

Fast VMM-based overlay networking for bridging the cloud and high performance computing

Lei Xia · Zheng Cui · John Lange · Yuan Tang · Peter Dinda · Patrick Bridges

Received: 1 October 2012 / Accepted: 1 May 2013
© Springer Science+Business Media New York 2013

Abstract A collection of virtual machines (VMs) interconnected with an overlay network with a layer 2 abstraction has proven to be a powerful, unifying abstraction for adaptive distributed and parallel computing on loosely-coupled environments. It is now feasible to allow VMs hosting high performance computing (HPC) applications to seamlessly bridge distributed cloud resources and *tightly-coupled* supercomputing and cluster resources. However, to achieve the application performance that the tightly-coupled resources are capable of, it is important that the overlay network not introduce significant overhead relative to the native hardware, which is not the case for current user-level tools, including our own existing VNET/U system. In response, we describe the design, implementation, and evaluation of a virtual networking system that has negligible latency and bandwidth overheads in 1–10 Gbps networks. Our sys-

tem, VNET/P, is directly embedded into our publicly available Palacios virtual machine monitor (VMM). VNET/P achieves native performance on 1 Gbps Ethernet networks and very high performance on 10 Gbps Ethernet networks. The NAS benchmarks generally achieve over 95 % of their native performance on both 1 and 10 Gbps. We have further demonstrated that VNET/P can operate successfully over more specialized tightly-coupled networks, such as Infiniband and Cray Gemini. Our results suggest it is feasible to extend a software-based overlay network designed for computing at wide-area scales into tightly-coupled environments.

Keywords Overlay networks · Virtualization · HPC · Scalability

L. Xia (✉) · P. Dinda
Northwestern University, Evanston, IL, USA
e-mail: lxia@northwestern.edu

P. Dinda
e-mail: pdinda@northwestern.edu

Z. Cui · P. Bridges
University of New Mexico, Albuquerque, NM, USA

Z. Cui
e-mail: cui Zheng@cs.unm.edu

P. Bridges
e-mail: bridges@cs.unm.edu

J. Lange
University of Pittsburgh, Pittsburgh, PA, USA

Y. Tang
University of Electronic Science and Technology of China,
Chengdu, China
e-mail: y tang@uestc.edu.cn

1 Introduction

Cloud computing in the “infrastructure as a service” (IaaS) model has the potential to provide economical and effective on-demand resources for high performance computing. In this model, an application is mapped into a collection of virtual machines (VMs) that are instantiated as needed, and at the scale needed. Indeed, for loosely-coupled applications, this concept has readily moved from research [8, 44] to practice [39]. As we describe in Sect. 3, such systems can also be adaptive, autonomously selecting appropriate mappings of virtual components to physical components to maximize application performance or other objectives. However, *tightly-coupled* scalable high performance computing (HPC) applications currently remain the purview of resources such as clusters and supercomputers. We seek to extend the adaptive IaaS cloud computing model into these regimes, allowing an application to dynamically span both kinds of environments.

The current limitation of cloud computing systems to *loosely-coupled* applications is not due to machine virtualization limitations. Current virtual machine monitors (VMMs) and other virtualization mechanisms present negligible overhead for CPU and memory intensive workloads [18, 37]. With VMM-bypass [34] or self-virtualizing devices [41] the overhead for direct access to network devices can also be made negligible.

Considerable effort has also gone into achieving low-overhead network virtualization and traffic segregation within an individual data center through extensions or changes to the network hardware layer [12, 25, 38]. While these tools strive to provide uniform performance across a cloud data center (a critical feature for many HPC applications), they do not provide the same features once an application has migrated outside the local data center, or spans multiple data centers, or involves HPC resources. Furthermore, they lack compatibility with the more specialized interconnects present on most HPC systems.

Beyond the need to support our envisioned computing model across today's and tomorrow's tightly-coupled HPC environments, we note that data center network design and cluster/supercomputer network design seem to be converging [1, 13]. This suggests that future data centers deployed for general purpose cloud computing will become an increasingly better fit for tightly-coupled parallel applications, and therefore such environments could potentially also benefit.

The current limiting factor in the adaptive cloud- and HPC-spanning model described above for tightly-coupled applications is the performance of the virtual networking system. Current adaptive cloud computing systems use software-based overlay networks to carry inter-VM traffic. For example, our VNET/U system, which is described in more detail later, combines a simple networking abstraction within the VMs with location-independence, hardware-independence, and traffic control. Specifically, it exposes a layer 2 abstraction that lets the user treat his VMs as being on a simple LAN, while allowing the VMs to be migrated seamlessly across resources by routing their traffic through the overlay. By controlling the overlay, the cloud provider or adaptation agent can control the bandwidth and the paths between VMs over which traffic flows. Such systems [43, 49] and others that expose different abstractions to the VMs [56] have been under continuous research and development for several years. Current virtual networking systems have sufficiently low overhead to effectively host loosely-coupled scalable applications [7], but their performance is insufficient for tightly-coupled applications [40].

In response to this limitation, we have designed, implemented, and evaluated VNET/P, which shares its model and vision with VNET/U, but is designed to achieve near-native performance in the 1 Gbps and 10 Gbps switched networks

common in clusters today, as well as to operate effectively on top of even faster networks, such as Infiniband and Cray Gemini. VNET/U and our model is presented in more detail in Sect. 3.

VNET/P is implemented in the context of our publicly available, open source Palacios VMM [30], which is in part designed to support virtualized supercomputing. A detailed description of VNET/P's design and implementation is given in Sect. 4. As a part of Palacios, VNET/P is publicly available. VNET/P could be implemented in other VMMs, and as such provides a proof-of-concept that overlay-based virtual networking for VMs, with performance overheads low enough to be inconsequential even in a tightly-coupled computing environment, is clearly possible.

The performance evaluation of VNET/P (Sect. 5) shows that it is able to achieve native bandwidth on 1 Gbps Ethernet with a small increase in latency, and very high bandwidth on 10 Gbps Ethernet with a similar, small latency increase. On 10 Gbps hardware, the kernel-level VNET/P system provides on average 10 times more bandwidth and 7 times less latency than the user-level VNET/U system can.

In a related paper from our group [5], we describe additional techniques, specifically optimistic interrupts and cut-through forwarding, that bring bandwidth to near-native levels for 10 Gbps Ethernet. Latency increases are predominantly due to the lack of selective interrupt exiting in the current AMD and Intel hardware virtualization extensions. We expect that latency overheads will be largely ameliorated once such functionality become available, or, alternatively, when software approaches such as ELI [11] are used.

Although our core performance evaluation of VNET/P is on 10 Gbps Ethernet, VNET/P can run on top of any device that provides an IP or Ethernet abstraction within the Linux kernel. The portability of VNET/P is also important to consider, as the model we describe above would require it to run on many different hosts. In Sect. 6 we report on preliminary tests of VNET/P running over Infiniband via the IPoIB functionality, and on the Cray Gemini via the IPoG virtual Ethernet interface. Running on these platforms requires few changes to VNET/P, but creates considerable flexibility. In particular, using VNET/P, existing, unmodified VMs running guest OSES with commonplace network stacks can seamlessly run on top of such diverse hardware. To the guest, a complex network of commodity and high-end networks looks like a simple Ethernet network. Also in Sect. 6, we describe a version of VNET/P that has been designed for use with the Kitten lightweight kernel as its host OS. Kitten is quite different from Linux—indeed the combination of Palacios and Kitten is akin to a “type-I” (unhosted) VMM—resulting in a different VNET/P architecture. This system and its performance provide evidence that the VNET/P model can be successfully brought to different host/VMM environments.

Our contributions are as follows:

- We articulate the benefits of extending virtual networking for VMs down to clusters and supercomputers with high performance networks. These benefits are also applicable to data centers that support IaaS cloud computing.
- We describe the design and implementation of a virtual networking system, VNET/P, that does so. The design could be applied to other VMMs and virtual network systems.
- We perform an extensive evaluation of VNET/P on 1 and 10 Gbps Ethernet networks, finding that it provides performance with negligible overheads on the former, and manageable overheads on the latter. VNET/P generally has little impact on performance for the NAS benchmarks.
- We describe our experiences with running the VNET/P implementation on Infiniband and Cray Gemini networks. VNET/P allows guests with commodity software stacks to leverage these networks.
- We describe the design, implementation, and evaluation of a version of VNET/P for lightweight kernel hosts, particularly the Kitten LWK.

Through the use of low-overhead overlay-based virtual networking in high-bandwidth, low-latency environments such as current clusters and supercomputers, and future data centers, we seek to make it practical to use virtual networking at all times, even when running tightly-coupled applications on such high-end environments. This would allow us to seamlessly and *practically* extend the already highly effective adaptive virtualization-based IaaS cloud computing model to such environments.

This paper is an extended version of a previous conference publication [57]. Compared to the conference paper, it provides an extended presentation of the core design aspects of VNET/P as well as descriptions of implementations of VNET/P for Infiniband and Cray Gemini. It also includes initial performance evaluations on these platforms.

2 Related work

VNET/P is related to NIC virtualization, overlays, and virtual networks, as we describe below.

NIC virtualization There is a wide range of work on providing VMs with fast access to networking hardware, where no overlay is involved. For example, VMware and Xen support either an emulated register-level interface [47] or a paravirtualized interface to guest operating system [36]. While purely software-based virtualized network interface has high overhead, many techniques have been proposed to support simultaneous, direct-access network I/O. For example, some work [34, 41] has demonstrated the use of self-virtualized network hardware that allows direct guest access, thus provides high performance to untrusted guests. Willmann et al

have developed a software approach that also supports concurrent, direct network access by untrusted guest operating systems [45]. In addition, VPPIO [59] can be applied on network virtualization to allow virtual passthrough I/O on non-self-virtualized hardware. Virtual WiFi [58] is an approach to provide the guest with access to wireless networks, including functionality specific to wireless NICs. In contrast with such work, VNET/P provides fast access to an overlay network, which includes encapsulation and routing. It makes a set of VMs appear to be on the same local Ethernet regardless of their location anywhere in the world and their underlying hardware. Our work shows that this capability can be achieved without significantly compromising performance when the VMs happen to be very close together.

Overlay networks: Overlay networks implement extended network functionality on top of physical infrastructure, for example to provide resilient routing (e.g. [3]), multicast (e.g. [17]), and distributed data structures (e.g., [46]) without any cooperation from the network core; overlay networks use end-systems to provide their functionality. VNET is an example of a specific class of overlay networks, namely virtual networks, discussed next.

Virtual networking: Virtual networking systems provide a service model that is compatible with an existing layer 2 or 3 networking standard. Examples include VIOLIN [21], ViNe [52], VINI [4], SoftUDC VNET [23], OCALA [22], WoW [10], and the emerging VXLAN standard [35]. Like VNET, VIOLIN, SoftUDC, WoW, and VXLAN are specifically designed for use with virtual machines. Of these, VIOLIN is closest to VNET (and contemporaneous with VNET/U), in that it allows for the dynamic setup of an arbitrary private layer 2 and layer 3 virtual network among VMs. The key contribution of VNET/P is to show that this model can be made to work with minimal overhead even in extremely low latency, high bandwidth environments.

Connections: VNET/P could itself leverage some of the related work described above. For example, effective NIC virtualization might allow us to push encapsulation directly into the guest, or to accelerate encapsulation via a split scatter/gather map. Mapping unencapsulated links to VLANs would enhance performance on environments that support them. There are many options for implementing virtual networking and the appropriate choice depends on the hardware and network policies of the target environment. In VNET/P, we make the choice of minimizing these dependencies.

3 VNET model and VNET/U

The VNET model was originally designed to support adaptive computing on distributed virtualized computing resources within the Virtuoso system [6], and in particular to support the adaptive execution of a distributed or parallel

computation executing in a collection of VMs potentially spread across multiple providers or supercomputing sites. The key requirements, which also hold for the present paper, were as follows.

- VNET would make within-VM network configuration the sole responsibility of the VM owner.
- VNET would provide location independence to VMs, allowing them to be migrated between networks and from site to site, while maintaining their connectivity, without requiring any within-VM configuration changes.
- VNET would provide hardware independence to VMs, allowing them to use diverse networking hardware without requiring the installation of specialized software.
- VNET would provide minimal overhead, compared to native networking, in the contexts in which it is used.

The VNET model meets these requirements by carrying the user's VMs' traffic via a configurable overlay network. The overlay presents a simple layer 2 networking abstraction: a user's VMs appear to be attached to the user's local area Ethernet network, regardless of their actual locations or the complexity of the VNET topology/properties. Further information about the model can be found elsewhere [49].

The VNET overlay is dynamically reconfigurable, and can act as a locus of activity for an adaptive system such as Virtuoso. Focusing on parallel and distributed applications running in loosely-coupled virtualized distributed environments (e.g., "IaaS Clouds"), we demonstrated that the VNET "layer" can be effectively used to:

1. monitor application communication and computation behavior [14, 15]),
2. monitor underlying network behavior [16],
3. formulate performance optimization problems [48, 51], and
4. address such problems through VM migration and overlay network control [50], scheduling [32, 33], network reservations [31], and network service interposition [27].

