doi:10.3772/j.issn.1002-0470.2021.12.005

基于半休眠模式的云计算任务卸载策略及性能研究①

宋 家② 余 靖③ 金顺福

(燕山大学信息科学与工程学院 秦皇岛 066004)

摘 要 云系统中的任务卸载是一种将密集的移动应用程序迁移到云端计算的机制。考虑云用户服务质量和云系统能耗水平,在云端引入周期性同步半休眠模式,提出一种云计算任务卸载策略。针对不同负载情况下的虚拟机服务速率,构建多服务台同步多重工作休假排队模型。基于拟生灭过程和矩阵几何解方法,得到云用户的平均响应时间与云系统的能源消耗。进行系统实验,评估卸载策略的系统性能,并验证其有效性。基于权重分配构建系统成本函数,改进海鸥优化算法,给出云计算任务卸载策略的智能优化方案,实现系统成本的最小化。

关键词 云计算;任务卸载;半休眠;工作休假排队;智能优化

0 引言

随着智能手机的大量普及和移动通信网络的迅 猛发展,各类网络应用程序,如交互式游戏、虚拟现 实、自然语言处理等层出不穷。这些应用程序通常 是计算密集型和高耗能的,不适合在计算能力和电 池供应有限的移动设备上工作。为了满足复杂应用 程序的资源需求,移动云计算的概念应运而生。如 何设计同时保证云用户的服务质量和云服务供应商 运营利润的云计算卸载策略成为热门研究问题^[12]。

过去的几年中,云数据中心的能源消耗持续增 长,云数据中心消耗的大量能量对环境造成了不利 影响。文献[3]提出了以节省能量消耗为目标的云 资源调度优化管理程序,利用虚拟机监控程序分配 资源,实现云计算资源分配的最小能耗和最高效率。 文献[4]为了解决云用户的服务质量和能耗等问题 提出了一种云计算环境下动态资源分配模型,使用 间距多目标蚁群算法将云系统状态的动态变化与任 务调度相结合,通过调整未使用的虚拟机,节省云系 统的能源消耗。上述文献均采用资源优化配置的方 式节省云系统能源消耗。

考虑到大部分云资源的使用率都较低,部分研 究通过应用休眠机制以实现降低能源消耗的目的。 文献[5]为提高云系统节能水平,同时满足云用户 响应性能,引入休眠延迟机制与休眠唤醒阈值,提出 了一种云资源调度策略。建立带有休假延迟与启动 过程的同步 N-策略多重休假排队模型,评估了所提 策略的系统性能。文献[6]以实现绿色云计算为目 的,将基于唤醒阈值的休眠模式引入到集群虚拟机 中,提出了一种云任务分配策略。建立异步 N-策略 多重休假排队模型,给出了任务平均延迟和系统节 能率的表达式。在以上的休眠机制中,处于休眠状 态的虚拟机完全不工作,虽然有效降低了云系统能 源消耗,但是云用户的服务质量得不到保证。

本文兼顾云用户的服务质量与云系统的能源消 耗水平,在云端引入周期性同步半休眠模式,提出一 种新型的云计算任务卸载策略。考虑不同负载情况 下的云端虚拟机服务速率,构建连续时间工作休假 的排队模型,研究云用户的随机行为。结合系统中 任务请求的数量、本地处理器状态以及云端虚拟机

① 国家自然科学基金(61872311, 61973261)资助项目。

② 女,1996 年生,硕士生;研究方向:计算机网络性能分析;E-mail: ysusj750626@163.com。

③ 通信作者, E-mail: xyyj@ysu.edu.cn。 (收稿日期:2020-11-24)

状态,建立二维连续时间的 Markov 链,进行云用户 平均响应时间与云系统能源消耗的理论分析。考虑 不同休眠参数和不同系统负载,基于 Matlab 进行系 统实验,给出系统性能指标的统计结果。构造系统 成本函数,利用 Logistic 映射混沌方法初始化种群位 置,改进海鸥优化算法。以最小化系统成本函数为 目标,优化任务请求本地分配概率,为云计算任务卸 载策略的实施提供理论依据。

1 云计算任务卸载策略及系统模型

1.1 云计算任务卸载策略

针对移动设备存储空间不足且处理能力有限等 问题,将部分任务请求迁移到云端^[7]。云服务通过 分布式软件依托于云计算环境。在云数据服务中心 上通常部署一个或多个物理机,而一个物理机上又 可以部署多个虚拟机。当一个物理机上没有需要被 执行的任务时,如果其上部署的虚拟机保持空闲状 态,会因此产生大量的能耗,让空闲的虚拟机进入休 眠状态是降低云计算能耗的有效方法之一。然而, 休眠模式也可能会降低响应性能。因此在云端引入 半休眠模式,提出一种新的云计算任务卸载策略。

产生于移动设备的任务请求首先汇聚于接入点, 然后在本地负载均衡器的调度下就近分配到本地处 理器接受服务,或卸载至云端接受服务。卸载到云 端的任务请求由云端负载均衡器分配给某一个物理 机。云计算任务卸载策略的系统架构如图1所示。



图1 云计算任务卸载策略系统架构

当一个物理机上的任务请求全部完成时,若缓 存区中也无等待的任务请求,则该物理机上部署的 全部虚拟机同时由常规状态进入半休眠状态。在常 规状态下的虚拟机有高速工作状态和高速空闲状态 两种状态。在半休眠状态下,虚拟机不是完全停止 工作,而是以低速率运转,此时包含有低速工作状态 和低速空闲状态。在所提出的云计算任务卸载策略 中,虚拟机共有4种状态:高速工作状态、高速空闲 状态、低速工作状态和低速空闲状态。虚拟机的状 态转换过程如图2所示。



