

Singapore Management University  
**Institutional Knowledge at Singapore Management University**

---

Research Collection School Of Information Systems

School of Information Systems

---

6-2013

# Mitigating Access-Driven Timing Channels in Clouds using StopWatch

Peng LI

Debin GAO

Singapore Management University, [dbgao@smu.edu.sg](mailto:dbgao@smu.edu.sg)

Michael K. Reiter

**DOI:** <https://doi.org/10.1109/DSN.2013.6575299>

Follow this and additional works at: [https://ink.library.smu.edu.sg/sis\\_research](https://ink.library.smu.edu.sg/sis_research)

 Part of the [Information Security Commons](#)

---

## Citation

LI, Peng; GAO, Debin; and Reiter, Michael K.. Mitigating Access-Driven Timing Channels in Clouds using StopWatch. (2013). *43rd Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN 2013)*. 1-12. Research Collection School Of Information Systems.

**Available at:** [https://ink.library.smu.edu.sg/sis\\_research/2038](https://ink.library.smu.edu.sg/sis_research/2038)

This Conference Proceeding Article is brought to you for free and open access by the School of Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email [libIR@smu.edu.sg](mailto:libIR@smu.edu.sg).

# Mitigating Access-Driven Timing Channels in Clouds using StopWatch

Peng Li

Department of Computer Science  
University of North Carolina  
Email: pengli@cs.unc.edu

Debin Gao

School of Information Systems  
Singapore Management University  
Email: dbgao@smu.edu.sg

Michael K. Reiter

Department of Computer Science  
University of North Carolina  
Email: reiter@cs.unc.edu

**Abstract**—This paper presents *StopWatch*, a system that defends against timing-based side-channel attacks that arise from coresidency of victims and attackers in infrastructure-as-a-service clouds. *StopWatch* triplicates each cloud-resident guest virtual machine (VM) and places replicas so that the three replicas of a guest VM are coresident with nonoverlapping sets of (replicas of) other VMs. *StopWatch* uses the timing of I/O events at a VM’s replicas collectively to determine the timings observed by each one or by an external observer, so that observable timing behaviors are similarly likely in the absence of any other individual, coresident VM. We detail the design and implementation of *StopWatch* in Xen, evaluate the factors that influence its performance, and address the problem of placing VM replicas in a cloud under the constraints of *StopWatch* so as to still enable adequate cloud utilization.

## I. INTRODUCTION

Implicit timing-based information flows threaten the use of clouds for very sensitive computations. In an “infrastructure as a service” (IaaS) cloud, such an attack could be mounted by an attacker submitting a virtual machine (VM) to the cloud that times the duration between events that it can observe, to make inferences about a *victim* VM with which it is running simultaneously on the same host but otherwise cannot access. Such “access-driven” attacks [1] were first studied in the context of timing-based *covert channels*, in which the victim VM is infected with a Trojan horse that intentionally signals information to the attacker VM by manipulating the timings that the attacker VM observes. Of more significance in modern cloud environments, however, are timing-based *side channels*, which leverage the same principles to attack an uninfected but oblivious victim VM (e.g., [2], [1]).

In this paper we propose an approach to defend against timing attacks and a system that implements this method for IaaS clouds. Our system, called *StopWatch*, alters timings observed by the attacker VM to “match” those of a *replica* attacker VM that is *not* coresident with the victim. Since *StopWatch* cannot identify attackers and victims *a priori*, realizing this intuition in practice requires replicating each VM on multiple hosts and enforcing that the replicas are coresident with nonoverlapping sets of (replicas of) other VMs. Moreover, two replicas is not enough: one might be coresident with its victim, and by symmetry, its timings would necessarily influence the timings imposed on the pair. *StopWatch* thus uses three replicas that coreside with nonoverlapping sets of (replicas of) other VMs and imposes the median timing of the three on all replicas. Even if the median timing of an event is that which occurred at an attacker replica that is coresident with a victim replica, timings both below and above the median occurred at attacker replicas that do not coreside with the

victim. The median can thus be viewed as “microaggregating” the timings to confound inferences from them (c.f., [3], [4]).

We detail the implementation of *StopWatch* in Xen, specifically to intervene on all real-time clocks and, notably, to enforce this median behavior on “clocks” available via the I/O subsystem (e.g., network interrupts). Moreover, for uniprocessor VMs, *StopWatch* enforces deterministic execution across all of a VM’s replicas, making it impossible for an attacker VM to utilize other internally observable clocks and ensuring the same outputs from the VM replicas. By applying the median principle to the timing of these outputs, *StopWatch* further interferes with inferences that an observer external to the cloud could make on the basis of output timings.

We evaluate the performance of our *StopWatch* prototype for supporting web service (file downloads) and various types of computations. Our analysis shows that the latency overhead of *StopWatch* is less than  $2.8\times$  even for network-intensive applications. We also identify adaptations to a service that can vastly increase its performance when run over *StopWatch*, e.g., making file download over *StopWatch* competitive with file download over unmodified Xen. For computational benchmarks, the latency induced by *StopWatch* is less than  $2.3\times$  and is directly correlated with their amounts of disk I/O.

We also study the impact of *StopWatch* on cloud utilization, i.e., how many guest VMs can be simultaneously executed on an infrastructure of  $n$  machines, each with a capacity of  $c$  guest VMs, under the constraint that the three replicas for each guest VM coreside with nonoverlapping sets of (replicas of) other VMs. We show that for any  $c \leq \frac{n-1}{2}$ ,  $\Theta(cn)$  guest VMs (three replicas of each) can be simultaneously executed; we also identify practical algorithms for placing replicas to achieve this bound. This distinguishes *StopWatch* from the alternative of simply running each guest VM on a separate computer, which permits simultaneous execution of only  $n$  guest VMs.

To summarize, our contributions are as follows: First, we introduce a novel approach for defending against access-driven timing side-channel attacks in “infrastructure-as-a-service” (IaaS) compute clouds that leverages replication of guest VMs with the constraint that the replicas of each guest VM coreside with nonoverlapping sets of (replicas of) other VMs. The median timings of I/O events across the three guest VM replicas are then imposed on these replicas to interfere with their use of event timings to extract information from a victim VM with which one is coresident. Second, we detail the implementation of this strategy in Xen, yielding a system called *StopWatch*, and evaluate the performance of *StopWatch* on a variety of workloads. This evaluation sheds light on the features of workloads that most impact the performance

of applications running on *StopWatch* and how they can be adapted for best performance. Third, we show how to place replicas under the constraints of *StopWatch* to utilize a cloud infrastructure more effectively than running each guest VM in isolation. Finally, in the Appendix we analyze the median as a microaggregation function and explain its benefits over the alternative of obscuring event timings with random noise (e.g., [5]).

## II. RELATED WORK

**Timing channel defenses.** Defenses against information leakage via timing channels are diverse, taking numerous different angles on the problem. Research on type systems and security-typed languages to eliminate timing attacks offers powerful solutions (e.g., [6], [7], [8]), but this work is not immediately applicable to our goal here, namely adapting an existing virtual machine monitor (VMM) to support practical mitigation of timing channels today. Other research has focused on the elimination of timing side channels within cryptographic computations (e.g., [9]) or as enabled by specific hardware components (e.g., [10], [11]), but we seek an approach that is comprehensive.

Askarov et al. [12] distinguish between *internal* timing channels that involve the implicit or explicit measurement of time from within the system, and *external* timing channels that involve measuring the system from the point of view of an external observer. Defenses for both internal (e.g., [5], [6], [7], [13]) and external (e.g., [14], [15], [12], [16], [17]) timing channels have received significant attention individually, though to our knowledge, *StopWatch* is novel in addressing timing channels through a combination of both techniques. *StopWatch* incorporates internal defenses to interfere with an attacker’s use of real-time clocks or “clocks” that it might derive from the I/O subsystem. In doing so, *StopWatch* imposes determinism on uniprocessor VMs and then uses this feature to additionally build an effective external defense against such attacker VMs.

*StopWatch*’s internal and external defense strategies also differ individually from prior work, in interfering with timing channels by allowing replicas (in the internal defenses) and external observers (in the external defenses) to observe only median I/O timings across the three replicas. The median offers several benefits over the alternative of obfuscating event timings by adding random noise (without replicating VMs): to implement random noise, a distribution from which to draw the noise must be chosen without reference to an execution in the absence of the victim—i.e., how the execution “should have” looked—and so ensuring that the chosen noise distribution is sufficient to suppress all timing channels can be quite difficult. *StopWatch* uses replication and careful replica placement (in terms of the other VMs with which each replica coresides) exactly to provide such a reference. Moreover, in the Appendix we show that the median permits the delays incurred by the system to scale better than uniformly random noise allows for the same protection, as the distinctiveness of victim behavior increases.

**Replication.** To our knowledge, *StopWatch* is novel in utilizing replication for timing channel defense. That said, replication has a long history that includes techniques similar

to those we use here. For example, state-machine replication to mask Byzantine faults [18] ensures that correct replicas return the same response to each request so that this response can be identified by “vote” (a technique related to one employed in *StopWatch*; see Sec. III and Sec. VI). To ensure that correct replicas return the same responses, these systems enforce the delivery of requests to replicas in the same order; moreover, they typically assume that replicas are deterministic and process requests in the order they are received. *Enforcing* replica determinism has also been a focus of research in (both Byzantine and benignly) fault-tolerant systems; most (e.g., [19], [20], [21]), but not all (e.g., [22]), do so at other layers of the software stack than *StopWatch* does.