These and other features that can be implemented within the VNET model have only marginal utility if carrying traffic via the VNET overlay has significant overhead compared to the underlying native network.

The VNET/P system described in this paper is compatible with, and compared to, our previous VNET implementation, VNET/U. Both support a dynamically configurable general overlay topology with dynamically configurable routing on a per MAC address basis. The topology and routing configuration is subject to global or distributed control (for example, by the VADAPT [50] part of Virtuoso). The overlay carries Ethernet packets encapsulated in UDP packets, TCP streams with and without SSL encryption, TOR privacy-preserving streams, and others. Because Ethernet packets are used, the VNET abstraction can also

easily interface directly with most commodity network devices, including virtual NICs exposed by VMMs in the host, and with fast virtual devices (e.g., Linux virtio network devices) in guests.

While VNET/P is implemented within the VMM, VNET/U is implemented as a user-level system. As a user-level system, it readily interfaces with VMMs such as VMware Server and Xen, and requires no host changes to be used, making it very easy for a provider to bring it up on a new machine. Further, it is easy to bring up VNET daemons when and where needed to act as proxies or way-points. A VNET daemon has a control port which speaks a control language for dynamic configuration. A collection of tools allows for the wholesale construction and teardown of VNET topologies, as well as dynamic adaptation of the topology and forwarding rules to the observed traffic and conditions on the underlying network.

The last reported measurement of VNET/U showed it achieving 21.5 MB/s (172 Mbps) with a 1 ms latency overhead communicating between Linux 2.6 VMs running in VMware Server GSX 2.5 on machines with dual 2.0 GHz Xeon processors [27]. A current measurement, described in Sect. 5, shows 71 MB/s with a 0.88 ms latency. VNET/U's speeds are sufficient for its purpose in providing virtual networking for wide-area and/or loosely-coupled distributed computing. They are not, however, sufficient for use within a cluster at gigabit or greater speeds. Making this basic VM-to-VM path competitive with hardware is the focus of this paper. VNET/U is fundamentally limited by the kernel/user space transitions needed to handle a guest's packet send or receive. In VNET/P, we move VNET directly into the VMM to avoid such transitions.

4 Design and implementation

We now describe how VNET/P has been architected and implemented in the context of Palacios as embedded in a Linux host. Section 6.3 describes how VNET/P is implemented in the context of a Kitten embedding. The nature of the embedding affects VNET/P primarily in how it interfaces to the underlying networking hardware and networking stack. In the Linux embedding, this interface is accomplished directly in the Linux kernel. In the Kitten embedding, the interface is done via a service VM.

4.1 Palacios VMM

VNET/P is implemented in the context of our Palacios VMM. Palacios is an OS-independent, open source, BSD-licensed, publicly available embeddable VMM designed as part of the V3VEE project (<http://v3vee.org>). The V3VEE project is a collaborative community resource development

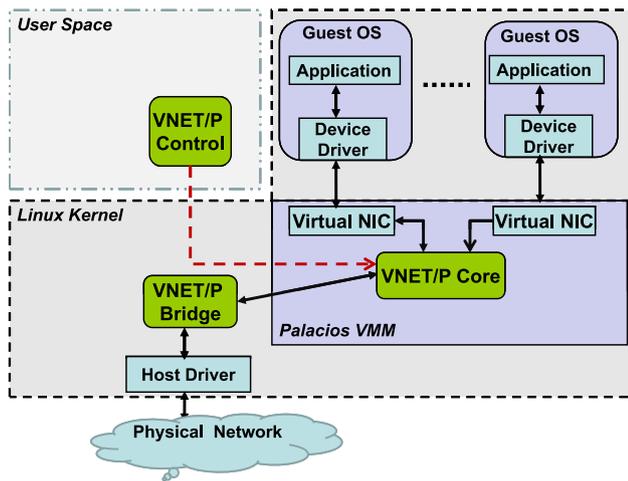


Fig. 1 VNET/P architecture

project involving Northwestern University, the University of New Mexico, Sandia National Labs, and Oak Ridge National Lab. Detailed information about Palacios can be found elsewhere [28, 30]. Palacios is capable of virtualizing large scale (4096+ nodes) with < 5 % overheads [29]. Palacios’s OS-agnostic design allows it to be embedded into a wide range of different OS architectures.

The Palacios implementation is built on the virtualization extensions deployed in current generation x86 processors, specifically AMD’s SVM [2] and Intel’s VT [53]. Palacios supports both 32 and 64 bit host and guest environments, both shadow and nested paging models, and a significant set of devices that comprise the PC platform. Due to the ubiquity of the x86 architecture Palacios is capable of operating across many classes of machines. Palacios has successfully virtualized commodity desktops and servers, high end Infiniband clusters, and Cray XT and XK supercomputers.

4.2 Architecture

Figure 1 shows the overall architecture of VNET/P, and illustrates the operation of VNET/P in the context of the Palacios VMM embedded in a Linux host. In this architecture, *guests* run in *application VMs*. Off-the-shelf guests are fully supported. Each application VM provides a virtual (Ethernet) NIC to its guest. For high performance applications, as in this paper, the virtual NIC conforms to the virtio interface, but several virtual NICs with hardware interfaces are also available in Palacios. The virtual NIC conveys Ethernet packets between the application VM and the Palacios VMM. Using the virtio virtual NIC, one or more packets can be conveyed from an application VM to Palacios with a single VM exit, and from Palacios to the application VM with a single VM exit+entry.

The *VNET/P core* is the component of VNET/P that is directly embedded into the Palacios VMM. It is responsible for routing Ethernet packets between virtual NICs on

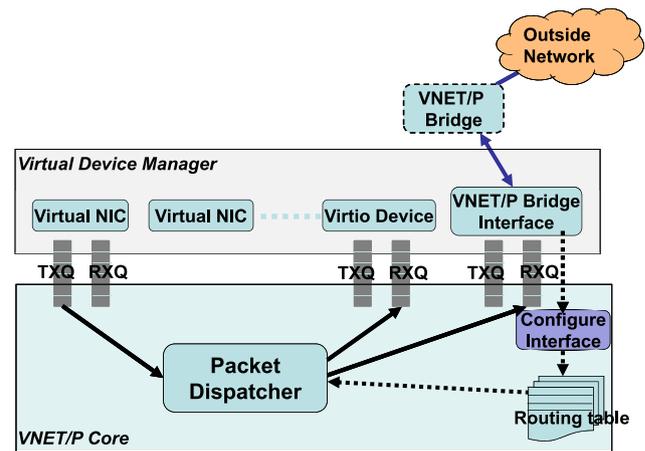


Fig. 2 VNET/P core’s internal logic

the machine and between this machine and remote VNET on other machines. The VNET/P core’s routing rules are dynamically configurable, through the control interface by the utilities that can be run in user space.

The VNET/P core also provides an expanded interface that the control utilities can use to configure and manage VNET/P. The *VNET/P control* component uses this interface to do so. It in turn acts as a daemon that exposes a TCP control port that uses the same configuration language as VNET/U. Between compatible encapsulation and compatible control, the intent is that VNET/P and VNET/U be interoperable, with VNET/P providing the “fast path”.

To exchange packets with a remote machine, the VNET/P core uses a *VNET/P bridge* to communicate with the physical network. The VNET/P bridge runs as a kernel module in the host kernel and uses the host’s networking facilities to interact with physical network devices and with the host’s networking stack. An additional responsibility of the bridge is to provide encapsulation. For performance reasons, we use UDP encapsulation, in a form compatible with that used in VNET/U. TCP encapsulation is also supported. The bridge selectively performs UDP or TCP encapsulation for packets destined for remote machines, but can also deliver an Ethernet packet without encapsulation. In our performance evaluation, we consider only encapsulated traffic.

The VNET/P core consists of approximately 2500 lines of C in Palacios, while the VNET/P bridge consists of about 2000 lines of C comprising a Linux kernel module. VNET/P is available via the V3VEE project’s public git repository, as part of the “devel” branch of the Palacios VMM.

4.3 VNET/P core

The VNET/P core is primarily responsible for routing and dispatching raw Ethernet packets. It intercepts all Ethernet packets from virtual NICs that are associated with VNET/P, and forwards them either to VMs on the same host machine

or to the outside network through the VNET/P bridge. Each packet is routed based on its source and destination MAC addresses. The internal processing logic of the VNET/P core is illustrated in Fig. 2.

Routing: To route Ethernet packets, VNET/P maintains routing tables indexed by source and destination MAC addresses. Although this table structure only provides linear time lookups, a hash table-based routing cache is layered on top of the table, and the common case is for lookups to hit in the cache and thus be serviced in constant time.

A routing table entry maps to a destination, which is either a *link* or an *interface*. A link is an overlay destination—it is the next UDP/IP-level (i.e., IP address and port) destination of the packet, on some other machine. A special link corresponds to the local network. The local network destination is usually used at the “exit/entry point” where the VNET overlay is attached to the user’s physical LAN. A packet routed via a link is delivered to another VNET/P core, a VNET/U daemon, or the local network. An interface is a local destination for the packet, corresponding to some virtual NIC.

For an interface destination, the VNET/P core directly delivers the packet to the relevant virtual NIC. For a link destination, it injects the packet into the VNET/P bridge along with the destination link identifier. The VNET/P bridge demultiplexes based on the link and either encapsulates the packet and sends it via the corresponding UDP or TCP socket, or sends it directly as a raw packet to the local network.

Packet processing: Packet forwarding in the VNET/P core is conducted by *packet dispatchers*. A packet dispatcher interacts with each virtual NIC to forward packets in one of two modes: *guest-driven mode* or *VMM-driven mode*. The operation of these modes is illustrated in the top two timelines in Fig. 3.

The purpose of guest-driven mode is to minimize latency for small messages in a parallel application. For example, a barrier operation would be best served with guest-driven mode. In the guest-driven mode, the packet dispatcher is invoked when the guest’s interaction with the NIC explicitly causes an exit. For example, the guest might queue a packet on its virtual NIC and then cause an exit to notify the VMM that a packet is ready. In guest-driven mode, a packet dispatcher runs at this point. Similarly, on receive, a packet dispatcher queues the packet to the device and then immediately notifies the device.

The purpose of VMM-driven mode is to maximize throughput for bulk data transfer in a parallel application. Unlike guest-driven mode, VMM-driven mode tries to handle multiple packets per VM exit. It does this by having VMM poll the virtual NIC. The NIC is polled in two ways. First, it is polled, and a packet dispatcher is run, if needed, in the context of the current VM exit (which is unrelated to

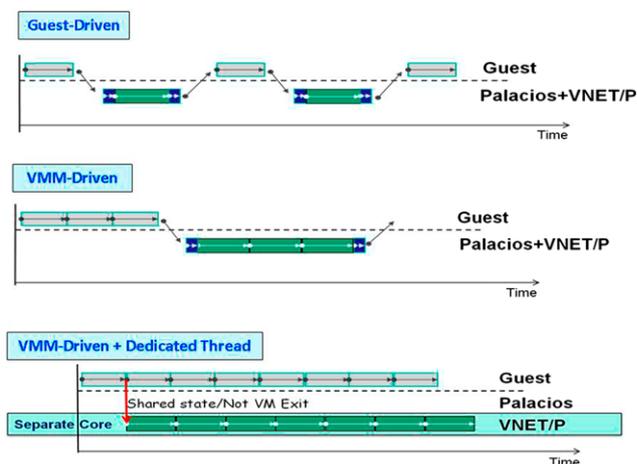


Fig. 3 The VMM-driven and guest-driven modes in the virtual NIC. Guest-driven mode can decrease latency for small messages, while VMM-driven mode can increase throughput for large messages. Combining VMM-driven mode with a dedicated packet dispatcher thread results in most send-related exits caused by the virtual NIC being eliminated, thus further enhancing throughput

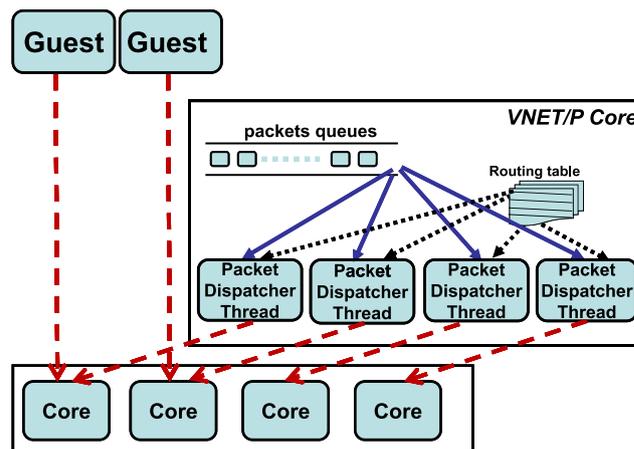


Fig. 4 VNET/P running on a multicore system. The selection of how many, and which cores to use for packet dispatcher threads is made dynamically

the NIC). Even if exits are infrequent, the polling and dispatch will still make progress during the handling of timer interrupt exits.

