linux, however, is supported across almost all ARM platforms and by integrating KVM/ARM with Linux, KVM/ARM is automatically available on any device running a recent version of the Linux kernel. This is in contrast to bare metal approaches such as Xen [32], which must actively support every platform on which they wish to install the Xen hypervisor. For example, for every new SoC that Xen needs to support, the developers must implement a new serial device driver in the core Xen hypervisor.

While KVM/ARM benefits from its integration with Linux in terms of portability and hardware support, a key problem we had to address was that the ARM hardware virtualization extensions were designed to support a standalone hypervisor design where the hypervisor is completely separate from any standard kernel functionality, as discussed in Section 2. In the following, we describe how KVM/ARM's novel design makes it possible to benefit from integration with an existing kernel and at the same time take advantage of the hardware virtualization features.

3.1 Split-mode Virtualization

Simply running a hypervisor entirely in ARM's Hyp mode is attractive since it is the most privileged level. However, since KVM/ARM leverages existing kernel infrastructure such as the scheduler, running KVM/ARM in Hyp mode implies running the Linux kernel in Hyp mode. This is problematic for at least two reasons. First, low-level architecture dependent code in Linux is written to work in kernel mode, and would not run unmodified in Hyp mode, because Hyp mode is a completely different CPU mode from normal kernel mode. The significant changes that would be required to run the kernel in Hyp mode would be very unlikely to be accepted by the Linux kernel community. More importantly, to preserve compatibility with hardware without Hyp mode and to run Linux as a guest OS, low-level code would have to be written to work in both modes, potentially resulting

in slow and convoluted code paths. As a simple example, a page fault handler needs to obtain the virtual address causing the page fault. In Hyp mode this address is stored in a different register than in kernel mode.

Second, running the entire kernel in Hyp mode would adversely affect native performance. For example, Hyp mode has its own separate address space. Whereas kernel mode uses two page table base registers to provide the familiar 3GB/1GB split between user address space and kernel address space, Hyp mode uses a single page table register and therefore cannot have direct access to the user space portion of the address space. Frequently used functions to access user memory would require the kernel to explicitly map user space data into kernel address space and subsequently perform necessary teardown and TLB maintenance operations, resulting in poor native performance on ARM.

These problems with running a Linux hypervisor using ARM Hyp mode do not occur for x86 hardware virtualization. x86 root mode is orthogonal to its CPU privilege modes. The entire Linux kernel can run in root mode as a hypervisor because the same set of CPU modes available in non-root mode are available in root mode. Nevertheless, given the widespread use of ARM and the advantages of Linux on ARM, finding an efficient virtualization solution for ARM that can leverage Linux and take advantage of the hardware virtualization support is of crucial importance.

KVM/ARM introduces split-mode virtualization, a new approach to hypervisor design that splits the core hypervisor so that it runs across different privileged CPU modes to take advantage of the specific benefits and functionality offered by each CPU mode. KVM/ARM uses split-mode virtualization to leverage the ARM hardware virtualization support enabled by Hyp mode, while at the same time leveraging existing Linux kernel services running in kernel mode. Split-mode virtualization allows KVM/ARM to be integrated with the Linux kernel without intrusive modifications to the existing code base.

This is done by splitting the hypervisor into two components, the lowvisor and the highvisor, as shown in Figure 2. The lowvisor is designed to take advantage of the hardware virtualization support available in Hyp mode to provide three key functions. First, the lowvisor sets up the correct execution context by appropriate configuration of the hardware, and enforces protection and isolation between different execution contexts. The lowvisor directly interacts with hardware protection features and is therefore highly critical and the code base is kept to an absolute minimum. Second, the lowvisor switches from a VM execution context to the host execution context and vice-versa. The host execution context is used to run the hypervisor and the host Linux kernel. We refer to an execution context as a world, and switching from one world to another as a world switch, because the entire state of the system is changed. Since the lowvisor is the only component that runs in Hyp mode, only it can be responsible for the hardware reconfiguration necessary to perform a world switch. Third, the lowvisor provides a virtualization trap handler, which handles interrupts and exceptions that must trap to the hypervisor. The lowvisor performs only the minimal amount of

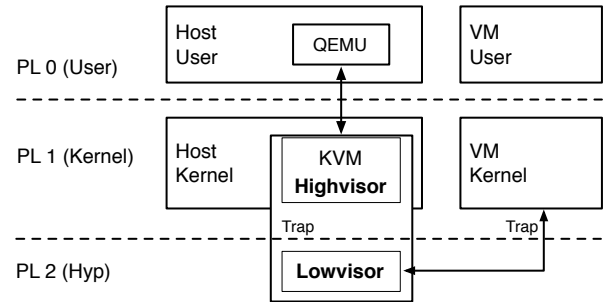


Figure 2: KVM/ARM System Architecture

processing required and defers the bulk of the work to be done to the highvisor after a world switch to the highvisor is complete.

The highvisor runs in kernel mode as part of the host Linux kernel. It can therefore directly leverage existing Linux functionality such as the scheduler, and can make use of standard kernel software data structures and mechanisms to implement its functionality, such as locking mechanisms and memory allocation functions. This makes higher-level functionality easier to implement in the highvisor. For example, while the lowvisor provides a low-level trap-handler and the low-level mechanism to switch from one world to another, the highvisor handles Stage-2 page faults from the VM and performs instruction emulation. Note that parts of the VM run in kernel mode, just like the highvisor, but with Stage-2 translation and trapping to Hyp mode enabled.

Because the hypervisor is split across kernel mode and Hyp mode, switching between a VM and the highvisor involves multiple mode transitions. A trap to the highvisor while running the VM will first trap to the lowvisor running in Hyp mode. The lowvisor will then cause another trap to run the highvisor. Similarly, going from the highvisor to a VM requires trapping from kernel mode to Hyp mode, and then switching to the VM. As a result, split-mode virtualization incurs a double trap cost in switching to and from the highvisor. On ARM, the only way to perform these mode transitions to and from Hyp mode is by trapping. However, as shown in Section 5, this extra trap is not a significant performance cost on ARM.

KVM/ARM uses a memory mapped interface to share data between the highvisor and lowvisor as necessary. Because memory management can be complex, we leverage the highvisor's ability to use the existing memory management subsystem in Linux to manage memory for both the highvisor and lowvisor. Managing the lowvisor's memory involves additional challenges though, because it requires managing Hyp mode's separate address space. One simplistic approach would be to reuse the host kernel's page tables and also use them in Hyp mode to make the address spaces identical. This unfortunately does not work, because Hyp mode uses a different page table format from kernel mode. Therefore, the highvisor explicitly manages the Hyp mode page tables to map any code executed in Hyp mode and any data structures shared between the highvisor and the lowvisor to the same virtual addresses in Hyp mode and in kernel mode.

