

Dynamic Virtual Resource Management In Cloud Computing Environment

D.Golden Jemi, C.S Soumiya

PG Student, Department of CSE Loyola Institute of Technology and Science ¹jemicse@gmail.com
Assistant Professor, Department of CSE Loyola Institute of Technology and Science ²soumiya.cs@gmail.com

Abstract

Cloud Computing is the latest technology used by many organizations in this competitive world. As many organizations are using this technology, the major issue is efficient resource provisioning and management because of the dynamic nature of the Cloud and the need to satisfy heterogeneous resource requirements. In this paper, we present the Dynamic Resource Allocation system that uses virtual machines on physical machines with automatic virtual machine migrations and without service interruption. This can be done with the help of Skewness Algorithm. The flexible resource provisioning and migrations of machine state improves the efficiency of resource usage and dynamic resource provisioning capabilities. We introduce the concept of data compression in resource provisioning for reducing the storage space and bandwidth. We also focus on ensuring data storage security in cloud data centers. The experimental results demonstrate that our system enhance the performance and security of the cloud computing environment.

Keywords— Virtual machines, Cloud Computing, dynamic resource allocation, data compression.

I. INTRODUCTION

One of the most significant benefits of cloud computing is reducing the operating cost of data center through virtualization to support cost reduction, the resource of physical machines (PM) in data center should be efficiently utilized. Virtualization technologies enable application computation and data to be hosted inside virtual containers (e.g., virtual machines, virtual disks) which are decoupled from the underlying physical resources. For example, server virtualization technologies such as VMware and Xen [8], [11] facilitate application computation to be packaged inside virtual machines (VMs) and enable multiple such virtual machines to be run alongside each other on a single physical machine. This allows extensive sharing of physical resources across applications. Additionally, the new live-migration advancements [2], [14] allow VMs to be migrated from one server to another without any downtime to the application running inside it.

However, if the provider only considers maximizing the utilization of data centers (i.e., maximizing the utilization level of physical machines), eventually, it has a bad influence upon the performance of virtual machines (VM) [1],[5] in data center due to high workload of each associated physical machine. To prevent such a performance degradation, an appropriate VM allocation scheme is needed for the overall performance of cloud computing. In addition, the provider needs to set an appropriate threshold of utilization level that does not

affect to the performance degradation of virtual machines in the data center. For better performance of VMs in terms of response time, Consider the location of each VM to provide worldwide services, the provider should have several data centers according to geographically locations. Since cloud computing services are delivered over the public internet which does not guaranteed reliability in general, there may be undesirable performance degradations such as slow response time.

Although the provider can designate the allocation for new VMs to a low utilized PM [1] that guarantees no performance degradation due to utilization level, a performance degrade is still possible to occur. If the location of a PM that is providing the user's VM is far from the location of a user, the geographical distance between the PM and the user affect to the response time of the VM.

Therefore, a cloud provider needs to consider not only the utilization level of PMs, but also the location of a PM to allocate the user request as a VM. To address these issues, this paper proposes a dynamic resource allocation model. The new model considers 1) the location of PMs, and 2) the dynamic utilization level of PMs.

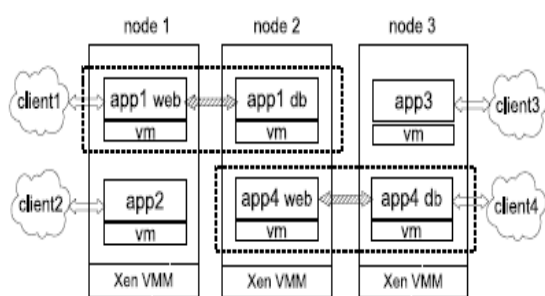


Figure 1: Virtualized infrastructure

Figure 1 shows a three-node subset of a virtualized infrastructure shared by multiple applications, where each tier of an application is hosted in a virtual machine (VM), and a multi-tier application (such as app1 or app4) may span multiple nodes [4]. Unlike the traditional hosting model where applications run on dedicated nodes, resulting in low resource utilization, this model allows applications to be consolidated onto fewer nodes, reducing capital expenditure on infrastructure as well as operating costs on power, cooling, maintenance, and support. It also leads to much higher resource utilization on the shared nodes.

