

Erasmus - HA

also, appeared in NSDI 2007.



ELSEVIER

Contents lists available at ScienceDirect

Computer Networks

journal homepage: www.elsevier.com/locate/comnet



SLA: service level agreement.

Sandpiper: Black-box and gray-box resource management for virtual machines

Timothy Wood ^{a,*}, Prashant Shenoy ^a, Arun Venkataramani ^a, Mazin Yousif ^b

^a University of Massachusetts, Dept. of Computer Science, 140 Governor's Drive, Amherst, MA 01003, United States

^b Avirtec, 1236 E. Grant Road, Tucson, AZ 85719, United States

hotspots
monitoring
detection
resize
migrate

ARTICLE INFO

Article history:
Available online 8 July 2009

Keywords:
Virtualization
Data centers
Migration
Dynamic provisioning

ABSTRACT

Virtualization can provide significant benefits in data centers by enabling dynamic virtual machine resizing and migration to eliminate hotspots. We present Sandpiper, a system that automates the task of monitoring and detecting hotspots, determining a new mapping of physical to virtual resources, resizing virtual machines to their new allocations, and initiating any necessary migrations. Sandpiper implements a black-box approach that is fully OS- and application-agnostic and a gray-box approach that exploits OS- and application-level statistics. We implement our techniques in Xen and conduct a detailed evaluation using a mix of CPU, network and memory-intensive applications. Our results show that Sandpiper is able to resolve single server hotspots within 20 s and scales well to larger, data center environments. We also show that the gray-box approach can help Sandpiper make more informed decisions, particularly in response to memory pressure.

© 2009 Elsevier B.V. All rights reserved.

service downtime < 1sec.

[Q.] e.g. SLAs? - avg. service time min throughput | R.T < 5sec, 99%

1. Introduction

Data centers—server farms that run networked applications—have become popular in a variety of domains such as web hosting, enterprise systems, and e-commerce sites. Server resources in a data center are multiplexed across multiple applications—each server runs one or more applications and application components may be distributed across multiple servers. Further, each application sees dynamic workload fluctuations caused by incremental growth, time-of-day effects, and flash crowds [1]. Since applications need to operate above a certain performance level specified in terms of a service level agreement (SLA),

effective management of data center resources while meeting SLAs is a complex task.

One possible approach for reducing management complexity is to employ *virtualization*. In this approach, applications run on virtual servers that are constructed using virtual machines, and one or more virtual servers are mapped onto each physical server in the system. Virtualization of data center resources provides numerous benefits. It enables application isolation since malicious or greedy applications can not impact other applications co-located on the same physical server. It enables server consolidation and provides better multiplexing of data center resources across applications. Perhaps the biggest advantage of employing virtualization is the ability to flexibly remap physical resources to virtual servers in order to handle workload dynamics. A workload increase can be handled by increasing the resources allocated to a virtual server if idle resources are available on the physical server, or by simply migrating the virtual server to a less loaded physical server. Migration is transparent to the applications and modern virtualization platforms support this capability [6,16]. How-

(i)
(ii)
(iii)

* A preliminary version of this paper, *Black-box and Gray-box Strategies for Virtual Machine Migration*, appeared in NSDI 2007.

* Corresponding author. Tel.: +1 4135454753.

E-mail addresses: twood@cs.umass.edu (T. Wood), shenoy@cs.umass.edu (P. Shenoy), arun@cs.umass.edu (A. Venkataramani), mazin@avirtec.net (M. Yousif).

[Q.]

metrics associated with resource mgmt.? ~ cost, resource utilization, SLAs met.

~ challenges/reqs.

(i)

ever, detecting workload hotspots and initiating a migration is currently handled manually. Manually-initiated migration lacks the agility to respond to sudden workload changes; it is also error-prone since each reshuffle might require migrations or swaps of multiple virtual servers to rebalance system load. Migration is further complicated by the need to consider multiple resources—CPU, network, and memory—for each application and physical server.

(ii) resize

To address this challenge, this paper studies automated black-box and gray-box strategies for virtual machine provisioning in large data centers. Our techniques automate the tasks of monitoring system resource usage, hotspot detection, allocating resources and initiating any necessary migrations. More importantly, our black-box techniques can make these decisions by simply observing each virtual machine from the outside and without any knowledge of the application resident within each VM. We also present a gray-box approach that assumes access to a small amount of OS-level statistics in addition to external observations to better inform the provisioning algorithm. Since a black-box approach is more general by virtue of being OS and application-agnostic, an important aspect of our research is to understand if a black-box approach alone is sufficient and effective for hotspot detection and mitigation. We have designed and implemented the Sandpiper system to support either black-box, gray-box, or combined techniques. We seek to identify specific limitations of the black-box approach and understand how a gray-box approach can address them.

(i)
(ii)
(iii)

Sandpiper implements a hotspot detection algorithm that determines when to resize or migrate virtual machines, and a hotspot mitigation algorithm that determines what and where to migrate and how many resources to allocate. The hotspot detection component employs a monitoring and profiling engine that gathers usage statistics on various virtual and physical servers and constructs profiles of resource usage. These profiles are used in conjunction with prediction techniques to detect hotspots in the system. Upon detection, Sandpiper grants additional resources to the overloaded servers if available. If necessary, Sandpiper's migration manager is invoked for further hotspot mitigation. The migration manager employs provisioning techniques to determine the resource needs of overloaded VMs and uses a greedy algorithm to determine a sequence of moves or swaps to migrate overloaded VMs to underloaded servers.

We have implemented our techniques using the Xen platform [3]. We conduct a detailed experimental evaluation on a testbed of two dozen servers using a mix of CPU-, network- and memory-intensive applications. Our results show that Sandpiper can alleviate single server hotspots in less than 20 s and more complex multi-server hotspots in a few minutes. Our results show that Sandpiper imposes negligible overheads and that gray-box statistics enable Sandpiper to make better migration decisions when alleviating memory hotspots.

The rest of this paper is structured as follows. Section 2 presents some background and Sections 3–6 present our design of Sandpiper. Section 7 presents our implementation and evaluation. Finally, Sections 8 and 9 present related work and our conclusions, respectively.

2. Background and system overview

Historically, approaches to dynamic provisioning have either focused on dynamic replication, where the number of servers allocated to an application is varied, or dynamic slicing, where the fraction of a server allocated to an application is varied. With the re-emergence of server virtualization, application migration has become an option for dynamic provisioning. Since migration is transparent to applications executing within virtual machines, our work considers this third approach—resource provisioning via dynamic migrations in virtualized data centers. We present *Sandpiper*,¹ a system for automated resource allocation and migration of virtual servers in a data center to meet application SLAs. Sandpiper assumes a large cluster of possibly heterogeneous servers. The hardware configuration of each server—its CPU, network interface, disk and memory characteristics—is assumed to be known to Sandpiper. Each physical server (also referred to as a physical machine or PM) runs a *virtual machine monitor* and one or more virtual machines. Each virtual server runs an application or an application component (the terms virtual servers and virtual machine are used interchangeably). Sandpiper currently uses Xen to implement such an architecture. Each virtual server is assumed to be allocated a certain slice of the physical server resources. In the case of CPU, this is achieved by assigning a subset of the host's CPUs to each virtual machine, along with a weight that the underlying Xen CPU scheduler uses to allocate CPU bandwidth. In case of the network interface, Xen is yet to implement a similar fair-share scheduler; a best-effort FIFO scheduler is currently used and Sandpiper is designed to work with this constraint. In case of memory, a slice is assigned by allocating a certain amount of RAM to each resident VM. All storage is assumed to be on a network file system or a storage area network, thereby eliminating the need to move disk state during VM migrations [6].

(i)
(ii)

Q.
e.g.c?
!

setup & elasticity

Sandpiper runs a component called the *nucleus* on each physical server; the *nucleus* runs inside a special virtual server (domain-0 in Xen) and is responsible for gathering resource usage statistics on that server (see Fig. 1). It employs a *monitoring engine* that gathers processor, network interface and memory swap statistics for each virtual server. For gray-box approaches, it implements a daemon within each virtual server to gather OS-level statistics and perhaps application logs.

The *nuclei* periodically relay these statistics to the *Sandpiper control plane*. The control plane runs on a distinguished node and implements much of the intelligence in Sandpiper. It comprises three components: a *profiling engine*, a *hotspot detector* and a *migration and resizing manager* (see Fig. 1). The *profiling engine* uses the statistics from the nuclei to construct resource usage profiles for each virtual server and aggregate profiles for each physical server. The hotspot detector continuously monitors these usage profiles to detect hotspots—informally, a hotspot is said to have occurred if the aggregate usage of any resource (processor, network or memory) exceeds a threshold or if SLA violations occur for a “sustained” period. Thus, the hotspot

⊗ hotspot defn.

¹ A migratory bird.

⊗ elastic resources via ballooning?

Sandpiper metrics :- #SLAs met - resource util.

- time to detect overcome hotspot.

