

On Node Virtualization for Scalable Network Emulation

Steffen Maier, Daniel Herrscher, Kurt Rothermel

University of Stuttgart, Institute of Parallel and Distributed Systems (IPVS)

Universitätsstr. 38, D-70569 Stuttgart, Germany

Phone: +49(711)7816-226, Fax: -424

maier@informatik.uni-stuttgart.de

Keywords: software performance evaluation, network emulation, mobile ad hoc networks, scalability, virtual routing.

Abstract

During the development of network protocols and distributed applications, their performance has to be analyzed in appropriate environments. Network emulation testbeds provide a synthetic, configurable network environment for comparative performance measurements of real implementations. Realistic scenarios have to consider hundreds of communicating nodes. Common network emulation approaches limit the number of nodes in a scenario to the number of computers in an emulation testbed. To overcome this limitation, we introduce a virtual node concept for network emulation. The key problem for node virtualization is a transparent, yet efficient separation of node resources. In this paper, we provide a brief survey of candidate node virtualization approaches to facilitate scalable network emulation. Based on the gathered insights, we propose a lightweight virtualization solution to achieve maximum scalability and discuss the main points regarding its implementation. We present extensive evaluations that show the scalability and transparency of our approach in both a traditional wired infrastructure-based, and in a wireless ad hoc network emulation scenario. The measurements indicate that our solution can push the upper limit of emulation scenario sizes by a factor of 20 to 30. Given our emulation testbed consisting of 64 computers, this translates to possible scenario sizes of up to 1920 nodes.

I. INTRODUCTION

During the design and implementation of distributed applications and network protocols, it is essential to analyze the impact of various network environments on their performance. While mathematical analysis and simulations are commonly used in early design stages, measurements are used to check the theoretical results as soon as implementations become available. Such measurements usually compare the performance of one implementation in different network environments or of different implementations in the same network environment.

Comparative performance measurements in real environments are considered problematic for two reasons. First, es-

pecially in scenarios with mobile nodes and wireless networking, it is hard to obtain multiple comparable measurement runs. Secondly, resource requirements prohibit measurements in larger scenarios. Therefore, there is strong demand for synthetic network environments that can be parametrized in order to reproduce an original or fictitious network.

The process of introducing network properties that differ from the actual properties of the hardware in use is called *network emulation*. A *network emulation tool* is software capable of altering network traffic in a specified way. A facility consisting of a combination of flexible networking hardware and suitable emulation tools is called *network emulation testbed*. Network protocols and distributed applications subjected to performance measurements in a network emulation testbed are called *software under test*.

Comparative performance measurements for mobile computing scenarios, e.g. the evaluation of an ad hoc routing protocol, typically require large scenarios with hundreds of nodes. The analysis of new applications for traditional infrastructure-based networks, e.g. a large-scale location service, may also require a high number of nodes, since both the end systems and the intermediate systems of the underlying infrastructure have to be considered.

Common network emulation systems assume that one communicating node in an emulation scenario corresponds to one physical computer in an emulation testbed. This severely limits the scalability, since testbeds with the required number of hundreds of computers are typically not available.

However, a number of applications aiming at resource-poor devices, e.g. in mobile computing scenarios, only need a fraction of the resources that a testbed node can provide. Therefore, we propose to run several instances of the software under test on a single testbed node ("physical node," *pnode*). Each instance of the software under test has to be provided a separate execution environment ("virtual node," *vnode*). In this paper, we provide a brief survey of candidate approaches for node virtualization. Based on these approaches, we present a transparent, yet lightweight and thus very scalable solution to node virtualization for network emulation testbeds. Our implementation not only supports scalable emulation of networks consisting of point to point links but also

shared media based networks such as mobile ad hoc networks.

The remainder of this paper is structured as follows. The Network Emulation Testbed, which we use as a basis for our scalable network emulation approach, is introduced in Section II. In Section III, we provide a brief survey of candidate node virtualization approaches. We choose one of the candidate approaches for our implementation, which we discuss in Section IV. In Section V, we provide extensive measurements showing the scalability of our approach for two important kinds of scenarios: emulation of infrastructure-based networks and MANETs (mobile ad hoc networks). Furthermore, we discuss the achievable degree of transparency for the software under test. We discuss related work in Section VI. Finally, we conclude the paper in Section VII.

II. OVERVIEW OF THE NETWORK EMULATION TESTBED

The Network Emulation Testbed (NET) [1] at the University of Stuttgart provides the basis for our scalable network emulation approach. It consists of 64 PC-nodes connected by a monolithic, programmable gigabit switch, and a separate administration network for setup and control (see Fig. 1). Using IEEE 802.1Q VLAN (virtual LAN) technology, the gigabit switch is able to create an arbitrary connection topology between the nodes. Each point-to-point link or shared media network segment in an emulation scenario, e.g. a WLAN (wireless LAN) channel, is mapped to a uniquely tagged VLAN.

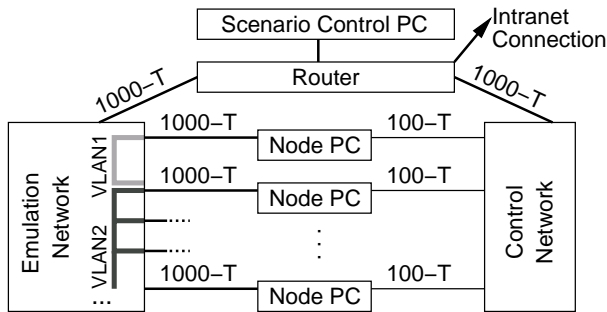


Figure 1. The Network Emulation Testbed.

Custom network emulation tools running on each node introduce the desired artificial network properties. They provide the service level abstraction of an unreliable datagram service to the software under test (see Fig. 2). This is the lowest possible emulation abstraction feasible to be implemented in software. The tool is implemented as a virtual network device driver, and therefore completely transparent to implementations on the network layer. As a result, the protocol

stack including the network layer and all higher layers can be considered as software under test.

On a testbed node, several VLANs represent several virtual network interfaces, each of which is assigned a separate instance of the emulation tool. The tool enables the configuration of arbitrary bandwidth limitations, delays, and frame error loss ratios. Additionally, to enable the realistic emulation of shared media networks, the effects of a MAC (media access control) layer can be reproduced. At the present time, this tool is capable of emulating IEEE 802.3 (Ethernet) [2]. We are currently extending the tool to allow the emulation of the ad hoc mode of IEEE 802.11 WLAN.

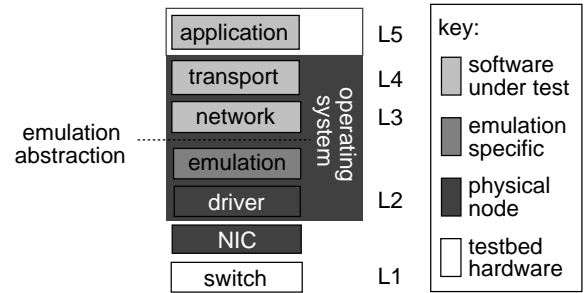


Figure 2. Software under test and network emulation tools on a physical node in NET.

