

# Bare-Metal Performance for Virtual Machines with Exitless Interrupts

By Nadav Amit, Abel Gordon, Nadav Har'El, Muli Ben-Yehuda, Alex Landau, Assaf Schuster, and Dan Tsafir

## Abstract

**Direct device assignment enhances the performance of guest virtual machines by allowing them to communicate with I/O devices without host involvement. But even with device assignment, guests are still unable to approach bare-metal performance, because the host intercepts all interrupts, including those generated by assigned devices to signal to guests the completion of their I/O requests. The host involvement induces multiple unwarranted guest/host context switches, which significantly hamper the performance of I/O intensive workloads. To solve this problem, we present ExitLess Interrupts (ELI), a software-only approach for handling interrupts within guest virtual machines *directly* and *securely*. By removing the host from the interrupt handling path, ELI improves the throughput and latency of unmodified, untrusted guests by 1.3×–1.6×, allowing them to reach 97–100% of bare-metal performance even for the most demanding I/O-intensive workloads.**

## 1. INTRODUCTION

I/O activity is a dominant factor in the performance of virtualized environments,<sup>17,25</sup> motivating *direct device assignment* where the host assigns physical I/O devices directly to guest virtual machines. Examples of such devices include disk controllers, network cards, and GPUs. Direct device assignment provides superior performance than alternative I/O virtualization approaches, because it almost entirely removes the host from the guest's I/O path. Without direct device assignment, I/O-intensive workloads might suffer unacceptable performance degradation.<sup>17,19,25</sup> Still, on x86 CPUs (the most popular platform for virtualization), direct assignment *alone* does not allow I/O-intensive workloads to approach bare-metal (nonvirtual) performance<sup>6,9,16,25</sup>; by our measurements, such workloads achieve only 60–65% of bare-metal performance. We find that nearly the *entire* performance difference is induced by interrupts of assigned devices.

I/O devices generate interrupts to notify the CPU of I/O operations' completion. In virtualized settings, each device interrupt triggers a costly *exit*,<sup>1,6</sup> causing the guest to be suspended and the host to be resumed, regardless of whether or not the device is assigned. The host first signals to the hardware the completion of the physical interrupt as mandated by the x86 specification. It then injects a corresponding (virtual) interrupt to the guest and resumes the guest's execution. The guest in turn handles the virtual interrupt and, like the host, signals completion, believing that it directly

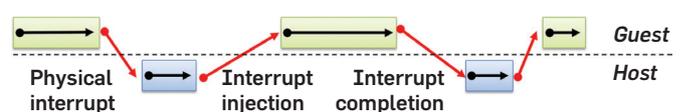
interacts with the hardware. This action triggers yet another exit, prompting the host to emulate the completion of the virtual interrupt and to resume the guest again. The chain of events for handling interrupts is illustrated in Figure 1.

The guest/host context switches caused by interrupts induce a tolerable overhead for non-I/O-intensive workloads, a fact that allowed some previous virtualization studies to claim they achieved bare-metal performance.<sup>5,14</sup> But our measurements indicate that this overhead quickly ceases to be tolerable, adversely affecting guests that require throughput of as little as 50 Mbps. Notably, previous studies improved virtual I/O by relaxing protection<sup>13,14</sup> or by modifying guests,<sup>5</sup> whereas we focus on the most challenging virtualization scenario of untrusted and unmodified guests.

Many previous studies identified interrupts as a major source of overhead,<sup>6,15</sup> and many proposed techniques to reduce it, both in bare-metal settings<sup>10,21,23,26</sup> and in virtualized settings<sup>3,9,16,25</sup>. In principle, it is possible to tune devices and their drivers to generate fewer interrupts, thereby reducing the related overhead. But doing so in practice is far from trivial<sup>22</sup> and can adversely affect both latency and throughput.

Our approach rests on the observation that the high interrupt rates experienced by a core running an I/O-intensive guest are mostly generated by devices assigned to the guest. Indeed, we measure rates of over 150K physical interrupts per second, even while employing standard techniques to reduce the number of interrupts, such as *interrupt coalescing*<sup>3,21,26</sup> and *hybrid polling*.<sup>10,23</sup> As noted, the resulting guest/host context switches are nearly exclusively responsible for the inferior performance relative to bare metal. To eliminate these switches, we propose ExitLess Interrupts (ELI), a software-only approach for handling physical interrupts directly within the guest in a secure manner.

Figure 1. Exits during interrupt handling.



A full version of this paper is available in *Proceedings of ACM Architectural Support for Programming Languages (ASPLOS) 2012*.

With ELI, physical interrupts are delivered directly to guests, allowing them to process their devices' interrupts without host involvement; ELI makes sure that each guest forwards all other interrupts to the host. With x86 hardware, interrupts are delivered using a software-controlled table of pointers to functions, such that the hardware invokes the  $k$ th function whenever an interrupt of type  $k$  fires. Instead of utilizing the guest's table, ELI maintains, manipulates, and protects a "shadow table," such that entries associated with assigned devices point to the guest's code, whereas the other entries are set to trigger an exit to the host.

We experimentally evaluate ELI with micro- and macro-benchmarks. Our baseline configuration employs standard techniques to reduce (coalesce) the number of interrupts, demonstrating ELI's benefit beyond the state-of-the-art. We show that ELI reduces CPU overheads that limit the attainable throughput, and thereby it improves the throughput and latency of guests by  $1.3\times$ – $1.6\times$ . Notably, whereas I/O-intensive guests were so far limited to 60–65% of bare-metal throughput, with ELI they reach performance that is within 97–100% of the optimum. Consequently, ELI makes it possible to, for example, consolidate traditional data-center workloads that nowadays remain nonvirtualized due to unacceptable performance loss.

## 2. MOTIVATION AND RELATED WORK

For the past several decades, interrupts have been the main method by which hardware devices can send asynchronous events to the operating system.<sup>7</sup> The main advantage of using interrupts to receive notifications from devices over polling them is that the processor is free to perform other tasks while waiting for an interrupt. This advantage applies when interrupts happen relatively infrequently, as was the case until high performance storage and network adapters came into existence. With these devices, the CPU can be overwhelmed with interrupts, leaving no time to execute code other than the interrupt handler.<sup>18</sup> When the operating system is run in a guest, interrupts have a higher cost, since every interrupt causes multiple exits.<sup>1,6</sup> ELI eliminates most of these exits and their associated overhead.