(1)高速工作状态。物理机及物理机上部署的 全部虚拟机正常工作,此状态下物理机中至少有一 个虚拟机正在服务任务请求。高速工作状态下到达 缓存区的任务请求依次在虚拟机接受服务。当服务 完成最后一个任务请求时,物理机启动休眠定时器, 开始一段半休眠期,物理机上的所有虚拟机同时进 入低速空闲状态。若处于高速工作状态的虚拟机服 务完成当前的任务请求,未分配到新的任务请求,且 同一物理机内有其他虚拟机还处于高速工作状态, 则该虚拟机转换到高速空闲状态。

(2)高速空闲状态。处于高速空闲状态的虚拟 机若分配到新的任务请求,则立即转换到高速工作 状态提供服务。若物理机服务完成最后一个任务请 求,则启动休眠定时器,其上的所有虚拟机同时进入 低速空闲状态。

(3)低速空闲状态。处于低速空闲状态的虚拟 机可以随时为新到达的任务请求提供服务,不受半 休眠定时器的控制。如果在休眠定时器所规定的范 围内没有任务请求到达,物理机将重新启动休眠定 时器,开始一次新的半休眠期,其上的虚拟机均保持 在低速空闲状态。当一个虚拟机处于低速空闲状态 时,若休眠定时器到期,且同一物理机内有其他虚拟 机未服务完成任务请求,则物理机结束半休眠,该虚 拟机转换到高速空闲状态;若分配到新的任务请求, 该虚拟机立即由低速空闲状态转换到低速工作状 态。

(4)低速工作状态。当一个虚拟机以较低的速 率完成当前的任务请求时,如果缓存区没有其他等 待的任务请求,该虚拟机将由低速工作状态转换到 低速空闲状态;否则,该虚拟机保持在低速工作状 态,以较低的服务速率为缓存区中等待的任务请求 提供服务。当休眠定时器到期时,如果有其他虚拟 机还在服务当前的任务请求,该物理机上的全部虚 拟机同时进入高速工作状态^[8]。

1.2 系统模型

将移动设备的 CPU 视为本地服务台,建立单一 服务台连续工作排队。云端每个虚拟机视为远程服 务台,建立多服务台同步多重工作休假排队。结合 单一服务台连续工作排队与同步多重工作休假排 队,建立系统模型。

假设移动设备任务请求的生成间隔是参数为 $\lambda(0 < \lambda < + \infty)$ 的指数分布。经过本地负载均衡器调度后以概率p分配到本地,以概率 $\bar{p} = (1 - p)$ 卸载至云端。

显然,本地服务台任务请求的到达服从参数 为 $\lambda_0 = p\lambda(0 < \lambda_0 < \lambda)$ 的 Poisson 过程。假设本地 服务台对一个任务请求的服务时间服从参数为 $\mu_0(0 < \mu_0 < + \infty)$ 的指数分布。假设本地只有一 个服务台,等待缓存区容量无穷大,假设任务请求的 到达过程与服务过程是彼此独立的。因此,本地服 务台执行任务请求的过程可看做 M/M/1 排队。

在云端负载均衡器的调度下,分配到云端的任 务请求以 $q_i(0 < i \leq n)(n$ 为云端部署的物理机个 数)的概率分配到第i个物理机接受服务。对于第i个物理机,任务请求的到达服从参数为 $\lambda_i = \bar{p}\lambda q_i$ 的 Poisson 过程。假设第i个物理机中部署了 c_i 个相互 独立并行工作的虚拟机,当任务请求到达一个物理 机时,若有高速空闲状态的虚拟机,则立即接受服务,否则排队等待。一旦物理机服务完成最后一个 任务请求, c_i 个虚拟机同时开始随机长度 V_i 的工作 休假,其中 V_i 服从参数为 $\theta_i(0 < \theta_i < + \infty)$ 的指数 分布。在非工作休假期内,任务请求的服务时间服 从参数为 $\mu_{bi}(0 < \mu_{bi} < + \infty)$ 的指数分布,工作休假 期服务时间服从参数为 $\mu_{ii}(0 < \mu_{ii} < \mu_{bi})$ 的指数分 布。假设所有物理机具有无限缓存,假设任务请求 的到达间隔、服务时间和工作休假期的长度相互独 立。因此,第*i*个物理机执行任务请求的过程可看 做同步多重工作休假 M/M/ c_i 排队^[9-10]。

令随机变量 $L_i(t) = l(l = 0, 1, ...)$ 表示时刻 t第 i 个物理机中任务请求的数量,称为系统水平。 令随机变量 $J_i(t) = j(j = 0, 1)$ 表示时刻 t 物理机 所处的状态: j = 0 表示物理机处于工作休假状态, j= 1 表示物理机处于正规忙状态。 $\{(L_i(t), J_i(t)), t \ge 0\}$ 构成一个二维连续时间 Markov 链, 其状态空间 Ω_i 为

 $\boldsymbol{\Omega}_{i} = \{(0,0)\} \cup \{(l,j): l \ge 1, j = 0,1\}$

令 $\pi_i(l, j)$ 表示稳态下系统水平为 l 且物理机 状态为 j 的概率分布。 $\pi_i(l, j)$ 定义为

 $\begin{aligned} \pi_i(l,j) &= \lim_{i \to \infty} P\{L_i(t) = l, J_i(t) = j\}, (l,j) \in \mathbf{\Omega}_i \\ &\Leftrightarrow \pi_i(0) = \pi_i(0,0) \text{ 表示为稳态下系统水平} \end{aligned}$

为0的概率向量,令 $\pi_i(l) = (\pi_i(l, 0), \pi_i(l, 1))$ 表示为稳态下系统水平为l的概率向量。二维连续 时间 Markov 链 { $(L_i(t), J_i(t)), t \ge 0$ } 的稳态概 率分布 Π_i 表示为

