

New Cloud Architectures For The Next Generation Internet

Fangfei Zhou
Northeastern University
Boston, Massachusetts
April, 2013

Abstract

Cloud computing has ushered in a new paradigm with the availability of computing as a service, letting customers share the same physical infrastructure and purchase computing resources on demand (e.g. Amazon EC2 and Windows Azure). The multi-tenancy property of cloud computing offers clients flexibility while creating a unique set of challenges in areas such as reliability and security.

In this thesis we study one challenge (SecureCloud) and three opportunities (DNSCloud, WebCloud and SamaritanCloud). In SecureCloud we explore how multi-tenancy, or the sharing of resources across users, can lead to the undermining of privacy and security. Taking Amazon EC2 as an example we identify an important scheduling vulnerability in Virtual Machine Monitors (VMMs). We create an attack scenario and demonstrate how it can be used to steal cycles in the cloud. We also discuss how attacks can be coordinated across the cloud on a collection of VMs. We present a general framework of solutions to combat such attacks. DNSCloud, WebCloud and SamaritanCloud are proposals for new architectures that improve delivery of existing infrastructural services and enable entirely new functionalities. The Domain Name System (DNS) has long been recognized as the Achilles' heel of the Internet and a variety of new (cache-poisoning and other) attacks surfacing over the past few years have only served to reinforce that notion. We present DNSCloud, a new architecture for providing a more robust DNS service that does not require a forklift upgrade (unlike DNSSEC). Today, content on Web sites like online social networks is created at the edge of network but distributed using a traditional client-server model. WebCloud is a novel cloud architecture that

leverages the burgeoning phenomenon of social networks for enabling a more efficient and scalable system for peer-to-peer content delivery. SamaritanCloud is a proposal for a new architecture that exploits the mobility of personal computing devices to share relevant locality-specific information. It allows people to offer physical help to each other remotely, in a secure and private way. Taken as a whole this thesis represents a synthesis of theory and practice that will hasten the ongoing transition to the era of cloud computing.

Introduction

1.1 Motivation

Over the past decades, data and applications have been slowly but steadily migrating to the Internet (or its more fashionable synonym, the “Cloud”), with services being delivered to end-users through a client-server model using standardized over-the-web protocols. This outsourcing to the Internet continues to be a work-in-progress with cloud-computing the phrase du jour.

Cloud computing [20] is built on top of server virtualization technology [21], which allows multiple instances of virtual machines running on a single physical server. A cloud computing model is comprised of a front end and a back end. The front end is the interface for users to interact with the system – for example, users can open and edit files through browser, or remotely connect to virtual machines to enjoy applications and other computing resources. The back end is the cloud itself, it maintains data and applications and delivers them to end users as a service through the Internet.

Some of the biggest cloud computing services include Web-based email (e.g. Hotmail, Gmail), social networking (e.g. Facebook, Twitter, LinkedIn), document hosting services, (e.g. Google Docs, Flickr, Picasa, YouTube) and back up services (e.g. Dropbox), etc. Major corporations including Amazon, Google, Microsoft and Oracle have invested in cloud computing services and offer individuals and businesses a range of cloud-based solutions, e.g. Amazon EC2 [1] and Microsoft

Azure [6]. In these services, users rent virtual machines for running their own applications and are charged by the amount of time their virtual machine is running (in hours or months).

Cloud computing is a new and different way to architect and remotely manage computing resources. Therefore, it is important to study both its shortcomings and its advantages which could help resolve current Internet issues. Here are some of the significant advantages.

- *Lower cost:* Cloud computing requires lower cost for data processing when compared with the older model of maintaining software on local machines. The use of the cloud removes the need for purchasing software and hardware and shifts the costs to a pay-per-use model.
- *Mobility and flexibility:* As the service is delivered through the Internet, users can access the resources no matter where they are located. Instant deployment and the pay-as-you-go business model offer users more flexibility than other computing services. Users do not need to worry about application installation and updates.
- *More capacity:* Cloud computing offers virtually limitless storage and computing power, thus enhancing the user's capability beyond the capacity of their local machines.
- *Device independence:* Finally, the ultimate cloud computing advantage is that users are no longer tethered to a single computer or network. Existing applications and data follow users through the cloud.

On the other hand, cloud computing has created its own challenges. One of the biggest concerns about cloud computing is security. Cloud computing also requires a constant-on and high-bandwidth Internet connection, which is a limitation for some clients. As data and application are stored on provider owned machines, clients lose their control over the software and become dependent on the provider to maintain, update and manage it. Cloud computing can also bring risks in the areas

of privacy and confidentiality. Clients storing their sensitive data and information on third-party servers may become compromised if the provider has inadequate security in terms of technology or processes.

Today's cloud computing services are typically categorized into the following three classes of service: Software as a Service (SaaS), Platform as a Service (PaaS) and Infrastructure as a Service (IaaS), each with their own security issues.

- *SaaS*: Cloud infrastructures providing software services often employ a multi-tenancy system architecture, where users access the same environment for different applications. Optimization of speed, security, availability, recovery are required in such services.
- *PaaS*: Cloud infrastructures allowing clients to develop cloud services and applications should provide stable programming environment, tools and configuration management.
- *IaaS*: Cloud infrastructures delivering virtualized resources on demand should meet growing or shrinking resource demand from clients, as well as guarantee fair share among users.

A major shortcoming in the above approaches is the lack of effective and efficient architectures that deliver the much-touted benefits of the cloud while overcoming its shortcomings.

1.2 Overview

This thesis proposes new architectures for cloud-computing that not only fix existing infrastructural deficiencies but also pave the way for newer and improved services. We present four new architectures.

SecureCloud [141]: Theft-free schedulers for multi-tenant architectures. Secure Cloud explores the security issues of cloud infrastructures. In cloud computing services, customers rent virtual machines (VMs) running on hardware owned and

managed by third-party providers, and pay for the amount of time their VM is running (in hours or months), rather than by the amount of CPU time used. Under this business model, the choice of scheduling algorithm of the hypervisor involves a trade-off between factors such as fairness, usage caps and scheduling latency. As in operating systems, a hypervisor scheduler may be vulnerable to behavior by virtual machines which results in inaccurate or unfair scheduling. However, cloud computing represents a new environment for such attacks for two reasons. First, the economic model of many services renders them vulnerable to theft-of-service attacks, which can be successful with far lower degrees of unfairness than required for strong denial-of-service attacks. In addition, the lack of detailed application knowledge in the hypervisor (e.g. to differentiate I/O wait from voluntary sleep) makes it more difficult to harden a hypervisor scheduler against malicious behavior. Therefore, verifying the scheduling fairness in cloud computing services is very important to both the providers and the customers. In Chapter 2, we study the scheduling algorithm used by Amazon EC2. We discover that the hypervisor scheduler in Amazon EC2 is vulnerable to behavior by virtual machines which results in inaccurate or unfair scheduling. We describe our attack scenario, provide a novel analysis of the necessary conditions for theft-of-service attacks, and propose scheduler modifications to eliminate the vulnerability. Moreover, our attack itself provides a mechanism for detecting the co-placement of VMs, which in conjunction with appropriate algorithms can be utilized to reveal this mapping. We present an algorithm that is provably optimal when the maximum partition size is bounded. In the unbounded case we show upper and lower bounds using the probabilistic method [18] in conjunction with a sieving technique.

DNSCloud [64]: DNS architecture protects itself from cache-poisoning attacks. DNSCloud explores the domain of new cloud architectures that allow for improved and robust Internet services. The Domain Name System or DNS [85] is a critical and integral part of the Internet. Kaminsky [52, 74, 100] showed that DNS caches throughout the world were susceptible to cache poisoning. This could lead to large-scale impersonation attacks by fake websites. Not only could it potentially

lead to the personal and financial ruin of individuals but successful poisoning of this system could also result in catastrophic consequences at a national level [13, 81] were any infrastructural networks such as the military network domains to be targeted by an enemy. Currently, there are only two ways to mitigate this vulnerability, patching the DNS server or patching the DNS protocol itself. Patching the DNS protocol, i.e. substituting DNSSEC [85] instead of DNS requires a forklift upgrade at both the client and resolving name-server ends and for now the world has settled for updating and patching the DNS server software (typically BIND [85]). The problem with the patch however, is that it requires upgrades to all client side resolving name servers, which is the more vulnerable and numerous end of things. In section 3, we propose DNSCloud, a new cloud-based architecture, to dramatically reduce the probability of suffering from a poisoned cache in the DNS.

