

Optimal Multiserver Configuration for Profit Maximization in Cloud Computing

Junwei Cao, *Senior Member, IEEE*, Kai Hwang, *Fellow, IEEE*,
Keqin Li, *Senior Member, IEEE*, and Albert Y. Zomaya, *Fellow, IEEE*

Abstract—As cloud computing becomes more and more popular, understanding the economics of cloud computing becomes critically important. To maximize the profit, a service provider should understand both service charges and business costs, and how they are determined by the characteristics of the applications and the configuration of a multiserver system. The problem of optimal multiserver configuration for profit maximization in a cloud computing environment is studied. Our pricing model takes such factors into considerations as the amount of a service, the workload of an application environment, the configuration of a multiserver system, the service-level agreement, the satisfaction of a consumer, the quality of a service, the penalty of a low-quality service, the cost of renting, the cost of energy consumption, and a service provider's margin and profit. Our approach is to treat a multiserver system as an M/M/m queueing model, such that our optimization problem can be formulated and solved analytically. Two server speed and power consumption models are considered, namely, the idle-speed model and the constant-speed model. The probability density function of the waiting time of a newly arrived service request is derived. The expected service charge to a service request is calculated. The expected net business gain in one unit of time is obtained. Numerical calculations of the optimal server size and the optimal server speed are demonstrated.

Index Terms—Cloud computing, multiserver system, pricing model, profit, queueing model, response time, server configuration, service charge, service-level agreement, waiting time

1 INTRODUCTION

CLOUD computing is quickly becoming an effective and efficient way of computing resources and computing services consolidation [10]. By centralized management of resources and services, cloud computing delivers hosted services over the Internet, such that accesses to shared hardware, software, databases, information, and all resources are provided to consumers on-demand. Cloud computing is able to provide the most cost-effective and energy-efficient way of computing resources management and computing services provision. Cloud computing turns information technology into ordinary commodities and utilities by using the pay-per-use pricing model [3], [5], [18]. However, cloud computing will never be free [8], and understanding the economics of cloud computing becomes critically important.

One attractive cloud computing environment is a three-tier structure [15], which consists of infrastructure vendors, service providers, and consumers. The three parties are also called cluster nodes, cluster managers, and consumers

in cluster computing systems [21], and resource providers, service providers, and clients in grid computing systems [19]. An infrastructure vendor maintains basic hardware and software facilities. A service provider rents resources from the infrastructure vendors, builds appropriate multiserver systems, and provides various services to users. A consumer submits a service request to a service provider, receives the desired result from the service provider with certain service-level agreement, and pays for the service based on the amount of the service and the quality of the service. A service provider can build different multiserver systems for different application domains, such that service requests of different nature are sent to different multiserver systems. Each multiserver system contains multiple servers, and such a multiserver system can be devoted to serve one type of service requests and applications. An application domain is characterized by two basic features, i.e., the workload of an application environment and the expected amount of a service. The configuration of a multiserver system is characterized by two basic features, i.e., the size of the multiserver system (the number of servers) and the speed of the multiserver system (execution speed of the servers).

Like all business, the pricing model of a service provider in cloud computing is based on two components, namely, the income and the cost. For a service provider, the income (i.e., the revenue) is the service charge to users, and the cost is the renting cost plus the utility cost paid to infrastructure vendors. A pricing model in cloud computing includes many considerations, such as the amount of a service (the requirement of a service), the workload of an application environment, the configuration (the size and the speed) of a multiserver system, the service-level agreement, the satisfaction of a consumer (the expected service time), the quality of a service (the task waiting time and the task response time),

- J. Cao is with the Research Institute of Information Technology, Tsinghua National Laboratory for Information Science and Technology, Tsinghua University, Beijing 100084, China. E-mail: jcao@tsinghua.edu.cn.
- K. Hwang is with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA 90089. E-mail: kaihawang@usc.edu.
- K. Li is with the Department of Computer Science, State University of New York, New Paltz, New York 12561. E-mail: lik@newpaltz.edu.
- A.Y. Zomaya is with the School of Information Technologies, University of Sydney, Sydney, NSW 2006, Australia. E-mail: albert.zomaya@sydney.edu.au.

Manuscript received 23 Feb. 2012; revised 18 June 2012; accepted 28 June 2012; published online 25 June 2012.

Recommended for acceptance by V.B. Misis, R. Buyya, D. Milojicic, and Y. Cui.

For information on obtaining reprints of this article, please send e-mail to: tpds@computer.org, and reference IEEECS Log Number TPDSSI-2012-02-0137.

Digital Object Identifier no. 10.1109/TPDS.2012.203.

the penalty of a low-quality service, the cost of renting, the cost of energy consumption, and a service provider's margin and profit. The profit (i.e., the net business gain) is the income minus the cost. To maximize the profit, a service provider should understand both service charges and business costs, and in particular, how they are determined by the characteristics of the applications and the configuration of a multiserver system.

The service charge to a service request is determined by two factors, i.e., the expected length of the service and the actual length of the service. The expected length of a service (i.e., the expected service time) is the execution time of an application on a standard server with a baseline or reference speed. Once the baseline speed is set, the expected length of a service is determined by a service request itself, i.e., the service requirement (amount of service) measured by the number of instructions to be executed. The longer (shorter, respectively) the expected length of a service is, the more (less, respectively) the service charge is. The actual length of a service (i.e., the actual service time) is the actual execution time of an application. The actual length of a service depends on the size of a multiserver system, the speed of the servers (which may be faster or slower than the baseline speed), and the workload of the multiserver system. Notice that the actual service time is a random variable, which is determined by the task waiting time once a multiserver system is established.

There are many different service performance metrics in service-level agreements [2]. Our performance metric in this paper is the task response time (or the turn around time), i.e., the time taken to complete a task, which includes task waiting time and task execution time. The service-level agreement is the promised time to complete a service, which is a constant times the expected length of a service. If the actual length of a service is (or, a service request is completed) within the service-level agreement, the service will be fully charged. However, if the actual length of a service exceeds the service-level agreement, the service charge will be reduced. The longer (shorter, respectively) the actual length of a service is, the more (less, respectively) the reduction of the service charge is. In other words, there is penalty for a service provider to break a service-level agreement. If the actual service time exceeds certain limit (which is service request dependent), a service will be entirely free with no charge. Notice that the service charge of a service request is a random variable, and we are interested in its expectation.

The cost of a service provider includes two components, i.e., the renting cost and the utility cost. The renting cost is proportional to the size of a multiserver system, i.e., the number of servers. The utility cost is essentially the cost of energy consumption and is determined by both the size and the speed of a multiserver system. The faster (slower, respectively) the speed is, the more (less, respectively) the utility cost is. To calculate the cost of energy consumption, we need to establish certain server speed and power consumption models.

To increase the revenue of business, a service provider can construct and configure a multiserver system with many servers of high speed. Since the actual service time (i.e., the task response time) contains task waiting time and task

execution time, more servers reduce the waiting time and faster servers reduce both waiting time and execution time. Hence, a powerful multiserver system reduces the penalty of breaking a service-level agreement and increases the revenue. However, more servers (i.e., a larger multiserver system) increase the cost of facility renting from the infrastructure vendors and the cost of base power consumption. Furthermore, faster servers increase the cost of energy consumption. Such increased cost may counterweight the gain from penalty reduction. Therefore, for an application environment with specific workload which includes the task arrival rate and the average task execution requirement, a service provider needs to decide an optimal multiserver configuration (i.e., the size and the speed of a multiserver system), such that the expected profit is maximized.