 $\boldsymbol{\Pi}_i = (\boldsymbol{\pi}_i(0), \boldsymbol{\pi}_i(1), \cdots)$

2 系统模型的稳态分析

2.1 稳态条件

根据云计算任务卸载策略建立的系统模型由 1个M/M/1排队和n(n为云端部署的物理机个数) 个同步多重工作休假 M/M/ c_i 排队构成,其中 c_i 为 云端部署的第 $i(0 < i \le n)$ 个物理机上虚拟机的个 数。为确保整个系统模型达到稳态,需要系统模型 中的每一个排队均达到稳态。

本地处理器所对应的 M/M/1 排队的稳态条件 — 1271 — 为

$$\rho_0 = \frac{\lambda_0}{\mu_0} < 1 \tag{1}$$

其中, ρ_0 表示为本地处理器的流通强度。

云端服务器上部署的第 $i(0 < i \leq n)$ 个物理机 所对应的同步多重工作休假 M/M/ c_i 排队的稳态条 件为

$$\rho_i = \frac{\lambda_i}{c_i \mu_{bi}} < 1 \tag{2}$$

其中, ρ_i 表示云端服务器上部署的第i个物理机的流通强度。

结合式(1)与式(2),给出系统模型处于稳定状态的充分必要条件为

$$\operatorname{Max}\{\rho_0, \rho_1, \cdots, \rho_n\} < 1 \tag{3}$$

2.2 状态转移

根据云端中每个物理机上任务请求的数量 l 与 物理机所处状态 j 之间的关系,给出 Markov 链 $\{(L_i(t), J_i(t)), t \ge 0\}$ 的状态转移过程,如图 3 所示。



由图 3 可知,状态转移只发生在相邻的系统水 平之间,表明二维 Markov 链 { $(L_i(t), J_i(t)), t \ge 0$ } 是拟生灭过程。

2.3 稳态分布

令 Q_i 表示 Markov 链 { $(L_i(t), J_i(t)), t \ge 0$ } 的一步转移率矩阵。根据系统水平划分 Q_i 为若干 个子阵。令 $Q_i(x, y)$ 表示系统水平由 $x(x = 0, 1, \dots)$ …) 到 $y(y = 0, 1, \dots)$ 的转移率子阵。为表述方 便,将 $Q_i(x, x - 1) Q_i(x, x)$ 和 $Q_i(x, x + 1)$ 分别 记为 $B_i(x) A_i(x)$ 和 $C_i(x)$ 。

(1) 当第*i*个物理机中任务请求的数量由 *x*变为*x*-1时,转移率子阵为*B_i*(*x*)。

当 x = 1 时,若以低速率服务完成任务请求,虚 拟机将由低速工作状态变为低速空闲状态,即由状 态(1,0)转换到状态(0,0),转移率为 μ_{ii} ;若以高 速率服务完成任务请求,虚拟机将由高速工作状态 变为低速空闲状态,即由状态(1,1)转换到状态 (0,0),转移率为 μ_{bi} 。 $B_i(1)$ 为2×1维矩阵,可表 示为

$$\boldsymbol{B}_{i}(1) = \begin{pmatrix} \boldsymbol{\mu}_{vi} \\ \boldsymbol{\mu}_{bi} \end{pmatrix}$$

当2 ≤ $x \le c_i$ 时,若以低速率服务完成任务请 求,虚拟机将由低速工作状态变为低速空闲状态,即 由状态 (x, 0)转换到状态 (x – 1, 0),转移率为 $x\mu_{xi}$;若以高速率服务完成任务请求,虚拟机将由高 速工作状态变为高速空闲状态,即由状态 (x, 1)转 换到状态 (x – 1, 1),转移率为 $x\mu_{bi}$ 。 $B_i(x)$ 为2×2 维矩阵,可表示为

$$\boldsymbol{B}_{i}(x) = \begin{pmatrix} x\mu_{vi} & 0\\ 0 & x\mu_{bi} \end{pmatrix}$$

当 $x > c_i$ 时,若以低速率服务完成任务请求,虚 拟机将继续保持低速工作状态,即由状态(x, 0)转 换到状态(x - 1, 0),转移率为 $c_i\mu_{ii}$;若以高速率服 务完成任务请求,虚拟机将继续保持高速工作状态, 即由状态(x, 1)转换到状态(x - 1, 1),转移率为 $c_i\mu_{bi}$ 。 $B_i(x)$ 为2×2维矩阵,可表示为

$$\boldsymbol{B}_{i}(x) = \begin{pmatrix} c_{i}\boldsymbol{\mu}_{vi} & 0\\ 0 & c_{j}\boldsymbol{\mu}_{bi} \end{pmatrix}$$

(2) 当第*i*个物理机中任务请求的数量 *x*保持不变时,转移率子阵为*A_i(x)*。

当 x = 0 时,第 i 个物理机处于半休眠状态,其 上所有虚拟机低速空转。若无任务请求到达,状态

— 1272 —

(0, 0)保持不变,转移率为 – $\lambda_i \circ A_i(0)$ 退化为一 个数,可表示为

 $A_i(0) = -\lambda_i$

当 $1 \le x \le c_i$ 时,若j = 0,第i个物理机处于半 休眠状态,其上部署的虚拟机部分低速运转部分低 速空转。若无任务请求到达也无任务请求离开且休 眠定时器未到期,状态 (x, 0)保持不变,转移率为 – ($\lambda_i + x\mu_{xi} + \theta_i$);若休眠定时器到期,物理机由状 态 (x, 0)转换到状态 (x, 1),转移率为 θ_i 。若j =1,第i个物理机处于常规状态,其上部署的虚拟机 部分高速运转部分高速空转。若无任务请求到达也 无任务请求离开,状态 (x, 1)保持不变,转移率为 – ($\lambda_i + x\mu_{ii}$)。 $A_i(x)$ 为 2×2 维矩阵,可表示为

$$\boldsymbol{A}_{i}(\boldsymbol{x}) = \begin{pmatrix} -(\boldsymbol{\lambda}_{i} + \boldsymbol{x}\boldsymbol{\mu}_{vi} + \boldsymbol{\theta}_{i}) & \boldsymbol{\theta}_{i} \\ 0 & -(\boldsymbol{\lambda}_{i} + \boldsymbol{x}\boldsymbol{\mu}_{bi}) \end{pmatrix}$$

