

基于 SMDP 的动态云计算资源优化管理系统*

梁宏斌¹⁺, 彭代渊², 刘燕³

¹(信息安全国家重点实验室(中国科学院 信息工程研究所), 北京 100093)

²(西南交通大学 信息科学与技术学院, 四川 成都 610031)

³(北京大学 软件与微电子学院, 北京 102600)

Dynamic Cloud Resources Allocation Based on SMDP

LIANG Hong-Bin¹⁺, PENG Dai-Yuan², LIU Yan³

¹(State Key Laboratory of Information Security (Institute of Information Engineering, The Chinese Academy of Sciences), Beijing 100093, China)

²(School of Information Science & Technology, Southwest Jiaotong University, Chengdu 610031, China)

³(School of Software and Microelectronics, Peking University, Beijing 102600, China)

+ Corresponding author: E-mail: liang.hongbi@gmail.com

Liang HB, Peng DY, Liu Y. Dynamic cloud resources allocation based on SMDP. *Journal of Software*, 2012, 23(Suppl. (1)): 25-37 (in Chinese). <http://www.jos.org.cn/1000-9825/12004.htm>

Abstract: With mobile cloud computing system architecture gradually replacing the traditional Client-Server model, one of the most critical issues that needs to be addressed is the resource allocation deciding problem to resolve how a cloud could efficiently handle the cloud resources to satisfy the requirements of mobile device for cloud services. Simultaneously, this is in effort to obtain the maximal utilization rate of cloud resources and system rewards of cloud. The study first proposes a novel cloud computing resource allocation model based on semi-Markov decision process (SMDP) to achieve the optimal resource allocation scheme in terms of maximal system reward, cloud resource utilization, and the QoS of mobile device while capturing the dynamics of resource request arrivals and departures. Finally, the performance of the proposed model is evaluated by the simulation results, which show that the obtained theoretic results are consistent with this simulation results.

Key words: SMDP, mobile cloud computing service domain, system reward, blocking probability

摘要: 随着移动云计算方式正在逐步替代传统的 Client-Server 方式,在移动云计算网络中,如何有效地分配云计算资源来尽量满足移动终端对云计算服务的需求,同时使得移动云计算网络的云计算资源利用率和系统收益最大,就成为当前云计算领域中一个重要的研究课题.首先提出了一种基于半马氏决策过程(SMDP)的移动云计算服务域动态云计算资源优化管理模型,通过该模型获得的云计算资源优化管理决策策略不仅能使移动云计算服务域的系统收益最大,同时也能提高云计算资源的利用率以及移动终端的用户满意度和服务质量(QoS),并能反映由于移动终端的服务请求到达移动云计算网络以及移动终端结束服务离开移动云计算网络而引起的云计算资源动态变化的真实情况.最后对提出的移动云计算服务域动态云计算资源优化管理模型的性能通过仿真进行了验证,实验仿真

* 基金项目: 国家重点基础研究发展计划(973)(2011CB302902, 2012CB316100); 国家科技重大专项(2010ZX03006-001-01); 中国科学院先导专项课题(XDA06040100);

收稿时间: 2012-05-05; 定稿时间: 2012-08-11

结果验证了理论分析的正确性.

关键词: 半马氏决策过程;移动云计算服务域;系统收益;阻塞率

云计算是一种以资源按需分配、pay-as-you-go 以及效能计算为特征的新的计算服务模式^[1].云计算不仅为云计算服务商同时也为个人用户提供了一种新的计算模式,它可以被广义的分为 Infrastructure-as-a-service (IaaS),Platform-as-a-service(PaaS)以及 Software-as-a-service(SaaS)三大类.随着无线通信技术以及互联网技术的发展,根据 Gartner 的预测,在 2013 年以后,移动终端将取代 PC 机成为全球最主要的互联网接入设备^[2].由于移动终端(MD)与传统有线终端相比具有更多的优势(例如移动性、灵活性以及感知能力等),因此将移动计算和云计算技术结合在一起自然就成为构建移动应用的新方法,目前无论在学术界还是工业界也吸引了越来越多的关注.因而,一个新的研究领域-移动云计算(mobile cloud computing)也就应运而生.

在以前关于移动云计算的研究中,主要的研究方向集中在计算任务的上传下载、远程运行以及动态组织等.Li 等人提出了一个在移动终端和云端都能运行移动应用的移动云计算模型,从而资源有限的移动终端能将计算、传输以及存储任务上传到云端运行^[3].Chun 和 Maniatis 通过增加执行次数来配置 CloneCloud 云资源,但是没有考虑到用户终端的实际运行状态^[4].Zhang 等人对移动终端通过云计算网络对弹性应用服务的资源管理做了一些初步的研究^[5].Huang 等人提出了移动云计算架构,该模型允许移动终端将相关应用上传到云端的虚拟机(virtual machine,简称 VM)运行^[6].Meng 等人提出了一种根据不同地域、不同流量来配置虚拟机,通过优化配置虚拟机的放置位置来提高网络利用率的新方法^[7].实际上,由于这些关于移动云计算的架构设施的研究探讨已经比较充分,因此,移动云计算的资源管理自然将成为下一个的主要研究方向.

在移动云计算网络中,基于服务器族群在地理位置上的分布式放置,系统的云计算资源(例如 CPU、内存以及存储等)分别由多个移动云计算服务域来负责管理.每一个移动云计算服务域由多个虚拟机组成,而每一个虚拟机则由能处理一个云计算服务的最小云计算资源组成.尽管与移动终端相比,移动云计算网络的云计算资源通常被认为是无限的,但是仍然非常有必要充分利用移动云计算服务域中的云计算资源来实现移动云计算网络的低成本运行.

尽管资源优化管理在无线通信网络中已经被广泛地进行了研究^[8-10,22],但是目前对云计算尤其是移动云计算的资源优化管理的研究还比较少.Liang 等人提出了一个经济型的移动云计算资源管理模型,该模型能在给定系统配置的情况下,通过在云端和移动终端之间优化分配移动应用来获得移动云计算网络的最大收益^[11].Wei 等人提出了一个基于博弈论的云计算资源分配模型^[12],该模型能根据移动终端对用户服务质量(QoS)的需求来分配云计算资源.另外,还有一些文献对云计算网络如何通过虚拟机或者是数据中心的服务器来优化管理云计算资源进行了研究.Lorincz 等人提出了一个新的云计算操作模型,该操作模型不仅能使用户在掌握云计算资源的情况下进行编程,同时也能实现云计算网络中,云计算服务重用云计算资源的管理模式^[13].他还对云计算网络中事件应用的云计算资源管理进行了研究^[14].Tesauro 等人提出了一个基于增强型自学习系统的资源管理模型来对云计算网络中的服务器进行动态分配,从而提高云计算网络的收益^[15].Boloor 等人提出了一个通用的对云计算服务请求进行分配和规划的方案,该方案在获得用户指定的服务质量的同时,提高了云计算服务提供商的收益^[16].