In the remainder of this section we introduce the existing approaches to reduce the overheads induced by interrupts, and we highlight the novelty of ELI in comparison to these approaches. We subdivide the approaches into two categories.

### 2.1. Generic interrupt handling approaches

We now survey approaches that apply equally to bare metal and virtualized environments.

**Polling** disables interrupts entirely and polls the device for new events at regular intervals. The benefit is that handling device events becomes synchronous, allowing the operating system to decide when to poll and thus limit the number of handler invocations. The drawbacks are added latency, increased power consumption (since the processor cannot enter an idle state), and wasted cycles when no events are pending. If polling is done on a different core, latency is improved, but a core is wasted.

A **hybrid** approach for reducing interrupt-handling overhead is to switch dynamically between using interrupts and

polling.<sup>10,18</sup> Linux uses this approach by default through the NAPI mechanism.<sup>23</sup> Switching between interrupts and polling does not always work well in practice, partly due to the complexity of predicting the number of interrupts a device will issue in the future.

Another approach is **interrupt coalescing**,<sup>3,21,26</sup> in which the OS programs the device to send one interrupt in a time interval or one interrupt per several events, as opposed to one interrupt per event. As with the hybrid approaches, coalescing delays interrupts and hence might increase latency<sup>15</sup> and burst TCP traffic.<sup>26</sup> Deciding on the right model and parameters for coalescing is particularly complex when the workload runs within a guest.<sup>9</sup> Getting it right for a wide variety of workloads is hard if not impossible.<sup>3,22</sup> Unlike coalescing, ELI does not reduce the number of interrupts; instead it streamlines the handling of interrupts targeted at virtual machines. Coalescing and ELI are therefore complementary, as we show in Section 5.4: coalescing reduces the number of interrupts, and ELI reduces their cost.

All evaluations in Section 5 were performed with the default Linux configuration, which combines the hybrid approach (via NAPI) and coalescing.

### 2.2. Virtualization-specific approaches

Using an emulated or paravirtual<sup>5</sup> device provides much flexibility on the host side, but its performance is much lower than that of device assignment, not to mention bare metal. Liu<sup>16</sup> shows that device assignment of SR-IOV devices can achieve throughput close to bare metal at the cost of as much as  $2\times$  higher CPU utilization. He also demonstrates that interrupts have a great impact on performance and are a major expense for both the transmit and receive paths.

There are software techniques<sup>2</sup> to reduce the number of exits by finding blocks of exiting instructions and exiting only once for the whole block. These techniques can increase the efficiency of running a virtual machine when the main reason for the overhead is in the guest code. When the reason is in external interrupts, such as for I/O intensive workloads with SR-IOV, such techniques do not alleviate the overhead.

Dong et al.<sup>9</sup> discuss a framework for implementing SR-IOV support in the Xen hypervisor. Their results show that SR-IOV can achieve line rate with a 10Gbps network interface controller (NIC). However, the CPU utilization is 148% of bare metal. In addition, this result is achieved using adaptive interrupt coalescing, which increases I/O latency.

Several studies attempted to reduce the aforementioned extra overhead of interrupts in virtual environments. vIC<sup>3</sup> discusses a method for interrupt coalescing in virtual storage devices and shows an improvement of up to 5% in a macro-benchmark. Their method uses the number of "commands in flight" to decide how many to coalesce. Therefore, as the authors say, this approach cannot be used for network devices due to the lack of information on commands (or packets) in flight. Dong et al.<sup>8</sup> use virtual interrupt coalescing via polling in the guest and receive side scaling to reduce network overhead in a paravirtual environment. Polling has its drawbacks, as discussed above, and ELI improves the more performance-oriented device assignment environment.

NoHype<sup>13</sup> argues that modern hypervisors are prone to attacks by their guests. In the NoHype model, the hypervisor is a thin layer that starts, stops, and performs other administrative actions on guests, but is not otherwise involved. Guests use assigned devices and interrupts are delivered directly to guests. No details of the implementation or performance results are provided. Instead, the authors focus on describing the security and other benefits of the model.

### 3. X86 INTERRUPT HANDLING

To put ELI's design in context, we begin with a short overview of how interrupt handling works on x86 today.

#### 3.1. Interrupts in bare-metal environments

x86 processors use interrupts and exceptions to notify system software about incoming events. Interrupts are asynchronous events generated by external entities such as I/O devices; exceptions are synchronous events—such as page faults—caused by the code being executed. In both cases, the currently executing code is interrupted and execution jumps to a pre-specified interrupt or exception handler.

x86 operating systems specify handlers for each interrupt and exception using an architected in-memory table, the Interrupt Descriptor Table (IDT). This table contains up to 256 entries, each entry containing a pointer to a handler. Each architecturally-defined exception or interrupt has a numeric identifier—an exception number or interrupt *vector*—which is used as an index to the table. The operating systems can use one IDT for all of the cores or a separate IDT per core. The operating system notifies the processor where each core's IDT is located in memory by writing the IDT's virtual memory address into the Interrupt Descriptor Table Register (IDTR). Since the IDTR holds the virtual (not physical) address of the IDT, the OS must always keep the corresponding address mapped in the active set of page tables. In addition to the table's location in memory, the IDTR holds the table's size.

When an external I/O device raises an interrupt, the processor reads the current value of the IDTR to find the IDT. Then, using the interrupt vector as an index to the IDT, the CPU obtains the virtual address of the corresponding handler and invokes it. Further interrupts may or may not be blocked while an interrupt handler runs.