A central scenario controller dynamically updates the parameters of the distributed network emulation tools. For MANET emulation, this includes changing connection quality and thus frame error rates between communicating nodes. The connection quality is automatically derived from the simulated node mobility [3].

Without node virtualization, each node in an emulation scenario is mapped to a physical node of NET, which limits the scenario size to 64.

III. APPROACHES TO NODE VIRTUALIZATION

In general, node virtualization provides a way to schedule formerly exclusive hardware resources to a number of consumers. With respect to network emulation, our consumers are execution environments for software under test, which is to be subjected to emulated network properties. We derive the following requirements from node virtualization and network emulation:

1. Our paramount goal is scalability. This requires *minimal virtualization overhead* in order to preserve resources for the software under test.
2. If two (or more) vnodes inside the same pnode communicate, they should make use of *efficient intra-pnode communication*. This requirement supports our paramount

goal 1 by minimizing overhead due to the virtualization process.

3. An execution environment introduced by node virtualization should be as *transparent* as possible for the software under test. This is important to support performance measurement of unmodified real implementations.

In the following, we present candidate node virtualization approaches and discuss their suitability regarding scalable network emulation. They all have in common that they allow multiple instances of software under test to be executed on top of the emulation abstraction interface shown in Fig. 2. For the discussion, we assume that the network stack is part of the kernel space, as is prevalent in commodity operating systems. Finally, we evaluate each approach for its suitability based on our requirements. The presented approaches can be classified into two main categories: virtual machines and virtual network stacks.

A. Virtual Machine

A straightforward way to introduce node virtualization is using a virtual machine (VM) approach. Instead of running an operating system (OS) directly on the bare hardware, a shim of software is inserted in between. This software – the virtual machine monitor (VMM) – schedules access of multiple guest operating systems to exclusive hardware resources managed by the VMM (Fig. 3).

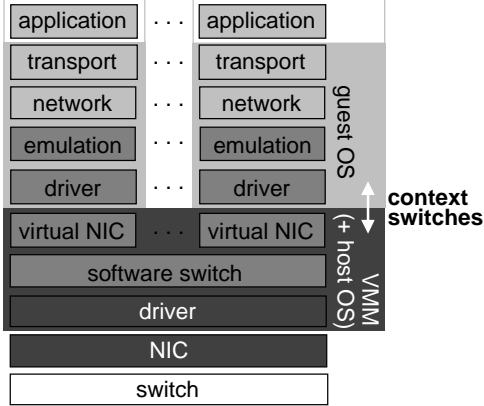


Figure 3. Virtual machine approach.

In order to support emulated network parameters, we need to insert our emulation tool on top of network interface drivers inside each guest operating system. For communication with other vnodes on the same and other pnodes, a software switch forwards frames correspondingly.

1. Classical Virtual Machine

Classical virtual machines such as Plex86 [4] or VMware Workstation [5] have in common that they support unmodified operating systems, and thus network stacks, in each guest

instance. Therefore, they provide transparency with respect to software under test. However, context switches between guest OS and VMM happen whenever privileged commands are trapped. Since network communication causes such context switches, classical VMs imply considerable virtualization overhead limiting scalability. This is especially an issue for VMs, that do not virtualize a certain kind of system hardware, but e.g. the system call interface of the host OS, e.g. User-Mode Linux (UML) [6] or UMLinux (now FAUmachine) [7]. Even with a modified host OS such VMs only show comparable performance to e.g. VMware [8].

2. Lightweight Virtual Machine

Lightweight virtual machines such as VMware ESX [9], Denali [10] or Xen [11] directly access the host hardware without a host OS in order to reduce virtualization overhead. However, they may require custom or ported guest operating systems and are thus only partly transparent.

B. Virtual Network Stack

The virtual machine approach described in the previous subsection actually virtualizes more than is needed for network emulation. It would be sufficient to provide virtual execution environments for just the software under test, i.e. for exactly those layers above the emulation abstraction interface (Fig. 2). This can be accomplished with virtual network stacks (Fig. 4). In order to extend the virtualization of network and transport layer also to the application layer, sets of processes get associated with a certain network stack instance. Consequently, a vnode consists of the following sets: a set of processes on application layer, a set of sockets on transport layer, and a routing table on network layer.

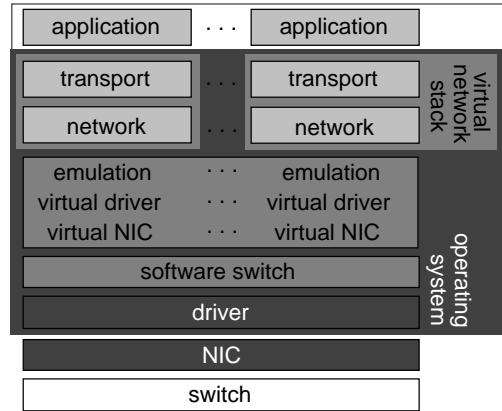


Figure 4. Virtual network stack approach.

In contrast to virtual machines, there is no more need for separate virtual network devices and their drivers. The emulation tool itself can appear as several instances of a virtual network device. In tight cooperation, a software switch for-

wards frames appropriately in order to allow communication between any vnodes. The virtualization overhead for virtual network stacks is as low as possible. Compared to virtual machine approaches, there are no redundant context switches and copy operations.

1. Duplicated Network Stack

Duplicated network stacks such as vimage [12] allow the flexible execution of different network stack implementations on the same pnode. However, they need extensive modifications to become fully virtualized and are thus hardly transparent.

2. Virtual Routing

Virtual routing [13], [14] requires only the essential variables, that have to be allocated separately for each stack instance, to be touched. Thus, virtual routing is more transparent than duplicated network stacks. Though multiple instances are supported, only one specific implementation of a stack can be executed on a single pnode at a time. Yet, using different implementations on different pnodes remedies this partial flexibility.

C. Summary and Selected Approach

Tab. 1 shows a summary of the discussion in the previous subsections. We rate each approach on a scale of three levels with plus denoting good fulfillment, a circle denoting partial fulfillment, and minus denoting restrictions with respect to our requirements.

Table 1: Comparison of candidate virtualization approaches.

virtualization approach	scalability	efficiency	transparency
A.1. classical virtual machine	–	–	+
A.2. lightweight virtual machine	◦	◦	◦
B.1. duplicated network stack	+	+	–
B.2. virtual routing	+	+	◦

For deriving the most suitable approach, we evaluate the fulfillment of our requirements in order of descending priority. While virtual machine approaches can be fully transparent and flexible, they do not fully comply with our paramount goal scalability. Virtual network stack approaches fulfill the requirement of low virtualization overhead and efficient intra-pnode communication. Of the two alternatives, virtual routing is more transparent. We thus consider virtual routing best suited for scalable network emulation.

IV. IMPLEMENTATION

Virtual routing as discussed in the previous section fits our requirements best. Hence, we choose virtual routing along with a custom software switch that enables communication between any vnodes in an emulation scenario. Linux 2.4 serves as operating system for the implementations. In the following, we describe the two main blocks of our approach traversing the layers from bottom to top.

