

RESEARCH ARTICLE



ISSN: 2321-7758

SLA-BASED OPTIMIZATION OF POWER AND MIGRATION COST IN CLOUD COMPUTING

S. JAYANTHI*¹, Dr. P. SRINIVASA BABU²

¹M.E Scholar, Department of Computer Science and Engineering, Adhiyamaan College of Engineering, Hosur, India.

² Professor, Department of Computer Science and Engineering, Adhiyamaan College of Engineering, Hosur, India.

Article Received: 30/01/2015

Article Revised on:10/02/2015

Article Accepted on:14/02/2015



S. JAYANTHI

ABSTRACT

Cloud computing systems (e.g., hosting datacenters) have attracted a lot of attention in recent years. Utility computing, reliable data storage, and infrastructure-independent computing are example applications of such systems. Operational cost in these systems is highly dependent on the resource management algorithms used to assign virtual machines (VMs) to physical servers and possibly migrate them in case of power and thermal emergencies. Energy non-proportionality of IT devices in a datacenter, cooling system inefficiency, and power delivery network constraints should be considered by the resource management algorithms in order to minimize the energy cost as much as possible. Scalability of the resource assignment solution is one of the biggest concerns in designing these algorithms. This thesis examines the resource management problem in datacenters. First a centralized datacenter resource management is proposed, which considers service level agreements (SLAs) in VM placement in order to minimize the total operational cost of the datacenter. Second, a hierarchical SLA-based resource management structure is proposed, which considers the peak power constraints and cooling-related power consumption in addition to the scalability issue. The proposed hierarchical structure fits the hierarchical resource distribution in datacenters. The proposed structure is suitable to track and react to dynamic changes inside the datacenter to satisfy SLA constraints and avoid emergencies.

Key Words—SLA, Data centers, virtual machine, power, cloud computing, resource management, resource allocation.

©KY Publications

INTRODUCTION

Operational cost and admission control policy in the cloud computing system are affected by its power and VM management policies. Power management techniques control the average and/or peak power dissipation in datacenters in a distributed or centralized manner. VM management techniques

[75, 76, 77, 78, 20] control the VM placement in physical servers as well as VM migration from a server to another one. In this chapter, we focus on the SLA-based VM management to minimize the operational cost in a cloud computing system. Optimal provisioning of the resources is crucial in order to reduce the cost incurred on the datacenter

operators as well as minimize the environmental impact of datacenters. The problem of optimal resource provisioning is challenging due to the diversity present in the clients (applications) that are hosted as well as in SLAs. For example: some applications may be compute-intensive while others may be memory intensive, some applications may run well together while others do not, etc. In this chapter, we focus on online service applications in cloud computing systems. Our goal in this chapter is to minimize the total cost of the cloud computing system under performance-related constraints—in particular, upper bounds on the response times (service latencies) for serving clients' requests. The operational cost in the cloud computing system includes power and migration cost and the SLA violation penalty of serving clients. A lower bound on the total operational cost is presented, and the average effectiveness of the presented algorithm is demonstrated by comparing with previous works' algorithms and lower bound value. Content of this chapter is presented in reference [70].

The outline of this chapter is as follows. In section II, cloud computing system model is presented. The optimization problem, Simulation results and conclusions are given in the sections III and IV .

II. SYSTEM MODEL

An SLA-aware resource allocation method for a cloud computing system is presented to minimize the total operational cost of the system. The structure of the datacenter, the VM manager (VMM), as well as performance model and type of SLA used by the clients are Service level agreements are, by their nature, "output" based – the result of the service as received by the customer is the subject of the "agreement." The (expert) service provider can demonstrate their value by organizing themselves with ingenuity, capability, and knowledge to deliver the service required, perhaps in an innovative way. Organizations can also specify the way the service is to be delivered, through a specification (a service level specification) and using subordinate "objectives" other than those related to the level of service. This type of agreement is known as an "input" SLA. This latter type of requirement is becoming obsolete as organizations become more demanding and shift the delivery methodology risk on to the service provider.