II. RELATED WORKS

There are numerous advantages of cloud computing, the most basic ones being lower costs, re-provisioning of resources and remote accessibility. Cloud computing lowers cost by avoiding the capital expenditure by the company in renting the physical infrastructure from a third party provider. Due to the flexible nature of cloud computing, we can quickly access more resources from cloud providers when we need to expand our business.

The remote accessibility enables us to access the cloud services from anywhere at any time. To gain the maximum degree of the above mentioned benefits, the services offered in terms of resources should be allocated optimally to the applications running in the cloud. The following section discusses the significance of resource allocation.

A. Significance of Dynamic Resource Allocation

In cloud computing, Dynamic Resource Allocation is the process of assigning available resources to the needed cloud applications over the internet. Resource allocation starves services if the allocation is not managed precisely. Resource provisioning solves that problem by allowing the service providers to manage the resources for each individual module.

Resource Allocation Strategy is all about integrating cloud provider activities for utilizing and allocating scarce resources within the limit of cloud environment so as to meet the needs of the cloud

application. It requires the type and amount of resources needed by each application in order to complete user requirements. The order and time of allocation of resources are also an input for resource allocation. From the perspective of a cloud provider, predicting the dynamic nature of users, user demands, and application demands are impractical. Cloud resources consist of physical and virtual resources [15].

The physical resources are shared across multiple compute requests through virtualization and provisioning. In Cloud environments, efficient resource provisioning and management present today a challenging issue because of the dynamic nature of the Cloud on one hand, and the need to satisfy heterogeneous resource requirements on the other hand.

In such dynamic environments where end-users can arrive and leave the Cloud at any time, a Cloud service provider (CSP) should be able to make accurate decisions for scaling up or down its data centers while taking into account several utility criteria, the delay of virtual resources setup, the migration of existing processes, the resource utilization, etc. In order to satisfy parties (the CSP and the end-users), an efficient and dynamic resource allocation strategy is mandatory.

Dynamic Resource Allocation dealing with virtualization machines on physical machines. The results confirmed that the virtual machine which loading becomes too high, it will automatically migrated to another low loading physical machine without service interrupting.

The remainder of this paper is organized as follows: Section 3 presents a system overview of the proposed model. Section 4 describes the proposed system of dynamic resource allocation model. Section 5 describes the objective of Skewness Algorithm. Section 6 describes the analysis of Skewness Algorithm and shows the performance evaluation of the proposed model in terms of the allocation outcomes (response time of user's VM). Finally, Section 7 concludes this paper.

III. SYSTEM OVERVIEW

The system is designed to provide flexible and on-demand service, the proposed system allows to user to request an arbitrary amount of resources at any time and from anywhere. For managing flexible user request, the proposed system is designed as a hybrid architecture that is combined with centralized and distributed resource management architectures.

Physical Machine is a real resource (i.e. it is combined different types of resources such CPU, memory, or network bandwidth) that can allocate many VMs. In a PM, there is a specialized layer for virtualization called a hypervisor [1]. The hypervisor

generally has the responsibility to allocate VMs and share its resources like traditional operating systems [2]. The cloud data center is designed by means of security.

Each PM runs the hypervisor which supports one or more applications such as Web server, remote desktop, DNS, Mail, Map/ Reduce, etc. Assume all PMs Share backend storage. The load predictor predicts the future resource demands of VMs and the future load of PMs based on past statistics. Compute the load of a PM by aggregating the resource usage of its VMs. The hot spot solver detects if the resource utilization of any PM is above the hot threshold (i.e., a hot spot).

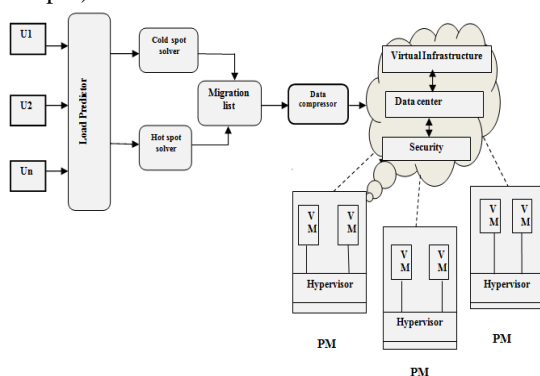


Figure 2: System Architecture

If some PM is present VMs running on them will be migrated away to reduce their load. The cold spot solver checks if the average utilization of actively used PMs (APMs) is below the green computing threshold. If so, some of those PMs could potentially be turned off to save energy. It identifies the set of PMs whose utilization is below the cold threshold (i.e., cold spots) and then attempts to migrate away all their VMs. It then compiles a migration list of VMs and passes it to the controller for execution and the compressed data is stored on cloud data centers.