A. Software Communication Switch

In the context of our network emulation testbed, each software switch introduces a “stacked” sub-switch using the emulation network connection as an uplink to the emulation switch (cp. to Fig. 1). A software switch resembles the functionality of a hardware Ethernet switch. It mediates both between vnodes located on the same pnode as well as between vnodes located on different pnodes. This provides transparent switching between any vnodes in a scenario.

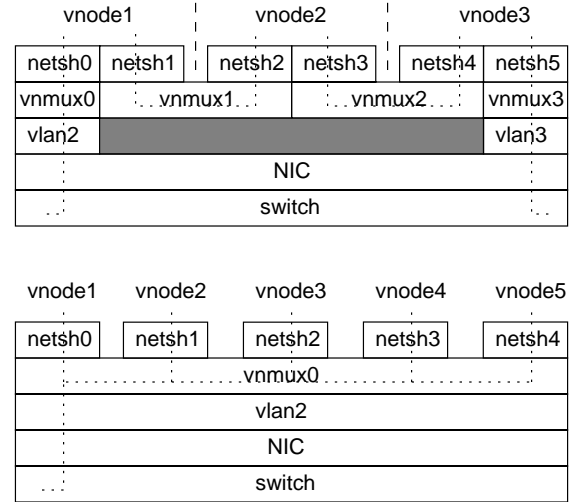


Figure 5. Pnode configuration examples: link based (top), shared media based, e.g. MANET, (bottom).

In contrast to the software bridge already existing in Linux, we need one uplink to a real device and multiple local ends. Therefore we designed a custom Linux kernel module providing instances of a virtual switch network device “vnmux” (virtual node multiplexer) (Fig. 5). In order to get an uplink, this device can be internally bound to the driver of a real network device (NIC). The latter could also be a tagged VLAN device which is in turn bound to a real network device. The bound device is put into promiscuous mode to be able to receive frames destined for local vnode devices. The emulation hardware switch takes care of filtering. It delivers only locally targeted frames after its learning phase, so that the software switch only has to process frames it really is responsible for.

Processing frames is done without extra copying of payload data. This is essential to fulfill our requirement 2 for efficient intra-pnode communication. Switching decisions work with constant destination lookup time resulting in a complexity of $O(1)$. At the upper interface of the switch, virtual network devices provided by emulation tool instances (“netshX”) register themselves to generate local switch “ports.”

B. Virtual Routing

Virtual routing instances are located on top of our emulation tool’s virtual network devices. Those virtual routing instances and applications on top of them represent possible software under test. We base our implementation on kernel patches for “Virtual Routing and Forwarding” (VRF) [15] version 0.100 by James R. Leu. VRF provides multiple instances of forwarding information bases as well as mechanisms to associate network devices, IPv4 UDP/TCP sockets, and processes with instances. User space tools exist for instance maintenance and for associating devices and processes with instances. Despite all these features, virtual routing is still not sufficiently transparent for application processes and common routing daemon implementations.

Therefore, we extend system interfaces that operate on routing tables – some IOCTLs and the protocol route-netlink – to work on the specific routing table of the VRF instance the calling process is associated with. Additionally, we extend the `ip_queue` feature of the protocol netlink-firewall to allow queueing of IPv4 packets to one process within *each* VRF instance. Thereby we gain full transparency for unmodified network applications including routing daemons, that potentially need to overhear certain packets with the help of `ip_queue`.

To implement all of the above mentioned functionality, only limited modifications to a standard Linux 2.4.24 source tree are necessary. The modifications comprise 1409 lines of code, which consist of 416 additions, 980 changes, and 13 deletions.

V. EVALUATION

In the following we provide an extensive evaluation of the building blocks as well as of the complete implemented system. All measurements are performed on pnodes in our test-bed equipped with an Intel Pentium 4 2.4 GHz processor, 133 MHz frontside bus, 512 MB main memory, and an Altima AC9100 Gigabit Ethernet adapter in a 32 bit, 33 MHz PCI bus. Passing through the different network stack layers from bottom to top, we start our evaluation with the software communication switch at the data link layer and show the accuracy of our network emulation tool in variably virtualized scenarios. Network and transport layer are treated twice due to two considered types of network requiring different routing algorithms: first for a wired infrastructure based network, sec-

only for a wireless ad-hoc network. The evaluation aims at showing the scalability of the system by comparing the non-virtualized cases to variably virtualized cases of the same scenarios. We would like to point out that the software under test is by no means limited to protocols on the network layer. After all, our load generators are processes on application layer communicating through sockets with the transport layer. Of course, more complex applications such as peer to peer systems can also be analyzed in our emulation environment without modification.

A. Software Communication Switch

Our software communication switch is a core component in our scalable emulation environment, since it has to switch frames quickly and at the lowest overhead possible. In order to show that it fulfills the expectations, we measure both the duration of switching decisions and the resulting throughput.

The scenario for measuring the duration of switching decisions consists of two pnodes connected by a point to point link provided by a tagged VLAN. One of the pnodes hosts one switch instance as the test subject. We vary the number of vnode devices attached to local ports of the switch instance between 1 and 64. The other pnode generates load by injecting randomly sized frames targeted at the software switch ports.

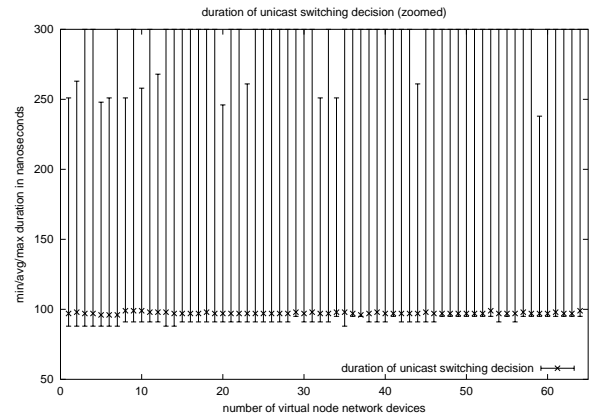


Figure 6. Duration of unicast switching decision versus number of vnode devices.

Fig. 6 shows the efficiency of unicast switching decisions. The average measured duration is about 98 ns, independent of the number of vnode devices per switch. The profiled machine code comprises 24 instructions, which take at least about 50 CPU clock cycles on the superscalar out-of-order core of an Intel Pentium 4, if all data is available in the first level cache [16]. At a frequency of 2.4 GHz, 50 clock cycles take about 20 ns. This marks a lower bound for the execution time. Taking cache misses into account, our measured average duration constitutes a reasonable value. A few spikes in

the maxima up to 1388 ns are due to cache effects having impact on such short measured intervals. Since the average is close to the minimum of 88 ns, all maxima appear rarely. We conclude from these measurements that our implementation of the switching decision, and therefore the core functionality of the software switch, is highly efficient.

The scenario for measuring switch throughput consists of one pnode with one switch instance and a varying number of vnode devices attached. As before, the switch also has an up-link to a tagged VLAN to cover all paths in the implementation. We vary frame sizes between 64 and 1500 Bytes. Fig. 7 shows constant throughput for unicast frames which only depends on the frame size. Small frames imply more overhead and thus less throughput. For comparison we measured a memory bandwidth of 1020 MByte/s with STREAM [17]. Obviously, frame handling overhead is the limiting factor in switch throughput. Nevertheless, a throughput of about 3 GBit/s can serve as an upper limit for aggregate link bandwidth inside one pnode and is 3 times larger than the external uplink over the Gigabit Ethernet network interface.