国内目前对云计算资源优化管理的研究还较少.崔云飞等人设计了基于多级递阶控制的云计算资源共享模型,该模型对如何构建云计算系统提供了一种可行的思路^[20].袁文成等人提出了一种面向虚拟资源的云计算资源管理机制,该机制通过对虚拟资源的划分、预留及调度策略,为用户提供有效的 IaaS 服务^[21].上述两篇文献的研究主要是为了提高并试图最大化云计算资源的利用率,保证云计算资源管理机制的有效性和可靠性,但是都没有考虑到云计算网络的系统收益.云计算系统收益对云计算网络的性能至关重要,因此云计算网络迫切需要一种既考虑云计算资源利用率又考虑云计算系统收益的新的云计算资源优化管理方案.

移动云计算网络的一个主要优势是允许移动终端在云端运行他们的移动应用服务.而一个云计算服务(移动应用服务)可以被分配多个 VM 的云计算资源来使移动终端获得更高的计算和存储能力.当移动云计算服务

域收到一个从移动终端发送过来的云计算服务请求时,系统需要分析当前可用的云计算资源,并基于分析结果决定是否接收该云计算服务请求;如果决定是接收,那么系统还需要进一步判断具体为该移动终端的云计算服务请求分配多少云计算资源(即 VM 的个数).如果移动云计算服务域中所有的云计算资源已经被占用,那么由于云计算资源的不足,系统会拒绝该移动终端的云计算服务请求(我们假设在移动云计算中,没有队列缓冲).对云计算服务请求的拒绝不仅对移动终端的用户满意度和服务质量带来了负面的影响,而且也极大地降低了系统的收益.

移动云计算服务域的系统收入通常随着被接收的云计算服务请求数量的增加而增加.但另一方面,随着系统接收的云计算服务请求越多,那么分配给每一个云计算服务的云计算资源也就越少,从而降低了正在接受服务的移动终端的用户满意度以及移动云计算服务域的系统性能.而现有关于云计算资源分配管理的方法大部分只考虑了系统的收入,没有考虑到云计算资源被占用所带来的支出,也没有考虑到移动终端的用户满意度和服务质量(QoS).因此,为了能得到移动云计算服务域全面的系统收益,在计算移动云计算服务域的系统收益时,不仅需要考虑移动云计算网络的收入,也需要考虑云计算资源被占用所带来的支出以及移动终端的用户满意度和服务质量(QoS).根据我们现有的调查,目前国内外还没有文献在这方面进行深入地研究.

在本文中,我们基于半马氏决策过程(SMDP)提出了新的移动云计算服务域动态云计算资源优化管理模型,通过该模型来获得移动云计算服务域的云计算资源的优化管理决策策略,并得到移动云计算服务域的最大收益,该收益不仅考虑了接收云计算服务请求所带来的收入,同时也考虑了因云计算服务占用云计算资源所带来的支出,以及移动终端的用户满意度和服务质量(QoS).

本文所提出的移动云计算服务域动态云计算资源优化管理模型的主要贡献有如下 3 点:

- 基于半马氏决策过程(SMDP)推导出了移动云计算服务域的动态云计算资源优化分配决策策略.
- 该模型能基于移动云计算服务域当前可用的云计算资源,为云计算服务请求自适应地分配不同的云计算资源,通过充分利用该移动云计算服务域的云计算资源来提高云计算资源利用率,并获得移动云计算服务域的最大整体收益.
- 该模型获得的移动云计算服务域的最大系统收益,既考虑了移动云计算服务域接收云计算服务请求所带来的收入,也考虑了因云计算资源被占用所带来的支出,还考虑了移动终端的用户满意度和服务质量(QoS).因此,通过该模型得到的系统收益是全面的整体收益.

本文所提出的移动云计算服务域动态云计算资源优化管理的系统模型将会在第 1 节进行阐述;在第 2 节中,我们将会对移动云计算服务域动态优化管理云计算资源的方案基于半马氏决策过程(SMDP)进行建模分析;根据本文所提出的移动云计算服务域动态云计算资源优化管理模型,第 3 节分析并推导了系统对云计算服务请求采取不同云计算资源分配方案的概率以及对云计算服务请求的拒绝概率;第 4 节对本文所提出的移动云计算服务域动态云计算资源优化管理模型的性能通过仿真实验进行了评估;在最后一节中,对本文内容进行了总结并且提出了以后的研究方向.

1 动态云计算资源优化管理的系统模型

移动云计算网络与传统的 Client-Server 服务模式相比的一个主要优势是:当移动终端将他们的应用服务上传到云端运行时,移动终端能获得更多的容量以及更好的性能(例如更少的处理时间以及移动终端电池电量的节约等).移动终端弹性应用任务的上传可以通过连接云端和移动终端的 Weblet 来实现.一个 Weblet 可以使用独立于平台的 Java 或者 .Net 或者 Python 语言,也可以使用平台编程语言.Chun 和 Maniatis 研究了将 Weblet 从移动终端上传到云端运行的算法^[4].通过 Weblet 上传弹性应用服务到云端运行,移动终端可以大幅提高自身的计算能力、存储能力以及网络带宽.通常,移动终端决定是否将任务上传到

云端运行取决于移动终端自身的状态(例如,移动终端的 CPU 处理能力、电池的电量、网络连接质量以及移动终端对安全的考虑等因素).在本文中,当移动终端决定将任务上传到云端运行时,它会首先给云端发送一个服务请求,如果云端接收了移动终端的服务请求,那么移动终端随后就会将任务上传到云端运行,运行结束

后,云端会将运行结果返回给移动终端.

如图 1 所示,一个 VM 负责管理 Weblet 在移动云计算网络中的上传、卸载以及处理.如前所述,在本文所提出的移动云计算服务域动态云计算资源优化管理模型中,一个 VM 负责管理移动云计算服务域中处理一个云计算服务所需的最小云计算资源(CPU、内存以及存储等),每一个 VM 所管理的云计算资源一次只能处理一个云计算服务请求.虽然我们可以认为移动云计算网络中的云计算资源是有限的,但是在移动云计算网络中,某个具体的移动云计算服务域的以 VM 个数来计算的云计算资源又是有限的.因此,在移动云计算服务域中,如果到达的云计算服务请求的数量超过了该服务域中可用的云计算资源 VM 个数,则随后到来的云计算服务请求将会被该服务域拒绝.另一方面,如果到达的云计算服务请求的数量远低于该服务域中可用的云计算资源的 VM 个数,那么该服务域就可以为每个云计算服务请求分配更多的 VM 个数来充分利用该服务域的云计算资源,以此来提高该移动云计算服务域的云计算资源利用率以及移动终端的用户满意度和服务质量(QoS).