Action	Nr.	State
Context Switch	38	General Purpose (GP) Registers
	26	Control Registers
	16	VGIC Control Registers
	4	VGIC List Registers
	2	Arch. Timer Control Registers
	32	64-bit VFP registers
Trap-and-Emulate	4	32-bit VFP Control Registers
	-	CP14 Trace Registers
	-	WFI Instructions
	-	SMC Instructions
	-	ACTLR Access
	-	Cache ops. by Set/Way
	-	L2CTLR / L2ECTLR Registers

Table 1: VM and Host State on a Cortex-A15

3.2 CPU Virtualization

To virtualize the CPU, KVM/ARM must present an interface to the VM which is essentially identical to the underlying real hardware CPU, while ensuring that the hypervisor remains in control of the hardware. This involves ensuring that software running in the VM must have persistent access to the same register state as software running on the physical CPU, as well as ensuring that physical hardware state associated with the hypervisor and its host kernel is persistent across running VMs. Register state not affecting VM isolation can simply be context switched by saving the VM state and restoring the host state from memory when switching from a VM to the host and vice versa. KVM/ARM configures access to all other sensitive state to trap to Hyp mode, so it can be emulated by the hypervisor.

Table 1 shows the CPU register state visible to software running in kernel and user mode, and KVM/ARM's virtualization method for each register group. The lowvisor has its own dedicated configuration registers only for use in Hyp mode, and is not shown in Table 1. KVM/ARM context switches registers during world-switches whenever the hardware supports it, because it allows the VM direct access to the hardware. For example, the VM can directly program the Stage-1 page table base register without trapping to the hypervisor, a fairly common operation in most guest OSes. KVM/ARM performs trap and emulate on sensitive instructions and when accessing hardware state that could affect the hypervisor or would leak information about the hardware to the VM that violates its virtualized abstraction. For example, KVM/ARM traps if a VM executes the WFI instruction, which causes the CPU to power down, because such an operation should only be performed by the hypervisor to maintain control of the hardware. KVM/ARM defers switching certain register state until absolutely necessary, which slightly improves performance under certain workloads.

The difference between running inside a VM in kernel or user mode and running the hypervisor in kernel or user mode is determined by how the virtualization extensions have been configured by Hyp mode during the world switch. A world switch from the host to a VM performs the following actions: (1) store all host GP

registers on the Hyp stack, (2) configure the VGIC for the VM, (3) configure the timers for the VM, (4) save all host-specific configuration registers onto the Hyp stack, (5) load the VM's configuration registers onto the hardware, which can be done without affecting current execution, because Hyp mode uses its own configuration registers, separate from the host state, (6) configure Hyp mode to trap floating-point operations for lazy context switching, trap interrupts, trap CPU halt instructions (WFI/WFE), trap SMC instructions, trap specific configuration register accesses, and trap debug register accesses, (7) write VM-specific IDs into shadow ID registers, (8) set the Stage-2 page table base register (VTTBR) and enable Stage-2 address translation, (9) restore all guest GP registers, and (10) trap into either user or kernel mode.

The CPU will stay in the VM world until an event occurs, which triggers a trap into Hyp mode. Such an event can be caused by any of the traps mentioned above, a Stage-2 page fault, or a hardware interrupt. Since the event requires services from the highvisor, either to emulate the expected hardware behavior for the VM or to service a device interrupt, KVM/ARM must perform another world switch back into the highvisor and its host. The world switch back to the host from a VM performs the following actions: (1) store all VM GP registers, (2) disable Stage-2 translation, (3) configure Hyp mode to not trap any register access or instructions, (4) save all VM-specific configuration registers, (5) load the host's configuration registers onto the hardware, (6) configure the timers for the host, (7) save VM-specific VGIC state, (8) restore all host GP registers, and (9) trap into kernel mode.

3.3 Memory Virtualization

KVM/ARM provides memory virtualization by enabling Stage-2 translation for all memory accesses when running in a VM. Stage-2 translation can only be configured in Hyp mode, and its use is completely transparent to the VM. The highvisor manages the Stage-2 translation page tables to only allow access to memory specifically allocated for a VM; other accesses will cause Stage-2 page faults which trap to the hypervisor. This mechanism ensures that a VM cannot access memory belonging to the hypervisor or other VMs, including any sensitive data. Stage-2 translation is disabled when running in the highvisor and lowvisor because the highvisor has full control of the complete system and directly manages the host physical addresses. When the hypervisor performs a world switch to a VM, it enables Stage-2 translation and configures the Stage-2 page table base register accordingly. Although both the highvisor and VMs share the same CPU modes, Stage-2 translations ensure that the highvisor is protected from any access by the VMs.

KVM/ARM uses split-mode virtualization to leverage existing kernel memory allocation, page reference counting, and page table manipulation code. KVM/ARM handles Stage-2 page faults by considering the IPA of the fault, and if that address belongs to normal memory in the VM memory map, KVM/ARM allocates a page for the VM by simply calling an existing kernel function, such as `get_user_pages`, and maps the allocated page to the VM in the Stage-2 page tables. In comparison, a bare metal

Ask about what each step is for ?

What happens when switch from VM to Hyp mode

How ??

hypervisor would be forced to either statically allocate memory to VMs or write an entire new memory allocation subsystem.

3.4 I/O Virtualization

KVM/ARM leverages existing QEMU and Virtio [29] user space device emulation to provide I/O virtualization. At a hardware level, all I/O mechanisms on the ARM architecture are based on load/store operations to MMIO device regions. With the exception of devices directly assigned to VMs, all hardware MMIO regions are inaccessible from VMs. KVM/ARM uses Stage-2 translations to ensure that physical devices cannot be accessed directly from VMs. Any access outside of RAM regions allocated for the VM will trap to the hypervisor, which can route the access to a specific emulated device in QEMU based on the fault address. This is somewhat different from x86, which uses x86-specific hardware instructions such as `inl` and `outl` for port I/O operations in addition to MMIO. As we show in Section 5, KVM/ARM achieves low I/O performance overhead with very little implementation effort.

3.5 Interrupt Virtualization

KVM/ARM leverages its tight integration with Linux to reuse existing device drivers and related functionality, including handling interrupts. When running in a VM, KVM/ARM configures the CPU to trap all hardware interrupts to Hyp mode. On each interrupt, it performs a world switch to the highvisor and the host handles the interrupt, so that the hypervisor remains in complete control of hardware resources. When running in the host and the highvisor, interrupts trap directly to kernel mode, avoiding the overhead of going through Hyp mode. In both cases, all hardware interrupt processing is done in the host by reusing Linux's existing interrupt handling functionality.

However, VMs must receive notifications in the form of virtual interrupts from emulated devices and multicore guest OSes must be able to send virtual IPIs from one virtual core to another. KVM/ARM uses the VGIC to inject virtual interrupts into VMs to reduce the number of traps to Hyp mode. As described in Section 2, virtual interrupts are raised to virtual CPUs by programming the list registers in the VGIC hypervisor CPU control interface. KVM/ARM configures the Stage-2 page tables to prevent VMs from accessing the control interface and to allow access only to the VGIC virtual CPU interface, ensuring that only the hypervisor can program the control interface and that the VM can access the VGIC virtual CPU interface directly. However, guest OSes will still attempt to access a GIC distributor to configure the GIC and to send IPIs from one virtual core to another. Such accesses will trap to the hypervisor and the hypervisor must emulate the distributor.

KVM/ARM introduces the virtual distributor, a software model of the GIC distributor as part of the highvisor. The virtual distributor exposes an interface to user space, so emulated devices in user space can raise virtual interrupts to the virtual distributor and exposes an MMIO interface to the VM identical to that of the physical GIC distributor. The virtual distributor keeps internal software state about the state of each interrupt and

uses this state whenever a VM is scheduled, to program the list registers to inject virtual interrupts. For example, if virtual CPU0 sends an IPI to virtual CPU1, the distributor will program the list registers for virtual CPU1 to raise a virtual IPI interrupt the next time virtual CPU1 runs.