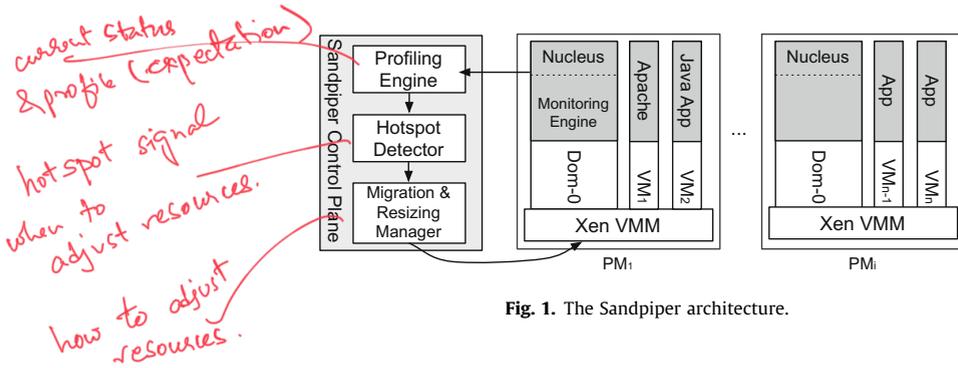


Fig. 1. The Sandpiper architecture.

detection component determines *when* to signal the need for resource adjustments and invokes the resource manager upon hotspot detection, which attempts hotspot mitigation via resizing or dynamic migrations. It implements algorithms that determine *how much* of a resource to allocate the virtual servers (i.e., determine a new resource allocation to meet the target SLAs), *what* virtual servers to migrate from the overloaded servers, and *where* to move them. The resource manager assumes that the virtual machine monitor implements a migration mechanism that is transparent to applications and uses this mechanism to automate migration decisions; Sandpiper currently uses Xen’s migration mechanisms that were presented in [6].

3. Monitoring and profiling in Sandpiper

Sandpiper supports both black- and gray-box monitoring techniques that are combined with profile generation tools to detect hotspots and predict VM resource requirements.

3.1. Unobtrusive black-box monitoring

The monitoring engine is responsible for tracking the processor, network and memory usage of each virtual server. It also tracks the total resource usage on each physical server by aggregating the usages of resident VMs. The monitoring engine tracks the usage of each resource over a measurement interval \mathcal{I} and reports these statistics to the control plane at the end of each interval.

In a pure black-box approach, all usages must be inferred solely from external observations and without relying on OS-level support inside the VM. Fortunately, much of the required information can be determined directly from the Xen hypervisor or by monitoring events within domain-0 of Xen. Domain-0 is a distinguished VM in Xen that is responsible for I/O processing; domain-0 can host device drivers and act as a “driver” domain that processes I/O requests from other domains [3,9]. As a result, it is possible to track network and disk I/O activity of various VMs by observing the driver activity in domain-0 [9]. Similarly, since CPU scheduling is implemented in the Xen hypervisor, the CPU usage of various VMs can be determined by tracking scheduling events in the hypervisor [10]. Thus, black-box monitoring can be implemented in the nucleus by tracking various domain-0 events and without modifying any virtual server. Next, we discuss CPU, network and memory monitoring using this approach.

CPU monitoring: By instrumenting the Xen hypervisor, it is possible to provide domain-0 with access to CPU scheduling events which indicate when a VM is scheduled and when it relinquishes the CPU. These events are tracked to determine the duration for which each virtual machine is scheduled within each measurement interval \mathcal{I} . The Xen 3 distribution includes a monitoring application called *XenMon* [10] that tracks the CPU usages of the resident virtual machines using this approach; for simplicity, the monitoring engine employs a modified version of XenMon to gather CPU usages of resident VMs over a configurable measurement interval \mathcal{I} . On a multi-cpu system, a VM may only be granted access to a subset of the total CPUs. However, the number of CPUs allocated to a virtual machine can be adjusted dynamically.

It is important to realize that these statistics do not capture the CPU overhead incurred for processing disk and network I/O requests; since Xen uses domain-0 to process disk and network I/O requests on behalf of other virtual machines, this processing overhead gets charged to the CPU utilization of domain-0. To properly account for this request processing overhead, analogous to proper accounting of interrupt processing overhead in OS kernels, we must apportion the CPU utilization of domain-0 to other virtual machines. We assume that the monitoring engine and the nucleus impose negligible overhead and that all of the CPU usage of domain-0 is primarily due to requests processed on behalf of other VMs. Since domain-0 can also track I/O request events based on the number of memory page exchanges between domains, we determine the number of disk and network I/O requests that are processed for each VM. Each VM is then charged a fraction of domain-0’s usage based on the proportion of the total I/O requests made by that VM. A more precise approach requiring a modified scheduler was proposed in [9].

Network monitoring: Domain-0 in Xen implements the network interface driver and all other domains access the driver via clean device abstractions. Xen uses a virtual firewall-router (VFR) interface; each domain attaches one or more virtual interfaces to the VFR [3]. Doing so enables Xen to multiplex all its virtual interfaces onto the underlying physical network interfaces.

Consequently, the monitoring engine can conveniently monitor each VM’s network usage in Domain-0. Since each virtual interface looks like a modern NIC and Xen uses Linux drivers, the monitoring engines can use the Linux /proc interface (in particular /proc/net/dev) to monitor

imp
+
approx.
does Tx & Rx Nw & Disk need same CPU? work.

the number of bytes sent and received on each interface. These statistics are gathered over interval \mathcal{I} and returned to the control plane.

Memory monitoring: Black-box monitoring of memory is challenging since Xen allocates a user-specified amount of memory to each VM and requires the OS within the VM to manage that memory; as a result, the memory utilization is only known to the OS within each VM. It is possible to instrument Xen to observe memory accesses within each VM through the use of shadow page tables, which is used by Xen’s migration mechanism to determine which pages are dirtied during migration. However, trapping each memory access results in a significant application slowdown and is only enabled during migrations[6]. Thus, memory usage statistics are not directly available and must be inferred.

The only behavior that is visible externally is swap activity. Since swap partitions reside on a network disk, I/O requests to swap partitions need to be processed by domain-0 and can be tracked. By tracking the reads and writes to each swap partition from domain-0, it is possible to detect memory pressure within each VM.

Our monitoring engine tracks the number of read and write requests to swap partitions within each measurement interval \mathcal{I} and reports it to the control plane. Since substantial swapping activity is indicative of memory pressure, our current black-box approach is limited to reactive decision making and can not be proactive.

3.2. Gray-box monitoring

Black-box monitoring is useful in scenarios where it is not feasible to “peek inside” a VM to gather usage statistics. Hosting environments, for instance, run third-party applications, and in some cases, third-party installed OS distributions. Amazon’s Elastic Computing Cloud (EC2) service, for instance, provides a “barebone” virtual server where customers can load their own OS images. While OS instrumentation is not feasible in such environments, there are environments such as corporate data centers where both the hardware infrastructure and the applications are owned by the same entity. In such scenarios, it is feasible to gather OS-level statistics as well as application logs, which can potentially enhance the quality of decision making in Sandpiper.

Sandpiper supports gray-box monitoring, when feasible, using a light-weight monitoring daemon that is installed inside each virtual server. In Linux, the

monitoring daemon uses the `/proc` interface to gather OS-level statistics of CPU, network, and memory usage. The memory usage monitoring, in particular, enables proactive detection and mitigation of memory hotspots. The monitoring daemon also can process logs of applications such as web and database servers to derive statistics such as request rate, request drops and service times. Direct monitoring of such application-level statistics enables explicit detection of SLA violations, in contrast to the black-box approach that uses resource utilizations as a proxy metric for SLA monitoring.

3.3. Profile generation

The profiling engine receives periodic reports of resource usage from each nucleus. It maintains a usage history for each server, which is then used to compute a profile for each virtual and physical server. A profile is a compact description of that server’s resource usage over a sliding time window W . Three black-box profiles are maintained per virtual server: CPU utilization, network bandwidth utilization, and swap rate (i.e., page fault rate). If gray-box monitoring is permitted, four additional profiles are maintained: memory utilization, service time, request drop rate and incoming request rate. Similar profiles are also maintained for each physical server, which indicate the aggregate usage of resident VMs.

Each profile contains a distribution and a time series. The distribution, also referred to as the distribution profile, represents the probability distribution of the resource usage over the window W . To compute a CPU distribution profile, for instance, a histogram of observed usages over all intervals \mathcal{I} contained within the window W is computed; normalizing this histogram yields the desired probability distribution (see Fig. 2).

While a distribution profile captures the variations in the resource usage, it does not capture temporal correlations. For instance, a distribution does not indicate whether the resource utilization increased or decreased within the window W . A time-series profile captures these temporal fluctuations and is simply a list of all reported observations within the window W . For instance, the CPU time-series profile is a list (C_1, C_2, \dots, C_k) of the k reported utilizations within the window W . Whereas time-series profiles are used by the hotspot detector to spot increasing utilization trends, distribution profiles are used by the migration manager to estimate peak resource requirements and provision accordingly.

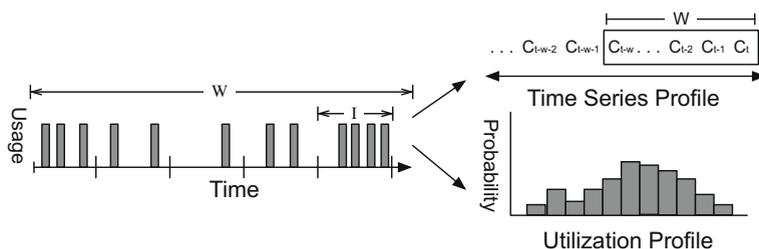


Fig. 2. Profile generation in Sandpiper

gray box memory mngmt.