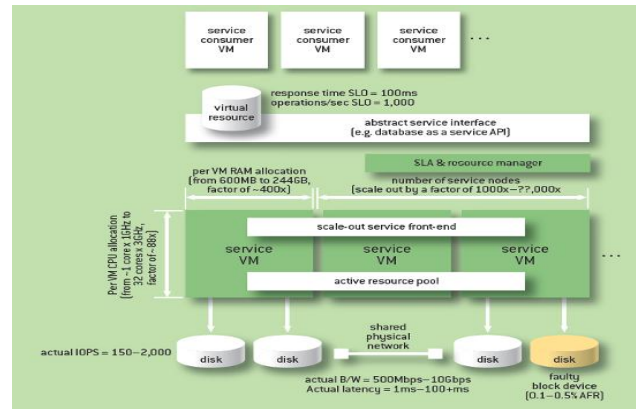


Figure 1. SLA System model

Service level agreements are also defined at different levels:

Customer-based SLA: An agreement with an individual customer group, covering all the services they use. For example, an SLA between a supplier (IT service provider) and the finance department of a large organization for the services such as finance system, payroll system, billing system, procurement/purchase system, etc.

Service-based SLA: An agreement for all customers using the services being delivered by the service provider. For example:

A car service station offers a routine service to all the customers and offers certain maintenance as a part of offer with the universal charging.

A mobile service provider offers a routine service to all the customers and offers certain maintenance as a part of offer with the universal charging

An email system for the entire organization. There are chances of difficulties arising in this type of SLA as level of the services being offered may vary for different customers (for example, head office staff may use high-speed LAN connections while local offices may have to use a lower speed leased line).

Multilevel SLA: The SLA is split into the different levels, each addressing different set of customers for the same services, in the same SLA.

Corporate-level SLA: Covering all the generic service level management (often abbreviated as SLM) issues appropriate to every customer throughout the organization. These issues are likely to be less volatile and so updates (SLA reviews) are less frequently required.

Customer-level SLA: covering all SLM issues relevant to the particular customer group, regardless of the services being used.

Service-level SLA: covering all SLM issue relevant to the specific services, in relation to this specific customer group.

2.1 Datacenter Configuration

We describe the type of the datacenter that we have assumed as well as our observations and key assumptions about where the performance bottlenecks are in the system and how we can account for the energy cost associated with a client’s VM running in a datacenter.

Data center transformation takes a step-by-step approach through integrated projects carried out over time. This differs from a traditional method of data center upgrades that takes a serial and siloed approach. The typical projects within a data center transformation initiative include standardization/consolidation, virtualization, automation and security.

Standardization/consolidation: The purpose of this project is to reduce the number of data centers a large organization may have. This project also helps to reduce the number of hardware, software platforms, tools and processes within a data center. Organizations replace aging data center equipment with newer ones that provide increased capacity and performance. Computing, networking and management platforms are standardized so they are easier to manage.

Virtualize: There is a trend to use IT virtualization technologies to replace or consolidate multiple data center equipment, such as servers. Virtualization helps to lower capital and operational expenses and reduce energy consumption. Virtualization technologies are also used to create virtual desktops, which can then be hosted in data centers and rented out on a subscription basis. Data released by investment bank Lazard Capital Markets reports that 48 percent of enterprise operations will be virtualized by 2012. Gartner views virtualization as a catalyst for modernization.

Automating: Data center automation involves automating tasks such as provisioning, configuration, patching, release management and compliance. As enterprises suffer from few skilled IT workers, automating tasks make data centers run more efficiently.

Securing: In modern data centers, the security of data on virtual systems is integrated with existing security of physical infrastructures. The security of a

modern data center must take into account physical security, network security, and data and user security.

2.2 VM Management System

Datacenter management is responsible for admitting the VMs into the datacenter, servicing them to satisfy SLAs, and minimizing the operational cost of the datacenter. We consider two main resource managers in the datacenter: VM manager (VMM) and power manager (PM). An exemplary architecture for the datacenter management system with emphasis on the VMM and per server PM is depicted in Figure 2.