Ideally, the virtual distributor only accesses the hardware list registers when necessary, since device MMIO operations are typically significantly slower than cached memory accesses. A complete context switch of the list registers is required when scheduling a different VM to run on a physical core, but not necessarily required when simply switching between a VM and the hypervisor. For example, if there are no pending virtual interrupts, it is not necessary to access any of the list registers. Note that once the hypervisor writes a virtual interrupt to a list register when switching to a VM, it must also read the list register back when switching back to the hypervisor, because the list register describes the state of the virtual interrupt and indicates, for example, if the VM has ACKed the virtual interrupt. The initial unoptimized version of KVM/ARM uses a simplified approach which completely context switches all VGIC state including the list registers on each world switch.

3.6 Timer Virtualization

Reading counters and programming timers are frequent operations in many OSes for process scheduling and to regularly poll device state. For example, Linux reads a counter to determine if a process has expired its time slice, and programs timers to ensure that processes don't exceed their allowed time slices. Application workloads also often leverage timers for various reasons. Trapping to the hypervisor for each such operation is likely to incur noticeable performance overheads, and allowing a VM direct access to the time-keeping hardware typically implies giving up timing control of the hardware resources as VMs can disable timers and control the CPU for extended periods of time.

KVM/ARM leverages ARM's hardware virtualization features of the generic timers to allow VMs direct access to reading counters and programming timers without trapping to Hyp mode while at the same time ensuring the hypervisor remains in control of the hardware. Since access to the physical timers is controlled using Hyp mode, any software controlling Hyp mode has access to the physical timers. KVM/ARM maintains hardware control by using the physical timers in the hypervisor and disallowing access to physical timers from the VM. The Linux kernel running as a guest OS only accesses the virtual timer and can therefore directly access timer hardware without trapping to the hypervisor.

Unfortunately, due to architectural limitations, the virtual timers cannot directly raise virtual interrupts, but always raise hardware interrupts, which trap to the hypervisor. KVM/ARM detects when a VM virtual timer expires, and injects a corresponding virtual interrupt to the VM, performing all hardware ACK and EOI operations in the highvisor. The hardware only provides a single virtual timer per physical CPU, and multiple virtual CPUs may be multiplexed across this single hardware instance. To support virtual timers in this scenario, KVM/ARM detects unexpired timers when a VM traps to the hypervisor

Important

so to call a virtual timer interrupt, the physical timer interrupts to run KVM/ARM which checks for virtual timer and then sends a virtual interrupt for guest OS to handle ?

and leverages existing OS functionality to program a software timer at the time when the virtual timer would have otherwise fired, had the VM been left running. When such a software timer fires, a callback function is executed, which raises a virtual timer interrupt to the VM using the virtual distributor described above.

4. Implementation and Adoption

We have successfully integrated our work into the Linux kernel and KVM/ARM is now the standard ARM hypervisor on Linux platforms, as it is included in every kernel beginning with version 3.9. Its relative simplicity and rapid completion was facilitated by specific design choices that allow it to leverage substantial existing infrastructure despite differences in the underlying hardware. We share some lessons we learned from our experiences in hopes that they may be helpful to others in getting research ideas widely adopted by the open source community.

Code maintainability is key. It is a common misconception that a research software implementation providing potential improvements or interesting new features can simply be open sourced and thereby quickly integrated by the open source community. An important point that is often not taken into account is that any implementation must be maintained. If an implementation requires many people and much effort to be maintained, it is much less likely to be integrated into existing open source code bases. Because maintainability is so crucial, reusing code and interfaces is important. For example, KVM/ARM builds on existing infrastructure such as KVM and QEMU, and from the very start we prioritized addressing code review comments to make our code suitable for integration into existing systems. An unexpected but important benefit of this decision was that we could leverage the community for help to solve hard bugs or understand intricate parts of the ARM architecture.

Be a known contributor. Convincing maintainers to integrate code is not just about the code itself, but also about who submits it. It is not unusual for researchers to complain about kernel maintainers not accepting their code into Linux only to have some known kernel developer submit the same idea and have it accepted. The reason is an issue of trust. Establishing trust is a catch-22: one must be well-known to submit code, yet one cannot become known without submitting code. One way to do this is to start small. As part of our work, we also made various small changes to KVM to prepare support for ARM, which included cleaning up existing code to be more generic and improving cross platform support. The KVM maintainers were glad to accept these small improvements, which generated goodwill and helped us become known to the KVM community.

Make friends and involve the community. Open source development turns out to be quite a social enterprise. Networking with the community helps tremendously, not just online, but in person at conferences and other venues. For example, at an early stage in the development of KVM/ARM, we traveled to ARM headquarters in Cambridge, UK to establish contact with both

ARM management and the ARM kernel engineering team, who both contributed to our efforts.

As another example, an important issue in integrating KVM/ARM into the kernel was agreeing on various interfaces for ARM virtualization, such as reading and writing control registers. Since it is an established policy to never break released interfaces and compatibility with user space applications, existing interfaces cannot be changed, and the community puts great effort into designing extensible and reusable interfaces. Deciding on the appropriateness of an interface is a judgment call and not an exact science. We were fortunate enough to receive help from well-known kernel developers such as Rusty Russell, who helped us drive both the implementation and communication about our interfaces, specifically for user space save and restore of registers, a feature useful for both debugging and VM migration. Working with an established developer like Rusty was a tremendous help because we benefited from both his experience and strong voice in the kernel community.

Involve the community early. An important issue in developing KVM/ARM was how to get access to Hyp mode across the plethora of available ARM SoC platforms supported by Linux. One approach would be to initialize and configure Hyp mode when KVM is initialized, which would isolate the code changes to the KVM subsystem. However, because getting into Hyp mode from the kernel involves a trap, early stage bootloader must have already installed code in Hyp mode to handle the trap and allow KVM to run. If no such trap handler was installed, trapping to Hyp mode could end up crashing the kernel. We worked with the kernel community to define the right ABI between KVM and the bootloader, but soon learned that agreeing on ABIs with SoC vendors had historically been difficult.

In collaboration with ARM and the open source community, we reached the conclusion that if we simply required the kernel to be booted in Hyp mode, we would not have to rely on fragile ABIs. The kernel then simply tests when it starts up whether it is in Hyp mode, in which case it installs a trap handler to provide a hook to re-enter Hyp mode at a later stage. A small amount of code must be added to the kernel boot procedure, but the result is a much cleaner and robust mechanism. If the bootloader is Hyp mode unaware and the kernel does not boot up in Hyp mode, KVM/ARM will detect this and will simply remain disabled. This solution avoids the need to design a new ABI and it turned out that legacy kernels would still work, because they always make an explicit switch into kernel mode as their first instruction. These changes were merged into the mainline Linux 3.6 kernel, and official ARM kernel boot recommendations were modified to recommend that all bootloaders boot the kernel in Hyp mode to take advantage of the new architecture features.