WebCloud [142]: P2P content delivery system for (wide-area) social networks. Over the past few years, we have witnessed the beginnings of a shift in the patterns of content creation and exchange over the web. Previously, content was primarily created by a small set of entities and was delivered to a large audience of web clients. Now, individual Internet clients are creating content that makes up a significant fraction of Internet traffic [17, 47]. The net result is that, content today is generated by a large number of clients located at the edge of the network, is of more uniform popularity, and exhibits a workload that is governed by the social network. Unfortunately, existing architectures for content distribution are ill-suited for these new patterns of content creation and exchange. In section 4.6.1, we propose WebCloud - a content delivery system designed for social networks. The key insight is to leverage the spatial and temporal locality of interest between social network clients.

SamaritanCloud: Location-based cyber-physical network for obtaining real-world assistance. The tremendous rise in the popularity of social networks has been a defining aspect of the online experience over the last few years. Social network services provide a platform allowing for the building of social relations among people that share similar interests and activities. More recently the phenomenon of

geosocial networking has come to the fore. These services use the GPS capabilities of mobile devices to enable additional social dynamics. In geosocial network services such as Foursquare, Gowalla, Facebook Places, and Yelp [3–5, 9], clients share their current locations as well as recommendations for events, food or other interests. In section 5 we propose not just a new architecture, but, in fact, a new service – SamaritanCloud – the goal of which is to provide a way for people to connect with others (possibly strangers) in a remote location and obtain (physical) help from them. It is possible that such a service may never take off but we believe there are several instances in which such a service can be of use to people, e.g., please tell the marathon runner with bib #123 I will wait at the finish line; did I leave my textbook in the restroom? is there a brown puppy roaming in the playground? Such a service will require efficient technical solutions to problems such as scalability, privacy, reputation etc., to overcome the social barriers of soliciting help from strangers. We focus primarily on the technical aspects of scalability – efficiently finding candidates for a request by comparing client profiles with request criteria; efficiency – improving server throughput on high dimensional profiles with locality sensitive hashing; and privacy – matching people with strangers in a secure and private way so that the need for help and the ability, desire to assist are disclosed in a safe and controlled manner.

CHAPTER 2

SecureCloud

2.1 Introduction

Server virtualization [21] enables multiple instances of an operating system and applications (*virtual machines* or VMs) to run on the same physical hardware, as if each were on its own machine. Recently server virtualization has been used to provide so-called *cloud computing* services, in which customers rent virtual machines running on hardware owned and managed by third-party providers. Two such cloud computing services are Amazon's Elastic Compute Cloud (EC2) service and Microsoft Windows Azure Platform; in addition, similar services are offered by a number of web hosting providers (e.g. Rackspace's Rackspace Cloud and ServePath Dedicated Hosting's GoGrid) and referred to as Virtual Private Servers (VPS). In each of these services customers are charged by the amount of time their virtual machine is running (in hours or months), rather than by the amount of CPU time used.

The operation of a hypervisor is in many ways similar to that of an operating system; as an operating system manages access by processes to underlying resources, so too a hypervisor must manage access by multiple virtual machines to a single physical machine. In either case the choice of scheduling algorithm will involve a trade-off between factors such as fairness, usage caps and scheduling latency.

As in operating systems, a hypervisor scheduler may be vulnerable to behavior by virtual machines which results in inaccurate or unfair scheduling. Such anomalies and their potential for malicious use have been recognized in the past in operating systems—McCanne and Torek [88] demonstrate a denial-of-service attack on 4.4BSD, and more recently Tsafirir [130] presents a similar attack against Linux 2.6 which was fixed only recently. Such attacks typically rely on the use of periodic sampling or a low-precision clock to measure CPU usage; like a train passenger hiding whenever the conductor checks tickets, an attacking process ensures it is never scheduled when a scheduling tick occurs.

Cloud computing represents a new environment for such attacks, however, for two reasons. First, the economic model of many services renders them vulnerable to theft-of-service attacks, which can be successful with far lower degrees of unfairness than required for strong denial-of-service attacks. In addition, the lack of detailed application knowledge in the hypervisor (e.g. to differentiate I/O wait from voluntary sleep) makes it more difficult to harden a hypervisor scheduler against malicious behavior.

The scheduler used by the Xen hypervisor (and with modifications by Amazon EC2) is vulnerable to such timing-based manipulation—rather than receiving its fair share of CPU resources, a VM running on unmodified Xen using our attack can obtain up to 98% of total CPU cycles, regardless of the number of other VMs running on the same core. In addition we demonstrate a kernel module allowing unmodified applications to readily obtain 80% of the CPU. The Xen scheduler also supports a non-work-conserving (NWC) mode where each VM’s CPU usage is “capped”; in this mode our attack is able to evade its limits and use up to 85% of total CPU cycles. The modified EC2 scheduler uses this to differentiate levels of service; it protects other VMs from our attack, but we still evade utilization limits (typically 40%) and consume up to 85% of CPU cycles.

We give a novel analysis of the conditions which must be present for such attacks to succeed, and present four scheduling modifications which will prevent this attack without sacrificing efficiency, fairness, or I/O responsiveness. We have

implemented these algorithms, and present experimental results on Xen 3.2.1.

Chapter outline: The rest of chapter is organized as follows. Section 2.2 provides a brief introduction to VMM architectures, Xen VMM and Amazon EC2 as background. Section 2.3 discusses related work. Section 2.4 describes the details of the Xen Credit scheduler. Section 2.5 explains our attacking scheme and Section 2.7 presents experimental results in the lab as well as on Amazon EC2. Section 2.8 details our scheduling modifications to prevent this attack, and evaluates their performance and overhead. Section 2.9 proofs how our attacking scheme could coordinate attacks, and we conclude in Section 2.10.

2.2 Background

We first provide a brief overview of hardware virtualization technology, and of the Xen hypervisor and Amazon’s Elastic Compute Cloud (EC2) service in particular.

2.2.1 Hardware virtualization

Hardware virtualization refers to any system which interposes itself between an operating system and the hardware on which it executes, providing an emulated or *virtualized* view of physical resources. Almost all virtualization systems allow multiple operating system instances to execute simultaneously, each in its own *virtual machine* (VM). In these systems a Virtual Machine Monitor (VMM), also known as a *hypervisor*, is responsible for resource allocation and mediation of hardware access by the various VMs.

Modern hypervisors may be classified by the methods of executing guest OS code without hardware access: (a) binary emulation and translation, (b) para-virtualization, and (c) hardware virtualization support. Binary emulation executes privileged guest code in software, typically with just-in-time translation for speed [15]. Hardware virtualization support [2] in recent x86 CPUs supports a privilege level beyond supervisor mode, used by the hypervisor to control guest

OS execution. Finally, para-virtualization allows the guest OS to execute directly in user mode, but provides a set of *hypercalls*, like system calls in a conventional operating system, which the guest uses to perform privileged functions.

2.2.2 The Xen hypervisor

Xen is an open source VMM for x86/x64 [33]. It introduced para-virtualization on the x86, using it to support virtualization of modified guest operating systems without hardware support (unavailable at the time) or the overhead of binary translation. Above the hypervisor there are one or more virtual machines or *domains* which use hypervisor services to manipulate the virtual CPU and perform I/O.

2.2.3 Amazon Elastic Compute Cloud (EC2)

Amazon EC2 is a commercial service which allows customers to run their own virtual machine instances on Amazon’s servers, for a specified price per hour each VM is running. Details of the different instance types currently offered, as well as pricing per instance-hour, are shown in Table 2.1.

Table 2.1: Amazon EC2 Instance Types and Pricing. (Spring 2012. Speed is given in “Amazon EC2 Compute Units”.)

Instance Type	Memory	Cores \times speed	\$/Hr
Small	1.7GB	1×1	0.085
Large	7.5	2×2	0.34
X-Large	15	4×2	0.68
Hi-CPU Med.	1.7	2×2.5	0.17
Hi-CPU X-Large	7	8×2.5	0.68

Amazon states that EC2 is powered by “a highly customized version of Xen, taking advantage of virtualization” [11]. The operating systems supported are Linux, OpenSolaris, and Windows Server 2003; Linux instances (and likely OpenSolaris) use Xen’s para-virtualized mode, and it is suspected that Windows instances do so as well [25].