System software needs to perform operations such as enabling and disabling interrupts, signaling the completion of interrupt handlers, configuring the timer interrupt, and sending interprocessor interrupts (IPIs). Software performs these operations through the Local Advanced Programmable Interrupt Controller (LAPIC) interface. The LAPIC has multiple registers used to configure, deliver, and signal completion of interrupts. Signaling the completion of interrupts, which is of particular importance to ELI, is done by writing to the end-of-interrupt (EOI) LAPIC register. The newest LAPIC interface, x2APIC,<sup>11</sup> exposes its registers using model specific registers (MSRs), which are accessed through “read MSR” and “write MSR” instructions. Previous LAPIC interfaces exposed the registers only in a predefined memory area which is accessed through regular load and store instructions.

#### 3.2. Interrupts in virtual environments

x86 hardware virtualization<sup>11</sup> provides two modes of operation, *guest mode* and *host mode*. The host, running in host mode, uses guest mode to create new contexts for running guest virtual machines. Once the processor starts running a guest, execution continues in guest mode until some sensitive event forces an exit back to host mode. The host handles any necessary events and then resumes the execution of the guest, causing an entry into guest mode. These exits and entries are the primary cause of virtualization overhead,<sup>1,6,19</sup> which is particularly pronounced in I/O intensive workloads.<sup>16,20,24</sup> It comes from the processor cycles spent switching between contexts, the time spent in host mode to handle the exit, and the resulting cache pollution.

This work focuses on running unmodified and untrusted operating systems. On the one hand, unmodified guests are not aware they run in a virtual machine, and they expect to control the IDT exactly as they do on bare metal. On the other hand, the host cannot easily give untrusted and unmodified guests control of each core's IDT. This is because having full control over the physical IDT implies total control of the core. Therefore, x86 hardware virtualization extensions use a different IDT for each mode. Guest mode execution on each core is controlled by the guest IDT and host mode execution is controlled by the host IDT. An I/O device can raise a physical interrupt when the CPU is executing either in host mode or in guest mode. If the interrupt arrives while the CPU is in guest mode, the CPU forces an exit and delivers the interrupt to the host through the host IDT.

Guests receive virtual interrupts, which are not necessarily related to physical interrupts. The host may decide to inject the guest with a virtual interrupt because the host received a corresponding physical interrupt, or the host may decide to inject the guest with a virtual interrupt manufactured by the host. The host injects virtual interrupts through the guest IDT. When the processor enters guest mode after an injection, the guest receives and handles the virtual interrupt.

During interrupt handling, the guest will access its LAPIC. Just like the IDT, full access to a core's physical LAPIC implies total control of the core, so the host cannot easily give untrusted guests access to the physical LAPIC. For guests using the first LAPIC generation, the processor forces an exit when the guest accesses the LAPIC memory area. For guests using x2APIC, the host traps LAPIC accesses according to an MSR bitmap, which specifies the sensitive MSRs that cannot be accessed directly by the guest. When the guest accesses sensitive MSRs, execution exits back to the host. In general, x2APIC registers are considered sensitive MSRs.

#### 3.3. Interrupts from assigned devices

The key to virtualization performance is for the CPU to spend most of its time in guest mode, running the guest, and not in the host, handling guest exits. I/O device emulation and paravirtualized drivers<sup>5</sup> incur significant overhead for I/O intensive workloads running in guests.<sup>6,16</sup> The overhead is incurred by the host's involvement in its guests' I/O paths for programmed I/O (PIO), memory-mapped I/O (MMIO), direct memory access (DMA), and interrupts.

Direct device assignment is the best performing approach for I/O virtualization<sup>9,16</sup> because it removes some of the host's involvement in the I/O path. With device assignment, guests are granted direct access to assigned devices. Guest I/O operations bypass the host and are communicated directly to devices. As noted, device DMA also bypasses the host; devices perform DMA accesses to and from guest memory directly. Interrupts generated by assigned devices, however, still require host intervention.

In theory, when the host assigns a device to a guest, it should also assign the physical interrupts generated by the device to that guest. Unfortunately, current x86 virtualization only supports two modes: either all physical interrupts on a core are delivered to the currently running guest, or all physical interrupts in guest mode cause an exit and are delivered to the host. An untrusted guest may handle its own interrupts, but it must not be allowed to handle the interrupts of the host and the other guests. Consequently, before ELI, the host had no choice but to configure the processor to force an exit when *any* physical interrupt arrives in guest mode. The host then inspected the interrupt and decided whether to handle it by itself or inject it to the associated guest.

Figure 1 describes the interrupt handling flow with baseline device assignment. Each physical interrupt from the guest's assigned device forces at least two exits from guest to host: when the interrupt arrives and when the guest signals completion of the interrupt handling. As we show in Section 5, interrupt-related exits are the foremost contributors to virtualization overhead for I/O intensive workloads.

#### 4. ELI: DESIGN AND IMPLEMENTATION

ELI enables unmodified and untrusted guests to handle interrupts directly and securely. ELI does not require any guest modifications, and thus should work with any operating system. It does not rely on any device-specific features, and thus should work with any assigned device.

##### 4.1. Exitless interrupt delivery

ELI's design was guided by the observation that *nearly* all physical interrupts arriving at a given core are targeted at the guest running on that core. This is due to several reasons. First, in high-performance deployments, guests usually have their own physical CPU cores (or else they would waste too much time context switching); second, high-performance deployments use device assignment with SR-IOV devices; and third, interrupt rates are usually proportional to execution time. The longer each guest runs, the more interrupts it receives from its assigned devices. Following this observation, ELI makes use of available hardware support to deliver *all* physical interrupts on a given core to the guest running on it, since most of them should be handled by that guest anyway, and forces the (unmodified) guest to reflect back to the host all those interrupts which should be handled by the host.

The guest OS continues to prepare and maintain its own IDT. Instead of running the guest with this IDT, ELI runs the guest in guest mode with a different IDT prepared by the host. We call this second guest IDT the *shadow* IDT. Just like shadow page tables can be used to virtualize the guest

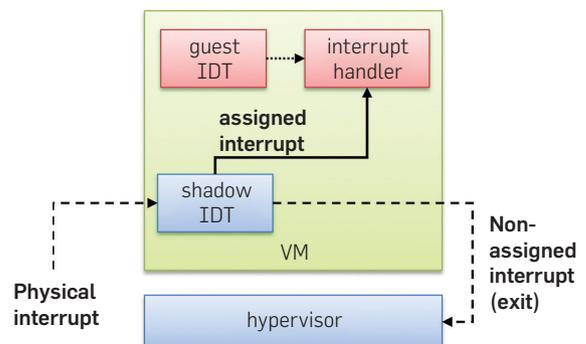
MMU,<sup>1,5</sup> IDT shadowing can be used to virtualize interrupt delivery. This mechanism, which is depicted in Figure 2 and described below, requires no guest cooperation.