Know the chain of command. There were multiple possible upstream paths for KVM/ARM. Historically, other architectures supported by KVM such as x86 and PowerPC were merged through the KVM tree directly into Linus Torvalds' tree with the appropriate approval of the respective architecture maintainers. KVM/ARM, however, required a few minor changes to

ARM-specific header files and the idmap subsystem, and it was therefore not clear whether the code would be integrated via the KVM tree with approval from the ARM kernel maintainer or via the ARM kernel tree. Russell King is the ARM kernel maintainer, and Linus pulls directly from his ARM kernel tree for ARM-related code. The situation was particularly interesting, because Russell King did not want to merge virtualization support in the mainline kernel [24] and he did not review our code. At the same time, the KVM community was quite interested in integrating our code, but could not do so without approval from the ARM maintainers, and Russell King refused to engage in a discussion about this procedure.

Be persistent. While we were trying to merge our code into Linux, a lot of changes were happening around Linux ARM support in general. The amount of churn in SoC support code was becoming an increasingly big problem for maintainers, and much work was underway to reduce board specific code and support a single ARM kernel binary bootable across multiple SoCs. In light of these ongoing changes, getting enough time from ARM kernel maintainers to review the code was challenging, and there was extra pressure on the maintainers to be highly critical of any new code merged into the ARM tree. We had no choice but to keep maintaining and improving the code, and regularly send out updated patch series that followed upstream kernel changes. Eventually Will Deacon, one of the ARM maintainers, made time for several comprehensive and helpful reviews, and after addressing his concerns, he gave us his approval of the code. After all this, when we thought we were done, we finally received some feedback from Russell King.

When MMIO operations trap to the hypervisor, the virtualization extensions populate a register which contains information useful to emulate the instruction (whether it was a load or a store, source/target registers, and the length of MMIO accesses). A certain class of instructions used by older Linux kernels do not populate such a register. KVM/ARM therefore loads the instruction from memory and decodes it in software. Even though the decoding implementation was well tested and reviewed by a large group of people, Russell King objected to including this feature. He had already implemented multiple forms of instruction decoding in other subsystems and demanded that we either rewrite significant parts of the ARM kernel to unify all instruction decoding to improve code reuse, or drop the MMIO instruction decoding support from our implementation. Rather than pursue a rewriting effort that could drag on for months, we abandoned the otherwise well-liked and useful code base. We can only speculate about the true motives behind this decision, as the ARM maintainer would not engage in a discussion about the subject.

After 15 main patch revisions and more than 18 months, the KVM/ARM code was successfully merged into Linus's tree via Russell King's ARM tree in February 2013. In getting all these things to come together in the end before the 3.9 merge window, the key was having a good relationship with many of the kernel developers to get their help, and being persistent in continuing to push to have the code merged in the face of various challenges.

5. Experimental Results

We present some experimental results that quantify the performance of KVM/ARM on multicore ARM hardware. We evaluate the virtualization overhead of KVM/ARM compared to native execution by running both microbenchmarks and real application workloads within VMs and directly on the hardware. We measure and compare the performance, energy, and implementation costs of KVM/ARM versus KVM x86 to demonstrate the effectiveness of KVM/ARM against a more mature hardware virtualization platform. These results provide the first real hardware measurements of the performance of ARM hardware virtualization support as well as the first comparison between ARM and x86.

5.1 Methodology

ARM measurements were obtained using an Insignal Arndale board [19] with a dual core 1.7GHz Cortex A-15 CPU on a Samsung Exynos 5250 SoC. This is the first and most widely used commercially available development board based on the Cortex A-15, the first ARM CPU with hardware virtualization support. Onboard 100Mb Ethernet is provided via the USB bus and an external 120GB Samsung 840 series SSD drive was connected to the Arndale board via eSATA. x86 measurements were obtained using both a low-power mobile laptop platform and an industry standard server platform. The laptop platform was a 2011 MacBook Air with a dual core 1.8GHz Core i7-2677M CPU, an internal Samsung SM256C 256GB SSD drive, and an Apple 100Mb USB Ethernet adapter. The server platform was a dedicated OVH SP 3 server with a dual core 3.4GHz Intel Xeon E3 1245v2 CPU, two physical SSD drives of which only one was used, and 1GB Ethernet connected to a 100Mb network infrastructure. x86 hardware with virtual APIC support was not yet available at the time of our experiments.

Given the differences in hardware platforms, our focus was not on measuring absolute performance, but rather the relative performance differences between virtualized and native execution on each platform. Since our goal is to evaluate hypervisors, not raw hardware performance, this relative measure provides a useful cross-platform basis for comparing the virtualization performance and power costs of KVM/ARM versus KVM x86.

To provide comparable measurements, we kept the software environments across all hardware platforms the same as much as possible. Both the host and guest VMs on all platforms were Ubuntu version 12.10. We used the mainline Linux 3.10 kernel for our experiments, with patches for huge page support applied on top of the source tree. Since the experiments were performed on a number of different platforms, the kernel configurations had to be slightly different, but all common features were configured similarly across all platforms. In particular, Virtio drivers were used in the guest VMs on both ARM and x86. We used QEMU version v1.5.0 for our measurements. All systems were configured with a maximum of 1.5GB of RAM available to the respective guest VM or host being tested. Furthermore, all multicore measurements were done using two physical cores and guest VMs with two virtual CPUs, and single-core measurements

were configured with SMP disabled in the kernel configuration of both the guest and host system; hyperthreading was disabled on the x86 platforms. CPU frequency scaling was disabled to ensure that native and virtualized performance was measured at the same clock rate on each platform.

For measurements involving the network and another server, 100Mb Ethernet was used on all systems. The ARM and x86 laptop platforms were connected using a Netgear GS608v3 switch, and a 2010 iMac with a 3.2GHz Core i3 CPU with 12GB of RAM running Mac OS X Mountain Lion was used as a server. The x86 server platform was connected to a 100Mb port in the OVH network infrastructure, and another identical server in the same data center was used as the server. While there are some differences in the network infrastructure used for the x86 server platform because it is controlled by someone else, we do not expect these differences to have any significant impact on the relative performance between virtualized and native execution.

We present results for four sets of experiments. **First, we measured the cost of various micro architectural characteristics of the hypervisors on multicore hardware using custom small guest OSes [11, 21] with some bugfix patches applied.** We further instrumented the code on both KVM/ARM and KVM x86 to read the cycle counter at specific points along critical paths to more accurately determine where overhead time was spent.

Second, we measured the cost of a number of common low-level OS operations using lmbench [25] v3.0 on both single-core and multicore hardware. When running lmbench on multicore configurations, we pinned each benchmark process to a separate CPU to measure the true overhead of interprocessor communication in VMs on multicore systems.

Third, we measured real application performance using a variety of workloads on both single-core and multicore hardware. Table 2 describes the eight application workloads we used.