The second manner in which the NIC can be polled is in the context of a packet dispatcher running in a kernel thread inside the VMM context, as shown in Fig. 4. The packet dispatcher thread can be instantiated multiple times, with these threads running on different cores in the machine. If a packet dispatcher thread decides that a virtual NIC queue is full, it forces the NIC’s VM to handle it by doing a cross-core IPI to force the core on which the VM is running to exit. The exit handler then does the needed event injection. Using this approach, it is possible to dynamically employ idle processor cores to increase packet forwarding bandwidth. Combined

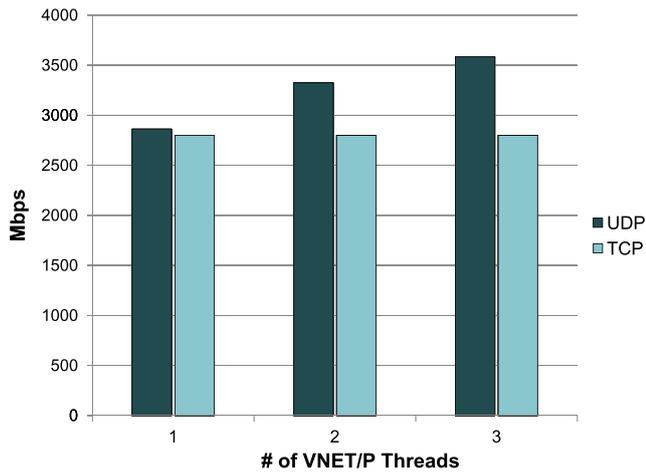


Fig. 5 Early example of scaling of receive throughput by executing the VMM-based components of VNET/P on separate cores, and scaling the number of cores used. The ultimate on-wire MTU here is 1500 bytes

with VMM-driven mode, VNET/P packet processing can then proceed in parallel with guest packet sends, as shown in the bottommost timeline of Fig. 3.

Influenced by Sidecore [26], an additional optimization we developed was to offload in-VMM VNET/P processing, beyond packet dispatch, to an unused core or cores, thus making it possible for the guest VM to have full use of its cores (minus the exit/entry costs when packets are actually handed to/from it). Figure 5 is an example of the benefits of doing so for small MTU communication.

VNET/P can be configured to statically use either VMM-driven or guest-driven mode, or *adaptive operation* can be selected. In adaptive operation, which is illustrated in Fig. 6, VNET/P switches between these two modes dynamically depending on the arrival rate of packets destined to or from the virtual NIC. For a low rate, it enables guest-driven mode to reduce the single packet latency. On the other hand, with a high arrival rate it switches to VMM-driven mode to increase throughput. Specifically, the VMM detects whether the system is experiencing a high exit rate due to virtual NIC accesses. It recalculates the rate periodically. The algorithm employs simple hysteresis, with the rate bound for switching from guest-driven to VMM-driven mode being larger than the rate bound for switching back. This avoids rapid switching back and forth when the rate falls between these bounds.

For a 1 Gbps network, guest-driven mode is sufficient to allow VNET/P to achieve the full native throughput. On a 10 Gbps network, VMM-driven mode is essential to move packets through the VNET/P core with near-native throughput. Generally, adaptive operation will achieve these limits.

```

num_packets = {number of exits caused by virtual NIC
               from last period of time window}
rate = num_packets/window_time

if (rate > upper_rate_limit && current-mode == Guest-driven)
do:
    switch_mode(VMM-Driven)
    current-mode= VMM-Driven
end
else if (rate < lower_rate_limit && current-mode == VMM-driven)
do:
    switch_mode(Guest-driven)
    current-mode= Guest-driven
end
else
    do-nothing
endif

```

Fig. 6 Adaptive mode dynamically selects between VMM-driven and guest-driven modes of operation to optimize for both throughput and latency

4.4 Virtual NICs

VNET/P is designed to be able to support any virtual Ethernet NIC device. A virtual NIC must, however, register itself with VNET/P before it can be used. This is done during the initialization of the virtual NIC at VM configuration time. The registration provides additional callback functions for packet transmission, transmit queue polling, and packet reception. These functions essentially allow the NIC to use VNET/P as its backend, instead of using an actual hardware device driver backend.

Linux virtio virtual NIC: Virtio [42], which was recently developed for the Linux kernel, provides an efficient abstraction for VMMs. A common set of virtio device drivers are now included as standard in the Linux kernel. To maximize performance, our performance evaluation configured the application VM with Palacios’s virtio-compatible virtual NIC, using the default Linux virtio network driver.

MTU: The maximum transmission unit (MTU) of a networking layer is the size of the largest protocol data unit that the layer can pass onwards. A larger MTU improves throughput because each packet carries more user data while protocol headers have a fixed size. A larger MTU also means that fewer packets need to be processed to transfer a given amount of data. Where per-packet processing costs are significant, larger MTUs are distinctly preferable. Because VNET/P adds to the per-packet processing cost, supporting large MTUs is helpful.

VNET/P presents an Ethernet abstraction to the application VM. The most common Ethernet MTU is 1500 bytes. However, 1 Gbit and 10 Gbit Ethernet can also use “jumbo frames”, with an MTU of 9000 bytes. Other networking technologies support even larger MTUs. To leverage the large MTUs of underlying physical NICs, VNET/P itself

supports MTU sizes of up to 64 KB.¹ The application OS can determine the virtual NIC's MTU and then transmit/receive accordingly. VNET/P advertises the appropriate MTU.

The MTU used by virtual NIC can result in encapsulated VNET/P packets that exceed the MTU of the underlying physical network. In this case, fragmentation has to occur, either in the VNET/P bridge or in the host NIC (via TCP Segmentation Offloading (TSO)). Fragmentation and reassembly is handled by VNET/P and is totally transparent to the application VM. However, performance will suffer when significant fragmentation occurs. Thus it is important that the application VM's device driver select an MTU carefully, and recognize that the desirable MTU may change over time, for example after a migration to a different host. In Sect. 5, we analyze throughput using different MTUs.

4.5 VNET/P bridge

The VNET/P bridge functions as a network bridge to direct packets between the VNET/P core and the physical network through the host NIC. It operates based on the routing decisions made by the VNET/P core which are passed along with the packets to be forwarded. It is implemented as a kernel module running in the host.

When the VNET/P core hands a packet and routing directive up to the bridge, one of two transmission modes will occur, depending on the destination. In a *direct send*, the Ethernet packet is directly sent. This is common for when a packet is exiting a VNET overlay and entering the physical network, as typically happens on the user's network. It may also be useful when all VMs will remain on a common layer 2 network for their lifetime. In an *encapsulated send* the packet is encapsulated in a UDP packet and the UDP packet is sent to the directed destination IP address and port. This is the common case for traversing a VNET overlay link. Similarly, for packet reception, the bridge uses two modes, simultaneously. In a *direct receive* the host NIC is run in promiscuous mode, and packets with destination MAC addresses corresponding to those requested by the VNET/P core are handed over to it. This is used in conjunction with direct send. In an *encapsulated receive* UDP packets bound for the common VNET link port are disassembled and their encapsulated Ethernet packets are delivered to the VNET/P core. This is used in conjunction with encapsulated send. Our performance evaluation focuses solely on encapsulated send and receive.

4.6 Control

The VNET/P control component allows for remote and local configuration of links, interfaces, and routing rules so

that an overlay can be constructed and changed over time. VNET/U already has user-level tools to support VNET, and, as we described in Sect. 3, a range of work already exists on the configuration, monitoring, and control of a VNET overlay. In VNET/P, we reuse these tools as much as possible by having the user-space view of VNET/P conform closely to that of VNET/U. The *VNET/P configuration console* allows for local control to be provided from a file, or remote control via TCP-connected VNET/U clients (such as tools that automatically configure a topology that is appropriate for the given communication pattern among a set of VMs [50]). In both cases, the VNET/P control component is also responsible for validity checking before it transfers the new configuration to the VNET/P core.

4.7 Performance-critical data paths and flows

Figure 7 depicts how the components previously described operate during packet transmission and reception. These are the performance critical data paths and flows within VNET/P, assuming that virtio virtual NICs (Sect. 4.4) are used. The boxed regions of the figure indicate steps introduced by virtualization, both within the VMM and within the host OS kernel. There are also additional overheads involved in the VM exit handling for I/O port reads and writes and for interrupt injection.

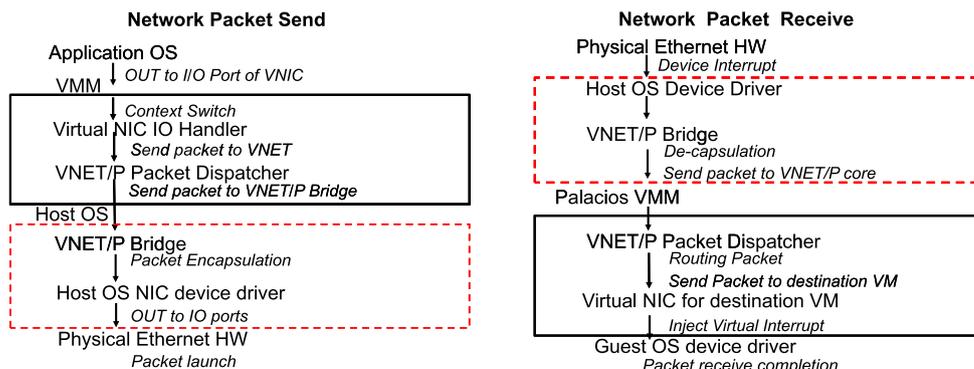
Transmission: The guest OS in the VM includes the device driver for the virtual NIC. The driver initiates packet transmission by writing to a specific virtual I/O port after it puts the packet into the NIC's shared ring buffer (TXQ). The I/O port write causes an exit that gives control to the virtual NIC I/O handler in Palacios. The handler reads the packet from the buffer and writes it to VNET/P packet dispatcher. The dispatcher does a routing table lookup to determine the packet's destination. For a packet destined for a VM on some other host, the packet dispatcher puts the packet into the receive buffer of the VNET/P bridge and notify it. Meanwhile, VNET/P bridge fetches the packet from the receive buffer, determines its destination VNET/P bridge, encapsulates the packet, and transmits it to the physical network via the host's NIC.

Note that while the packet is handed off multiple times, it is copied only once inside the VMM, from the send buffer (TXQ) of the receive buffer of the VNET/P bridge. Also note that while the above description and the diagram suggest sequentiality, packet dispatch can occur on a separate kernel thread running on a separate core, and the VNET/P bridge itself introduces additional concurrency. From the guest's perspective, the I/O port write that initiated transmission returns essentially within a VM exit/entry time.

Reception: The path for packet reception is essentially symmetric to that of transmission. The host NIC in the host machine receives a packet using its standard driver and delivers it to the VNET/P bridge. The bridge unencapsulates

¹This may be expanded in the future. Currently, it has been sized to support the largest possible IPv4 packet size.

Fig. 7 Performance-critical data paths and flows for packet transmission and reception. Solid boxed steps and components occur within the VMM itself, while dashed boxed steps and components occur in the host OS



the packet and sends the payload (the raw Ethernet packet) to the VNET/P core. The packet dispatcher in VNET/P core determines its destination VM and puts the packet into the receive buffer (RXQ) of its virtual NIC.

Similar to transmission, there is considerably concurrency in the reception process. In particular, packet dispatch can occur in parallel with the reception of the next packet.

4.8 Performance tuning parameters

VNET/P is configured with the following parameters:

- Whether guest-driven, VMM-driven, or adaptive mode is used.
- For adaptive operation, the upper and lower rate bounds (α_l, α_u) for switching between guest-driven, and VMM-driven modes, as well as the window size ω over which rates are computed.
- The number of packet dispatcher threads $n_{dispatchers}$ that are instantiated.
- The yield model and parameters for the bridge thread, the packet dispatch threads, and the VMM’s halt handler.

The last of these items requires some explanation, as it presents an important tradeoff between VNET/P latency and CPU consumption. In both the case of packets arriving from the network and packets arriving from the guest’s virtual NIC, there is, conceptually, a thread running a receive operation that could block or poll. First, consider packet reception from the physical network. Suppose a bridge thread is waiting for UDP packet to arrive. If no packets have arrived recently, then blocking would allow the core on which the thread is running to be yielded to another thread, or for the core to be halted. This would be the ideal for minimizing CPU consumption, but on the next arrival, there will be a delay, even if the core is idle, in handling the packet since the thread will need to be scheduled. There is a similar tradeoff in the packet dispatcher when it comes to packet transmission from the guest.