Power manager is responsible for minimizing the average power consumption and satisfying the peak power constraints (thermal or peak power capacity limitation) subject to providing the required performance to VMs. Power management system in datacenter includes hierarchical power provisioners and a power manager for each server. Power provisioners distribute the peak power allowance between lower level power consumers and make sure that these power budget constraints are met. Servers are located at the lowest level of this hierarchy. Power manager in each server tries to minimize the average power consumption subject to satisfying the peak power constraint and performance requirements of the assigned VMs. This manager uses different dynamic power management techniques such as DVFS and clock throttling to minimize the power consumption.

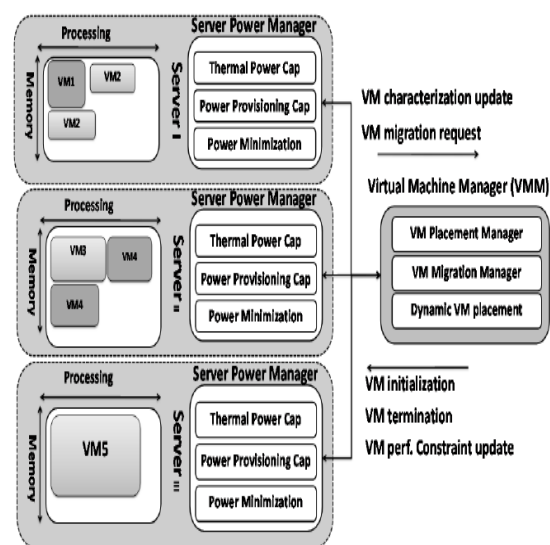


Figure 2. VM management structure in a datacenter

VMM is responsible for assigning VMs to servers, determining their performance requirements and migrating them if needed. VMM performs these tasks based on two optimization procedures: periodic and reactive. In contrast to periodic optimization procedure, reactive optimization procedure is performed when it is needed.

In the periodic optimization procedure, VMM considers the whole active set of VMs, the previous assignment solution, feedbacks generated from power, thermal and performance sensors, and workload prediction to generate the best VM placement solution for the next epoch.

The length of the epoch depends on the type and size of the datacenter and its workload. In reactive optimization procedure, VMM finds a temporary VM placement.

2.3 Performance Modeling

Performance of each client in the cloud computing system should be monitored and necessary decisions should be taken to satisfy SLA requirements. We focus on the online service applications that are sensitive to latency. A client in this system is application software that can produce a number of requests in each time unit. To model the response time of clients, we assume that the inter-arrival times of the requests for each client follow an exponential distribution function similar to the inter-arrival times of the requests in e-commerce applications [22]. Streams of requests generated by each client (application) may be decomposed into a number of different VMs. In case of more than one VM serving client, requests are assigned probabilistically portion of the incoming requests are forwarded to the server s (host of a VM) for execution, independently of the past or future forwarding decisions. Based on this assumption, the request arrival rate for each application in each server follows the Poisson distribution function.

There are different resources in the servers that are used by VMs such as processing units, memory, communication bandwidth, and secondary storage. These resources can be allocated to VMs by a fixed or round-robin scheduling policy. In this work, we consider the processing unit and memory to have fixed allocation policy whereas others are allocated by round-robin scheduling. Our algorithm determines the portion of processing unit and

memory allocated to each VM, which is assigned to a physical server. The amount of memory allocated to a VM does not significantly affect performance of the VM under different workloads as long as it is not less than a certain value [16].

However, these values can be changed in each server as a function of workload changes or power/performance optimization at the server. VMM considers the clients' workload to determine the resource allocation parameters to control the wait time of the processing queue for different applications based on SLA requirements.

A multi-class single server queue exists in servers that have more than one VM (from different clients). We consider *generalized processor sharing* (GPS) model at each queue; GPS model approximates the scheduling policy used by most operating systems, e.g., weighted fair queuing and the CPU time sharing in Linux. Using this scheduling policy, multi-class single server queue can be replaced by multiple single-server queues.

2.4 Initial Solution

To find an initial solution for P1, a constructive approach is used to assign clients to servers and allocate resources to them based on the assignment solution in the previous epoch. For this purpose, clients are divided into four groups. Clients that were served in the previous epoch are placed in one of the first three groups. The first group includes clients that leave the datacenter in the new epoch. The second group includes clients whose request arrival rates drop in the new epoch and the third group includes clients whose request arrival rates rise in the new epoch. Finally, the fourth group includes clients that were not served in the previous epoch.