Fourth, we measured energy efficiency using the same eight application workloads used for measuring application performance. ARM power measurements were performed using an ARM Energy Probe [3] which measures power consumption over a shunt attached to the power supply of the Arndale board. Power to the external SSD was delivered by attaching a USB power cable to the USB ports on the Arndale board thereby factoring storage power into the total SoC power measured at the power supply. x86 power measurements were performed using the powerstat tool, which reads ACPI information. powerstat measures total system power draw from the battery, so power measurements on the x86 system were run from battery power and could only be run on the x86 laptop platform. Although we did not measure the power efficiency of the x86 server platform, it is expected to be much less efficient than the x86 laptop platform, so using the x86 laptop platform provides a conservative comparison of energy efficiency against ARM. The display and wireless features of the x86 laptop platform were turned off to ensure a fair comparison. Both tools reported instantaneous power draw in watts with a 10Hz interval. These

apache	Apache v2.2.22 Web server running ApacheBench v2.3 on the local server, which measures number of handled requests per seconds serving the index file of the GCC 4.4 manual using 100 concurrent requests
mysql	MySQL v14.14 (distrib 5.5.27) running the SysBench OLTP benchmark using the default configuration
memcached	memcached v1.4.14 using the memslap benchmark with a concurrency parameter of 100
kernel compile	kernel compilation by compiling the Linux 3.6.0 kernel using the vexpress.defconfig for ARM using GCC 4.7.2 on ARM and the GCC 4.7.2 arm-linux-gnueabi- cross compilation toolchain on x86
untar	untar extracting the 3.6.0 Linux kernel image compressed with bz2 compression using the standard tar utility
curl 1K	curl v7.27.0 downloading a 1KB randomly generated file 1,000 times from the respective iMac or OVH server and saving the result to /dev/null with output disabled, which provides a measure of network latency
curl 1G	curl v7.27.0 downloading a 1GB randomly generated file from the respective iMac or OVH server and saving the result to /dev/null with output disabled, which provides a measure of network throughput
hackbench	hackbench [26] using unix domain sockets and 100 process groups running with 500 loops

Table 2: Benchmark Applications

measurements were averaged and multiplied by the duration of the test to obtain an energy measure.

5.2 Performance and Power Measurements

Table 3 presents various micro-architectural costs of virtualization using KVM/ARM on ARM and KVM x86 on x86. Measurements are shown in cycles instead of time to provide a useful comparison across platforms with different CPU frequencies. We show two numbers for the ARM platform where possible, with and without VGIC and virtual timers support.

Hypercall is the cost of two world switches, going from the VM to the host and immediately back again without doing any work in the host. **KVM/ARM takes three to four times as many cycles for this operation versus KVM x86 due to two main factors. First, saving and restoring VGIC state to use virtual interrupts is quite expensive on ARM;** available x86 hardware does not yet provide such mechanism. The ARM no VGIC/vtimers measurement does not include the cost of saving and restoring VGIC state, showing that this accounts for over half of the cost of a world switch on ARM. **Second, x86 provides hardware support to save and restore state on the world switch, which is much faster. ARM requires software to explicitly save and restore state, which provides greater flexibility, but higher costs.** Nevertheless, without the VGIC state, the hypercall costs are only about 600 cycles more than the hardware accelerated hypercall cost on the x86 server platform. The ARM world switch costs have not been optimized and can be reduced further. For example, a small patch eliminating unnecessary atomic operations reduces the hypercall cost by roughly 300 cycles, but did not make it into the mainline kernel until after v3.10 was released. As another example, if parts of the VGIC state were lazily context switched instead of

Micro Test	ARM	ARM no VGIC/vtimers	x86 laptop	x86 server
Hypercall	5,326	2,270	1,336	1,638
Trap	27	27	632	821
I/O Kernel	5,990	2,850	3,190	3,291
I/O User	10,119	6,704	10,985	12,218
IPI	14,366	32,951	17,138	21,177
EOI+ACK	427	13,726	2,043	2,305

Table 3: Micro-Architectural Cycle Counts

being saved and restored on each world switch, this may also reduce the world switch costs.

Trap is the cost of switching the hardware mode from the VM into the respective CPU mode for running the hypervisor, Hyp mode on ARM and root mode on x86. **ARM is much faster than x86 because it only needs to manipulate two registers to perform this trap, whereas the cost of a trap on x86 is roughly the same as the cost of a world switch because the same amount of state is saved by the hardware in both cases. The trap cost on ARM is a very small part of the world switch costs, indicating that the double trap incurred by split-mode virtualization on ARM does not add much overhead.**

I/O Kernel is the cost of an I/O operation from the VM to a device, which is emulated inside the kernel. I/O User shows the cost of issuing an I/O operation to a device emulated in user space, adding to I/O Kernel the cost of transitioning from the kernel to a user space process and doing a small amount of work in user space on the host for I/O. This is representative of the cost of using QEMU. **Since these operations involve world switches, saving and restoring VGIC state is again a significant cost on ARM. KVM x86 is faster than KVM/ARM on I/O Kernel, but slightly slower on I/O User. This is because the hardware optimized world switch on x86 constitutes the majority of the cost of performing I/O in the kernel, but transitioning from kernel to a user space process on the host side is more expensive on x86 because x86 KVM saves and restores additional state lazily when going to user space. Note that the added cost of going to user space includes saving additional state, doing some work in user space, and returning to the kernel and processing the KVM_RUN ioctl call for KVM.**

IPI is the cost of issuing an IPI to another virtual CPU core when both virtual cores are running on separate physical cores and both are actively running inside the VM. **IPI measures time starting from sending an IPI until the other virtual core responds and completes the IPI. It involves multiple world switches and sending and receiving a hardware IPI. Despite its higher world switch cost, ARM is faster than x86 because the underlying hardware IPI on x86 is expensive, x86 APIC MMIO operations require KVM x86 to perform instruction decoding not needed on ARM, and completing an interrupt on x86 is more expensive. ARM without VGIC/vtimers is significantly slower than with VGIC/vtimers even though it has lower world switch costs because sending, EOling and ACKing interrupts trap to the hypervisor and are handled by QEMU in user space.**

EOI+ACK is the cost of completing a virtual interrupt on both platforms. **It includes both interrupt acknowledgment and completion on ARM, but only completion on the x86 platform. ARM requires an additional operation, the acknowledgment, to the interrupt controller to determine the source of the interrupt. x86 does not because the source is directly indicated by the interrupt descriptor table entry at the time when the interrupt is raised.** However, the operation is roughly 5 times faster on ARM than x86 because there is no need to trap to the hypervisor on ARM because of VGIC support for both operations. On x86, the EOI operation must be emulated and therefore causes a trap to the hypervisor. This operation is required for every virtual interrupt including both virtual IPIs and interrupts from virtual devices.

Figures 3 to 7 show virtualized execution measurements normalized relative to their respective native execution measurements, with lower being less overhead. Figures 3 and 4 show normalized performance for running Imbench in a VM versus running directly on the host. Figure 3 shows that KVM/ARM and KVM x86 have similar virtualization overhead in a single core configuration. For comparison, we also show KVM/ARM performance without VGIC/vtimers. **Overall, using VGIC/vtimers provides slightly better performance except for the pipe and ctxsw workloads where the difference is substantial. The reason for the high overhead in this case is caused by updating the run-queue clock in the Linux scheduler every time a process blocks, since reading a counter traps to user space without vtimers on the ARM platform. We verified this by running the workload with VGIC support, but without vtimers, and we counted the number of timer read exits when running without vtimers support.**

Figure 4 shows more substantial differences in virtualization overhead between KVM/ARM and KVM x86 in a multicore configuration. KVM/ARM has less overhead than KVM x86 fork and exec, but more for protection faults. Both systems have the worst overhead for the pipe and ctxsw workloads, though KVM x86 is more than two times worse for pipe. This is due to the cost of repeatedly sending an IPI from the sender of the data in the pipe to the receiver for each message and the cost of sending an IPI when scheduling a new process. **x86 not only has higher IPI overhead than ARM, but it must also EOI each IPI, which is much more expensive on x86 than on ARM because this requires trapping to the hypervisor on x86 but not on ARM. Without using VGIC/vtimers, KVM/ARM also incurs high overhead comparable to KVM x86 because it then also traps to the hypervisor to ACK and EOI the IPIs.**

Figures 5 and 6 show normalized performance for running application workloads in a VM versus running directly on the host. Figure 5 shows that KVM/ARM and KVM x86 have similar virtualization overhead across all workloads in a single core configuration except for the MySQL workloads, but Figure 6 shows that there are more substantial differences in performance on multicore. On multicore, KVM/ARM has significantly less virtualization overhead than KVM x86 on Apache and MySQL. Overall on multicore, KVM/ARM performs within 10% of running directly on the hardware for all application workloads,

What is the ACK used for ?