More fundamentally, to our knowledge all prior systems that enforce timing determinism across replicas permit one replica to dictate timing-related events for the others, which does not suffice for our goals: that replica could be the one coresident with the victim, and so permitting it to dictate timing related events would simply “copy” the information it gleans from the victim to the other replicas, enabling that information to then be leaked out of the cloud. Rather, by forcing the timing of events to conform to the median timing across three VM replicas, at most one of which is coresident with the victim, the enforced timing of each event is either the timing of a replica not coresident with the victim or else between the timing of two replicas that are not coresident with the victim. This strategy is akin to ones used for Byzantine fault-tolerant clock synchronization (e.g., see [23, Sec. 5.2]) or sensor replication (e.g., see [18, Sec. 5.1]), though we use it here for information hiding (versus integrity).

Aside from replication for fault tolerance, replication has been explored to detect server penetration [24], [25], [26], [27]. These approaches purposely employ diverse replica codebases or data representations so as to reduce the likelihood of a single exploit succeeding on multiple replicas. Divergence of replica behavior in these approaches is then indicative of an exploit succeeding on one but not others. In contrast to these approaches, *StopWatch* leverages (necessarily) *identical* guest VM replicas to address a different class of attacks (timing side channels) than replica compromise.

Research on VM execution *replay* (e.g., [28], [29]) focuses on recording nondeterministic events that alter VM execution and then coercing these events to occur the same way when the VM is replayed. The replayed VM is a replica of the original, albeit a temporally delayed one, and so this can also be viewed as a form of replication. *StopWatch* similarly coerces VM replicas to observe the same event timings, but again, unlike these timings being determined by one replica (the original), they are determined collectively using median calculations, so as to interfere with one attacker VM replica that is coresident with the victim from simply propagating its timings to all replicas. That said, the state-of-the-art in VM replay (e.g., [29]) addresses multiprocessor VM execution, which our present implementation of *StopWatch* does not. *StopWatch* could be extended to support multiprocessor execution with techniques for deterministic multiprocessor scheduling (e.g., [30]). Mechanisms for enforcing deterministic execution through O/S-level modifications (e.g., [31]) are less relevant to our goals, as they are not easily used by an IaaS cloud provider that accepts arbitrary VMs to execute.

### III. DESIGN

Our design is focused on “infrastructure as a service” (IaaS) clouds that accept virtual machine images, or “guest VMs,” from customers to execute. Amazon EC2 (<http://aws.amazon.com/ec2/>) and Rackspace (<http://www.rackspace.com/>) are example providers of public IaaS clouds. Given the concerns associated with side-channel attacks in cloud environments (e.g., [2], [1]), we seek to develop virtualization software that would enable a provider to construct a cloud that offers substantially stronger assurances against leakage via timing channels. This cloud might be a higher assurance offering that a provider runs alongside its normal cloud (while presumably charging more for the greater assurance it offers) or a private cloud with substantial assurance needs (e.g., run by and for an intelligence or military community).

**Threat model.** Our threat model is a customer who submits *attacker VMs* for execution that are designed to employ timing side channels. We presume that the attacker VM is designed to extract information from a particular victim VM, versus trying to learn general statistics about the cloud such as its average utilization. We assume that access controls prevent the attacker VMs from accessing victim VMs directly or from escalating their own privileges in a way that would permit them to access victim VMs. The cloud’s virtualization software (in our case, Xen and our extensions thereof) is trusted.

According to Wray [32], to exploit a timing channel, the attacker VM measures the timing of observable events using a *clock* that is independent of the timings being measured. While the most common such clock is real time, a clock can be any sequence of observable events. With this general definition of a “clock,” a timing attack simply involves measuring one clock using another. Wray identified four possible clock sources in conventional computers [32]:

- TL: the “CPU instruction-cycle clock” (e.g., a clock constructed by executing a simple timing loop);
- Mem: the memory subsystem (e.g., data/instruction fetches);
- IO: the I/O subsystem (e.g., network, disk, and DMA interrupts); and
- RT: real-time clocks provided by the hardware platform (e.g., time-of-day registers).

**Defense strategy.** *StopWatch* is designed to interfere with the use of IO and RT clocks and, for uniprocessor VMs, TL or Mem clocks, for timing attacks. (As discussed in Sec. II, extension to multiprocessor VMs is a topic of future work.) IO and RT (especially RT) clocks are an ingredient in every timing side-channel attack in the research literature that we have found, undoubtedly because real time is the most intuitive, independent and reliable reference clock for measuring another clock. So, intervening on these clocks is of paramount importance. Moreover, the way *StopWatch* does so forces the scheduler in a uniprocessor guest VM to behave deterministically, interfering with attempts to use TL or Mem clocks.

More specifically, to interfere with IO clocks, *StopWatch* replicates each attacker VM (i.e., every VM, since we do not presume to know which ones are attacker VMs) threefold so that the three replicas of a guest VM are coresident with

nonoverlapping sets of (replicas of) other VMs. Then, when determining the timing with which an event is made available to each replica, the median timing value of the three is adopted. *StopWatch* addresses RT clocks by replacing a VM’s view of real time with a virtual time that depends on the VM’s own progress, an idea due to Popek and Kline [33].

A side effect of how *StopWatch* addresses IO and RT clocks is that it enforces deterministic execution of uniprocessor attacker VM replicas, also disabling its ability to use TL or Mem clocks. These mechanisms thus deal effectively with internal observations of time, but it remains possible that an external observer could glean information from the real-time duration between the arrival of packets that the attacker VM sends. To interfere with this timing channel, we emit packets to an external observer with timing dictated by, again, the median timing of the three VM replicas.

**Justification for the median.** Permitting only the median timing of an IO event to be observed limits the information that an attacker VM can glean from being co-located with a victim VM of interest, because the distribution of the median timings substantially dampens the visibility of a victim’s activities. We formally justify this assertion in the Appendix. Here we simply provide an example illustration of how the median does so.

Consider a victim VM that induces observable timings that are exponentially distributed with rate  $\lambda'$ , versus a baseline (i.e., non-victim) exponential distribution with rate  $\lambda > \lambda'$ .<sup>1</sup> Fig. 1(a) plots example distributions of the attacker VMs’ observations under *StopWatch* when an attacker VM is coresident with the victim (“Median of two baselines, one victim”) and when attacker VM is not (“Median of three baselines”). This figure shows that these median distributions are quite similar, even when  $\lambda$  is substantially larger than  $\lambda'$ ; e.g.,  $\lambda = 1$  and  $\lambda' = 1/2$  in the example in Fig. 1(a). In this case, to even reject the null hypothesis that the attacker VM is not coresident with the victim using a  $\chi$ -square test, the attacker can do so with high confidence in the absence of *StopWatch* with only a single observation, but doing so under *StopWatch* requires almost two orders of magnitude more (Fig. 1(b)). This improvement becomes even more pronounced if  $\lambda$  and  $\lambda'$  are closer; the case  $\lambda = 1$ ,  $\lambda' = 10/11$  is shown in Fig. 1(c).

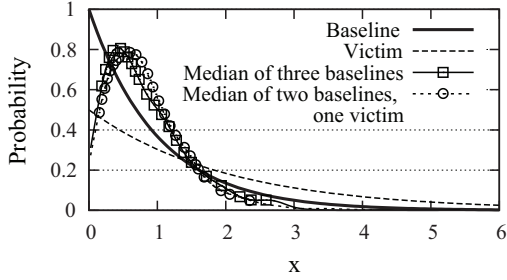
In terms of the number of observations needed to extract meaningful information from the victim VM, this assessment is very conservative, since the attacker would face numerous pragmatic difficulties that we have not modeled here [1]. But even this simple example shows the power of disclosing only median timings of three VM replicas, and in Sec. V-B we will repeat this illustration using actual message traces. Please see the Appendix for a formal analysis.

### IV. RT CLOCKS

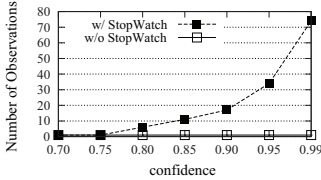
Real-time clocks provide reliable and intuitive reference clocks for measuring the timings of other events. In this section, we describe the high-level strategy taken in *StopWatch* to interfere with their use for timing channels and detail the implementation of this strategy in Xen with hardware-assisted virtualization (HVM).

---

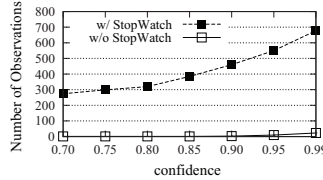
<sup>1</sup>It is not uncommon to model packet inter-arrival time, for example, using an exponential distribution (e.g., [34]).