By shadowing the guest's IDT, the host has explicit control over the interrupt handlers invoked by the CPU on interrupt delivery. The host can configure the shadow IDT to deliver assigned interrupts directly to the guest's interrupt handler or force an exit for nonassigned interrupts. The simplest method to cause an exit is to force the CPU to generate an exception, because exceptions can be selectively trapped by the host and can be easily generated if the host intentionally misconfigures the shadow IDT. For our implementation, we decided to force exits primarily by generating not-present (NP) exceptions. Each IDT entry has a present bit. Before invoking an entry to deliver an interrupt, the processor checks whether that entry is present (has the present bit set). Interrupts delivered to NP entries raise a NP exception. ELI configures the shadow IDT as follows: for exceptions and physical interrupts belonging to devices assigned to the guest, the shadow IDT entries are copied from the guest's original IDT and marked as present. Every other entry in the shadow IDT should be handled by the host and is therefore marked as not present to force a NP exception when the processor tries to invoke the handler. Additionally, the host configures the processor to force an exit from guest mode to host mode whenever a NP exception occurs.

Any physical interrupt reflected to the host appears in the host as a NP exception and must be converted back to the original interrupt vector. The host inspects the cause for this exception. If the exit was actually caused by a physical interrupt, the host raises a software interrupt with the same vector as the physical interrupt, which causes the processor to invoke the appropriate IDT entry. If the exit was not caused by a physical interrupt, then it is a true guest NP exception and should be handled by the guest. In this case, the host injects the exception back into the guest. True NP exceptions are rare in normal execution.

The host sometimes also needs to inject into the guest virtual interrupts raised by devices that are emulated by the host (e.g., the keyboard). These interrupt vectors will have their entries in the shadow IDT marked NP. To deliver such virtual interrupts through the guest IDT handler, ELI enters a special *injection mode* by configuring the

Figure 2. ELI interrupt delivery flow.



processor to cause an exit on any physical interrupt and running the guest with the original guest IDT. ELI then injects the virtual interrupt into the guest for handling, similarly to how it is usually done (Figure 1). After the guest signals completion of the injected virtual interrupt, ELI leaves injection mode by reconfiguring the processor to let the guest handle physical interrupts directly and resuming the guest with the shadow IDT. As we later show in Section 5, the number of injected virtual interrupts is orders of magnitude smaller than the number of physical interrupts generated by the assigned device. Thus, the number of exits due to physical interrupts while running in injection mode is negligible.

Even when all the interrupts require exits, ELI is not slower than baseline device assignment. The number of exits never increases and cost per exit remains the same. Common OS rarely modify the IDT content after system initialization. Entering and leaving injection mode requires only two memory writes, one to change the IDT pointer and the other to change the CPU execution mode.

#### 4.2. Placing the shadow IDT

There are several requirements on where in guest memory to place the shadow IDT. First, it should be hidden from the guest, that is, placed in memory not normally accessed by the guest. Second, it must be placed in a guest physical page that is always mapped in the guest's kernel address space. This is an x86 architectural requirement, since the IDTR expects a virtual address. Third, since the guest is unmodified and untrusted, the host cannot rely on any guest cooperation for placing the shadow IDT. ELI satisfies all three requirements by placing the shadow IDT in an extra page of a device's PCI Base Address Register (BAR).

PCI devices which expose their registers to system software as memory do so through BAR registers. BARs specify the location and sizes of device registers in physical memory. Linux and Windows drivers will map the full size of their devices' PCI BARs into the kernel's address space, but they will only access specific locations in the mapped BAR that are known to correspond to device registers. Placing the shadow IDT in an additional memory page tacked onto the end of a device's BAR causes the guest to (1) map it into its address space, (2) keep it mapped, and (3) not access it during normal operation. All of this happens as part of normal guest operation and does not require any guest awareness or cooperation. To detect runtime changes to the guest IDT, the host also write-protects the shadow IDT page.

#### 4.3. Configuring guest and host vectors

Neither the host nor the guest have absolute control over precisely when an assigned device interrupt fires. Since the host and the guest may run at different times on the core receiving the interrupt, both must be ready to handle the same interrupt. (The host handles the interrupt by injecting it into the guest.) Interrupt vectors also control that interrupt's priority relatively to other interrupts. Therefore, ELI makes sure that for each device interrupt,

the respective guest and host interrupt handlers are assigned to the same vector.

#### 4.4. Exitless interrupt completion

Although ELI IDT shadowing delivers hardware interrupts to the guest without host intervention, signaling interrupt completion still forces an exit to host mode. This exit is caused by the guest signaling the completion of an interrupt. As explained in Section 3.2, guests signal completion by writing to the EOI LAPIC register. This register is exposed to the guest either as part of the LAPIC area (older LAPIC interface) or as an x2APIC MSR (the new LAPIC interface). With the old interface, every LAPIC access causes an exit, whereas with the new one, the host can decide on a per-x2APIC-register basis which register accesses cause exits.

Before ELI, the host configured the CPU's MSR bitmap to force an exit when the guest accessed the EOI MSR. ELI exposes the x2APIC EOI register directly to the guest by configuring the MSR bitmap to not cause an exit when the guest writes to the EOI register. Combining this interrupt completion technique with ELI IDT shadowing eliminates the exits on the critical interrupt handling path.

Guests are not aware of the distinction between physical and virtual interrupts. They signal the completion of all interrupts the same way, by writing the EOI register. When the host injects a virtual interrupt, the corresponding completion should go to the host for emulation and not to the physical EOI register. Thus, during injection mode (described in Section 4.1), the host temporarily traps accesses to the EOI register. Once the guest signals the completion of all pending virtual interrupts, the host leaves injection mode.

