

A Light-Weight Virtual Machine Monitor for Blue Gene/P

Jan Stoess^{§¶1}
Jonathan Appavoo^{†1}

Udo Steinberg^{‡1}
Amos Waterland^{|1}

Volkmar Uhlig^{¶1}
Jens Kehne[§]

[§]Karlsruhe Institute of Technology

[‡]Technische Universität Dresden

[¶]HStreaming LLC

[|]Harvard School of Engineering and Applied Sciences

[†]Boston University

ABSTRACT

In this paper, we present a light-weight, micro-kernel-based virtual machine monitor (VMM) for the Blue Gene/P Supercomputer. Our VMM comprises a small μ -kernel with virtualization capabilities and, atop, a user-level VMM component that manages virtual BG/P cores, memory, and interconnects; we also support running native applications directly atop the μ -kernel. Our design goal is to enable compatibility to standard OSes such as Linux on BG/P via virtualization, but to also keep the amount of kernel functionality small enough to facilitate shortening the path to applications and lowering OS noise.

Our prototype implementation successfully virtualizes a BG/P version of Linux with support for Ethernet-based communication mapped onto BG/P's collective and torus network devices. First experiences and experiments show that our VMM still shows a substantial performance hit; nevertheless, our approach poses an interesting OS alternative for Supercomputers, providing the convenience of a fully-featured commodity software stack, while also promising to deliver the scalability and low latency of an HPC OS.

1. INTRODUCTION

A substantial fraction of supercomputer programmers today write software using a parallel programming run-time such as MPI on top of a customized light-weight kernel. For Blue Gene/P (BG/P) machines in production, IBM provides such a light-weight kernel called Compute Node Kernel (CNK) [9]. CNK runs tasks massively parallel, in a single-thread-per-core fashion. Like other light-weight kernels, CNK supports a subset of a standardized application interface (POSIX), facilitating the development of dedicated (POSIX-like) applications for a supercomputer. However, CNK is not fully POSIX-compatible: it lacks, for instance, comprehensive scheduling or memory management as well as standards-compatible networking or support for standard

¹This research was mostly conducted by the authors while at IBM Watson Research Center, Yorktown Heights, NY.

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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ROSS '11, May 31, 2011, Tucson, AZ, USA.

Copyright 2011 ACM 978-1-4503-0761-1/11/05...\$10.00.

debugging tools. CNK also supports I/O only via function-shipping to I/O nodes.

CNK's lightweight kernel model is a good choice for the current set of BG/P HPC applications, providing low operating system (OS) noise and focusing on performance, scalability, and extensibility. However, today's HPC application space is beginning to scale out towards Exascale systems of *truly* global dimensions, spanning companies, institutions, and even countries. The restricted support for standardized application interfaces of light-weight kernels in general and CNK in particular renders porting the sprawling diversity of scalable applications to supercomputers more and more a bottleneck in the development path of HPC applications.

In this paper, we explore an alternative, hybrid OS design for BG/P: a μ -kernel-based virtual machine monitor (VMM). At the lowest layer, in kernel mode, we run a μ -kernel that provides a small set of basic OS primitives for constructing customized HPC applications and services at user level. We then construct a user-level VMM that fully virtualizes the BG/P platform and allows arbitrary Blue Gene OSes to run in virtualized compartments. In this paper, we focus on the *virtualization layer* of our OS architecture; the field of research on native application frameworks for μ -kernels has been rich [10, 17], and has even been explored for HPC systems [21].

The benefits of a μ -kernel-based VMM architecture are twofold: on the one hand, it provides compatibility to BG/P hardware, allowing programmers to ship the OS they require for their particular applications along, like a library. For instance, our VMM successfully virtualizes a Blue Gene version of Linux with support for Ethernet-based communication, allowing virtually any general-purpose Linux application or service to run on BG/P. On the other hand, our μ -kernel also resembles the light-weight kernel approach in that it reduces the amount of kernel functionality to basic resource management and communication. Those mechanisms are available to native applications running directly on top of the μ -kernel, and programmers can use them to customize their HPC applications for better efficiency and scalability, and to directly exploit the features of the tightly-interconnected BG/P hardware. However, μ -kernel and VMM architectures also imply potential penalty to efficiency, as they increase kernel-user interaction and add another layer of indirection to the system software stack. Nevertheless, the need for standardized application interfaces is becoming more prevalent, and we expect our work to be an insightful step towards supporting such standardization on supercomputer platforms.

The idea of a virtualization layer for HPC systems is not new [21], nor is the idea of using a decomposed VMM architecture to deliver predictable application performance [12]. However, to our knowledge, this is the first approach to providing a commodity system software stack and hiding the hardware peculiarities of such a highly-customized architecture, while still being able to run hand-optimized code side-by-side. Initial results are promising: Our prototype based on the L4 μ -kernel fully virtualizes BG/P compute nodes and their high-performance network interconnects, and successfully runs multiple instances of a BG/P version of Linux.

The rest of the paper is structured as follows: Section 2 presents the basic architecture of our μ -kernel-based virtualization approach for BG/P. Section 3 presents details of the μ -kernel, followed by Section 4 presenting details about our user-level VMM component. Finally, Section 5 presents an initial evaluation of our prototype, followed by related work in Section 6 and a summary in Section 7.