(a) Distribution of median;  $\lambda' = 1/2$



(b) Observations needed to detect victim;  $\lambda' = 1/2$



(c) Observations needed to detect victim;  $\lambda' = 10/11$

Fig. 1. Justification for median; baseline distribution  $\text{Exp}(\lambda)$ ,  $\lambda = 1$ , and victim distribution  $\text{Exp}(\lambda')$

## A. Strategy

The strategy adopted in *StopWatch* to interfere with a VM’s use of real-time clocks is to virtualize these real-time clocks so that their values observed by a VM are a deterministic function of the VM’s instructions executed so far [33]. That is, after the VM executes  $instr$  instructions, the virtual time observed from within the VM is

$$virt(instr) \leftarrow slope \times instr + start \quad (1)$$

To determine  $start$  at the beginning of VM replica execution, the VMMs hosting the VM’s replicas exchange their current real times;  $start$  is initially set to the median of these values.  $slope$  is initially set to a constant determined by the tick rate of the machines on which the replicas reside.

Optionally, the VMMs can adjust  $start$  and  $slope$  periodically, e.g., after the replicas execute an “epoch” of  $I$  instructions, to coarsely synchronize  $virt$  and real time. For example, after the  $k$ -th epoch, each VMM can send to the others the duration  $D_k$  over which its replica executed those  $I$  instructions and its real time  $R_k$  at the end of that duration. Then, the VMMs can select the median real time  $R_k^*$  and the duration  $D_k^*$  from that same machine and reset

$$start_{k+1} \leftarrow virt_k(I)$$

$$slope_{k+1} \leftarrow \arg \min_{v \in [\ell, u]} \left| \frac{R_k^* - virt_k(I) + D_k^*}{I} - v \right|$$

for a preconfigured constant range  $[\ell, u]$ , to yield the formula for  $virt_{k+1}$ .<sup>2</sup> The use of  $\ell$  and  $u$  ensures that  $slope_{k+1}$  is not too extreme and, if  $\ell > 0$ , that  $slope_{k+1}$  is positive. In this way,  $virt_{k+1}$  should approach real time on the computer contributing the median real time  $R_k^*$  over the next  $I$  instructions, assuming that the machine and VM workloads stay roughly

<sup>2</sup>In other words, if  $(R_k^* - virt_k(I) + D_k^*)/I \in [\ell, u]$  then this value becomes  $slope_{k+1}$ . Otherwise, either  $\ell$  or  $u$  does, whichever is closer to  $(R_k^* - virt_k(I) + D_k^*)/I$ .

the same. Of course, the smaller  $I$ -values are, the more  $virt$  follows real time and so poses the risk of becoming useful in timing attacks. So,  $virt$  should be adjusted only for tasks for which coarse synchronization with real time is important and then only with large  $I$  values.

## B. Implementation in Xen

Real-time clocks on a typical x86 platform include timer interrupts and various hardware counters. Closely related to these real-time clocks is the time stamp counter register, which is accessed using the `rdtsc` instruction and stores a count of processor ticks since reset.

**Timer interrupts.** Operating systems typically measure the passage of time by counting timer interrupts; i.e., the operating system sets up a hardware device to interrupt periodically at a known rate, such as 100 times per second [35]. There are various such hardware devices that can be used for this purpose. Our current implementation of *StopWatch* assumes the guest VM uses a Programmable Interval Timer (PIT) as its timer interrupt source, but our implementation for other sources would be similar. The *StopWatch* VMM generates timer interrupts for a guest on a schedule dictated by that guest’s *virtual* time  $virt$  as computed in Eqn. 1. To do so, it is necessary for the VMM to be able to track the instruction count  $instr$  executed by the guest VM.

In our present implementation, *StopWatch* uses the guest *branch count* for  $instr$ , i.e., keeping track only of the number of branches that the guest VM executes. Several architectures support hardware branch counters, but these are not sensitive to the multiplexing of multiple guests onto a single hardware processor and so continue to count branches regardless of the guest that is currently executing. So, to track the branch count for a guest, *StopWatch* implements a *virtualized* branch counter for each guest.

A question is when to inject each timer interrupt. Intel VT augments IA-32 with two new forms of CPU operations: virtual machine extensions (VMX) root operation and VMX non-root operation [36]. While the VMM uses root operation, guest VMs use VMX non-root operation. In non-root operation, certain instructions and events cause a *VM exit* to the VMM, so that the VMM can emulate those instructions or deal with those events. Once completed, control is transferred back to the guest VM via a *VM entry*. The guest then continues running as if it had never been interrupted.

VM exits give the VMM the opportunity to inject timer interrupts into the guest VM as the guest’s virtual time advances. However, so that guest VM replicas observe the same timer interrupts at the same points in their executions, *StopWatch* injects timer interrupts only after VM exits that are caused by guest execution. Other VM exits can be induced by events external to the VM, such as hardware interrupts on the physical machine; these would generally occur at different points during the execution of the guest VM replicas but will not be visible to the guest [37, Sec. 29.3.2]. For VM exits caused by guest VM execution, the VMM injects any needed timer interrupts on the next VM entry.

**rdtsc calls and CMOS RTC values.** Another way for a guest VM to measure time is via `rdtsc` calls. Xen already

emulates the return values to these calls. More specifically, to produce the return value for a `rdtsc` call, the Xen hypervisor computes the time passed since guest reset using its real-time clock, and then this time value is scaled by a constant factor. *StopWatch* replaces this use of a real-time clock with the guest’s virtual clock (Eqn. 1).

A virtualized real-time clock (RTC) is also provided to HVM guests in Xen; this provides time to the nearest second for the guest to read. The virtual RTC gets updated by Xen using its real-time clock. *StopWatch* responds to requests to read the RTC using the guest’s virtual time.

**Reading counters.** The guest can also observe real time from various hardware counters, e.g., the PIT counter, which repeatedly counts down to zero (at a pace dictated by real time) starting from a constant. These counters, too, are already virtualized in modern VMMs such as Xen. In Xen, these return values are calculated using a real-time clock; *StopWatch* uses the guest virtual time, instead.

## V. IO CLOCKS

IO clocks are typically network, disk and DMA interrupts. (Other device interrupts, such as keyboards, mice, graphics cards, etc., are typically not relevant for guest VMs in clouds.) We outline our strategy for mitigating their use to implement timing channels in Sec. V-A, and then in Sec. V-B we describe our implementation of this strategy in *StopWatch*.

### A. Strategy

The method described in Sec. IV for dealing with RT clocks by introducing virtual time provides a basis for addressing sources of IO clocks. A component of our strategy for doing so is to synchronize I/O events across the three replicas of each guest VM in virtual time, so that every I/O interrupt occurs at the same virtual time at all replicas. Among other things, this synchronization will force uniprocessor VMs to execute deterministically, but it alone will not be enough to interfere with IO clocks; it is also necessary to prevent the timing behavior of one replica’s machine from imposing I/O interrupt synchronization points for the others, as discussed in Sec. II–III. This is simpler to accomplish for disk accesses and DMA transfers since replica VMs initiate these themselves, and so we will discuss this case first. The more difficult case of network interrupts, where we explicitly employ median calculations to dampen the influence of any one machine’s timing behavior on the others, will then be addressed.

**Disk and DMA interrupts.** The replication of each guest VM at start time includes replicating its entire disk image, and so any disk blocks available to one VM replica will be available to all. By virtue of the fact that (uniprocessor) VMs execute deterministically in *StopWatch*, replicas will issue disk and DMA requests at the same virtual time. Upon receiving such a request from a replica at time  $V$ , the VMM adds an offset  $\Delta_d$  to determine a “delivery time” for the interrupt, i.e., at virtual time  $V + \Delta_d$ , and initiates the corresponding I/O activities (disk access or DMA transfer). The offset  $\Delta_d$  must be large enough to ensure that the data transfer completes by the virtual delivery time. Once the virtual delivery time has been determined, the VMM simply waits for the first VM exit caused by the guest VM (as in Sec. IV-B) that occurs at a

virtual time at least as large as this delivery time. The VMM then injects the interrupt prior to the next VM entry of the guest. This interrupt injection also includes copying the data into the address space of the guest, so as to prevent the guest VM from polling for the data in advance of the interrupt to create a form of clock (e.g., see [5, Sec 4.2.2]).

**Network interrupts.** Unlike the initiation of disk accesses and DMA transfers, the activity giving rise to a network interrupt, namely the arrival of a network packet that is destined for the guest VM, is not synchronized in virtual time across the three replicas of the guest VM. So, the VMMs on the three machines hosting these replicas must coordinate to synchronize the delivery of each network interrupt to the guest VM replicas. To prevent the timing of one from dictating the delivery time at all three, these VMMs exchange proposed delivery times and select the median, as discussed in Sec. III. To solicit proposed timings from the three, it is necessary, of course, that the VMMs hosting the three replicas all observe each network packet. So, *StopWatch* replicates every network packet to all three computers hosting replicas of the VM for which the packet is intended. This is done by a logically separate “ingress node” that we envision residing on a dedicated computer in the cloud. (Of course, there need not be only one such ingress for the whole cloud.)

When a VMM observes a network packet to be delivered to the guest, it sends its proposed virtual time — i.e., in the guest’s virtual time, see Sec. IV — for the delivery of that interrupt to the VMMs on the other machines hosting replicas of the same guest VM. (We stress that these proposals are not visible to the guest VM replicas.) Each VMM generates its proposed delivery time by adding a constant offset  $\Delta_n$  to the virtual time of the guest VM at its last VM exit.  $\Delta_n$  must be large enough to ensure that once the three proposals have been collected and the median determined at all three replica VMMs, the chosen median virtual time has not already been passed by any of the guest VMs. The virtual-time offset  $\Delta_n$  is thus determined using an assumed upper bound on the real time it takes for each VMM to observe the interrupt and to propagate its proposal to the others,<sup>3</sup> as well as the maximum allowed difference between the fastest two replicas’ virtual times. This difference can be limited by slowing the execution of the fastest replica.