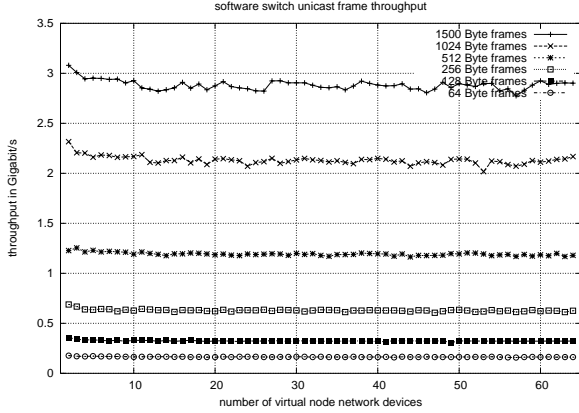


Figure 7. Unicast frame throughput versus number of vnode devices.

For broadcast frames their administration structure `sk_buff` (not the payload) has to be cloned on delivery for each local recipient. This is necessary since the receive path assumes exclusive administrative frame data structures. We also evaluate switch throughput for broadcast frames. The throughput for the starting value of 2 vnode network devices is slightly lower than for the unicast case because an additional frame clone has to be transmitted on the uplink (Fig. 8). With an increasing number of vnode devices per pnode, throughput decreases due to the overhead of cloning. Yet, aggregate switch throughput stays significantly above the memory bandwidth. However, in order to avoid any decrease, we plan to investi-

gate possible improvements by sharing administration structures of broadcast frames between vnodes on the same pnode.

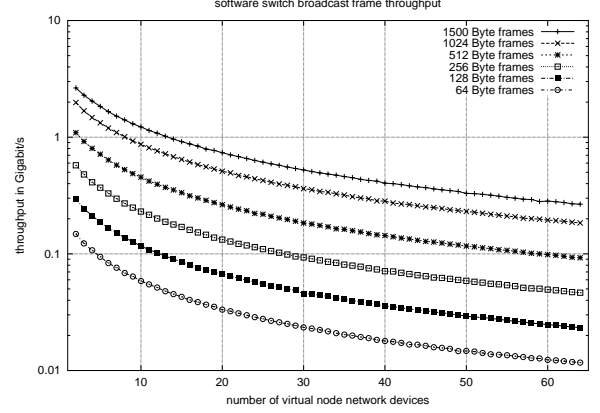


Figure 8. Broadcast frame throughput versus number of vnode devices.

B. Network Emulation Tool

Our network emulation tool is able to accurately enforce specified network properties consisting of bandwidth limitation, delay, and frame loss ratio [1]. In this section, we show that our tool remains accurate in the virtualized case up to a machine dependent limit for the degree of virtualization.

The scenario for measuring the accuracy of emulated network properties consists of a varying number of vnodes on a single pnode. n vnodes are interconnected in a chain of $n-1$ full duplex links having either limited bandwidth or specific delay in each direction (similar to Fig. 12). To measure loss ratio only one direction of each link is configured to lose frames.

In order to measure the accuracy of bandwidth limitation, we put load on each link by measuring maximum TCP throughput concurrently. Fig. 9 shows the results consisting of the measured average link bandwidth with minimum and maximum over all links, i.e. TCP flows. Depending on the number of vnodes, the specified bandwidth is accurately enforced by our network emulation tool. Up to an emulated bandwidth of 5 MBit/s, at least 64 vnodes can be hosted on a single pnode without loss of accuracy. 8 to 16 vnodes can be safely interconnected at 54 MBit/s and at least 4 vnodes can be hosted on a pnode in a Fast-Ethernet scenario with 100 MBit/s.

We measure ICMP round trip times (RTT) on each link concurrently to investigate the accuracy of delay emulation. Since the full duplex links are symmetric, the actual delay results from half the measured RTT. The results in Fig. 10 indicate that delay is emulated accurately independent of the number of vnodes per pnode. Thus, the emulation of delay scales perfectly with the degree of virtualization. The measured deviations from the average delay values stay within

bounds of 5 ms and are due to the granularity of the timer used to introduce the delay [1].

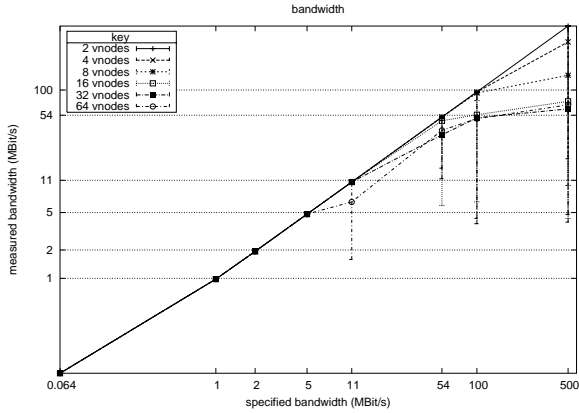


Figure 9. Enforced versus specified bandwidth for different numbers of vnodes.

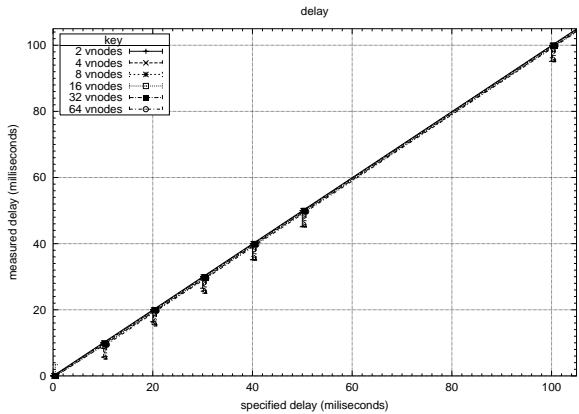


Figure 10. Enforced versus specified delay for different numbers of vnodes.

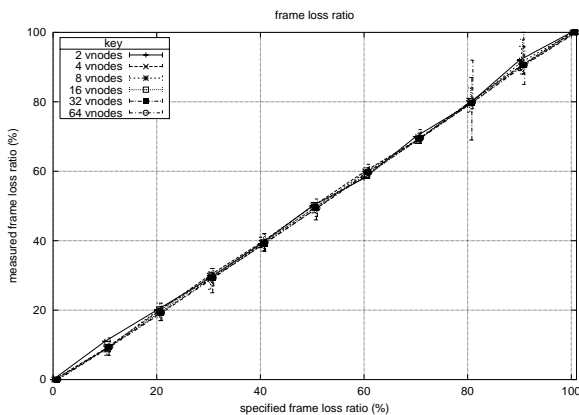


Figure 11. Enforced versus specified loss ratio for different numbers of vnodes.