What statistic about memory is measured?

(*)

✓

*

*

(i)

(i)

4. Hotspot detection

The hotspot detection algorithm is responsible for signaling a need for VM resizing whenever SLA violations are detected implicitly by the black-box approach or explicitly by the gray-box approach. Hotspot detection is performed on a per-physical server basis in the black-box approach—a hotspot is flagged if the aggregate CPU or network utilizations on the physical server exceed a threshold or if the total swap activity exceeds a threshold. In contrast, explicit SLA violations must be detected on a per-virtual server basis in the gray-box approach—a hotspot is flagged if the memory utilization of the VM exceeds a threshold or if the response time or the request drop rate exceed the SLA-specified values.

To ensure that a small transient spike does not trigger needless migrations, a hotspot is flagged only if thresholds or SLAs are exceeded for a sustained time. Given a time-series profile, a hotspot is flagged if at least k out of the n most recent observations as well as the next predicted value exceed a threshold. With this constraint, we can filter out transient spikes and avoid needless migrations. The values of k and n can be chosen to make hotspot detection aggressive or conservative. For a given n , small values of k cause aggressive hotspot detection, while large values of k imply a need for more sustained threshold violations and thus a more conservative approach. Similarly, larger values of n incorporate a longer history, resulting in a more conservative approach. In the extreme, $n = k = 1$ is the most aggressive approach that flags a hotspot as soon as the threshold is exceeded. Finally, the threshold itself also determines how aggressively hotspots are flagged; lower thresholds imply more aggressive migrations at the expense of lower server utilizations, while higher thresholds imply higher utilizations with the risk of potentially higher SLA violations.

Sandpiper employs time-series prediction techniques to predict future values [4]. Specifically, Sandpiper relies on the auto-regressive family of predictors, where the n th order predictor $AR(n)$ uses n prior observations in conjunction with other statistics of the time series to make a prediction. To illustrate the first-order $AR(1)$ predictor, consider a sequence of observations: u_1, u_2, \dots, u_k . Given this time series, we wish to predict the demand in the $(k + 1)$ th interval. Then the first-order $AR(1)$ predictor makes a prediction using the previous value u_k , the mean of the time series values μ , and the parameter ϕ which captures the variations in the time series [4]. The prediction \hat{u}_{k+1} is given by:

$$\hat{u}_{k+1} = \mu + \phi(u_k - \mu). \quad (1)$$

As new observations arrive from the nuclei, the hot spot detector updates its predictions and performs the above checks to flag new hotspots in the system.

5. Resource provisioning

A hotspot indicates a resource deficit on the underlying physical server to service the collective workloads of resident VMs. Before the hotspot can be resolved, Sandpiper

must first estimate how much additional resources are needed by the overloaded VMs to fulfill their SLAs; these estimates are then used to determine if local resource allocation adjustments or migrations are required to resolve the hotspot.

5.1. Black-box provisioning

The provisioning component needs to estimate the peak CPU, network and memory requirement of each overloaded VM; doing so ensures that the SLAs are not violated even in the presence of peak workloads.

Estimating peak CPU and network bandwidth needs: Distribution profiles are used to estimate the peak CPU and network bandwidth needs of each VM. The tail of the usage distribution represents the peak usage over the recent past and is used as an estimate of future peak needs. This is achieved by computing a high-percentile (e.g., the 95th percentile) of the CPU and network bandwidth distribution as an initial estimate of the peak needs.

Since both the CPU scheduler and the network packet scheduler in Xen are work-conserving, a VM can use more than its fair share, provided that other VMs are not using their full allocations. In case of the CPU, for instance, a VM can use a share that exceeds the share determined by its weight, so long as other VMs are using less than their weighted share. In such instances, the tail of the distribution will exceed the guaranteed share and provide insights into the actual peak needs of the application. Hence, a high-percentile of the distribution is a good first approximation of the peak needs.

However, if all VMs are using their fair shares, then an overloaded VM will not be allocated a share that exceeds its guaranteed allocation, even though its peak needs are higher than the fair share. In such cases, the observed peak usage (i.e., the tail of the distribution) will equal its fair share. In this case, the tail of the distribution will underestimate the actual peak need. To correct for this underestimate, the provisioning component must scale the observed peak to better estimate the actual peak. Thus, whenever the CPU or the network interface on the physical server are close to saturation, the provisioning component first computes a high-percentile of the observed distribution and then adds a constant Δ to scale up this estimate.

Example. Consider two virtual machines that are assigned CPU weights of 1:1 resulting in a fair share of 50% each. Assume that VM₁ is overloaded and requires 70% of the CPU to meet its peak needs. If VM₂ is underloaded and only using 20% of the CPU, then the work-conserving Xen scheduler will allocate 70% to VM₁. In this case, the tail of the observed distribution is a good indicator of VM₁'s peak need. In contrast, if VM₂ is using its entire fair share of 50%, then VM₁ will be allocated exactly its fair share. In this case, the peak observed usage will be 50%, an underestimate of the actual peak need. Since Sandpiper can detect that the CPU is fully utilized, it will estimate the peak to be 50 + Δ .

The above example illustrates a fundamental limitation of the black-box approach—it is not possible to estimate

Q: what is tail of usage?

Q: what is work conserving?

approx

the true peak need when the underlying resource is fully utilized. The scale-up factor Δ is simply a guess and might end up over- or under-estimating the true peak.

Estimating peak memory needs: Xen allows an adjustable amount of physical memory to be assigned to each resident VM; this allocation represents a hard upper-bound that can not be exceeded regardless of memory demand and regardless of the memory usage in other VMs. Consequently, our techniques for estimating the peak CPU and network usage do not apply to memory. The provisioning component uses observed swap activity to determine if the current memory allocation of the VM should be increased. If swap activity exceeds the threshold indicating memory pressure, then the current allocation is deemed insufficient and is increased by a constant amount Δ_m . Observe that techniques such as Geiger and hypervisor level caches that attempt to infer working set sizes by observing swap activity [11,14] can be employed to obtain a better estimate of memory needs; however, our current prototype uses the simpler approach of increasing the allocation by a fixed amount Δ_m whenever memory pressure is observed.

5.2. Gray-box provisioning

Since the gray-box approach has access to application-level logs, information contained in the logs can be utilized to estimate the peak resource needs of the application. Unlike the black-box approach, the peak needs can be estimated even when the resource is fully utilized.

To estimate peak needs, the peak request arrival rate is first estimated. Since the number of serviced requests as well as the number of dropped requests are typically logged, the incoming request rate is the summation of these two quantities. Given the distribution profile of the arrival rate, the peak rate is simply a high-percentile of the distribution. Let λ_{peak} denote the estimated peak arrival rate for the application.

Estimating peak CPU needs: An application model is necessary to estimate the peak CPU needs. Applications such as web and database servers can be modeled as G/G/1 queuing systems [24]. The behavior of such a G/G/1 queuing system can be captured using the following queuing theory result [13]:

$$\lambda_{cap} \geq \left[s + \frac{\sigma_a^2 + \sigma_b^2}{2 \cdot (d - s)} \right]^{-1}, \tag{2}$$

where d is the mean response time of requests, s is the mean service time, and λ_{cap} is the request arrival rate. σ_a^2 and σ_b^2 are the variance of inter-arrival time and the variance of service time, respectively. Note that response time includes the full queuing delay, while service time only reflects the time spent actively processing a request.

While the desired response time d is specified by the SLA, the service time s of requests as well as the variance of inter-arrival and service times σ_a^2 and σ_b^2 can be determined from the server logs. By substituting these values into Eq. 2, a lower bound on request rate λ_{cap} that can be serviced by the virtual server is obtained. Thus, λ_{cap} represents the current capacity of the VM.

To service the estimated peak workload λ_{peak} , the current CPU capacity needs to be scaled by the factor $\frac{\lambda_{peak}}{\lambda_{cap}}$. Ob-

serve that this factor will be greater than 1 if the peak arrival rate exceeds the currently provisioned capacity. Thus, if the VM is currently assigned a CPU weight w , its allocated share needs to be scaled up by the factor $\frac{\lambda_{peak}}{\lambda_{cap}}$ to service the peak workload.

Estimating peak network needs: The peak network bandwidth usage is simply estimated as the product of the estimated peak arrival rate λ_{peak} and the mean requested file size b ; this is the amount of data transferred over the network to service the peak workload. The mean request size can be computed from the server logs.

Estimating memory needs: Using OS-level information about a virtual machine's memory utilization allows the gray-box approach to more accurately estimate the amount of memory required by a virtual machine. The gray-box approach can proactively adjust memory allocations when the OS reports that it is low on memory (but before swapping occurs). This data is also used to safely reduce the amount of memory allocated to VMs which are not using their full allotment, something which is impossible to do with only black-box information about swapping.

6. Hotspot mitigation

Once a hotspot has been detected, Sandpiper must determine if the hotspots can be resolved with local resource adjustments, or if migrations are required to balance load between hosts.

6.1. VM resizing

While large changes in resource needs may require migration between servers, some hot spots can be handled by adjusting the resource allocation of the overloaded VM. Sandpiper first attempts to increase the resource allocation for an overloaded VM by either adding additional CPUs, network interfaces, or memory depending on which resource utilizations exceeded the warning thresholds.