A related tradeoff, for packet reception in the guest, exists in the VMM’s model for handling the halt state. If the

guest is idle, it will issue an HLT or related instruction. The VMM’s handler for this case consists of waiting until an interrupt needs to be injected into the guest. If it polls, it will be able to respond as quickly as possible, thus minimizing the packet reception latency, while if it blocks, it will consume less CPU.

For all three of these cases VNET/P’s and the VMM’s *yield strategy* comes into play. Conceptually, these cases are written as polling loops, which in their loop bodies can yield the core to another thread, and optionally put the thread to sleep pending a wake-up after a signaled event or the passage of an interval of time. Palacios currently has selectable yield strategy that is used in these loops, and the strategy has three different options, one of which is chosen when the VM is configured:

- Immediate yield. Here, if there is no work, we immediately yield the core to any competing threads. However, if there are no other active threads that the core can run, the yield immediately returns. A poll loop using the immediate yield has the lowest latency while still being fair to competing threads.
- Timed yield: Here, if there is no work, we put the thread on a wake up queue and yield CPU. Failing any other event, the thread will be awakened by the passage of time T_{sleep} . A poll loop using the timed yield strategy minimizes CPU usage, but at the cost of increased latency.
- Adaptive yield: Here, the poll loop reports how much time has passed since it last did any work. Until that time exceeds a threshold T_{nowork} , the immediate yield strategy is used, and afterwards the timed yield strategy is used. By setting the threshold, different tradeoffs between latency and CPU usage are made.

In our performance evaluations, we use the parameters given in Table 1 to focus on the performance limits of VNET/P.

Table 1 Parameters used to configure VNET/P in our performance evaluation

Parameter	Value
Mode	Adaptive
α_l	10^3 packets/s
α_u	10^4 packets/s
ω	5 ms
$n_{dispatchers}$	1
Yield strategy	Immediate yield
T_{sleep}	Not used
T_{nowork}	Not used

5 Performance evaluation

The purpose of our performance evaluation is to determine how close VNET/P comes to native throughput and latency in the most demanding (lowest latency, highest throughput) hardware environments. We consider communication between two machines whose NICs are directly connected in most of our detailed benchmarks.

In the virtualized configuration the guests and performance testing tools run on top of Palacios with VNET/P carrying all traffic between them using encapsulation. In the native configuration, the same guest environments run directly on the hardware.

Our evaluation of communication performance in this environment occurs at three levels. First, we benchmark the TCP and UDP bandwidth and latency. Second, we benchmark MPI using a widely used benchmark. Finally, we evaluated the performance of the HPCC and NAS application benchmarks in a cluster to see VNET/P's impact on the performance and scalability of parallel applications.

5.1 Testbed and configurations

Most of our microbenchmark tests are focused on the end-to-end performance of VNET/P. Therefore our testbed consists of two physical machines, which we call host machines. Each machine has a quadcore 2.4 GHz X3430 Intel Xeon processor, 8 GB RAM, a Broadcom NetXtreme II 1 Gbps Ethernet NIC (1000BASE-T), and a NetEffect NE020 10 Gbps Ethernet fiber optic NIC (10GBASE-SR) in a PCI-e slot. The Ethernet NICs of these machines are directly connected with twisted pair and fiber patch cables.

A range of our measurements are made using the cycle counter. We disabled DVFS control on the machine's BIOS, and in the host Linux kernel. We also sanity-checked the measurement of larger spans of time by comparing cycle counter-based timing with a separate wall clock.

All microbenchmarks included in the performance section are run in the testbed described above. The HPCC and

NAS application benchmarks are run on a 6-node test cluster described in Sect. 5.4.

We considered the following two software configurations:

- *Native*: In the native configuration, neither Palacios nor VNET/P is used. A minimal BusyBox-based Linux environment based on an unmodified 2.6.30 kernel runs directly on the host machines. We refer to the 1 and 10 Gbps results in this configuration as *Native-1G* and *Native-10G*, respectively.
- *VNET/P*: The VNET/P configuration corresponds to the architectural diagram given in Fig. 1, with a single guest VM running on Palacios. The guest VM is configured with one virtio network device, 2 cores, and 1 GB of RAM. The guest VM runs a minimal BusyBox-based Linux environment, based on the 2.6.30 kernel. The kernel used in the VM is identical to that in the Native configuration, with the exception that the virtio NIC drivers are loaded. The virtio MTU is configured as 9000 Bytes. We refer to the 1 and 10 Gbps results in this configuration as *VNET/P-1G* and *VNET/P-10G*, respectively.

To assure accurate time measurements both natively and in the virtualized case, our guest is configured to use the CPU's cycle counter, and Palacios is configured to allow the guest direct access to the underlying hardware cycle counter. Our 1 Gbps NIC only supports MTUs up to 1500 bytes, while our 10 Gbps NIC can support MTUs of up to 9000 bytes. We use these maximum sizes unless otherwise specified.

5.2 TCP and UDP microbenchmarks

Latency and throughput are the fundamental measurements we use to evaluate the VNET/P system performance. First, we consider these at the IP level, measuring the round-trip latency, the UDP goodput, and the TCP throughput between two nodes. We measure round-trip latency using *ping* by sending ICMP packets of different sizes. UDP and TCP throughput are measured using *ttcp-1.10*.

UDP and TCP with a standard MTU: Fig. 8 shows the TCP throughput and UDP goodput achieved in each of our configurations on each NIC. For the 1 Gbps network, host MTU is set to 1500 bytes, and for the 10 Gbps network, host MTUs of 1500 bytes and 9000 bytes are both tested. For 1 Gbps, we also compare with VNET/U running on the same hardware with Palacios. Compared to previously reported results (21.5 MB/s, 1 ms), the combination of the faster hardware we use here, and Palacios, leads to VNET/U increasing its bandwidth by 330 %, to 71 MB/s, with a 12 % reduction in latency, to 0.88 ms. We also tested VNET/U with VMware, finding that bandwidth increased by 63 % to 35 MB/s, with no change in latency. The difference in performance of VNET/U on the two VMMs is due to a custom

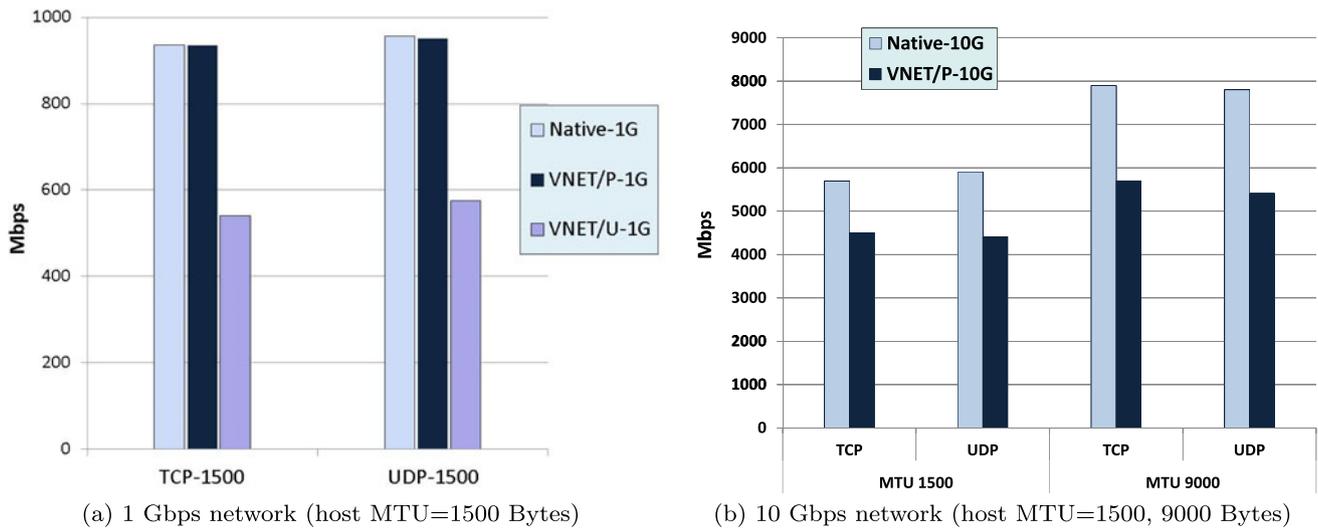


Fig. 8 End-to-end TCP throughput and UDP goodput of VNET/P on 1 and 10 Gbps network. VNET/P performs identically to the native case for the 1 Gbps network and achieves 74–78 % of native throughput for the 10 Gbps network

tap interface in Palacios, while on VMware, the standard host-only tap is used. Even with this optimization, VNET/U cannot saturate a 1 Gbps link.

We begin by considering UDP goodput when a standard host MTU size is used. For UDP measurements, `ttcp` was configured to use 64000 byte writes sent as fast as possible over 60 seconds. For the 1 Gbps network, VNET/P easily matches the native goodput. For the 10 Gbps network, VNET/P achieves 74 % of the native UDP goodput.

For TCP throughput, `ttcp` was configured to use a 256 KB socket buffer, and to communicate 40 MB writes were made. Similar to the UDP results, VNET/P has no difficulty achieving native throughput on the 1 Gbps network. On the 10 Gbps network, using a standard Ethernet MTU, it achieves 78 % of the native throughput. The UDP goodput and TCP throughput that VNET/P is capable of, using a standard Ethernet MTU, are approximately 8 times those we would expect from VNET/U given the 1 Gbps results.

UDP and TCP with a large MTU: We now consider TCP and UDP performance with 9000 byte jumbo frames our 10 Gbps NICs support. We adjusted the VNET/P MTU so that the ultimate encapsulated packets will fit into these frames without fragmentation. For TCP we configure `ttcp` to use writes of corresponding size, maximize the socket buffer size, and do 4 million writes. For UDP, we configure `ttcp` to use commensurately large packets sent as fast as possible for 60 seconds. The results are also shown in the Fig. 8. We can see that performance increases across the board compared to the 1500 byte MTU results. Compared to the VNET/U performance we would expect in this configuration, the UDP goodput and TCP throughput of VNET/P are over 10 times higher.

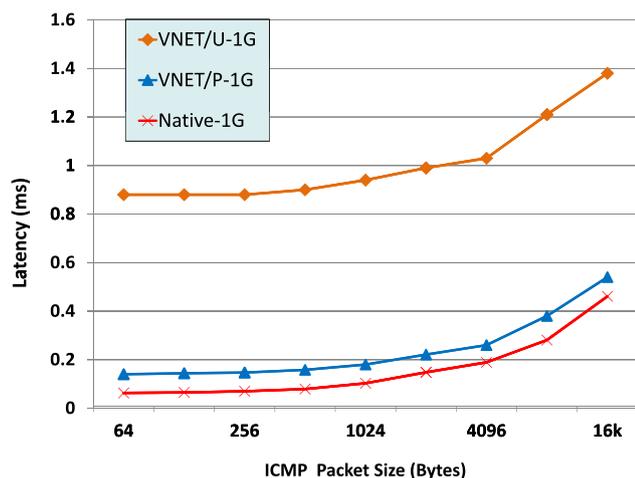
Latency: Fig. 9 shows the round-trip latency for different packet sizes, as measured by ping. The latencies are the average of 100 measurements. While the increase in latency of VNET/P over Native is significant in relative terms ($2\times$ for 1 Gbps, $3\times$ for 10 Gbps), it is important to keep in mind the absolute performance. On a 10 Gbps network, VNET/P achieves a 130 μ s round-trip, end-to-end latency. The latency of VNET/P is almost seven times lower than that of VNET/U.

5.3 MPI microbenchmarks

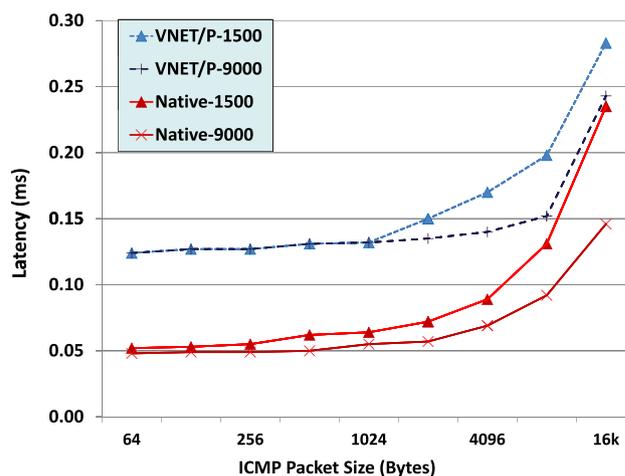
Parallel programs for distributed memory computers are typically written to the MPI interface standard. We used the OpenMPI 1.3 [9] implementation in our evaluations. We measured the performance of MPI over VNET/P by employing the widely-used Intel MPI Benchmark Suite (IMB 3.2.2) [20], focusing on the point-to-point messaging performance. We compared the basic MPI latency and bandwidth achieved by VNET/P and natively.