Fig. 11 depicts measurement results showing the fidelity of emulated frame loss ratio. We put load on each link concurrently using adaptive ping. On average, frame loss emulation scales well with the number of vnodes per pnode. Minima and maxima – especially at a loss ratio of 80 or 90 % – are outliers and the deviation is below or equal to 2.7 for all measuring points.

We conclude from our measurements that our network emulation tool is able to accurately enforce specified network properties over a wide range of virtualization degree.

C. Wired Infrastructure Based Network Emulation

Having treated the data link layer in the previous section, we now continue our evaluation of the network and transport layer in a wired infrastructure based network emulation scenario. The system model is described first. Afterwards, we present measurement results for the network and then the transport layer. Additionally, we report on the system utilization caused by executing multiple vnodes on the same pnode.

Fig. 12 shows the network topology of the scenario. It consists of a linear chain with a varying number of router nodes. Point to point links connecting the routers are full duplex and have an emulated limited bandwidth of 100 MBit/s in each direction. Each router uses static routing table entries to reach its predecessors and successors in the chain. We conduct two types of experiments for a scenario. First, we measure in a scenario with real pnodes to obtain reference values. Secondly, we place all routers inside vnodes on a single pnode except the last router, which resides on a separate pnode without any virtualization. Thereby we show that communication over the software switch works transparently, and mixing of arbitrarily configured pnodes is possible. Note that the layers of the real network stack implementation are always traversed on communication even if the network traffic does not leave the left pnode but for the last hop. For each packet forwarding on each vnode, the network layer is traversed once on the input and once on the output path. If appropriate, traversal reaches up to the application layer on each vnode.



Figure 12. Infrastructure emulation scenario, virtualized case.

On network layer, we measure ICMP round trip time delays. The ping utility executed on the leftmost router sends ICMP echo requests through the router chain to the rightmost router. Fig. 13 shows linear increase of the mean ICMP round trip time delay with an increasing number of hops in the rout-

ing chain. The left y-axis corresponds to the results of this infrastructure based scenario. The figure also contains measurement results for the wireless ad-hoc scenario discussed in the next section. For the variant with pnodes only, we had 48 of 64 nodes available at the time of the experiment. For the virtualized variant of the scenario, the slope is more flat than with pnodes only. This is because the software switch has lower communication delay than the hardware emulation switch. The emulation tool could compensate for that, if a particular scenario requires inter-node delays to be exactly the same. Delays occur within time bounds depicted in Fig. 14. A comparison to the mean values in Fig. 13 shows that mean and minima fall close together, i.e. maxima occur rarely and are due to route cache misses.

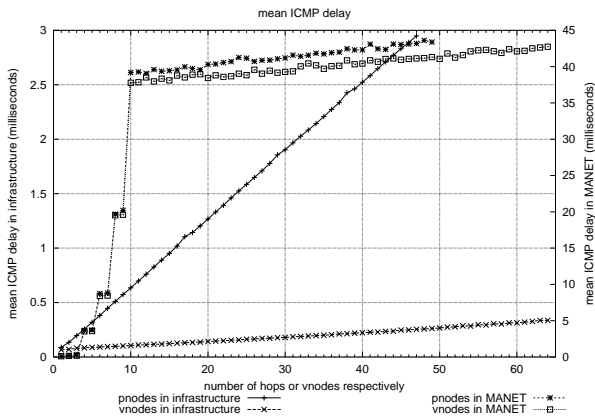


Figure 13. Round trip time mean delay versus number of vnodes.

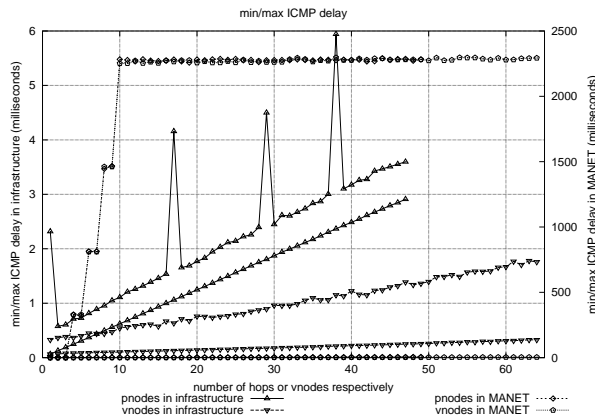


Figure 14. Round trip time minima and maxima delay versus number of vnodes.

On transport layer, we measure TCP throughput over the router chain with a block size of 1 kByte (Fig. 15). For the pnode-only chain, the throughput should stay constant over different numbers of hops. However, it starts decreasing at 41 hops. This might be due to TCP behavior on paths with a

pathologically large number of hops. For the virtualized scenario, we observe different behavior in each direction. On the reverse direction (rx) from the last router on a pnode to the first router on a vnode, throughput starts dropping at 46 vnodes due to resource contention. On the forward path, throughput drops earlier at 21 vnodes. This is due to the fact that the TCP source is virtualized: Both the sending TCP protocol with its timers in the leftmost vnode and all the other vnodes compete for the same resources of their shared pnode. This confirms that TCP adheres to the smart sender/dumb receiver protocol design rule [18].

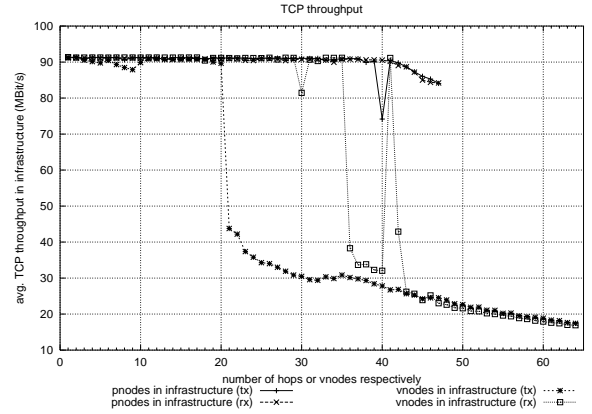


Figure 15. Throughput versus number of vnodes.

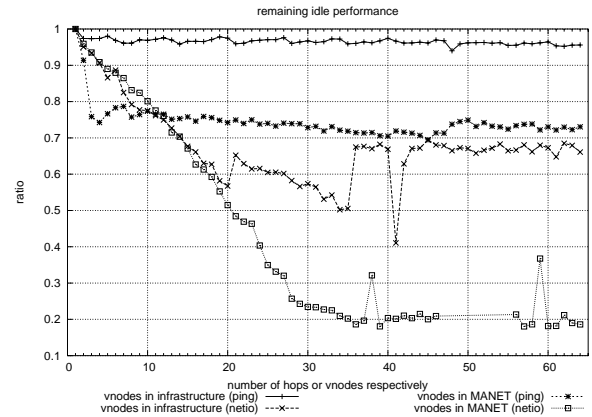


Figure 16. Remaining system compute performance versus number of vnodes.

Measurements for throughput already showed deviation from the reality, if too many nodes are hosted on the same pnode. In order to gain insight into system utilization, we measured the remaining idle performance on the pnode hosting vnodes while executing the two previously mentioned experiments. The results of all measurement runs are normalized to the result for one vnode per pnode resembling the unvirtualized case. Fig. 16 shows only one plot per throughput measurement, since the TCP load generator mea-

sures both directions (rx and tx in Fig. 15) back-to-back. In general, the remaining idle performance decreases with an increasing number of vnodes per pnode. This is an indicator for possible resource contention due to virtualization.