当 $x > c_i$ 时,若j = 0,第i个物理机处于半休眠 状态,其上部署的虚拟机部分低速运转部分低速 空转。若无任务请求到达也无任务请求离开且休 眠定时器未到期,状态(x,0)保持不变,转移率为 -($\lambda_i + c_i\mu_{ii} + \theta_i$);若休眠定时器到期,物理机由状 态(x, 0)转换到状态(x, 1),转移率为 θ_i 。若j =1,第i个物理机处于常规状态,其上部署的虚拟机 部分高速运转部分高速空转。若无任务请求到达也 无任务请求离开,状态(x,1)保持不变,转移率为 -($\lambda_i + c_i\mu_{ii}$)。 $A_i(x)$ 为2×2维矩阵,可表示为

$$\boldsymbol{A}_{i}(\boldsymbol{x}) = \begin{pmatrix} -(\boldsymbol{\lambda}_{i} + c_{i}\boldsymbol{\mu}_{vi} + \boldsymbol{\theta}_{i}) & \boldsymbol{\theta}_{i} \\ 0 & -(\boldsymbol{\lambda}_{i} + c_{i}\boldsymbol{\mu}_{bi}) \end{pmatrix}$$

(3) 当第*i*个物理机中任务请求的数量由 *x*变为*x*+1时,转移率子阵为*C_i(x)*。

当x = 0时,第i个物理机处于半休眠状态,其 上所有虚拟机低速空转。若有任务请求达到,其中 一个虚拟机将变为低速工作状态,即由状态(0,0) 转换到状态(1,0),转移率为 λ_i 。 $C_i(0)$ 为1×2 维 矩阵,可表示为

 $\boldsymbol{C}_i(0) = (\boldsymbol{\lambda}_i \quad 0)$

当 $x \ge 1$ 时,若j = 0,第i个物理机处于半休眠 状态,其上部署的虚拟机既有低速空转又有低速运 转。若有任务请求到达,物理机由状态(x, 0)转换 到状态(x+1, 0),转移率为 λ_i 。若j = 1,第i个物 理机处于常规状态,其上部署的虚拟机部分高速运转部分高速空转。若有任务请求到达,物理机由状态 (x, 1)转换到状态 (x + 1, 1),转移率为 λ_i 。 $C_i(x)$ 为2×2 维矩阵,可表示为

$$\boldsymbol{C}_{i}(\boldsymbol{x}) = \begin{pmatrix} \boldsymbol{\lambda}_{i} \\ & \boldsymbol{\lambda}_{i} \end{pmatrix}$$

结合以上给出的一步转移率子阵,发现 $B_i(x)(0 < i \le n)$ 和 $A_i(x)(0 < i \le n)$ 从系统水平 c_i 起开始重复,转移率子阵 $C_i(x)(0 < i \le n)$ 从系 统水平1起开始重复。将重复的 $B_i(x) \land A_i(x)$ 和 $C_i(x)$ 分别用 $B_i \land A_i$ 和 C_i 表示,则 Markov 链 $\{(L_i(t), J_i(t)), t \ge 0\}$ 的一步转移率矩阵 Q_i 可表 示为分块三对角形式

$$\boldsymbol{Q}_{i} = \begin{pmatrix} \boldsymbol{A}_{i}(0) & \boldsymbol{C}_{i}(0) \\ \boldsymbol{B}_{i}(1) & \boldsymbol{A}_{i}(1) & \boldsymbol{C}_{i} \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{B}_{i}(c_{i}-1) & \boldsymbol{A}_{i}(c_{i}-1) & \boldsymbol{C}_{i} \\ & \boldsymbol{B}_{i} & \boldsymbol{A}_{i} & \boldsymbol{C}_{i} \\ & \vdots & \ddots & \ddots \end{pmatrix}$$

拟生灭过程 { $(L_i(t), J_i(t)), t \ge 0$ } 正常返的 充分必要条件是矩阵二次方程式

$$\boldsymbol{R}_i^2 \boldsymbol{B}_i + \boldsymbol{R}_i \boldsymbol{A}_i + \boldsymbol{C}_i = \boldsymbol{0}$$
 (4)

有最小非负解 R_i ,且谱半径 Sp $(R_i) < 1$ 。最小非负 解 R_i 被称为率阵,可以推出其精确解。

由于转移率子阵 B_i 、 A_i 和 C_i 均为上三角矩阵, 因此,率阵 R_i 也为上三角矩阵,表示为

$$\boldsymbol{R}_i = \begin{pmatrix} r_i^{11} & r_i^{12} \\ 0 & r_i^{22} \end{pmatrix}$$

将 R_i 、 B_i 、 A_i 和 C_i 代入矩阵二次方程式(4)中, 可得方程组

$$\begin{cases} c_{i}\mu_{vi}(r_{i}^{11})^{2} - (\lambda_{i} + c_{i}\mu_{vi} + \theta_{i})r_{i}^{11} + \lambda_{i} = 0\\ c_{i}\mu_{bi}(r_{i}^{11}r_{i}^{12} + r_{i}^{12}r_{i}^{22}) + \theta_{i}r_{i}^{11} - (\lambda_{i} + c_{i}\mu_{bi})r_{i}^{12} = 0\\ c_{i}\mu_{bi}(r_{i}^{22})^{2} - (\lambda_{i} + c_{i}\mu_{bi})r_{i}^{22} + \lambda_{i} = 0 \end{cases}$$
(5)