#### 4.5. Protection

Full details of the considered threat model are available in the full paper. Here we briefly describe possible attacks and the mechanisms ELI employs to prevent them.

A malicious guest may try to steal CPU time by disabling interrupts forever. To prevent such an attack, ELI uses the *preemption timer* feature of x86, which triggers an unconditional exit after a configurable period of time elapses.

A misbehaving guest may refrain from signaling interrupt completion and thereby mask host interrupts. To prevent it, ELI signals interrupt completion for any assigned interrupt still in service after an exit. To maintain correctness, when ELI detects that the guest did not complete any previously delivered interrupts, it falls back to injection mode until the guest signals completions of all in-service interrupts. Since all of the registers that control CPU interruptibility are reloaded upon exit, the guest cannot affect host interruptibility.

A malicious guest can try to block or consume critical physical interrupts, such as a thermal interrupt. To protect against such an attack, ELI uses one of the following mechanisms. If there is a core which does not run any ELI-enabled guests, ELI redirects critical interrupts there. If no such core is available, ELI uses a combination of Non-maskable-Interrupts (NMIs) and IDT limiting.

NMIs trigger unconditional exits; they cannot be blocked by guests. ELI redirects critical interrupts to the core's single NMI handler. All critical interrupts are registered with this

handler, and whenever an NMI occurs, the handler calls all registered interrupt vectors to discern which critical interrupt occurred. NMI sharing has a negligible run-time cost (since critical interrupts rarely happen). However, some devices and device drivers may lock up or otherwise misbehave if their interrupt handlers are called when no interrupt was raised.

For critical interrupts whose handlers must only be called when an interrupt actually occurred, ELI uses a complementary coarse grained *IDT limit* mechanism. The IDT limit is specified in the IDTR register, which is protected by ELI and cannot be changed by the guest. IDT limiting reduces the limit of the shadow IDT, causing all interrupts whose vector is above the limit to trigger the usually rare general protection exception (GP). A GP is intercepted and handled by the host similarly to the NP exception. No events take precedence over the IDTR limit check,<sup>11</sup> and all handlers above the limit are therefore guaranteed to trap to the host when called.

## 5. EVALUATION

We implement ELI within the KVM hypervisor. This section evaluates the performance of our implementation.

### 5.1. Methodology and experimental setup

We measure and analyze ELI's effect on high-throughput network cards assigned to a guest virtual machine. Network devices are the most common use-case of device assignment, due to their high throughput and because SR-IOV network cards make it easy to assign one physical network card to multiple guests. We use throughput and latency to measure performance, and we contrast the results achieved by virtualized and bare-metal settings to demonstrate that the former can approach the latter. As noted earlier, performance-minded applications would typically dedicate whole cores to guests. We limit our evaluation to this case.

Our test machine is an IBM System x3550 M2 server, equipped with Intel Xeon X5570 CPUs, 24GB of memory, and an Emulex OneConnect 10Gbps NIC. We use another similar remote server (connected directly by 10Gbps fiber) as a workload generator and a target for I/O transactions. Guest mode and bare-metal configurations execute with a single core; 1GB of memory is assigned for each. All setups run Ubuntu 9.10 with Linux 2.6.35.

We run all guests on the KVM hypervisor (which is part of Linux 2.6.35) and QEMU-KVM 0.14.0. To check that ELI functions correctly in other setups, we also deploy it in an environment that uses a different device (BCM5709 1Gbps NIC) and a different OS (Windows 7); we find that ELI indeed operates correctly. We evaluate and compare the performance using baseline device assignment (i.e., unmodified KVM), ELI, and bare-metal system without virtualization.

We configure the hypervisor to back the guest's memory with 2MB *huge pages* and two-dimensional page tables. Huge pages minimize two-dimensional paging overhead and reduce TLB pressure. We note that only the host uses huge pages; in all cases the guest still operates with the default 4KB page size. We quantify the performance without huge pages, finding that they improve performance of both baseline and ELI runs similarly (data not shown).

Recall that ELI makes use of the x2APIC hardware to avoid exits on interrupt completions. x2APIC is available in every Intel x86 CPU since Sandy Bridge microarchitecture. Alas, the hardware we used for evaluation does not support x2APIC. To nevertheless measure the benefits of ELI utilizing x2APIC hardware, we slightly modify our Linux guest to emulate the x2APIC behavior. Specifically, we expose the physical LAPIC and a control flag to the guest, such that the guest may perform an EOI on the virtual LAPIC (forcing an exit) or the physical LAPIC (no exit), according to the flag. We verified that our approach conforms to the published specifications.

### 5.2. Throughput

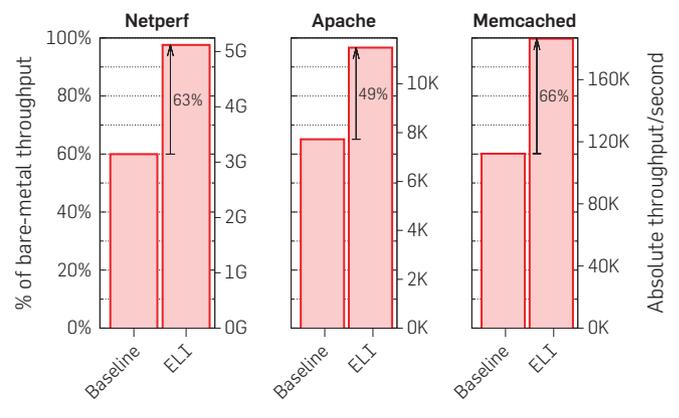
I/O virtualization performance suffers the most with workloads that are I/O intensive and which incur many interrupts. We start our evaluation by measuring three well-known examples of network-intensive workloads, and show that for these benchmarks ELI provides a significant (49–66%) throughput increase over baseline device assignment, and that it nearly (to 0–3%) reaches bare-metal performance.