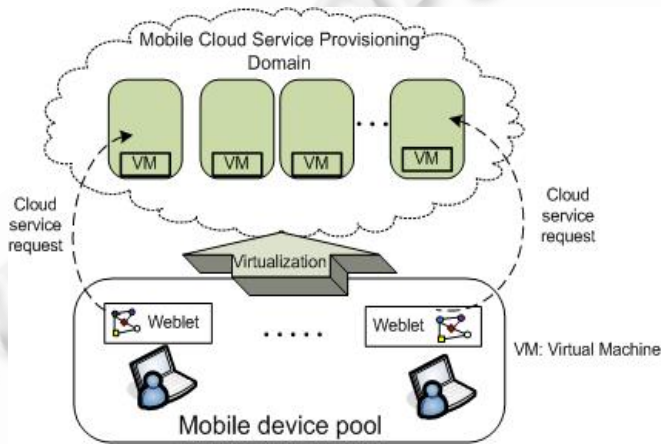


Fig.1 Service model of mobile cloud computing
图 1 移动云计算网络的服务模型

因此,本文所提出的移动云计算服务域动态云计算资源优化管理模型的目标就是通过充分利用云计算资源,使得移动云计算服务域既能获得最大整体系统收益,也能提高该服务域的云计算资源利用率以及用户满意度和服务质量(QoS).

在本文中,我们考虑一个只包含一个云计算服务域的移动云计算网络,设该移动云计算服务域的云计算资源总共为 K 个虚拟机(VM).用 r 表示分配给一个移动终端的云计算资源的 VM 个数,其中 r 是一个正整数,并且满足条件 $1 \leq r \leq K$.此外,我们可以用不同的效能函数来度量移动云计算用户的满意度.例如,可以用类似 Sigmoidal^[17]的函数来描述移动云计算网络中移动终端的用户满意度,

$$U(r) = 1 - \exp\left(-\frac{\omega_2 r^2}{\omega_1 + r}\right) \tag{1}$$

其中, $U(r)$ 表示移动云计算网络中移动终端的用户满意度, r 是移动云计算服务域分配给移动终端的云计算资源的 VM 个数, ω_1 和 ω_2 是用来调节 $U(r)$ 波形的参数,其函数的波形如图 2 所示.

通常来讲,参数 ω_1 和 ω_2 的选择是由移动云计算服务域和最终用户对服务质量(QoS)的需求来决定的,有效的指数选择对移动云计算服务域的云计算资源管理具有显著的影响.从公式(1)可知,为了提高移动云计算网络的用户满意度,需要给移动终端用户分配尽可能多的云计算资源.但是另一方面,为了能提高移动云计算网络的总体系统收益,根据该移动云计算服务域总的可用云计算资源和移动用户使用云计算服务时对云计算资源的需求,系统又不可能单独给每一个移动终端用户都分配最大的云计算资源.因此,这就成为本文建立的移动云计

算服务域动态云计算资源优化管理模型需要解决的一个矛盾.为了能对移动云计算网络的云计算服务对云计算资源的动态需求建立优化模型,我们假设移动终端请求接入移动云计算网络并使用云计算服务的过程服从泊松分布(Poisson),其均值为 λ ,同时也假设移动云计算终端用户的在网时间服从指数分布,其均值为 $\frac{1}{\mu}$.在以下的章节中,我们将会定义本文所建立的移动云计算服务域动态云计算资源优化管理模型的系统状态、决策以及收益模型.

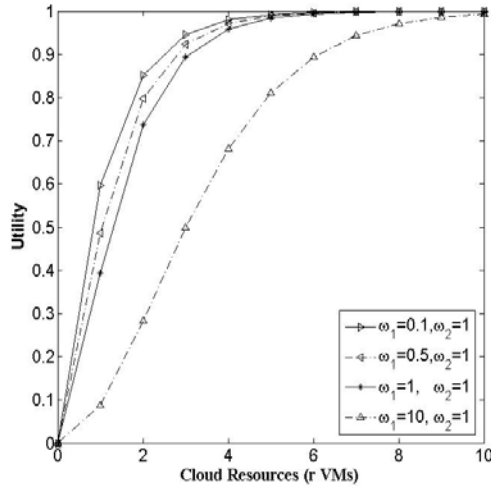


Fig.2 Utility function of cloud service

图 2 云计算服务的效能函数

1.1 系统状态

为了能使用半马氏决策过程来表征移动云计算服务域云计算资源的优化管理模型,我们假设用户对移动云计算网络的满意度与其服务请求被处理的时间成反比,即用户的云计算服务请求被处理的时间越短则用户对系统的满意度越高.显然,如果分配给该用户的云计算资源(即 VM 个数)越多,则该用户所请求的云计算服务被处理的时间则越短(云计算服务通常可以由多个 VM 的云计算资源来并行处理,从而提高运行速度).由此可知,用户对移动云计算网络的满意度与分配给该用户的云计算资源成正比,即分配给用户的云计算资源越多,则该用户的满意度越高.我们将移动云计算网络的用户满意度根据移动云计算服务域分配给移动终端用户的云计算资源的 VM 个数分为 N 类.因此我们定义本文基于半马氏过程的移动云计算服务域动态云计算资源优化管理模型的系统状态为每一个用户满意度下正在运行的云计算服务的数量以及在该移动云计算服务域中所发生的事件的集合.不失一般性,我们定义 k_i 为分配给用户满意度为 i 的用户的云计算资源的 VM 个数,这里 $i=1, 2, \dots, N$, 并且 $0 < k_1 < \dots < k_N \leq K$, 同时用 n_i 表示移动云计算服务域中用户满意度为 i 的正在运行的云计算服务的数量.

在移动云计算服务域中,总共有两种类型的事件:

- 1) 一个新到的云计算服务请求,用 R 来表示;
- 2) 一个用户满意度为 i 的云计算服务完成了运行,并释放了其所占用的云计算资源,用 D_i 来表示.

因此在该移动云计算服务域内任何事件 e 都可以用集合 $e \in \{R, D_1, D_2, \dots, D_N\}$ 来表示,从而移动云计算服务域的系统状态可用如下公式来表示:

$$S = \{s | s = \langle n_1, n_2, \dots, n_N, e \rangle\} \tag{2}$$

1.2 行动集合

当一个终端用户请求接入移动云计算服务域并申请使用云计算服务时($e=R$),该移动云计算服务域需要决定是否接收这个云计算服务请求,如果接收的话,那么应该给该云计算服务分配多少个 VM 的云计算资源.为简单计,我们用 $A(s)=0$ 表示该移动云计算服务域拒绝一个云计算服务请求;用 $A(s)=i$ 表示移动云计算服务域接收了该云计算服务请求,并且分配给这个云计算服务请求的云计算资源为 k_i 个 VM,以期能使移动用户对该云计算服务的满意度达到 i ,这里 s 表示当前的系统状态,而另一方面,我们用 $A(s)=-1$ 来表示在一个云计算服务结束运行并释放所占用云计算资源的状态(这里事件 $e=D_i$)下的行动,即统计现有可用的云计算资源的 VM 个数并等待下一个事件的发生.因此,该模型的行动集合总结如下:

$$A(s) = \begin{cases} -1, & e \in \{D_1, D_2, \dots, D_N\} \\ \{0, 1, \dots, N\}, & e = R \end{cases} \quad (3)$$