2.3 Related work

To the best of our knowledge, our work is the first to show how a deliberately misbehaving VM can unfairly monopolize CPU resources in a virtualized environment, with important implications for Cloud Computing environment. However, the concept of a timing attack long predates computers. Tsafirir *et al.* [130] designed a cheat attack based on a similar scenario in context of processes on the Linux 2.6 scheduler, allowing an attacking process to appear to consume no CPU and receive higher priority. McCanne and Torek [88] present the same cheat attack on 4.4BSD, and develop a uniform randomized sampling clock to estimate CPU utilization. They describe sufficient conditions for this estimate to be accurate, but unlike section 2.8.2 they do not examine conditions for a theft-of-service attack. Cherkasova and Gupta *et al.* [31, 32] have done an extensive performance analysis of scheduling in the Xen VMM. They studied I/O performance for the three schedulers: BVT, SEDF and Credit scheduler. Their work showed that both the CPU scheduling algorithm and the scheduler parameters drastically impact the I/O performance. Furthermore, they stressed that the I/O model on Xen remains an issue in resource allocation and accounting among VMs. Since Domain-0 is indirectly involved in servicing I/O for guest domains, I/O intensive domains may receive excess CPU resources by focusing on the processing resources used by Domain-0 on behalf of I/O bound domains. To tackle this problem, Gupta *et al.* [40] introduced the SEDF-DC scheduler, derived from Xen's SEDF scheduler, that charges guest domains for the time spent in Domain-0 on their behalf. Govindan *et al.* [57] proposed a CPU scheduling algorithm as an extension to Xen's SEDF scheduler that preferentially schedules I/O intensive domains. The key idea behind their algorithm is to count the number of packets flowing into or out of each domain and to schedule the one with highest count that has not yet consumed its entire slice. However, Ongaro *et al.* [101] pointed out that this scheme is problematic when bandwidth-intensive and latency-sensitive domains run concurrently on the same host - the bandwidth-intensive domains are likely to take priority over any latency-sensitive

domains with little I/O traffic. They explored the impact of VMM scheduler on I/O performance using multiple guest domains concurrently running different types of applications and evaluated 11 different scheduler configurations within Xen VMM with both the SEDF and Credit schedulers. They also proposed multiple methods to improve I/O performance. Weng *et al.* [134] found from their analysis that Xen's asynchronous CPU scheduling strategy wastes considerable physical CPU time. To fix this problem, they presented a hybrid scheduling framework that groups VMs into high-throughput type and concurrent type and determines processing resource allocation among VMs based on type. In a similar vein Kim *et al.* [77] presented a task-aware VM scheduling mechanism to improve the performance of I/O-bound tasks within domains. Their approach employs gray-box techniques to peer into VMs and identify I/O-bound tasks in mixed workloads. There are some configurations and policies [70, 95, 111] proposed to improve Xen security in other perspectives but not scheduling attacks. Very recently, Ristenpart *et al.* [109] instantiate new VMs on EC2 until one is placed co-resident with the target; they then show that known cross-VM side-channel attacks can extract information from the target. However the side-channel attacks were carried out on carefully controlled machines in the lab, not on EC2, and are likely to be significantly more difficult in practice [125], while, our exploit is directly applicable to EC2. We show that our exploit can be used to infer the mapping of VMs to physical hosts (Partition) and this information can then be used to amplify our scheduling attack into a coordinated siege by a collection of VMs. Although there is significant literature on similar problems for contention resolution in multiple-access channels [59, 60, 79], to the best of our knowledge Partition has not been studied earlier.

2.4 Xen Credit Scheduler analysis

In Xen (and other hypervisors) a single virtual machine consists of one or more virtual CPUs (VCPUs); the goal of the scheduler is to determine which VCPU to execute on each physical CPU (PCPU) at any instant. To do this it must determine which VCPUs are idle and which are active, and then from the active VCPUs choose one for each PCPU.

In a virtual machine, a VCPU is idle when there are no active processes running on it and the scheduler on that VCPU is running its *idle task*. On early systems the idle task would loop forever; on more modern ones it executes a HALT instruction, stopping the CPU in a lower-power state until an interrupt is received. On a fully-virtualized system this HALT traps to the hypervisor and indicates the VCPU is now idle; in a para-virtualized system a direct hypervisor call is used instead. When an exception (e.g. timer or I/O interrupt) arrives, that VCPU becomes active until HALT is invoked again.

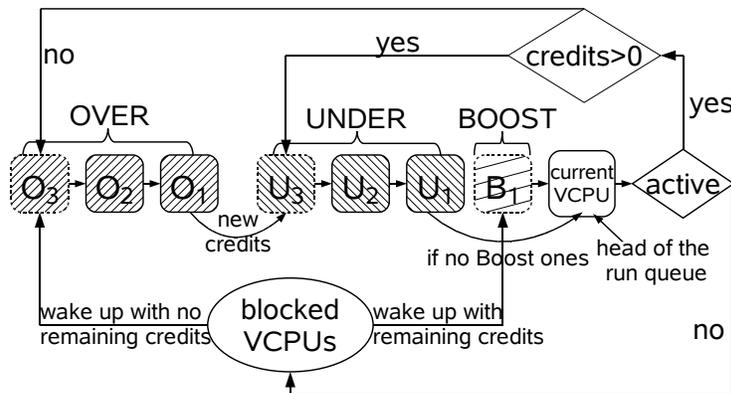
By default Xen uses the Credit scheduler [31], an implementation of the classic *token bucket* algorithm in which credits arrive at a constant rate, are conserved up to a maximum, and are expended during service. Each VCPU receives credits at an administratively determined rate, and a periodic scheduler tick debits credits from the currently running VCPU. If it has no more credits, the next VCPU with available credits is scheduled. Every 3 ticks the scheduler switches to the next runnable VCPU in round-robin fashion, and distributes new credits, capping the credit balance of each VCPU at 300 credits. Detailed parameters (assuming even weights) are:

Fast tick period:	10ms
Slower (rescheduling) tick:	30ms
Credits debited per fast tick:	100
Credit arrivals per fast tick:	100/N
Maximum credits:	300

where N is the number of VCPUs per PCPU. The fast tick decrements the running VCPU by 100 credits every 10 ms, giving each credit a value of $100 \mu\text{s}$ of CPU time; the cap of 300 credits corresponds to 30 ms, or a full scheduling quantum. Based on their credit balance, VCPUs are divided into three states: *UNDER*, with a positive credit balance, *OVER*, or out of credits, and *BLOCKED* or halted.

The VCPUs on a PCPU are kept in an ordered list, with those in *UNDER* state ahead of those in *OVER* state; the VCPU at the head of the queue is selected for execution. In work conserving mode, when no VCPUs are in the *UNDER* state, one in the *OVER* state will be chosen, allowing it to receive more than its share of CPU. In non-work-conserving (NWC) mode, the PCPU will go idle instead.

Figure 2.1: Per-PCPU Run Queue Structure



The executing VCPU leaves the run queue head in one of two ways: by going idle, or when removed by the scheduler while it is still active. VCPUs which go idle enter the *BLOCKED* state and are removed from the queue. Active VCPUs are enqueued after all other VCPUs of the same state—*OVER* or *UNDER*—as shown in Figure 2.1.

The basic credit scheduler accurately distributes resources between CPU-intensive workloads, ensuring that a VCPU receiving k credits per 30 ms epoch will receive at least $k/10$ ms of CPU time within a period of $30N$ ms. This fairness comes at the expense of I/O performance, however, as events such as packet reception may wait as long as $30N$ ms for their VCPU to be scheduled.

To achieve better I/O latency, the Xen Credit scheduler attempts to prioritize such I/O. When a VCPU sleeps waiting for I/O it will typically have remaining credits; when it wakes with remaining credits it enters the BOOST state and may immediately preempt running or waiting VCPUs with lower priorities. If it goes idle again with remaining credits, it will wake again in BOOST priority at the next I/O event.

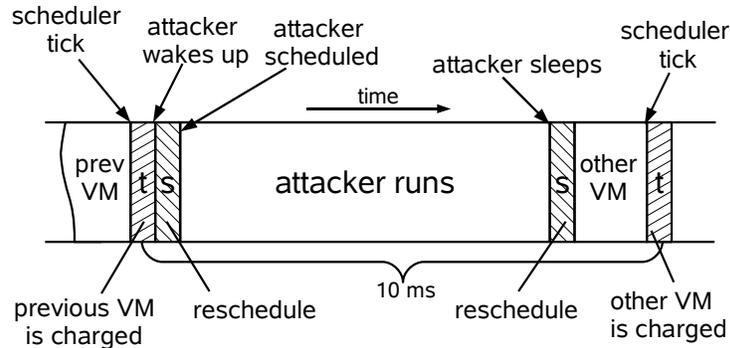
This allows I/O-intensive workloads to achieve very low latency, consuming little CPU and rarely running out of credits, while preserving fair CPU distribution among CPU-bound workloads, which typically utilize all their credits before being preempted. However, as we describe in the following section, it also allows a VM to “steal” more than its fair share of CPU time.