2. SYSTEM OVERVIEW

In this section, we present the design of our decomposed VMM for BG/P. We start with a very brief overview of the BG/P supercomputer: the basic building block of BG/P is a compute node, which is composed of an embedded quad-core PowerPC, five networks, a DDR2 controller, and either 2 or 4 GB of RAM, integrated into a system-on-a-chip (Figure 1). One of the large BG/P configurations is a 256-rack system of 2 mid-planes each, which, in turn, comprise 16 node cards with 32 nodes each, totaling over 1 million cores and 1 Petabyte of RAM. BG/P features three key communication networks, a torus and a collective interconnect, and an external 10GE Ethernet network on I/O nodes.

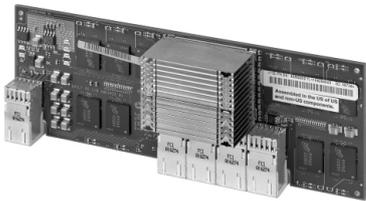


Figure 1: Blue Gene/P Compute Node Card

The basic architecture of our VMM is illustrated in Figure 2. We use an enhanced version of the L4 μ -kernel as privileged supercomputer kernel [24]. Our BG/P implementation uses a recent L4 version named L4Ka:Pistachio [19]. Traditional hypervisors such as Xen [29] or VMware [2] virtualization-only approaches in the sense that they provide virtual hardware – virtual CPUs, memory, disks, networks, etc. – as *first class* abstractions. L4 is different in that it offers a limited set of OS abstractions to enforce safe and secure execution: threads, address spaces and inter-process communication (IPC). Those abstractions are low-level enough for constructing an efficient virtualization layer atop, as has been demonstrated on commodity systems [11, 23].

While L4 provides core primitives, the actual VMM functionality is implemented as a user-level application outside of the privileged kernel. The basic mechanics of such an L4-based VMM is as follows: L4 merely acts as a safe messaging system propagating sensitive guest instructions to a

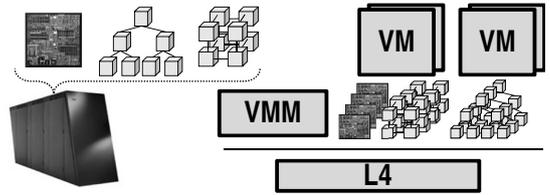


Figure 2: A μ -kernel based VMM virtualizing BG/P cores and interconnects

user-level VMM. That VMM, in turn, decodes the instruction and emulates it appropriately, and then responds with a fault reply message that instructs L4 to update the guest VM’s context and then to resume guest VM execution.

Our decomposed, μ -kernel-based VMM architecture has benefits that are particularly interesting on HPC systems: Since L4 provides minimal yet sufficiently generic abstractions, it also supports native applications that can bypass the complete virtualization stack whenever they need performance. Hybrid configurations are also possible: An application can be started within a guest VM, with access to all legacy services of the guest kernel; for improved efficiency, it can later choose to employ some native HPC library (e.g. an MPI library running directly atop L4). In the following, we will first describe how L4 facilitates running a VMM atop; we will then describe the user-level VMM part in the subsequent section.

3. A MICRO-KERNEL WITH VIRTUALIZATION CAPABILITIES

In our architecture, L4 acts as the privileged part of the VMM, responsible for partitioning processors, memory, device memory, and interrupts. Our virtualization model is largely common to L4’s normal execution model: we use i) L4 threads to virtualize the processor; ii) L4 memory mapping mechanisms to provide and manage guest-physical memory, iii) an IPC to allow emulation of sensitive instructions through the user-level VMM. However, in virtualization mode, threads and address spaces have access to an extended ISA and memory model (including a virtualized TLB), and have restricted access to L4-specific features, as we will describe in the following.

3.1 Virtual PowerPC Processor

L4 virtualizes cores by mapping each virtual CPU (vCPU) to a dedicated thread. Threads are schedulable entities, and vCPUS are treated equally: They are dispatched regularly from a CPU-local scheduling queue, and they can be moved and load-balanced among individual physical processors through standard L4 mechanisms.

In contrast to recent x86 and PowerPC based architectures, BG/P’s cores are based on the widely used embedded 440 architecture, which has no dedicated support to facilitate or accelerate virtualization. However, also in contrast to x86 architecture, PowerPC is much more virtualization-friendly in the first place: The ISA supports trap-based virtualization, and fixed instruction lengths simplify decoding and emulating sensitive instructions. L4 employs such a trap-and-emulate style virtualization method by compressing PowerPC privilege levels. L4 itself runs in supervisor mode, while user and guest both run in user mode (although

with different address space IDs, as described in Section 3.2). Guest application code runs undisturbed, but whenever the guest kernel issues a sensitive instruction, the processor causes a trap into L4.

IPC-based Virtualization. Unlike traditional VMMs, L4 itself does not emulate all sensitive instructions itself. Unless the instruction is related to the virtual TLB or can quickly be handled such as the modification of guest shadow register state, it hands off emulation to the user-level VMM component. Moving the virtualization service out of the kernel makes a fast guest-VMM interaction mechanism a prerequisite for efficient execution. We therefore rely on L4 IPCs to implement the virtualization protocol (i.e. the guest trap – VMM emulation – guest resume cycle). In effect, L4 handles guest traps the same way it handles normal page faults and exceptions, by synthesizing a fault IPC message on behalf of the guest to a designated per-vCPU exception handler.

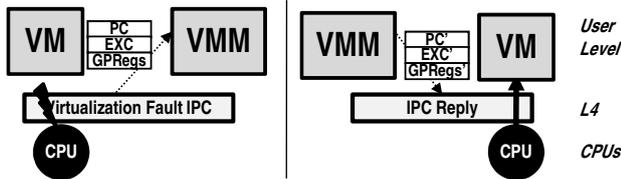


Figure 3: vCPU exits are propagated to the VMM as IPC message; the user-level VMM responds by sending back a reply IPC resuming the guest.