Clients within these groups are picked in the order of their average minimum processing requirement for VMs (biggest VM first) but the groups are processed in increasing order of their IDs. For clients in the first group, VMM releases their resources and updates the resource availabilities. Resource availability in each server is defined as the amount of processing and memory allocated to the existing VMs.

2.5 Turn OFF under-utilized servers

To decrease the total cost in the system, it may be possible to turn off some of the under-utilized servers (after finding the initial solution) to reduce

the idle energy cost of the servers at the expense of more migration cost (for clients that were assigned to these under-utilized servers in the previous epoch) or more SLA violation penalty.

An iterative method is presented to find the minimum cost solution based on the results of the previous steps. In each iteration, a server with utilization less than a threshold (e.g., 20%) is chosen and its VMs are removed. To assign the removed VMs to 50 other servers, DPRA method is used. Considering the high energy cost for inactive servers, the DPRA method encourages the VMM to choose more SLA violation penalty or pay for the migration cost instead of turning on a server. Note that these iterations do not always decrease the total cost in the system; therefore, the global lowest total cost is compared to the total cost after turning off a server, and the move is rejected if it is not beneficial.

This iterative method is continued until all servers with low utilization have been examined.

III. SIMULATION RESULTS

To evaluate the effectiveness of the presented VM placement algorithm, a simulation framework is implemented. Simulation setups, baseline heuristics and numerical results of this implementation are explained next.

3.1 Simulation Setup

For simulations, model parameters are chosen based on true-to-life cloud computing systems. The number of server types is set to 10. For each server type, an arbitrary number of servers are placed in datacenter. Processors in server types are selected from a set of Intel processors (e.g. Atom and Xeon) [80] with different number of cores, cache, power consumptions and working frequencies. Active power consumptions for different server types (excluding processor power consumption) are set to vary uniformly between three to six times the power consumption of their fully-utilized processor. Memory capacities of the servers are selected based on cache size of the processors with a constant scaling factor of 1,500. Energy cost is assumed to be 15 cents per KWhr at all times. Request arrival rates of the client are chosen uniformly between 0.1 and 1 request per second. The memory requirements for clients are also selected uniformly between 256MB and 4GB. These parameters are borrowed from the simulation setup of [27].

In each simulation, five different client classes are considered. Each client is randomly picked from one of the client classes. The amount of penalty for different client classes is selected based on the on-demand rates of Amazon EC2 cloud service [81]. Migration costs are set to be equal to downtime penalty of 65ms for each client. In addition, set based on the highest clock frequency for the servers.

Each simulation is repeated at least 1000 times to generate acceptable average results for each case. In each simulation, a number of clients are assigned to the servers for the first decision epoch. At the end of each epoch, an arbitrary number of clients leave the datacenter while an arbitrary number of clients join the datacenter. Less than 10% of current clients join or leave the datacenter at the beginning of each epoch. Moreover, inter-arrival rate of the remaining clients in the system are chosen uniformly between 0.1 and 1 request per second for the next epoch. To account for the physical infrastructure overhead, energy cost of the servers in the datacenter is multiplied by a factor of 1.3 as a typical power usage effectiveness of current datacenters [11].

3.2 Heuristics for Comparison

We implemented a slightly modified version of the FFD [79] for VM placement, called FFDP, and PMaP heuristic [7] as baseline. These approaches consider VMs that 52 have fixed processing size. The expected violation rate of the SLA response time constraints for each client. From (23) the amount of processing units required for different VMs on different physical servers were calculated.

The FFDP method picks clients based on the size of their VM (highest to lowest) and assigns them to the first server with available resources from the server type that has the lowest execution time for the client's requests. The PMaP method is a VM placement heuristic that tries to minimize the power and migration cost. PMaP computes the amount of resources that VMs need, determines the active servers and place the VMs on the servers. After these steps, a power and migration-aware local search is done to find the final solution. Details of PMaP may be found in [7].