2.5 Anatomy of an attack

Although the Credit scheduler provides fairness and low I/O latency for well-behaved virtual machines, poorly-behaved ones can evade its fairness guarantees. In this section we describe the features of the scheduler which render it vulnerable to attack, formulate an attack scheme, and present results showing successful theft of service both in the lab and in the field on EC2 instances. Our attack relies on periodic sampling as used by the Xen scheduler, and is shown as a timeline in Figure 2.2. Every 10 ms the scheduler tick fires and schedules the attacking VM, which runs for $10 - \epsilon$ ms and then calls *Halt()* to briefly go idle, ensuring that another VM will be running at the next scheduler tick. Our attack is self-synchronizing due to wake-up after a scheduling tick. In theory the efficiency of this attack increases as ϵ approaches 0; however in practice some amount of timing jitter is found, and overly small values of ϵ increase the risk of the VM being found executing when the scheduling tick arrives.

When perfectly executed on the non-BOOST credit scheduler, this ensures that the attacking VM will never have its credit balance debited. If there are N VMs with equal shares, then the $N-1$ victim VMs will receive credits at a total rate of

Figure 2.2: Attack Timing



$\frac{N-1}{N}$, and will be debited at a total rate of 1.

This vulnerability is due not only to the predictability of the sampling, but to the granularity of the measurement. If the time at which each VM began and finished service were recorded with a clock with the same 10 ms resolution, the attack would still succeed, as the attacker would have a calculated execution time of 0 on transition to the next VM.

This attack is more effective against the actual Xen scheduler because of its BOOST priority mechanism. When the attacking VM yields the CPU, it goes idle and waits for the next timer interrupt. Due to a lack of information at the VM boundary, however, the hypervisor is unable to distinguish between a VM waking after a deliberate sleep period—a non-latency-sensitive event—and one waking for e.g. packet reception. The attacker thus wakes in BOOST priority and is able to preempt the currently running VM, so that it can execute for $10 - \epsilon$ ms out of every 10 ms scheduler cycle.

2.6 Implementation

2.6.1 User-level

To examine the performance of our attack scenario in practice, we implement it using both user-level and kernel-based code and evaluate them in the lab and on Amazon EC2. In each case we test with two applications: a simple loop we refer to as “Cycle Counter” described below, and the Dhrystone 2.1 [133] CPU benchmark. Our attack described in Section 2.5 requires millisecond-level timing in order to sleep before the debit tick and then wake again at the tick; it performs best either with a tick-less Linux kernel [117] or with the kernel timer frequency set to 1000 Hz.

```
prev=rdtsc()  
loop:  
    if (rdtsc() - prev) > 9ms  
        prev = rdtsc()  
        usleep(0.5ms)
```

2.6.2 Kernel-level

To implement kernel-level attack, the most basic way to do is that to add files to the kernel source tree and modify kernel configuration to include the new files in compilation. However, most Unix kernels are monolithic. The kernel itself is a piece of compact code, in which all functions share a common space and are tightly related. When the kernel needs to be updated, all the functions must be relinked and reinstalled and the system rebooted before the change can take affect. This makes modifying kernel level code difficult, especially when the change is made to a kernel which is not in user’s control. Loadable kernel module (LKM) allows developers to make change to kernel when it is running. LKM is an object file that contains code to extend the running kernel of an operating system. LKMs are typically used for one of the three functionalities : device drivers, filesystem drivers

and system calls. Some of the advantages of using loadable kernel modules rather than modifying the base kernel are: (1) No kernel recompilation is required. (2) New kernel functionality can be added without root privileges. (3) New functionality takes effect right after module installation, no reboot is required. (4) LKM can be reloaded at run time, it does not add to the size of kernel permanently.

The attack is implemented as a kernel thread which invokes an OS sleep for $10 - \varepsilon$ ms, allowing user applications to run, and then invokes the `SCHED_block` hypercall via the safe halt function. In practice ε must be higher than for the user-mode attack, due to timing granularity and jitter in the kernel.

```
loop:
msleep(8);
safe_halt();
```

In theory the less ε is in kernel module implementation, the more CPU cycles the user application could consume. However due to our evaluation of `msleep`, there is a 1ms timing jitter in kernel. `msleep(8)` allows the kernel thread to wake up on time and temporarily pause all user applications. Then the VM is considered as idle and gets swapped out before debit tick happens. Since `msleep(9)` does not guarantee the thread wake up before the debit tick every time, thus the attacking VM may not be swapped in/swapped out as expected. According to this observation, $\varepsilon = 2$ is a safe choice in implementation for kernel 2.6.

2.7 Evaluation

2.7.1 User-level

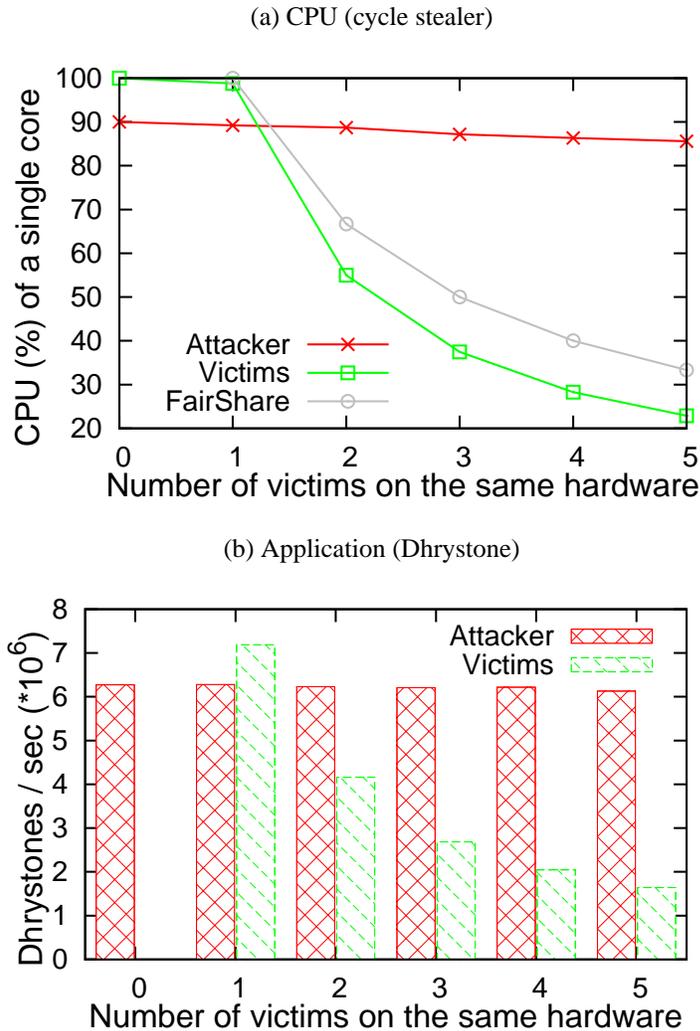
Experiments in the lab

Our first experiments evaluate our attack against unmodified Xen in the lab in work-conserving mode, verifying the ability of the attack to both deny CPU resources to

competing “victim” VMs, and to effectively use the “stolen” CPU time for computation. All experiments were performed on Xen 3.2.1, on a 2-core 2.7 GHz Intel Core2 CPU. Virtual machines were 32-bit, para-virtualized, single-VCPU instances with 192 MB memory, each running Suse 11.0 Core kernel 2.6.25 with a 1000 Hz kernel timer frequency. To test our ability to steal CPU resources from other VMs, we implement a “Cycle Counter”, which performs no useful work, but rather spins using the RDTSC instruction to read the timestamp register and track the time during which the VM is scheduled. The attack is performed by a variant of this code, “Cycle Stealer”, which tracks execution time and sleeps once it has been scheduled for $10 - \varepsilon$ (here $\varepsilon = 1$ ms). Note that the sleep time is slightly less than ε , as the process will be woken at the next OS tick after timer expiration and we wish to avoid over-sleeping. In Figure 2.3a we see attacker and victim performance on our 2-core test system. As the number of victims increases, attacker performance remains almost constant at roughly 90% of a single core, while the victims share the remaining core.