In this paper, we study the problem of optimal multiserver configuration for profit maximization in a cloud computing environment. Our approach is to treat a multiserver system as an M/M/m queueing model, such that our optimization problem can be formulated and solved analytically. We consider two server speed and power consumption models, namely, the idle-speed model and the constant-speed model. Our main contributions are as follows. We derive the probability density function (pdf) of the waiting time of a newly arrived service request. This result is significant in its own right and is the base of our discussion. We calculate the expected service charge to a service request. Based on these results, we get the expected net business gain in one unit of time, and obtain the optimal server size and the optimal server speed numerically. To the best of our knowledge, there has been no similar investigation in the literature, although the method of optimal multicore server processor configuration has been employed for other purposes, such as managing the power and performance tradeoff [17].

One related research is user-centric and market-based and utility-driven resource management and task scheduling, which have been considered for cluster computing systems [7], [20], [21] and grid computing systems [4], [12], [19]. To compete and bid for shared computing resources through the use of economic mechanisms such as auctions, a user can specify the value (utility, yield) of a task, i.e., the reward (price, profit) of completing the task. A utility function, which measures the value and importance of a task as well as a user's tolerance to delay and sensitivity to quality of service, supports market-based bidding, negotiation, and admission control. By taking an economic approach to providing service-oriented and utility computing, a service provider allocates resources and schedules tasks in such a way that the total profit earned is maximized. Instead of traditional system-centric performance optimization such as minimizing the average task response time, the main concern in such computational economy is user-centric performance optimization, i.e., maximizing the total utility delivered to the users (i.e., the total user-perceived value).

2 A MULTISERVER MODEL

Throughout the paper, we use $P[e]$ to denote the probability of an event e . For a random variable x , we use $f_x(t)$ to represent the probability density function of x , and $F_x(t)$ to

represent the cumulative distribution function (cdf) of x , and \bar{x} to represent the expectation of x .

A cloud computing service provider serves users' service requests by using a multiserver system, which is constructed and maintained by an infrastructure vendor and rented by the service provider. The architecture detail of the multiserver system can be quite flexible. Examples are blade servers and blade centers where each server is a server blade [16], clusters of traditional servers where each server is an ordinary processor [7], [20], [21], and multicore server processors where each server is a single core [17]. We will simply call these blades/processors/cores as *servers*. Users (i.e., customers of a service provider) submit service requests (i.e., applications and tasks) to a service provider, and the service provider serves the requests (i.e., run the applications and perform the tasks) on a multiserver system.

Assume that a multiserver system S has m identical servers. In this paper, a multiserver system is treated as an M/M/ m queuing system which is elaborated as follows. There is a Poisson stream of service requests with arrival rate λ , i.e., the interarrival times are independent and identically distributed (i.i.d.) exponential random variables with mean $1/\lambda$. A multiserver system S maintains a queue with infinite capacity for waiting tasks when all the m servers are busy. The first-come-first-served (FCFS) queuing discipline is adopted. The task execution requirements (measured by the number of instructions to be executed) are i.i.d. exponential random variables r with mean \bar{r} . The m servers (i.e., blades/processors/cores) of S have identical execution speed s (measured by the number of instructions that can be executed in one unit of time). Hence, the task execution times on the servers of S are i.i.d. exponential random variables $x = r/s$ with mean $\bar{x} = \bar{r}/s$.

Notice that although an M/G/ m queuing system has been considered (see, e.g., [13]), the M/M/ m queuing model is the only model that accommodates an analytical and closed-form expression of the probability density function of the waiting time of a newly arrived service request.

Let $\mu = 1/\bar{x} = s/\bar{r}$ be the average service rate, i.e., the average number of service requests that can be finished by a server of S in one unit of time. The server utilization is $\rho = \lambda/m\mu = \lambda\bar{x}/m = \lambda/m \cdot \bar{r}/s$, which is the average percentage of time that a server of S is busy. Let p_k denote the probability that there are k service requests (waiting or being processed) in the M/M/ m queuing system for S . Then, we have ([14, p. 102])

$$p_k = \begin{cases} p_0 \frac{(m\rho)^k}{k!}, & k \leq m; \\ p_0 \frac{m^m \rho^k}{m!}, & k \geq m, \end{cases}$$

where

$$p_0 = \left(\sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} + \frac{(m\rho)^m}{m!} \cdot \frac{1}{1-\rho} \right)^{-1}.$$

The probability of queuing (i.e., the probability that a newly submitted service request must wait because all servers are busy) is

$$P_q = \sum_{k=m}^{\infty} p_k = \frac{p_m}{1-\rho} = p_0 \frac{(m\rho)^m}{m!} \cdot \frac{1}{1-\rho}.$$

The average number of service requests (in waiting or in execution) in S is

$$\bar{N} = \sum_{k=0}^{\infty} k p_k = m\rho + \frac{\rho}{1-\rho} P_q.$$

Applying Little's result, we get the average task response time as

$$\bar{T} = \frac{\bar{N}}{\lambda} = \bar{x} \left(1 + \frac{P_q}{m(1-\rho)} \right) = \bar{x} \left(1 + \frac{p_m}{m(1-\rho)^2} \right).$$

The average waiting time of a service request is

$$\bar{W} = \bar{T} - \bar{x} = \frac{p_m}{m(1-\rho)^2} \bar{x}.$$

The waiting time is the source of customer dissatisfaction. A service provider should keep the waiting time to a low level by providing enough servers and/or increasing server speed, and be willing to pay back to a customer in case the waiting time exceeds certain limit.

3 POWER CONSUMPTION MODELS

Power dissipation and circuit delay in digital CMOS circuits can be accurately modeled by simple equations, even for complex microprocessor circuits. CMOS circuits have dynamic, static, and short-circuit power dissipation; however, the dominant component in a well-designed circuit is dynamic power consumption P (i.e., the switching component of power), which is approximately $P = aCV^2f$, where a is an activity factor, C is the loading capacitance, V is the supply voltage, and f is the clock frequency [6]. In the ideal case, the supply voltage and the clock frequency are related in such a way that $V \propto f^\phi$ for some constant $\phi > 0$ [22]. The processor execution speed s is usually linearly proportional to the clock frequency, namely, $s \propto f$. For ease of discussion, we will assume that $V = bf^\phi$ and $s = cf$, where b and c are some constants. Hence, we know that power consumption is $P = aCV^2f = ab^2Cf^{2\phi+1} = (ab^2C/c^{2\phi+1})s^{2\phi+1} = \xi s^\alpha$, where $\xi = ab^2C/c^{2\phi+1}$ and $\alpha = 2\phi + 1$. For instance, by setting $b = 1.16$, $aC = 7.0$, $c = 1.0$, $\phi = 0.5$, $\alpha = 2\phi + 1 = 2.0$, and $\xi = ab^2C/c^\alpha = 9.4192$, the value of P calculated by the equation $P = aCV^2f = \xi s^\alpha$ is reasonably close to that in [11] for the Intel Pentium M processor.