Hosting too many vnodes per pnode leads to severe resource contention which can lead to measurement artifacts. Since we are interested in realistic results, the number of hosted vnodes is limited. For the above measurements, the earliest undesirable deviation from the unvirtualized case happens for TCP throughput at a number of 20 vnodes (Fig. 15). We conclude that our approach supports 20 IP-router instances per pnode connected at a bandwidth of 100 MBit/s on each link. Given our testbed hardware with 64 pnodes and the scenario above, we thus can support scenario sizes of up to 1280 nodes.

D. Wireless Ad hoc Network Emulation

Wireless ad hoc networks typically consist of a large number of communicating devices, which are often resource-poor. By using node virtualization, wireless ad hoc scenarios can be emulated with a meaningful number of devices on an affordable smaller number of computers in an emulation testbed. Hence, we evaluate the scalability of our approach for the emulation of such scenarios. As before, we describe the system model, followed by evaluation results for the network and transport layer as well as for system utilization.

Fig. 17 shows the emulation scenario. For comparison with the infrastructure scenario, we configured the virtual node position and the emulated wireless network transmission range – depicted by dotted circles – such that the connectivity of the nodes resembles a chain. This is accomplished by a frame loss ratio for ingress traffic of zero for frames from reachable neighbors, and one for all others, as described in [3]. The wireless links between nodes are full duplex and have a limited bandwidth of 11 MBit/s. Here, we do not emulate the effects of a MAC layer, i.e. there are no frame collisions. Incorporating a MAC layer emulation as mentioned in Section II requires more resources and could reduce scalability. In this scenario, we use AODV-UU [19] version 0.7.2 as software under test. AODV-UU is an implementation of the ad hoc on-demand distance vector routing protocol. One instance of the routing daemon is executed on each node. In contrast to most routing daemon implementations, AODV-UU uses a small kernel module to overhear and queue relevant packets for accessing them in user space, and re-route packets after route establishment. In order to maintain the VRF ID on each routing decision, we had to add one line of code to the module specifying the ID as additional key component on route lookup. We consider this to be almost completely transparent for the software under test. Similar to the infrastructure scenario, we measure this scenario once with

only pnodes and once with all vnodes on a single pnode, except the last node, which resides on a separate pnode.

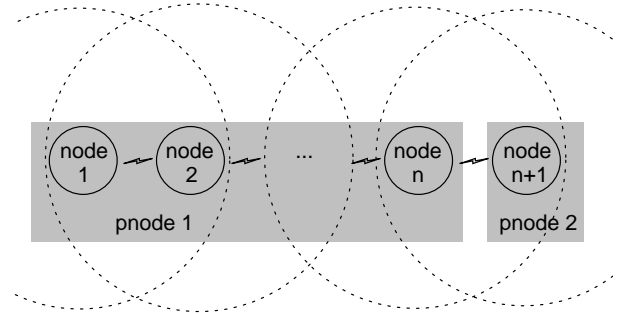


Figure 17. Mobile ad hoc network emulation scenario, virtualized case.

On network layer, we measure ICMP round trip time delays between the leftmost and the rightmost node in the scenario. Results are averaged over 10^4 echo/reply cycles for each hop count. The right y-axis in Fig. 13 corresponds to the measurement results for mean delay times. Starting with one hop we observe expanding ring search in combination with binary exponential backoff for outgoing route requests as described in [20], which is implemented by AODV-UU. Beyond the default time to live threshold, route requests work without expanding ring search leading to ICMP delays with linear increase starting at ten hops. Vnodes within the same pnode communicate efficiently observing shorter delays and thus leading to a more flat slope. The right y-axis in Fig. 13 corresponds to minima and maxima in ICMP delay times. Maxima resemble multiples of the mean values due to route cache misses on route establishment. Minima show very flat linear increase being observed for established routes with route cache hits.

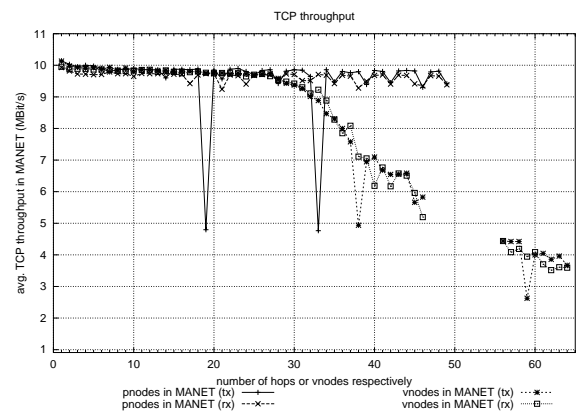


Figure 18. Throughput versus number of vnodes.

Measurement results for the transport layer are depicted in Fig. 18. TCP throughput starts deviating from the reference

values in the unvirtualized scenario variant at about 30 hops. As a result of the lower limited bandwidth compared to the infrastructure scenario in the previous subsection, more virtual nodes can be executed on a pnode without interference. Also TCP timing does not seem to be as critical so that both transmission directions behave similarly in the virtualized variant.

Remaining idle performance for the virtualized MANET scenario is shown in Fig. 16. In addition to the forwarding operation on network layer as per the previous infrastructure based scenario, one ad hoc routing daemon is executed on application layer on each vnode. This results in higher system utilization, the more vnodes are hosted on a pnode. The consequence is remaining idle performance of about 20 % at a maximum of 30 vnodes.

Similar to the infrastructure based scenario, the earliest undesirable deviation from the realistic reference values happens for TCP throughput (Fig. 18). We conclude that our approach supports 30 ad hoc routing instances per pnode connected at a bandwidth of 11 MBit/s on each link. Given our testbed hardware with 64 pnodes, we thus can support scenario sizes of up to 1920 nodes for similar wireless scenarios.

VI. RELATED WORK

In this section, we review existing scalable network emulation approaches. All approaches have in common that they place components of the software under test to certain, different positions within an emulated network scenario. We classify the architectures to build the emulated network in centralized and distributed.

A. Centralized Network Emulation

Centralized approaches emulate a whole scenario within a single instance of a network emulation tool. The traffic that can be handled by the central instance constitutes the upper limit for the scalability of these approaches.

Ns-e [21] is an emulation extension of the well-known network simulator ns-2 [22]. The scalability of ns-e depends on the amount of traffic in the scenario. For a typical MANET experiment, a scenario size of about 50 nodes is possible [23]. To some extent, this can be alleviated by extending the discrete event simulation into a parallel engine [24].