Once the median proposed virtual time for a network interrupt has been determined at a VMM, the VMM simply waits for the first VM exit caused by the guest VM (as in Sec. IV-B) that occurs at a virtual time at least as large as that median value.<sup>4</sup> The VMM then injects the interrupt prior to the next VM entry of the guest. As with disk accesses and DMA transfers, this interrupt injection also includes copying the data into the address space of the guest, so as to prevent the guest VM from polling for the data in advance of the interrupt to create a form of clock (e.g., see [5, Sec. 4.2.2]).

The process of determining the delivery time of a network

<sup>3</sup>In distributed computing parlance, we thus assume a *synchronous* system, i.e., there are known bounds on processor execution rates and message delivery times.

<sup>4</sup>If the median time determined by a VMM has already passed, then our synchrony assumption was violated by the underlying system. In this case, that VMM’s replica has diverged from the others and so must be recovered by, e.g., copying the state of another replica.

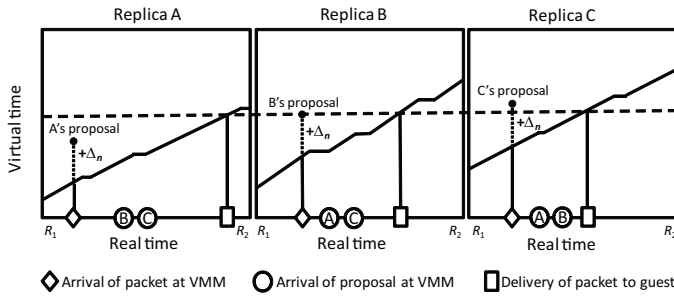


Fig. 2. Delivering a packet to guest VM replicas.

packet to a guest VM’s replicas is pictured in Fig. 2. This figure depicts a real-time interval  $[R_1, R_2]$  at the three machines at which a guest VM is replicated, showing at each machine: the arrival of a packet at the VMM, the proposal made by each VMM, the arrival of proposals from other replica machines, the selection of the median, and the delivery of the packet to the guest replica. Each stepped diagonal line shows the progression of virtual time at that machine.

### B. Implementation in Xen

Xen presents to each HVM guest a virtualized platform that resembles a classic PC/server platform with a network card, disk, keyboard, mouse, graphics display, etc. This virtualized platform support is provided by virtual I/O devices (device models) in Dom0, a domain in Xen with special privileges. QEMU (<http://fabrice.bellard.free.fr/qemu>) is used to implement device models. One instance of the device models is run in Dom0 per HVM domain.

**Network card emulation.** In the case of a network card, the device model running in Dom0 receives packets destined for the guest VM. Without *StopWatch* modification, the device model copies this packet to the guest address space and asserts a virtual network device interrupt via the virtual Programmable Interrupt Controller (vPIC) exposed by the VMM for this guest. HVM guests cannot see real external hardware interrupts since the VMM controls the platform’s interrupt controllers [37, Sec. 29.3.2].

In *StopWatch*, we modify the network card device model so as to place each packet destined for the guest VM into a buffer hidden from the guest, rather than delivering it to the guest. The device model then reads the current virtual time of the guest (as of the guest’s last VM exit), adds  $\Delta_n$  to this virtual time to create its proposed delivery (virtual) time for this packet, and multicasts this proposal to the other two replicas (step 1 in Fig. 3). A memory region shared between Dom0 and the VMM allows device models in Dom0 to read guest virtual time.

Once the network device model receives the two proposals in addition to its own, it takes the median proposal as the delivery time and stores this delivery time in the memory it shares with the VMM. The VMM compares guest virtual time to the delivery time stored in the shared memory upon every guest VM exit caused by guest VM execution. Once guest virtual time has passed the delivery time, the network device model copies the packet into the guest address space (step 2

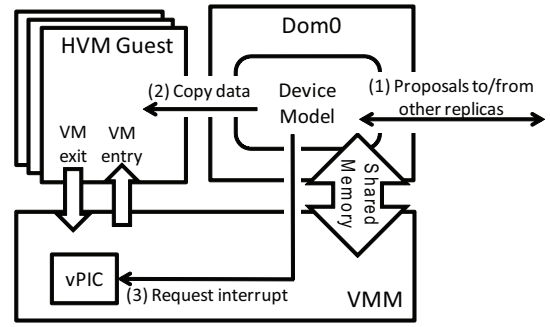


Fig. 3. Emulation of network I/O device in *StopWatch*.

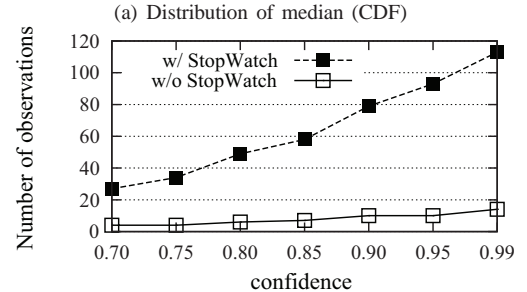
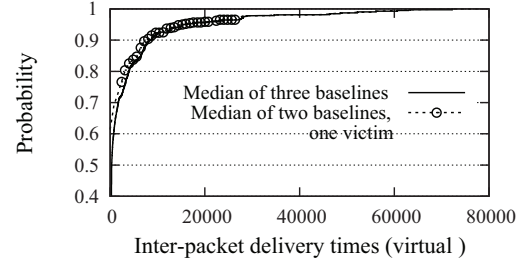


Fig. 4. Virtual inter-packet delivery times to attacker VM replicas with coresident victim (“two baselines, one victim”) and in a run where no replica was coresident with a victim (“three baselines”)

in Fig. 3) and asserts a virtual network interrupt on the vPIC prior to the next VM entry (step 3).

Fig. 4(a) shows the CDF of virtual inter-packet delivery times to replicas of an attacker VM in an actual run where one replica is coresident with a victim VM continuously serving a file, in comparison to the virtual delivery times with no victim present. This plot is directly analogous to that in Fig. 1(a) but is generated from a real *StopWatch* run and shows the distribution as a CDF for ease of readability. Fig. 4(b) shows the number of observations needed to distinguish the victim and no-victim distributions in Fig. 4(a) using a  $\chi$ -squared test, as a function of the desired confidence. This figure is analogous to Fig. 1(b) and confirms that *StopWatch* strengthens defense against timing attacks by an order of magnitude in this scenario. Again, the absolute number of observations needed to distinguish these distributions is likely quite conservative, owing to numerous practical challenges to gathering these observations [1].

**Disk and DMA emulation.** The emulation of the IDE disk and DMA devices is similar to the network card emulation above. *StopWatch* controls when the disk and DMA device

models complete requests and notify the guest. Instead of copying data read to the guest address space, the device model in *StopWatch* prepares a buffer to receive this data. In addition, rather than asserting an appropriate interrupt via the vPIC to the guest as soon as the data is available, the *StopWatch* device model reads the current guest virtual time from memory shared with the VMM, adds  $\Delta_d$ , and stores this value as the interrupt delivery time in the shared memory. Upon the first VM exit caused by guest execution at which the guest virtual time has passed this delivery time, the device model copies the buffered data into the guest address space and asserts an interrupt on the vPIC. Disk writes are handled similarly, in that the interrupt indicating write completion is delivered as dictated by adding  $\Delta_d$  to the virtual time at which the write was initiated.

## VI. EXTERNAL OBSERVERS

The mechanisms described in Sec. IV–V intervene on two significant sources of clocks; though VM replicas can measure the progress of one relative to the other, for example, their measurements will be the same and will reflect the median of their timing behaviors. Moreover, by forcing each guest VM to execute (and, in particular, schedule its internal activities) on the basis of virtual time and by synchronizing I/O events across replicas in virtual time, uniprocessor guest VMs execute deterministically, stripping them of the ability to leverage TL and Mem clocks, as well. (More specifically, the progress of TL and Mem clocks are functionally determined by the progress of virtual time and so are not independent of it.) There nevertheless remains the possibility that an external observer, on whose real-time clock we cannot intervene, could discern information on the basis of the real-time behavior of his attacker VM. In this section we describe our approach to addressing this form of timing channel.

Because guest VM replicas will run deterministically, they will output the same network packets in the same order. *StopWatch* uses this property to interfere with a VM’s ability to exfiltrate information on the basis of its real-time behavior as seen by an external observer. *StopWatch* does so by adopting the median timing across the three guest VM replicas for each output packet. The median is selected at a separate “egress node” that is dedicated for this purpose (c.f., [38]), analogous to the “ingress node” that replicates every network packet destined to the guest VM to the VM’s replicas (see Sec. V). Like the ingress node, there need not be only one egress node for the whole cloud.

To implement this scheme in Xen, every packet sent by a guest VM replica is tunneled by the network device model on that machine to the egress node over TCP. The egress node forwards an output packet to its destination after receiving the second copy of that packet (i.e., the same packet from two guest VM replicas). Since the second copy of the packet it receives exhibits the median output timing of the three replicas, this strategy ensures that the timing of the output packet sent toward its destination is either the timing of a guest replica not coresident with the victim VM or else a timing that falls between those of guest replicas not coresident with the victim.