If the profiling engine detects that a VM is experiencing an increasing usage of CPU, Sandpiper will attempt to allocate an additional virtual CPU to the VM. Xen and other virtualization platforms support dynamic changes in the number of CPUs a VM has access to by exploiting hot-swapping code that already exists in many operating system kernels. A similar approach can be used to add network interfaces to a VM, although this is not currently supported by Sandpiper.

In many cases, memory hotspots can also be resolved through local provisioning adjustments. When a VM has insufficient memory as detected by either swapping (black-box) or OS statistics (gray-box), Sandpiper will first attempt to increase the VM's memory allocation on its current host. Only if there is insufficient spare memory will the VM be migrated to a different host.

6.2. Load balancing with migration

If there are insufficient spare resources on a host, the migration manager invokes its hotspot mitigation algorithm to determine which virtual servers to migrate and

Q. ballooning or memory hotplug?
? web server.
how much to add?/remove?

of memory? on vho? or disk IO? (is approx. only CPU)

where in order to dissipate the hotspot. Determining a new mapping of VMs to physical servers that avoids threshold violations is NP-hard—the multi-dimensional bin packing problem can be reduced to this problem, where each physical server is a bin with dimensions corresponding to its resource constraints and each VM is an object that needs to be packed with size equal to its resource requirements. Even the problem of determining if a valid packing exists is NP-hard.

Consequently, our hotspot mitigation algorithm resorts to a heuristic to determine which overloaded VMs to migrate and where *such that migration overhead is minimized*. Reducing the migration overhead (i.e., the amount of data transferred) is important, since Xen's live migration mechanism works by iteratively copying the memory image of the VM to the destination while keeping track of which pages are being dirtied and need to be resent. This requires Xen to intercept all memory accesses for the migrating domain, which significantly impacts the performance of the application inside the VM. By reducing the amount of data copied over the network, Sandpiper can minimize the total migration time, and thus, the performance impact on applications. Note that network bandwidth available for application use is also reduced due to the background copying during migrations; however, on a gigabit LAN, this impact is small.

Capturing multi-dimensional loads: Once the desired resource allocations have been determined by either our black-box or gray-box approach, the problem of finding servers with sufficient idle resource to house overloaded VMs is identical for both. The migration manager employs a greedy heuristic to determine which VMs need to be migrated. The basic idea is to move load from the most overloaded servers to the least-overloaded servers, while attempting to minimize data copying incurred during migration. Since a VM or a server can be overloaded along one or more of three dimensions—CPU, network and memory—we define a new metric that captures the combined CPU-network-memory load of a virtual and physical server. The *volume* of a physical or virtual server is defined as the product of its CPU, network and memory loads:

$$Vol = \frac{1}{1 - cpu} * \frac{1}{1 - net} * \frac{1}{1 - mem}, \quad (3)$$

where *cpu*, *net* and *mem* are the corresponding utilizations of that resource normalized by the number of CPUs and network interfaces allocated to the virtual or physical server.² The higher the utilization of a resource, the greater the volume; if multiple resources are heavily utilized, the above product results in a correspondingly higher volume. The volume captures the degree of (over)load along multiple dimensions in a unified fashion and can be used by the mitigation algorithms to handle all resource hotspots in an identical manner.

² If a resource is fully utilized, its utilization is set to $1 - \epsilon$, rather than one, to avoid infinite volume servers. Also, since the black-box approach is oblivious of the precise memory utilization, the value of *mem* is set to 0.5 in the absence of swapping and to $1 - \epsilon$ if memory pressure is observed; the precise value of *mem* is used in the gray-box approach.

Migration phase: To determine which VMs to migrate, the algorithm orders physical servers in decreasing order of their volumes. Within each server, VMs are considered in decreasing order of their *volume-to-size ratio* (VSR); where VSR is defined as *Volume/Size*; size is the memory footprint of the VM. By considering VMs in VSR order, the algorithm attempts to migrate the maximum volume (i.e., load) per unit byte moved, which has been shown to minimize migration overhead [21].

The algorithm proceeds by considering the highest VSR virtual machine from the highest volume server and determines if it can be housed on the least volume (least loaded) physical server. The move is feasible only if that server has sufficient idle CPU, network and memory resources to meet the desired resource allocation of the candidate VM as determined by the provisioning component (Section 5). Since we use VSR to represent three resource quantities, the least loaded server may not necessarily “fit” best with a particular VM's needs. If sufficient resources are not available, then the algorithm examines the next least loaded server and so on, until a match is found for the candidate VM. If no physical server can house the highest VSR VM, then the algorithm moves on to the next highest VSR VM and attempts to move it in a similar fashion. The process repeats until the utilizations of all resources on the physical server fall below their thresholds.

The algorithm then considers the next most loaded physical server that is experiencing a hotspot and repeats the process until there are no physical servers left with a hotspot. The output of this algorithm is a list of overloaded VMs and a new destination server for each; the actual migrations are triggered only after all moves have been determined.

Swap phase: In cases where there are not sufficient idle resources on less loaded servers to dissipate a hotspot, the migration algorithm considers VM swaps as an alternative. A swap involves exchanging a high VSR virtual machine from a loaded server with one or more low VSR VMs from an underloaded server. Such a swap reduces the overall utilization of the overloaded server, albeit to a lesser extent than a one-way move of the VM. Our algorithm considers the highest VSR VM on the highest volume server with a hotspot; it then considers the lowest volume server and considers the *k* lowest VSR VMs such that these VMs collectively free up sufficient resources to house the overloaded VM. The swap is considered feasible if the two physical servers have sufficient resources to house the other server's candidate VM(s) without violating utilization thresholds. If a swap cannot be found, the next least loaded server is considered for a possible swap and so on. The process repeats until sufficient high VSR VMs have been swapped with less loaded VMs so that the hotspot is dissipated. Although multi-way swaps involving more than two servers can also be considered, our algorithm presently does not implement such complex swaps. The actual migrations to perform the swaps are triggered only after a list of all swaps is constructed. Note that a swap may require a third server with “scratch” RAM to temporarily house a VM before it moves to its final destination. An alternative is to (i) suspend one of the VMs on disk, (ii) use the freed up RAM to accommodate the other VM,

⊛ select PM in hotspot and move VM with highest VSR to least loaded PM! ~ VM that is not in hotspot may have to move!

and (iii) resume the first VM on the other server; doing so is not transparent to the temporarily suspended VM.

7. Implementation and evaluation

The implementation of Sandpiper is based on Xen. The Sandpiper *control plane* is implemented as a daemon that runs on the control node. It listens for periodic usage reports from the various nuclei, which are used to generate profiles. The profiling engine currently uses a history of the past 200 measurements to generate virtual and physical server profiles. The hotspot detector uses these profiles to detect hotspots; currently a hotspot is triggered when 3 out of 5 past readings and the next predicted value exceed a threshold. The default threshold is set to 75%. The migration manager implements our provisioning and hotspot mitigation algorithms; it notifies the nuclei of any desired migrations, which then trigger them. In all, the control plane consists of less than 750 lines of Python code.

The Sandpiper *nucleus* is a Python application that extends the XenMon CPU monitor to also acquire network and memory statistics for each VM. The monitoring engine in the nucleus collects and reports measurements once every 10 s—the default measurement interval. The nucleus uses Xen's Python management API to trigger migrations and adjust resource allocations as directed by the control plane. While black-box monitoring only requires access to domain-0 events, gray-box monitoring employs two additional components: a Linux OS daemon and an Apache module.

The gray-box linux daemon runs on each VM that permits gray-box monitoring. It currently gathers memory statistics via the `/proc` interface—the memory utilization, the number of free pages and swap usage are reported to the monitoring engine in each interval. The gray-box Apache module comprises of a real-time log analyzer and a dispatcher. The log-analyzer processes log-entries as they are written to compute statistics such as the service time, request rate, request drop rate, inter-arrival times, and request/file sizes. The dispatcher is implemented as a kernel module based on Linux IP Virtual server (IPVS) ver 1.2.1; the goal of the kernel module is to accurately estimate the request arrival rate during overload periods, when a high fraction of requests may be dropped. Since requests can be dropped at the TCP layer as well as at the HTTP layer

during overloads, the use of a transport-level dispatcher such as IPVS is necessary for accurately estimating the drop (and hence arrival) rates. Ordinarily, the kernel dispatcher simply forwards incoming requests to Apache for processing. In all, the nucleus comprises 650 lines of Python code.

Our evaluation of Sandpiper is based on a prototype data center consisting of twenty 2.4Ghz Pentium-4 servers connected over a gigabit Ethernet. These servers run Linux 2.6 and Xen 3.0.2-3 and are equipped with at least 1 GB of RAM. Experiments involving multi-core systems run on Intel Quad-Core servers with 4 GB of RAM and Xen 3.1. A cluster of Pentium-3 Linux servers is used to generate workloads for our experiments. One node in the cluster is designated to run the Sandpiper control plane, while the rest host one or more VMs, all of which run the Sandpiper nucleus in domain-0. In the following experiments, our VMs run Apache 2.0.54, PHP 4.3.10, and MySQL 4.0.24.

7.1. VM resizing

While migrations are necessary for large changes in resource allocations, it is less expensive if resources can be adjusted locally without the overhead of migration. This experiment demonstrates Sandpiper's ability to detect increasing CPU requirements and respond by allocating additional CPU cores to the virtual machine.