Figures 10 and 11(a) illustrate the latency and bandwidth reported by Intel MPI PingPong benchmark for our 10 Gbps configuration. Here the latency measured is the one-way, end-to-end, application-level latency. That is, it is the time from when an MPI send starts on one machine to when its matching MPI receive call completes on the other machine. For both Native and VNET/P, the host MTU is set to 9000 bytes.

VNET/P's small message MPI latency is about 55 μ s, about 2.5 times worse than the native case. However, as the message size increases, the latency difference decreases. The measurements of end-to-end bandwidth as a function of message size show that native MPI bandwidth is slightly



(a) 1 Gbps network (Host MTU=1500 Bytes)



(b) 10 Gbps network (Host MTU=1500, 9000 Bytes)

Fig. 9 End-to-end round-trip latency of VNET/P as a function of ICMP packet size. Small packet latencies on a 10 Gbps network in VNET/P are $\sim 130 \mu\text{s}$

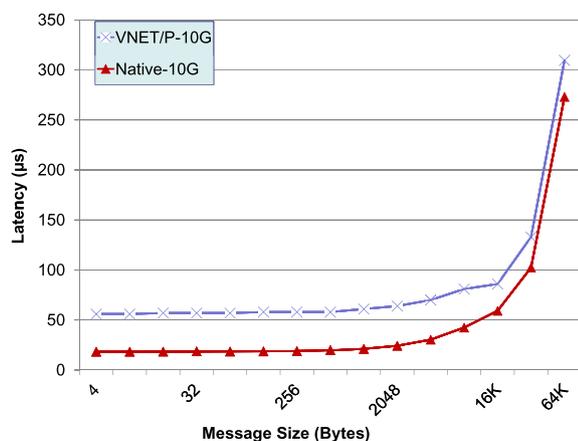


Fig. 10 One-way latency on 10 Gbps hardware from Intel MPI Ping-Pong microbenchmark

lower than raw UDP or TCP throughput, and VNET/P performance tracks it similarly. The bottom line is that the current VNET/P implementation can deliver an MPI latency of $55 \mu\text{s}$ and bandwidth of 510 MB/s on 10 Gbps Ethernet hardware.

Figure 11(b) shows the results of the MPI SendRecv microbenchmark in which each node simultaneously sends and receives. There is no reduction in performance between the bidirectional case and the unidirectional case. For both figures, the absolute numbers are shown, but it is important to note that the overhead is quite stable in percentage terms once the message size is large: beyond 256K, the one-way bandwidth of VNET/P is around 74 % of native, and the two-way bandwidth is around 62 % of native. A possible interpretation is that we become memory copy bandwidth limited.

5.4 HPCC benchmarks on more nodes

To test VNET/P performance on more nodes, we ran the HPCC benchmark [19] suite on a 6 node cluster with 1 Gbps and 10 Gbps Ethernet. Each node was equipped with two quad-core 2.3 GHz 2376 AMD Opterons, 32 GB of RAM, an nVidia MCP55 Forthdeth 1 Gbps Ethernet NIC and a Net-Effect NE020 10 Gbps Ethernet NIC. The nodes were connected via a Fujitsu XG2000 10 Gb Ethernet Switch. Similar to our A range of our measurements are made using the cycle counter. We disabled DVFS control on the machine's BIOS, and in the host Linux kernel. We also used the cycle counter on these machines, and we applied the same techniques (deactivated DVFS, sanity-checking against an external clock for larger time spans) described in Sect. 5.1 to ensure accurate timing.

The VMs were all configured exactly as in previous tests, with 4 virtual cores, 1 GB RAM, and a virtio NIC. For the VNET/P test case, each host ran one VM. We executed tests with 2, 3, 4, 5, and 6 VMs, with 4 HPCC processes per VM (one per virtual core). Thus, our performance results are based on HPCC with 8, 12, 16, 20 and 24 processes for both VNET/P and Native tests. In the native cases, no VMs were used, and the processes ran directly on the host. For 1 Gbps testing, the host MTU was set to 1500, while for the 10 Gbps cases, the host MTU was set to 9000.

Latency-bandwidth benchmark: This benchmark consists of the ping-pong test and the ring-based tests. The ping-pong test measures the latency and bandwidth between all distinct pairs of processes. The ring-based tests arrange the processes into a ring topology and then engage in collective communication among neighbors in the ring, measuring bandwidth and latency. The ring-based tests model

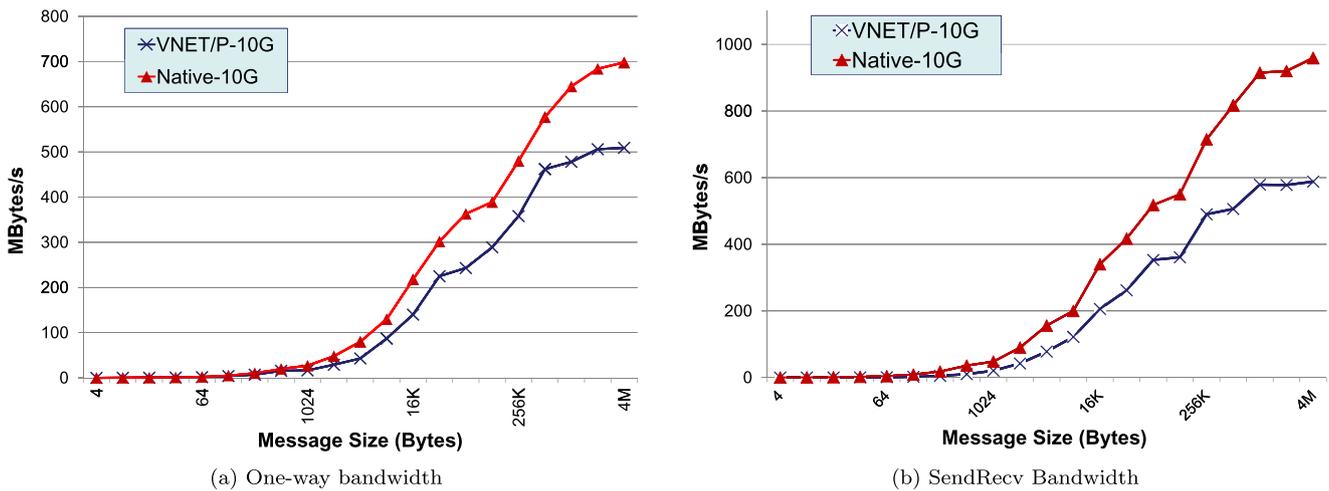


Fig. 11 Intel MPI PingPong microbenchmark showing (a) one-way bandwidth and (b) bidirectional bandwidth as a function of message size on the 10 Gbps hardware

the communication behavior of multi-dimensional domain-decomposition applications. Both naturally ordered rings and randomly ordered rings are evaluated. Communication is done with MPI non-blocking sends and receives, and MPI SendRecv. Here, the bandwidth per process is defined as total message volume divided by the number of processes and the maximum time needed in all processes. We reported the ring-based bandwidths by multiplying them with the number of processes in the test.

Figure 12 shows the results for different numbers of test processes. The ping-pong latency and bandwidth results are consistent with what we saw in the previous microbenchmarks: in the 1 Gbps network, bandwidth are nearly identical to those in the native cases while latencies are 1.2–2 times higher. In the 10 Gbps network, bandwidths are within 60–75 % of native while latencies are about 2 to 3 times higher. Both latency and bandwidth under VNET/P exhibit the same good scaling behavior of the native case.

5.5 Application benchmarks

We evaluated the effect of a VNET/P overlay on application performance by running two HPCC application benchmarks and the whole NAS benchmark suite on the cluster described in Sect. 5.4. Overall, the performance results from the HPCC and NAS benchmarks suggest that VNET/P can achieve high performance for many parallel applications.

HPCC application benchmarks: We considered the two application benchmarks from the HPCC suite that exhibit the large volume and complexity of communication: MPI-RandomAccess and MPIFFT. For 1 Gbps networks, the difference in performance is negligible so we focus here on 10 Gbps networks.

In MPIRandomAccess, random numbers are generated and written to a distributed table, with local buffering. Performance is measured by the billions of updates per second

(GUPs) that are performed. Figure 13(a) shows the results of MPIRandomAccess, comparing the VNET/P and Native cases. VNET/P achieves 65–70 % application performance compared to the native cases, and performance scales similarly.

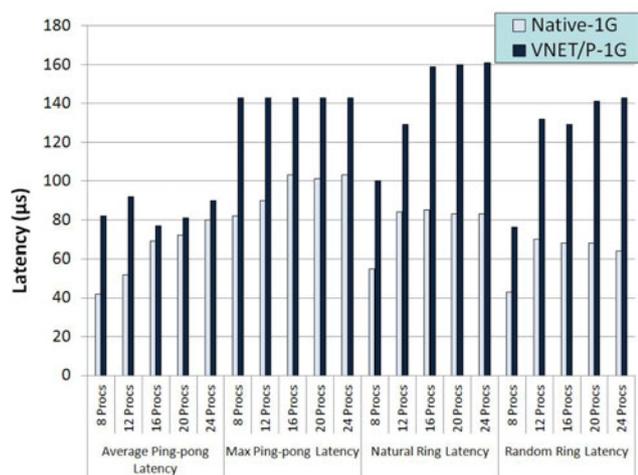
MPIFFT implements a double precision complex one-dimensional Discrete Fourier Transform (DFT). Figure 13(b) shows the results of MPIFFT, comparing the VNET/P and Native cases. It shows that VNET/P’s application performance is within 60–70 % of native performance, with performance scaling similarly.

NAS parallel benchmarks: The NAS Parallel Benchmark (NPB) suite [54] is a set of five kernels and three pseudo-applications that is widely used in parallel performance evaluation. We specifically use NPB-MPI 2.4 in our evaluation. In our description, we name executions with the format “name.class.procs”. For example, *bt.B.16* means to run the BT benchmark on 16 processes with a class B problem size.

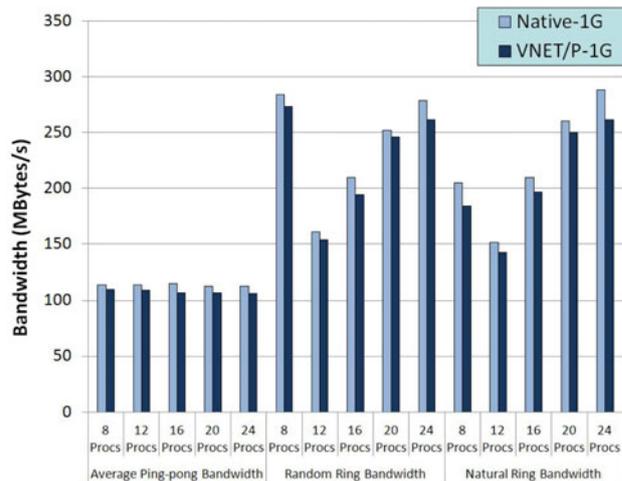
We run each benchmark with at least two different scales and one problem size, except FT, which is only run with 16 processes. One VM is run on each physical machine, and it is configured as described in Sect. 5.4. The test cases with 8 processes are running within 2 VMs and 4 processes started in each VM. The test cases with 9 processes are run with 4 VMs and 2 or 3 processes per VM. Test cases with 16 processes have 4 VMs with 4 processes per VM. We report each benchmark’s *Mop/s total* result for both native and with VNET/P.

Figure 14 shows the NPB performance results, comparing the VNET/P and Native cases on both 1 Gbps and 10 Gbps networks. The upshot of the results is that for most of the NAS benchmarks, VNET/P is able to achieve in excess of 95 % of the native performance even on 10 Gbps networks. We now describe the results for each benchmark.