— 1273 —

$$\begin{array}{cccc} \ddots & \ddots & \ddots \\ \boldsymbol{B}_i(c_i-1) & \boldsymbol{A}_i(c_i-1) & \boldsymbol{C}_i \\ \boldsymbol{B}_i & \boldsymbol{R}_i \boldsymbol{B}_i + \boldsymbol{A}_i \end{array}$$

由矩阵几何解中的平衡方程和归一化条件可 知,拟生灭过程 { $(L_i(t), J_i(t)), t \ge 0$ } 的平稳分 布满足方程组

$$\begin{cases} (\boldsymbol{\pi}_{i}(0), \boldsymbol{\pi}_{i}(1), \cdots, \boldsymbol{\pi}_{i}(c_{i})) B[\boldsymbol{R}_{i}] = \boldsymbol{0} \\ (\boldsymbol{\pi}_{i}(0), \boldsymbol{\pi}_{i}(1), \cdots, \boldsymbol{\pi}_{i}(c_{i}-1)) \boldsymbol{e}_{1} \\ + \boldsymbol{\pi}_{i}(c_{i}) (\boldsymbol{I} - \boldsymbol{R}_{i})^{-1} \boldsymbol{e}_{2} = 1 \end{cases}$$

$$(6)$$

其中, e_1 是 ($2c_i - 1$) × 1 的全 1 列向量, e_2 是 2 × 1 的全 1 列向量, I 是 2 × 2 的单位阵。

引入增广矩阵
$$P_i = \left(B[R_i] \middle| \begin{array}{c} e_1 \\ (I - R_i)^{-1} e_2 \end{array} \right),$$
方

程组式(6)可等价表示为

$$(\boldsymbol{\pi}_{i}(0),\boldsymbol{\pi}_{i}(1),\cdots,\boldsymbol{\pi}_{i}(c_{i}))\boldsymbol{P}_{i} = (0,0,\cdots,0,1)$$
(7)

利用高斯-塞德尔法求解方程式(7),可以给出 $\pi_i(k)(0 \le k \le c_i)$ 的数值解。

基于矩阵几何解方法,由 $\pi_i(c_i)$ 进一步给出 $\pi_i(k)$ ($k \ge c_i + 1$)如下:

 $\boldsymbol{\pi}_i(k) = \boldsymbol{\pi}_i(c_i) \boldsymbol{R}_i^{k-c_i}$

3 系统性能指标

关注云用户响应时间和云系统能耗水平,评估 基于半休眠模式的云计算任务卸载策略的性能指 标。 任务请求在本地处理器接受服务的平均响应时间 T_{MD}包括任务请求在本地缓存中的等待时间与在本地处理器上执行的时间。根据 M/M/1 排队模型,可以得到 T_{MD}的表达式为

$$T_{\rm MD} = \frac{1}{\mu_0 - \lambda_0} \tag{8}$$

任务请求卸载到云端服务器接受服务的平均响 应时间 *T*_c包括任务请求在云端缓存中的等待时间 与在云端服务器上执行的时间。卸载到云端后第 *i* 个物理机上任务请求的平均响应时间 *T_i*的表达式 为

$$T_i = \frac{1}{\lambda_i} \sum_{l=c_i+1}^{\infty} (l - c_i) \boldsymbol{\pi}_i(l) \boldsymbol{e} + T_{si}$$

其中, *T_{si}* 表示任务请求在第*i* 个物理机上的平均服务时间, *e* 为2×1的全1列向量。

根据卸载到第 i 个物理机上的任务请求到达云 端时的状态,分以下 4 种情况讨论 $T_{sio}(1)$ 若第 i 个物理机处于常规状态,当前任务请求接受高速率 服务,平均服务时间为 $\frac{1}{\mu_{i}}$; (2) 若第*i*个物理机处 于半休眠状态,休眠定时器到期之前缓存中排队等 待的所有任务请求及当前任务请求均完成服务,当 前任务请求接受低速率服务,平均服务时间为 1;; (3) 若第*i*个物理机处于半休眠状态,休眠定时器 到期之前缓存中排队等待的所有任务请求均完成服 务,但当前任务请求未服务完成,当前任务请求经历 一段低速率服务和一段高速率服务,平均服务时间 为 $\frac{1}{\mu_i}$ + $\frac{1}{\theta_i}$; (4) 若第 *i* 个物理机处于半休眠状态, 休眠定时器到期时,缓存中有排队等待的任务请求 未服务完成,当前任务请求接受高速率服务,平均服 务时间为 $\frac{1}{\mu_{i}}$ 。因此,任务请求在第*i*个物理机上的 平均服务时间 T_{st} 表示为

$$T_{si} = \frac{d_i^1}{\mu_{bi}} + \frac{d_i^0 \mu_{vi}}{\mu_{vi} + \theta_i} \left(\frac{\mu_{vi}}{\mu_{vi} + \theta_i}\right)^a \frac{1}{\mu_{vi}} \\ + \frac{d_i^0 \theta_i}{\mu_{vi} + \theta_i} \left(\frac{\mu_{vi}}{\mu_{vi} + \theta_i}\right)^a \left(\frac{1}{\theta_i} + \frac{1}{\mu_{bi}}\right) \\ + d_i^0 \left(1 - \left(\frac{\mu_{vi}}{\mu_{vi} + \theta_i}\right)^a\right) \frac{1}{\mu_{bi}}$$

— 1274 —

其中, $d_i^1 = \sum_{l=1}^{\infty} \pi_i(l,1)$ 表示云端第 *i* 个物理机处于 常规状态的概率, $d_i^0 = \sum_{l=0}^{\infty} \pi_i(l,0)$ 表示云端第 *i* 个 物理机处于半休眠状态的概率, $a = \sum_{l=c_i+1}^{\infty} (l - c_i)\pi_i(l)e$ 表示第 *i* 个物理机上缓存中任务请求的 数量。