1.3 收益模型

基于系统状态和对应的行动,我们可以预估一个移动云计算服务域能获得的收益(用 $r(s,a)$ 来表示),这个收益由两部分组成,一部分是系统的收入,另一部分是系统的支出,可以用下式来表示,

$$r(s,a) = x(s,a) - \tau(s,a)y(s,a) \quad (4)$$

$x(s,a)$ 是系统在状态为 s ,选择的行动为 a 时,系统所获得的总收入,用如下的效能函数来表示为

$$x(s,a) = \begin{cases} -1, & e = R, a = 0 \\ U(k_i), & e = R, a = i \end{cases} \quad (5)$$

$\tau(s,a)$ 表示在当前系统状态为 s ,当选取的行动为 a 时,转移到下一个系统状态 j 所预期的服务时间; $y(s,a)$ 表示在当前系统状态为 s ,选取的行动为 a 时的支出, $y(s,a)$ 可以用正在运行的云计算服务所占用的云计算资源 VM 的总个数来度量,由下式表示为

$$y(s,a) = \sum_{i=1}^N n_i k_i \quad (6)$$

2 基于半马尔可夫决策系统的建模

通常一个半马氏决策过程模型包含有 6 个要素:

1) 系统状态;2) 事件;3) 行动集合;4) 收益;5) 决策时间点;6) 状态转移概率.前 4 个要素我们已经在上一节里进行了讨论.本节主要讨论后 2 个要素.

决策时间点是指当任何一个事件发生并需要做决策的时间点,例如一个云计算服务请求到达移动云计算服务域,或者是一个已经完成云计算服务的移动终端用户离开该云计算服务域并且释放所占用的云计算资源.在我们的系统模型中,由于在两个决策点之间的时间 $\tau(s,a)$ 均服从指数分布,因此,所有事件发生的平均速率 $\gamma(s,a)$ 可以表示为

$$\gamma(s,a) = \tau(s,a)^{-1} = \begin{cases} \lambda + \sum_{i=1}^N n_i \mu, & e = R, a = 0 | e = D_i \\ \lambda + \left(\sum_{i=1}^N n_i + 1 \right) \mu, & e = R, a = i \end{cases} \quad (7)$$

应用半马氏决策过程(SMDP)的折扣收益模型^[18,19],在时间 $\tau(s,a)$ 之间的期望折扣收益 $z(s,a)$ 可表示为

$$z(s,a) = x(s,a) - y(s,a) E_s^\alpha \left\{ \int_0^{\tau} e^{-\alpha t} dt \right\} = x(s,a) - y(s,a) E_s^\alpha \left\{ \frac{1 - e^{-\alpha \tau}}{\alpha} \right\} = x(s,a) - \frac{y(s,a)}{\alpha + \gamma(s,a)} \quad (8)$$

其中, $x(s,a)$ 和 $y(s,a)$ 分别在公式(5)和公式(6)已经定义, α 为连续时间下的折扣率.由此,我们可以得到最大长期折扣收益为

$$v(s) = \max_{a \in A} \left\{ z(s, a) + \eta \sum_{j \in S} p(j | s, a) v(j) \right\} \quad (9)$$

其中, $\eta = \frac{\gamma(s, a)}{\alpha + \gamma(s, a)}$, $p(j | s, a)$ 表示系统在状态 s , 当选取的行动为 a 时, 系统转移到状态 j 的状态转移概率. 以下将

推导出所有的状态转移概率. 为了简化书写, 我们定义如下的符号:

$$\begin{aligned} \hat{n}_1 &= \langle n_1, n_2, \dots, n_i, \dots, n_N \rangle \\ \hat{n}_{2,i} &= \langle n_1, \dots, n_i - 1, \dots, n_N \rangle \\ \hat{n}_{3,i} &= \langle n_1, \dots, n_i + 1, \dots, n_N \rangle \\ \hat{n}_{4,i,m} &= \langle n_1, \dots, n_i + 1, \dots, n_m - 1, \dots, n_N \rangle \end{aligned} \quad (10)$$

当一个新的云计算服务请求到达移动云计算服务域时, 如果这时系统的决策是拒绝, 那么有 $s = \langle \hat{n}_1, R \rangle$, 同时有 $a=0$; 或者当一个用户满意度为 i 的云计算服务结束运行, 该移动终端用户离开此移动云计算服务域并释放所占用的云计算资源时, 这时系统里面正在接受云计算服务的用户数量减少并且可用的云计算资源增加, 因此有 $s = \langle \hat{n}_1, D_i \rangle$.

在这两种情况下, 我们可以得到状态转移概率为

$$p(j | s, a) = \begin{cases} \frac{\lambda}{\gamma(s, a)}, & j = \langle \hat{n}_1, R \rangle \\ \frac{n_i \mu}{\gamma(s, a)}, & j = \langle \hat{n}_{2,i}, R \rangle, n_i \geq 1 \end{cases} \quad (11)$$

当一个新的云计算服务请求到达移动云计算服务域时, 如果系统这时的决策是同意接入并且准备为该云计算服务请求分配的云计算资源为 k_i 个 VM, 那么此时 $s = \langle \hat{n}_1, R \rangle$, 同时有 $a=i, i=1, 2, \dots, N$, 在这种情况下, 我们可以得到状态转移概率为

$$p(j | s, a) = \begin{cases} \frac{(n_i + 1) \mu}{\gamma(s, a)}, & j = \langle \hat{n}_1, D_i \rangle \\ \frac{\lambda}{\gamma(s, a)}, & j = \langle \hat{n}_{3,i}, R \rangle \\ \frac{n_m \mu}{\gamma(s, a)}, & j = \langle \hat{n}_{4,i,m}, D_m \rangle, n_m \geq 1, m \neq i \end{cases} \quad (12)$$

图 3 给出了基于我们所提出的移动云计算服务域动态云计算资源优化管理决策模型, $N=2$ 时的状态转移图. 图 3 中, 箭头线上第 1 项表示在当前状态下所采取的行动, 箭头线上第 2 项表示在当前状态下, 采取相应行动后, 转移到下一个状态的状态转移概率.

从公式(4)可知, 收益模型的支出是一个连续时间函数. 为了能够应用离散折扣马尔可夫决策模型, 需要对收益函数进行归一化处理并获得归一化后的长期期望收益. 如果我们能找到一个常数 ω , 使其满足 $[1 - p(s | s, a)] \gamma(s, a) \leq \omega < \infty$, 那么我们就可以得到归一化后的期望收益^[19].