CONCLUSION

In this chapter we presented a centralized VM placement to minimize the power and migration cost in a cloud system. Soft SLA constraints on

response time were considered for the clients in this system. We presented an algorithm based on convex solutions. Based on the results of this chapter, it can be seen that considering SLA with effective VM placement can help to minimize the operational cost in the cloud computing system

REFERENCES

- [1]. J. Koomey, "Growth in data center electricity use 2005 to 2010," Analytics Press, 2011.
- [2]. ENERGY STAR, "Report to Congress on Server and Datacenter Energy Efficiency Public Law 109-431," U.S.Environmental Protection Agency, Washington, D.C., 2007.
- [3]. D. Meisner, B. Gold and T. Wenischn, "PowerNap: eliminating server idle power," in *Proceedings of the ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, Washington, DC, 2009.
- [4]. S. Pelley, D. Meisner, T. F. Wenischn and J. VanGilder, "Understanding and abstracting total datacenter power," in *workshop on Energy-Efficient Design*, 2009.
- [5]. "EPA confenrece on Enterprise Servers and Datacenters: Opportunities for Energy Efficiency," EPA, Lawrence Berkeley National Laboratory, 2006.
- [6]. 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," in *19th ACM Symposium on Operating Systems Principles*, 2003.
- [7]. A. Verrna, P. Ahuja and A. Neogi, "pMapper: Power and migration cost aware application placement in virtualized systems," in *ACM/IFIP/USENIX 9th International Middleware Conference*, 2008.
- [8]. M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica and M. Zaharia, "A view of cloud computing," *Commun ACM*, vol. 53, no. 4, pp. 50-58, 2010.
- [9]. R. Buyya, "Market-oriented cloud computing: Vision, hype, and reality of delivering computing as the 5th utility," in *9th IEEE/ACM International Symposium on Cluster Computing and the Grid, CCGRID*, 2009.
- [10]. S. Ghemawat, H. Gobioff and S.-T. Leung, "The Google file system," in *The 19th ACM Symposium on Operating Systems Principles*, Lake George, NY, 2003.
- [11]. L. A. Barroso and U. Holzle, *The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines*, Morgan & Claypool Publishers, 2009.
- [12]. C. Belady, "Green Grid Datacenter Power Efficiency Metrics: PUE and DCiE," [Online]. Available: Available at http://www.thegreengrid.org/gg_content/TGG_Data_Center_Power_Efficiency_Metrics_PUE_and DCiE.pdf.
- [13]. L. A. Barroso and U.Hölzle, "The Case for Energy-Proportional Computing," *IEEE Computer*, vol. 40, 2007.
- [14]. X. Fan, W. Weber and L. A. Barroso, "Power provisioning for a warehouse-sized computer," in *in Proceedings of the 34th Annual International symposium on Computer Architecture*, San Diego, CA, 2007.
- [15]. A. Karve, T. Kimbre, G. Pacifici, M. Spreitzer, M. Steinder, M. Sviridenko and A. Tantawi, "Dynamic placement for clustered web applications," in *15th International Conference on World Wide Web, WWW'06*, 2006.
- [16]. C. Tang, M. Steinder, M. Spreitzer and G. Pacifici, "A scalable application placement controller for enterprise datacenters," in *16th International World Wide Web Conference, WWW2007*, 2007.
- [17]. F. Chang, J. Ren and R. Viswanathan, "Optimal resource allocation in clouds," in *3rd IEEE International Conference on Cloud Computing, CLOUD 2010*, 2010.
- [18]. 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 *18th ACM Symposium on Operating Systems Principles (SOSP'01)*, 2001.
- [19]. E. Pakbaznia, M. GhasemAzar and M. Pedram, "Minimizing datacenter cooling