Initially, a VM running a CPU intensive web application is allocated a single CPU core. During the experiment, the number of clients accessing the web server increases. Sandpiper responds by increasing the number of virtual CPUs allocated to the VM. The VM starts on a dual core host; as the load continues to rise, a migration is required to move the VM to a host with four available CPUs as shown in Fig. 3.

Result: Resizing a VM's resource allocation incurs little overhead. When additional resources are not available locally, migrations are required.

7.2. Migration effectiveness

Our next experiment exercises Sandpiper's hotspot detection and migration algorithms; we subject a set of black-box servers to a series of workloads that repeatedly place the system in overload. Our experiment uses three physical servers and five VMs with memory allocations

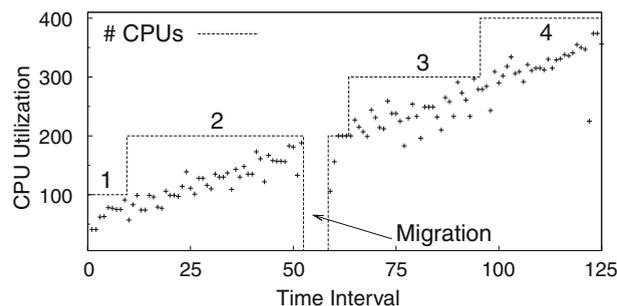


Fig. 3. Sandpiper increases the number of virtual CPU cores allocated to a VM. A migration is required to move from a 2 to 4 core PM.

Table 1

Workload in requests/second, memory allocations, and initial placement.

VM	Peak 1	Peak 2	Peak 3	RAM (MB)	Start PM
1	200	130	130	256	1
2	90	90	90	256	1
3	60	200	60	256	2
4	60	90	90	256	2
5	10	10	130	128	3

as shown in Table 1. All VMs run Apache serving dynamic PHP web pages. The PHP scripts are designed to be CPU intensive so that a low client request rate places a large CPU load on a server without significant network or memory utilization. We use *httperf* to inject a workload that goes through three phases, each of which causes a hotspot on a different physical machine. The peak request rates for each phase are shown in Table 1.

Fig. 4 presents a time series of the load placed on each VM along with the triggered migrations. In the first phase, a large load is placed on VM₁, causing the CPU utilization on PM₁ to exceed the CPU threshold. The system detects a hotspot at $t = 166$ s. The migration manager examines candidates for migration in VSR order. VM₁ has the highest VSR, so it is selected as a candidate. Since PM₃ has sufficient spare capacity to house VM₁, it is migrated there, thereby eliminating the hotspot less than 20 s after detection. This represents the ideal case for our algorithm: if possible, we try to migrate the most loaded VM from an overloaded PM to one with sufficient spare capacity.

In the second phase, PM₂ becomes overloaded due to increasing load on VM₃. However, the migration manager is unable to migrate this VM because there is insufficient capacity on the other PMs. As a result, at $t = 362$ s, the VM on PM₂ with the second highest VSR VM₄, is migrated to PM₁ that now has spare capacity. This demonstrates a more typical case where none of the underloaded PMs have sufficient spare capacity to run the overloaded PM's highest VSR VM, so instead we migrate less overloaded VMs that can fit elsewhere.

In the final phase, PM₃ becomes overloaded when both its VMs receive *identical* large loads. Unlike the previous two cases where candidate VMs had equal memory footprints, VM₅ has half as much RAM as VM₁, so it is chosen for migration.

Result: To eliminate hotspots while minimizing the overhead of migration, our placement algorithm tries to move the highest VSR VM to the least loaded PM. This maximizes the amount of load displaced from the hotspot per megabyte of data transferred.

7.3. Mixed resource workloads

Sandpiper can consolidate applications that stress different resources to improve the overall multiplexing of server resources. Our setup comprises two servers with two VMs each. Both VMs on the first server are network-intensive, involving large file transfers, while those on the second server are CPU-intensive running Apache with dynamic PHP scripts. All VMs are initially allocated 256 MB of memory. VM₂ additionally runs a main-memory database that stores its tables in memory, causing its memory usage to grow over time.

Figs. 5a and b shows the resource utilization of each PM over time. Since PM₁ has a network hotspot and PM₂ has a CPU hotspot, Sandpiper swaps a network-intensive VM for a CPU-intensive VM at $t = 130$. This results in a lower CPU and network utilization on both servers. Fig. 5(d) shows the initial utilizations of each VM; after the swap, the aggregate CPU and network utilizations on both servers falls below 50%.

In the latter half, memory pressure increases on VM₂ due to its main-memory database application. As shown in Fig. 5c, Sandpiper responds by increasing the RAM allocation in steps of 32 MB every time swapping is observed; when no additional RAM is available, the VM is swapped to the second physical server at $t = 430$. This is feasible because two cpu-intensive jobs are swapped, leaving CPU and network utilization balanced, and the second physical server has more RAM than the first. Memory allocations

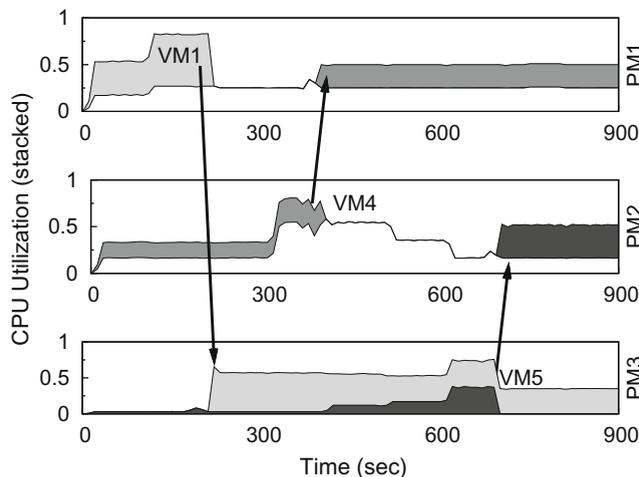


Fig. 4. A series of migrations resolve hotspots. Different shades are used for each migrating VM.

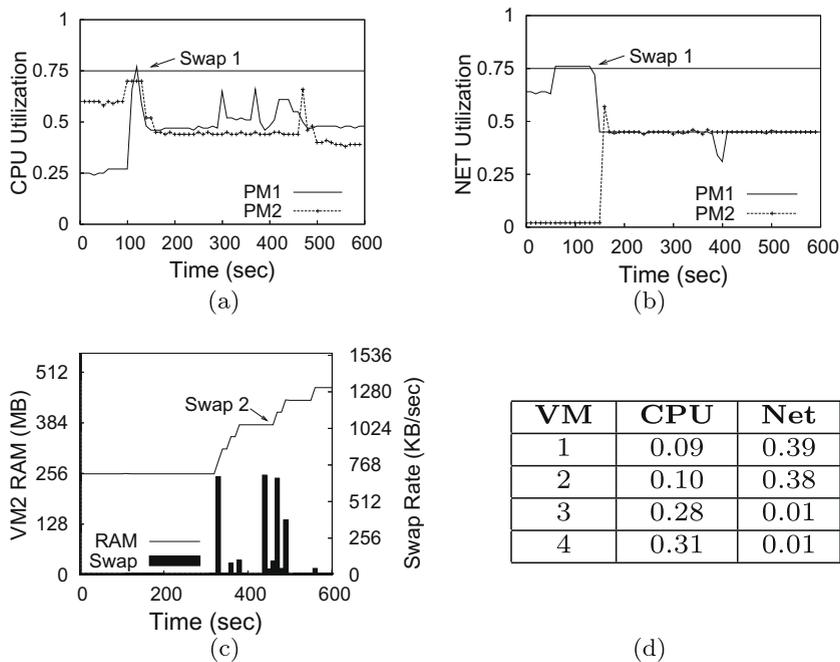


Fig. 5. Swaps and migrations to handle network- and memory-intensive loads. Initially, VM₁ and VM₂ are on PM₁, the rest on PM₂. After two swaps, PM₁ hosts VM₁ and VM₄.

are reactive since only black-box stats are available. Next we demonstrate how a gray-box approach can proactively respond to memory pressure.

Result: Sandpiper can respond to network, CPU, or memory hotspots and can collocate VMs that stress different resources to improve overall system utilization.

7.4. Gray v. Black: memory allocation

We compare the effectiveness of the black- and gray-box approaches in mitigating memory hotspots using the SPECjbb 2005 benchmark. SPECjbb emulates a three-tier web application based on J2EE servers. We use SPECjbb to apply an increasingly intense workload to a single VM. The workload increases every two minutes, causing a significant increase in memory usage. After twenty minutes, the application reaches its peak intensity, after which the workload decreases at a similar rate.

The VM is initially assigned 256 MB of RAM, and resides on a physical machine with 384 MB total RAM. We also run a second, idle physical server which has 1 GB RAM. We run the experiment with two separate pairs of servers, Black and Gray, that correspond to the black- and gray-box approaches, respectively. The Gray system is configured to signal a hotspot whenever the amount of free RAM in the virtual machine falls below 32 MB.

Fig. 6a plots the memory allocation of the VM over time. Both systems gradually increase the VM's memory until all unused RAM is exhausted. Since Black can only respond to swapping, it lags in responsiveness. At $t = 380$ s, Gray determines that there is insufficient RAM for the VM and migrates it to the second PM with 1 GB RAM; Black initiates the same migration shortly afterward. Both continue

to increase the VM's memory as the load rises. Throughout the experiment, Black writes a total of 32 MB to swap, while Gray only writes 2 MB. Note that a lower memory hotspot threshold in Gray can prevent swapping altogether, while Black can not eliminate swapping due to its reactive nature.