To measure the ability of our attack to effectively use stolen CPU cycles, we embed the attack within the Dhrystone benchmark. By comparing the time required for the attacker and an unmodified VM to complete the same number of Dhrystone iterations, we can determine the *net* amount of work stolen by the attacker. Our baseline measurement was made with one VM running unmodified Dhrystone, with no competing usage of the system; it completed 1.5×10^9 iterations in 208.8 seconds. When running 6 unmodified instances, three for each core, each completed the same 1.5×10^9 iterations in 640.8 seconds on average—32.6% the baseline speed, or close to the expected fair share performance of 33.3%. With one modified attacker instance competing against 5 unmodified victims, the attacker completed in 245.3 seconds, running at a speed of 85.3% of baseline, rather than 33.3%, with a corresponding decrease in victim performance. Full results for experiments with 0 to 5 unmodified victims and the modified Dhrystone attacker are shown in Figure 2.3b. In the modified Dhrystone attacker the TSC register is sampled once for each iteration of a particular loop, as described in the appendix; if this

Figure 2.3: Lab experiments (User Level) - CPU and application performance for attacker and victims.



sampling occurs too slowly the resulting timing inaccuracy might affect results. To determine whether this might occur, lengths of the compute and sleep phases of the attack were measured. Almost all (98.8%) of the compute intervals were found to lie within the bounds 9 ± 0.037 ms, indicating that the Dhrystone attack was able to attain timing precision comparable to that of Cycle Stealer. As described in Section 2.5, the attacker runs for a period of length $10 - \varepsilon$ ms and then briefly goes to sleep to avoid the sampling tick. A smaller value of ε increases the CPU time stolen by the attacker; however, too small an ε increases the chance of being charged due to

timing jitter. To examine this trade-off we tested values of $10 - \epsilon$ between 7 and 9.9 ms. Figure 2.7 shows that under lab conditions the peak value was 98% with an execution time of 9.8 ms and a requested sleep time of 0.1 ms. When execution time exceeded 9.8 ms the attacker was seen by sampling interrupts with high probability. In this case it received only about 21% of one core, or even less than the fair share of 33.3%.

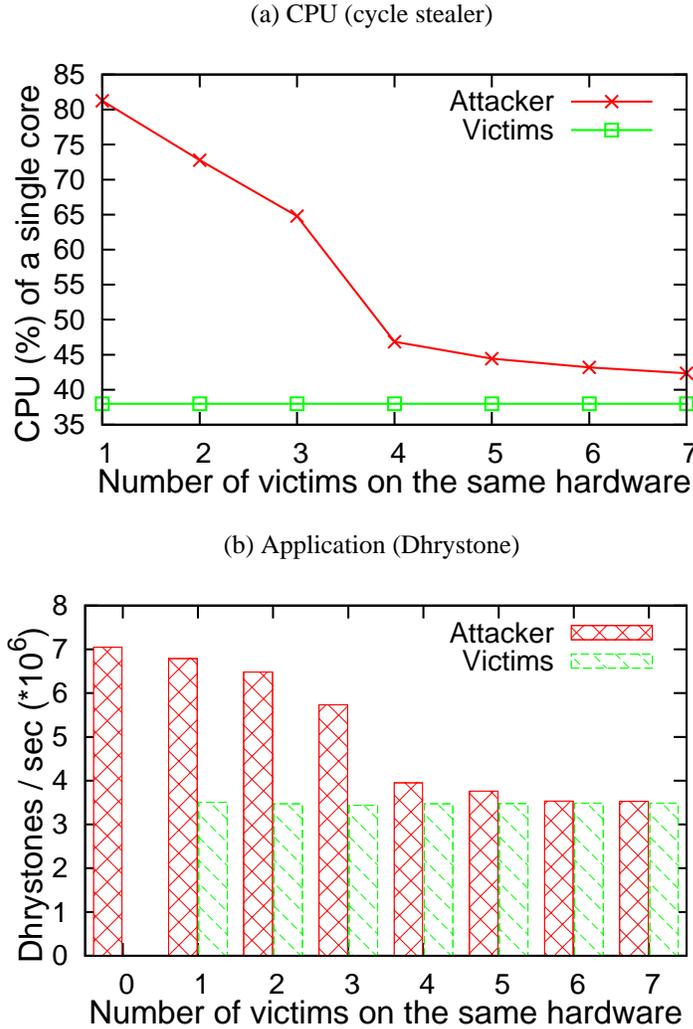
Experiments on Amazon

We evaluate our attacking using Amazon EC2 Small instances with the following attributes: 32-bit, 1.7 GB memory, 1 VCPU, running Amazon’s Fedora Core 8 kernel 2.6.18, with a 1000 Hz kernel timer. We note that the VCPU provided to the Small instance is described as having “1 EC2 Compute Unit”, while the VCPUs for larger and more expensive instances are described as having 2 or 2.5 compute units; this indicates that the scheduler is being used in non-work-conserving mode to throttle Small instances. To verify this hypothesis, we ran Cycle Stealer in measurement (i.e. non-attacking) mode on multiple Small instances, verifying that these instances are capped to less than $\frac{1}{2}$ of a single CPU core—in particular, approximately 38% on the measured systems. We believe that the nominal CPU cap for 1-unit instances on the measured hardware is 40%, corresponding to an unthrottled capacity of 2.5 units. Additional experiments were performed on a set of 8 Small instances co-located on a single 4-core 2.6 GHz physical system provided by our partners at Amazon.¹ The Cycle Stealer and Dhrystone attacks measured in the lab were performed in this configuration, and results are shown in Figure 2.4a and Figure 2.4b, respectively. We find that our attack is able to evade the CPU cap of 40% imposed by EC2 on Small instances, obtaining up to 85% of one core in the absence of competing VMs. When co-located CPU-hungry “victim” VMs were present, however, EC2 performance diverged from that of unmodified Xen. As seen in Figures 2.4a and 2.4b, co-located VM performance was virtually unaffected by

¹ This configuration allowed direct measurement of attack impact on co-located “victim” VMs, as well as eliminating the possibility of degrading performance of other EC2 customers.

our attack. Although this attack was able to steal cycles from EC2, it was unable to steal cycles from other EC2 customers.

Figure 2.4: Amazon EC2 experiments - CPU and application performance for attacker and victims.



2.7.2 Kernel-level

Lab experiments were performed with one attacking VM, which loads the kernel module, and up to 5 victims running a simple CPU-bound loop. In this case the fair CPU share for each guest instance on our 2-core test system would be $\frac{200}{N}\%$ of a single core, where N is the total number of VMs. Due to the granularity of

kernel timekeeping, requiring a larger ε , the efficiency of the kernel-mode attack is slightly lower than that of the user-mode attack. We evaluate the kernel level attacking with Cycle Stealer in measurement (non-attacking) mode and unmodified Dhrystone. In our tests, the attacker consumed up to 80.0% of a single core, with the remaining CPU time shared among the victim instances; results are shown in Figure 2.5a. The average amount of CPU stolen by the attacker decreases slightly (from 80.0% to 78.2%) as the number of victims increases; we speculate that this may be due to increased timing jitter causing the attacker to occasionally be charged by the sampling tick. The current implementation of the kernel module does not succeed in stealing cycles on Amazon EC2. We dig into Amazon EC2 kernel synchronization and our analysis of timing traces indicates a lack of synchronization of the attacker to the hypervisor tick, as seen for NWC mode in the lab, below; in this case, however, synchronization was never achieved. The user-level attack displays strong self-synchronizing behavior, aligning almost immediately to the hypervisor tick; we are investigating approaches to similarly strengthen self-synchronizing in the kernel module.

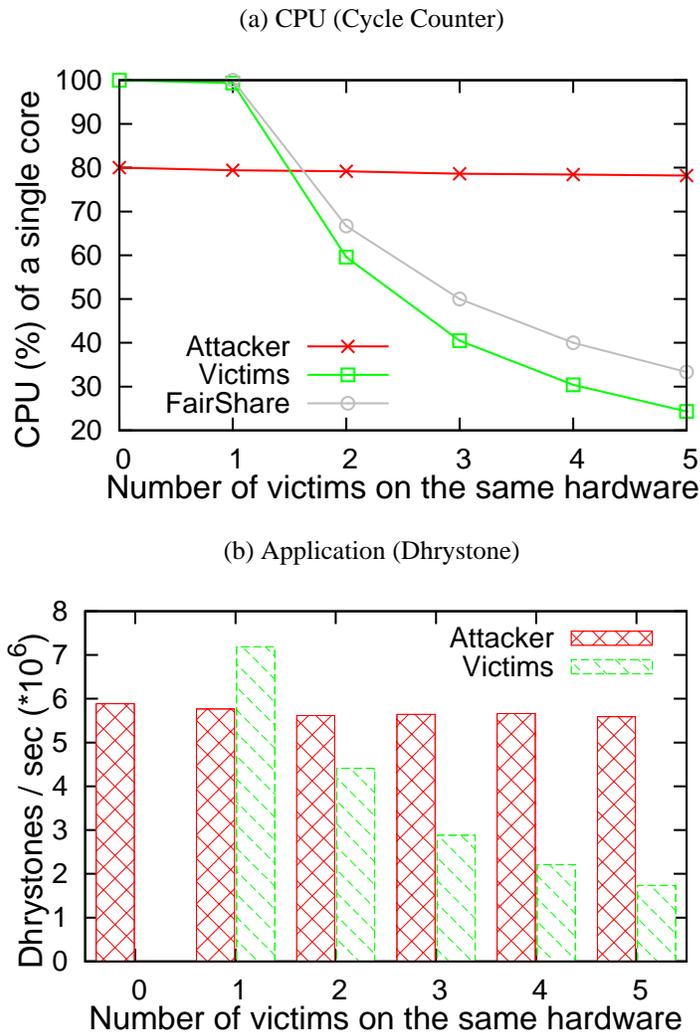
Table 2.2: Lab: Attack performance in non-work-conserving mode with 33.3% limit.