令任务请求分配到第i个物理机上的概率为 $q_i(0 < q_i < 1)$,则卸载到云端接受服务的任务请求 的平均响应时间 T_c 的表达式为

$$T_{\rm C} = \sum_{i=1}^{n} q_i T_i \tag{9}$$

云用户的平均响应时间 T 定义为任务请求从 进入系统开始直到离开系统所经历的时间长度的平 均值。结合式(8)和式(9),给出平均响应时间 T 的 表达式为

 $T = pT_{MD} + \bar{p}T_{C}$ (10) 其中, p(0 为任务请求在本地处理器接受 $服务的概率, <math>\bar{p} = 1 - p$ 为任务请求卸载到云端服务 器接受服务的概率。

本地处理器处于工作状态时的平均运行功率表示为 P_{busy},处于空闲状态时的平均运行功率表示为 P_{idle}。任务请求在本地处理器接受服务的平均功率 消耗 E_{MD} 为

$$E_{\rm MD} = \frac{\lambda_0}{\mu_0} P_{\rm busy} + \left(1 - \frac{\lambda_0}{\mu_0}\right) P_{\rm idle}$$
(11)

云端上部署的第 *i* 个物理机处于半休眠状态 时,该物理机维持其上一个虚拟机空转的平均运行 功率表示为 *Pⁱ*_{vf};维持其上一个虚拟机运转的平均 运行功率表示为 *Pⁱ*_{vb}。云端上部署的第 *i* 个物理机处 于常规状态时,该物理机维持其上一个虚拟机空转 的平均运行功率表示为 *Pⁱ*_{bf};维持其上一个虚拟机 运转的平均运行功率表示为 *Pⁱ*_{bb}^[11-12]。任务请求在 云端服务器接受服务的平均功率消耗 *Eⁱ*_c为

$$E_{\rm C}^{i} = \sum_{k=0}^{c_{i}-1} (c_{i} - k)\pi_{i}(k,0)P_{\rm vf}^{i} + \sum_{k=1}^{c_{i}-1} (c_{i} - k)\pi_{i}(k,1)P_{\rm bf}^{i} + \left(\sum_{k=0}^{c_{i}-1} k\pi_{i}(k,0) + \sum_{k=c_{i}}^{\infty} c_{i}\pi_{i}(k,0)\right)P_{\rm vb}^{i} + \left(\sum_{k=1}^{c_{i}-1} k\pi_{i}(k,1) + \sum_{k=c_{i}}^{\infty} c_{i}\pi_{i}(k,1)\right)P_{\rm bb}^{i}$$
(12)

结合式(11)和式(12),给出系统平均功率 E 的 表达式为

$$E = pE_{\rm MD} + \bar{p}\sum_{i=1}^{n} q_i E_{\rm C}^{i}$$
(13)

4 系统实验

为了进一步进行系统实验研究任务请求初始分 配概率与休眠参数对云计算任务卸载策略的影响, 进行性能指标变化趋势的系统实验。系统实验的计 算机硬件环境为 Intel(R) Core(TM),i7-4790 CPU@ 3.60 GHz,8.00 GB RAM。软件环境为 Matlab R2016a。

在保证系统模型稳定的条件下的参数设置如 表1所示。

表1 实验参数设置

参数	值		
云端物理机数量 n	3个		
不同物理机上部署的虚拟机 个数 c_1, c_2, c_3	(2,2,2)个		
任务请求分配到不同物理机上 的概率为 q1, q2, q3	(0.49,0.22,0.29)		
本地处理器的服务率 μ_0	16 个/s		
不同物理机的高速服务率 μ_{b1} , μ_{b2} , μ_{b3}	(20,25,18)个/s		
不同物理机的低速服务率 $\mu_{v1}, \mu_{v2}, \mu_{v3}$	(7,8,5)个/s		
本地处理器处于工作状态时 的运行功率 P _{busy}	70 mW		
本地处理器处于空闲状态时 的运行功率 P _{idle}	20 mW		
物理机维持一个虚拟机在半休眠 状态下空转的运行功率 $P_{vf}^1, P_{vf}^2, P_{vf}^3$	(15,17,19) mW		
物理机维持一个虚拟机在半休眠 状态下低速工作的运行功率 P ¹ _{vb} , P ² _{vb} , P ³ _{vb}	(25,27,29) mW		
物理机维持一个虚拟机在常规状态 下空转的运行功率 P ¹ _{bf} , P ² _{bf} , P ³ _{bf}	(40,42,44) mW		
物理机维持一个虚拟机在常规状态 下高速工作的运行功率 $P^1_{\rm bb}, P^2_{\rm bb}, P^3_{\rm bb}$	(60,62,64) mW		

根据表1设定的实验参数,针对不同的休眠参数 θ和不同的任务请求到达率λ,给出云用户平均 响应时间 T 随任务请求本地分配概率 p 的变化趋势 如图4 所示。

— 1275 —



图 4 云用户平均响应时间的变化趋势

对比图4 (a)和4 (b)可知,当任务请求本地分 配概率 *p*和任务请求到达率 λ 一定时,随着休眠参 数 θ 的增大,云用户平均响应时间 *T* 减小。休眠参 数 θ 较大时,云端物理机处于半休眠状态的时间长 度缩短。这种情况下,云端服务器可以更及时地返 回常规状态以高速服务率为任务请求提供服务。因 此,云用户平均响应时间 *T* 减小。