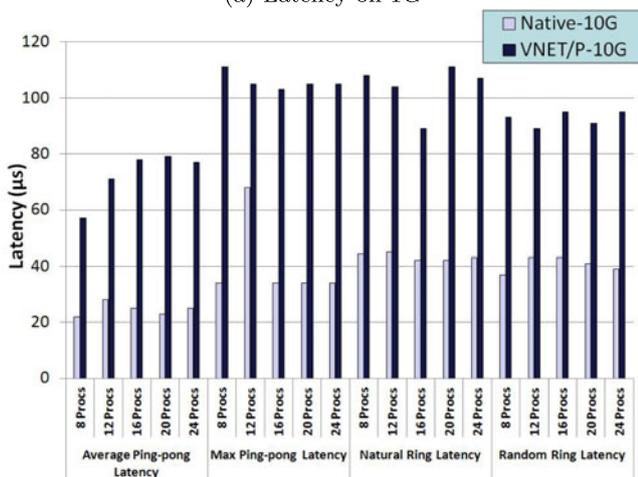
EP is an “embarrassingly parallel” kernel that estimates the upper achievable limits for floating point performance. It



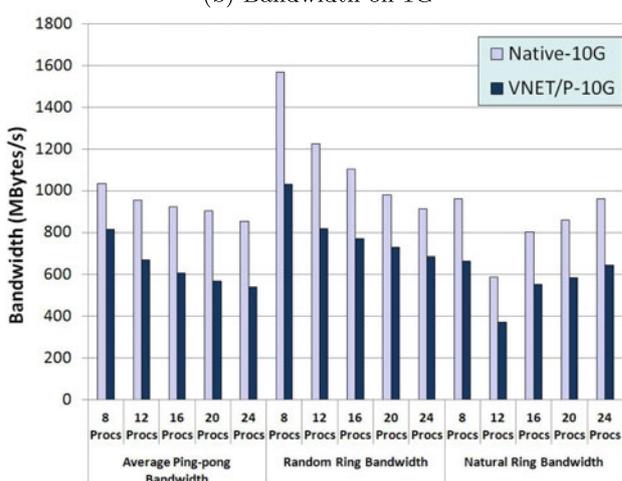
(a) Latency on 1G



(b) Bandwidth on 1G



(c) Latency on 10G



(d) Bandwidth on 10G

Fig. 12 HPC Latency-bandwidth benchmark for both 1 Gbps and 10 Gbps. Ring-based bandwidths are multiplied by the total number of processes in the test. The ping-pong latency and bandwidth tests show

results that are consistent with the previous microbenchmarks, while the ring-based tests show that latency and bandwidth of VNET/P scale similarly to the native cases

does not require a significant interprocessor communication. VNET/P achieves native performance in all cases.

MG is a simplified multigrid kernel that requires highly structured long distance communication and tests both short and long distance data communication. With 16 processes, MG achieves native performance on the 1 Gbps network, and 81 % of native performance on the 10 Gbps network.

CG implements the conjugate gradient method to compute an approximation to the smallest eigenvalue of a large sparse symmetric positive definite matrix. It is typical of unstructured grid computations in that it tests irregular long distance communication, employing unstructured matrix vector multiplication. With 16 processes, CG achieves native performance on the 1 Gbps network and 94 % of native performance on the 10 Gbps network.

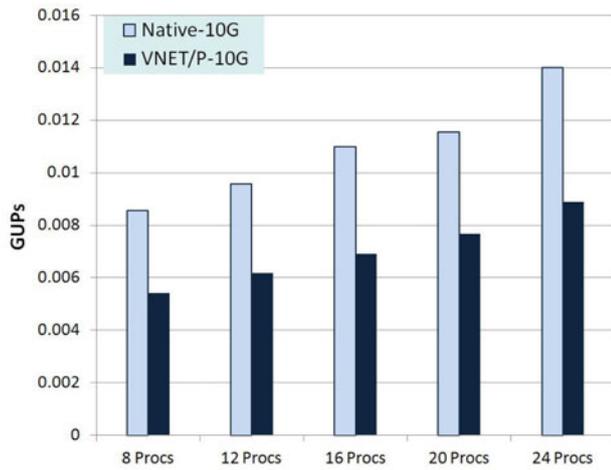
FT implements the solution of partial differential equations using FFTs, and captures the essence of many spectral

codes. It is a rigorous test of long-distance communication performance. With 16 nodes, it achieves 82 % of native performance on 1 Gbps and 86 % of native performance on 10 Gbps.

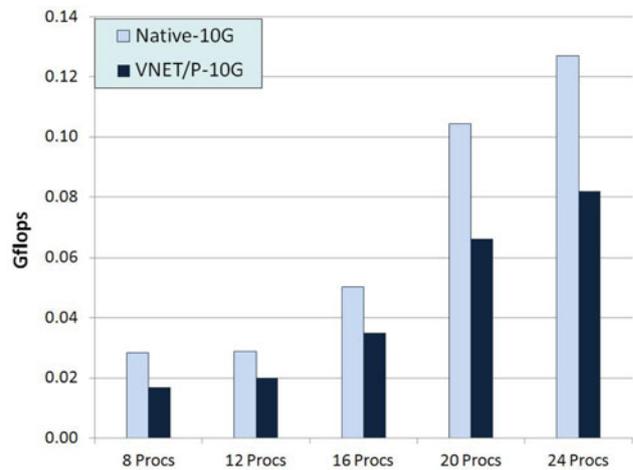
IS implements a large integer sort of the kind that is important in particle method codes and tests both integer computation speed and communication performance. Here VNET/P achieves native performance in all cases.

LU solves a regular-sparse, block (5×5) lower and upper triangular system, a problem associated with implicit computational fluid dynamics algorithms. VNET/P achieves 75 %–85 % of native performance on this benchmark, and there is no significant difference between the 1 Gbps and 10 Gbps network.

SP and BT implement solutions of multiple, independent systems of non diagonally dominant, scalar, pentadiagonal equations, also common in computational fluid dynamics.



(a) MPIRandomAccess



(b) MPIFFT

Fig. 13 HPC application benchmark results. VNET/P achieves reasonable and scalable application performance when supporting communication-intensive parallel application workloads on 10 Gbps networks. On 1 Gbps networks, the difference is negligible

Mop/s	Native-1G	VNET/P-1G	$\frac{VNET/P-1G}{Native-1G}(\%)$	Native-10G	VNET/P-10G	$\frac{VNET/P-10G}{Native-10G}(\%)$
ep.B.8	103.15	101.94	98.8 %	102.18	102.12	99.9 %
ep.B.16	204.88	203.9	99.5 %	208	206.52	99.3 %
ep.C.8	103.12	102.1	99.0 %	103.13	102.14	99.0 %
ep.C.16	206.24	204.14	99.0 %	206.22	203.98	98.9 %
mg.B.8	4400.52	3840.47	87.3 %	5110.29	3796.03	74.3 %
mg.B.16	1506.77	1498.65	99.5 %	9137.26	7405	81.0 %
cg.B.8	1542.79	1319.43	85.5 %	2096.64	1806.57	86.2 %
cg.B.16	160.64	159.69	99.4 %	592.08	554.91	93.7 %
ft.B.16	1575.83	1290.78	81.9 %	1432.3	1228.39	85.8 %
is.B.8	78.88	74.61	94.6 %	59.15	59.04	99.8 %
is.B.16	35.99	35.78	99.4 %	23.09	23	99.6 %
is.C.8	89.54	82.15	91.7 %	132.08	131.87	99.8 %
is.C.16	84.76	82.22	97.0 %	77.77	76.94	98.9 %
lu.B.8	6818.52	5495.23	80.6 %	7173.65	6021.78	83.9 %
lu.B.16	7847.99	6694.12	85.3 %	12981.86	9643.21	74.3 %
sp.B.9	1361.38	1215.85	89.3 %	2634.53	2421.98	91.9 %
sp.B.16	1489.32	1399.6	94.0 %	3010.71	2916.81	96.9 %
bt.B.9	3423.52	3297.04	96.3 %	5229.01	4076.52	78.0 %
bt.B.16	4599.38	4348.99	94.6 %	6315.11	6105.11	96.7 %

Fig. 14 NAS Parallel Benchmark performance with VNET/P on 1 Gbps and 10 Gbps networks. VNET/P can achieve native performance on many applications, while it can get reasonable and scalable per-

formance when supporting highly communication-intensive parallel application workloads

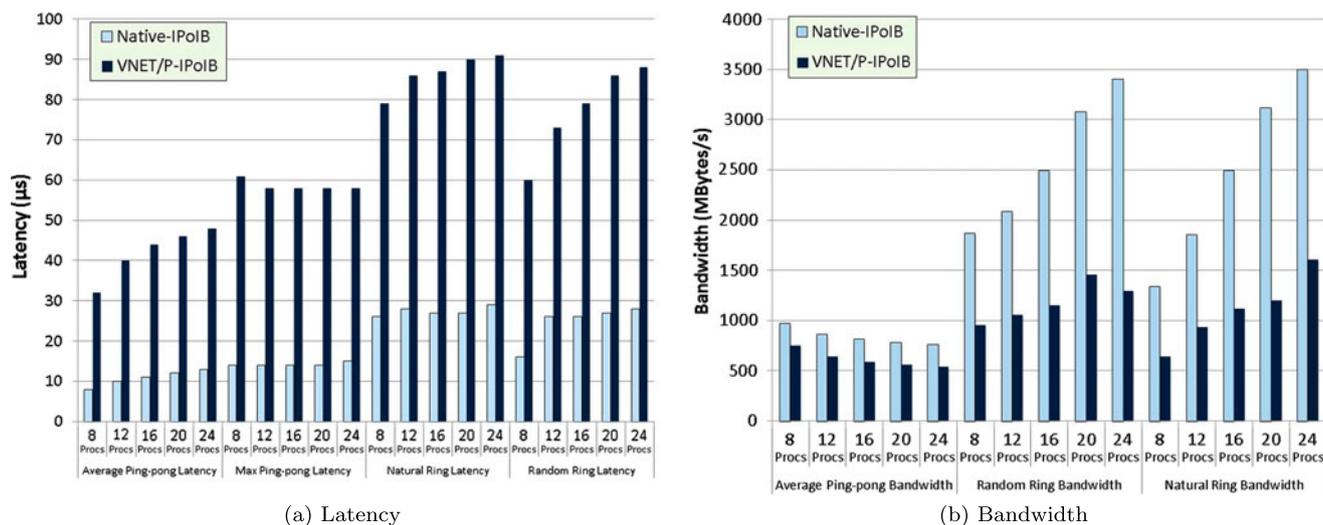


Fig. 15 Preliminary results for the HPCC latency-bandwidth benchmark, comparing the fully virtualized environment using VNET/P running on IPoIB with the purely native environment using IPoIB. Ring-based bandwidths are multiplied by the total number of processes in the test

The salient difference between the two is the communication to computation ratio. For SP with 16 processes, VNET/P achieves 94 % of native performance on 1 Gbps and around 97 % of native on 10 Gbps. For BT at the same scale, 95 % of native at 1 Gbps and 97 % of native at 10 Gbps are achieved.

It is worth mentioning that in microbenchmarking we measured throughput and latency separately. During the throughput microbenchmarking, large packets were continuously sent, saturating VNET/P. On the other hand, for the application benchmarks, the network communication consisted of a mixture of small and large packets, and so their performance was determined both by throughput and latency. Recall that the latency overhead of VNET/P is on the order of 2–3 \times . This may explain why some application benchmarks cannot achieve native performance even in the 1 Gbps case (e.g. LU, FT) despite VNET/P achieving native throughput in the microbenchmarks. This also suggests that in order to gain fully native application performance, we need to further reduce the small packet latency. In a separate paper from our group [5], we have described additional techniques to do so, resulting in microbenchmark latency overheads of 1.2–1.3 \times . This results in the latency-limited application benchmarks achieving > 95 % of native performance.

6 VNET/P portability

We now describe our experiences in running the Linux-based VNET/P on alternative hardware, and a second implementation of VNET/P for a different host environment, the Kitten lightweight kernel. These experiences are in support of hardware independence, the 3rd goal of VNET artic-

ulated in Sect. 3. Regardless of the underlying networking hardware or host OS, the guests see a simple Ethernet LAN.

6.1 Infiniband

Infiniband in general, and our Mellanox hardware in particular, supports IP over Infiniband (IPoIB) [24]. The IPoIB functionality of the device driver, running in the host, allows the TCP/IP stack of the host to use the NIC to transport IP packets. Because VNET/P sends UDP packets, it is trivial to direct it to communicate over the IB fabric via this mechanism. No code changes to VNET/P are needed. The primary effort is in bringing up the Infiniband NICs with IP addresses, and establishing routing rules based on them. The IP functionality of Infiniband co-exists with native functionality.

It is important to point out that we have not yet tuned VNET/P for this configuration, but we describe our current results below. On our testbed the preliminary, “out of the box” performance of VNET/P includes a ping latency of 155 μ s and a TTCP throughput of 3.6 Gbps. We also ran the HPCC benchmarks on a 6-node cluster with the same configuration as in Sect. 5.4, comparing the performance of VNET/P on IPoIB with native performance on IPoIB.

Figure 15 shows the results for the HPCC latency-bandwidth benchmark with different numbers of test processes. In the pingpong test, VNET/P on IPoIB can achieve 70–75 % of native bandwidth, while its latencies are about 3 to 4 times higher. In the ring-based tests, VNET/P on IPoIB can achieve on average 50–55 % of the native bandwidth, while the latency is on average about 2.5 to 3 times higher.