We will consider two types of server speed and power consumption models. In the *idle-speed model*, a server runs at zero speed when there is no task to perform. Since the power for speed s is ξs^α , the average amount of energy consumed by a server in one unit of time is $\rho \xi s^\alpha = \frac{\lambda}{m} \bar{r} \xi s^{\alpha-1}$, where we notice that the speed of a server is zero when it is idle. The average amount of energy consumed by an m -server system S in one unit of time, i.e., the power supply to the multiserver system S , is $P = m\rho \xi s^\alpha = \lambda \bar{r} \xi s^{\alpha-1}$, where $m\rho = \lambda\bar{x}$ is the average number of busy servers in S . Since a server still consumes some amount of power P^* even when it is idle (assume that an idle server consumes certain base power P^* , which includes static power dissipation, short-circuit power dissipation, and other leakage and wasted power [1]), we will include P^* in P , i.e., $P = m(\rho \xi s^\alpha + P^*) = \lambda \bar{r} \xi s^{\alpha-1} + mP^*$. Notice that when $P^* = 0$, the above P is independent of m .

In the *constant-speed model*, all servers run at the speed s even if there is no task to perform. Again, we use P to represent the power allocated to multiserver system S . Since the power for speed s is ξs^α , the power allocated to multiserver system S is $P = m(\xi s^\alpha + P^*)$.

4 WAITING TIME DISTRIBUTION

Let W denote the waiting time of a new service request that arrives to a multiserver system. In this section, we find the pdf $f_W(t)$ of W . To this end, we consider W in different situations, depending on the number of tasks in the queuing system when a new service request arrives. Let W_k denote the waiting time of a new task that arrives to an M/M/m queuing system under the condition that there are k tasks in the queuing system when the task arrives.

We define a *unit impulse function* $u_z(t)$ as follows:

$$u_z(t) = \begin{cases} z, & 0 \leq t \leq \frac{1}{z}; \\ 0, & t > \frac{1}{z}. \end{cases}$$

The function $u_z(t)$ has the following property:

$$\int_0^\infty u_z(t) dt = 1,$$

namely, $u_z(t)$ can be treated as a pdf of a random variable with expectation

$$\int_0^\infty t u_z(t) dt = z \int_0^{1/z} t dt = \frac{1}{2z}.$$

Let $z \rightarrow \infty$ and define $u(t) = \lim_{z \rightarrow \infty} u_z(t)$. It is clear that any random variable whose pdf is $u(t)$ has expectation 0.

The following theorem gives the pdf of the waiting time of a newly arrived service request:

Theorem 1. *The pdf of the waiting time W of a newly arrived service request is*

$$f_W(t) = (1 - P_q)u(t) + m\mu p_m e^{-(1-\rho)m\mu t},$$

where $P_q = p_m/(1 - \rho)$ and $p_m = p_0(m\rho)^m/m!$.

Sketch of the Proof. Let W_k be the waiting time of a new service request if there are k tasks in the queuing system when the service request arrives. We find the pdf of W_k for all $k \geq 0$. Then, we have

$$f_W(t) = \sum_{k=0}^{\infty} p_k f_{W_k}(t).$$

Actually, W_k can be found for two cases, i.e., when $k < m$ and when $k \geq m$. A complete proof of the theorem is given in Section 9. \square

Notice that a multiserver system with multiple identical servers has been configured to serve requests from certain application domain. Therefore, we will only focus on task waiting time in a waiting queue and do not consider other sources of delay, such as resource allocation and provision, virtual machine instantiation and deployment, and other overhead in a complex cloud computing environment.

5 SERVICE CHARGE

If all the servers have a fixed speed s , the execution time of a service request with execution requirement r is known as $x = r/s$. The response time to the service request is $T = W + x = W + r/s$. The response time T is related to the service charge to a customer of a service provider in cloud computing.

To study the expected service charge to a customer, we need a complete specification of a service charge based on the amount of a service, the service-level agreement, the satisfaction of a consumer, the quality of a service, the penalty of a low-quality service, and a service provider's margin and profit.

Let s_0 be the baseline speed of a server. We define the *service charge function* for a service request with execution requirement r and response time T to be

$$C(r, T) = \begin{cases} ar, & \text{if } 0 \leq T \leq \frac{c}{s_0}r; \\ ar - d\left(T - \frac{c}{s_0}r\right), & \text{if } \frac{c}{s_0}r < T \leq \left(\frac{a}{d} + \frac{c}{s_0}\right)r; \\ 0, & \text{if } T > \left(\frac{a}{d} + \frac{c}{s_0}\right)r. \end{cases}$$

The above function is defined with the following rationals:

- If the response time T to process a service request is no longer than $(c/s_0)r = c(r/s_0)$ (i.e., a constant c times the task execution time with speed s_0), where the constant c is a parameter indicating the service-level agreement, and the constant s_0 is a parameter indicating the expectation and satisfaction of a consumer, then a service provider considers that the service request is processed successfully with high quality of service and charges a customer ar , which is linearly proportional to the task execution requirement r (i.e., the amount of service), where a is the service charge per unit amount of service (i.e., a service provider's margin and profit).
- If the response time T to process a service request is longer than $(c/s_0)r$ but no longer than $(a/d + c/s_0)r$, then a service provider considers that the service request is processed with low quality of service and the charge to a customer should decrease linearly as T increases. The parameter d indicates the degree of penalty of breaking the service-level agreement.
- If the response time T to process a service request is longer than $(a/d + c/s_0)r$, then a service provider considers that the service request has been waiting too long, so there is no charge and the service is free.

Notice that the task response time T is compared with the task execution time on a server with speed s_0 (i.e., the baseline or reference speed). The actual speed s of a server can be decided by a service provider, which can be either lower or higher than s_0 , depending on the workload (i.e., λ and \bar{r}) and system parameters (e.g., m , α , and P^*) and the service charge function (i.e., a , c , and d), such that the net business gain to be defined below is maximized.

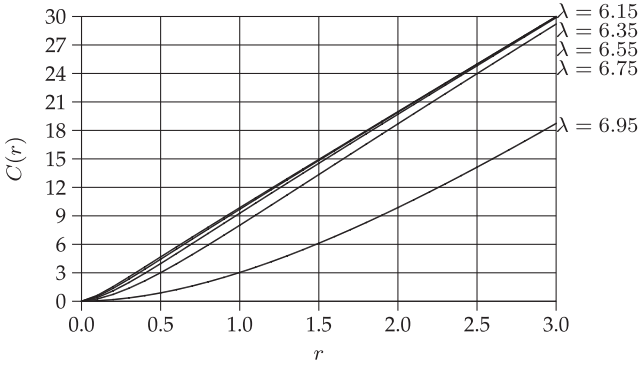


Fig. 1. Service charge $C(r)$ versus r and λ .

To build our discussion upon our earlier result on task waiting time, we notice that the service charge function can be rewritten equivalently in terms of r and W as

$$C(r, W) = \begin{cases} ar, & \text{if } 0 \leq W \leq \left(\frac{c}{s_0} - \frac{1}{s}\right)r; \\ \left(a + \frac{cd}{s_0} - \frac{d}{s}\right)r - dW, & \text{if } \left(\frac{c}{s_0} - \frac{1}{s}\right)r < W \leq \left(\frac{a}{d} + \frac{c}{s_0} - \frac{1}{s}\right)r; \\ 0, & \text{if } W > \left(\frac{a}{d} + \frac{c}{s_0} - \frac{1}{s}\right)r. \end{cases}$$

The following theorem gives the expected charge to a service request:

Theorem 2. *The expected charge to a service request is*

$$C = a\bar{r} \left(1 - \frac{P_q}{((ms - \lambda\bar{r})(c/s_0 - 1/s) + 1)} \times \frac{1}{((ms - \lambda\bar{r})(a/d + c/s_0 - 1/s) + 1)} \right),$$

where $P_q = p_m/(1 - \rho)$ and $p_m = p_0(m\rho)^m/m!$.