IV. PROPOSED WORK

A. Resource Management:

Resource management poses particular challenges in large-scale systems, such as server clusters that simultaneously process requests from a large number of clients. Dynamic resource management in large scale cloud environment includes the physical infrastructure and associated control functionality that enables the provisioning and management of cloud services. Cloud resources consist of physical and virtual resources. The physical resources are shared across multiple compute requests through virtualization and provisioning.

The request for virtualized resources is described through a set of parameters detailing the

processing, memory and disk needs. Provisioning satisfies the request by mapping virtualized resources to physical ones. The hardware and software resources are allocated to the cloud applications on-demand basis. The Load Predictor is the entity responsible for optimizing resource allocation. When it receives a resource request, the Load Predictor iterates through the possible subsets of available resources

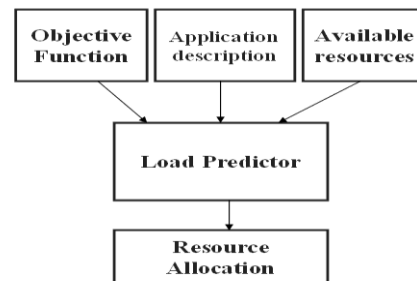


Figure No 3: Resource Allocation

Objective Function: The objective function defines the metric that should optimize. For example, given the increasing cost and scarcity of power in the data center, an objective function might measure the increase in power usage due to a particular allocation.

Application Description: The application description consists of three parts: i) the framework type that identifies the framework model to use ii) workload specific parameters that describe the particular application's resource usage and iii) a request for resources including the number of VMs, storage.

Available Resources: The final input required by the Load Predictor is a resource snapshot of the IaaS data centre. This includes information derived from both the virtualization layer and the IaaS monitoring service.

Load Predictor: The Load Predictor maps resource allocation candidates to the user with respect to a given objective function.

By using Resource Management, the Cloud providers provide resources for more number of users with less response time and can share their resources over the internet with high performance.

B. Virtualization:

In virtualization based Cloud Computing model the user not need to know the specific location of each physical device, operating system, memory, number of processor cores, middleware technology and so on. The users simply send their requests to the cloud.

The heart of virtualization is the "virtual machine" (VM), a tightly isolated software container with an operating system and application inside. Because each VM is completely separate and

independent, many of them can run simultaneously on a single computer. A thin layer of software called a hypervisor decouples the VMs from the host, and dynamically allocates computing resources to each VM as needed. This architecture redefines computing equation, to deliver:

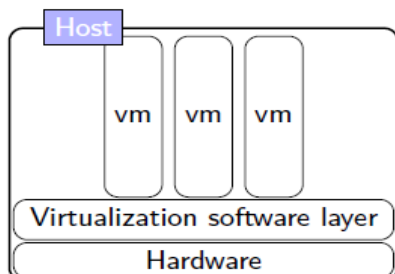


Figure No 4: Virtualization

Many applications on each server. As each VM encapsulates an entire machine, many applications and operating systems can be run on one host at the same time.

Maximum server utilization, minimum server count. Every physical machine is used to its full capacity, allowing you to significantly reduce costs by deploying fewer servers overall.

Faster, easier application and resources provisioning. As self-contained software files, VMs can be manipulated with copy-and-paste ease. This brings unprecedented simplicity, speed, and flexibility to IT provisioning and management. VMs can even be transferred from one physical server to another while running, via a process known as live migration.

C. Migration:

Migration of VMs has been investigated as a mean to adjust data-centers utilization. The VMs are periodically reallocated using migration according to their current resource demand. Virtual machine migration takes a running virtual machine and moves it from one physical machine to another.

When a VM is running a live service it is important that this transfer occurs in a manner that balances the requirements of minimizing both downtime and total migration time. The only perceived change should be a brief slowdown during the migration and a possible improvement in performance after the migration because the VM was moved to a machine with more available resources.

In terms of VM migration [5], there have been a few proposals to enable a VM to migrate from one physical machine to another for different purposes such as improving the power efficiency and satisfying performance requirements.