ModelNet [25] is a parallel network emulator. Its primary design is to emulate a given network topology of point-to-point links. The topology is partitioned among a cluster of emulation computers. Each cluster node processes network packets through internal arbitrarily connected links and routing instances. Computers running the software under test have to be externally connected to the central emulation cluster. Several instances of the software under test can run on each such edge node, which makes the approach scalable. However, the interface to the emulated network is based on

socket calls, which restricts the software under test to the application layer. Existing implementations of network protocols cannot be analyzed but have to be specifically re-implemented for the cluster nodes. This is also true with MobiNet [26], which is an extension to emulate MANET scenarios. The presented 802.11 MAC emulation in MobiNet is completely centralized and only works on a single core node. For a MANET scenario the authors report scalability up to 200 emulated nodes with re-implemented MANET routing protocol.

vBET [27] is an approach designed for emulating a network scenario on a single computer. It makes use of User Mode Linux (UML) [6] in order to provide virtual machines as execution environment for multiple virtual nodes on one computer. In combination with additional network emulation tools, it is possible to connect software under test to an emulated scenario. Connecting multiple of such vBET computers could allow larger scenarios. However, the use of UML's virtual machine concept introduces considerable overhead and thus limits the number of virtual nodes per computer. The authors report a maximum throughput for their software switch between vnodes of 128 MBit/s, which is an order of magnitude below our approach. vBET is more suitable for qualitative analysis than comparative performance analysis.

B. Distributed Network Emulation

Distributed approaches connect several instances of a network emulator together to form a comprehensive scenario.

Empower [28] allows the emulation of multiple routing instances on one computer, making up a link-based or wireless network topology. Each connection to the emulated network is mapped to a physical link of an existing hardware network interface. The authors equip each testbed node with several network cards to increase scalability. The number of network interfaces per pnode limits the number to a few vnodes per pnode.

Entrapid [29] and Alpine [30] virtualize the network stack in user space and thus provide multiple execution environments for software under test on a single computer. In combination with network emulation tools connecting such virtualized stacks, the emulation of network scenarios is possible. However, the software under test has to be adapted in order to interact with the user space network stacks. The packet processing in user space also introduces considerable timing inaccuracies, compared to real network stacks. Thus, these approaches are more suitable for testing than performance evaluation.

Vimage [12] virtualizes the network stack in the kernel. While common operating systems support one single instance of a network stack, vimage supports multiple independent instances. To accomplish this, the stack is modified to have all formerly global instance variables independently

available for each stack instance. Processes are associated with a certain network stack instance. Thus, the virtualization is transparent for software under test on the application layer. In combination with the network emulation tool dummynet [31], it is possible to emulate link-based scenarios in a scalable way. However, modifying all instance variables and access to them incorporates substantial changes to the network stack and thus the software under test. In [32] the authors report TCP throughput of 420 MBytes/s over 15 routing hops on a single machine with a slightly faster processor than used in our evaluation. Though scaling significantly better than a VMware based virtualization approach, the throughput was measured in a best case without any introduction of emulated network properties such as bandwidth limitation or delay and is thus hardly comparable to our results. Emulated network properties are however essential for network emulation and imply emulation overhead due to timer management reducing the accumulated throughput that can be realistically emulated.

The emulation testbed Netbed [33] supports scalable network emulation by introducing virtual nodes [34] on the basis of BSD jails [35] and multiple routing tables [13]. Netbed aims at emulating scenarios with fixed links. While it is possible to link *real* wireless nodes to an emulated scenario [36], there is no support for the *emulation* of wireless networks.

VII. SUMMARY AND CONCLUSION

Network emulation testbeds provide a synthetic, configurable network environment for comparative performance measurements of distributed applications and protocols. Common approaches limit the scenario size to the number of computers in the testbed, whereas meaningful emulation scenarios often require hundreds of communicating nodes. Testbeds of such sizes are hardly available.

In this paper, we propose to execute multiple instances of the software under test on a single testbed computer. Therefore, we introduce virtual nodes providing the software under test with a virtual execution environment with respect to the network stack. From a set of candidate node virtualization approaches, we choose the most lightweight approach fulfilling our paramount goal for scalability. In addition to our emulation software tools, we implemented an efficient software communication switch and extensions to “Virtual Routing and Forwarding” for Linux by James R. Leu. We provide an extensive evaluation of the implemented network emulation system. For a wired infrastructure-based and a wireless mobile ad hoc network emulation scenario, our measurement results show that node virtualization can increase the possible scenario size by the factor 20 or 30, respectively. Given our testbed hardware with 64 physical emulation nodes, this translates to scenario sizes of up to 1920 nodes.

Clearly, for scenario sizes of several hundred nodes, it is no more possible for an experimenter to manually map the nodes from an emulation scenario to the available testbed computers. Thus, our next step is an automated mapping based on constraints that evolve from the requirements of a scenario description and offerings of the testbed hardware. While we investigated the possible degree of virtualization by comparing measurements to the non-virtualized variant of a scenario, this procedure is not always desired or even possible. Therefore, we will introduce quality criteria for realistic network emulation which can be monitored while executing an experiment. In case of undesired resource contention due to virtualization, the experimenter will be informed and may decide to modify the mapping of vnodes to pnodes in order to prevent contention in another emulation run.

ACKNOWLEDGMENTS

This work is partially funded by the Deutsche Forschungsgemeinschaft (German Research Foundation) under grant DFG-GZ RO 1086/9-1.

AVAILABILITY

The source code of our implementations is publicly available under the terms of the GNU general public license at our NET-project web page: <http://net.informatik.uni-stuttgart.de>.

REFERENCES

- [1] Herrscher, D. and K. Rothermel. 2002. “A Dynamic Network Scenario Emulation Tool.” In *Proceedings of the 11th International Conference on Computer Communications and Networks (ICCCN '02)*, (Miami, October), 262–267.
- [2] Herrscher, D., S. Maier, and K. Rothermel. 2003. “Distributed Emulation of Shared Media Networks.” In *Proceedings of the 2003 International Symposium on Performance Evaluation of Computer and Telecommunication Systems (SPECTS '03)*, (Montréal, Quebec, Canada, July), 226–233.
- [3] Herrscher, D., S. Maier, J. Tian, and K. Rothermel. 2004. “A Novel Approach to Evaluating Implementations of Location-Based Software.” In *Proceedings of the 2004 International Symposium on Performance Evaluation of Computer and Telecommunication Systems (SPECTS '04)*, (San Jose, CA, USA, July), 484–490.
- [4] Lawton, K. 2000. “plex86: an i80x86 virtual machine.” In *Proceedings of the 4th Annual Linux Showcase & Conference*, (Atlanta, Georgia, USA, October).
- [5] Sugerman, J., G. Venkitachalam, and B.-H. Lim. 2001. “Virtualizing I/O Devices on VMware Workstation’s Hosted Virtual Machine Monitor.” In *Proceedings of the 2001 USENIX Annual Technical Conference*, (Boston, Massachusetts, USA, June), 1–14.
- [6] Dike, J. 2000. “A user-mode port of the Linux kernel.” In *Proceedings of the 4th Annual Linux Showcase & Conference*, (Atlanta, Georgia, USA, October).

