

基于云计算的可反馈负载均衡策略的研究

孟 蒙^{1,2,3}, 茅 苏^{1,2,3}

- (1. 南京邮电大学 计算机学院, 江苏 南京 210023;
2. 江苏省无线传感网高技术研究重点实验室, 江苏 南京 210023;
3. 宽带无线通信与传感网技术教育部重点实验室, 江苏 南京 210023)

摘 要:现有负载均衡算法存在仅考虑单一云环境资源利用率情况或缺少云环境多维资源负载监控机制问题。文中基于现有的云环境中的负载均衡技术,量化了云环境中多种资源的使用情况,定义了云计算中不同的负载均衡类型,提出了一种可反馈的负载均衡策略及双监控器策略,均衡负载,避免反复创建、释放虚拟机导致时间、资源、性能的损耗,从而提高资源的利用率。CloudSim 仿真实验结果显示,文中提出的可反馈的负载均衡监控机制,能够更有效地监控资源使用情况以及提高资源的利用率。

关键词:负载均衡;云计算;虚拟机;资源利用率;CloudSim

中图分类号:TP301

文献标识码:A

文章编号:1673-629X(2014)10-0135-05

doi:10.3969/j.issn.1673-629X.2014.10.032

Study on Feedback Load Balancing Strategy Based on Cloud Computing

MENG Meng^{1,2,3}, MAO Su^{1,2,3}

- (1. School of Computer, Nanjing University of Posts and Telecommunications,
Nanjing 210023, China;
2. Jiangsu High Technology Research Key Laboratory for Wireless Sensor Networks,
Nanjing 210023, China;
3. Key Lab of Broadband Wireless Communication and Sensor Network Technology of
Ministry of Education, Nanjing 210023, China)

Abstract: Existing load balancing algorithm has considered a single cloud resource utilization only or a lack of cloud environments multi-dimensional load monitoring mechanism. Based on existing load balancing technology in cloud computing environment, quantify the use of a variety of resources, define different types of load balancing in cloud computing, present a new feedback load balancing strategy and dual monitor strategy, which is for balance load, avoiding repeated creating and releasing the virtual machine to result in loss of time, resources, performance, improving resource utilization. The CloudSim simulation results show that the feedback load balancing strategy can more effectively monitor the use of resources and improve resource utilization.

Key words: load balancing; cloud computing; virtual machines; resource utilization; CloudSim

0 引 言

近年来,具有抽象和封装特性的虚拟化技术成为云资源和调度管理中的关键技术^[1]。传统的针对物理实体资源的并行计算^[2]、分布式计算^[3]和网格计算^[4]中负载均衡技术并不能很好地适用于云计算负载均衡

的需要。

为减少云计算环境中虚拟机 HotPots^[5]的现象,提高用户满意度,提高云资源利用率和降低能耗减少碳排放^[6],云计算中的负载均衡技术成为了研究热点。文献[7-10]针对不同的问题和应用场景提出了多种负载均衡机制,但文献中提到的负载均衡算法存在仅

收稿日期:2013-12-09

修回日期:2014-03-12

网络出版时间:2014-07-28

基金项目:江苏省高校优势学科建设工程资助项目(yx002001)

作者简介:孟 蒙(1989-),男,硕士研究生,研究方向为基于网络的计算机软件应用技术;茅 苏,高级工程师,研究方向为基于网络的计算机软件应用技术。

网络出版地址: <http://www.cnki.net/kcms/detail/61.1450.TP.20140728.1228.047.html>

考虑单一的云环境资源利用率情况或缺少云环境负载监控机制问题。文中首先量化了云环境中多种资源的使用情况,定义了负载不均衡的不同状态,提出了一种可反馈负载均衡监控策略,均衡负载,避免因反复创建、释放、迁移虚拟机导致的时间、资源、性能的损耗,从而提高资源的利用率。通过 CloudSim 仿真实验结果表明,文中提出的可反馈的负载均衡监控机制,能够更有效地监控资源使用情况以及提高资源的利用率。