During the second phase of the trial, Gray is able to detect the decreasing memory requirements and is able to safely reduce the VM's memory allocation. Since the black-box system can only detect swapping, it cannot reduce the memory allocation without fear of causing swapping and worse performance.

Result: A key weakness of the black-box approach is its inability to infer memory usage. Using this information, the gray-box system can reduce or eliminate swapping and can safely decrease a VM's memory allocation.

7.5. Gray v. Black: Apache performance

Recall from Section 5 that when resources are fully utilized, they hamper the black-box approach from accurately determining the needs of overloaded VMs. This experiment demonstrates how a black-box approach may incur extra migrations to mitigate a hotspot, whereas a gray-box approach can use application-level knowledge for faster hotspot mitigation.

Our experiment employs three physical servers and four VMs. Initially, VM₁ through VM₃ reside on PM₁, VM₄ resides on PM₂, and PM₃ is idle. We use httpperf to generate requests for CPU intensive PHP scripts on all VMs. At $t = 80$ s, we rapidly increase the request rates on VM₁ and VM₂ so that actual CPU requirement for each virtual machine reaches 70%, creating an extreme hotspot on

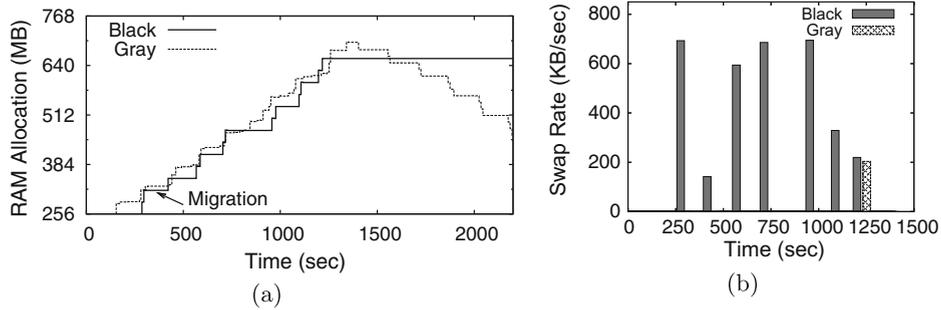


Fig. 6. The black-box system lags behind the gray-box system in allocating memory. The gray-box approach proactively increases memory and safely reduces the VM's memory allocation when demand falls.

PM₁. The request rates for VM₃ and VM₄ remain constant, requiring 33% and 7% CPU respectively. We use an aggressive 6 s measurement interval so that Sandpiper can respond quickly to the increase in workload.

Without accurate estimates of each virtual machine's resource requirements, the black-box system falters in its decision making as indicated in Fig. 7. Since the CPU on PM₁ is saturated, each virtual machine receives an equal portion of processing time and appears equivalent to Sandpiper. Sandpiper must select a VM at random, and in the worst case, tries to eliminate the hotspot by migrating VM₃ to PM₃. Since VM₁ and VM₂ continue to reside on PM₁, the hotspot persists even after the first migration. Next, the black-box approach assumes that VM₂ requires only 50% of the CPU and migrates it to PM₂. Unfortunately, this results in PM₂ becoming overloaded, so a final migration must be performed to move VM₄ to PM₃.

We repeat this scenario with the Apache gray-box module running inside of each virtual machine. Since the gray-box monitor can precisely measure the incoming request rates, Sandpiper can accurately estimate the CPU needs of VM₁ and VM₂. By using this information, Sandpiper is able to efficiently respond to the hotspot by immediately migrating VM₃ to PM₂ and VM₂ to PM₃. Fig. 8 depicts the improved performance of the gray-box approach. Note that since Sandpiper requires the hotspot to persist for k out of n intervals before it acts, it is not until $t = 98$ s that either system considers itself overloaded. Once a hotspot is flagged, the gray-box approach can mitigate it within 40 s with just two migrations, while the black-box approach requires 110 s and three migrations to do so. Although response time increases equally under both systems, the

gray-box approach is able to reduce response times to an acceptable level 61% faster than the black-box system, producing a corresponding reduction in SLA violations.

Result: Application-level statistics enable the gray-box approach to better infer resource needs and improves the quality of migration decisions, especially in scenarios where resource demands exceed server capacity.

7.6. Prototype data center evaluation

Next we conduct an experiment to demonstrate how Sandpiper performs under realistic data center conditions. We deployed a prototype data center on a cluster of 16 servers that run a total of 35 VMs. An additional node runs the control plane and one node is reserved as a scratch node for swaps. The virtual machines run a mix of data center applications ranging from Apache and streaming servers to LAMP servers running Apache, PHP, and MySQL within a single VM. We run RUBiS on our LAMP servers—RUBiS is an open-source multi-tier web application that implements an eBay-like auction web site and includes a workload generator that emulates users browsing and bidding on items.

Of the 35 deployed VMs, 5 run the RUBiS application, 5 run streaming servers, 5 run Apache serving CPU-intensive PHP scripts, 2 run main memory database applications, and the remaining 15 serve a mix of PHP scripts and large HTML files. We use the provided workload generators for the RUBiS applications and use httperf to generate requests for the other servers.

To demonstrate Sandpiper's ability to handle complex hotspot scenarios, we orchestrate a workload that causes

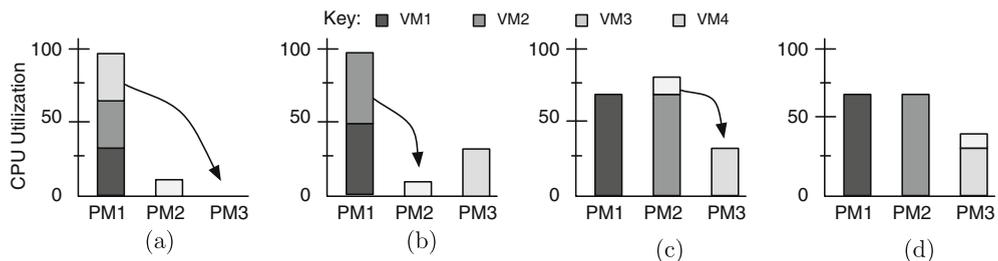


Fig. 7. The black-box system incorrectly guesses resource requirements since CPU usage is saturated, resulting in an increased resolution time. The gray-box system infers usage requirements and transitions directly from (a) to (d).

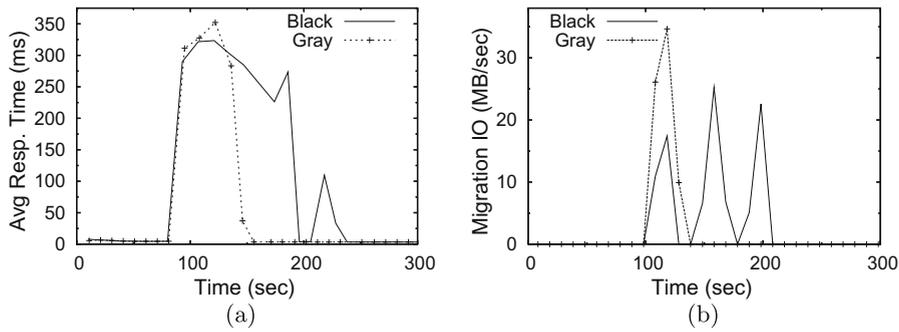


Fig. 8. The gray-box system balances the system more quickly due to more informed decision making. The black-box system must perform migrations sequentially and incurs an additional migration.

multiple network and CPU hotspots on several servers. Our workloads causes six physical servers running a total of 14 VMs to be overloaded—four servers see a CPU hotspot and two see a network hotspot. Of the remaining PMs, 4 are moderately loaded (greater than 45% utilization for at least one resource) and 6 have lighter loads of between 25% and 40% utilization. We compare Sandpiper to a statically allocated system with no migrations.

Fig. 9 demonstrates that Sandpiper eliminates hotspots on all six servers by interval 60. These hotspots persist in the static system until the workload changes or a system administrator triggers manual migrations. Due to Xen’s migration overhead, there are brief periods where Sandpiper causes more physical servers to be overloaded than in the static case. Despite this artifact, even during periods where migrations are in progress, Sandpiper reduces the number of intervals spent in sustained overload by 61%. In all, Sandpiper performs seven migrations and two swaps to eliminate all hotspots over a period of 237 s after hotspot detection.

Result: Sandpiper is capable of detecting and eliminating simultaneous hotspots along multiple resource dimensions. Despite Xen’s migration overhead, the number of servers experiencing overload is decreased even while migrations are in progress.

7.7. System overhead and scalability

Sandpiper’s CPU and network overhead is dependent on the number of PMs and VMs in the data center. With only

black-box VMs, the type of application running in the VM has no effect on Sandpiper’s overhead. If gray-box modules are in use, the overhead may vary depending on the size of application-level statistics gathered.

Nucleus overheads: Sandpiper’s nucleus sends reports to the Control Plane every measurement interval (10 s by default). The table in Fig. 10a gives a breakdown of overhead for each report type. Since each report uses only 288 bytes per VM, the resulting overhead on a gigabit LAN is negligible. To evaluate the CPU overhead, we compare the performance of a CPU benchmark with and without our resource monitors running. Even on a single physical server running 24 concurrent VMs, our monitoring overheads only reduce the CPU benchmark performance by approximately one percent. This is comparable to the overheads reported by XenMon, which much of our code is based on [10].