- [7] Buchacker, K. and V. Sieh. 2001. "Framework for Testing the Fault-Tolerance of Systems Including OS and Network Aspects." In *Proceedings of the Third IEEE International High-Assurance System Engineering Symposium (HASE 2001)*, (Boca Raton, Florida), 95–105.
- [8] King, S.T., G.W. Dunlap, and P.M. Chen. 2003. "Operating System Support for Virtual Machines." In *Proceedings of the 2003 USENIX Annual Technical Conference*, (San Antonio, Texas, June), 71–84.
- [9] Waldspurger, C.A. 2002. "Memory Resource Management in VMware ESX Server." In *Proceedings of the 5th Symposium on Operating Systems Design and Implementation (OSDI'02)*, (Boston, MA, USA, December), 181–194.
- [10] Whitaker, A., M. Shaw, and S.D. Gribble. 2002. "Scale and Performance in the Denali Isolation Kernel." In *Proceedings of the 5th Symposium on Operating Systems Design and Implementation (OSDI'02)*, (Boston, MA, USA, December).
- [11] Barham, P., B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, and A. Warfield. 2003. "Xen and the Art of Virtualization." In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP 03)*, (Bolton Landing, New York, USA, October), 164–177.
- [12] Zec, M. 2003. "Implementing a Clonable Network Stack in the FreeBSD Kernel." In *Proceedings of the 2003 USENIX Annual Technical Conference*, (San Antonio, Texas, June), 137–150.
- [13] Scandariato, R. and F. Risso. 2002. "Advanced VPN support on FreeBSD systems." In *Proceedings of the 2nd European BSD Conference (BSDCon Europe 2002)*, (Amsterdam, the Netherlands, November).
- [14] Kourai, K., T. Hirotsu, K. Sato, O. Akashi, K. Fukuda, T. Sugawara, and S. Chiba. 2003. "Secure and Manageable Virtual Private Networks for Endusers." In *Proceedings of the 28th Annual IEEE International Conference on Local Computer Networks (LCN '03)*, (Bonn/Königswinter, Germany, October), 385–394.
- [15] Leu, J.R. 2004. Linux Virtual Routing and Forwarding, <http://linux-vrf.sourceforge.net/>.
- [16] Intel. 2004. IA-32 Intel Architecture Optimization Reference Manual, (USA).
- [17] McCalpin, J.D. 1995. "Memory Bandwidth and Machine Balance in Current High Performance Computers." *IEEE Technical Committee on Computer Architecture (TCCA) Newsletter*, (December), 19–25.
- [18] Peterson, L.L. and B.S. Davie. 1996. *Computer Networks: A Systems Approach*. Morgan Kaufmann Publishers.
- [19] Lundgren, H., E. Nordström, and C. Tschudin. 2002. "Coping with Communication Gray Zones in IEEE 802.11b based Ad hoc Networks." In *Proceedings of the 5th ACM International Workshop on Wireless Mobile Multimedia (WoWMoM'02)*, (Atlanta, Georgia, USA, September), 49–55.
- [20] Perkins, C.E., E.M. Belding-Royer, and S.R. Das. 2003. "Ad hoc On-Demand Distance Vector (AODV) Routing (work in progress)." Internet Draft, Internet Engineering Task Force, (February).
- [21] Fall, K. 1999. "Network Emulation in the Vint/NS Simulator." In *Proceedings of the Fourth IEEE Symposium on Computers and Communications (ISCC'99)*, (Red Sea, Egypt, July), 244–250.
- [22] Breslau, L., D. Estrin, K. Fall, S. Floyd, J. Heidemann, A. Helmy, P. Huang, S. McCanne, K. Varadhan, Y. Xu, and H. Yu. 2000. "Advances in Network Simulation." *IEEE Computer*, 33(5), (May): 59–67.
- [23] Ke, Q., D.A. Maltz, and D.B. Johnson. 2000. "Emulation of Multi-Hop Wireless Ad Hoc Networks." In *Proceedings of the 7th International Workshop on Mobile Multimedia Communications (MoMuC 2000)*, (Tokyo, Japan, October).
- [24] Riley, G.F., R.M. Fujimoto, and M.H. Ammar. 1999. "A Generic Framework for Parallelization of Network Simulations." In *Proceedings of the 7th International Symposium on Modelling, Analysis and Simulation of Computer and Telecommunications Systems (MASCOTS 99)*, (College Park, Maryland, March), 128–135.
- [25] Vahdat, A., K. Yocum, K. Walsh, P. Mahadevan, D. Kostic, J. Chase, and D. Becker. 2002. "Scalability and Accuracy in a Large-Scale Network Emulator." In *Proceedings of the 5th Symposium on Operating Systems Design and Implementation (OSDI'02)*, (Boston, MA, USA, December).
- [26] Mahadevan, P., A. Rodriguez, D. Becker, and A. Vahdat. 2004. "MobiNet: A Scalable Emulation Infrastructure for Ad Hoc and Wireless Networks." Technical Report CS2004-0792, Department of Computer Science, University of California, San Diego, (June).
- [27] Jiang, X. and D. Xu. 2003. "vBET: a VM-Based Emulation Testbed." In *Proceedings of the ACM SIGCOMM 2003 Workshops*, (Karlsruhe, Germany, August), 95–104.
- [28] Zheng, P. and L.M. Ni. 2002. "EMPOWER: A Scalable Framework for Network Emulation." In *Proceedings of the 2002 International Conference on Parallel Processing (ICPP'02)*, (Vancouver, B.C., Canada, August), 185–192.
- [29] Huang, X.W., R. Sharma, and S. Keshav. 1999. "The ENTRAPID Protocol Development Environment." In *Proceedings of the Conference on Computer Communications (INFOCOM 99)*, volume 3, (New York, NY, USA, March), 1107–1115.
- [30] Ely, D., S. Savage, and D. Wetherall. 2001. "Alpine: A User-Level Infrastructure for Network Protocol Development." In *Proceedings of the 3rd USENIX Symposium on Internet Technologies and Systems (USITS 2001)*, (San Francisco, California, USA, March).
- [31] Rizzo, L. 1997. "Dummynet: A simple approach to the evaluation of network protocols." *ACM Computer Communication Review*, 27(1), (January), 31–41.
- [32] Zec, M. and M. Mikuc. 2004. "Operating System Support for Integrated Network Emulation in IMUNES." In *Proceedings of the 1st Workshop on Operating System and Architectural Support for the on demand IT Infrastructure (2004 OASIS)*, (Boston, MA, October), 3–12.
- [33] White, B., J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar. 2002. "An Integrated Experimental Environment for Distributed Systems and Networks." In *Proceedings of the 5th Symposium on Operating Systems Design and Implementation (OSDI'02)*, (Boston, MA, USA, December), 255–270.
- [34] Guruprasad, S., L. Stoller, M. Hibler, and J. Lepreau. 2003. "Scaling Network Emulation with Multiplexed Virtual Resources." SIGCOMM 2003 Poster Abstract, (August).
- [35] Kamp, P.-H. and R.N.M. Watson. 2000. "Jails: Confining the omnipotent root." In *Proceedings of the 2nd International System Administration and Networking Conference (SANE 2000)*, (Maastricht, The Netherlands, May).
- [36] White, B., J. Lepreau, and S. Guruprasad. 2002. "Lowering the Barrier to Wireless and Mobile Experimentation." In *Proceedings of the First Workshop on Hot Topics in Networks (Hotnets-I)*, (Princeton, New Jersey, USA, October).