1 负载均衡策略理论模型

1.1 资源量化定义

首先对云环境中的资源做如下量化定义。

定义1: $V_{i \text{ conf}}(R_{i \text{ cpu}}, R_{i \text{ mem}}, R_{i \text{ sto}}, R_{i \text{ band}})$, $V_i(R_{i \text{ cpu}}, R_{i \text{ mem}}, R_{i \text{ sto}}, R_{i \text{ band}})$ 分别表示虚拟机 V_i 产生时资源的配置情况和使用时的实时利用率。其中, $R_{i \text{ cpu}}$, $R_{i \text{ mem}}$, $R_{i \text{ sto}}$, $R_{i \text{ band}}$ 依次表示虚拟机 V_i 的 CPU、内存、存储和带宽资源的配置情况和虚拟机 V_i 的 CPU、内存、存储和带宽的实时利用率情况。因此,可以根据 $V_{i \text{ conf}}$ 的各项资源配置值设定计算 V_{ih} 和 V_{il} 的各项高低告警阈值。计算公式如下所示。

$$\begin{aligned} V_{ih}(R_{i \text{ cpu}}, R_{i \text{ mem}}, R_{i \text{ sto}}, R_{i \text{ band}}) &= \\ \alpha_h * V_{i \text{ conf}}(R_{i \text{ cpu}}, R_{i \text{ mem}}, R_{i \text{ sto}}, R_{i \text{ band}}) \\ V_{il}(R_{i \text{ cpu}}, R_{i \text{ mem}}, R_{i \text{ sto}}, R_{i \text{ band}}) &= \\ \alpha_l * V_{i \text{ conf}}(R_{i \text{ cpu}}, R_{i \text{ mem}}, R_{i \text{ sto}}, R_{i \text{ band}}) \\ R_{i \text{ <cpu mem sto band> h}} &= \\ \alpha_h * R_{i \text{ <cpu mem sto band> h}} &= \\ \alpha_l * R_{i \text{ <cpu mem sto band> l}} &= \end{aligned} \quad (1)$$

其中, α_h, α_l 为常数, $0 < \alpha_l < \alpha_h < 100\%$, 大小根据系统的负载情况实时动态调整设定。

定义2: $H_{i \text{ conf}}(\text{HR}_{i \text{ cpu}}, \text{HR}_{i \text{ mem}}, \text{HR}_{i \text{ sto}}, \text{HR}_{i \text{ band}})$ 表示当前物理机 H_i 配置的资源情况。其中, $\text{HR}_{i \text{ cpu}}$, $\text{HR}_{i \text{ mem}}$, $\text{HR}_{i \text{ sto}}$ 和 $\text{HR}_{i \text{ band}}$ 表示物理机 H_i 当前配置的 CPU、内存、存储和带宽资源。

同样可以定义 H_i 为物理机当前资源利用情况和 H_{ih} , H_{il} 为物理机 H_i 的资源利用率达到的高低告警阈值。

在云计算环境中,为了充分利用资源和确保用户程序的安全性,通常一个或多个虚拟机同时映射到单个物理实体机上^[11],因此,物理机和虚拟机需要满足条件(2):

$$\sum_{i=0}^n V_{i \text{ conf}} \leq H_{i \text{ conf}} \quad (2)$$

其中, $\sum_{i=0}^n R_{i \text{ (cpu mem sto band)}} \leq \text{HR}_{i \text{ <cpu mem sto band>}}$; n 表示物理机 H 上面运行的虚拟机数量。