因此, 我们可以找到一个常数 ω , 使其满足条件 $\omega = \lambda + \left[\frac{B}{u} \right] \cdot u < \infty$, 并且让 $\tilde{p}(j | s, a)$ 、 $\tilde{v}(s)$ 和 $\tilde{z}(s, a)$ 分别代表归一化后的状态转移概率、长期期望收益和收益函数. 从而我们可以得到归一化后的状态转移概率为

$$\tilde{p}(j | s, a) = \begin{cases} 1 - \frac{[1 - p(s | s, a)] \gamma(s, a)}{\omega}, & j = s \\ \frac{\gamma(s, a) p(j | s, a)}{\omega}, & j \neq s \end{cases} \quad (13)$$

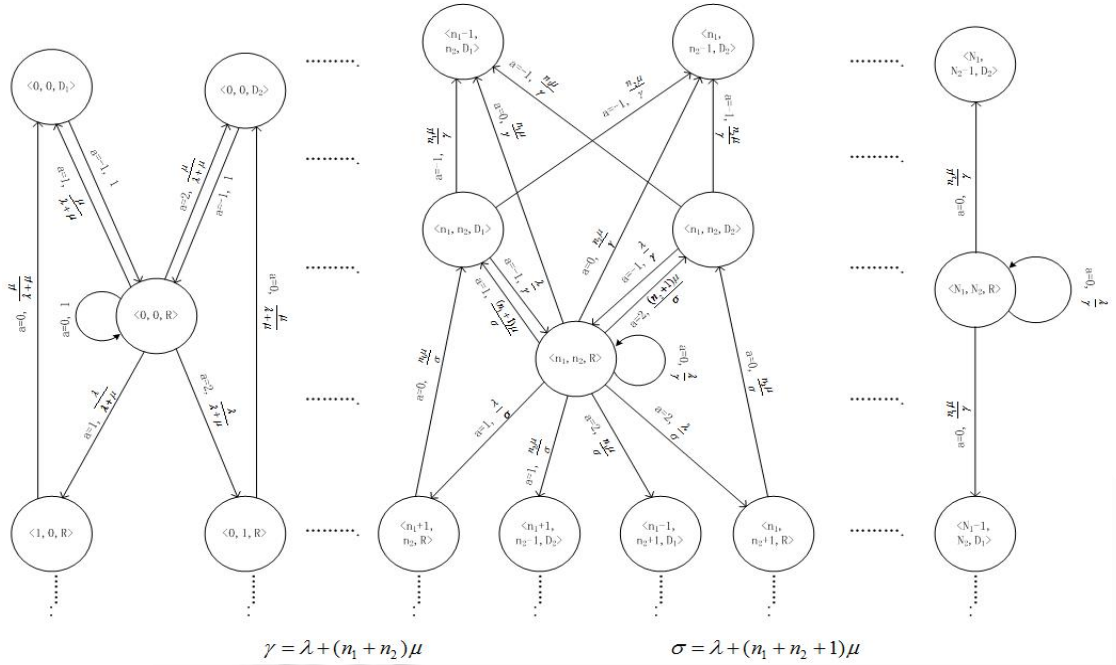


Fig.3 State transition diagram ($N=2$)

图 3 状态转移图($N=2$)

由此,在 $s = \langle \hat{n}_1, R \rangle$, 同时 $a=0$ 的状态下,状态转移概率 $\tilde{p}(j|s,a)$ 可以简化为

$$\tilde{p}(j|s,a) = \begin{cases} \frac{(\omega + \lambda - \gamma(s,a))}{\omega}, & j = \langle \hat{n}_1, R \rangle \\ \frac{n_i \mu}{\omega}, & j = \langle \hat{n}_{2,i}, D_i \rangle, n_i \geq 1 \end{cases} \quad (14)$$

同理,在 $s = \langle \hat{n}_1, R \rangle$, 同时 $a=i (i=1,2,\dots,N)$ 的状态下,状态转移概率 $\tilde{p}(j|s,a)$ 可以重写为

$$\tilde{p}(j|s,a) = \begin{cases} \frac{(n_i + 1)\mu}{\omega}, & j = \langle \hat{n}_1, D_i \rangle \\ \frac{\lambda}{\omega}, & j = \langle \hat{n}_{3,i}, R \rangle \\ \frac{n_m \mu}{\omega}, & j = \langle \hat{n}_{4,i,m}, D_m \rangle, n_m \geq 1, m \neq i \\ \frac{(\omega - \gamma(s,a))}{\omega}, & j = s \end{cases} \quad (15)$$

另外,对于状态 $s = \langle \hat{n}_1, D_i \rangle$, 归一化的状态转移概率为

$$\tilde{p}(j|s,a) = \begin{cases} \frac{\lambda}{\omega}, & j = \langle \hat{n}_1, R \rangle \\ \frac{n_i \mu}{\omega}, & j = \langle \hat{n}_{2,i}, D_i \rangle, n_i \geq 1 \\ \frac{(\omega - \gamma(s,a))}{\omega}, & j = s \end{cases} \quad (16)$$

因此,在归一化处理后,可以得到最大化的长期期望收益为

$$\tilde{v}(s) = \max_{a \in A(s)} \left\{ \tilde{z}(s,a) + \tilde{\eta} \sum_{j \in S} \tilde{p}(j|s,a) \tilde{v}(j) \right\} \quad (17)$$

其中, $\tilde{z}(s, a) = z(s, a) \frac{\gamma(s, a) + \alpha}{(\alpha + \omega)}$ 是归一化后的收益函数, 并且 $\tilde{\eta} = \frac{\omega}{\omega + \alpha}$.

3 模型的性能分析

对于移动云计算网络来说, 阻塞率是一个非常重要的服务质量(QoS)性能指标. 在本节里, 我们推导出了本文提出的基于 SMDP 的移动云计算服务域动态云计算资源优化管理模型对云计算服务请求的阻塞率和为移动终端分配不同云计算资源 VM 个数的概率.

我们让 π_s 表示移动云计算服务域中状态 s 的稳态概率. 由图 3 我们可以推导出移动云计算网络的稳态概率. 基于公式(11)和公式(12), 在计算稳态概率 $\pi_{\langle \hat{n}_1, e \rangle}$ 时, 必须考虑如下 3 种可能的情况:

- 1) 在一个云计算服务被完成后, 新到一个云计算服务请求;
- 2) 拒绝一个新到的云计算服务请求;
- 3) 接收一个新到的云计算服务请求.