We consider the following three benchmarks: **Netperf** TCP stream, which opens a single TCP connection to the remote machine, and makes as many rapid `write()` calls of a given size as possible; **Apache** HTTP server, measured using remote *ApacheBench* which repeatedly requests a static page from several concurrent threads; and **Memcached**, a high-performance in-memory key-value storage server, measured using the *Memslap* benchmark which sends a random sequence of `get` (90%) and `set` (10%) requests.

We configure each benchmark with parameters that fully load the tested machine's core (so that throughput can be compared), but do not saturate the tester machine. We configure Netperf to do 256-byte writes, ApacheBench to request 4KB static pages from 4 concurrent threads, and Memslap to make 64 concurrent requests from 4 threads.

Figure 3 illustrates how ELI improves the throughput of these three benchmarks. Each of the benchmarks was run on bare metal and under two virtualized setups: baseline device assignment, and device assignment with ELI.

**Figure 3. Performance of I/O intensive workloads relatively to bare-metal.**



The figure shows that baseline device assignment performance is still considerably below bare-metal performance: Netperf throughput on a guest is at 60% of bare-metal throughput, Apache is at 65%, and Memcached at 60%. With ELI, Netperf achieves 98% of the bare-metal throughput, Apache 97%, and Memcached 100%. It is evident that using ELI gives a significant throughput increase, 63%, 49%, and 66% for Netperf, Apache, and Memcached, respectively.

### 5.3. Execution breakdown

Breaking down the execution time to host, guest, and overhead components allows us to better understand how and why ELI improves the guest's performance. Table 1 shows this breakdown for the Apache benchmark (Netperf and Memcached appear in the full paper). We summarize here the results of the three benchmarks.

Guest performance should be better with ELI because the guest gets a larger fraction of the CPU (the host uses less), and/or because the guest runs more efficiently when it gets to run. With baseline device assignment, only 60–69% of the CPU time is spent in the guest; the rest is spent in the host, handling exits. ELI eliminates most of the exits, and thereby reduces both the fraction of time spent in the host (down to 1–2%) and the number of exits (down to 764–1118 per second).

In baseline device assignment, all interrupts arrive at the host and are then injected to the guest. The injection rate is slightly higher than the interrupt rate because the host injects additional virtual interrupts, such as timer interrupts. The number of interrupts “handled in host” is very low (103–207) when ELI is used, because the fraction of the time that the CPU is running the host is much lower.

Baseline device assignment is further slowed down by “IRQ window” exits: on bare metal, when a device interrupt occurs while interrupts are blocked, the interrupt will be delivered by the LAPIC hardware some time later. But when a guest is running, an interrupt always causes an immediate exit. The host wishes to inject this interrupt to the guest (if it is an interrupt from the assigned device), but if the guest has interrupts blocked, it cannot. The x86 architecture solution is to run the guest with an “IRQ window” enabled, requesting an exit as soon as the guest enables interrupts. We see 7801–9069 of these exits every second in the baseline device assignment run. ELI mostly eliminates IRQ window overhead, by eliminating most injections. Consequently, as expected, ELI slashes the number of exits, from 90,506 to 123,134 in the baseline device assignment runs, to just 764–1118.

**Table 1. Apache benchmark execution breakdown.**

| Statistics      | Baseline | ELI    | Bare-metal |
|-----------------|----------|--------|------------|
| Exits/s         | 90,506   | 1118   |            |
| Time in guest   | 67%      | 98%    |            |
| Interrupts/s    | 36,418   | 66,546 | 68,851     |
| Handled in host | 36,418   | 195    |            |
| Injections/s    | 36,671   | 458    |            |
| IRQ windows/s   | 7801     | 192    |            |
| Requests/s      | 7729     | 11,480 | 11,875     |
| Avg response ms | 0.518    | 0.348  | 0.337      |

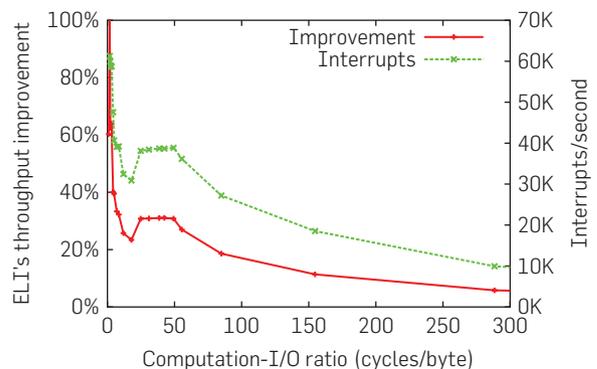
### 5.4. Impact of interrupt rate

The benchmarks in the previous section demonstrated that ELI significantly improves throughput over baseline device assignment for I/O intensive workloads. But as the workload spends less of its time on I/O and more of its time on computation, it seems likely that ELI's improvement will be less pronounced. Nonetheless, counterintuitively, we shall now show that ELI continues to provide relatively large improvements until we reach some fairly high computation-per-I/O ratio (and some fairly low throughput). To this end, we modify the Netperf benchmark to perform a specified amount of extra computation per byte written to the stream. This resembles many useful server workloads, where the server does some computation before sending its response.

A useful measure of the ratio of computation to I/O is *cycles/byte*, the number of CPU cycles spent to produce one byte of output; this ratio is easily measured as the quotient of CPU frequency (in cycles/second) and workload throughput (in bytes/second). Note that cycles/byte is inversely proportional to throughput. Figure 4 depicts ELI's improvement and the interrupt rate as a function of this ratio. As shown, until after 60 cycles/byte—which corresponds to throughput of only 50Mbps—ELI's improvement stays over 25% and the interrupt rate remains between 30K and 60K. As will be shown below, interrupt rates are kept in this range due to the NIC (which coalesces interrupts) and the Linux driver (which employs NAPI), and they would have been higher if it were not for these mechanisms. Since ELI lowers the overhead of handling interrupts, its benefit is proportional to their rate, *not* to throughput, a fact that explains why the improvement is similar over a range of computation-I/O values.