Sketch of the Proof. The proof is actually a detailed calculation of C , which contains two steps. In the first step, we calculate $C(r)$, i.e., the expected charge to a service request with execution requirement r , based on the pdf of W obtained from Theorem 1. In the second step, we calculate C based on the pdf of r . A complete proof of the theorem is given in Section 9. \square

In Fig. 1, we consider the expected charge to a service request with execution requirement r , i.e.,

$$C(r) = ar - \frac{dP_q}{(1 - \rho)m\mu} \left(e^{-(1-\rho)m\mu(c/s_0 - 1/s)r} - e^{-(1-\rho)m\mu(a/d + c/s_0 - 1/s)r} \right).$$

We assume that $\bar{r} = 1$ billion instructions, $m = 7$ servers, $s_0 = 1$ billion instructions per second, $s = 1$ billion instructions per second, $a = 10$ cents per one billion instructions (Note: The monetary unit “cent” in this paper may not be identical but should be linearly proportional to the real cent in US dollars.), $c = 3$, and $d = 1$ cents per second. (Note: These parameters are chosen only for illustration and should be scaled to any values.) For $\lambda = 6.15, 6.35, 6.55, 6.75, 6.95$

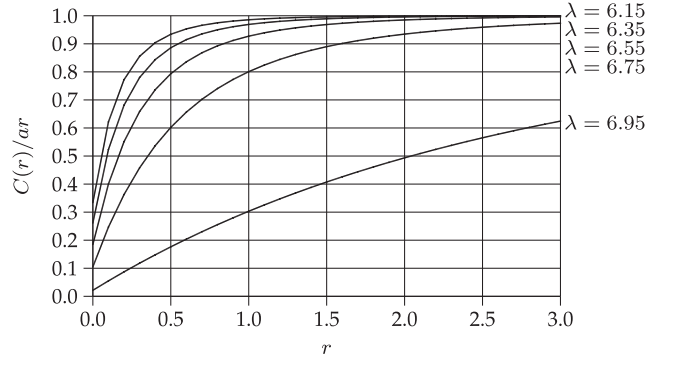


Fig. 2. Normalized service charge $C(r)/ar$ versus r and λ .

service requests per second, we show $C(r)$ for $0 \leq r \leq 3$. It can be seen that the service charge is a decreasing function of λ , since the waiting time and lateness penalty increase as λ increases. It can also be seen that the service charge is an increasing function of r , i.e., large service requests generate more revenue than small service requests.

In Fig. 2, we further display $C(r)/ar$ using the same parameters in Fig. 1. Since ar is the ideal (maximum) charge to a service request with execution requirement r , $C(r)/ar$ is considered as the *normalized service charge*. For $\lambda = 6.15, 6.35, 6.55, 6.75, 6.95$ service requests per second, we show $C(r)/ar$ for $0 \leq r \leq 3$. It can be seen that the normalized service charge is a decreasing function of λ , since the waiting time and lateness penalty increase as λ increases. It can also be seen that the normalized service charge is an increasing function of r , i.e., the percentage of lost service charge due to waiting time decreases as service requirement r increases. In other words, it is more likely to make profit from large service requests and it is more likely to give free services to small service requests. It can be verified that as r approaches 0, the normalized service charge is $\lim_{r \rightarrow 0} \frac{C(r)}{ar} = 1 - P_q$, where P_q increases (and $1 - P_q$ decreases) as λ increases. It can also be verified that as r approaches infinity, the normalized service charge is $\lim_{r \rightarrow \infty} \frac{C(r)}{ar} = 1$, for all λ .

6 NET BUSINESS GAIN

Since the number of service requests processed in one unit of time is λ in a stable M/M/m queuing system, the expected service charge in one unit of time is λC , which is actually the expected *revenue* of a service provider. Assume that the rental cost of one server for unit of time is β . Also, assume that the cost of energy is γ per Watt. The *cost* of a service provider is the sum of the cost of infrastructure renting and the cost of energy consumption, i.e., $\beta m + \gamma P$. Then, the expected *net business gain* (i.e., the net profit) of a service provider in one unit of time is $G = \lambda C - (\beta m + \gamma P)$, which is defined as the revenue minus the cost. The above equation is $G = \lambda C - (\beta m + \gamma(\lambda \bar{r} \xi s^{\alpha-1} + mP^*))$, for the idle-speed model, and $G = \lambda C - (\beta m + \gamma m(\xi s^{\alpha} + P^*))$, for the constant-speed model.

In Figs. 3 and 4, we demonstrate the revenue λC and the net business gain G in one unit of time as a function of λ for the two power consumption models, respectively, using the same parameters in Figs. 1 and 2. Furthermore, we assume

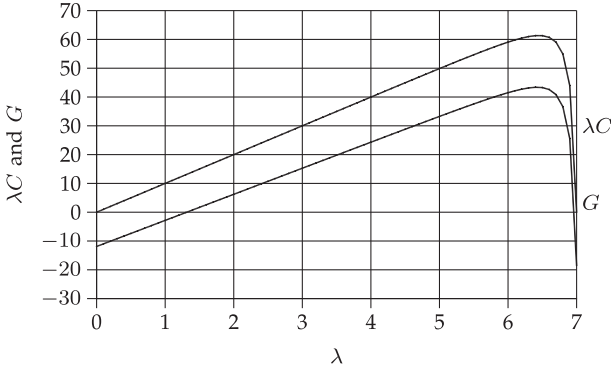


Fig. 3. Revenue and net business gain versus λ (idle-speed model).

that $P^* = 2$ Watts, $\alpha = 2.0$, $\xi = 9.4192$, $\beta = 1.5$ cents per second, and $\gamma = 0.1$ cents per Watt. For $0 \leq \lambda \leq 7$, we show λC and G . The cost of infrastructure renting is $\beta m = 14$ cents per second, and the cost of energy consumption is $0.5\lambda + 7$ cents per second for the idle-speed model and 10.5 cents per second for the constant-speed model. We observe that both λC and G increase with λ almost linearly and drop sharply after certain point. In other words, more service requests bring more revenue and net business gain; however, after the number of service requests per unit of time reaches certain point, the excessive waiting time causes increased lateness penalty, so that there is no revenue and negative business gain.

There are two situations that cause negative business gain. In the first case, there is not enough business (i.e., service requests). In this case, a service provider should consider reducing the number of servers m and/or server speed s , so that the cost of infrastructure renting and the cost of energy consumption can be reduced. In the second case, there is too much business (i.e., service requests). In this case, a service provider should consider increasing the number of servers and/or server speed, so that the waiting time can be reduced and the revenue can be increased. However, increasing the number of servers and/or server speed also increases the cost of infrastructure renting and the cost of energy consumption. Therefore, we have the problem of selecting the optimal server size and/or server speed so that the profit is maximized.