Figure 16 shows the results of two HPCC application benchmarks comparing the VNET/P and native cases. For

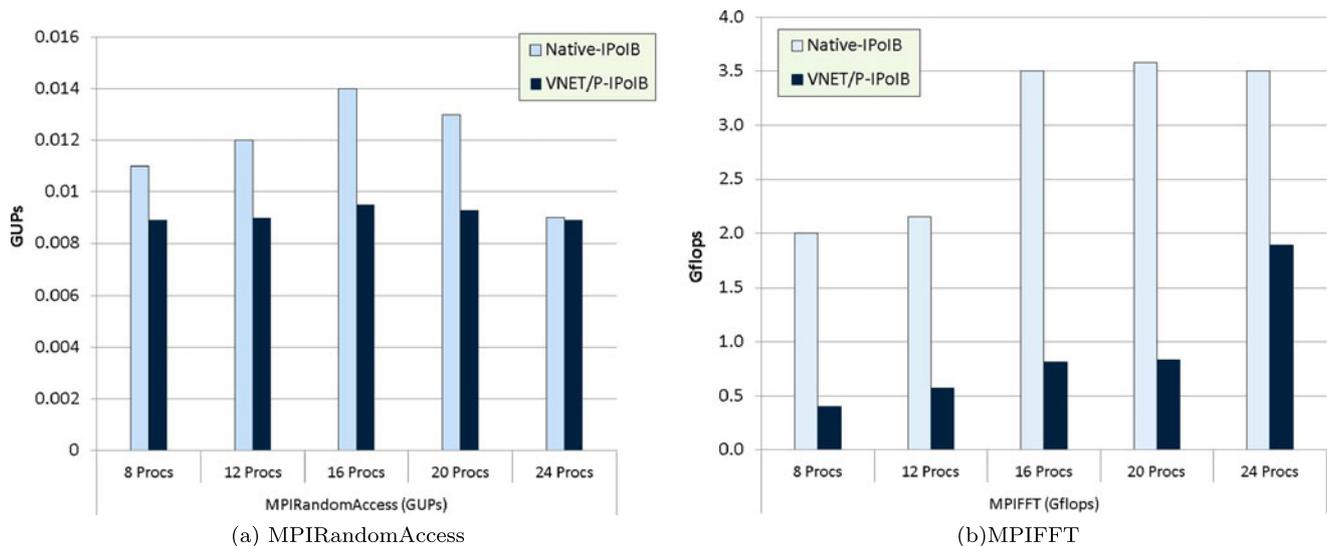


Fig. 16 Preliminary results for HPC application benchmarks, comparing the fully virtualized environment using VNET/P over IPoIB with the purely native environment using IPoIB

MPIRandomAccess, VNET/P on IPoIB achieves 75–80 % of the native application performance, and performance scales similarly. On the other hand, the results of MPIFFT show that VNET/P's achieve about 30–45 % of the native performance, similar performance scaling.

Note also that there are two overheads involved in this form of mapping, the VNET/P overhead, and the IPoIB overhead. We have not yet determined the relative sizes of these overheads.

6.2 Cray Gemini

As a proof-of-concept we have brought Palacios and VNET/P up on a Cray XK6, the Curie testbed at Sandia National Labs. Curie consists of 50 compute nodes that each combine an AMD Opteron 6272, an nVidia Tesla M2090, and 32 GB of RAM. The nodes communicate using the Gemini network, which is a three dimensional torus. The Gemini NIC connects to a node via a HyperTransport 3 link. The theoretical peak inter-node latency and bandwidth are 1.5 μ s and 5 GB/s (40 Gbps). Measured MPI throughput for large (64 KB) messages is about 4 GB/s (32 Gbps). Vaughan et al [55] give a more detailed description of Gemini in the context of a large installation.

Although the XK6 runs Cray Compute Node Linux as its host OS, the changes to Palacios needed to support this configuration have proven to be relatively minor. Gemini supports an “IPoG” layer. This creates a virtual Ethernet NIC abstraction on top of the underlying Gemini NIC, and this abstraction then supports the host OS’s TCP/IP stack. Given this, VNET/P maps its encapsulated UDP traffic straightforwardly, simply by using the IP addresses assigned to the relevant Gemini NICs. Except for talking to the IPoG layer

instead of talking to an Ethernet NIC, the architecture of VNET/P on Gemini is identical to that shown in Fig. 1.

We are in the process of tuning VNET/P in this environment. Our preliminary results show VNET/P achieves a TCP throughput of 1.6 GB/s (13 Gbps). We are currently addressing a likely precision-timing problem that is limiting performance. We have not yet determined the relative overheads of VNET/P and IPoG.

6.3 VNET/P for Kitten

The VNET/P model can be implemented successfully in host environments other than Linux. Specifically, we have developed a version of VNET/P for the Kitten Lightweight Kernel. This version focuses on high performance Infiniband support and leverages precision timing and other functionality made possible by Kitten.

Figure 17 shows the architecture of VNET/P for Kitten and can be compared and contrasted with the VNET/P for Linux architecture shown in Fig. 1. While the architectures are substantially different, the abstraction that the guest VMs see is identical: the guests still believe they are connected by a simple Ethernet network.

Kitten has, by design, a minimal set of in-kernel services. For this reason, the VNET/P Bridge functionality is not implemented in the kernel, but rather in a privileged service VM called the Bridge VM that has direct access to the physical Infiniband device. In place of encapsulating Ethernet packets in UDP packets for transmission to a remote VNET/P core, VNET/P's for InfiniBand on Kitten simply maps Ethernet packets to InfiniBand frames. These frames are then transmitted through an InfiniBand queue pair accessed via the Linux IPoIB framework.

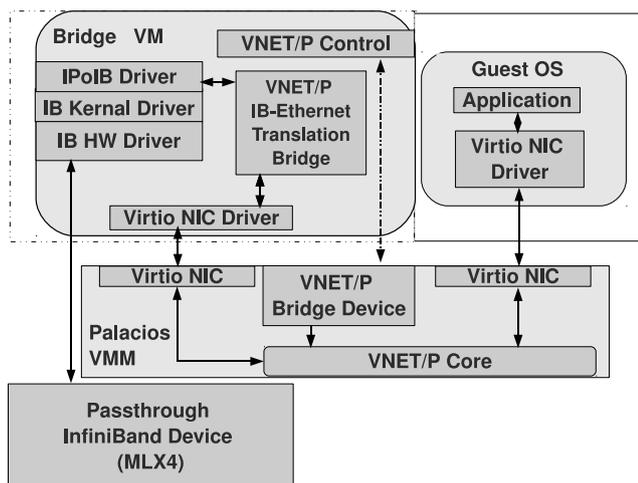


Fig. 17 The architecture of VNET/P for Kitten running on InfiniBand. Palacios is embedded in the Kitten lightweight kernel and forwards packets to/receives packets from a service VM called the bridge VM

We conducted preliminary performance tests of VNET/P for Kitten on InfiniBand using 8900 byte TCP payloads running on `ttcp` on a testbed similar to the one described in Sect. 5.1. Here, each node was a dual quad-core 2.3 GHz 2376 AMD Opteron machine with 32 GB of RAM and a Mellanox MT26428 InfiniBand NIC in a PCI-e slot. The InfiniBand NICs were connected via a Mellanox MTS 3600 36-port 20/40 Gbps InfiniBand switch. This version of VNET/P is able to achieve 4.0 Gbps end-to-end TCP throughput, compared to 6.5 Gbps when run natively on top of IP-over-InfiniBand in Reliable Connected (RC) mode.

The Kitten version of VNET/P described here is the basis of VNET/P+, which is described in more detail elsewhere [5]. VNET/P+ uses two techniques, optimistic interrupts and cut-through forwarding, to increase the throughput and lower the latency compared to VNET/P. Additionally, by leveraging the low-noise environment of Kitten, it is able to provide very little jitter in latency compared to the Linux version. When run with a 10 Gbps Ethernet network, native throughput is achievable. We are currently back-porting the VNET/P+ techniques into the Linux version.

7 Conclusion and future work

We have described the VNET model of overlay networking in a distributed virtualized computing environment and our efforts in extending this simple and flexible model to support tightly-coupled high performance computing applications running on high-performance networking hardware in current supercomputing environments, future data centers, and future clouds. VNET/P is our design and implementation of VNET for such environments. Its design goal is to achieve near-native throughput and latency on 1 and

10 Gbps Ethernet, InfiniBand, Cray Gemini, and other high performance interconnects.

To achieve performance, VNET/P relies on several key techniques and systems, including lightweight virtualization in the Palacios virtual machine monitor, high-performance I/O, and multicore overlay routing support. Together, these techniques enable VNET/P to provide a simple and flexible level 2 Ethernet network abstraction in a large range of systems no matter what the actual underlying networking technology is. While our VNET/P implementation is tightly integrated into our Palacios virtual machine monitor, the principles involved could be used in other environments as well.

We are currently working to further enhance VNET/P's performance through its guarded privileged execution directly in the guest, including an uncooperative guest.

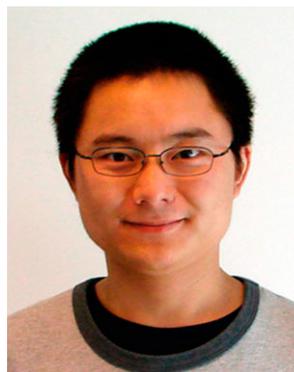
Acknowledgements We would like to thank Kevin Pedretti and Kyle Hale for their efforts in bringing up Palacios and VNET/P under CNL on the Cray XK6. This project is made possible by support from the United States National Science Foundation (NSF) via grants CNS-0709168 and CNS-0707365, and by the Department of Energy (DOE) via grant DE-SC0005343. Yuan Tang's visiting scholar position at Northwestern was supported by the China Scholarship Council.

References

1. Abu-Libdeh, H., Costa, P., Rowstron, A., O'Shea, G., Donnelly, A.: Symbiotic routing in future data centers. In: Proceedings of SIGCOMM, August (2010)
2. AMD Corporation: AMD64 Virtualization Codenamed "Pacific" Technology: Secure Virtual Machine Architecture Reference Manual, May (2005)
3. Andersen, D., Balakrishnan, H., Kaashoek, F., Morris, R.: Resilient overlay networks. In: Proceedings of SOSP, March (2001)
4. Bavier, A.C., Feamster, N., Huang, M., Peterson, L.L., Rexford, J.: In vini veritas: realistic and controlled network experimentation. In: Proceedings of SIGCOMM, September (2006)
5. Cui, Z., Xia, L., Bridges, P., Dinda, P., Lange, J.: Optimizing overlay-based virtual networking through optimistic interrupts and cut-through forwarding. In: Proceedings of the ACM/IEEE International Conference on High Performance Computing, Networking, Storage and Analysis (SC'12—Supercomputing), November (2012)
6. Dinda, P., Sundararaj, A., Lange, J., Gupta, A., Lin, B.: Methods and Systems for Automatic Inference and Adaptation of Virtualized Computing Environments, March 2012. United States patent number 8,145,760
7. Evangelinos, C., Hill, C.: Cloud computing for parallel scientific HPC applications: feasibility of running coupled atmosphere-ocean climate models on Amazon's EC2. In: Proceedings of Cloud Computing and Its Applications (CCA), October (2008)
8. Figueiredo, R., Dinda, P.A., Fortes, J.: A case for grid computing on virtual machines. In: Proceedings of the 23rd International Conference on Distributed Computing Systems (ICDCS 2003), May (2003)
9. Gabriel, E., Fagg, G.E., Bosilca, G., Angskun, T., Dongarra, J.J., Squyres, J.M., Sahay, V., Kambadur, P., Barrett, B., Lumsdaine, A., Castain, R.H., Daniel, D.J., Graham, R.L., Woodall, T.S.: open MPI