We now proceed to investigate the dependence of ELI's improvement on the amount of coalescing done by the NIC, which immediately translates to the amount of generated interrupts. Our NIC imposes a configurable cap on coalescing, allowing its users to set a time duration  $T$ , such that the NIC will not fire more than one interrupt per  $T \mu\text{s}$  (longer  $T$  implies less interrupts). We set the NIC's coalescing cap to the following values: 16  $\mu\text{s}$ , 24  $\mu\text{s}$ , 32  $\mu\text{s}$ , ..., 96  $\mu\text{s}$ . Figure 5 plots the results of the associated experiments (the data along the curve denotes values of  $T$ ). Higher interrupt rates

**Figure 4. Throughput improvement and baseline interrupt rate of modified-Netperf workloads with various computation-I/O ratios.**



imply higher savings due to ELI. Even with the maximal coalescing ELI still provides a 10% performance improvement over the baseline. ELI achieves at least 99% of bare-metal throughput in all of the experiments described in this section. These results indicate that when ELI is used, coalescing has lesser effect on throughput. The granularity of coalescing can therefore be made finer, so as to refrain from the increased latency that coarse coalescing induces.

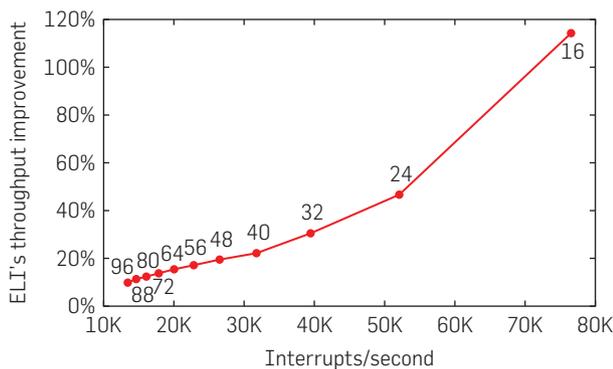
### 5.5. Latency

By removing the exits caused by external interrupts, ELI substantially reduces the time it takes to deliver interrupts to the guest. This period of time is critical for latency-sensitive workloads. We measure ELI's latency improvement using Netperf UDP request-response, which sends a UDP packet and waits for a reply before sending the next. To simulate a busy guest that has work to do alongside a latency-sensitive application, we run a busy-loop within the guest. As the results in Table 2 show, baseline device assignment increases bare metal latency by 8.21  $\mu$ s and that ELI reduces this gap to only 0.58  $\mu$ s, which is within 98% of bare-metal latency.

## 6. AFTERMATH

In our original ASPLOS 2012 paper, we urged hardware vendors to add hardware support that would simplify implementing direct interrupt delivery to guest virtual machines. We made the case that the substantial performance improvement demonstrated by ELI merits the effort to add such support. We are happy to report that, since then, a few positive steps has been taken in this direction.

**Figure 5. Throughput improvement and interrupt rate for Netperf benchmark with different interrupt coalescing intervals (shown in labels).**



**Table 2. Latency measured by Netperf UDP request-response benchmark.**

| Configuration | Latency ( $\mu$ s) | % Overhead |
|---------------|--------------------|------------|
| Baseline      | 36.14              | 29         |
| ELI           | 28.51              | 2          |
| Bare-metal    | 27.93              | 0          |

To mitigate some of the overheads caused by interrupts delivery, hypervisors can now use the Intel “virtual APIC” (APICv) feature. Assume that a guest currently runs on core  $C_1$ . The hypervisor can arrange things such that the relevant (physical) interrupts are triggered on a different core  $C_2$  that runs in host mode. When an interrupt reaches  $C_2$ , APICv then allows the hypervisor to “forward” the corresponding (virtual) interrupts to  $C_1$  to the guest *without* inducing an exit. Although such a scheme eliminates unwarranted exits for the guest, it is inferior to ELI due to two reasons. First, it requires dedicating special host cores (like  $C_2$ ) for redirecting guest interrupts. Second, it increases interrupts delivery latency, as they must first be processed by the hypervisor at  $C_2$  and only then can they be delivered to the guest at  $C_1$ .

Both Intel and AMD indicate that they intend to support direct ELI-like delivery in hardware.<sup>12</sup> Some ARM chips already support such delivery.<sup>4</sup> It is still unclear, however, whether this hardware support would live up to its promise. The first generation of Intel implementation, for instance, would deliver each guest interrupt to a certain core. As a result, this implementation may not be usable for multi-core guests whose OS spreads interrupts across the guest cores.

## 7. CONCLUSION

The key to high virtualization performance is for the CPU to spend most of its time in guest mode, running the guest, and not in the host, handling guest exits. Yet current approaches to x86 virtualization induce multiple exits by requiring host involvement in the critical interrupt handling path. The result is that I/O performance suffers. We propose to eliminate the unwarranted exits by introducing ELI, an approach that lets guests handle interrupts directly and securely. Building on many previous efforts to reduce virtualization overhead, ELI finally makes it possible for untrusted and unmodified virtual machines to reach nearly bare-metal performance, even for the most I/O-intensive workloads. Considering, it seems that the next logical step for chip vendors is extend the posted interrupts architecture so as to support the ELI paradigm in hardware, thereby simplifying its implementation.

### Acknowledgments

The research leading to the results presented in this paper is partially supported by the European Community's Seventh Framework Programme ([FP7/2001–2013]) under grant agreements #248615 (IOLanes) and #248647 (ENCORE). ■

### References

- Adams, K., Agesen, O. A comparison of software and hardware techniques for x86 virtualization. In *ACM Architectural Support for Programming Languages & Operating Systems (ASPLOS)* (2006).
- Agesen, O., Mattson, J., Rugina, R., Sheldon, J. Software techniques for avoiding hardware virtualization exits. In *USENIX Annual Technical Conference (ATC)* (2012), 373–385.
- Ahmad, I., Gulati, A., Mashtizadeh, A. vIC: Interrupt coalescing for virtual machine storage device IO. In *USENIX Annual Technical Conference (ATC)* (2011).
- ARM Ltd. *Arm Generic Interrupt Controller Architecture Version 2.0*. ARM IHI 0048B, 2011.
- Barham, P., Dragovic, B., Fraser, K., Hand, S., Harris, T., Ho, A., Neugebauer, R., Pratt, I., Warfield, A. Xen and the art of virtualization. In *ACM Symposium on Operating Systems Principles (SOSP)* (2003).
- Ben-Yehuda, M., Day, M.D., Dubitzky, Z., Factor, M., Har'El, N., Gordon, A., Liguori, A., Wasserman, O., Yassour, B.-A. The turtles project: Design and implementation of nested virtualization. In *USENIX Symposium on*

*Operating Systems Design & Implementation (OSDI)* (2010).