7 PROFIT MAXIMIZATION

To formulate and solve our optimization problems analytically, we need a closed-form expression of C . To this end, let us use the following closed-form approximation, $\sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} \approx e^{m\rho}$, which is very accurate when m is not too small and ρ is not too large [17]. We also need Stirling's approximation of $m!$, i.e., $m! \approx \sqrt{2\pi m} \left(\frac{m}{e}\right)^m$. Therefore, we get the following closed-form approximation of p_m :

$$p_m \approx \frac{1 - \rho}{\sqrt{2\pi m}(1 - \rho)(e^\rho/e\rho)^m + 1},$$

and the following closed-form approximation of P_q :

$$P_q \approx \frac{1}{\sqrt{2\pi m}(1 - \rho)(e^\rho/e\rho)^m + 1}.$$

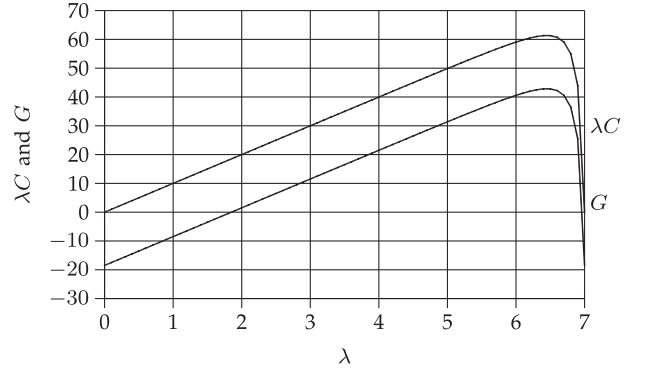


Fig. 4. Revenue and net business gain versus λ (constant-speed model).

By using the above-closed-form expression of P_q , we get a closed-form approximation of the expected service charge to a service request as

$$C \approx a\bar{r} \left(1 - \frac{1}{(\sqrt{2\pi m}(1 - \rho)(e^\rho/e\rho)^m + 1)} \times \frac{1}{((ms - \lambda\bar{r})(c/s_0 - 1/s) + 1)} \times \frac{1}{((ms - \lambda\bar{r})(a/d + c/s_0 - 1/s) + 1)} \right).$$

For convenience, we rewrite C as $C = a\bar{r}(1 - \frac{1}{D_1 D_2 D_3})$, where

$$\begin{aligned} D_1 &= \sqrt{2\pi m}(1 - \rho)(e^\rho/e\rho)^m + 1, \\ D_2 &= (ms - \lambda\bar{r})(c/s_0 - 1/s) + 1, \\ D_3 &= (ms - \lambda\bar{r})(a/d + c/s_0 - 1/s) + 1. \end{aligned}$$

Our discussion in this section is based on the above-closed-form expression of C .

7.1 Optimal Size

Given $\lambda, \bar{r}, s, P^*, \alpha, \beta, \gamma, a, c$, and d , our first problem is to find m such that G is maximized. To maximize G , we need to find m such that

$$\frac{\partial G}{\partial m} = \lambda \frac{\partial C}{\partial m} - (\beta + \gamma P^*) = 0,$$

for the idle-speed model, and

$$\frac{\partial G}{\partial m} = \lambda \frac{\partial C}{\partial m} - (\beta + \gamma(\xi s^\alpha + P^*)) = 0,$$

for the constant-speed model, where

$$\begin{aligned} \frac{\partial C}{\partial m} &= \frac{a\bar{r}}{(D_1 D_2 D_3)^2} \\ &\times \left(D_2 D_3 \frac{\partial D_1}{\partial m} + D_1 D_3 \frac{\partial D_2}{\partial m} + D_1 D_2 \frac{\partial D_3}{\partial m} \right). \end{aligned}$$

To continue the calculation, we rewrite D_1 as $D_1 = \sqrt{2\pi m}(1 - \rho)R + 1$, where $R = (e^\rho/e\rho)^m$. Notice that $\ln R = m \ln(e^\rho/e\rho) = m(\rho - \ln \rho - 1)$. Since

$$\frac{\partial \rho}{\partial m} = -\frac{\lambda\bar{r}}{m^2 s} = -\frac{\rho}{m},$$

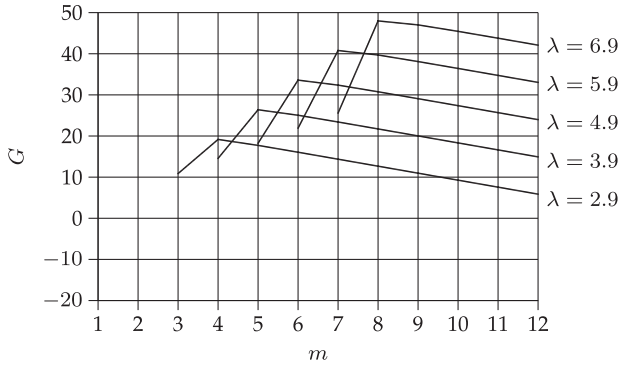


Fig. 5. Net business gain G versus m and λ (idle-speed model).

we get

$$\frac{1}{R} \frac{\partial R}{\partial m} = (\rho - \ln \rho - 1) + m \left(1 - \frac{1}{\rho}\right) \frac{\partial \rho}{\partial m} = -\ln \rho,$$

and

$$\frac{\partial R}{\partial m} = -R \ln \rho.$$

Now, we have

$$\begin{aligned} \frac{\partial D_1}{\partial m} &= \sqrt{2\pi} \left(\frac{1}{2\sqrt{m}} (1 - \rho) R + \sqrt{m} \left(-\frac{\partial \rho}{\partial m} \right) R \right. \\ &\quad \left. + \sqrt{m} (1 - \rho) \frac{\partial R}{\partial m} \right) \\ &= \sqrt{2\pi} \left(\frac{1}{2\sqrt{m}} (1 - \rho) R + \sqrt{m} \frac{\rho}{m} R - \sqrt{m} (1 - \rho) R \ln \rho \right) \\ &= \sqrt{2\pi} \left(\frac{1}{2\sqrt{m}} (1 - \rho) R + \frac{1}{\sqrt{m}} \rho R - \sqrt{m} (1 - \rho) (\ln \rho) R \right) \\ &= \sqrt{2\pi} \left(\frac{1}{2\sqrt{m}} (1 + \rho) R - \sqrt{m} (1 - \rho) (\ln \rho) R \right). \end{aligned}$$

Furthermore, we have

$$\frac{\partial D_2}{\partial m} = cs/s_0 - 1,$$

and

$$\frac{\partial D_3}{\partial m} = as/d + cs/s_0 - 1.$$

Although there is no closed-form solution to m , we notice that $\partial G/\partial m$ is a decreasing function of m . Therefore, m can be found numerically by using the standard bisection method.

In Figs. 5 and 6, we demonstrate the net business gain G in one unit of time as a function of m and λ for the two power consumption models, respectively, using the same parameters in Figs. 1, 2, 3, and 4. For $\lambda = 2.9, 3.9, 4.9, 5.9, 6.9$, we display G for m large enough such that $\rho < 1$. We notice that there is an optimal choice of m such that G is maximized. Using our analytical results, we can find m such that $\partial G/\partial m = 0$. The optimal value of m is 3.67479, 4.79218, 5.89396, 6.98457, 8.06655, respectively, for the idle-speed model, and 3.54842, 4.64834, 5.73478, 6.81160, 7.88104, respectively, for the constant-speed model.