## VII. PERFORMANCE EVALUATION

In this section we evaluate the performance of our *StopWatch* prototype. We present additional implementation details

that impact performance in Sec. VII-A, our experimental setup in Sec. VII-B, and our tests and their results in Sec. VII-C–VII-D.

### A. Selected implementation details

Our prototype is a modification of Xen version 4.0.2-rc1-pre, amounting to insertions or changes of roughly 1500 source lines of code (SLOC) in the hypervisor. There were also about 2000 SLOC insertions and changes to the QEMU device models distributed with that Xen version. In addition to these changes, we incorporated OpenPGM (<http://code.google.com/p/openpgm/>) into the network device model in Dom0. OpenPGM is a high-performance reliable multicast implementation, specifically of the Pragmatic General Multicast (PGM) specification [39]. In PGM, reliable transmission is accomplished by receivers detecting loss and requesting retransmission of lost data. OpenPGM is used in *StopWatch* for replicating packets destined to a guest VM to all of that VM’s replicas and for communication among the VMMs hosting guest VM replicas.

Recall from Sec. V that each VMM proposes (via an OpenPGM multicast) a virtual delivery time for each network interrupt, and the VMMs adopt the median proposal as the actual delivery time. As noted there, each VMM generates its proposal by adding a constant offset  $\Delta_n$  to the current virtual time of the guest VM.  $\Delta_n$  must be large enough to ensure that by the time each VMM selects the median, that virtual time has not already passed in the guest VM. However, subject to this constraint,  $\Delta_n$  should be minimized since the real time to which  $\Delta_n$  translates imposes a lower bound on the latency of the interrupt delivery. (Note that because  $\Delta_n$  is specified in virtual time and virtual time can vary in its relationship to real time, the exact real time to which  $\Delta_n$  translates can vary during execution.) We selected  $\Delta_n$  to accommodate timing differences in the arrivals of packets destined to the guest VM at its three replicas’ VMMs, the delays for delivering each VMM’s proposed virtual delivery time to the others, and the maximum allowed difference in progress between the two fastest guest VM replicas (which *StopWatch* enforces by slowing the fastest replica, if necessary). For the platform used in our experiments (see Sec. VII-B) and under diverse networking workloads, we found that a value of  $\Delta_n$  that typically translates to a real-time delay in the vicinity of 7–12ms sufficed to meet the above criteria. The analogous offset  $\Delta_d$  for determining the virtual delivery time for disk and DMA interrupts was determined based on the maximum observed disk access times and translates to roughly 8–15ms.

### B. Experimental setup

Our “cloud” consisted of three machines with the same hardware configuration: 4 Intel Core2 Quad Q9650 3.00GHz CPUs, 8GB memory, and a 70GB rotating hard drive. Dom0 was configured to run Linux kernel version 2.6.32.25. Each HVM guest had one virtual CPU, 2GB memory and 16GB disk space. Each guest ran Linux kernel 2.6.32.24 and was configured to use the Programmable Interrupt Controller (PIC) as its interrupt controller and a Programmable Interrupt Timer (PIT) of 250Hz as its clock source. The Advanced Programmable Interrupt Controller (APIC) was disabled. An emulated ATA QEMU disk and a QEMU Realtek RTL-8139/8139C/8139C+



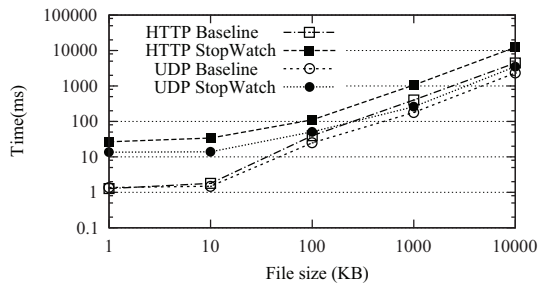


Fig. 5. HTTP and UDP file-retrieval latency.

were provided to the guest as its disk and network card. In each of our tests, we installed an application (e.g., a web server or other program) in the guest VM, as will be described later.

After the guest VM was configured, we copied it to our three machines and restored the VM at each. In this way, our three replicas started running from the same state. In addition, we copied the disk file to all three machines to provide identical disk state to the three replicas.

Once the guest VM replicas were started, inbound packets for this guest VM were replicated to all three machines for delivery to their replicas as discussed in Sec. V. These three machines were attached to a /24 subnet within the UNC campus network, and so broadcast traffic on the network (e.g., ARP requests) was replicated for delivery as in Sec. V. These broadcasts averaged roughly 50-100 packets per second. As such, this background activity was present throughout our experiments and is reflected in our numbers.

### C. Network Services

In this section we describe tests involving network services deployed on the cloud. In all of our tests, our client that interacted with the cloud-resident service was a Lenovo T400 laptop with a dual-core 2.8GHz CPU and 2GB memory attached to an 802.11 wireless network on the UNC campus.

**File downloads.** Our first experiments tested the performance of file download by the client from a web server in the cloud. The total times for the client to retrieve files of various sizes over HTTP are shown in Fig. 5. This figure shows tests in which our guest VM ran Apache version 2.2.14, and the file retrieval was from a cold start (and so file-system caches were empty). The “HTTP Baseline” curve in Fig. 5 shows the average latency for the client to retrieve a file from an unmodified Xen guest VM. The “HTTP *StopWatch*” curve shows the average cost of file retrieval from our *StopWatch* implementation. Every average is for ten runs. Note that both axes are log-scale.

Fig. 5 shows that for HTTP download, a service running on our current *StopWatch* prototype loses less than  $2.8\times$  in download speed for files of 100KB or larger. Diagnosing this cost reveals that the bottleneck, by an order of magnitude or more, was the network transmission delay (vs. disk access delay) in both the baseline and for *StopWatch*. Moreover, the performance cost of *StopWatch* in comparison to the baseline was dominated by the time for delivery of *inbound* packets to the web-server guest VM, i.e., the TCP SYN and ACK messages in the three-way handshake, and then additional

acknowledgments sent by the client. Enforcing a median timing on output packets (Sec. VI) adds modest overhead in comparison.

This combination of insights, namely the detriment of inbound packets (mostly acknowledgments) to *StopWatch* file download performance and the fact that these costs so outweigh disk access costs, raises the possibility of recovering file download performance using a transport protocol that minimizes packets inbound to the web server, e.g., using negative acknowledgments or forward error correction. Alternatively, an unreliable transport protocol with no acknowledgments, such as UDP, could be used; transmission reliability could then be enforced at a layer above UDP using negative acknowledgments or forward error correction. Though TCP does not define negative acknowledgments, transport protocols that implement reliability using them are widely available, particularly for *multicast* where positive acknowledgments can lead to “ack implosion.” Indeed, recall that the PGM protocol specification [39], and so the OpenPGM implementation that we use, ensures reliability using negative acknowledgments.

To illustrate this point, in Fig. 5 we repeat the experiments using UDP to transfer the file.<sup>5</sup> The “UDP Baseline” curve shows the performance using unmodified Xen; the “UDP *StopWatch*” curve shows the performance using *StopWatch*. Not surprisingly, baseline UDP shows performance comparable to (but slightly more efficient than, by less than a factor of two) baseline TCP, but rather than losing an order of magnitude, UDP over *StopWatch* is *competitive* with these baseline numbers for files of 100KB or more.

**NFS.** We also set up a Network File System (NFSv4) server in our guest VM. On our client machine, we installed an NFSv4 client; remotely mounted the filesystem exported by the NFS server; performed file operations manually; and then ran `nfsstat` on the NFS server to print its server-side statistics, including the mix of operations induced by our activity. We then used the `nfsstone` benchmarking utility to evaluate the performance of the NFS server with and without *StopWatch*. `nfsstone` generates an artificial load with a specified mix of NFS operations. The mix of NFS operations used in our tests was the previously extracted mix file.<sup>6</sup> In each test, the client machine ran five processes using the mounted file system, making calls at a constant rate ranging from 25 to 400 per second in total across the five client processes.

The average latency per operation is shown in Fig. 6(a). In this figure, the horizontal axis is the rate at which operations were submitted to the server; note that this axis is log-scale. Fig. 6(a) suggests that an NFS server over *StopWatch* incurs a less than  $2.7\times$  increase in latency over an NFS server running over unmodified Xen. Since the NFS implementation used TCP, in some sense this is unsurprising in light of the file download results in Fig. 5. That said, it is also perhaps surprising that *StopWatch*’s cost increased only roughly log-

<sup>5</sup>We are not advocating UDP for file retrieval generally but rather are simply showing the advantages for *StopWatch* of a protocol that minimizes client-to-server packets. We did not use OpenPGM in these tests since the web site (as the “multicast” originator) would need to initiate the connection to the client; this would have required more substantial modifications. This “directionality” issue is not fundamental to negative acknowledgments, however.

<sup>6</sup>This mix was 11.37% `setattr`, 24.07% `lookup`, 11.92% `write`, 7.93% `getattr`, 32.34% `read` and 12.37% `create`.