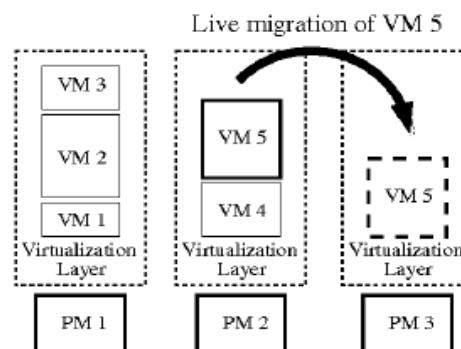


Figure No 5: Migration

Three physical machines with a virtualization layer are used to execute five virtual machines. In the initial configuration, PM 3 can be in low power state or powered down because it is not hosting any VMs. In response to demand change PM 3 is activated and VM 5 (denoted by solid lines) is migrated from PM 2 to PM 3. The migration occurs without service interruption.

D. Green Computing:

Green computing is the term used to denote efficient use of resources in computing. Core objectives of Green Computing Strategy is to Minimizing energy consumption, Purchasing green energy, Reducing the paper and other consumables used, Minimizing equipment disposal requirements and Reducing travel requirements for employees/customers. It also used for reduce costs

Computing Power Consumption has reached a Critical Point in Data centers have run out of usable power and cooling due to high densities. Computer virtualization is the process of running two or more logical computer systems on one set of physical hardware.

When the resource utilization of active servers [3] is too low, some of them can be turned off to save energy. This is handled by Green Computing. The challenge is to reduce the number of active servers during low load without sacrificing performance either now or in the future.

E. Data Compression

Data Compression is used to storing data in a format that requires less space than usual. Data compression is particularly useful in communications because it enables devices to transmit or store the same amount of data in fewer bits. A simple characterization of data compression is that it involves transforming a string of characters in some representation into a new string which contains the same information but whose length is as small as possible. Data compression has important application in the areas of data transmission and data storage.

When the amount of data to be transmitted is reduced, the effect is that of increasing the capacity of the communication channel bandwidth. Similarly, compressing a file to half of its original size is equivalent to doubling the capacity of the storage medium. It may then become feasible to store the data at a higher, thus faster, level of the storage hierarchy and reduce the load on the input/output channels of the computer system.

F. Security

The challenge is that data protection and security concerns will always be with us as we evolve to newer IT technologies and as long as there are those who try to hack through security measures and access data. In this paper, we focus on Ensuring data storage security in cloud computing, which is an important aspect of Quality of Service (QoS). Because cloud-based services use the Internet, storing data in the cloud can be risky and can mean less control over your data. The security issues can be solved by providing security for the user resources while allocating resources dynamically. This will improve the performance and security of Data centers.

V. OBJECTIVE OF SKEWNESS ALGORITHM

In proposed Dynamic Resource Allocation, the resources are allocated using Skewness Algorithm. The concept of skewness Algorithm is to quantify the unevenness in the utilization of multiple resources on a server. Let n be the number of resources consider is and r_i be the utilization of the i^{th} resource. Define the resource skewness of a server p as

$$skewness(p) = \sqrt{\sum_{i=1}^n \left(\frac{r_i}{\bar{r}} - 1\right)^2},$$

where \bar{r} is the average utilization of all resources for server p . In practice, not all types of resources are performance critical and hence only need to consider bottleneck resources in the above calculation. By minimizing the skewness, different types of workloads can combine nicely and improve the overall utilization of server resources.

A. Hot and Cold Spots

Algorithm is executed periodically to evaluate the resource allocation status based on the predicted future resource demands of VMs. The server as a hot spot if the utilization of any of its resources is above a hot threshold. This indicates that the server is overloaded and hence some VMs running on it should be migrated away. The temperature of a hot spot p is defined as the square

sum of its resource utilization beyond the hot threshold.

$$temperature(p) = \sum_{r \in R} (r - r_t)^2,$$

Where R is the set of overloaded resources in server p and r_t is the hot threshold for resource r . The temperature of a hot spot reflects its degree of overload. If a server is not a hot spot, its temperature is zero. A server is defined as a cold spot if the utilizations of all its resources are below a cold threshold.