- and server power costs," in *Proc. of Design Automation and Test in Europe*, 2010.
- [20]. S. Srikantaiah, A. Kansal and F. Zhao, "Energy aware consolidation for cloud computing," in *Conference on Power aware computing and systems (HotPower'08)*, 2008.
- [21]. B. Urgaonkar, P. Shenoy and T. Roscoe, "Resource Overbooking and Application Profiling in Shared Hosting Platforms," in *Symposium on Operating Systems Design and Implementation*, 2002.
- [22]. Z. Liu, M. S. Squillante and J. L. Wolf, "On maximizing service-level-agreement profits," in *Third ACM Conference on Electronic Commerce*, 2001.
- [23]. K. Le, R. Bianchini, T. D. Nguyen, O. Bilgir and M. Martonosi, "Capping the brown energy consumption of internet services at low cost," in *International Conference on Green Computing (Green Comp)*, 2010.
- [24]. L. Zhang and D. Ardagna, "SLA based profit optimization in autonomic computing systems," in *Proceedings of the Second International Conference on Service Oriented Computing*, 2004.
- [25]. D. Ardagna, M. Trubian and L. Zhang, "SLA based resource allocation policies in autonomic environments," *Journal of Parallel and Distributed Computing*, vol. 67, no. 3, pp. 259-270, 2007.
- [26]. H. Goudarzi and M. Pedram, "Maximizing profit in the cloud computing system via resource allocation," in *Proc. of international workshop on Datacenter Performance*, 2011.
- [27]. D. Ardagna, B. Panicucci, M. Trubian and L. Zhang, "Energy-Aware Autonomic Resource Allocation in Multi-Tier Virtualized Environments," *IEEE Transactions on Services Computing*, vol. 99, 2010.
- [28]. H. Goudarzi and M. Pedram, "Multi-dimensional SLA- based resource allocation for multi-tier cloud computing systems," in *proceeding of 4th IEEE conference on cloud computing (Cloud 2011)*, 2011.
- [29]. G. Tesauro, N. K. Jong, R. Das and M. N. Bennani, "A hybrid reinforcement learning approach to autonomic resource allocation," in *Proceedings of International Conference on Autonomic Computing (ICAC '06)*, 2006.
- [30]. D. Kusic, J. O. Kephart, J. E. Hanson, N. Kandasamy and G. Jiang, "Power and performance management of virtualized computing environments via lookahead control," in *Proceedings of International Conference on Autonomic Computing (ICAC '08)*, 2008.
- [31]. A. Chandra, W. Gongt and P. Shenoy, "Dynamic resource allocation for shared datacenters using online measurements," in *International Conference on Measurement and Modeling of Computer Systems ACM SIGMETRICS*, 2003.
- [32]. N. Bobroff, A. Kochut and K. Beaty, "Dynamic Placement of Virtual Machines for Managing SLA Violations," in *Proceedings of the 10th IFIP/IEEE International Symposium on Integrated Management (IM2007)*, 2007.
- [33]. M. N. Bennani and D. A. Menasce, "Resource allocation for autonomic datacenters using analytic performance models," in *Second International Conference on Autonomic Computing*, 2005.
- [34]. B. Urgaonkar, G. Pacifici, P. Shenoy, M. Spreitzer and A. Tantawi, "An analytical model for multi-tier internet services and its applications," in *SIGMETRICS 2005: International Conference on Measurement and Modeling of Computer Systems*, 2005.
- [35]. M. Pedram and I.Hwang, "Power and performance modeling in a virtualized server system," in *39th International Conference on Parallel Processing workshops (ICPPW)*, 2010.
- [36]. R. Raghavendra, P. Ranganathan, V. Talwar, Z. Wang and X. Zhu, "No "power" struggles: Coordinated multi-level power management for the datacenter," *ACM SIGPLAN Notices*, vol. 43, no. 3, pp. 48-59, 2008.
- [37]. S. Pelley, D. Meisner, P. Zandevakili, T. F. Wenisch and J. Underwood, "Power Routing : Dynamic Power Provisioning in the Datacenter," in *ASPLOS '10*:

- Architectural Support for Programming Languages and Operating Systems*, 2010.
- [38]. M. Srivastava, A. Chandrakasan and R. Brodersen, "Predictive system shutdown and other architectural techniques for energy efficient programmable computation," *IEEE Trans. on VLSI*, 1996.
- [39]. Q. Qiu and M. Pedram, "Dynamic Power Management Based on Continuous-Time Markov Decision Processes," in *ACM design automation conference (DAC'99)*, 1999.
- [40]. G. Dhiman and T. S. Rosing, "Dynamic power management using machine learning," in *ICCAD '06*, 2006.
- [41]. D. Meisner, C. Sadler, L. Barroso, W. Weber and T. Wenisch, "Power Management of Online Data-Intensive Services," in *Proceedings of the 38th Annual International symposium on Computer Architecture*, 2011.
- [42]. Y. Chen, A. Das, W. Qin, A. Sivasubramaniam, Q. Wang and N. Gautam, "Managing server energy and operational costs in hosting centers," in *ACM SIGMETRICS '05*, 2005.
- [43]. X. Wang and Y. Wang, "Co-con: Coordinated control of power and application performance for virtualized server clusters," in *IEEE 17th International workshop on Quality of Service (IWQoS)*, 2009.
- [44]. E. Elnozahy, M. Kistler and R. Rajamony, "Energy-Efficient Server Clusters," in *Proc. 2nd workshop Power-Aware Computing Systems*, 2003.
- [45]. R. Buyya and A. Beloglazov, "Energy efficient resource management in virtualized cloud datacenters," in *10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing (CCGrid)*, 2010.
- [46]. L. Liu, H. Wang, X. Liu, X. Jin, W. He, Q. Wang and Y. Chen, "Greencloud: A new architecture for green datacenter," in *6th International Conference Industry Session on Autonomic Computing and Communications Industry Session, ICAC-INDST'09*, 2009.
- [47]. R. Nathuji and K. Schwan, "VirtualPower: Coordinated power management in virtualized enterprise systems," *Operating Systems Review*, vol. 41, no. 6, pp. 265-278, 2007.
- [48]. N. Rasmussen, "Calculating Total Cooling Requirements for Datacenters," *American Power Conversion*, 2007.
- [49]. E. Pakbaznia and M. Pedram, "Minimizing datacenter cooling and server power costs," in *Proceedings of the International Symposium on Low Power Electronics and Design*, 2009.
- [50]. R. Sharma, C. Bash, C. Patel, R. Friedrich and J. Chase, "Balance of power: dynamic thermal management for Internet datacenters," *IEEE Internet Computing*,
- [51]. J. Moore, J. Chase, P. Ranganathan and R. Sharma, "Making scheduling "cool": temperature-aware workload placement in datacenters," in *Proceedings of the Annual Conference on USENIX Annual Technical Conference*, 2005.
- [52]. Q. Tang, S. Gupta and G. Varsamopoulos, "Thermal-Aware Task Scheduling for Datacenters through Minimizing Heat Recirculation," in *Proc. IEEE Cluster*, 2007.
- [53]. Q. Tang, S. Gupta and G. Varsamopoulos, "Energy-Efficient Thermal-Aware Task Scheduling for Homogeneous High-Performance Computing Datacenters: A Cyber-Physical Approach," *IEEE Transactions on Parallel and Distributed Systems*, 2008.
- [54]. S. Biswas, M. Tiwari, T. Sherwood, L. Theogarajan and F. T. Chong, "Fighting fire with fire: modeling the datacenter-scale effects of targeted superlattice thermal management," in *Proceedings of the 38th Annual International symposium on Computer Architecture*, 2011.
- [55]. C. Patel, R. Sharma, C. Bash and A. Beitelmal, "Thermal considerations in cooling large scale high compute density datacenters," in *Proceedings of the Eighth Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems*, 2002.