这样, 我们可以计算出到达状态的稳态概率 $\pi_{\langle \hat{n}_1, R \rangle}$ 为

$$\pi_{\langle \hat{n}_1, R \rangle} = \hat{a}_{\langle \hat{n}_1, R \rangle} \cdot \frac{\lambda}{\gamma(s, a)} \cdot \pi_{\langle \hat{n}_1, R \rangle} + \frac{\lambda}{\gamma(s, a)} \cdot \sum_{i=1}^N \pi_{\langle \hat{n}_1, D_i \rangle} + \frac{\lambda}{\gamma(s, a)} \cdot \sum_{i=1}^N \hat{a}_{\langle \hat{n}_2, j, R \rangle} \cdot \pi_{\langle \hat{n}_2, j, R \rangle} \quad (18)$$

其中, $\hat{a}_{\langle \hat{n}_1, R \rangle}$ 和 $\hat{a}_{\langle \hat{n}_2, j, R \rangle}$ 分别可由下式得到

$$\hat{a}_{\langle \hat{n}_1, R \rangle} = \begin{cases} 1, & \hat{a}_{\langle \hat{n}_1, R \rangle} = 0 \\ 0, & \text{otherwise} \end{cases}$$

和

$$\hat{a}_{\langle \hat{n}_2, j, R \rangle} = \begin{cases} 1, & \hat{a}_{\langle \hat{n}_2, j, R \rangle} = i, \quad i = 1, 2, \dots, N \\ 0, & \text{otherwise} \end{cases}$$

同理, 我们可以计算出离去状态的稳态概率 $\pi_{\langle \hat{n}_1, D_i \rangle}$ 为

$$\begin{aligned} \pi_{\langle \hat{n}_1, D_i \rangle} &= \hat{a}_{\langle \hat{n}_3, j, R \rangle} \cdot \frac{(n_i + 1)\mu}{\gamma(s, a)} \cdot \pi_{\langle \hat{n}_3, j, R \rangle} + \frac{(n_i + 1)\mu}{\gamma(s, a)} \cdot \sum_{l=1}^N \pi_{\langle \hat{n}_3, j, D_l \rangle} \\ &+ \hat{a}_{\langle \hat{n}_1, R \rangle} \cdot \frac{(n_i + 1)\mu}{\gamma(s, a)} \cdot \pi_{\langle \hat{n}_1, R \rangle} + \frac{(n_i + 1)\mu}{\gamma(s, a)} \cdot \sum_{l=1, l \neq i}^N \hat{a}_{\langle \hat{n}_4, j, l, R \rangle} \cdot \pi_{\langle \hat{n}_4, j, l, R \rangle} \end{aligned} \quad (19)$$

其中, $\hat{a}_{\langle \hat{n}_3, j, R \rangle}$, $\hat{a}_{\langle \hat{n}_1, R \rangle}$ 和 $\hat{a}_{\langle \hat{n}_4, j, l, R \rangle}$ 可以分别由下式得到:

$$\begin{aligned} \hat{a}_{\langle \hat{n}_3, j, R \rangle} &= \begin{cases} 1, & \hat{a}_{\langle \hat{n}_3, j, R \rangle} = 0 \\ 0, & \text{otherwise} \end{cases} \\ \hat{a}_{\langle \hat{n}_1, R \rangle} &= \begin{cases} 1, & \hat{a}_{\langle \hat{n}_1, R \rangle} = i, \quad i = 1, 2, \dots, N \\ 0, & \text{otherwise} \end{cases} \\ \hat{a}_{\langle \hat{n}_4, j, l, R \rangle} &= \begin{cases} 1, & \hat{a}_{\langle \hat{n}_4, j, l, R \rangle} = l, \quad l \neq i \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

由于总概率等于 1, 我们可以得到:

$$\sum_s \pi_s = 1 \quad (20)$$

通过对公式(18)~公式(20)求解, 我们可以得到系统在任何状态下的稳态概率. 因此, 给移动云计算服务域中新到的云计算服务请求分配 k_i 个 VM 云计算资源的概率(用 $P_T^{k_i}$ 表示)可以表示为

$$P_T^{k_i} = \frac{\sum \pi_{\langle \hat{n}_1, R \rangle}}{\sum \pi_{\langle \hat{n}_1, R \rangle}}, \quad i = 1, 2, \dots, N \quad (21)$$

相应地,可以得到移动云计算服务域中云计算服务请求的阻塞率为

$$P_{blocking} = 1 - \sum_k P_T^k \tag{22}$$

5 模型的性能评估

在本节中,我们通过用 Matlab 编程的事件触发模拟器,对我们所提出的基于半马氏决策过程(SMDP)的移动云计算服务域动态云计算资源优化管理模型的性能进行了仿真和评估.在该仿真中,共划分了 3 个服务质量(QoS)级别,分别对应着分配云计算资源 VM 的数量,即 $k_1=1, k_2=2$ 和 $k_3=3$.如果没有另外说明,本章仿真中所使用的移动云计算服务域的云计算资源 VM 的总数为 10(即 $K=10$),云计算服务请求到达移动云计算服务域的均值速率为 $\lambda=7$,移动终端用户结束云计算服务并且释放所占用的云计算资源的均值速率分别为 $\mu=5$ 和 $\mu=10$.为了保证计算长期收益能收敛,我们设置模型中的折扣因子 α 为 0.1.另外,为了保证仿真的准确度,我们对每一个所测试性能的仿真时间定为 1800s.

3.1 优化行动

表 1 和表 2 列出了在不同效能函数的参数 ω_1 下,每个状态的云计算资源分配管理的优化决策.

Table 1 Optimal Decision Table I

表 1 优化决策表 I

$\omega_1 = 0.5, b = 1, \lambda = 7, \mu = 10, n_3 = 0$

$n_1 \setminus n_2$	0	1	2	3	4	5
0	3	3	3	2	2	0
1	3	3	3	2	1	0
2	3	3	2	2	0	0
3	3	3	2	1	0	0
4	3	2	2	0	0	0
5	3	2	1	0	0	0
6	2	2	0	0	0	0
7	2	1	0	0	0	0
8	2	0	0	0	0	0
9	1	0	0	0	0	0
10	0	0	0	0	0	0

Table 2 Optimal Decision Table II

表 2 优化决策表 II

$\omega_1 = 0.1, b = 1, \lambda = 7, \mu = 10, n_3 = 0$

$n_1 \setminus n_2$	0	1	2	3	4	5
0	2	2	2	2	2	0
1	2	2	2	2	1	0
2	2	2	2	2	0	0
3	2	2	2	1	0	0
4	2	2	2	0	0	0
5	2	2	1	0	0	0
6	2	2	0	0	0	0
7	2	1	0	0	0	0
8	2	0	0	0	0	0
9	1	0	0	0	0	0
10	0	0	0	0	0	0

在表中的数字表示在系统状态 $\langle n_1, n_2, n_3, R \rangle$ 时,本文所提出的移动云计算服务域动态云计算资源优化管理模型所做的优化决策.例如,在移动云计算服务域中,如果当前没有任何云计算服务在使用云计算资源,当一个新的云计算服务请求到达移动云计算服务域时,系统这时所做的决策为 $a=3$,即系统将会分配 k_3 个 VM 的云计算资源给新到的云计算服务请求.如果在系统中已经有 $n_2=3$ 个云计算服务在占用云计算资源,其中每个云计算服务占用 k_2 个 VM 的云计算资源,这就意味着该移动云计算服务域内还有 4 个 VM 未占用的云计算资源,那么当一个新的云计算服务请求到达该移动云计算服务域时,系统所采取的决策是 $a=2$,即分配 k_2 个 VM 的云计算资源给这个新到的云计算服务请求.当该移动云计算服务域中可用的云计算资源很充足时,为了获得更高的效能收益,我们的优化模型将会选择决策 $a=3$ 而不会选择决策 $a=1$ 和 $a=2$.另一方面来讲,当移动云计算服务域内可用的云计算资源减少时,我们的优化模型在做决策时会更保守一些.另外,从表中我们也可以注意到,当我们选取不同的效能函数的参数 ω_1 时,系统所做的决策也会因此而显著变化.例如,当 $\omega_1=0.1$ 时,对新到的云计算服务请求分配 k_1, k_2 和 k_3 个 VM 的云计算资源所对应获得的效能收益分别为 0.74427, 0.94257 和 0.98716.由于分配 k_2 和 k_3 个 VM 的云计算资源所带来的效能收益差别非常小,而另一方面系统选取决策 $a=3$ 时,会占用更多的云计算资源,从而增加系统的支出,因此我们的优化模型宁愿选择决策 $a=2$,而不会选择决策 $a=3$.