During a virtualization fault IPC, the trapping guest automatically blocks waiting for a reply message. To facilitate proper decoding and emulation of sensitive instructions, the virtualization IPC contains the vCPU’s current execution state such as instruction and stack pointer and general-purpose registers. The reply from the VMM may contain updated state, which L4 then transparently installs into the vCPU’s execution frame. To retain efficient transfer of guest vCPU state, the kernel allows the VMM to configure the particular state to be transferred independently for each class of faults. For instance, a VMM may choose to always transfer general-purpose registers, but to only transfer TLB-related guest state on TLB-related faults. L4 also offers, via a separate system call, to inspect and modify all guest state asynchronously from the VMM. That way, we keep the common virtualization fault path fast, while deferring loading of additional guest state to a non-critical path.

3.2 Virtualized Memory Management

PowerPC 440 avoids die costs for the page table walker and does not imply any page table format (or even construction). Instead, address translation is based solely on a translation look-aside buffer (TLB) managed in software [13]. On BG/P cores, the TLB has 64 entries managed by the kernel. Each entry maps 32-bit logical addresses to their physical counterparts, and can cover a configurable size ranging from 1 KByte to 1 GByte. For translation, the TLB further uses an 8-bit process identifier (PID) and an 1-bit address space identifier (AS) that can be freely set by the kernel, effectively implementing a tagged TLB (Figure 4).

A normal OS only needs to provide a single level of memory translation – from virtual to physical. A VMM, in contrast, must support two such translation levels, from guest-

virtual to guest-physical, and from guest-physical to host-physical memory (details on virtualized translations can be found, e.g., in [2, 18]). As PowerPC 440 hardware only supports a single translation level, the VMM must merge the two levels when inserting translations into the hardware TLB, effectively translating guest-virtual to host-physical addresses. Both L4 and the user-level VMM are involved in providing a two-level memory translation: To provide guest-physical memory, we use L4’s existing memory management abstractions based on external pagers [24]. We then provide a second, in-kernel virtual TLB that caters for the additional notion of guest virtual memory. We will describe the L4-specific extensions in the following two paragraphs.

Virtual Physical Memory. Providing guest-physical memory is largely identical to the provisioning of normal virtual memory: L4 treats a guest VM’s physical address space like a regular address space and exports establishing of translations within that address-space to a user-level pager (which is, in this case, the user-level VMM). Whenever a guest suffers a TLB miss, L4’s miss handler inspects its kernel data structures to find out whether the miss occurred because of a missing guest-physical to host-physical translation. This is the case if the VMM has not yet mapped the physical memory into the VM’s guest-physical address space. If so, L4 synthesizes a page fault IPC message to the user-level pager VMM on behalf of the faulting guest, requesting to service the fault.

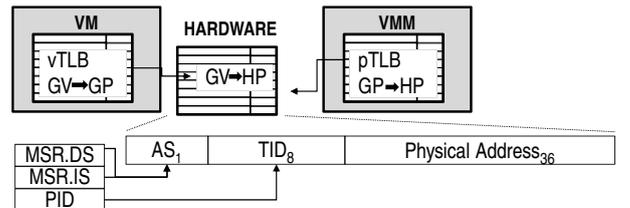


Figure 4: Merging two levels of virtualized address translation into a single-level hardware TLB.

When the VMM finds the guest-physical page fault to be valid, it responds with a mapping message, which will cause L4 to insert a valid mapping into the TLB and then to resume guest execution. Otherwise, the VMM may terminate the VM for an invalid access to non-existing physical memory or inject a hardware exception. To keep track of a guest’s mappings independent of the actual state of the hardware TLB, L4 maintains them with an in-kernel database. Should the user-level VMM revoke a particular mapping, L4 flushes the corresponding database and hardware TLB entries.

Virtual TLB. Emulating virtual address translations and the TLB is extremely critical for the overall performance of a VMM. For that reason, most virtualization-enabled hardware platforms have specific support, like the recent x86 [6] or embedded PowerPC [14] processors. On BG/P, we lack such support and reverted to a software solution using a table shadowing the hardware TLB, which we hence call virtual TLB. In order to further reduce the cost of TLB updates, we employ a number of heuristics and tracking methods. From a security perspective we only need to ensure that the hardware TLB always contains a subset of the virtual TLB.

L4 provides the notion of a virtual TLB, which the guest has access to via normal (trapped) hardware instructions for TLB management. While management of guest-physical memory involves the user-level VMM, our solution for guest-virtual memory is L4-internal: Whenever the guest kernel accesses the virtual TLB, L4’s internal instruction emulator stores those entries into a per-VM virtual TLB data structure. On a hardware TLB miss, L4’s miss handler parses that data structure to find out whether the guest has inserted a valid TLB mapping into its virtual TLB for the given fault address. If not, it injects a TLB miss fault into the guest VM to have the miss handled by the guest kernel. If the virtual TLB indeed contains a valid entry, L4 checks its mapping database to find out whether the miss occurred at the second stage, from guest-physical to host-physical. If that translation is valid as well, L4 inserts the resulting guest-virtual to host-physical mapping into the hardware TLB and resumes the VM; if it turns out to be missing, L4 synthesizes a page fault IPC to the VMM, as discussed in the previous paragraph.