根据文献[9]中的链式负载,可以把某台物理机 H_i 的负载均衡度用公式(3)表示:

$$\omega_i = \frac{\sum_{j=0}^n (V_{ij} - \bar{V})}{n} \quad (3)$$

其中, $\bar{V} = \frac{\sum_{j=0}^n V_{ij}}{n}$, n 表示物理机上运行的虚拟机个数。

ω_i 值越小,说明物理机 H_i 的负载均衡度越好。

也可以把整个云环境的某一资源负载均衡度用公式(4)表示:

$$\mu = \frac{\sum_{i=0}^m (\omega_i - \bar{\omega})}{m} \quad (4)$$

其中, $\bar{\omega} = \frac{\sum_{i=0}^m \omega_i}{m}$, m 表示云环境中运行的物理机个数。

μ 值越小,说明云环境负载均衡度越好。

1.2 采样方法说明

为了避免瞬时负载峰值^[12]现象,文中采用平均采样法来解决这一问题。平均采样法的采样公式如公式(5)所示:

$$R_{i \text{ (cpu mem sto band)}} = \frac{\sum_{i=0}^t R_{i \text{ (cpu mem sto band)}}}{n} \quad (5)$$

其中, n 为采样个数。

在文献[13]中,建议一般可将采集周期设置在 1 ~ 11 s 之间。

1.3 负载状态分类

根据云环境中二级链式负载^[9],可以把云环境中负载均衡状态分为两大类,虚拟机负载状态和物理机负载状态。其中,虚拟机负载状态根据不同的负载状态可以分为以下5类。

(1) 虚拟机单个资源负载不均衡状态:在虚拟机集群中,存在一个虚拟机 V_i 的 CPU、内存、存储和带宽资源利用率之一小于或大于虚拟机 V_i 初始化时设定的低或高告警阈值,并且其他虚拟机 V_j 的资源利用率都满足均衡条件。

(2) 虚拟机多个资源负载不均衡状态:在虚拟机集群中,一个虚拟机 V_i 至少有一种资源满足虚拟单个负载不均衡状态。

针对虚拟机单个资源负载不均衡现象,可以采用文献[14]中提出的自发调节和全局调节之间协作的算法解决虚拟机内存资源不均衡的问题。在虚拟机对应的实体机资源充足的情况下,动态添加虚拟机的该

项资源。在虚拟机对应的实体机资源不足的情况下,可以通过文中提出的动态降低 R_l 或者提高 R_h 的值,使虚拟机该项资源达到相对均衡的状态。

针对虚拟机多个资源负载不均衡现象,同样采取动态地添加虚拟机的各项资源,降低资源利用率或降低 R_l ,提高 R_h 的值使资源利用达到均衡状态。

(3) 单个虚拟机负载不均衡状态:在虚拟机集群中,存在虚拟机 V_i 的 CPU、内存、存储和带宽资源同时小于虚拟机初始化时设定的低告警阈值或大于虚拟机初始化时设定的高告警阈值,且同一物理机部署的其他任意虚拟机 V_j 满足均衡条件。

针对单个虚拟机 V_i 负载过低现象,可以通过提高虚拟机 V_i 对用户请求响应的权限,降低其他虚拟机 V_j 对用户请求响应的权限。另外,还可以通过降低虚拟机 V_i 的 R_l ,提高其他虚拟机 V_j 的 R_h 来达到资源利用率的相对均衡状态。

针对单个虚拟机 V_i 负载过高现象,可以通过使用与单个虚拟机 V_i 过低相反的方法来达到资源利用率的相对均衡状态。

(4) 虚拟机集群整体负载不均衡状态:在虚拟机集群中,多个虚拟机 V_i 的 CPU、内存、存储和带宽资源小于或大于虚拟机初始化时设定的低或高告警阈值,称为虚拟机集群达到整体负载过低或过高状态。

虚拟机集群过低状态虽然属于一种负载均衡状态,但是云资源没有得到充分利用,为了提高资源利用率,可以选择关闭虚拟机,回收虚拟机的资源,降低能耗。

虚拟机集群过高状态也属于一种负载均衡状态,云资源得到了充分利用,但是会带来响应用户请求时间过长,致使云系统整体性能下降。如果在这种状态的基础上,实体机资源充足的情况下,可以选择动态映射新的虚拟机,降低平均负载。