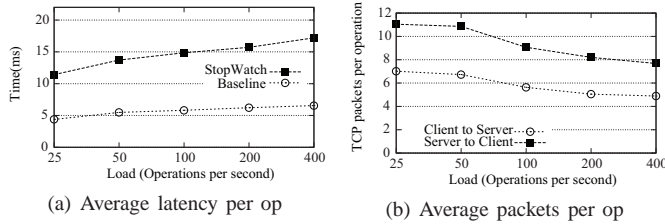


Fig. 6. Tests of NFS server using `nhfsstone`

arithmically as a function of the offered rate of operations. This modest growth is in part because *StopWatch* schedules packets for delivery to guest VM replicas independently — the scheduling of one does not depend on the delivery of a previous one, and so they can be “pipelined” — and because the number of TCP packets from the client to the server actually decreases per operation, on average, as the offered load grows (Fig. 6(b)).

#### D. Computations

In this section we evaluate the performance of various computations on *StopWatch* that may be representative of future cloud workloads. For this purpose, we employ the PARSEC benchmarks [40]. PARSEC is a diverse set of benchmarks that covers a wide range of computations that are likely to become important in the near future (see <http://parsec.cs.princeton.edu/overview.htm>). Here we take PARSEC as representative of future cloud workloads.

We utilized the following five applications from the PARSEC suite (version 2.1), providing each the “native” input designated for it. `ferret` is representative of next-generation search engines for non-text document data types. In our tests, we configured the application for image similarity search. `blackscholes` calculates option pricing with Black-Scholes partial differential equations and is representative of financial analysis applications. `canneal` is representative of engineering applications and uses simulated annealing to optimize routing cost of a chip design. `dedup` represents next-generation backup storage systems characterized by a combination of global and local compression. `streamcluster` is representative of data mining algorithms for online clustering problems. Each of these applications involves various activities, including initial configuration, creating a local directory for results, unpacking input files, performing its computation, and finally cleaning up temporary files.

We ran each benchmark ten times in one guest VM over unmodified Xen, and then ten more times with three guest VM replicas over *StopWatch*. Fig. 7(a) shows the average runtimes of these applications in both cases. In this figure, each application is described by two bars; the black bar on the left shows its performance over unmodified Xen, and the gray bar on the right shows its performance over *StopWatch*. *StopWatch* imposed an overhead of at most  $2.3\times$  (for `blackscholes`) to the average running time of the applications. Owing to the dearth of network traffic involved in these applications, the overhead imposed by *StopWatch* is mostly due to the overhead involved in intervening on disk I/O (see Sec. V). As shown in Fig. 7(b), there is a direct correlation between the number of disk interrupts to deliver during the application run and the performance penalty (in absolute terms) that *StopWatch*

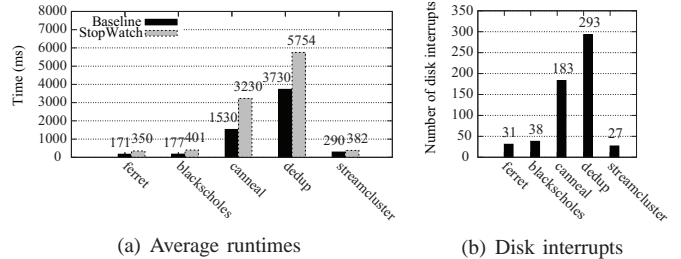


Fig. 7. Tests of PARSEC applications

imposes. If the computers in our experiments used solid-state drives (versus hard disks), we conjecture that their reduced access times would permit us to shrink  $\Delta_d$  and so improve the performance of *StopWatch* for these applications.

### VIII. REPLICA PLACEMENT IN THE CLOUD

*StopWatch* requires that the three replicas of each guest VM are coresident with nonoverlapping sets of (replicas of) other VMs. This constrains how a cloud operator places guest VM replicas on its machines. In this section we clarify the significance of these placement constraints in terms of the provider’s ability to best utilize its infrastructure. After all, if under these constraints, the provider were able to simultaneously run a number of guest VMs that scales, say, only linearly in the number of cloud nodes, then the provider should forgo *StopWatch* and simply run each guest VM (non-replicated) in isolation on a separate node. Here we show that the cloud operator is not limited to such poor utilization of its machines.

If the cloud has  $n$  machines, then consider the complete, undirected graph (clique)  $K_n$  on  $n$  vertices, one per machine. For every guest VM, the placement of its three replicas forms a *triangle* in  $K_n$  consisting of the vertices for the machines on which the replicas are placed and the edges between those vertices. The placement constraints of *StopWatch* can be expressed by requiring that the triangles representing VM replica placements be pairwise *edge-disjoint*. As such, the number of guest VMs that can simultaneously be run on a cloud of  $n$  machines is the same as the number of edge-disjoint triangles that can be *packed* into  $K_n$ . A corollary of a result due to Horsley [41, Thm. 1.1] is:

*Theorem 1:* A maximum packing of  $K_n$  with pairwise edge-disjoint triangles has exactly  $k$  triangles, where: (i) if  $n$  is odd, then  $k$  is the largest integer such that  $3k \leq \binom{n}{2}$  and  $\binom{n}{2} - 3k \notin \{1, 2\}$ ; and (ii) if  $n$  is even, then  $k$  is the largest integer such that  $3k \leq \binom{n}{2} - \frac{n}{2}$ .

So, a cloud of  $n$  machines using *StopWatch* can simultaneously execute  $k = \Theta(n^2)$  guest VMs. The existence of such a placement, however, does not guarantee an efficient algorithm to find it. Moreover, this theorem ignores machine capacities. Below we address both of these shortcomings.

Under the constraints of *StopWatch*, one node in a cloud of  $n$  nodes can simultaneously execute up to  $\frac{n-1}{2}$  guest VMs, since the other replicas of the guest VMs that it executes (two per VM) must occupy distinct nodes. If each node has resources to simultaneously execute  $c \leq \frac{n-1}{2}$  guest VMs, then

the following theorem provides for an algorithm to efficiently place them subject to the per-machine capacity constraint  $c$ :

*Theorem 2:* Let  $n \equiv 3 \pmod 6$  and  $c \leq \frac{n-1}{2}$ . If  $c \equiv 0$  or  $1 \pmod 3$ , then there is an efficient algorithm to place  $k \leq \frac{1}{3}cn$  guest VMs. If  $c \equiv 2 \pmod 3$ , then there is an efficient algorithm to place  $k \leq \frac{1}{3}(c-1)n + \frac{n-3}{6}$  guest VMs.

*Proof:* Following Bose’s construction of a Steiner Triple System [42], let  $n = 6v + 3$  and let  $(Q, \circ)$  be a multiplicative idempotent commutative quasigroup of order  $2v + 1$ , where  $Q = \{a_0, a_1, \dots, a_{2v}\}$ . An idempotent commutative quasigroup of size  $2v + 1$  has the property that its multiplication table is an idempotent, commutative matrix of size  $(2v + 1) \times (2v + 1)$ , and each element of  $Q$  appears exactly once in each row and in each column. Let  $Q \times \{0, 1, 2\}$  denote the  $n$  nodes, and consider the following sets  $G_t$ ,  $0 \leq t \leq v$ , of triangles:

$$G_0 = \bigcup_{0 \leq i \leq 2v} \{(a_i, 0), (a_i, 1), (a_i, 2)\}$$

and for  $1 \leq t \leq v$ ,

$$G_t = \bigcup_{\substack{0 \leq i \leq 2v \\ 0 \leq \ell \leq 2}} \{(a_i, \ell), (a_j, \ell), (a_i \circ a_j, \ell + 1 \pmod 3)\}$$

where  $j = i + t \pmod{2v + 1}$ .

There are  $2v+1$  triangles in  $G_0$  and  $(2v+1) \times 3 = 6v+3 = n$  triangles in  $G_t$  for each  $1 \leq t \leq v$ . Moreover, all of these triangles are edge-disjoint. Triangles in  $G_0$  visit each of the  $n$  nodes exactly once. Triangles in any  $G_t$ ,  $1 \leq t \leq v$ , visit each node  $(a^*, \ell^*)$  exactly three times: when  $a^* = a_i$  and  $\ell^* = \ell$ ; when  $a^* = a_j$  for  $j = i + t \pmod{2v + 1}$  and  $\ell^* = \ell$ ; and when  $a^* = a_i \circ a_j$  for  $j = i + t \pmod{2v + 1}$  and  $\ell^* = \ell + 1 \pmod 3$ . So, collectively the triangles in  $G_0, \dots, G_v$  visit each node  $3v + 1 = \frac{n-1}{2} \geq c$  times.