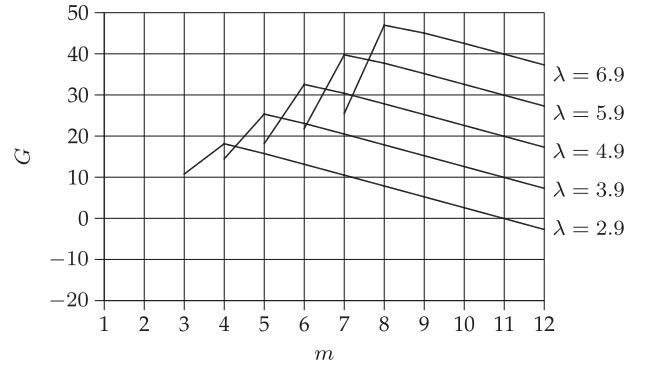


Fig. 6. Net business gain G versus m and λ (constant-speed model).

Such server size optimization has clear physical interpretation. When m is small such that ρ is close to 1, the waiting times of service requests are excessively long, and the service charges and the net business gain are low. As m increases, the waiting times are significantly reduced, and the service charges and the net business gain are increased. However, as m further increases, there will be no more increase in the expected services charge which has an upper bound $a\bar{r}$; on the other hand, the cost of a service provider (i.e., the rental cost and base power consumption) increases, so that the net business gain is actually reduced. Hence, there is an optimal choice of m which maximizes the profit.

7.2 Optimal Speed

Given $\lambda, \bar{r}, m, P^*, \alpha, \beta, \gamma, a, c$, and d , our second problem is to find s such that G is maximized. To maximize G , we need to find s such that

$$\frac{\partial G}{\partial s} = \lambda \frac{\partial C}{\partial s} - \gamma \lambda \bar{r} \xi (\alpha - 1) s^{\alpha-2} = 0,$$

for the idle-speed model, and

$$\frac{\partial G}{\partial s} = \lambda \frac{\partial C}{\partial s} - \gamma m \xi \alpha s^{\alpha-1} = 0,$$

for the constant-speed model, where

$$\begin{aligned} \frac{\partial C}{\partial s} &= \frac{a\bar{r}}{(D_1 D_2 D_3)^2} \\ &\quad \times \left(D_2 D_3 \frac{\partial D_1}{\partial s} + D_1 D_3 \frac{\partial D_2}{\partial s} + D_1 D_2 \frac{\partial D_3}{\partial s} \right). \end{aligned}$$

Similar to the calculation in the last section, we have

$$\frac{\partial \rho}{\partial s} = -\frac{\lambda \bar{r}}{ms^2} = -\frac{\rho}{s},$$

and

$$\frac{1}{R} \frac{\partial R}{\partial s} = m \left(1 - \frac{1}{\rho}\right) \frac{\partial \rho}{\partial s} = \frac{m}{s} (1 - \rho),$$

and

$$\frac{\partial R}{\partial s} = \frac{m}{s} (1 - \rho) R.$$

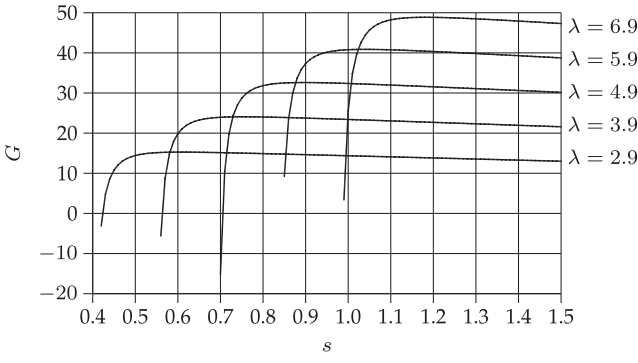


Fig. 7. Net business gain G versus s and λ (idle-speed model).

Now, we have

$$\begin{aligned}\frac{\partial D_1}{\partial s} &= \sqrt{2\pi m} \left(\left(-\frac{\partial \rho}{\partial s} \right) R + (1 - \rho) \frac{\partial R}{\partial s} \right) \\ &= \sqrt{2\pi m} \left(\frac{\rho}{s} R + \frac{m}{s} (1 - \rho)^2 R \right) \\ &= \sqrt{2\pi m} (\rho + m(1 - \rho)^2) \frac{R}{s}.\end{aligned}$$

Furthermore, we have

$$\begin{aligned}\frac{\partial D_2}{\partial s} &= m \left(\frac{c}{s_0} - \frac{1}{s} \right) + (ms - \lambda \bar{r}) \frac{1}{s^2} \\ &= \frac{mc}{s_0} - \frac{\lambda \bar{r}}{s^2},\end{aligned}$$

and

$$\begin{aligned}\frac{\partial D_3}{\partial s} &= m \left(\frac{a}{d} + \frac{c}{s_0} - \frac{1}{s} \right) + (ms - \lambda \bar{r}) \frac{1}{s^2} \\ &= m \left(\frac{a}{d} + \frac{c}{s_0} \right) - \frac{\lambda \bar{r}}{s^2}.\end{aligned}$$

Although there is no closed-form solution to s , we notice that $\partial G/\partial s$ is a decreasing function of s . Therefore, s can be found numerically by using the standard bisection method.

In Figs. 7 and 8, we demonstrate the net business gain G in one unit of time as a function of s and λ for the two power consumption models, respectively, using the same parameters in Figs. 1, 2, 3, 4, 5, and 6. For $\lambda = 2.9, 3.9, 4.9, 5.9, 6.9$, we display G for s large enough such that $\rho < 1$. We notice that there is an optimal choice of s such that G is maximized. Using our analytical results, we can find s such that

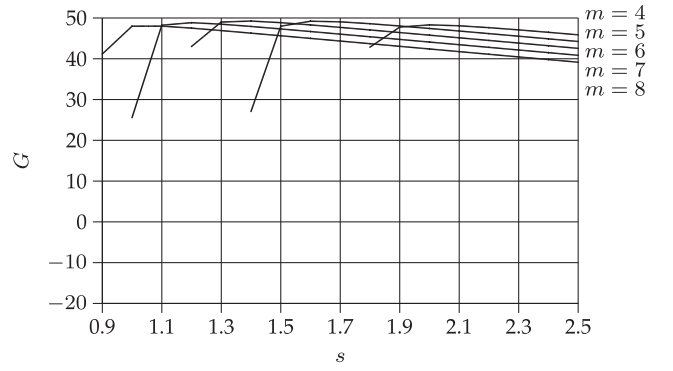


Fig. 9. Net business gain G versus s and m (idle-speed model).

$\partial G/\partial s = 0$. The optimal value of s is 0.63215, 0.76982, 0.90888, 1.04911, 1.19011, respectively, for the idle-speed model, and 0.57015, 0.71009, 0.85145, 0.99348, 1.13584, respectively, for the constant-speed model.