	Cycle Stealer (% of 1 core achieved)	Dhrystones per second	% of baseline
attacker	81.0	5749039	80.0
victims	23.4	1658553	23.1

2.7.3 Non-work conserving mode

As with most hypervisors, Xen CPU scheduling can be configured in two modes: work-conserving mode and non-work-conserving mode [42]. In order to optimize scheduling and allow near native performance, work-conserving scheduling schemes are efficient and do not waste any processing cycles. As long as there are instructions to be executed and there is enough CPU capacity, work-conserving schemes assign instructions to physical CPU to be carried out. Otherwise the in-

Figure 2.5: Lab experiments (Kernel Level) - CPU and application performance for attacker and victims.



structions will be queued and will be executed based on their priority. In contrast, non-work-conserving schemes allow CPU resources to go unused. In such schemes, there is no advantage to execute an instruction sooner. Usually, the CPU resources are allocated to VMs in proportion to their weights, e.g. two VMs with equal weight will own 50% each of CPU. When one VM goes idle, instead of obtaining the free cycles, the other VM is capped to its 50% sharing. Our attack program helps a VM to evade its capacity cap and obtain more, no matter if other VMs are busy or idle. Additional experiments were performed to examine our attack per-

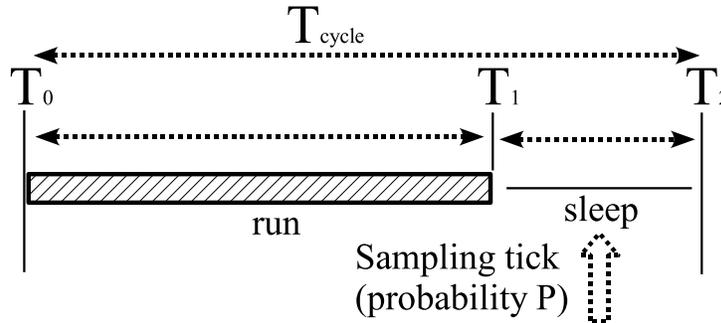
formance against the Xen scheduler in non-work-conserving mode. One attacker and 5 victims were run on our 2-core test system, with a CPU cap of 33% set for each VM; results are presented in Table 2.2 for both Cycle Stealer and Dhrystone attackers, including user-level and kernel-level implementation. With a per-VM CPU cap of 33%, the attacker was able to obtain 80% of a single core, as well. In each case the primary limitation on attack efficiency is the granularity of kernel timekeeping; we speculate that by more directly manipulating the hypervisor timer it would be possible to increase efficiency. In non-work-conserving case, when a virtual machine running attacking program yields CPU, there is a chance (based on our experiment results, not often though) that reschedule does not happen, the attacking VM is still considered as the scheduled VM when debit tick happens and gets debited. In other words, non-work-conserving mode does not stop our attack, but adds some interferences. In addition we note that it often appeared to take seconds or longer for the attacker to synchronize with the hypervisor tick and evade resource caps, while the user-level attack succeeds immediately.

2.8 Theft-resistant Schedulers

The class of theft-of-service attacks on schedulers which we describe is based on a process or virtual machine voluntarily sleeping when it could have otherwise remained scheduled. As seen in Figure 2.6, this involves a tradeoff—the attack will only succeed if the expected benefit of sleeping for T_{sleep} is greater than the guaranteed cost of yielding the CPU for that time period. If the scheduler is attempting to provide each user with its fair share based on measured usage, then sleeping for a duration t must reduce measured usage by more than t in order to be effective. Conversely, a scheduler which ensures that yielding the CPU will never reduce measured usage more than the sleep period itself will be resistant to such attacks.

This is a broader condition than that of maintaining an unbiased estimate of CPU usage, which is examined by McCanne and Torek [88]. Some theft-resistant schedulers, for instance, may over-estimate the CPU usage of attackers and give

Figure 2.6: Attacking trade-offs. The benefit of avoiding sampling with probability P must outweigh the cost of forgoing T_{sleep} CPU cycles.



them less than their fair share. In addition, for schedulers which do not meet our criteria, if we can bound the ratio of sleep time to measurement error, then we can establish bounds on the effectiveness of a timing-based theft-of-service attack.

2.8.1 Exact scheduler

The most direct solution is the *Exact scheduler*: using a high-precision clock (in particular, the TSC) to measure actual CPU usage when a scheduler tick occurs or when a VCPU yields the CPU and goes idle, thus ensuring that an attacking VM is always charged for exactly the CPU time it has consumed. In particular, this involves adding logic to the Xen scheduler to record a high-precision timestamp when a VM begins executing, and then calculate the duration of execution when it yields the CPU. This is similar to the approach taken in e.g. the recent tickless Linux kernel [117], where timing is handled by a variable interval timer set to fire when the next event is due rather than using a fixed-period timer tick.²

2.8.2 Randomized schedulers

An alternative to precise measurement is to sample as before, but on a random schedule. If this schedule is uncorrelated with the timing of an attacker, then over

²Although Kim et al. [77] use TSC-based timing measurements in their modifications to the Xen scheduler, they do not address theft-of-service vulnerabilities.

sufficiently long time periods we will be able to estimate the attacker's CPU usage accurately, and thus prevent attack. Assuming a fixed charge per sample, and an attack pattern with period T_{cycle} , the probability P of the sampling timer falling during the sleep period must be no greater than the fraction of the cycle $\frac{T_{sleep}}{T_{cycle}}$ which it represents.

Poisson Scheduler: This leads to a Poisson arrival process for sampling, where the expected number of samples during an interval is exactly proportional to its duration, regardless of prior history. This leads to an exponential arrival time distribution,

$$\Delta T = \frac{-\ln U}{\lambda}$$

where U is uniform on $(0,1)$ and λ is the rate parameter of the distribution. We approximate such Poisson arrivals by choosing the inter-arrival time according to a truncated exponential distribution, with a maximum of 30 ms and a mean of 10 ms, allowing us to retain the existing credit scheduler structure. Due to the possibility of multiple sampling points within a 10 ms period we use a separate interrupt for sampling, rather than re-using or modifying the existing Xen 10 ms interrupt.

Bernoulli Scheduler: The discrete-time analog of the Poisson process, the *Bernoulli* process, may be used as an approximation of Poisson sampling. Here we divide time into discrete intervals, sampling at any interval with probability p and skipping it with probability $q = 1-p$. We have implemented a Bernoulli scheduler with a time interval of 1 ms, sampling with $p = \frac{1}{10}$, or one sample per 10 ms, for consistency with the unmodified Xen Credit scheduler. Rather than generate a timer interrupt with its associated overhead every 1 ms, we use the same implementation strategy as for the Poisson scheduler, generating an inter-arrival time variate and then setting an interrupt to fire after that interval expires. By quantizing time at a 1 ms granularity, our Bernoulli scheduler leaves a small vulnerability, as an attacker may avoid being charged during any 1 ms interval by sleeping before the end of the interval. Assuming that (as in Xen) it will not resume until the beginning of the next 10 ms period, this limits an attacker to gaining no more than 1 ms every 10 ms above its fair share, a relatively insignificant theft of service.

Uniform Scheduler: The final randomized scheduler we propose is the *Uniform* scheduler, which distributes its sampling uniformly across 10 ms scheduling intervals. Rather than generating additional interrupts, or modifying the time at which the existing scheduler interrupt fires, we perform sampling within the virtual machine switch code as we did for the exact scheduler. In particular, at the beginning of each 10 ms interval (time t_0) we generate a random offset Δ uniformly distributed between 0 and 10 ms. At each VCPU switch, as well as at the 10 ms tick, we check to see whether the current time has exceeded $t_0 + \Delta$. If so, then we debit the currently running VCPU, as it was executing when the “virtual interrupt” fired at $t_0 + \Delta$. Although in this case the sampling distribution is not memoryless, it is still sufficient to thwart our attacker. We assume that sampling is undetectable by the attacker, as it causes only a brief interruption indistinguishable from other asynchronous events such as network interrupts. In this case, as with Poisson arrivals the expected number of samples within any interval in a 10 ms period is exactly proportional to the duration of the interval. Our implementation of the uniform scheduler quantizes Δ with 1 ms granularity, leaving a small vulnerability as described in the case of the Bernoulli scheduler. As in that case, however, the vulnerability is small enough that it may be ignored. We note also that this scheduler is not theft-proof if the attacker is able to observe the sampling process. If we reach the 5 ms point without being sampled, for instance, the probability of being charged 10 ms in the remaining 5 ms is 1, while avoiding that charge would only cost 5 ms.