Important

Dont understand the reason

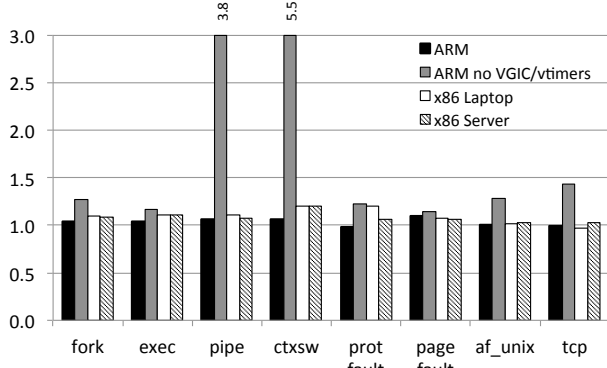


Figure 3: UP VM Normalized Imbench Performance

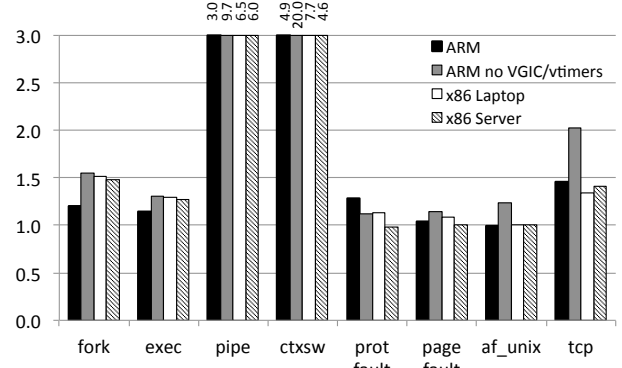


Figure 4: SMP VM Normalized Imbench Performance

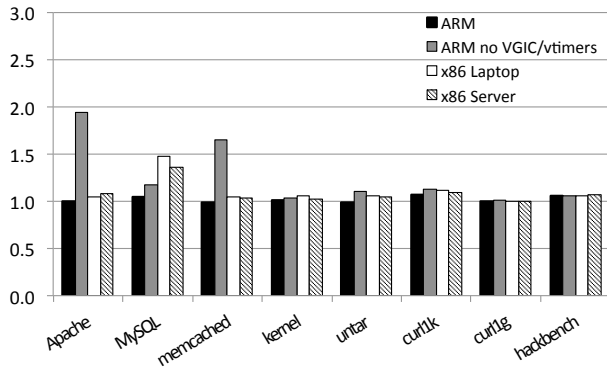


Figure 5: UP VM Normalized Application Performance

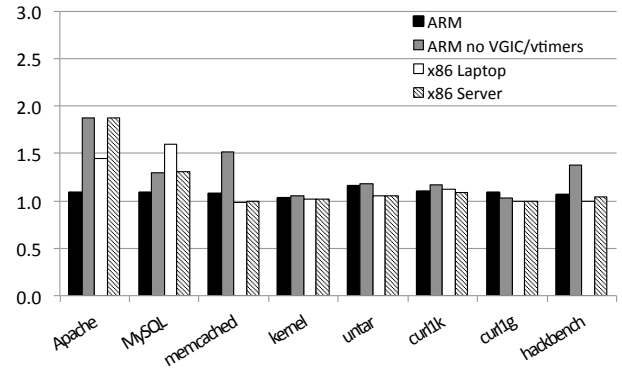


Figure 6: SMP VM Normalized Application Performance

while the more mature KVM x86 system has significantly higher virtualization overheads for Apache and MySQL. KVM/ARM's split-mode virtualization design allows it to leverage ARM hardware support with comparable performance to a traditional hypervisor using x86 hardware support. The measurements also show that KVM/ARM performs better overall with ARM VGIC/vtimers support than without.

Figure 7 shows normalized power consumption of using virtualization versus direct execution for various application workloads on multicore. We only compared KVM/ARM on ARM against KVM x86 on x86 laptop. The Intel Core i7 CPU used in these experiments is one of Intel's more power optimized processors, and we expect that server power consumption would be even higher. The measurements show that KVM/ARM using VGIC/vtimers is more power efficient than KVM x86 virtualization in all cases except memcached and untar. Both workloads are not CPU bound on both platforms and the power consumption is not significantly affected by the virtualization layer. However, due to ARM's slightly higher virtualization overhead for these workloads, the energy virtualization overhead is slightly higher on ARM for the two workloads. While a more detailed study of energy aspects of virtualization is beyond the scope of this paper, these measurements nevertheless provide useful data comparing ARM and x86 virtualization energy costs.

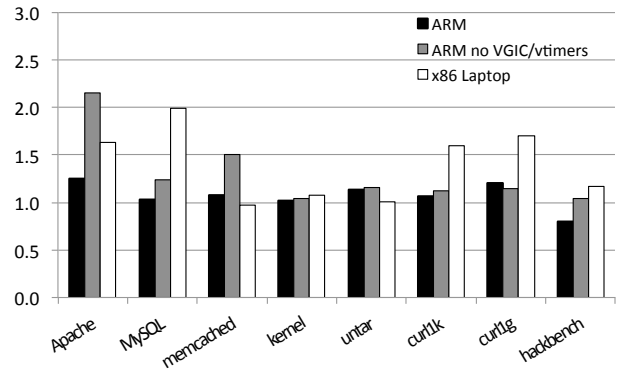


Figure 7: SMP VM Normalized Energy Consumption

5.3 Implementation Complexity

We compare the code complexity of KVM/ARM to its KVM x86 counterpart in Linux 3.10. KVM/ARM is 5,812 lines of code (LOC), counting just the architecture-specific code added to Linux to implement it, of which the lowvisor is a mere 718 LOC. As a conservative comparison, KVM x86 is 25,367 LOC, excluding guest performance monitoring support, not yet supported by KVM/ARM, and 3,311 LOC required for AMD support. These numbers do not include KVM's architecture-generic code, 7,071 LOC, which is shared by all systems. Table 4 shows a

Component	KVM/ARM	KVM x86 (Intel)
Core CPU	2,493	16,177
Page Fault Handling	738	3,410
Interrupts	1,057	1,978
Timers	180	573
Other	1,344	1,288
Architecture-specific	5,812	25,367

Table 4: Code Complexity in Lines of Code (LOC)

breakdown of the total hypervisor architecture-specific code into its major components.

By inspecting the code we notice that the striking additional complexity in the x86 implementation is mainly due to the five following reasons: (1) Since EPT was not supported in earlier hardware versions, KVM x86 must support shadow page tables. (2) The hardware virtualization support have evolved over time, requiring software to conditionally check for support for a large number of features such as EPT. (3) A number of operations require software decoding of instructions on the x86 platform. KVM/ARM’s out-of-tree MMIO instruction decode implementation was much simpler, only 462 LOC. (4) The various paging mode on x86 requires more software logic to handle page faults. (5) x86 requires more software logic to support interrupts and timers than ARM, which provides VGIC/vtimers hardware support that reduces software complexity.