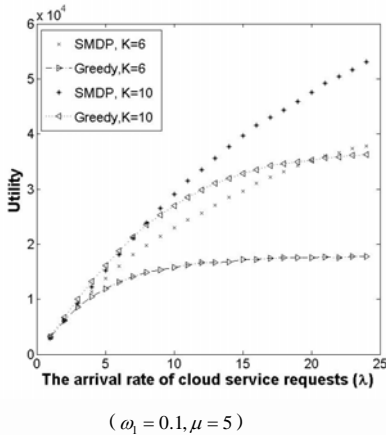
3.2 系统效能收益和阻塞率

为了进一步验证本文所提出的移动云计算服务域动态云计算资源优化管理模型的性能,我们将本文所提出的移动云计算服务域动态云计算资源优化管理模型的网络总效能收益和阻塞率分别与贪婪分配算法

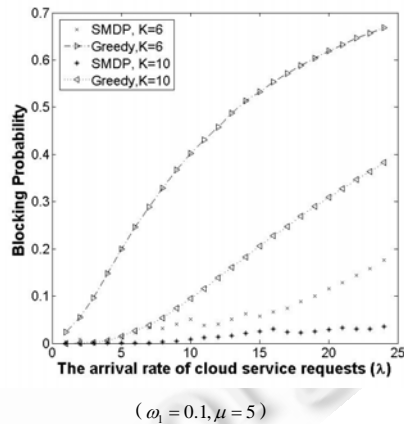
(Greedy Scheme)进行了比较.只要有足够的云计算资源,贪婪分配算法总是给云计算服务请求分配最大的云计算资源,以此来获得更高的系统效能收益.

图 4 比较了我们所提出的移动云计算服务域动态云计算资源优化管理模型与贪婪分配算法的网络系统效能收益的区别.从图中我们可以看到,随着云计算服务请求到达率 λ 的增加,系统效能收益也随之而增加,并且我们所提出的优化模型的性能明显优于贪婪分配算法的性能.这是因为,当一个云计算服务请求到达移动云计算服务域时,贪婪分配算法总是为该云计算服务请求分配最大的云计算资源,因而就存在这样的风险:当下一个云计算服务请求到达移动云计算服务域时,由于该移动云计算服务域云计算资源的不足,贪婪分配算法只能拒绝下一个到来的云计算服务请求.而我们所提出的移动云计算服务域动态云计算资源优化管理模型在一个新的云计算服务请求到达移动云计算网络时,即考虑了接收当前到来的云计算服务请求所带来的收入,同时也考虑了长期的系统预期收益,因此在给该云计算服务分配云计算资源时,会相对保守一些.

在图 5 中,我们可以看到,随着云计算服务请求到达率 λ 的增加,贪婪分配算法的云计算服务请求的阻塞率迅速增加,而我们所提出的移动云计算服务域动态云计算资源优化管理模型对云计算服务请求的阻塞率的增长却相对缓慢.这从云计算服务请求的阻塞率上表明了我们所提出的移动云计算服务域动态云计算资源优化管理模型的性能远远优于贪婪分配算法.



($\omega_1 = 0.1, \mu = 5$)
Fig.4 System utility
图 4 系统效能收益



($\omega_1 = 0.1, \mu = 5$)
Fig.5 Blocking probability of cloud service
图 5 云计算服务请求的阻塞率

我们在图 6 和图 7 中更进一步分析了当一个云计算服务请求到达移动云计算服务域时,本文所提出的移动云计算服务域动态云计算资源优化管理模型采取不同决策的概率.当云计算服务请求的到达率较低时,这样就有更多可用的云计算资源可以预留下来.因此在这种情况下,移动云计算服务域接收一个新到的云计算服务请求的概率就会较高.例如,当云计算服务请求的到达率为 $\lambda=3$ 时,系统分配 k_3 个 VM 云计算资源给该云计算服务请求的概率高达 96%.而当云计算服务的用户量增加时,更多的云计算资源被占用,从而为以后到来的云计算服务请求可预留的云计算资源就大为减少,这样就会导致云计算服务请求的阻塞率显著提高.

由于我们所提出的移动云计算服务域动态云计算资源优化管理模型需要考虑整个系统的长期收益,因此在为新到的云计算服务请求分配云计算资源就会更谨慎一些,所以在图 6 中可以看到,随着云计算服务请求到达率的增加,系统对新到的云计算服务请求采用决策 $a=3$ 的概率在减少,而采用决策 $a=1$ 和 $a=2$ 的概率则相应增加.图 7 则表示了移动云计算网络采取各个决策的概率与系统总的云计算资源之间的关系.从图中可以看出,随着云计算资源 VM 数量的增加,系统分配给新到云计算服务更多云计算资源的概率也随之而增加.如图 7 所示,当移动云计算服务域的云计算资源从 1 个 VM 增加到 10 个 VM 时,系统对新到云计算服务请求采用 $a=3$ 的决策的概率从 0 增加到了 96%.