Virtual Address Space Protection. Finally, a VMM must virtualize not only the translation engine of the TLB but also its protection features. Again, the virtualization logically requires two levels, allowing the guest to use the virtual TLB’s protection bits and identifiers in the same manner as on native hardware, but, at the second level, also permitting L4 and its user-level address spaces to shield their data from being accessed by guest kernel and applications. The TLB of the PowerPC 440 gives us great help. The 440 can hold up to 256 address space mappings in the TLB, (via the TID field, see Figure 4). The particular mapping is chosen through a processor register. Address space translation 0 is always accessible independently of which particular mapping is active. The 440 additionally features a 1-bit translation space, with the active translation space being selected via processor register bits, one for instruction and one for data fetches.

To facilitate trap-and-emulate virtualization, both guest kernel and applications run in user mode. L4 puts guest kernel and user into the second translation space, and keeps the first translation space reserved for itself, the VMM, and native L4 applications. The processor automatically switches to the first translation space when an interrupt or trap occurs, directly entering L4 with trap and interrupt handlers in place. To read guest memory when decoding sensitive instructions, L4 temporarily switches the translation space for data fetches, while retaining the space for instructions. Altogether, our solution allows the guest to receive a completely empty address space, but on any exception, the processor switches to a hypervisor-owned address space. That way, we keep virtualization address spaces clean without the need for ring compression as done on x86 systems lacking hardware-virtualization [31].

Our ultimate goal is to reduce the number of TLB flushes to enforce protection on user-to-kernel switches. We observed the following usage scenario for standard OSes (e.g., Linux) and implemented our algorithm to mimic the behavior: Common OSes map application code and data as user and kernel accessible, while kernel code and data is only accessible in privileged mode. We strive is to make the transition from user to kernel and back fast to achieve good system call performance; we therefore disable the address

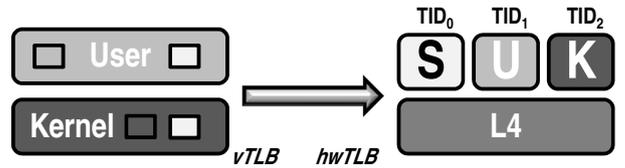


Figure 5: Virtualized TLB Protection. vTLB protection bits are mapped to TID in the hardware TLB, with TID=0 for shared pages.

space mappings completely and flush the hardware TLB on each guest address space switch (Figure 5). We then use address space 0 for mappings that are accessible to guest user and guest kernel. We use address space 1 for mappings that are only guest-user accessible and address space 2 for mappings that are only guest-kernel accessible. A privilege level switch from guest-user to guest-kernel mode therefore solely requires updating the address space identifier from 1 (user-mode mappings) to 2 (kernel-mode mappings).

Our virtual TLB effectively compresses guest user/kernel protection bits into address space identifiers; as a result, it requires hardware TLB entries to be flushed whenever the guest kernel switches guest application address spaces. Also, our scheme requires that TLB entries are flushed during world switches between different guests. It does not require, however, any TLB flushes during guest system calls or other switches from guest user to kernel, or during virtualization traps and resumes within the same VM. Thus, we optimize for frequent kernel/user and kernel/VMM switches rather than for address space or world switches, as the former occur more frequently.

3.3 Interrupt Virtualization

BG/P provides a custom interrupt controller called Blue Gene Interrupt Controller (BIC), which gathers and delivers device signals to the cores as interrupts or machine check exceptions. The BIC supports a large number (several hundreds) of interrupts, which are organized in groups and have different types for prioritization and routing purposes. The BIC supports delivering interrupts to different cores, as it supports core-to-core interrupts.

The original L4 version already provides support for user-level interrupt management, mapping interrupt messages and acknowledgments onto IPC. L4 further permits user software to migrate interrupts to different cores. Our user-level VMM uses those L4 interrupt features to receive and acknowledge interrupts for BG/P devices. To inject virtual interrupts into the guest, the VMM modifies guest vCPU state accordingly, either using L4’s state modification system call (Section 3.1), or by piggybacking the state update onto a virtualization fault reply, in case the guest VM is already waiting for the VMM when the interrupt occurs.

4. USER-LEVEL VMM

Our user-level VMM component runs as a native L4 program, and provides the virtualization service based on L4’s core abstractions. It can be described as an interface layer that translates virtualization API invocations (i.e. sensitive instructions) into API invocations of the underlying L4 architecture. As described, L4 facilitates virtualization by means of virtualization fault IPC. The user-level VMM

mainly consists of a server executing an IPC loop, waiting for any incoming IPC message from a faulting guest VM. Upon reception, it retrieves the VM register context that L4 has sent along, emulates the sensitive instruction accordingly, and finally responds with a reply IPC containing updated vCPU state such as result registers of the given sensitive instruction and an incremented program counter. Before resuming the VM, L4 installs the updated context transparently into the VM, while the VMM waits for the next message to arrive.

4.1 Emulating Sensitive Instructions

Our user-level VMM largely resembles typical other VMMs such as VMware or Xen: It contains a virtual CPU object and a map translating guest physical memory pages into memory pages owned by the VMM. To emulate sensitive instructions upon a virtualization fault IPC, the VMM decodes the instruction and its parameters based on the program counter pointer and general-purpose register file of the guest VM, which are stored within the IPC message that was sent from L4 on behalf of the trapping VM. For convenience, L4 also passes along the *value* of the program counter, that is, the trapping instruction. In comparison to x86 processors, which have variable-sized instructions of lengths up to 15 bytes, fetching and decoding sensitive instructions on embedded PowerPC are rather trivial tasks, as instructions have a fixed size of 32 bits on PowerPC.