(5) 虚拟机集群负载均衡状态:在虚拟机集群中,任意虚拟机 V_i 的 CPU、内存、存储和带宽资源利用率都处于虚拟机 V_i 初始化时设定的高、低告警阈值之间。

物理机的负载均衡状态分类与虚拟机类似,在此不再赘述。

2 可反馈的负载均衡策略设计

2.1 总体架构设计

参考文献[15]中提出的一种基于绿色计算资源池策略的云环境弹性负载均衡机制,和文献[10]提出的基于阈值的动态资源分配调度策略中的阈值的方法,同时,基于前文论述的云环境中负载不均衡的分类,文中提出了一种可反馈的负载均衡策略,通过监控云环境中的资源利用率情况,判断云节点属于什么类

型负载均衡状态,选择相应的策略,使云环境达到相对负载均衡的状态。

可反馈的负载均衡策略模型如图 1 所示,主要由配置管理器、监控机、信息收集器三个主要模块组成。

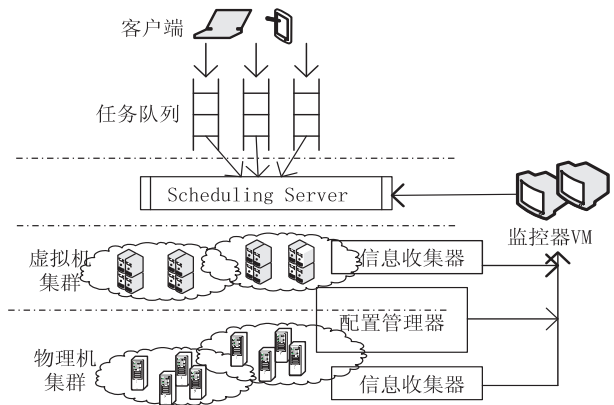


图 1 可反馈负载均衡策略模型

下面对各个模块进行定义。

(1) 配置管理器:主要负责管理各物理机和虚拟机当前配置信息,根据当前负载情况修改各虚拟机的高低阈值和采样频率,备份一份各物理机和虚拟机最近的资源利用率清单,启动和停止备份监控机。

(2) 监控器:主要负责根据采样频率,定期收集各物理机和虚拟机当前资源利用率,根据平均采样法得出的资源利用率和配置管理器中的高低阈值做比较,把比较结果反馈到调度服务器,定期把各物理机和虚拟机的资源利用率备份到配置管理器模块。

(3) 信息收集器:主要负责收集虚拟机和物理机的资源利用率情况,把收集结果反馈到 VM 监控器中。

2.2 可反馈负载均衡模型算法过程

可反馈负载均衡模型主要通过配置管理器、监控器、信息收集器配合实现,其模型流程图如图 2 所示。

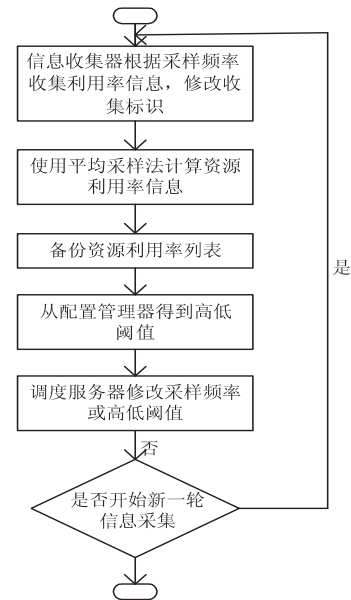


图 2 可反馈负载均衡流程图

可反馈的负载均衡策略模型具体过程如下:

Step1:位于每个虚拟机节点上的信息收集器根据内部信息是否已经被监控器采集标识,进行新一轮的信息收集,同时修改收集标识,以备下一轮信息收集;

Step2:监控器每隔一段时间从信息收集器中采集负载信息,同时修改存在信息收集器中的标识,对采集的信息利用采样平均法,得到一组最新的资源利用率列表,定期把资源利用率列表更新到配置管理器中,并从配置管理器中得到最新的高低阈值;