Control plane scalability: The main source of computational complexity in the control plane is the computation of a new mapping of virtual machines to physical servers after detecting a hotspot. Although the problem is NP-hard, we only require an approximate solution, and our heuristics make the problem tractable for reasonable system sizes. For data centers with up to 500 virtual servers, the algorithm completes in less than 5 s as shown in Fig. 10b. For very large data centers with thousands of virtual machines, the computation could be split up across multiple nodes, or the center’s servers can be broken up into pools, each controlled independently by its own control plane.

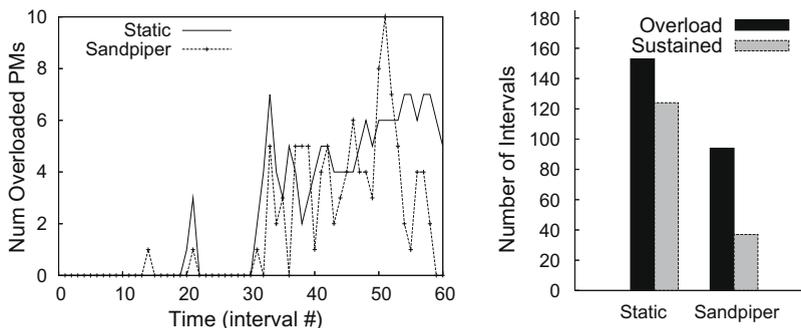


Fig. 9. Sandpiper eliminates all hotspots and reduces the number of intervals experiencing sustained overload by 61%.

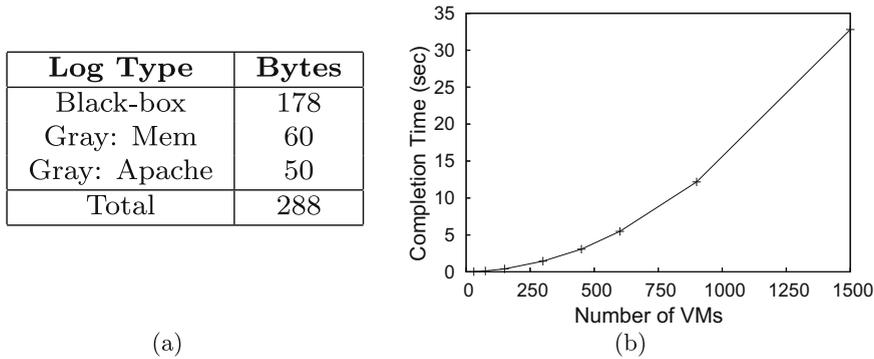


Fig. 10. Sandpiper overhead and scalability.

7.8. Stability during overloads

This section demonstrates how Sandpiper ensures stable system behavior by avoiding “thrashing” migrations. First, Sandpiper avoids migrations to physical machines with rising loads, since this can trigger additional migrations if the load rises beyond the threshold; time-series predictions are used to determine future load trends when selecting a physical server. Thus, Fig. 11a shows that when a migration decision is required at $t = 140$ s, Sandpiper will prefer PM_2 over PM_1 as a target. Even though PM_2 has a higher current load, the 120 s prediction window indicates a rising load on PM_1 .

Next, we demonstrate Sandpiper’s behavior in the presence of increasing number of hotspots. We simulate a data center with fifty physical servers, each with three virtual servers. We increase the number of simultaneous hotspots from 20 to 45; the mean utilizations are set to 85% and 45% for servers with and without hotspots. Fig. 11b depicts the mean number of migrations performed to resolve these hotspots over multiple runs. If fewer than half of the servers are overloaded, then all hotspots can typically be resolved with one migration per overloaded server. After this threshold, swaps are required and it is increasingly difficult to fully resolve overload until it becomes infeasible. With 35 overloaded servers, Sandpiper was able to eliminate all hotspots 73% of the time (over multiple runs); with

40 overloaded servers, a complete solution was found only 3% of the time. In the extreme case, Sandpiper is still able to resolve 22 of the 45 hotspots before giving up. In all cases, Sandpiper first finds a solution before initiating migrations or swaps; when no feasible solutions are found, Sandpiper either implements a partial solution or gives up entirely rather than attempting wasteful migrations. This bounds the number of migrations which will ever be performed and explains the decrease in migrations beyond 40 overloaded servers, where there is no feasible solution.

7.9. Tuning Sandpiper

Sandpiper has several parameters which the system administrator can tune to make hotspot detection and mitigation more or less aggressive. Our experiments suggest the following rules of thumb.

Setting thresholds: If overload thresholds are set too high, then the additional overhead during migration can cause additional SLA violations. Our experiments show that the average throughput of a CPU-intensive Apache server can drop by more than 50% during a migration. We suggest a CPU threshold of 75% to absorb the CPU overhead of migration while maximizing server utilization. We also suggest a 75% threshold for network utilization based on experiments in [6] which indicate that the network throughput of a highly loaded server can drop by about

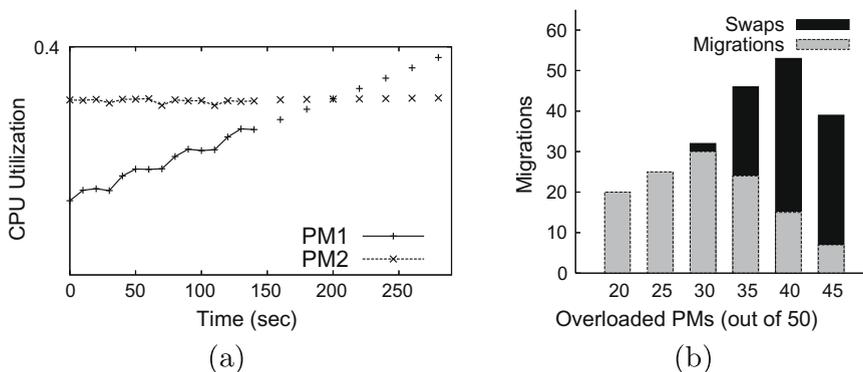


Fig. 11. (a) Using time-series predictions (the dotted lines) allows Sandpiper to better select migration destinations, improving stability. (b) Higher levels of overload requires more migrations until there is no feasible solution.

20% during portions of a migration (due to network copying overheads).

Sustained overload requirement: Our experiments (not reported here) reveal that Sandpiper is not sensitive to a particular choice of the measurement interval \mathcal{I} so long as it is between a few seconds and a few tens of seconds. For a measurement interval of 10 s, we suggest $k = 3$ and $n = 5$ for the “ k out of n ” check; this corresponds to requiring the time period of about 3 migrations to exceed the resource threshold before we initiate a migration. The Δ parameter is used in the black-box system to increase resource allocations when utilization is saturated. This should be set equal to the maximum increase in resource requirements that a service is likely to see during a measurement interval and may vary based on workload; we use 10% in our experiments. Using more advanced time series forecasting techniques would allow Sandpiper to dynamically determine Δ .

8. Related work

Our work draws upon recent advances in virtual machines and dynamic provisioning in data centers to address a question of increasing research and commercial interest: can virtual machine migration enable robust and highly responsive provisioning in data centers? The Xen migration work [6] alludes to this motivation. What is missing is a convincing validation and algorithms to effect migration, which is the focus of this paper.

The idea of process migration was first investigated in the 80s [23]. Support for migrating groups of processes across OSES was presented in [17], but applications had to be suspended and it did not address the problem of maintaining open network connections. Virtualization support for commodity operating systems in [7] led towards techniques for virtual machine migration over long time spans, suitable for WAN migration [20]. More recently, Xen [6] and VMWare [16] have implemented “live” migration of VMs that involve extremely short downtimes ranging from tens of milliseconds to a second. VM migration has been used for dynamic resource allocation in Grid environments [19,22,8]. A system employing automated VM migrations for scientific nano-technology workloads on federated grid environments was investigated in [19]. The Shirako system provides infrastructure for leasing resources within a federated cluster environment and was extended to use virtual machines for more flexible resource allocation in [8]. Shirako uses migrations to enable dynamic placement decisions in response to resource broker and cluster provider policies. In contrast, we focus on data center environments with stringent SLA requirements that necessitate highly responsive migration algorithms for online load balancing. VMware’s Distributed Resource Scheduler [25] uses migration to perform automated load balancing in response to CPU and memory pressure. DRS uses a userspace application to monitor memory usage similar to Sandpiper’s gray-box monitor, but unlike Sandpiper, it cannot utilize application logs to respond directly to potential SLA violations or to improve placement decisions.

Dedicated hosting is a category of dynamic provisioning in which each physical machine runs at most one application and workload increases are handled by spawning a new replica of the application on idle servers. Physical server granularity provisioning has been investigated in [1,18]. Techniques for modeling and provisioning multi-tier Web services by allocating physical machines to each tier are presented in [24]. Although dedicated hosting provides complete isolation, the cost is reduced responsiveness – without virtualization, moving from one physical machine to another takes on the order of several minutes [24] making it unsuitable for handling flash crowds. Our current implementation does not replicate virtual machines, implicitly assuming that PMs are sufficiently provisioned.