B. Hot Spot Mitigation

Sort the list of hot spots in the system in descending temperature. Our goal is to eliminate all hot spots if possible. Otherwise, keep their temperature as low as possible. For each server p , first decide which of its VMs should be migrated away. Sort its list of VMs based on the resulting temperature of the server if that VM is migrated away.

The aim is to migrate away the VM that can reduce the server's temperature the most. In case of ties, select the VM whose removal can reduce the skewness of the server the most. For each VM in the list, the algorithm finds a destination server to accommodate it. The server must not become a hot spot after accepting this VM. Among all such servers, algorithm select one whose skewness can be reduced the most by accepting this VM. Note that this reduction can be negative which means algorithm selects the server whose skewness increases the least.

If a destination server is found, the algorithm records the migration of the VMs to that server and updates the predicted load of related servers. Otherwise, move onto the next VM in the list and try to find a destination server for it. As long as find a destination server for any of its VMs, Consider this run of the algorithm a success and then move onto the next hot spot. Note that each run of the algorithm migrates away at most one VM from the overloaded server.

C. Green Computing

When the resource utilization of active servers is too low, some of them can be turned off to save energy. This is handled in our green computing algorithm. The challenge is to reduce the number of active servers during low load without sacrificing performance either now or in the future.

Green computing algorithm is invoked when the average utilizations of all resources on active servers are below the green computing threshold. Sort the list of cold spots in the system based on the ascending order of their memory size. Since algorithm needs to migrate away all its VMs before shut down an under-utilized server, define the memory size of a cold spot as the aggregate memory size of all VMs running on it. For a cold spot p ,

check if all of its VMs are migrated somewhere else. For each VM on p , try to find a destination server to accommodate it. The resource utilizations of the server after accepting the VM must be below the warm threshold.

The energy can be saved by consolidating under-utilized servers, overdoing it may create hot spots in the future. The warm threshold is designed to prevent that. All things being equal, then select a destination server whose skewness can be reduced the most by accepting this VM. If the destination server is finding for all VMs on a cold spot, then record the sequence of migrations and update the predicted load of related servers. Otherwise, do not migrate any of its VMs. Restrict the number of cold spots that can be eliminated in each run of the algorithm to be no more than a certain percentage of active servers in the system. This is called the consolidation limit.

D. Consolidated Movements

The movements generated in each step above are not executed until all steps have finished. The list of movements are then consolidated so that each VM is moved at most once to its final destination. For example, hot spot mitigation may dictate a VM to move from PM A to PM B, while green computing dictates it to move from PM B to PM C. In the actual execution, the VM is moved from A to C directly.

VI. ANALYSIS OF SKEWNESS ALGORITHM

The skewness algorithm consists of three parts: load prediction, hot spot mitigation, and green computing. Let n and m be the number of PMs and VMs respectively. The number of resources (CPU, memory, I/O, etc.) that need to be considered is usually a small constant. Thus the computation of the skewness and the temperature metrics for a single server takes a constant amount of time. During load prediction, The time complexity is $O(n+m)$.

A. Complexity of Hot Spot Mitigation

For hot spot mitigation, let n_h be the number of hot spots in the system during a decision. Sorting them based on their temperature takes $O(n_h \log(n_h))$. Hence, the sorting takes a constant amount of time. For each VM, need to scan the rest of the PMs to find a suitable destination for it, which takes $O(n)$. The overall complexity of this phase is thus $O(n_h * n)$.

B. Complexity of Green Computing

For the green computing phase, let n_c be the number of cold spots in the system during a decision run. Sorting them based on their memory sizes takes $O(n_c \log(n_c))$. For each VM in a cold spot, it takes

$O(n)$ time to find a destination server for it. The overall complexity of this phase is $O(n_c * n)$.

VII. CONCLUSION

Dynamic migrations of virtual machines are becoming an interesting opportunity to allow cloud infrastructures to accommodate changing demands for different types of processing with heterogeneous workloads and time constraints. The proposed system multiplexes virtual to physical resources adaptively based on the changing demand.

The skewness algorithm metric is used to combine VMs with different resource characteristics appropriately so that the capacities of servers are well utilized. The algorithm achieves both overload avoidance and green computing for systems with multi resource constraints. This improves the performance of the Cloud Data Centers.