- [56]. J. Choi, Y. Kim, A. Sivasubramaniam, J. Srebric, Q. Wang and J. Lee, "Modeling and Managing Thermal Profiles of Rack-mounted Servers with ThermoStat," in *Proceedings of International Symposium on High Performance Computer Architecture*, 2007.
- [57]. A. Ipakchi and F. Albuyeh, "Grid of the future," *IEEE Power and Energy Magazine*, vol. 7, no. 2, pp. 52-62, 2009.
- [58]. "<http://www.google.com/green/energy/>," [Online].
- [59]. R. Miller, "Facebook installs solar panels at new data center," *DatacenterKnowledge*, 16 April 2011. [Online].
- [60]. L. Rao, X. Liu, L. Xie and W. Liu, "Minimizing electricity cost: Optimization of distributed internet data centers in a multi-electricity market environment," in *IEEE INFOCOM*, 2010.
- [61]. R. Stanojevic and R. Shorten, "Distributed dynamic speed scaling," in *IEEE INFOCOM*, 2010.
- [62]. X. Wang and M. Chen, "Cluster-level feedback power control for performance optimization," in *IEEE HPCA*, 2008.
- [63]. M. Lin, Z. Liu, A. Wierman and L. L. Andrew, "Online algorithms for geographical load balancing," in *Proc. Int. Green Computing Conf.*, San Jose, CA, 2012.
- [64]. Z. Liu, M. Lin, A. Wierman, S. H. Low and L. L. H. Andrew, "Geographical load balancing with renewables," in *Proc. ACM GreenMetrics*, 2011.
- [65]. Z. Liu, M. Lin, A. Wierman, S. H. Low and L. L. H. Andrew, "Greening geographical load balancing," in *Proc. ACM SIGMETRICS*, San Jose, CA, 2011.
- [66]. K. Le, R. Bianchini, M. Martonosi and T. D. Nguyen, "Cost and energy-aware load distribution across data centers," in *HotPower'09*, Big Sky, MT, 2009.
- [67]. M. A. Adnan, R. Sugihara and R. Gupta, "Energy Efficient Geographical Load Balancing via Dynamic Deferral of Workload," in *proceeding of 5th IEEE conference on cloud computing (Cloud 2012)*, Honolulu, HI, 2012.
- [68]. Z. Liu, Y. Chen, C. Bash, A. Wierman, D. Gmach, Z. Wang and C. Hyser, "Renewable and cooling aware workload management for sustainable data centers," *ACM SIGMETRICS Performance Evaluation Review*, vol. 40, no. 1, 2012.
- [69]. H. Goudarzi and M. Pedram, "Profit-maximizing resource allocation for multi-tier cloud computing systems under service level agreements," in *Large Scale Network-Centric Distributed Systems*, Wiley-IEEE Computer Society Press, 2013.
- [70]. H. Goudarzi, M. Ghasemazar and M. Pedram, "SLA-ased Optimization of Power and Migration Cost in Cloud Computing," in *12th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing (CCGrid)*, 2012.
- [71]. H. Goudarzi and M. Pedram, "Hierarchical SLA-Driven Resource Management for Peak Power-Aware and Energy-Efficient Operation of a Cloud Datacenter," *Submitted to IEEE transaction on computers*.
- [72]. H. Goudarzi and M. Pedram, "Energy-efficient Virtual Machine Replication and Placement in a Cloud Computing System," in *EEE international conference on cloud computing (CLOUD 2012)*, Honolulu, 2012.
- [73]. H. Goudarzi and M. Pedram, "Geographical Load Balancing for Online Service Applications in Distributed Datacenters," in *IEEE international conference on cloud computing (CLOUD 2013)*, Santa Clara, 2013.
- [74]. H. Goudarzi and M. Pedram, "Force-directed Geographical Load Balancing and Scheduling for Batch Jobs in Distributed Datacenters," in *IEEE international conference on cluster computing (CLUSTER 2013)*, Indianapolis, 2013.
- [75]. A. Beloglazov and R. Buyya, "Energy efficient resource management in virtualized cloud datacenters," in *Proceeding of 2010 10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing (CCGrid)*, 2010.
- [76]. I. Goiri, J. Fitó, F. Julià, R. Nou, J. Berral, J. Guitart and J. Torres, "Multifaceted

- resource management for dealing with heterogeneous workloads in virtualized data centers," in *Proceeding of IEEE/ACM International Conference on Grid Computing (GRID)*, 2010.
- [77]. B. Sotomayor, R. S. Montero, I. M. Llorente and I. Foster, "Capacity Leasing in Cloud Systems using the OpenNebula Engine," in *Workshop on Cloud Computing and its Applications*, 2008.
- [78]. R. Buyya, Y. S. Chee and S. Venugopal, "Market-Oriented Cloud Computing: Vision, Hype and Reality for Delivering IT Services as Computing Utilities," in *IEEE International Conference on High Performance Computing and Communications*, 2008.
- [79]. S. Martello and P. Toth, *Knapsack Problems: Algorithms and Computer Implementations*, Wiley, 1990.
- [80]. "<http://ark.intel.com/>," [Online].
-