Step3:监控器备份资源利用列表;

Step4:得到最新各物理机和虚拟机高低阈值列表后,监控器把最新资源利用率列表与阈值高低列表进行对比,把对比结果发送给调度服务器;

Step5:调度服务器根据对比结果动态调整配置管理器中的高低阈值,采样频率,关闭或初始化新虚拟机以及添加或关闭物理机等相关措施;

Step6:是否开始新一轮的信息采集? 是,重复 Step1 ~ Step5 过程,使云系统维持在相对负载平衡的状态。

2.3 可反馈负载算法优化策略

2.3.1 动态采集周期

在文献[13]中,建议一般可将采集周期设置在 1 ~ 11 s 之间,过大的采集间隔会影响精确性,较小的采集间隔会给系统带来额外开销。文中采取动态设置的方式控制采集周期,调度服务器根据用户请求任务队列长度大小,和监控器发来的比较结果列表动态地调整采集间隔的大小。当用户任务队列长度小,说明负载低时,缩短采集周期;当用户任务队列长度大,说明负载高时,延长采集周期。

2.3.2 可反馈高低阈值策略

文献[10]同样提出了使用阈值对于资源分配过程中出现的负载不均衡的现象的解决方案,但是由于阈值的计算粒度过大,而且阈值不可动态调控,极易造成不恰当的重新分配,最终导致资源的浪费。文中提出的可反馈的高低阈值,可以根据系统的负载状况动态调整物理机或虚拟机的高低阈值,对于第一部分的几种负载不均衡定义都可以动态地调整高低阈值,使云环境达到认可下的“负载不均衡”,避免不必要的虚拟机迁移所带来的问题。

2.3.3 双 VM 监控器机制

为了避免 VM 监控器成为负载均衡策略的瓶颈或者故障的发生,提出了双 VM 监控器机制,一个为主 VM 监控器,另一个为辅 VM 监控器。为了解决系统瓶颈问题,二者也可以随时从不同节点的信息收集器收集性能信息,计算后反馈给调度服务器,并且共享配置管理器中的配置,高低阈值,采样频率等信息。另

外,在负载较小的情况下,如果同时使用两台 VM 监控器,对于资源利用率信息采集的精度比单台 VM 监控器更具优势。假设一个监控机发生故障的概率为 p ,那么两台监控机同时发生故障的概率就是 p^2 。

3 仿真与结果分析

3.1 实验参数设置

文中提出的负载均衡监控策略使用 CloudSim^[16] 模拟器进行仿真。10 个数据中心,每个数据中心里主机数目为 15 台,主机处理器均为单核处理器,处理速度分别为 1 000,2 000 或 3 000 MIPS,内存 5 G,外存 1 TB,带宽 10 M。每台主机需要部署的虚拟机数目为 9 个,每个虚拟机需要的 CPU 处理速度为 250,500,750 或 1 000 MIPS,内存 256 M,外存 1 GB。2 个用户,每个用户提交的任务数目为 300 个不同类型的任务。 α_h, α_v 参数分别为 75% 和 15%,随着云环境中物理机和虚拟机负载的动态变化相应地依次把 α_v, α_h 分别在 5% ~ 50% 和 50% ~ 100% 之间进行调整,使云环境中物理机和虚拟机达到相对的负载均衡状态。

3.2 实验结果及分析

文献[17]提出利用过载拒绝、容错率、预测的精确性、稳定性、进程迁移、资源利用率等因素来评价一个负载均衡策略的性能。文中提出的可反馈的负载均衡策略更注重对于资源利用率以及各节点的负载均衡情况。因此,仿真实验对比了使用可反馈的负载均衡策略与未使用可反馈的负载均衡策略的资源利用率和各节点负载均衡度的对比。