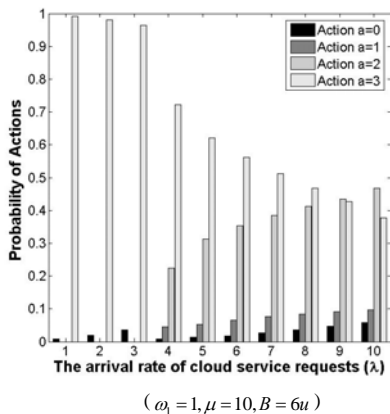


Fig.6 Probability of actions for cloud service

图6 对云计算服务请求采取各个决策的概率

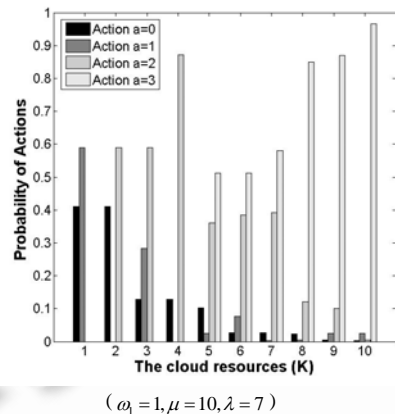


Fig.7 Probability of actions for cloud service

图7 对云计算服务请求采取各个决策的概率

4 结论与展望

本文主要研究了在移动云计算网络中,对云计算服务进行动态分配云计算资源的优化问题.我们基于半马尔可夫决策过程(SMDP)来对动态云计算资源优化分配管理问题进行建模,并且新提出了移动云计算服务域动态云计算资源优化管理模型.根据系统现有可用的云计算资源,该模型在既考虑网络效能收益,又考虑因占用云计算资源所带来的开销的基础上,对新到的云计算服务请求做出了一个优化决策策略,使其能保证移动云计算服务域的系统长期收益最大.同时通过动态分配不同的云计算资源给新到的云计算服务请求,我们所提出的移动云计算服务域动态云计算资源优化管理模型能为云计算服务提供多个不同的服务质量(QoS)级别.我们对所提出的移动云计算服务域动态云计算资源优化管理模型进行了进一步推导,得出了该优化管理模型的重要服务指标——阻塞率,同时也推导出了系统对新到云计算服务请求采取不同决策的概率.

我们以后的工作将主要集中在如下两方面:

1) 通过运用不同类型的效能函数,从而构建更加复杂的决策模型.

2) 在给定云计算服务质量(QoS)的条件下,基于本文的研究成果,构建一个新的动态云计算资源优化管理模型,使其能用最小的云计算资源来实现移动云计算网络的最大整体收益.

References:

- [1] Armbrust M, Fox A, Griffith R, Joseph A, Katz R, Konwinski A, Lee G, Patterson D, Rabkin A, Stoica I, *et al.* Above the clouds: A Berkeley view of cloud computing. Tech. Rep. UCB/EECS-2009-28, EECS Department, University of California, Berkeley, 2009.
- [2] Walshy M. Gartner: Mobile to outpace desktop Web by 2013.
- [3] Li XH, Zhang H, Zhang YF. Deploying mobile computation in cloud service. In: Proc. of the 1st Int'l Conf. for Cloud Computing (CloudCom). 2009. 301.
- [4] Chun B, Maniatis P. Augmented smartphone applications through clone cloud execution. In: Proc. of the USENIX HotOS XII. 2009.
- [5] Zhang X, Schiffman J, Gibbs S, Kunjithapatham A, Jeong S. Securing elastic applications on mobile devices for cloud computing. In: Proc. of the 2009 ACM Workshop on Cloud Computing Security. 2009. 127-134.
- [6] Huang D, Zhang X, Kang M, Luo J. Mobicloud: A secure mobile cloud framework for pervasive mobile computing and communication. In: Proc. of the 5th IEEE Int'l Symp. on Service-Oriented System Engineering. 2010.
- [7] Meng X, Pappas V, Zhang L. Improving the scalability of data center networks with traffic-aware virtual machine placement. In: Proc. of the IEEE INFOCOM. 2010.
- [8] Cai LX, Cai L, Shen X, Mark JW. Resource management and QoS provisioning for IPTV over mmWave-based WPANs with directional antenna. ACM Mobile Networks and Applications (MONET), 2009,14(2):210-219.

- [9] Cheng HT, Zhuang W. Novel packet-level resource allocation with effective QoS provisioning for wireless mesh networks. *IEEE Trans. on Wireless Communications*, 2009,8(2):694–700.
- [10] Cai LX, Shen X, Mark JW. Efficient MAC protocol for ultrawideband networks. *IEEE Communications Magazine*, 2009,47(6):179–185.
- [11] Liang H, Huang D, Peng D. On economic mobile cloud computing model. In: *Proc. of the Int'l Workshop on Mobile Computing and Clouds (MobiCloud in conjunction with MobiCASE)*. 2010.
- [12] Wei G, Vasilakos AV, Zheng Y, Xiong N. A game-theoretic method of fair resource allocation for cloud computing services. 2009.
- [13] Lorincz K, Chen BR, Waterman J, Werner-Allen G, Welsh M. Resource aware programming in the pixie OS. In: *Proc. of the SenSys 2008*. 2008.
- [14] Lorincz K, Chen B., Waterman J, Werner-Allen G, Welsh M. A stratified approach for supporting high throughput event processing applications. In: *Proc. of the DEBS 2009*. 2009.
- [15] Tesauro G, Jong NK, Das R, Bannani MN. A hybrid reinforcement learning approach to autonomic resource allocation. In: *Proc. of the ICAC 2006*. 2006.
- [16] Bloor K, Chirkova R, Viniotis R, Salo T. Dynamic request allocation and scheduling for context aware applications subject to a percentile response time sla in a distributed cloud. In: *Proc. of the 2nd IEEE Int'l Conf. on Cloud Computing Technology and Science*. 2010.
- [17] Lee JW, Mazumdar RR, Shroff NB. Non-Convex optimization and rate control for multi-class services in the Internet. *IEEE/ACM Trans. on Networking*, 2005,13(4):827–840.
- [18] Mine H, Osaki S, Puterman ML. *Markovian Decision Process*. Amsterdam: Elsevier, 1970.
- [19] Puterman ML. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2005.
- [20] Cui YF, Wang S, Li Y, Dong KW. Cloud computing resource management research based on multi-level hierarchical control structure. *Journal of Academy of Equipment Command & Technology*, 2010,21(2) (in Chinese with English abstract).
- [21] Yuan WC, Zhu YA, Lu W. Exploring virtualized resource management mechanism for cloud computing. *Journal of Northwestern Polytechnical University*, 2010,28(5) (in Chinese with English abstract).
- [22] Jiang Q, Xi HS, Yin BQ. An online adaptive bandwidth allocation optimization algorithm for wireless multimedia communication networks. *Journal of Software*, 2007,18(6):1491–1500 (in Chinese with English abstract). <http://www.jos.org.cn/1000-9825/18/1491.htm> [doi: 10.1360/jos181491]

中文参考文献

- [20] 崔云飞,王帅,李艺,董可为.基于多级递阶控制结构的云计算资源管理研究.装备指挥技术学院学报,2010,21(2).
- [21] 袁文成,朱怡安,陆伟.面向虚拟资源的云计算资源管理机制.西北工业大学学报,2010,28(5).
- [22] 江琦,奚宏生,殷保群.无线多媒体通信网适应带宽配置在线优化算法.软件学报,2007,18(6):1491–1500. <http://www.jos.org.cn/1000-9825/18/1491.htm> [doi: 10.1360/jos181491]



梁宏斌(1972—),男,四川成都人,博士,高级工程师,主要研究领域为移动云计算架构,云安全,云存储,资源优化管理.



刘燕(1971—),女,博士,副教授,主要研究领域为计算机网络,软件工程.



彭代渊(1955—),男,教授,博士生导师,主要研究领域为信息与编码,信息安全.