KVM/ARM’s LOC is less than partially complete bare-metal microvisors written for Hyp mode [31], with the lowvisor LOC almost an order of magnitude smaller. Unlike standalone hypervisors, KVM/ARM’s code complexity is so small because lots of functionality simply does not have to be implemented because it is already provided by Linux. Table 4 does not include other non-hypervisor architecture-specific Linux code, such as basic bootstrapping, which is significantly more code. Porting a standalone hypervisor such as Xen from x86 to ARM is much more complicated because all of that ARM code for basic system functionality needs to be written from scratch. In contrast, since Linux is dominant on ARM, KVM/ARM just leverages existing Linux ARM support to run on every platform supported by Linux.

6. Recommendations

From our experiences building KVM/ARM, we offer a few recommendations for hardware designers to simplify and optimize future hypervisor implementations.

Share kernel mode memory model. The hardware mode to run a hypervisor should use the same memory model as the hardware mode to run OS kernels. Software designers then have greater flexibility in deciding how tightly to integrate a hypervisor with existing OS kernels. ARM Hyp mode unfortunately did not do this, preventing KVM/ARM from simply reusing the kernel’s page tables in Hyp mode. This reuse would have simplified the implementation and allowed for performance critical emulation code to run in Hyp mode, avoiding a complete world switch in some cases. Some might argue that this recommendation makes

for more complicated standalone hypervisor implementations, but this is not really true. For example, ARM kernel mode already has a simple option to use one or two page table base registers to unify or split the address space. Our recommendation is different from the x86 virtualization approach, which does not have a separate and more privileged, hypervisor CPU mode. Having a separate CPU mode potentially improves stand-alone hypervisor performance and implementation, but not sharing the kernel memory model complicates the design of hypervisors integrated with host kernels.

Make VGIC state access fast, or at least infrequent. While VGIC support can improve performance especially on multicore systems, our measurements also show that access to VGIC state adds substantial overhead to world switches. This is caused by slow MMIO access to the VGIC control interface in the critical path. Improving the MMIO access time is likely to improve VM performance, but if this is not possible or cost-effective, MMIO accesses to the VGIC could at least be made less frequent. For example, a summary register could be introduced describing the state of each virtual interrupt. This could be read when performing a world switch from the VM to the hypervisor to get information which can currently only be obtained by reading all the list registers (see Section 3.5) on each world switch.

Completely avoid IPI traps. Hardware support to send virtual IPIs directly from VMs without the need to trap to the hypervisor would improve performance. Hardware designers may underestimate how frequent IPIs are on modern multicore OSes, and our measurements reveal that sending IPIs adds significant overhead for some workloads. The current VGIC design requires a trap to the hypervisor to emulate access to the IPI register in the distributor, and this emulated access must be synchronized between virtual cores using a software locking mechanism, which adds significant overhead for IPIs. Current ARM hardware supports receiving the virtual IPIs, which can be ACKed and EOled without traps, but unfortunately does not address the also important issue of sending virtual IPIs.

7. Related Work

Virtualization has a long history [27], but has enjoyed a resurgence starting in the late 1990s. Most efforts have almost exclusively focused on virtualizing the x86 architecture. While systems such as VMware [1, 10] and Xen [8] were originally based on software-only approaches before the introduction of x86 hardware virtualization support, all x86 virtualization platforms, VMware [2], Xen, and KVM [22], now leverage x86 hardware virtualization support. Because x86 hardware virtualization support differs substantially from ARM in the ability to completely run the hypervisor in the same mode as the kernel, x86 virtualization approaches do not lend themselves directly to take advantage of ARM hardware virtualization support.

Some x86 approaches also leverage the host kernel to provide functionality for the hypervisor. VMware Workstation’s hypervisor creates a VMM separate from the host kernel, but this

approach is different from KVM/ARM in a number of important ways. First, the VMware VMM is not integrated into the host kernel source code and therefore cannot reuse existing host kernel code, for example, for populating page tables relating to the VM. Second, since the VMware VMM is specific to x86 it does not run across different privileged CPU modes, and therefore does not use a design similar to KVM/ARM. Third, most of the emulation and fault-handling code required to run a VM executes at the most privileged level inside the VMM. KVM/ARM executes this code in the less privileged kernel mode, and only executes a minimal amount of code in the most privileged mode. In contrast, KVM benefits from being integrated with the Linux kernel like KVM/ARM, but the x86 design relies on being able to run the kernel and the hypervisor together in the same hardware hypervisor mode, which is problematic on ARM.

Full-system virtualization of the ARM architecture is a relatively unexplored research area. Most approaches are software only. A number of standalone bare metal hypervisors have been developed [16, 17, 28], but these are not widespread, are developed specifically for the embedded market, and must be modified and ported to every single host hardware platform, limiting their adoption. An abandoned port of Xen for ARM [18] requires comprehensive modifications to the guest kernel, and was never fully developed. An earlier prototype for KVM on ARM [12, 15] used an automated lightweight paravirtualization approach to automatically patch kernel source code to run as a guest kernel, but had poor performance. VMware Horizon Mobile [9] uses hosted virtualization to leverage Linux's support for a wide range of hardware platforms, but requires modifications to guest OSes and its performance is unproven. None of these paravirtualization approaches could run unmodified guest OSes.

An earlier study attempted to estimate the performance of ARM hardware virtualization support using a software simulator and a simple hypervisor lacking important features like SMP support and use of storage and network devices by multiple VMs [31]. Because of the lack of hardware or a cycle-accurate simulator, no real performance evaluation was possible. In contrast, we present the first evaluation of ARM virtualization extensions using real hardware, provide a direct comparison with x86, and present the design and implementation of a complete hypervisor using ARM virtualization extensions, including SMP support.

A newer version of Xen exclusively targeting servers [32] is being developed using ARM hardware virtualization support. Because Xen is a bare metal hypervisor that does not leverage kernel functionality, it can be architected to run entirely in Hyp mode rather than using split-mode virtualization. At the same time, this requires a substantial commercial engineering effort. Since Xen is a standalone hypervisor, porting Xen from x86 to ARM is difficult in part because all ARM-related code must be written from scratch. Even after getting Xen to work on one ARM platform, it must be manually ported to each different ARM device that Xen wants to support. Because of Xen's custom I/O model using hypercalls from VMs for device emulation on ARM, Xen unfortunately cannot run guest OSes unless they have been configured

to include Xen's hypercall layer and include support for XenBus paravirtualized drivers. In contrast, KVM/ARM uses standard Linux components to enable faster development, full SMP support, and the ability to run unmodified OSes. KVM/ARM is easily supported on new devices with Linux support, and we spent almost no effort to support KVM/ARM on ARM's Versatile Express boards, the Arndale board, and hardware emulators. While Xen can potentially reduce world switch times for operations that can be handled inside the Xen hypervisor, switching to Dom0 for I/O support or switching to other VMs would involve context switching the same state as KVM/ARM.