当用户提交的任务全部为计算型时,随着用户任务提交数量的增加,CPU 负载不均衡 VM 告警数情况如图 3 所示。随着任务数的增加,运行可反馈负载均衡机制的云环境,VM 告警数由于进行了阈值调整,动态地添加或减少 VM 的 CPU 资源,VM 的低告警数在逐渐降低,随着任务数的稳定,VM 低告警数趋于 0,说明云环境中没有空闲的虚拟机,而且虚拟机中的资源

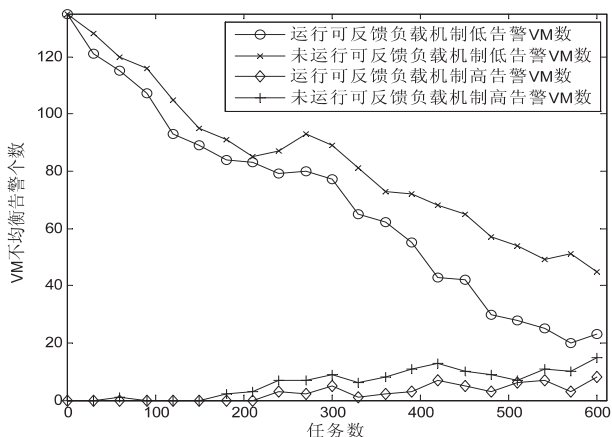


图 3 CPU 负载均衡 VM 告警数图

得到了充分的利用;相反,未运行可反馈负载均衡机制的云环境,随着任务数的增加,VM 低告警数趋于 0 的幅度小于运行可反馈负载机制的幅度,说明在运行的过程中存在虚拟机资源浪费的情况。另外,运行可反馈负载均衡机制的云环境中,出现高告警的 VM 数量维持在 0 到 1 个,而未运行可反馈负载均衡机制的云环境中,随着任务数量的增加,高告警的 VM 数量递增,同时,也出现了低负载告警的 VM,说明存在 VM 出现负载过重的同时有些 VM 出现负载过低现象,造成了资源的浪费。

根据公式(3)度量负载均衡度公式,某一资源负载均衡度越小,说明负载越均衡。在云环境运行的过程中选取了 21 个时间点,计算其 135 台 VM 的 CPU 资源利用率平均负载均衡度,如图 4 所示。随着任务数和时间的增加,运行可反馈负载机制的 VM 负载均衡度要小于未运行可反馈负载机制的 VM 负载均衡度,因此,可以得知运行可反馈负载机制的每个 VM 的负载均衡,很好地提高了云环境资源利用率。

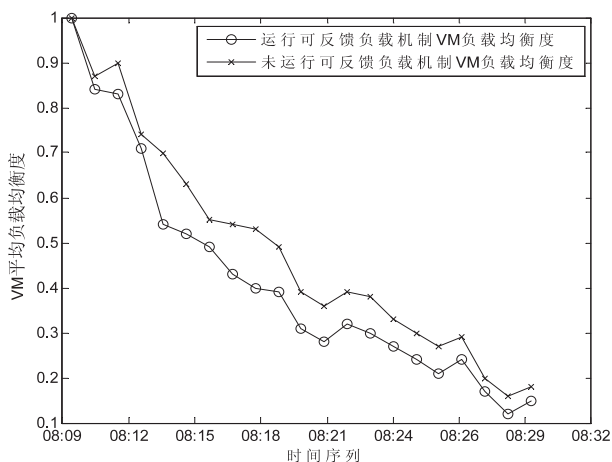


图 4 135 台 VM 平均负载均衡度图

4 结束语

文中提出的可反馈的负载均衡策略可以避免反复创建、释放虚拟机带来的时间、资源、性能的损耗,从而提高资源的利用率,一定程度上解决了云环境中的负载均衡问题。然而,随着云环境中物理或虚拟机的数量增加,云环境会变得更加复杂。文中未来的研究工作将从以下方面进行:首先,继续优化可反馈的双监控器负载均衡策略,解决因需要通信而造成的通信瓶颈问题;其次,在任务调度或者资源调度时,考虑到可反馈的双监控器负载均衡策略,从调度管理的角度实现负载均衡。

参考文献:

[1] Foster I, Zhao Yong, Raicu I, et al. Cloud computing and grid computing 360-degree compared[C]//Proc of grid computing

environments workshop. Austin, TX: IEEE, 2008: 1-10.