2.8.3 Evaluation

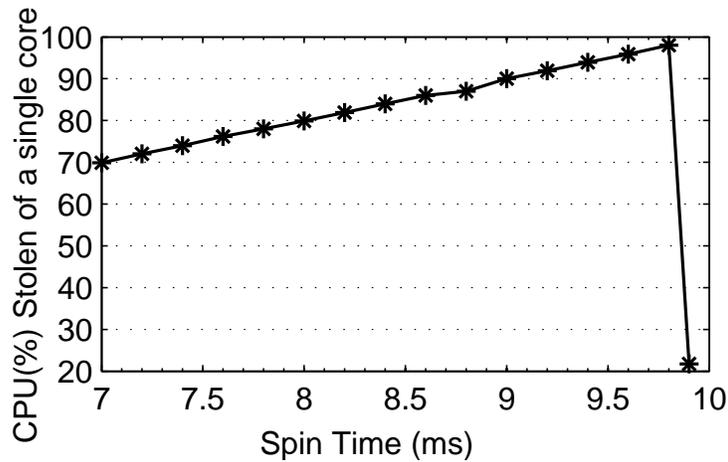
We have implemented each of the four modified schedulers on Xen 3.2.1. Since the basic credit and priority boosting mechanisms have not been modified from the original scheduler, our modified schedulers should retain the same fairness and I/O performance properties of the original in the face of well-behaved applications. To verify performance in the face of ill-behaved applications we tested attack performance against the new schedulers; in addition measuring overhead and I/O perfor-

mance.

Performance against attack

In Table 2.3 we see the performance of our Cycle Stealer on the Xen Credit scheduler and the modified schedulers. All four of the schedulers were successful in thwarting the attack: when co-located with 5 victim VMs on 2 cores on the unmodified scheduler, the attacker was able to consume 85.6% of a single-CPU with user-level attacking and 80% with kernel-level attacking, but no more than its fair share on each of the modified ones. (Note that the 85.6% consumption with user-level attacking in the unmodified case was limited by the choice of $\varepsilon = 1$ ms, and can increase with suitably reduced values of ε as shown in Figure 2.7.)

Figure 2.7: Lab: User-level attack performance vs. execute times. (note - sleep time ≤ 10 - spin time)



In Table 2.4 we see similar results for the modified Dhrystone attacker. Compared to the baseline, the unmodified scheduler allows the attacker to steal about 85.3% CPU cycles with user-level attacking and 77.8% with kernel-level attacking; while each of the improved schedulers limits the attacker to approximately its fair share.

Table 2.3: Performance of the schedulers against cycle stealer

Scheduler	CPU(%) obtained by the attacker	
	(user-level)	(kernel-level)
Xen Credit	85.6	78.8
Exact	32.9	33.1
Uniform	33.1	33.3
Poisson	33.0	33.2
Bernoulli	33.1	33.2

Table 2.4: Performance of the schedulers against Dhrystone

Scheduler	CPU(%) obtained by the attacker	
	(user-level)	(kernel-level)
Xen Credit	85.3	77.8
Exact	32.2	33.0
Uniform	33.0	33.4
Poisson	32.4	33.1
Bernoulli	32.5	33.3

Overhead measurement

To quantify the impact of our scheduler modifications on normal execution (i.e. in the absence of attacks) we performed a series of measurements to determine whether application or I/O performance had been degraded by our changes. Since the primary modifications made were to interrupt-driven accounting logic in the Xen scheduler, we examined overhead by measuring performance of a CPU-bound application (unmodified Dhrystone) on Xen while using the different scheduler. To reduce variance between measurements (e.g. due to differing cache line alignment [97]) all schedulers were compiled into the same binary image, and the desired scheduler selected via a global configuration variable set at boot or compile time.

Our modifications added overhead in the form of additional interrupts and/or accounting code to the scheduler, but also eliminated other accounting code which had performed equivalent functions. To isolate the effect of new code from that of the removal of existing code, we also measured versions of the Poisson and

Bernoulli schedulers (Poisson-2 and Bernoulli-2 below) which performed all accounting calculations of both schedulers, discarding the output of the original scheduler calculations.

Results from 100 application runs for each scheduler are shown in Table 2.5. Overhead of our modified schedulers is seen to be low—well under 1%—and in the case of the Bernoulli and Poisson schedulers is negligible. Performance of the Poisson and Bernoulli schedulers was unexpected, as each incurs an additional 100 interrupts per second; the overhead of these interrupts appears to be comparable to or less than the accounting code which we were able to remove in each case. We note that these experiments were performed with Xen running in para-virtualized mode; the relative cost of accounting code and interrupts may be different when using hardware virtualization.

We analyzed the new schedulers’ I/O performances by testing the I/O latency between two VMs in two configurations. In configuration 1, two VMs executed on the same core with no other VMs active, while in configuration 2 a CPU-bound VM was added on the other core. From the first test, we expected to see the performance of well-behaved I/O intensive applications on different schedulers; from the second one, we expected to see that the new schedulers retain the priority boosting mechanism.

In Table 2.6 we see the results of these measurements. Differences in performance were minor, and as may be seen by the overlapping confidence intervals, were not statistically significant.

Table 2.5: Scheduler CPU overhead, 100 data points per scheduler.

Scheduler	CPU overhead (%)	95% CI
Exact	0.50	0.24 – 0.76
Uniform	0.44	0.27 – 0.61
Poisson	0.04	-0.17 – 0.24
Bernoulli	-0.10	-0.34 – 0.15
Poisson-2	0.79	0.60 – 0.98
Bernoulli-2	0.79	0.58 – 1.00

Table 2.6: I/O latency by scheduler, with 95% confidence intervals.

Scheduler	Round-trip delay (μs)	
	(config. 1)	(config. 2)
Unmodified Xen Credit	53 ± 0.66	96 ± 1.92
Exact	55 ± 0.61	97 ± 1.53
Uniform	54 ± 0.66	96 ± 1.40
Poisson	53 ± 0.66	96 ± 1.40
Bernoulli	54 ± 0.75	97 ± 1.49

2.8.4 Additional discussion

A comprehensive comparison of our proposed schedulers is shown in Table 2.7. The *Poisson* scheduler seems to be the best option in practice, as it has no performance overhead nor vulnerability. Even though it has short-period variance, it guarantees exactly fair share in the long run. The *Bernoulli* scheduler would be an alternative if the vulnerability of up to 1ms is not a concern. The *Uniform* scheduler has similar performance to the Bernoulli one, and the implementation of sampling is simpler, but it has more overhead than Poisson and Bernoulli. Lastly, the *Exact* scheduler is the most straight-forward strategy to prevent cycle stealing, with a relatively trivial implementation but somewhat higher overhead.

Table 2.7: Comparison of the new schedulers

Schedulers	Short-run fairness	Long-run fairness	Low overhead	Ease of implementation	Deterministic	Theft-proof
Exact	✓	✓		✓	✓	✓
Uniform		✓		✓		
Poisson		✓	✓			✓
Bernoulli		✓	✓			

2.9 Coordinated cloud attacks

We have shown how an old vulnerability has manifested itself in the modern context of virtualization. Our attack [141] enables an attacking VM to steal cycles by gaming a vulnerability in the hypervisor’s scheduler. We now explain how, given a

collection of VMs in a cloud, attacks may be coordinated so as to steal the maximal number of cycles.

Consider a customer of a cloud service provider with a collection of VMs. Their goal is to extract the maximal number of cycles possible for executing their computational requirements. Note that the goal of the customer is not to inflict the highest load on the cloud but to extract the maximum amount of useful work for themselves. To this end they wish to steal as many cycles as possible. As has already been explained attacking induces an overhead and having multiple VMs in attack mode on the same physical host leads to higher total overhead reducing the total amount of useful work. Ideally, therefore, the customer would like to have exactly one VM in attack mode per physical host with the other VMs functioning normally in non-attack mode.