Microkernel approaches for hypervisors [16, 30] have been used to reduce the hypervisor TCB and run other hypervisor services in user mode. These approaches differ both in design and rationale from split-mode virtualization, which splits hypervisor functionality across privileged modes to leverage virtualization hardware support. Split-mode virtualization also provides a different split of hypervisor functionality. KVM/ARM's lowvisor is a much smaller code base that implements only the lowest level hypervisor mechanisms. It does not include higher-level functionality present in the hypervisor TCB used in these other approaches.

8. Conclusions

KVM/ARM is the mainline Linux ARM hypervisor and the first system that can run unmodified guest operating systems on ARM multicore hardware. KVM/ARM's split-mode virtualization makes it possible to use ARM hardware virtualization extensions while leveraging Linux kernel mechanisms and hardware support. Our experimental results show that KVM/ARM (1) incurs minimal performance impact from the extra traps incurred by split-mode virtualization, (2) has modest virtualization overhead and power costs, within 10% of direct native execution on multicore hardware for real application workloads, and (3) achieves comparable or lower virtualization overhead and power costs on multicore hardware compared to widely-used KVM x86 virtualization. Based on our experiences integrating KVM/ARM into the mainline Linux kernel, we provide some hints on getting research ideas and code adopted by the open source community, and recommendations for hardware designers to improve future hypervisor implementations.

9. Acknowledgments

Marc Zyngier helped with development, implemented VGIC and vtimers support, and assisted us with hardware bring up. Rusty Russell worked on the coprocessor user space interface and assisted with upstreaming. Will Deacon and Avi Kivity provided numerous helpful code reviews. Peter Maydell helped with QEMU support and debugging. Gleb Natapov helped us better understand KVM x86 performance. Marcelo Tosatti and Nicolas Viennot helped with resolving what became known as the voodoo bug. Keith Adams, Hani Jamjoom, and Emmett Witchel provided helpful comments on earlier drafts of this paper. This work was supported in part by ARM and NSF grants CNS-1162447, CCF-1162021, and CNS-1018355.

References

- [1] K. Adams and O. Agesen. A Comparison of Software and Hardware Techniques for x86 Virtualization. In *Proceedings of the 12th International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 2–13, Oct. 2006.
- [2] O. Agesen, J. Mattson, R. Rugina, and J. Sheldon. Software Techniques for Avoiding Hardware Virtualization Exits. In *Proceedings of the 2012 USENIX Annual Technical Conference*, pages 373–385, June 2012.
- [3] ARM Ltd. ARM Energy Probe. <http://www.arm.com/products/tools/arm-energy-probe.php>.
- [4] ARM Ltd. ARM Cortex-A15 Technical Reference Manual ARM DDI 0438C, Sept. 2011.
- [5] ARM Ltd. ARM Generic Interrupt Controller Architecture version 2.0 ARM IHI 0048B, June 2011.
- [6] ARM Ltd. ARM Architecture Reference Manual ARMv7-A DDI0406C.b, July 2012.
- [7] ARM Ltd. ARM Architecture Reference Manual ARMv8-A DDI0487A.a, Sept. 2013.
- [8] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, and A. Warfield. Xen and the Art of Virtualization. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles*, pages 164–177, Oct. 2003.
- [9] K. Barr, P. Bungale, S. Deasy, V. Gyuris, P. Hung, C. Newell, H. Tuch, and B. Zoppis. The VMware Mobile Virtualization Platform: is that a hypervisor in your pocket? *SIGOPS Operating Systems Review*, 44(4):124–135, Dec. 2010.
- [10] E. Bugnion, S. Devine, M. Rosenblum, J. Sugerman, and E. Y. Wang. Bringing Virtualization to the x86 Architecture with the Original VMware Workstation. *ACM Transactions on Computer Systems*, 30(4):12:1–12:51, Nov. 2012.
- [11] C. Dall and A. Jones. KVM/ARM Unit Tests. <https://github.com/columbia/kvm-unit-tests>.
- [12] C. Dall and J. Nieh. KVM for ARM. In *Proceedings of the Ottawa Linux Symposium*, pages 45–56, July 2010.
- [13] C. Dall and J. Nieh. Supporting KVM on the ARM architecture. LWN.net, July 2013. <http://lwn.net/Articles/557132/>.
- [14] David Brash, Architecture Program Manager, ARM Ltd. Personal communication, Nov. 2012.
- [15] J.-H. Ding, C.-J. Lin, P.-H. Chang, C.-H. Tsang, W.-C. Hsu, and Y.-C. Chung. ARMvisor: System Virtualization for ARM. In *Proceedings of the Ottawa Linux Symposium (OLS)*, pages 93–107, July 2012.
- [16] General Dynamics. OKL4 Microvisor. <http://www.ok-labs.com/products/okl4-microvisor>.
- [17] Green Hills Software. INTEGRITY Secure Virtualization. http://www.ghs.com/products/rtos/integrity_virtualization.html.
- [18] J. Hwang, S. Suh, S. Heo, C. Park, J. Ryu, S. Park, and C. Kim. Xen on ARM: System Virtualization using Xen Hypervisor for ARM-based Secure Mobile Phones. In *Proceedings of the 5th Consumer Communications and Network Conference*, Jan. 2008.
- [19] InSignal Co. ArndaleBoard.org. <http://arndaleboard.org>.
- [20] Intel Corporation. Intel 64 and IA-32 Architectures Software Developers Manual, 325462-044US, Aug. 2012.
- [21] A. Kivity. KVM Unit Tests. <https://git.kernel.org/cgit/virt/kvm/kvm-unit-tests.git>.
- [22] A. Kivity, Y. Kamay, D. Laor, U. Lublin, and A. Liguori. **kvm**: The Linux Virtual Machine Monitor. In *Proceedings of the Ottawa Linux Symposium (OLS)*, volume 1, pages 225–230, June 2007.
- [23] KVM/ARM Mailing List. <https://lists.cs.columbia.edu/cuclists/listinfo/kvmarm>.
- [24] Linux ARM Kernel Mailing List. A15 H/W Virtualization Support, Apr. 2011. <http://archive.arm.linux.org.uk/lurker/message/20110412.204714.a36702d9.en.html>.
- [25] L. McVoy and C. Staelin. Imbench: Portable Tools for Performance Analysis. In *Proceedings of the 1996 USENIX Annual Technical Conference*, pages 279–294, Jan. 1996.
- [26] I. Molnar. Hackbench. <http://people.redhat.com/mingo/cfs-scheduler/tools/hackbench.c>.
- [27] G. J. Popek and R. P. Goldberg. Formal Requirements for Virtualizable Third Generation Architectures. *Communications of the ACM*, 17(7):412–421, July 1974.
- [28] Red Bend Software. vLogix Mobile. <http://www.redbend.com/en/mobile-virtualization>.
- [29] R. Russell. virtio: Towards a De-Facto Standard for Virtual I/O Devices. *SIGOPS Operating Systems Review*, 42(5):95–103, July 2008.
- [30] U. Steinberg and B. Kauer. Nova: A Microhypervisor-Based Secure Virtualization Architecture. In *Proceedings of the 5th European Conference on Computer Systems*, pages 209–222, Apr. 2010.
- [31] P. Varanasi and G. Heiser. Hardware-Supported Virtualization on ARM. In *Proceedings of the Second Asia-Pacific Workshop on Systems*, pages 11:1–11:5, July 2011.
- [32] Xen.org. Xen ARM. http://xen.org/products/xen_arm.html.