4.2 Virtual Physical Memory

To the user-level VMM, paging a guest with virtualized physical memory is similar to regular user-level paging in L4 systems [11]: Whenever the guest suffers a physical TLB miss, L4 sends a page fault IPC containing the faulting instruction and address and other (virtual) TLB state necessary to service the fault. In its present implementation, the VMM organizes guest-physical memory in linear segments. Thus, when handling a fault, the VMM checks whether the accessed guest-physical address is within the segment limits. If so, it responds with a mapping IPC message that causes L4 to insert the corresponding mapping into its database and into the hardware TLB.

4.3 Device Virtualization

Besides virtualization of BG/P cores, the main task of the user-level VMM is to virtualize BG/P hardware devices. L4 traps and propagates sensitive device instructions such as moves to or from system registers (`mfocr`, `mtocr`), as well as generic load/store instructions to sensitive device memory regions. It is up to the VMM to back those transactions with appropriate device models and state machines that give the guest the illusion of real devices, and to multiplex them onto actual physical hardware shared among all guests. The VMM currently provides virtual models for the BIC and for the collective and the torus network devices. Emulation of the BIC is a rather straightforward task: The VMM intercepts all accesses to the memory-mapped BIC device and emulates them using L4's mechanisms for external interrupt handling and event injection. The following paragraphs detail the emulation of the collective and torus network.

Collective Network. BG/P's collective network is an over-connected binary tree that spans the whole installation. The collective is a one-to-all medium for broadcast or reduction

operations, with support for node-specific filtering and a complex routing scheme that allows to partition the network to increase the total bandwidth. Collective link bandwidth is 6.8 Gbit/s; hardware latencies are below $6\ \mu\text{s}$ for a 72-rack system [15]. Software transmits data over the collective network via packets injected into two memory-mapped virtual channels. Each packet consists of a header and 16 128-bit data words. The packet header is written and read using general-purpose register to (device-) memory instructions, while packet data is written and read through the floating-point unit.

Our VMM provides a fully virtualized version of BG/P's collective network device. The VMM leaves device memory unmapped, so that each device register access leads to a virtualization trap. The VMM furthermore emulates device control registers (DCRs) used to configure the collective network device. The corresponding instructions (`mtocr`, `mfocr`) are sensitive and directly trap to L4 and the VMM. For emulation, the VMM provides a per-VM shadow collective network interface model that contains virtual DCRs and, per channel, virtual injection and reception FIFOs with a virtual header and 16 virtual FPU words per packet. The VMM registers itself to L4 as an interrupt handler for all collective network device interrupts, which will cause L4 to emit an interrupt IPC message to the VMM whenever the physical device fires one of its hardware interrupts.

Whenever the guest causes a packet to be sent on its virtual collective network interface, the VMM loads the corresponding virtual registers into the physical collective network device. Receiving data is slightly more complex: Whenever the VMM receives an interrupt message from L4, it reads the packet header and data from the physical device into a private buffer, and then delivers a virtual interrupt to the corresponding guest VM. Subsequent VM accesses to packet header and data are then served out of the private buffer into the VM's general-purpose or floating-point registers. Since copying packets induces substantial software overhead to the high-performance collective network path (Section 5), we are also considering optimized virtual packet handling by means of device pass-through and/or para-virtualization techniques.

Torus. BG/P's torus network is the most important data transport network with respect to bisectional bandwidth, latency, and software overhead [1]. Each compute node is part of the 3D torus network spanning the whole installation. On each node, the torus device has input and output links for each of its six neighbors, each with a bandwidth of 3.4 Gbit/s, for a total node bandwidth of 40.8 Gbit/s. Worst-case end-to-end latency in a 64 k server system is below $5\ \mu\text{s}$. Nodes act as forwarding routers without software intervention.

The torus provides two transmission interfaces, a regular buffer-based one and one based on remote direct memory access (rDMA). For the latter, the torus DMA engine provides an advanced put/get-based interface to read or write segments of memory from remote nodes, based on memory descriptors inserted into its injection and reception FIFOs. Each memory descriptor denotes a contiguous region in physical memory on the local node. To facilitate rDMA transactions (e.g. a direct-put operation copying data directly to a remote node's memory), software identifies the corresponding remote memory segments via a selector/offset

tuple that corresponds to a descriptor in the receive FIFO on the remote node (Figure 6). To transmit data via the torus, software inserts one or more packets into a torus injection FIFO, with the packet specifying the destination using its X,Y,Z coordinates. For non-DMA transfer, the payload is embedded in the packet; for DMA transfers (local as well as remote), the corresponding sender and receiver memory segment descriptors and selectors are instead appended to the packet.

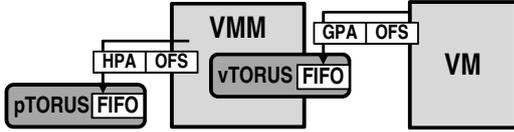


Figure 6: Virtualized Torus Interconnect. The VMM passes through guest VM descriptors, translating addresses from guest- to host-physical.

As with the collective network, our VMM provides a virtualized version of Blue Gene’s torus device, and traps and emulates all accesses to torus device registers. Again, DCR instructions directly trap into L4, while device memory accesses trap by means of invalid host TLB entries. Again, the VMM registers itself for physical torus interrupts and delivers them to the guests as needed.