固定任务请求到达率 λ 和休眠参数 θ, 分别观 察图 4(a) 和图 4(b) 可知,随着任务请求本地分配 概率 p 的增加,云用户平均响应时间 T 呈现出先下 降后上升的趋势。在曲线最低点左侧,任务请求本 地分配概率 p 较小,任务请求到达系统后大概率地 卸载到云端服务器接受服务,任务请求在云端的响 应时间成为云用户平均响应时间 T 的主导因素。随 着任务请求本地分配概率 p 的增加,分配到本地的 — 1276 — 任务请求逐渐增多,堆积在云端的任务请求减少,云 用户平均响应时间 T 随之减小。在曲线最低点右 侧,任务请求本地分配概率p 较大,任务请求到达系 统后大概率地选择在本地处理器接受服务,任务请 求在本地处理器接受服务的响应时间成为云用户平 均响应时间 T 的主导因素。随着任务请求本地分配 概率p 的增加,分配到本地的任务请求越来越多,大 量的任务请求将会堆积在本地处理器的缓存区,云 用户平均响应时间 T 随之增大。

固定休眠参数 θ 和任务请求本地分配概率 p, 分别观察图 4(a) 和图 4(b) 可知,随着任务请求到 达率 λ 的增加,云用户平均响应时间 T 增加。无论 是在本地还是云端,随着任务请求到达率 λ 的值增 大,系统内到达的任务请求数量也就越多,显然云用 户平均响应时间 T 将会增大。

针对不同的休眠参数 θ 和不同的任务请求到达率 λ ,给出云系统平均功率 E 随任务请求本地分配 概率 p 的变化趋势如图 5 所示。

对比图 5(a)和 5(b)可知,当任务请求到达率 λ 一定而任务请求本地分配概率 p较小时,随着休 眠参数 θ 的增大,云系统平均功率 E略有增大。休 眠参数 θ 较大时,云端物理机处于半休眠状态的时 间长度缩短,物理机可以更及时地转换到常规状态。 处于常规状态的物理机所消耗的功率大于处于半休 眠状态的物理机所消耗的功率。因此,云系统平均 功率 E增加。当任务请求到达率 λ 一定而任务请 求本地分配概率 p较大时,卸载到云端的任务请求 数量较少,云端物理机有更多的机会持续处于半休 眠状态。云系统平均功率 E与休眠参数 θ 的关系减 弱。

固定休眠参数 θ 和任务请求到达率 λ, 分别观 察图 5(a) 和图 5(b) 可知,随着任务请求本地分配 概率 p 的增加,云系统平均功率 E 呈现出先下降后 上升的趋势。在曲线最低点左侧,任务请求本地分 配概率 p 较小,任务请求到达系统后大概率地卸载 到云端服务器接受服务。随着任务请求本地分配概 率 p 的增加,卸载到云端的任务请求减少,云端物理 机有更多的机会进入半休眠状态,云系统平均功率 E 随之减小。在曲线最低点右侧,任务请求本地分 配概率 p 较大,任务请求到达系统后大概率地选择 在本地处理器接受服务。随着任务请求本地分配概 率 p 的增加,本地处理器提供服务需要的功耗增加, 云系统平均功率 E 随之变大。



固定休眠参数 θ 和任务请求本地分配概率 p, 分别观察图 5(a) 和图 5(b)还可知,随着任务请求 到达率 λ 的增加,云系统平均功率 E 增加。无论是 在本地还是云端,随着任务请求到达率 λ 的增大, 系统内到达的任务请求数量增多,执行任务请求需 要的功耗更多,因此云系统平均功率 E 将会增大。

综合图 4 和图 5 的实验结果可以发现,任务请 求本地分配概率是影响系统性能的重要因素。兼顾 云用户和云系统的不同要求,在不同的任务请求到 达率和相同的休眠参数情况下,优化任务请求本地 分配概率。

5 系统优化

为了权衡云用户平均响应时间与云系统平均功 率之间的折衷关系,构造系统成本函数 F 为

 $F = \beta_1 T + \beta_2 E$ (14) 其中, β_1 表示云用户平均响应时间 T 对系统成本的 影响因子, β_2 表示云系统平均功率 E 对系统成本的 影响因子。在实际应用中, 对于高响应需求的云系 统, 可将参数 β_1 的值设置得大些; 而对于高节能需 求的云系统,则可将参数 β_2 的值设置得大些。

与其他智能优化算法相比,海鸥智能优化算法 容易实现,且不需要调整过多的参数^[13]。为了加快 算法搜索速度,利用 Logistic 映射混沌方法进行海鸥 位置初始化,改进海鸥优化算法,优化任务请求本地 分配概率。算法中海鸥的位置表示任务请求本地分 配概率 *p*,海鸥的适应值表示系统成本函数 *F* 的值。 算法的主要步骤如下。

步骤1 初始化海鸥数量 N,最大迭代次数 $Max_{iterations}$,控制海鸥在给定搜索空间中运动行为 A 的频率 f_c ,螺旋攻击形状的相关常数 u、v,海鸥飞行 位置的上边界 Max_p 与下边界 Min_p 。初始化当前迭 代次数 t = 0。

步骤2 利用 Logistic 映射方法产生混沌变量 初始化海鸥的位置。

 $p_{1} = rand(Max_{p} - Min_{p}) + Min_{p}$ for m = 2:N $p_{m} = \varepsilon \times p_{m-1} \times (1 - p_{m-1})$ % $\varepsilon = 4$ 时完全处于混沌状态

endfor

步骤3 初始化当前最佳海鸥位置 *p** 和适应 值 *F**。

$$p^{*} = p_{1}$$

$$F^{*} = p_{1}(\beta_{1}T_{MD}(p_{1}) + \beta_{2}E_{MD}(p_{1}))$$

$$+ (1 - p_{1})(\beta_{1}T_{C}(p_{1}) + \beta_{2}\sum_{i=1}^{n}q_{i}E_{C}^{i}(p_{1}))$$