7. Codd, E.F. *Advances in Computers*, Volume 3, New York: Academic Press, 1962, 77–153.
8. Dong, Y., Xu, D., Zhang, Y., Liao, G. Optimizing network I/O virtualization with efficient interrupt coalescing and virtual receive side scaling. In *IEEE International Conference on Cluster Computing (CLUSTER)* (2011).
9. Dong, Y., Yang, X., Li, X., Li, J., Tian, K., Guan, H. High performance network virtualization with SR-IOV. In *IEEE International Symposium on High Performance Computer Architecture (HPCA)* (2010).
10. Dovrolis, C., Thayer, B., Ramanathan, P. HIP: Hybrid interrupt-polling for the network interface. *ACM SIGOPS Operat. Syst. Rev.* 35 (2001), 50–60.
11. Intel Corporation. *Intel 64 and IA-32 Architectures Software Developer's Manual*, 2014.
12. Intel Corporation. Intel virtualization technology for directed I/O architecture specification, 2014.
13. Keller, E., Szefer, J., Rexford, J., Lee, R.B. NoHype: Virtualized cloud infrastructure without the virtualization. In *ACM/IEEE International Symposium on Computer Architecture (ISCA)* (2010), ACM.
14. Lange, J.R., Pedretti, K., Dinda, P., Bridges, P.G., Bae, C., Soltero, P., Merritt, A. Minimal-overhead virtualization of a large scale supercomputer. In *ACM/USENIX International Conference on Virtual Execution Environments (VEE)* (2011).
15. Larsen, S., Sarangam, P., Huggahalli, R., Kulkarni, S.

Architectural breakdown of end-to-end latency in a TCP/IP network. *Int. J. Parallel Prog.* 37, 6 (2009), 556–571.

16. Liu, J. Evaluating standard-based self-virtualizing devices: A performance study on 10 GbE NICs with SR-IOV support. In *IEEE International Parallel & Distributed Processing Symposium (IPDPS)* (2010).
17. Liu, J., Huang, W., Abali, B., Panda, D.K. High performance VMM-bypass I/O in virtual machines. In *USENIX Annual Technical Conference (ATC)* (2006).
18. Mogul, J.C., Ramakrishnan, K.K. Eliminating receive livelock in an interrupt-driven kernel. *ACM Trans. Comput. Syst.* 15 (1997), 217–252.
19. Raj, H., Schwan, K. High performance and scalable I/O virtualization via self-virtualized devices. In *International Symposium on High Performance Distributed Computer (HPDC)* (2007).
20. Ram, K.K., Santos, J.R., Turner, Y., Cox, A.L., Rixner, S. Achieving 10Gbps using safe and transparent network interface virtualization. In *ACM/USENIX International Conference on Virtual Execution Environments (VEE)* (2009).
21. Salah, K. To coalesce or not to coalesce. *Int. J. Electron. Commun.* 61, 4 (2007), 215–225.
22. Salah, K., Qahtan, A. Boosting throughput of Snort NIDS under Linux. In *International Conference on Innovations in Information Technology (IIT)* (2008).
23. Salim, J.H., Olsson, R., Kuznetsov, A. Beyond softnet. In *Annual Linux Showcase & Conference* (2001).

24. Santos, J.R., Turner, Y., Janakiraman, G.J., Pratt, I. Bridging the gap between software and hardware techniques for I/O virtualization. In *USENIX Annual Technical Conference (ATC)* (2008).
25. Willmann, P., Shafer, J., Carr, D., Menon, A., Rixner, S., Cox, A.L., Zwaenepoel, W. Concurrent direct network access for virtual machine monitors. In *IEEE International*

*Symposium on High Performance Computer Architecture (HPCA)* (2007).

26. Zec, M., Mikuc, M., Žagar, M. Estimating the impact of interrupt coalescing delays on steady state TCP throughput. In *International Conference on Software, Telecommunications and Computer Networks (SoftCOM)* (2002).

**Nadav Amit, Assaf Schuster, and Dan Tsafir** ((namit, assaf, dan)@cs.technion.ac.il), Technion, Haifa, Israel.

**Nadav Har'El** (nadav@harel.org.il), Cloudius Systems, Herzliya Pituach, Israel.

**Abel Gordon and Muli Ben-Yehuda** (abel@stratoscale.com, mulix@mulix.org), Stratoscale, Haifa, Israel.

**Alex Landau** (landau.alex@gmail.com), Facebook, Seattle, WA.

Authors Gordon, Har'El, Ben-Yehuda, and Landau conducted the research discussed in the paper while employed at IBM.

Copyright held by authors.  
Publication rights licensed to ACM. \$15.00.

## ACM Transactions on Parallel Computing

*Solutions to Complex Issues in Parallelism*

Editor-in-Chief: Phillip B. Gibbons, Intel Labs, Pittsburgh, USA



*ACM Transactions on Parallel Computing* (TOPC) is a forum for novel and innovative work on all aspects of parallel computing, including foundational and theoretical aspects, systems, languages, architectures, tools, and applications. It will address all classes of parallel-processing platforms including concurrent, multithreaded, multicore, accelerated, multiprocessor, clusters, and supercomputers.

### Subject Areas

- Parallel Programming Languages and Models
- Parallel System Software
- Parallel Architectures
- Parallel Algorithms and Theory
- Parallel Applications
- Tools for Parallel Computing



Association for Computing Machinery

Advancing Computing as a Science & Profession

For further information or to submit your manuscript, visit [topc.acm.org](http://topc.acm.org)

Subscribe at [www.acm.org/subscribe](http://www.acm.org/subscribe)