However, in contrast to our virtual version of the collective network device, our virtual torus device only holds the DMA descriptor registers, but does not copy the actual data around during DMA send or receive operations. Instead, it passes along guest VM memory segment descriptors from virtual to physical FIFOs, merely validating that they reference valid guest-physical memory before translating them into host-physical addresses. As the VMM currently uses simple linear segmentation (see Section 4.2), that translation is merely an offset operation and the descriptors in physical FIFOs always reference valid guest memory.

In contrast to general-purpose virtualization environments, which typically provide virtualized Ethernet devices including virtual MAC addresses and a virtual switch architecture, or supercomputer VMM presently does not cater for extra naming or multiplexing of multiple virtual torus devices. Instead, it preserves the naming of the real world and maps virtual torus coordinates and channel identifiers idempotently to their physical counterparts. As a result, our VMM does not fully emulate the torus network, but merely provides safe and *partitioned* access to the physical torus network for individual VMs. As BG/P’s torus features four independent DMA send and receive groups, our VMM can supply up to four VMs with a different DMA channel holding all its descriptors, without having to multiplex descriptors from different VMs onto a single FIFO.

At present, our VMM intercepts all accesses to the virtual DMA groups and multiplexes them among the physical ones. However, since DMA groups are located on different physical pages and thus have different TLB entries, we plan, for future versions, to directly map torus channels into guest applications, effectively allowing them to bypass the guest OS *and hypervisor*. While such a pass-through approach will require a para-virtual torus driver (to translate guest-physical to host-physical addresses), it can still be made safe without requiring interception, as the torus DMA supports a set of range check registers that enable containment of

valid DMA addresses into the guest’s allowed allotment of physical memory.

5. INITIAL EVALUATION

Our approach is in an prototypical stage, and we did not optimize any of the frequently executed trap-and-emulate paths yet. For evaluation, we thus focused mostly on functionality rather than performance. Nevertheless, we have run some initial performance benchmarks to find out whether our approach is generally viable and where the most important bottlenecks and possibilities of optimization reside.

Guest OS Support. Our L4-based VMM generally supports running arbitrary guest OSes on BG/P, like CNK [9], or ZeptoOS [5]. With respect to the the implementation of the virtualized TLB, L4 is currently limited in that it reserves one of the two translation spaces. We have verified our implementation with Kittyhawk Linux, a BG/P version of the Linux Kernel [4] with support for BG/P’s hardware devices. It also provides an overlay that maps standard Linux Ethernet communication onto Blue Gene’s high-speed collective and torus interconnects [3]. Our VMM allows one or more instances of Kittyhawk Linux to run in a VM. Kittyhawk Linux runs unmodified, that is, the same Kittyhawk Linux binary that runs on BG/P also runs on top of our VMM. We currently support uni-processor guests only; however, as L4 itself supports multi-processing, individual guest vCPUs can be scheduled on each of the four physical cores of the Blue Gene node.

Initial Benchmark Results. For initial evaluation, we ran two experiments. In the first experiment, we compiled a small source-code project (about 1000 lines of code) under virtualized Kittyhawk Linux. We then used a debug build of L4 with an internal event tracing facility to find out frequently executed VM-related code paths. In this configuration, compilation took about 126s compared to 3s when running on native Kittyhawk Linux. Table 1 lists the results. We note that the number of IPCs – that is, the number of VM exits that involve the user-level VMM – is relatively low, meaning that L4 handles most guest exits internally. Also the experiment shows a high number of TLB misses and TLB-related instructions, indicating that virtualized memory subsystem is a bottleneck of our implementation.

Trace Point	Count	Trace Point	Count
SYSCALL_IPC	106 K	ITLB_MISS	4756 K
EXCEPT_DECR	66 K	EMUL_RFI	5021 K
EMUL_MTMSR	102 K	DTLB_MISS	6514 K
EMUL_WRTEE	117 K	EMUL_TLBWE	14745 K
EMUL_MFMSR	228 K	EMUL_MFSPR	29964 K
EMUL_WRTEEI	403 K	EMUL_MTSPR	30036 K

Table 1: Execution frequency of VMM-related L4 code paths for a compilation job in a VM.

In the second experiment, we measured Ethernet network throughput and latency between two compute nodes. Packets are delivered to the torus interconnect, by means of Kittyhawk Linux’ Ethernet driver module. For comparison, we ran the experiments for both a native and virtualized Kittyhawk Linux each running on a compute node.

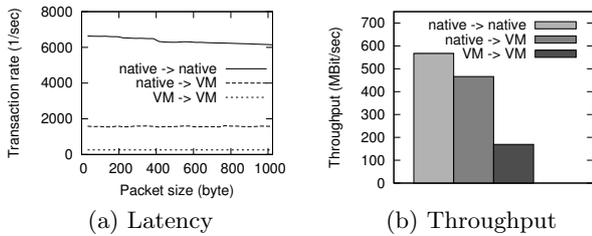


Figure 7: Performance of Torus-Routed Virtualized Networks.

For benchmarking we used `netperf`'s TCP stream test for throughput and the TCP request/response test for latency measurements [28]. Figure 7 shows the results. Our virtualization layer poses significant overhead to the Ethernet network performance (which is already off the actual performance that the torus and collective hardware can deliver [3]). We are confident, however, that optimizations such as para-virtual device support can render virtualization substantially more efficient. A recent study reports that VMware's engineers faced, and eventually addressed, similarly dismal performance with prototypical versions of their VMM [2]. We conclude from the preliminary evaluation and our own experiences, that our L4-based VMM is a promising approach for successfully deploying, running, and using virtualized OSes and applications on BG/P, warranting further exploration and optimization.