Unfortunately, cloud providers such as Amazon's EC2 not only control the mapping of VMs to physical hosts but also withhold information about the mapping from their customers. As infrastructure providers this is the right thing for them to do. By doing so they make it harder for malicious customers to snoop on others. [109] show the possibility of targeting victim VMs by mapping the internal cloud infrastructure though this is much harder to implement successfully in the wild [125]. Interestingly, though our attack can be turned on its head not just to steal cycles but also to identify co-placement. Assume without loss of generality for the purposes of the rest of this discussion that hosts are single core machines. Recall that a regular (i.e. non-attacking) VM on a single host is capped at 38% whereas an attacking VM obtains 87% of the cycles. Now, if we have two attacking VMs on the same host then they obtain about 42% each. Thus by exclusively (i.e without running any other VMs) running a pair of VMs in attack mode one can determine whether they are on the same physical host or not (depending on whether they get 42% or 87%). Thus a customer could potentially run every pair of VMs and determine which VMs are on the same physical host. (Of course, this is under the not-unreasonable assumption that only the VMs of this particular customer can be in attack mode.) Note that if there are n VMs then this scheme requires $\binom{n}{2}$ tests

for each pair. This leads to a natural question: what is the most efficient way to discover the mapping of VMs to physical host? This problem has a very clean and beautiful formulation.

Observe that the VMs get partitioned among (an unknown number of) the physical hosts. And we have the flexibility to activate some subset of the VMs in attack mode. We assume (as is the case of Amazon's EC2) that there is no advantage to running any of the other VMs in normal (non-attack) mode since they get their fair share (recall that in EC2 attackers only get additional *spare* cycles). When we activate a subset of the VMs in attack mode then we get back one bit of information for each VM in the subset - namely, whether that VM is the only VM from the subset on its physical host or whether there are 2 or more VMs from the subset on the same host. Thus the question now becomes: what is the fastest way to discover the unknown partition (of VMs among physical hosts)? We think of each subset that we activate as a query that takes unit time and we wish to use the fewest number of queries. More formally:

Partition. You are given an unknown partition $\mathfrak{P} = \{S_1, S_2, \dots, S_k\}$ of a ground set $[1 \dots n]$. You are allowed to ask queries of this partition. Each query is a subset $Q = \{q_1, q_2, \dots\} \subset [1 \dots n]$. When you ask the query Q you get back $|Q|$ bits, a bit b_i for each element $q_i \in Q$; let S_{q_i} denote the set of the partition \mathfrak{P} containing q_i ; then $b_i = 1$ if $|Q \cap S_{q_i}| = 1$ (and otherwise, $b_i = 0$ if $|Q \cap S_{q_i}| \geq 2$). The goal is to determine the query complexity of Partition, i.e., you have to find \mathfrak{P} using the fewest queries.

Recall that a partition is a collection of disjoint subsets whose union is the ground set, i.e., $\forall 1 \leq i, j \leq k, S_i \cap S_j = \emptyset$ and $\bigcup_{i=1}^k S_i = [1 \dots n]$. Observe that one can define both adaptive (where future queries are dependent on past answers) and oblivious variants of Partition. Obviously, adaptive queries have at least as much power as oblivious queries. In the general case we are able to show nearly tight bounds:

Theorem 1. *The oblivious query complexity of Partition is $O(n^{\frac{3}{2}}(\log n)^{\frac{1}{4}})$ and the adaptive query complexity of Partition is $\Omega(\frac{n}{\log^2 n})$.*

Proof Sketch: The upper bound involves the use of the probabilistic method in conjunction with sieving [18]. We are able to derandomize the upper bound to produce deterministic queries, employing expander graphs and pessimal estimators. The lower bound uses an adversarial argument based on the probabilistic method as well. \square

In practice, partitions cannot be arbitrary as there is a bound on the number of VMs that a cloud service provider will map to a single host. This leads to the problem B-Partition where all the sets in the partition have a size at most B . In the special case we are able to prove tight bounds:

Theorem 2. *The query complexity of B-Partition is $\theta(\log n)$.*

Proof Sketch: The lower bound follows directly from the information-theoretic argument that there are $\Omega(2^{n \log n})$ partitions while each query returns only n bits of information. The upper bound involves use of the probabilistic method (though the argument is simpler than for the general case) and can be derandomized to provide deterministic constructions of the queries. \square

Proof Sketch: We prove an upper bound of $O(\log n)$ on the query complexity of B-Partition, where the size of any set in the partition is at most B . Due to constraints of space we exhibit a randomized construction and leave the details of a deterministic construction to [141]. First, we query the entire ground set $[1 \dots n]$ and this allows us to identify any singleton sets in the partition. So now we assume that every set in our partitions has size at least 2. Consider a pair of elements x and y belonging to different sets S_x and S_y . Fix any other element $y' \in S_y$, arbitrarily (note that such a y' exists because we assume there are no singleton sets left). A query Q is said to be a *separating witness* for the tuple $\mathcal{T} = \langle x, y, y', S_x \rangle$ iff $Q \cap S_x = x$ and $y, y' \in Q$. (Observe that for any such query it will be the case that $b_x = 1$ while $b_y = 0$, hence the term “separating witness”. Observe that any partition \mathfrak{P} is completely determined by a set of queries with separating witnesses for all the tuples the partition contains. Now, form a query Q by picking each element independently and uniformly with probability $\frac{1}{2}$. Consider any particular

tuple $\mathbb{T} = \langle x, y, y', S_x \rangle$. Then

$$Pr_Q(Q \text{ is a separating witness for } \mathbb{T}) \geq \frac{1}{2^{B+2}}$$

Recall that $|S_x| \leq B$ since we are only considering partitions whose set sizes are at most B . Now consider a collection \mathcal{Q} of such queries each chosen independently of the other. Then for a given tuple \mathbb{T} we have that

$$Pr_{\mathcal{Q}}(\forall Q \in \mathcal{Q} Q \text{ is not a separating witness for } \mathbb{T}) \leq \left(1 - \frac{1}{2^{B+2}}\right)^{|\mathcal{Q}|}$$

But there are at most $\binom{n}{B+2} * B^3 \leq n^{B+2}$ such tuples. Hence,

$$Pr_{\mathcal{Q}}(\exists \text{tuple } \mathbb{T} \forall Q \in \mathcal{Q} Q \text{ is not a separating witness for } \mathbb{T}) \leq n^{B+2} * \left(1 - \frac{1}{2^{B+2}}\right)^{|\mathcal{Q}|}$$

Thus by choosing $|\mathcal{Q}| = O((B+2) * 2^{B+2} * \log n) = O(\log n)$ queries we can ensure that the above probability is below 1, which means that the collection \mathcal{Q} contains a separating witness for every possible tuple. This shows that the query complexity is $O(\log n)$. \square

2.10 Conclusion

Scheduling has a significant impact on the fair sharing of processing resources among virtual machines and on enforcing any applicable usage caps per virtual machine. This is specially important in commercial services like computing cloud services, where customers who pay for the same grade of service expect to receive the same access to resources and providers offer pricing models. However, the Xen hypervisor (and perhaps others) uses a scheduling mechanism which may fail to detect and account for CPU usage by poorly-behaved virtual machines, allowing malicious customers to obtain enhanced service at the expense of others.

We have demonstrated this vulnerability in the lab and on Amazon's Elastic Compute Cloud (EC2). Under laboratory conditions, we found that the applications exploiting this vulnerability are able to utilize up to 98% of a CPU core, regardless

of competition from other virtual machines. Amazon EC2 uses a patched version of Xen, which prevents the capped amount of CPU resources of other VMs from being stolen. However, our attack scheme can steal idle CPU cycles to increase its share, and obtain up to 85% of CPU resources (as mentioned earlier, we have been in discussions with Amazon about the vulnerability reported in this chapter and our recommendations for fixes; they have since implemented a fix that we have tested and verified). We describe four approaches to eliminating this vulnerability, and demonstrate their effectiveness and negligible overhead. Finally, we give an algorithm in conjunction with our attack to discover the co-placement of VMs, this important mechanism can be utilized to conduct coordinated attacks in Cloud.

Reference

- [1] Amazon Elastic Compute Cloud. <http://aws.amazon.com/ec2/>.
- [2] AMD Virtualization Technology. <http://www.amd.com/>.
- [3] Facebook places. <http://www.facebook.com/facebookplaces>.
- [4] Foursquare. <https://foursquare.com/>.
- [5] Gowalla. <http://gowalla.com/>.
- [6] Microsoft Azure Services Platform. <http://www.microsoft.com/azure/default.aspx>.
- [7] Microsoft Windows Azure Platform. <http://www.microsoft.com/azure/default.aspx>.
- [8] Strangers helping strangers. <http://www.facebook.com/SHStrangers>.
- [9] Yelp. <http://www.yelp.com/>.
- [10] Method and system for protecting websites from public internet threats. United States Patent 7,260,639, August 2007.
- [11] Amazon Web Services: Overview of Security Processes, 2008. <http://developer.amazonwebservices.com>.
- [12] Method and system for providing on-demand content delivery for an origin server. United States Patent 7,376,736, May 2008.