Such server speed optimization also has clear physical interpretation. When s is small such that ρ is close to 1, the waiting times of service requests are excessively long, and the service charges and the net business gain are low. As s increases, the waiting times are significantly reduced, and the service charges and the net business gain are increased. However, as s further increases, there will be no more increase in the expected services charge which has an upper bound $a\bar{r}$; on the other hand, the cost of a service provider (i.e., the cost of energy consumption) increases, so that the net business gain is actually reduced. Hence, there is an optimal choice of s which maximizes the profit.

7.3 Optimal Size and Speed

Given $\lambda, \bar{r}, P^*, \alpha, \beta, \gamma, a, c$, and d , our third problem is to find m and s such that G is maximized. To maximize G , we need to find m and s such that $\partial G/\partial m = 0$ and $\partial G/\partial s = 0$, where $\partial G/\partial m$ and $\partial G/\partial s$ have been derived in the last two sections. The two equations can be solved by a nested bisection search procedure.

In Figs. 9 and 10, we demonstrate the net business gain G in one unit of time as a function of s and m for the two power consumption models, respectively, using the same parameters in Figs. 1, 2, 3, 4, 5, 6, 7, and 8, where $\lambda = 6.9$. For $m = 4, 5, 6, 7, 8$, we display G for s large enough such that $\rho < 1$. Using our analytical results, we can find m and s such that $\partial G/\partial m = 0$ and $\partial G/\partial s = 0$. For the idle-speed

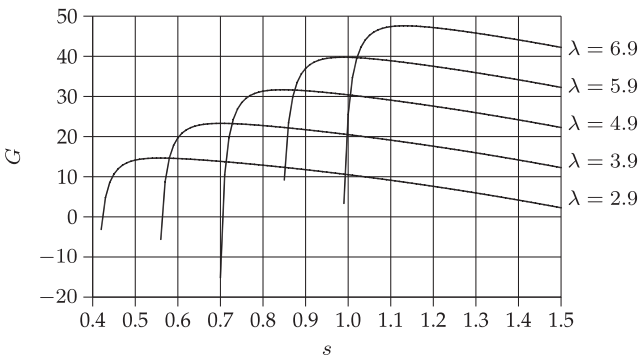


Fig. 8. Net business gain G versus s and λ (constant-speed model).

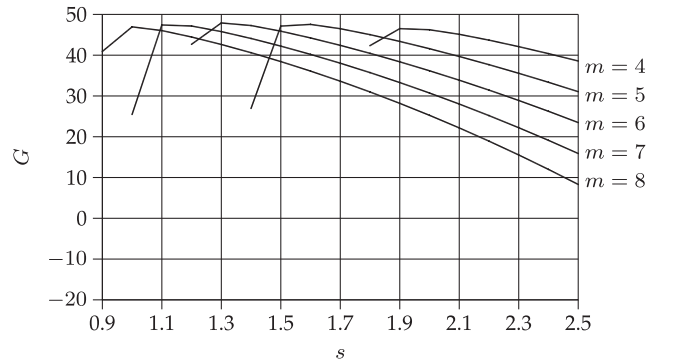


Fig. 10. Net business gain G versus s and m (constant-speed model).

model, the theoretically optimal values are $m = 5.56827$ and $s = 1.46819$, which result in the maximum $G = 49.25361$ by using the closed-form approximation of C . Practically, m can be either 5 or 6. When $m = 5$, the optimal value of s is 1.62236, which results in the maximum $G = 49.16510$. When $m = 6$, the optimal value of s is 1.37044, which results in the maximum $G = 49.18888$. Hence, the practically optimal setting is $m = 6$ and $s = 1.37044$, and the maximum net business gain in one unit of time is $G = 49.29273$ by using the exact expression of C . For the constant-speed model, the theoretically optimal values are $m = 5.79074$ and $s = 1.35667$, which result in the maximum $G = 47.80769$ by using the closed-form approximation of C . Practically, m can be either 5 or 6. When $m = 5$, the optimal value of s is 1.55839, which results in the maximum $G = 47.63979$. When $m = 6$, the optimal value of s is 1.31213, which results in the maximum $G = 47.78640$. Hence, the practically optimal setting is $m = 6$ and $s = 1.31213$, and the maximum net business gain in one unit of time is $G = 47.91830$ by using the exact expression of C .

8 SIMULATION SETTINGS AND RESULTS

See Section 1 of the supplementary material, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPDS.2012.203>.

9 PROOFS OF THEOREMS 1 AND 2

See Section 2 of the supplementary material, available online.

10 CONCLUDING REMARKS

We have proposed a pricing model for cloud computing which takes many factors into considerations, such as the requirement r of a service, the workload λ of an application environment, the configuration (m and s) of a multiserver system, the service level agreement c , the satisfaction (r and s_0) of a consumer, the quality (W and T) of a service, the penalty d of a low-quality service, the cost (β and m) of renting, the cost (α , γ , P^* , and P) of energy consumption, and a service provider's margin and profit a . By using an M/M/m queueing model, we formulated and solved the problem of optimal multiserver configuration for profit maximization in a cloud computing environment. Our discussion can be easily extended to other service charge functions. Our methodology can be applied to other pricing models.

ACKNOWLEDGMENTS

The authors are grateful to three anonymous reviewers for their constructive comments. Part of the work were performed while K. Hwang, K. Li, and A.Y. Zomaya were visiting Tsinghua National Laboratory for Information Science and Technology, Tsinghua University, in the winter of 2011 and the summer of 2012 as Intellectual Ventures endowed visiting chair professors. This work is partially supported by Ministry of Science and Technology of China under National 973 Basic Research Grants No. 2011CB302805 and No. 2011CB302505, and National 863 High-tech Program Grant No. 2011AA040501; Ministry of

Industry and Information Technology of China under the Internet of Things program; the Innovation R/D Team Program of Guangdong Province, China, under contract No. 201001D0104726115; Australian Research Grant DP1097110.