So, if  $c \equiv 0 \pmod 3$ , then we can place  $k \leq \frac{1}{3}cn$  VMs using the  $\frac{1}{3}cn$  triangles in groups  $G_1, \dots, G_{c/3}$ . If  $c \equiv 1 \pmod 3$ , then we can place  $k \leq \frac{1}{3}cn$  VMs by first using the  $2v + 1 = \frac{n}{3}$  triangles in  $G_0$  and then the  $\frac{1}{3}(c-1)n$  triangles in  $G_1, \dots, G_{(c-1)/3}$ . If  $c \equiv 2 \pmod 3$ , then we can place  $k \leq \frac{1}{3}(c-1)n + \frac{n-3}{6}$  VMs by first using the  $2v + 1 = \frac{n}{3}$  triangles in  $G_0$ , then  $\frac{1}{3}(c-2)n$  triangles in  $G_1, \dots, G_{(c-2)/3}$ , and finally any  $v = \frac{n-3}{6}$  triangles from  $G_v$  that visit each node at most one time (e.g.,  $\{(a_i, 0), (a_j, 0), (a_i \circ a_j, 1)\}$  for  $0 \leq i \leq v-1$  and  $j = i + v$ ). ■

## IX. COLLABORATING ATTACKER VMs

Our discussion so far has not explicitly addressed the possibility of attacker VMs collaborating to mount timing attacks. The apparent risks of such collaboration can be seen in the following possibility: replicas of one attacker VM (“VM1”) reside on machines A, B, and C; one replica of another attacker VM (“VM2”) resides on machine A; and a replica of the victim VM resides on machine C. If VM2 induces significant load on its machines, then this may slow the replica of VM1 on machine A to an extent that marginalizes its impact on median calculations among its replicas’ VMMs. The replicas of VM1 would then observe timings influenced by the larger

of the replicas on B and C — which may well reflect timings influenced by the victim.

Mounting such an attack, or any collaborative attack involving multiple attacker VMs on one machine, appears to be difficult, however. Just as the reasoning in Fig. 1 and its confirmation in Fig. 4 suggest that an attacker VM detecting its coresidence with a victim VM is made much harder by *StopWatch*, one attacker VM detecting coresidence with another using timing covert channels would also be impeded by *StopWatch*. If the cloud takes measures to avoid disclosing coresidence of one VM with another by other channels, it should be difficult for the attacker to even detect when he is in a position to mount such an attack or to interpret the results of mounting such an attack indiscriminately.

If such attacks are nevertheless feared, they can be made harder still by increasing the number of replicas of each VM. If the number were increased from three to, say, five, then inducing sufficient load to marginalize one attacker replica from its median calculations would not substantially increase the attacker’s ability to mount attacks on a victim. Rather, the attacker would need to marginalize multiple of its replicas, along with accomplishing the requisite setup to do so.

## X. CONCLUSION

We proposed a new method to address timing side channels in IaaS compute clouds that employs three-way replication of guest VMs and placement of these VM replicas so that they are coresident with nonoverlapping sets of (replicas of) other VMs. By permitting these replicas to observe only virtual (vs. real) time and the median timing of network events across the three replicas, we suppress their ability to glean information from a victim VM with which one is coresident. We described an implementation of this technique in Xen, yielding a system called *StopWatch*, and we evaluated the performance of *StopWatch* on a variety of workloads. Though the performance cost for our current prototype ranges up to  $2.8\times$  for networking applications, we used our evaluation to identify the sources of costs and alternative application designs (e.g., reliable transmission using negative acknowledgments, to support serving files) that can enhance performance considerably. We showed that clouds with  $n$  machines capable of each running  $c \leq \frac{n-1}{2}$  guest VMs simultaneously can efficiently schedule  $\Theta(cn)$  guest VMs under the constraints of *StopWatch*, a clear improvement over the alternative of running guest VMs in isolation. Finally, in the Appendix we analyze the median as a microaggregation function and explain its benefits over the alternative of obscuring event timings through the addition of random noise. We envision *StopWatch* as a basis for a high-security cloud, e.g., suitable for military, intelligence, or financial communities with high assurance needs.

**Acknowledgments.** This work was supported in part by NSF grant 0910483, the Science of Security Lablet at North Carolina State University, and a grant from VMWare.

## REFERENCES

- [1] Y. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart, “Cross-VM side channels and their use to extract private keys,” in *19th ACM CCS*, 2012.

- [2] T. Ristenpart, E. Tromer, H. Shacham, and S. Savage, "Hey, you, get off of my cloud: Exploring information leakage in third-party compute clouds," in *16th ACM CCS*, 2009.
- [3] J. Domingo-Ferrer and V. Torra, "Median-based aggregation operators for prototype construction in ordinal scales," *Intern. J. Intel. Sys.*, vol. 18, no. 6, 2003.
- [4] M. E. Kabir and H. Wang, "Microdata protection method through microaggregation: A median-based approach," *Information Security J.: A Global Perspective*, vol. 20, 2011.
- [5] W.-M. Hu, "Reducing timing channels with fuzzy time," in *1991 IEEE Symp. Security & Privacy*, 1991.
- [6] J. Agat, "Transforming out timing leaks," in *27th ACM POPL*, 2000.
- [7] S. Zdancewic and A. C. Myers, "Observational determinism for concurrent program security," in *16th IEEE CSFW*, 2003.
- [8] D. Zhang, A. Askarov, and A. C. Myers, "Language-based control and mitigation of timing channels," in *33rd ACM PLDI*, 2012.
- [9] E. Tromer, D. A. Osvik, and A. Shamir, "Efficient cache attacks on AES, and countermeasures," *J. Cryptology*, vol. 23, no. 1, 2010.
- [10] H. Raj, R. Nathuji, A. Singh, and P. England, "Resource management for isolation enhanced cloud services," in *ACM CCSW*, 2009.
- [11] T. Kim, M. Peinado, and G. Mainar-Ruiz, "STEALTH-MEM: System-level protection against cache-based side channel attacks in the cloud," in *21st USENIX Security Symp.*, 2012.
- [12] A. Askarov, A. C. Myers, and D. Zhang, "Predictive black-box mitigation of timing channels," in *17th ACM CCS*, 2010.
- [13] B. C. Vattikonda, S. Das, and H. Shacham, "Eliminating fine grained timers in Xen," in *ACM CCSW*, 2011.
- [14] M. H. Kang and I. S. Moskowitz, "A pump for rapid, reliable, secure communication," in *ACM CCS*, 1993.
- [15] J. Giles and B. Hajek, "An information-theoretic and game-theoretic study of timing channels," *IEEE TOIT*, vol. 48, no. 9, 2002.
- [16] A. Haeberlen, B. C. Pierce, and A. Narayan, "Differential privacy under fire," in *20th USENIX Security Symp.*, 2011.
- [17] D. Zhang, A. Askarov, and A. C. Myers, "Predictive mitigation of timing channels in interactive systems," in *18th ACM CCS*, 2011.
- [18] F. B. Schneider, "Implementing fault-tolerant services using the state machine approach: A tutorial," *ACM Comp. Surv.*, vol. 22, no. 4, 1990.
- [19] A. Borg, W. Blau, W. Graetsch, F. Herrmann, and W. Oberle, "Fault tolerance under UNIX," *ACM TOCS*, vol. 7, no. 1, 1989.
- [20] P. Narasimhan, L. E. Moser, and P. M. Melliar-Smith, "Enforcing determinism for the consistent replication of multithreaded CORBA applications," in *IEEE SRDS*, 1999.
- [21] C. Basile, Z. Kalbarczyk, and R. K. Iyer, "Active replication of multithreaded applications," *IEEE TPDS*, vol. 17, no. 5, 2006.
- [22] T. C. Bressoud and F. B. Schneider, "Hypervisor-based fault-tolerance," *ACM TOCS*, vol. 14, no. 1, 1996.
- [23] F. B. Schneider, "Understanding protocols for Byzantine clock synchronization," Department of Computer Science, Cornell University, Tech. Rep. 87-859, 1987.
- [24] D. Gao, M. K. Reiter, and D. Song, "Behavioral distance for intrusion detection," in *8th RAID*, 2005.
- [25] B. Cox, D. Evans, A. Filipi, J. Rowanhill, W. Hu, J. Davidson, J. Knight, A. Nguyen-Tuong, and J. Hiser, "N-variant systems: A secretless framework for security through diversity," in *15th USENIX Security Symp.*, 2006.
- [26] A. Nguyen-Tuong, D. Evans, J. C. Knight, B. Cox, and J. W. Davidson, "Security through redundant data diversity," in *38th DSN*, 2008.
- [27] D. Gao, M. K. Reiter, and D. Song, "Beyond output voting: Detecting compromised replicas using HMM-based behavioral distance," *IEEE TDSC*, vol. 6, no. 2, 2009.
- [28] M. Xu, V. Malyugin, J. Sheldon, G. Venkitachalam, and B. Weissman, "ReTrace: Collecting execution trace with virtual machine deterministic replay," in *3rd Workshop on Modeling, Benchmarking and Simulation*, 2007.
- [29] G. W. Dunlap, D. G. Lucchetti, P. M. Chen, and M. A. Fetterman, "Execution replay of multiprocessor virtual machines," in *4th ACM VEE*, 2008.
- [30] J. Devietti, B. Lucia, L. Ceze, and M. Oskin, "DMP: Deterministic shared memory multiprocessing," *IEEE Micro*, vol. 30, pp. 41–49, 2010.
- [31] A. Aviram, S.-C. Weng, S. Hu, and B. Ford, "Efficient system-enforced deterministic parallelism," in *9th USENIX OSDI*, 2010.
- [32] J. C. Wray, "An analysis of covert timing channels," in *1991 IEEE Symp. Security & Privacy*, 1991.
- [33] G. Popek and C. Kline, "Verifiable secure operating system software," in *AFIPS National Comp. Conf.*, 1974.
- [34] T. Karagiannis, M. Molle, M. Faloutsos, and A. Broido, "A nonstationary Poisson view of Internet traffic," in *INFOCOM*, 2004.
- [35] *Timekeeping in VMware Virtual Machines*, VMware Inc., 2010.
- [36] R. Uhlig, G. Neiger, D. Rodgers, A. L. Santoni, F. C. M. Martins, A. V. Anderson, S. M. Bennett, A. Kagi, F. H. Leung, and L. Smith, "Intel virtualization technology," *IEEE Comp.*, vol. 38, no. 3, 2005.
- [37] *Intel 64 and IA-32 Architectures Software Developer's Manual*, Intel Corporation, 2011.
- [38] J. Yin, A. Venkataramani, J.-P. Martin, L. Alvisi, and M. Dahlin, "Byzantine fault-tolerant confidentiality," in *Inter. Workshop on Future Directions in Distributed Computing*, 2002.
- [39] T. Speakman, et al., "PGM reliable transport protocol specification," Request for Comments 3208, Internet Engineering Task Force, 2001.
- [40] C. Bienia, "Benchmarking modern multiprocessors," Ph.D. dissertation, Princeton University, 2011.
- [41] D. Horsley, "Maximum packing of the complete graph with uniform length cycles," *J. Graph Theory*, vol. 68, no. 1, 2011.
- [42] C. C. Lindner and C. A. Rodger, *Design Theory*. CRC Press, 2008, ch. 1.
- [43] E. Deza and M. Deza, *Dictionary of Distances*. Elsevier, 2006.
- [44] M. Güngör, Y. Bulut, and S. Çalık, "Distributions of order statistics," *Appl. Math. Sci.*, vol. 3, no. 16, 2009.