- [2] 李鸿健, 豆育升, 唐红, 等. 一种可变周期反馈的动态负载均衡算法[J]. 电子测量与仪器学报, 2011, 25(11): 952-958.
- [3] Alakeel A M. A guide to dynamic load balancing in distributed computer systems[J]. International Journal of Computer Science and Information Security, 2010, 10(6): 153-160.
- [4] 蒋江, 张民选, 廖湘科. 基于多种资源的负载平衡算法的研究[J]. 电子学报, 2002, 30(8): 1148-1152.
- [5] Wood T, Shenoy P, Venkataramani A, et al. Sandpiper: black-box and gray-box resource management for virtual machines[J]. Computer Networks, 2009, 53(17): 2923-2938.
- [6] Kansal N J, Chana I. Cloud load balancing techniques: a step towards green computing[J]. International Journal of Computer Science Issues, 2012, 9(1): 238-246.
- [7] Hu Jinhua, Gu Jianhua, Sun Guofei, et al. A scheduling strategy on load balancing of virtual machine resources in cloud computing environment[C]//Proc of third international symposium on parallel architectures, algorithms and programming. Los Alamitos: IEEE, 2010: 89-96.
- [8] Clark C, Fraser K, Hand S, et al. Live migration of virtual machines[C]//Proceedings of the 2nd symposium on networked systems design & implementation. [s. l.]: USENIX Association, 2005: 273-286.
- [9] 刘之家. 一种基于云计算的负载均衡技术的研究[J]. 广西师范学院学报: 自然科学版, 2011, 28(2): 93-96.
- [10] Lin W, Wang J Z, Liang C, et al. A threshold-based dynamic resource allocation scheme for cloud computing[J]. Procedia Engineering, 2011, 23: 695-703.
- [11] Fang Y, Wang F, Ge J. A task scheduling algorithm based on load balancing in cloud computing[M]//Web information systems and mining. Berlin: Springer, 2010: 271-277.
- [12] Zhao Yi, Huang Wenlong. A daptive distributed load balancing algorithm based on live migration of virtual machines in cloud[C]//Proc of fifth international joint conference on INC, IMS and IDC. [s. l.]: [s. n.], 2009: 170-175.
- [13] 陈亮. 集群负载均衡关键技术研究[D]. 长沙: 中南大学, 2009.
- [14] 张伟哲, 张宏莉, 张迪, 等. 云计算平台中多虚拟机内存协同优化策略研究[J]. 计算机学报, 2011, 34(12): 2265-2277.
- [15] 杜垚, 郭涛, 陈俊杰. 云环境下机群弹性负载均衡机制[J]. 计算机应用, 2013, 33(3): 830-833.
- [16] Calheiros R N, Ranjan R, de Rose C A F, et al. CloudSim: a novel framework for modeling and simulation of cloud computing infrastructures and services[R]. Parkville, VIC: Grid Computing and Distributed Systems Laboratory, the University of Melbourne Australia, 2009.
- [17] Sharma S, Singh S, Sharma M. Performance analysis of load balancing algorithms[J]. World Academy of Science, Engineering and Technology, 2008, 38: 269-272.

基于云计算的可反馈负载均衡策略的研究

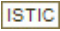
作者：

孟蒙，茅苏，[MENG Meng](#)，[MAO Su](#)

作者单位：

[南京邮电大学 计算机学院，江苏 南京210023；江苏省无线传感网高技术研究重点实验室，江苏 南京210023；宽带无线通信与传感网技术教育部重点实验室，江苏 南京210023](#)

刊名：

[计算机技术与发展](#)

英文刊名：

[Computer Technology and Development](#)

年，卷(期)：

2014(10)

本文链接：http://d.wanfangdata.com.cn/Periodical_wjfz201410033.aspx