6. RELATED WORK

There exists a plethora of VMM efforts, including Xen [29], VMware [2], and, directly related, micro-hypervisor-based systems [11, 12, 30]; those approaches mostly address the embedded or server rather than HPC space. The studies in [7, 25] identified virtualization as a system-level alternative to address the development challenges of HPC systems. Gavrilowska et al. explored virtualized HPC for x86/InfiniBand-based hardware [8]; P.R.O.S.E. explored a partitioning hypervisor architecture for PowerPC and x86 based HPC systems [32]. However, both approaches focus mostly on hypervisor infrastructure rather than on decomposed OS designs or on support for native applications. Finally, there exists a port of the KVM monitor to PowerPC Book E cores [18]; although designed towards embedded systems, it shares some implementation details with our PowerPC version of L4 and the VMM.

Arguably the most related effort to our approach is Palacios and Kitten [21], a light-weight kernel/VMM combination striving to achieve high performance and scalability on HPC machines. Palacios runs as a module extension within Kitten. Like L4, Kitten also provides a native environment that can be used to develop customizations. The most notable differences are: that Palacios and Kitten are being developed for x86-based HPC systems rather than for a highly-specialized Supercomputer platform such as BG/P; that Palacios runs as a kernel module in Kitten's privileged domain, whereas our VMM runs completely decomposed and deprived as user-level process; and that their approach to device virtualization introduces a scheme relying on guest cooperation in order to achieving high-performance virtualization of HPC network devices [20], whereas our ap-

proach currently fully intercepts guest device accesses (although it could be adapted to support a similar scheme).

Examples of traditional light-weight kernel approaches for Supercomputers are CNK [9, 27] and Sandia's Catamount [16]. Our approach strives to enhance such light-weight kernel approaches in that it provides the ability to virtualize guest OSes. Finally, research has explored whether a more fully-fledged OS such as Plan9 [26] or Linux [4, 5, 22] may be a more viable alternative than light-weight kernels for Supercomputers. Our approach complements those efforts with the alternative idea of a decomposed OS architecture with support for virtualization.

7. CONCLUSION

In this paper, we have presented a virtualization-capable μ -kernel OS architecture for BG/P. Our architecture consists of a virtualization-capable μ -kernel and a user-level VMM component running atop. The μ -kernel also supports running applications natively. Our L4-based prototype successfully virtualizes Kittyhawk Linux with support for virtualized collective and torus network devices. First experiences and experiments show that our VMM still takes a substantial performance hit. However, we believe that, pending optimizations, our approach poses an interesting OS alternative for Supercomputers, providing the convenience of a fully-featured commodity OS software stack, while also promising to satisfy the need for low latency and scalability of a HPC system. This research used resources of the Argonne Leadership Computing Facility at Argonne National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under contract DE-AC02-06CH11357. The infrastructure presented herein is part of the the projects L4Ka and Kittyhawk. Software is open source and can be downloaded from l4ka.org and kittyhawk.bu.edu.

8. REFERENCES

- [1] N. R. Adiga, M. A. Blumrich, D. Chen, P. Coteus, A. Gara, M. E. Giampapa, P. Heidelberger, S. Singh, B. D. Steinmacher-Burow, T. Takken, M. Tsao, and P. Vranas. Blue Gene/L torus interconnection network. *IBM Journal of Research and Development*, 49(2/3):265–276, June 2005.
- [2] O. Agesen, A. Garthwaite, J. Sheldon, and P. Subrahmanyam. The evolution of an x86 virtual machine monitor. *ACM Operating Systems Review*, 44(4):3–18, December 2010.
- [3] J. Appavoo, V. Uhlig, J. Stoess, A. Waterland, B. Rosenburg, R. Wisniewski, D. D. Silva, E. van Hensbergen, and U. Steinberg. Providing a cloud network infrastructure on a supercomputer. In *Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing*, pages 385–394, Chicago, IL, USA, June 2010.
- [4] J. Appavoo, V. Uhlig, A. Waterland, B. Rosenburg, D. D. Silva, and J. E. Moreira. Kittyhawk: Enabling cooperation and competition in a global, shared computational system. *IBM Journal of Research and Development*, 53(4):1–15, July 2009.
- [5] P. Beckman, K. Iskra, K. Yoshii, S. Coghlan, and A. Nataraj. Benchmarking the effects of operating