APPENDIX

Here we justify the use of the median as a microaggregation function in *StopWatch*. Let  $X_{r:m}$  denote the random variable that takes on the value of the  $r$ -th smallest of the  $m$  values obtained by sampling random variables  $X_1 \dots X_m$ . Let  $F_i(x)$  denote the CDF of  $X_i$  (i.e.,  $F_i(x) = \mathbb{P}(X_i \leq x)$ ) and let  $F_{r:m}(x)$  denote the CDF of  $X_{r:m}$ .

**Utility of the median.** The security of *StopWatch* (with  $m = 3$  replicas per VM) hinges on the distribution of the median  $X_{2:3}$  of three random variables  $X_1, X_2, X_3$ . In the case of delivering a packet to the attacker VM (Sec. V),  $X_1, X_2, X_3$  correspond to the proposed virtual delivery times of the packet to the three replicas, less the actual virtual delivery time of the previous packet. In the case of the attacker VM sending a packet to an external observer (Sec. VI),  $X_1, X_2, X_3$  correspond to the emission (real) time of the packet from each attacker VM replica, less the emission time of the preceding packet to the external observer. In either case,  $X_1, X_2, X_3$  are independent.

The attack considered in this paper is one in which the adversary learns information due to the difference between (i) the CDF  $F_{2:3}(x)$  for random variables  $X_1, X_2, X_3$  corresponding to attacker VM replicas that are *not* coresident with a victim VM of interest, and (ii) the CDF  $F'_{2:3}(x)$  for random variables  $X'_1, X_2, X_3$  where  $X'_1$  corresponds to an attacker VM that *is* coresident with the victim VM of interest. An example measure of the distance between two CDFs  $F(x)$  and  $\hat{F}(x)$  is their Kolmogorov-Smirnov distance [43, p. 179], defined as  $D(F, \hat{F}) = \max_x |F(x) - \hat{F}(x)|$ . The following theorem shows that adopting the median microaggregation function can only interfere with the adversary's goal:

*Theorem 3:* If the distributions of  $X_2$  and  $X_3$  are overlapping (i.e., for no  $x$  is  $F_2(x) = 0$  and  $F_3(x) = 1$ , or  $F_2(x) = 1$  and  $F_3(x) = 0$ ), then  $D(F_{2:3}, F'_{2:3}) < D(F_1, F'_1)$ .

*Proof:* Due to well-known results in order statistics (e.g., see GÜNGÖR et al. [44, Result 2.4]):<sup>7</sup>

$$F_{r:m}(x) = \sum_{\ell=r}^m (-1)^{\ell-r} \binom{\ell-1}{r-1} \sum_{\substack{I \subseteq \{1 \dots m\}: \\ |I|=\ell}} \prod_{i \in I} F_i(x)$$

In particular,

$$F_{2:3}(x) = F_1(x)F_2(x) + F_1(x)F_3(x) + F_2(x)F_3(x) - 2F_1(x)F_2(x)F_3(x)$$

$$F'_{2:3}(x) = F'_1(x)F_2(x) + F'_1(x)F_3(x) + F_2(x)F_3(x) - 2F'_1(x)F_2(x)F_3(x)$$

where  $F'_1(x)$  represents the CDF of  $X'_1$ . So,

$$D(F_{2:3}, F'_{2:3}) = \max_x |[F_2(x) + F_3(x) - 2F_2(x)F_3(x)][F_1(x) - F'_1(x)]|$$

Noting that  $D(F_1, F'_1) = \max_x |F_1(x) - F'_1(x)|$ , it suffices to show that  $|F_2(x) + F_3(x) - 2F_2(x)F_3(x)| < 1$  for all  $x$ . However, since  $F_2(x) \in [0, 1]$  and  $F_3(x) \in [0, 1]$  for all  $x$ ,

$|F_2(x) + F_3(x) - 2F_2(x)F_3(x)| \leq 1$  and, moreover, equals 1 only if for some  $x$ , one of  $F_2(x)$  and  $F_3(x)$  is 1 and the other is 0. This last case is precluded by the theorem. ■

In the limit, when the distributions of  $X_2$  and  $X_3$  overlap exactly, we get a much stronger result:

*Theorem 4:* If  $X_2$  and  $X_3$  are identically distributed, then  $D(F_{2:3}, F'_{2:3}) \leq \frac{1}{2}D(F_1, F'_1)$ .

*Proof:* In this case,  $F_2 = F_3$  and so

$$|F_2(x) + F_3(x) - 2F_2(x)F_3(x)|$$

reaches its maximum value of  $\frac{1}{2}$  at the value  $x$  yielding  $F_2(x) = F_3(x) = \frac{1}{2}$ . ■

**Comparison to uniformly random noise.** An alternative to *StopWatch* is simply adding random noise (without replicating VMs) to confound timing attacks. For simplicity, suppose that  $X_1$  and  $X'_1$  are exponentially distributed with rate parameters  $\lambda$  and  $\lambda'$ , respectively, as in the example of Fig. 1. For the random variable  $X_N$  representing added noise, assume that  $X_N$  is drawn uniformly from  $[0, b]$  (i.e.,  $X_N \sim U(0, b)$ ), a common choice to mitigate timing channels (e.g., [5], [15]).

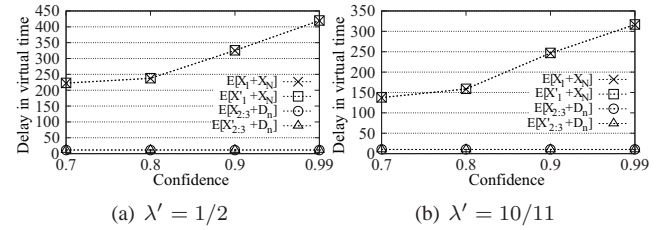


Fig. 8. Expected delay induced by *StopWatch* vs. by uniform noise, as a function of confidence with which attacker distinguishes the two distributions (coresident victim or not) after the same number of observations; baseline distribution  $\text{Exp}(\lambda)$ ,  $\lambda = 1$ ; victim distribution  $\text{Exp}(\lambda')$

We calculated expected delay imposed by *StopWatch* and by adding uniformly distributed noise. To make a fair comparison, we configured both approaches to provide the same strength of defense against timing attacks. Specifically, after calculating the number of observations the attacker requires in the case of *StopWatch* to distinguish, for a fixed confidence level, the distributions  $X_{2:3} + \Delta_n$  and  $X'_{2:3} + \Delta_n$  using a  $\chi$ -squared test, we calculated the minimum  $b$  that would give the attacker the same confidence in distinguishing  $X_1 + X_N$  and  $X'_1 + X_N$  after that number of observations. Fig. 8 shows the resulting expected delays in each case.

This figure indicates that *StopWatch* scales much better as the attacker's required confidence and the distinctiveness of the victim grows (as represented by  $\lambda'$  dropping). The delay of the *StopWatch* approach is tied most directly to  $\Delta_n$ , which is added to ensure that the replicas of each VM remain synchronized (see Section V-A); here we calculated it so that  $\Pr[|X_1 - X'_1| \leq \Delta_n] \geq 0.9999$ . That is, the probability of a desynchronization at this event is less than 0.0001. Note that  $E[X_{2:3} + \Delta_n]$  and  $E[X'_{2:3} + \Delta_n]$  are nearly the same in Fig. 8, since their difference is how the attacker differentiates the two, and similarly for  $E[X_1 + X_N]$  and  $E[X'_1 + X_N]$ .

<sup>7</sup>This equation assumes each  $F_i(x)$  is continuous. See GÜNGÖR et al. [44] for the case when some  $F_i(x)$  is not continuous.