10. Ganguly, A., Agrawal, A., Boykin, P.O., Figueiredo, R.: IP over P2P: enabling self-configuring virtual IP networks for grid computing. In: Proceedings of the 20th IEEE International Parallel and Distributed Processing Symposium (IPDPS), April (2006)
11. Gordon, A., Amit, N., Har'El, N., Ben-Yehuda, M., Landau, A., Schuster, A., Tsafir, D.: ELI: bare-metal performance for I/O virtualization. In: Proceedings of the 17th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 2012), March (2012)
12. Greenberg, A., Hamilton, J.R., Jain, N., Kandula, S., Kim, C., Lahiri, P., Maltz, D.A., Patel, P., Sengupta, S.: VL2: a scalable and flexible data center network. In: Proceedings of SIGCOMM, August (2009)
13. Guo, C., Lu, G., Li, D., Wu, H., Zhang, X., Shi, Y., Tian, C., Zhang, Y., Lu, S.: Bcube: a high performance, server-centric network architecture for modular data centers. In: Proceedings of SIGCOMM, August (2009)
14. Gupta, A.: Black Box Methods for Inferring Parallel Applications' Properties in Virtual Environments. PhD thesis, Northwestern University, May 2008. Technical report NWU-EECS-08-04, Department of Electrical Engineering and Computer Science
15. Gupta, A., Dinda, P.A.: Inferring the topology and traffic load of parallel programs running in a virtual machine environment. In: Proceedings of the 10th Workshop on Job Scheduling Strategies for Parallel Processing (JSSPP), June (2004)
16. Gupta, A., Zangrilli, M., Sundararaj, A., Huang, A., Dinda, P., Lowekamp, B.: Free network measurement for virtual machine distributed computing. In: Proceedings of the 20th IEEE International Parallel and Distributed Processing Symposium (IPDPS) (2006)
17. Hua Chu, Y., Rao, S., Sheshan, S., Zhang, H.: Enabling conferencing applications on the Internet using an overlay multicast architecture. In: Proceedings of ACM SIGCOMM (2001)
18. Huang, W., Liu, J., Abali, B., Panda, D.: A case for high performance computing with virtual machines. In: Proceedings of the 20th ACM International Conference on Supercomputing (ICS), June–July (2006)
19. Innovative Computing Laboratory: HPC challenge benchmark. <http://icl.cs.utk.edu/hpc/>
20. Intel: Intel cluster toolkit 3.0 for Linux. <http://software.intel.com/en-us/articles/intel-mpi-benchmarks/>
21. Jiang, X., Xu, D.: Violin: virtual internetworking on overlay infrastructure. Tech. rep. CSD TR 03-027, Department of Computer Sciences, Purdue University, July (2003)
22. Joseph, D.A., Kannan, J., Kubota, A., Lakshminarayanan, K., Stolica, I., Wehrle, K.: Ocala: an architecture for supporting legacy applications over overlays. In: Proceedings of the 3rd Symposium on Networked Systems Design and Implementation (NSDI), May (2006)
23. Kallahalla, M., Uysal, M., Swaminathan, R., Lowell, D.E., Wray, M., Christian, T., Edwards, N., Dalton, C.I., Gittler, F.: Softudc: a software-based data center for utility computing. *Computer* **37**(11), 38–46 (2004)
24. Kashyap, V.: IP over Infiniband (IPoIB) Architecture. IETF Network Working Group Request for Comments. RFC 4392, April (2006). Current expiration: August 2012
25. Kim, C., Caesar, M., Rexford, J.: Floodless in Seattle: a scalable Ethernet architecture for large enterprises. In: Proceedings of SIGCOMM, August (2008)
26. Kumar, S., Raj, H., Schwan, K., Ganev, I.: Re-architecting VMMS for multicore systems: the sidecore approach. In: Proceedings of the 2007 Workshop on the Interaction Between Operating Systems and Computer Architecture, June (2007)
27. Lange, J., Dinda, P.: Transparent network services via a virtual traffic layer for virtual machines. In: Proceedings of the 16th IEEE International Symposium on High Performance Distributed Computing (HPDC), June (2007)
28. Lange, J., Dinda, P., Hale, K., Xia, L.: An introduction to the Palacios virtual machine monitor—release 1.3. Tech. rep. NWU-EECS-11-10, Department of Electrical Engineering and Computer Science, Northwestern University, October (2011)
29. Lange, J., Pedretti, K., Dinda, P., Bae, C., Bridges, P., Soltero, P., Merritt, A.: Minimal-overhead virtualization of a large scale supercomputer. In: Proceedings of the 2011 ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments (VEE), March (2011)
30. Lange, J., Pedretti, K., Hudson, T., Dinda, P., Cui, Z., Xia, L., Bridges, P., Gocke, A., Jaconette, S., Levenhagen, M., Brightwell, R.: Palacios and Kitten: new high performance operating systems for scalable virtualized and native supercomputing. In: Proceedings of the 24th IEEE International Parallel and Distributed Processing Symposium (IPDPS), April (2010)
31. Lange, J., Sundararaj, A., Dinda, P.: Automatic dynamic run-time optical network reservations. In: Proceedings of the 14th International Symposium on High Performance Distributed Computing (HPDC), July (2005)
32. Lin, B., Dinda, P.: Vsched: mixing batch and interactive virtual machines using periodic real-time scheduling. In: Proceedings of ACM/IEEE SC (Supercomputing), November (2005)
33. Lin, B., Sundararaj, A., Dinda, P.: Time-sharing parallel applications with performance isolation and control. In: Proceedings of the 4th IEEE International Conference on Autonomic Computing (ICAC), June (2007)
34. Liu, J., Huang, W., Abali, B., Panda, D.: High performance VMM-Bypass I/O in virtual machines. In: Proceedings of the USENIX Annual Technical Conference, May (2006)
35. Mahalingam, M., Dutt, D., Duda, K., Agarwal, P., Kreeger, L., Sridhar, T., Bursell, M., Wright, C.: VXLAN: A framework for overlaying virtualized layer 2 networks over layer 3 networks. IETF Network Working Group Internet Draft, February (2012). Current expiration: August 2012
36. Menon, A., Cox, A.L., Zwaenepoel, W.: Optimizing network virtualization in Xen. In: Proceedings of the USENIX Annual Technical Conference (USENIX), May (2006)
37. Mergen, M.F., Uhlig, V., Krieger, O., Xenidis, J.: Virtualization for high-performance computing. *Oper. Syst. Rev.* **40**(2), 8–11 (2006)
38. Mysore, R.N., Pamboris, A., Farrington, N., Huang, N., Miri, P., Radhakrishnan, S., Subramanya, V., Vahdat, A.: Portland: a scalable fault-tolerant layer 2 data center network fabric. In: Proceedings of SIGCOMM, August (2009)
39. Nurmi, D., Wolski, R., Grzegorzczak, C., Obertelli, G., Soman, S., Youseff, L., Zagorodnov, D.: The Eucalyptus open-source cloud-computing system. In: Proceedings of the 9th IEEE/ACM International Symposium on Cluster Computing and the Grid (CCGrid), May (2009)
40. Ostermann, S., Iosup, A., Yigitbasi, N., Prodan, R., Fahringer, T., Epema, D.: An early performance analysis of cloud computing services for scientific computing. Tech. Rep. PDS2008-006, Delft University of Technology, Parallel and Distributed Systems Report Series, December (2008)
41. Raj, H., Schwan, K.: High performance and scalable i/o virtualization via self-virtualized devices. In: Proceedings of the 16th IEEE International Symposium on High Performance Distributed Computing (HPDC), July (2007)
42. Russell, R.: Virtio: towards a de-facto standard for virtual I/O devices. *Oper. Syst. Rev.* **42**(5), 95–103 (2008)
43. Ruth, P., Jiang, X., Xu, D., Goasguen, S.: Towards virtual distributed environments in a shared infrastructure. *Computer* **38**(5), 63–69 (2005)
44. Ruth, P., McGachey, P., Jiang, X., Xu D.: Viocluster: virtualization for dynamic computational domains. In: Proceedings of the IEEE International Conference on Cluster Computing (Cluster), September (2005)

45. Shafer, J., Carr, D., Menon, A., Rixner, S., Cox, A.L., Zwaenepoel, W., Willmann, P.: Concurrent direct network access for virtual machine monitors. In: Proceedings of the 13th International Symposium on High Performance Computer Architecture (HPCA), February (2007)
46. Stoica, I., Morris, R., Karger, D., Kaashoek, F., Balakrishnan, H.: Chord: a scalable peer-to-peer lookup service for Internet applications. In: Proceedings of ACM SIGCOMM 2001, pp. 149–160 (2001)
47. Sugerma, J., Venkitachalan, G., Lim, B.-H.: Virtualizing I/O devices on VMware workstation's hosted virtual machine monitor. In: Proceedings of the USENIX Annual Technical Conference, June (2001)
48. Sundararaj, A.: Automatic Run-time, and Dynamic Adaptation of Distributed Applications Executing in Virtual Environments. PhD thesis, Northwestern University, December (2006). Technical report NWU-EECS-06-18, Department of Electrical Engineering and Computer Science
49. Sundararaj, A., Dinda, P.: Towards virtual networks for virtual machine grid computing. In: Proceedings of the 3rd USENIX Virtual Machine Research and Technology Symposium (VM 2004), May (2004). Earlier version available as technical report NWU-CS-03-27, Department of Computer Science, Northwestern University
50. Sundararaj, A., Gupta, A., Dinda, P.: Increasing application performance in virtual environments through run-time inference and adaptation. In: Proceedings of the 14th IEEE International Symposium on High Performance Distributed Computing (HPDC), July (2005)
51. Sundararaj, A., Sanghi, M., Lange, J., Dinda, P.: An optimization problem in adaptive virtual environments. In: Proceedings of the Seventh Workshop on Mathematical Performance Modeling and Analysis (MAMA), June (2005)
52. Tsugawa, M.O., Fortes, J.A.B.: A virtual network (vine) architecture for grid computing. In: 20th International Parallel and Distributed Processing Symposium (IPDPS), April (2006)
53. Uhlig, R., Neiger, G., Rodgers, D., Santoni, A., Martin, F., Anderson, A., Bennett, S., Kagi, A., Leung, F., Smith, L.: Intel virtualization technology. *IEEE Computer*, 48–56 (2005)
54. Van der Wijngaart, R.: NAS parallel benchmarks version 2.4. Tech. rep. NAS-02-007, NASA Advanced Supercomputing (NAS Division), NASA Ames Research Center, October (2002)
55. Vaughan, C., Rajan, M., Barrett, R., Doerfler, D., Pedretti, K.: Investigating the impact of the Cielo Cray XE6 architecture on scientific application codes. In: Proceedings of the 2011 IEEE International Symposium on Parallel and Distributed Processing Workshops and PhD Forum (IPDPSW 2011) (2011)
56. Wolinsky, D., Liu, Y., Juste, P.S., Venkatasubramanian, G., Figueiredo, R.: On the design of scalable, self-configuring virtual networks. In: Proceedings of 21st ACM/IEEE International Conference of High Performance Computing, Networking, Storage, and Analysis (SuperComputing—SC), November (2009)
57. Xia, L., Cui, Z., Lange, J., Tang, Y., Dinda, P., Bridges, P.: VNET/P: bridging the cloud and high performance computing through fast overlay networking. In: Proceedings of the 21st ACM Symposium on High-Performance Parallel and Distributed Computing (HPDC 2012), June (2012)
58. Xia, L., Kumar, S., Yang, X., Gopalakrishnan, P., Liu, Y., Schoenberg, S., Guo, X.: Virtual WiFi: bring virtualization from wired to wireless. In: Proceedings of the 7th International Conference on Virtual Execution Environments (VEE'11) (2011)
59. Xia, L., Lange, J., Dinda, P., Bae, C.: Investigating virtual passthrough I/O on commodity devices. *Oper. Syst. Rev.* **43**(3), 83–94 (July 2009). Initial version appeared at WIOV 2008



Lei Xia is currently a PhD candidate in the Department of Electrical Engineering and Computer Science at Northwestern University. He holds a B.S. in Geophysics and M.S. in Computer Science, both from Nanjing University, China. His research focuses in virtualization for high performance parallel and distributed systems, as well as cloud systems.



Zheng Cui is a PhD candidate in Computer Science at University of New Mexico. She received the M.S.(2007) in Computer Science from University of New Mexico, and B.E. in Computer Science from Zhengzhou University (2003). Her research interests include large-scale distributed and parallel systems, system software, virtualization, high-performance virtual networking, and networking for HPC and cloud computing.



John Lange is currently an assistant professor in the Department of Computer Science at the University of Pittsburgh. Prior to joining the faculty at the University of Pittsburgh, he received a BSc, MSc and PhD in Computer Science as well as a BSc in Computer Engineering from Northwestern University. His research focuses in HPC and operating systems, as well as networking, virtualization and distributed systems. His current focus lies in the area of multi-instance system software architectures for high performance computing and cloud environments, as well as operating system optimizations for consolidated virtual environments.



Yuan Tang is an assistant professor of computer science at Chengdu University of Technology (CDUT) in China. He graduated from Hebei Finance University, China, in 1999, got his master degree from Kunming University of Science and Technology (KMUST), China, in 2007, and received his Ph.D. from University of Electronic Science and Technology of China (UESTC) in 2012. During 2008 to 2010 he was a visiting scholar in department of EECS, Northwestern University. His research has mainly focused on operating systems, virtualization and cloud computing.



Peter Dinda is a professor in the Department of Electrical Engineering and Computer Science at Northwestern University, and head of its Computer Engineering and Systems division, which includes 17 faculty members. He holds a B.S. in electrical and computer engineering from the University of Wisconsin and a Ph.D. in computer science from Carnegie Mellon University. He works in experimental computer systems, particularly parallel and distributed systems. His research currently involves virtualization

for distributed and parallel computing, programming languages for parallel computing, and empathic systems for bridging individual user satisfaction and systems-level decision-making. You can find out more about him at pdinda.org.



Patrick Bridges is an associate professor in the Department of Computer Science at the University of New Mexico. He received his BSc in Computer Science from Mississippi State University and his PhD in Computer Science from the University of Arizona. His overall research interests are in system software for large-scale computing systems, particularly large parallel and distributed systems. His recent research in this area has focused on techniques for leveraging virtualization for large-scale systems and on

fault tolerance and resilience issues in next-generation supercomputing systems.