The amount of data to be transmitted is reduced, the effect is that of increasing the capacity of the communication channel bandwidth. Security is the important factor that everyone thinks before choosing a cloud provider. The security issues can be solved by providing security for the user resources while allocating resources dynamically. This will improve the performance and security of Data centers.

REFERENCES

- [1] Zhen Xiao, Senior Member, IEEE, Weijia Song, and Qi Chen-“Dynamic Resource Allocation Using Virtual Machines for Cloud Computing Environment”- IEEE Transaction on parallel and Distributed System, vol. 24, no. 6, June 2013.
- [2] C. Clark, K. Fraser, S. Hand, J. G. Hansen, E. Jul, C. Limpach, I. Pratt, and A. Warfield, “Live migration of virtual machines,” in Proc. of the Symposium on Networked Systems Design and Implementation (NSDI’05), May 2005.
- [3] G. Chen, H. Wenbo, J. Liu, S. Nath, L. Rigas, L. Xiao, and F. Zhao, “Energy aware server provisioning and load dispatching for connection-intensive internet services,” in Proc. of the USENIX Symposium on Networked Systems Design and Implementation (NSDI’08), Apr. 2008.
- [4] P. Padala, K.-Y. Hou, K. G. Shin, X. Zhu, M. Uysal, Z. Wang, S. Singhal, and A. Merchant, “Automated control of multiple virtualized resources,” in Proc. of the ACM European conference on Computer systems (EuroSys’09), 2009.
- [5] N. Bobroff, A. Kochut, and K. Beaty, “Dynamic placement of virtual machines for managing sla violations,” in Proc. of the

- IFIP/IEEE International Symposium on Integrated Network Management (IM'07), 2007.
- [6] J. S. Chase, D. C. Anderson, P. N. Thakar, A. M. Vahdat, and R. P. Doyle, "Managing energy and server resources in hosting centers," in Proc. Of the ACM Symposium on Operating System Principles (SOSP'01), Oct. 2001.
- [7] C. Tang, M. Steinder, M. Spreitzer, and G. Pacifici, "A scalable application placement controller for enterprise data centers," in Proc. Of the International World Wide Web Conference (WWW'07), May 2007.
- [8] A. Singh, M. Korupolu, and D. Mohapatra, "Server-storage virtualization: integration and load balancing in data centers," in Proc. of the ACM/IEEE conference on Supercomputing, 2008.
- [9] T. Das, P. Padala, V. N. Padmanabhan, R. Ramjee, and K. G. Shin, "Litegreen: saving energy in networked desktops using virtualization," in Proc. of the USENIX Annual Technical Conference, 2010.
- [10] Y. Agarwal, S. Savage, and R. Gupta, "Sleepserver: a software-only approach for reducing the energy consumption of pcs within enterprise environments," in Proc. of the USENIX Annual Technical Conference, 2010.
- [11] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, and A. Warfield, "Xen and the Art of Virtualization," Proc. ACM Symp. Operating Systems Principles (SOSP '03), Oct. 2003.
- [12] T. Wood, P. Shenoy, A. Venkataramani, and M. Yousif, "Black-Box and Gray-Box Strategies for Virtual Machine Migration," Proc. Symp. Networked Systems Design and Implementation (NSDI '07), Apr. 2007.
- [13] N. Bila, E.d. Lara, K. Joshi, H.A. Lagar-Cavilla, M. Hiltunen, and M. Satyanarayanan, "Jettison: Efficient Idle Desktop Consolidation with Partial VM Migration," Proc. ACM European Conf. Computer Systems (EuroSys '12), 2012.
- [14] M. Nelson, B.-H. Lim, and G. Hutchins, "Fast Transparent Migration for Virtual Machines," Proc. USENIX Ann. Technical Conf., 2005.
- [15] V.Vinothina, R.Sridaran, P.Ganapathi "A survey on Resource Allocation Strategies in Cloud Computing" IJCSA Vol 3, No 6,2012.
- [16] Y. Toyoda, "A Simplified Algorithm for Obtaining Approximate Solutions to Zero-One Programming Problems," Management Science, vol. 21, pp. 1417-1427, Aug. 1975.
- [17] J.S. Chase, D.C. Anderson, P.N. Thakar, A.M. Vahdat, and R.P. Doyle, "Managing Energy and Server Resources in Hosting Centers," Proc. ACM Symp. Operating System Principles (SOSP '01), Oct. 2001.