Shared hosting is the second variety of dynamic provisioning, and allows a single physical machine to be shared across multiple services. Various economic and resource models to allocate shared resources have been presented in [5]. Mechanisms to partition and share resources across services include [2,5]. A dynamic provisioning algorithm to allocate CPU shares to VMs on a single physical machine (as opposed to a cluster) was presented and evaluated through simulations in [15]. In comparison to the above systems, our work assumes a shared hosting platform and uses VMs to partition CPU, memory, and network resources, but additionally leverages VM migration to meet SLA objectives.

Estimating the resources needed to meet an application’s SLA requires a model that inspects the request arrival rates for the application and infers its CPU, memory, and network bandwidth needs. Developing such models is not the focus of this work and has been addressed by several previous efforts such as [12,1].

9. Conclusions

This paper argued that virtualization provides significant benefits in data centers by enabling virtual machine migration to eliminate hotspots. We presented Sandpiper, a system that automates the task of monitoring and detecting hotspots, determining a new mapping of physical to virtual resources, and resizing or migrating VMs to eliminate the hotspots. We discussed a black-box strategy that is fully OS- and application-agnostic as well as a gray-box approach that can exploit OS- and application-level statistics. An evaluation of our Xen-based prototype showed that VM migration is a viable technique for rapid hotspot elimination in data center environments. Using solely black-box methods, Sandpiper is capable of eliminating simultaneous hotspots involving multiple resources. We found that utilizing gray-box information can improve the responsiveness of our system, particularly by allowing for proactive memory allocations and better inferences about resource requirements.

Acknowledgements

We would like to thank our anonymous reviewers for their helpful comments. This research was supported by

NSF Grants EEC-0313747, CNS-0720271, CNS-0720616, CNS-0325868 and a gift from Intel.

References

- [1] K. Appleby, S. Fakhouri, L. Fong, M. Goldszmidt, S. Krishnakumar, D. Pazel, J. Pershing, B. Rochwerger, Oceanos-sla-based management of a computing utility, in: Proceedings of the IFIP/IEEE Symposium on Integrated Network Management, May 2001.
- [2] M. Aron, P. Druschel, W. Zwaenepoel, Cluster reserves: a mechanism for resource management in cluster-based network servers, in: Proceedings of the ACM SIGMETRICS Conference, Santa Clara, CA, June 2000.
- [3] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, A. Warfield, Xen and the art of virtualization, in: Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP'03), Bolton Landing, NY, October 2003, pp. 164–177.
- [4] G.P. Box, G.M. Jenkins, G.C. Reinsel, Time Series Analysis Forecasting and Control, third ed., Prentice Hall, 1994.
- [5] J. Chase, D. Anderson, P. Thakar, A. Vahdat, R. Doyle, Managing energy and server resources in hosting centers, in: Proceedings of the 18th ACM Symposium on Operating Systems Principles (SOSP), October 2001, pp. 103–116.
- [6] C. Clark, K. Fraser, S. Hand, J. Hansen, E. Jul, C. Limpach, I. Pratt, A. Warfield, Live migration of virtual machines, in: Proceedings of Usenix Symposium on Network Systems Design and Implementation (NSDI), May 2005.
- [7] K. Govil, D. Teodosiu, Y. Huang, M. Rosenblum, Cellular disco: resource management using virtual clusters on shared-memory multiprocessors, in: Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'99), Kiawah Island Resort, SC, December 1999, pp. 154–169.
- [8] L. Grit, D. Irwin, A. Yumerefendi, J. Chase, Virtual machine hosting for networked clusters: building the foundations for autonomic orchestration, in: The First International Workshop on Virtualization Technology in Distributed Computing (VTDC), November 2006.
- [9] D. Gupta, L. Cherkasova, R. Gardner, A. Vahdat, Enforcing performance isolation across virtual machines in xen, in: Proceedings of the ACM/IFIP/USENIX Seventh International Middleware Conference (Middleware'2006), Melbourne, Australia, November 2006.
- [10] D. Gupta, R. Gardner, L. Cherkasova, Xenmon: Qos monitoring and performance profiling tool, Tech. Rep. HPL-2005-187, HP Labs, 2005.
- [11] S. Jones, A. Arpaci-Dusseau, R. Arpaci-Dusseau, Geiger: monitoring the buffer cache in a virtual machine environment, in: Proceedings of the 12th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'00), Cambridge, MA, October 2006, pp. 13–23.
- [12] A. Kamra, V. Misra, E. Nahum, Yaksha: a self-tuning controller for managing the performance of 3-tiered web sites, in: International Workshop on Quality of Service (IWQoS), June 2004.
- [13] L. Kleinrock, Queueing Systems, Computer Applications, vol. 2, John Wiley and Sons Inc., 1976.
- [14] P. Lu, K. Shen, Virtual machine memory access tracing with hypervisor exclusive cache, in: Usenix Annual Technical Conference, June 2007.
- [15] D.A. Menasce, M.N. Bennani, Autonomic virtualized environments, in: IEEE International Conference on Autonomic and Autonomous Systems, July 2006.
- [16] M. Nelson, B.-H. Lim, G. Hutchins, Fast transparent migration for virtual machines, in: USENIX Annual Technical Conference, 2005.
- [17] S. Osman, D. Subhraveti, G. Su, J. Nieh, The design and implementation of zap: a system for migrating computing environments, in: Proceedings of the Fifth Symposium on Operating Systems Design and Implementation (OSDI), 2002.
- [18] S. Ranjan, J. Rolia, H. Fu, E. Knightly, Qos-driven server migration for internet data centers, in: Proceedings of IWQoS 2002, Miami Beach, FL, May 2002.
- [19] P. Ruth, J. Rhee, D. Xu, R. Kennell, S. Goasguen, Autonomic live adaptation of virtual computational environments in a multi-domain infrastructure, in: IEEE International Conference on Autonomic Computing (ICAC), June 2006.
- [20] C.P. Sapuntzakis, R. Chandra, B. Pfaff, J. Chow, M.S. Lam, M. Rosenblum, Optimizing the migration of virtual computers, in: Proceedings of the Fifth Symposium on Operating Systems Design and Implementation, December 2002.
- [21] V. Sundaram, T. Wood, P. Shenoy, Efficient data migration in self-managing storage systems, in: Proceedings of the Third IEEE International Conference on Autonomic Computing (ICAC-06), Dublin, Ireland, June 2006.
- [22] A. Sundararaj, A. Gupta, P. Dinda, Increasing application performance in virtual environments through run-time inference and adaptation, in: Fourteenth International Symposium on High Performance Distributed Computing (HPDC), July 2005.
- [23] M. Theimer, K. Lantz, D. Cheriton, Preemptable remote execution facilities for the v-system, in: Proceedings of the 10th SOSP, Operating Systems Review, 1985.
- [24] B. Urgaonkar, P. Shenoy, A. Chandra, P. Goyal, Dynamic provisioning for multi-tier internet applications, in: Proceedings of the Second IEEE International Conference on Autonomic Computing (ICAC-05), Seattle, WA, June 2005.
- [25] Resource Management with VMware DRS, VMware Whitepaper.



Timothy Wood is a computer science Ph.D. student at the University of Massachusetts Amherst. He received his B.S. in Electrical and Computer Engineering from Rutgers University in 2005 and his M.S. in Computer Science from the University of Massachusetts in 2009. His research focuses on improving the reliability, performance, and management of modern data centers, with an emphasis on exploiting the benefits of virtualization.



Prashant Shenoy received the B.Tech. degree in Computer Science and Engineering from the Indian Institute of Technology, Bombay in 1993, and the M.S. and Ph.D. degrees in Computer Science from the University of Texas, Austin, in 1994 and 1998, respectively. He is currently an Associate Professor of Computer Science at the University of Massachusetts Amherst. His research interests are in operating and distributed systems, sensor networks, Internet systems and pervasive multimedia. He has been the recipient of the

National Science Foundation Career Award, the IBM Faculty Development Award, the Lilly Foundation Teaching Fellowship, the UT Computer Science Best Dissertation Award and an IIT Silver Medal. He is a senior member of the IEEE and the ACM.



Arun Venkataramani has been an assistant professor at University of Massachusetts Amherst since 2005 after receiving his Ph.D. from University of Texas at Austin by way of University of Washington. His research interests are in networked and distributed systems.



Mazin Yousif is currently the Chief Technology Officer for Avirtec, Inc., where he leads technical and business related issues for the company. Before that, he held technical executive positions at Numonyx, Intel and IBM corporations. He held adjunct and research professor positions at various universities including Arizona, the Oregon Graduate Institute (OGI), Duke and North Carolina State Universities. He finished his M.S. and Ph.D. degrees from the Pennsylvania State University in 1987 and 1992, respectively. His

research interests include computer architecture, cloud computing, autonomic computing, workload profiling/prediction and workload-dri-

ven platform architectures. He has published 55+ articles in his areas of research. He chaired several conferences and workshops and served in the program committee of many others. He is an Associate Editor in IEEE ToC, in the advisory board of the Journal of Pervasive Computing and Communications (JPCC); editor in the Journal of Autonomic and Trusted Computing (JoATC), and is an editor in Cluster Computing, the Journal of Networks, Software Tools and Applications. He chairs the Advisory Committee of ERCIM (The European Research Consortium for Informatics and Mathematics) and is an IEEE Distinguished Visitor program (2008 10210). He is an IEEE senior member.