REFERENCES

- [1] <http://en.wikipedia.org/wiki/CMOS>, 2012.
- [2] http://en.wikipedia.org/wiki/Service_level_agreement, 2012.
- [3] M. Armbrust et al., "Above the Clouds: A Berkeley View of Cloud Computing," Technical Report No. UCB/EECS-2009-28, Feb. 2009.
- [4] R. Buyya, D. Abramson, J. Giddy, and H. Stockinger, "Economic Models for Resource Management and Scheduling in Grid Computing," *Concurrency and Computation: Practice and Experience*, vol. 14, pp. 1507-1542, 2007.
- [5] R. Buyya, C.S. Yeo, S. Venugopal, J. Broberg, and I. Brandic, "Cloud Computing and Emerging IT Platforms: Vision, Hype, and Reality for Delivering Computing as the Fifth Utility," *Future Generation Computer Systems*, vol. 25, no. 6, pp. 599-616, 2009.
- [6] A.P. Chandrakasan, S. Sheng, and R.W. Brodersen, "Low-Power CMOS Digital Design," *IEEE J. Solid-State Circuits*, vol. 27, no. 4, pp. 473-484, Apr. 1992.
- [7] B.N. Chun and D.E. Culler, "User-Centric Performance Analysis of Market-Based Cluster Batch Schedulers," *Proc. Second IEEE/ACM Int'l Symp. Cluster Computing and the Grid*, 2002.
- [8] D. Durkee, "Why Cloud Computing Will Never be Free," *Comm. ACM*, vol. 53, no. 5, pp. 62-69, 2010.
- [9] R. Ghosh, K.S. Trivedi, V.K. Naik, and D.S. Kim, "End-to-End Performability Analysis for Infrastructure-as-a-Service Cloud: An Interacting Stochastic Models Approach," *Proc. 16th IEEE Pacific Rim Int'l Symp. Dependable Computing*, pp. 125-132, 2010.
- [10] K. Hwang, G.C. Fox, and J.J. Dongarra, *Distributed and Cloud Computing*. Morgan Kaufmann, 2012.
- [11] "Enhanced Intel SpeedStep Technology for the Intel Pentium M Processor," White Paper, Intel, Mar. 2004.
- [12] D.E. Irwin, L.E. Grit, and J.S. Chase, "Balancing Risk and Reward in a Market-Based Task Service," *Proc. 13th IEEE Int'l Symp. High Performance Distributed Computing*, pp. 160-169, 2004.
- [13] H. Khazaei, J. Misic, and V.B. Misic, "Performance Analysis of Cloud Computing Centers Using M/G/m/m+r Queueing Systems," *IEEE Trans. Parallel and Distributed Systems*, vol. 23, no. 5, pp. 936-943, May 2012.
- [14] L. Kleinrock, *Queueing Systems: Theory*, vol. 1. John Wiley and Sons, 1975.
- [15] Y.C. Lee, C. Wang, A.Y. Zomaya, and B.B. Zhou, "Profit-Driven Service Request Scheduling in Clouds," *Proc. 10th IEEE/ACM Int'l Conf. Cluster, Cloud and Grid Computing*, pp. 15-24, 2010.
- [16] K. Li, "Optimal Load Distribution for Multiple Heterogeneous Blade Servers in a Cloud Computing Environment," *Proc. 25th IEEE Int'l Parallel and Distributed Processing Symp. Workshops*, pp. 943-952, May 2011.
- [17] K. Li, "Optimal Configuration of a Multicore Server Processor for Managing the Power and Performance Tradeoff," *J. Supercomputing*, vol. 61, no. 1, pp. 189-214, 2012.
- [18] P. Mell and T. Grance, "The NIST Definition of Cloud Computing," Nat'l Inst. of Standards and Technology, <http://csrc.nist.gov/groups/SNS/cloud-computing/>, 2009.
- [19] F.I. Popovici and J. Wilkes, "Profitable Services in an Uncertain World," *Proc. ACM/IEEE Conf. Supercomputing*, 2005.
- [20] J. Sherwani, N. Ali, N. Lotia, Z. Hayat, and R. Buyya, "Libra: A Computational Economy-Based Job Scheduling System for Clusters," *Software - Practice and Experience*, vol. 34, pp. 573-590, 2004.
- [21] C.S. Yeo and R. Buyya, "A Taxonomy of Market-Based Resource Management Systems for Utility-Driven Cluster Computing," *Software - Practice and Experience*, vol. 36, pp. 1381-1419, 2006.
- [22] B. Zhai, D. Blaauw, D. Sylvester, and K. Flautner, "Theoretical and Practical Limits of Dynamic Voltage Scaling," *Proc. 41st Design Automation Conf.*, pp. 868-873, 2004.



Junwei Cao received the bachelor's and master's degrees in control theories and engineering in 1996 and 1998, respectively, both from Tsinghua University, Beijing, China. He received the PhD degree in computer science from the University of Warwick, Coventry, United Kingdom, in 2001. He is currently a professor and vice director, Research Institute of Information Technology, Tsinghua University, Beijing, China. He is also the director of Common Platform

and Technology Division, Tsinghua National Laboratory for Information Science and Technology. Before joining Tsinghua University in 2006, he was a research scientist at MIT LIGO Laboratory and NEC Laboratories Europe for about five years. He has published more than 130 papers and cited by international scholars for over 3,000 times. He is the book editor of *Cyberinfrastructure Technologies and Applications*, published by Nova Science in 2009. His research is focused on advanced computing technologies and applications. He is a senior member of the IEEE and IEEE Computer Society, and a member of the ACM and CCF.



Kai Hwang received the PhD degree from the University of California, Berkeley in 1972. He is a professor of EE/CS at the University of Southern California. He also chairs the IV-endowed visiting chair professor group at Tsinghua University in China. He has published eight books and more than 218 scientific papers in computer architecture, parallel processing, distributed systems, cloud computing, network security, and Internet applications. His popular books

have been adopted worldwide and translated into four foreign languages. His published papers have been cited more than 11,000 times by early 2012. His latest book *Distributed and Cloud Computing: from Parallel Processing to the Internet of Things* (with G. Fox and J. Dongarra) was just published by Kaufmann in 2011. He received the 2004 CFC Outstanding Achievement Award, and the Founders Award for his pioneering work in parallel processing from IEEE IPDPS in 2011. He has served as a founding editor-in-chief of the *Journal of Parallel and Distributed Computing* for 28 years. He has delivered 34 keynote addresses on advanced computing systems and cutting-edge information technologies in major IEEE/ACM Conferences. He has performed advisory, consulting and collaborative work for IBM, Intel, MIT Lincoln Lab, JPL at Caltech, ETL in Japan, ITRI in Taiwan, GMD in Germany, INRIA in France, and Chinese Academy of Sciences. He is a fellow of the IEEE (1986).



Keqin Li is a SUNY distinguished professor in computer science and an Intellectual Ventures endowed visiting chair professor at Tsinghua University, China. His research interests are mainly in design and analysis of algorithms, parallel and distributed computing, and computer networking. He has contributed extensively to processor allocation and resource management; design and analysis of sequential/parallel, deterministic/probabilistic, and approximation algorithms; parallel and distributed computing systems performance analysis, prediction, and evaluation; job scheduling, task dispatching, and load balancing in heterogeneous distributed systems; dynamic tree embedding and randomized load distribution in static networks; parallel computing using optical interconnections; dynamic location management in wireless communication networks; routing and wavelength assignment in optical networks; energy-efficient power management and performance optimization. He has published more than 240 research publications and has received several Best Paper Awards for his highest quality work. He is currently on the editorial board of *IEEE Transactions on Parallel and Distributed Systems*, *IEEE Transactions on Computers*, *International Journal of Parallel, Emergent and Distributed Systems*, *International Journal of High Performance Computing and Networking*, and *Optimization Letters*. He is a senior member of the IEEE.



Albert Y. Zomaya is currently the chair professor of high performance computing and networking and Australian Research Council Professorial fellow in the School of Information Technologies, The University of Sydney. He is also the director of the Centre for Distributed and High Performance Computing which was established in late 2009. He is the author/coauthor of seven books, more than 400 papers, and the editor of nine books and 11 conference proceedings. He is the

editor-in-chief of the *IEEE Transactions on Computers* and serves as an associate editor for 19 leading journals, such as, the *IEEE Transactions on Parallel and Distributed Systems* and *Journal of Parallel and Distributed Computing*. He is the recipient of the Meritorious Service Award (in 2000) and the Golden Core Recognition (in 2006), both from the IEEE Computer Society. Also, he received the IEEE Technical Committee on Parallel Processing Outstanding Service Award and the IEEE Technical Committee on Scalable Computing Medal for Excellence in Scalable Computing, both in 2011. He is a chartered engineer, a fellow of AAAS, IEEE, and IET (United Kingdom).

► **For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.**