- system interference on extreme-scale parallel machines. *Cluster Computing*, 11(1):3–16, Jan. 2008.
- [6] R. Bhargava, B. Serebrin, F. Spadini, and S. Manne. Accelerating two-dimensional page walks for virtualized systems. In *Proceedings of the 13th International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 48–59, Seattle, WA, USA, Mar. 2008.
- [7] C. Engelmann, S. Scott, H. Ong, G. Vallée, and T. Naughton. Configurable virtualized system environments for high performance computing. In *1st Workshop on System-level Virtualization for High Performance Computing*, Lisbon, Portugal, Mar. 2007.
- [8] A. Gavrilovska, S. Kumar, H. Raj, K. Schwan, V. Gupta, R. Nathuji, R. Niranjana, A. Ranadive, and P. Saraiya. High-performance hypervisor architectures: Virtualization in HPC systems. In *1st Workshop on System-level Virtualization for High Performance Computing*, Lisbon, Portugal, Mar. 2007.
- [9] M. Giampapa, T. Gooding, T. Inglett, and R. W. Wisniewski. Experiences with a lightweight supercomputer kernel: Lessons learned from Blue Gene’s CNK. In *Proceedings of the 2010 International Conference on Supercomputing*, pages 1–10, New Orleans, LA, USA, Nov. 2010.
- [10] H. Härtig, M. Hohmuth, N. Feske, C. Helmuth, A. Lackorzynski, F. Mehnert, and M. Peter. The Nizza secure-system architecture. In *Proceedings the 1st International Conference on Collaborative Computing: Networking, Applications and Worksharing*, San Jose, CA, USA, Dec. 2005.
- [11] H. Härtig, M. Hohmuth, J. Liedtke, and S. Schönberg. The performance of μ -kernel based systems. In *Proceedings of the 16th Symposium on Operating System Principles*, pages 66–77, Saint Malo, France, Oct. 1997.
- [12] G. Heiser and B. Leslie. The OKL4 microvisor: convergence point of microkernels and hypervisors. In *Proceedings of the first ACM asia-pacific workshop on Workshop on systems*, pages 19–24, Aug. 2010.
- [13] IBM. *IBM Power ISA Version 2.03*. IBM Corporation, 2006.
- [14] IBM. *IBM Power ISA Version 2.06*. IBM Corporation, 2009.
- [15] IBM Blue Gene team. Overview of the IBM Blue Gene/P Project. *IBM Journal of Research and Development*, 52(1/2):199–220, Jan./Mar. 2008.
- [16] S. M. Kelly and R. Brightwell. Software architecture of the light weight kernel, Catamount. In *Proceedings of the 2005 Cray User Group Annual Technical Conference*, Albuquerque, NM, USA, May 2005.
- [17] I. Kuz, Y. Liu, I. Gorton, and G. Heiser. Camkes: A component model for secure microkernel-based embedded systems. *Journal of Systems and Software*, 80(5):687 – 699, May 2007.
- [18] KVM Team. KVM for PowerPC. <http://www.linux-kvm.org/page/PowerPC/>.
- [19] L4 Development Team. *L4 X.2 Reference Manual*. University of Karlsruhe, Germany, May 2009.
- [20] J. R. Lange, K. Pedretti, P. Dinda, P. G. Bridges, C. Bae, P. Soltero, and A. Merritt. Minimal-overhead virtualization of a large scale supercomputer. In *Proceedings of the 7th ACM SIGPLAN/SIGOPS international conference on Virtual execution environments*, pages 169–180, Newport Beach, California, USA, Mar. 2011.
- [21] J. R. Lange, K. T. Pedretti, T. Hudson, P. A. Dinda, Z. Cui, L. Xia, P. G. Bridges, A. Gocke, S. Jaconette, M. Levenhagen, and R. Brightwell. Palacios and Kitten: New high performance operating systems for scalable virtualized and native supercomputing. In *Proceedings of the 24th IEEE International Symposium on Parallel and Distributed Processing*, pages 1–12, Atlanta, GA, USA, Apr. 2010.
- [22] Laurence S. Kaplan. Lightweight Linux for High-performance Computing. [LinuxWorld.com, http://www.linuxworld.com/news/2006/120406-lightweight-linux.html](http://www.linuxworld.com/news/2006/120406-lightweight-linux.html), 2006.
- [23] J. LeVasseur, V. Uhlig, J. Stoess, and S. Götz. Unmodified device driver reuse and improved system dependability via virtual machines. In *Proceedings of the 6th Symposium on Operating Systems Design and Implementation*, pages 17–30, San Fransisco, CA, USA, Dec. 2004.
- [24] J. Liedtke. On μ -Kernel construction. In *Proceedings of the 15th Symposium on Operating System Principles*, pages 237–250, Copper Mountain, CO, USA, Dec. 1995.
- [25] M. F. Mergen, V. Uhlig, O. Krieger, and J. Xenidis. Virtualization for high-performance computing. *ACM Operating Systems Review*, 40:8–11, April 2006.
- [26] R. Minnich and J. McKie. Experiences porting the Plan 9 research operating system to the IBM Blue Gene supercomputers. *Computer Science - Research and Development*, 23:117–124, May 2009.
- [27] J. Moreira, M. Brutman, J. Castanos, T. Gooding, T. Inglett, D. Lieber, P. McCarthy, M. Mundy, J. Parker, B. Wallenfelt, M. Giampapa, T. Engelsiepen, and R. Haskin. Designing a highly-scalable operating system: The Blue Gene/L story. In *Proceedings of the 2006 International Conference on Supercomputing*, pages 53–63, Tampa, FL, USA, Nov. 2006.
- [28] Netperf Team. Netperf. <http://www.netperf.org/netperf/>.
- [29] I. Pratt, K. Fraser, S. Hand, C. Limpach, A. Warfield, D. Magenheimer, J. Nakajima, and A. Malick. Xen 3.0 and the art of virtualization. In *Proceedings of the 2005 Ottawa Linux Symposium*, pages 65–78, Ottawa, Canada, July 2005.
- [30] U. Steinberg and B. Kauer. NOVA: A microhypervisor-based secure virtualization architecture. In *Proceedings of the 5th ACM SIGOPS EuroSys conference*, pages 209–221, Paris, France, Apr. 2010.
- [31] R. Uhlig, G. Neiger, D. Rodgers, A. L. Santoni, F. C. M. Martins, A. V. Anderson, S. M. Bennett, A. Kägi, F. H. Leung, and L. Smith. Intel virtualization technology. *IEEE Computer*, 38(5):48–56, May 2005.
- [32] E. Van Hensbergen. PROSE: partitioned reliable operating system environment. *ACM Operating Systems Review*, 40(2):12–15, 2006.