步骤4 根据海鸥迁徙行为和攻击行为更新海 鸥位置。

for
$$m = 1:N$$

 $A = f_c - \left(\frac{tf_c}{Max_{\text{iterations}}}\right)$

— 1277 —

$$D_{s}(m) = \begin{vmatrix} \left(f_{c} - \left(\frac{tf_{c}}{Max_{\text{iterations}}}\right)\right)p_{m} \\ + (2\operatorname{rand} f_{c} - f_{c})(p_{1} - p_{m}) \end{vmatrix}$$
$$p_{m} = D_{s}(m)(ue^{\omega v})^{3}\omega \cos\omega \sin\omega + p_{1}$$
% ω $\mathfrak{W}(0, 2\pi)$ $\mathfrak{M} \mathfrak{M} \mathfrak{M}$

endfor

步骤5 使用系统成本函数 F 计算海鸥的适应 值 $F(p_m), m = 1, 2, \dots, N_{\circ}$

for
$$m = 1:N$$

 $F(p_m) = p_m (\beta_1 T_{MD} + \beta_2 E_{MD})$
 $+ (1 - p_m) (\beta_1 T_C + \beta_2 \sum_{i=1}^n q_i E_C^i)$

endfor

for

步骤6 找出当前最佳海鸥位置 p^* 和适应值 F^* 。

$$m = 1:N$$

if $F(p_m) < F^*$
 $p^* = p_m$
 $F^* = F(p_m)$

endif

```
endfor
```

步骤7 更新迭代次数。

t = *t* + 1 % 更新当前迭代次数

- if $t < Max_{\text{iterations}}$
 - 跳转到步骤4

endif

步骤8 输出最佳海鸥位置为最优任务请求本 地分配概率 *p**;输出最佳适应值为最小系统成本 *F**。

沿用第4节系统实验的参数,设定 $\beta_1 = 2, \beta_2 = 0.2, N = 30, Max_{iterations} = 1000, f_c = 2, u = 1, v = 1, Max_p = 0.99999, Min_p = 0.00001。利用改进的海鸥 优化算法,给出任务请求本地分配概率的优化结果 如表 2 所示。$

由表 2 的优化结果可知,对于相同的休眠参数 θ ,当任务请求到达率 λ 增大时,受本地处理器处理 能力的限制,会有更多的任务请求卸载到云端服务 器,因此,最优任务请求本地分配概率 p^* 变小。对 于相同的任务请求到达率 λ ,休眠参数 θ 越大,云端

表 2 任务请求本地分配概率的优化结果

任务请求	休眠参数	最优任务请求	最小系统
到达率 λ	heta	本地分配概率 p^*	成本 F^*
11.0	0.5	0.37364	7.8140
11.0	1.0	0.38067	7.8386
11.5	0.5	0.36295	7.8832
11.5	1.0	0.37031	7.9113
12.0	0.5	0.35294	7.9500
12.0	1.0	0.36060	7.9817
12.5	0.5	0.34354	8.0145
12.5	1.0	0.35148	8.0499
13.0	0.5	0.33467	8.0770
13.0	1.0	0.34290	8.1161

服务器处于半休眠状态的时间越短,云端服务器产 生的功率消耗越高。为了实现系统成本的最小化, 更多的任务请求将在本地处理器接受服务,因此,最 优任务请求本地分配概率 *p** 越大越好。

6 结论

综合考虑云用户低响应时间与云系统低功率消 耗的要求,根据系统负载调整云端服务器的工作速 率,基于半休眠模式,提出了一种云计算任务卸载策 略。通过建立本地单一服务台连续工作排队和云端 多服务台同步多重工作休假排队,给出了稳态下云 用户平均响应时间和云系统平均功率的表达式,并 通过实验揭示出不同性能指标之间的折衷关系。从 经济学角度出发,构建了系统成本函数,改进传统的 海鸥优化算法,给出了任务请求本地分配概率的优 化方案。

本文所引入的周期性同步半休眠模式具有良好 的节能效果。在未来的研究中,将考虑具有多级适 应性的半休眠模式,以进一步提高用户的响应性能。

参考文献

[1] Zhang Y Z, Dong X S, Zhao Y N. Decentralized computation offloading over wierless-powered mobile-edge computing networks [C] // Proceedings of IEEE International Conference on Artificial Intelligence and Information Systems, Fukuoka, Japan, 2020: 137-140

[2] Kuang Z K, Shi Y W, Guo S T, et al. Multi-user offload-

— 1278 —

ing game strategy in OFDMA mobile cloud computing system[J]. IEEE Transactions on Vehicular Technology, 2019, 68(12); 12190-12201

- [3] Malarvizhi N, Priyatharsini S G, Koteeswaran S. Cloud resource scheduling optimal hypervisor for dynamic cloud computing environment[J]. Wireless Personal Communications, 2020, 115(1): 27-42
- [4] Belgacem A, Beghdad-Bey K, Nacer H, et al. Efficient dynamic resource allocation method for cloud computing environment[J]. *Cluster Computing*, 2020, 23: 2871-2889
- [5] 王秀双,金顺福.基于新型休眠机制的云任务调度策略的研究[J].高技术通讯,2018,28(11-12):907-914
- [6] Jin S F, Qie X C, Zhao W J, et al. A clustered virtual machine allocation strategy based on a sleep-mode with wake-up threshold in a cloud environment [J]. Annals of Operations Research, 2020, 293(1): 193-212
- [7] Zhou W C, Fang W W, Li Y Y, et al. Markov approximation for task offloading and computation scaling in mobile edge computing [J]. Mobile Information Systems,

2019, 2019: 1-12

- [8]李吉良,秦兵,李文江,等.融合唤醒阈值与半休眠 模式的云虚拟机调度策略[J].燕山大学学报,2020, 44(4):370-378
- [9] 郜燕,刘文芬,张建辉.同步多重工作休假 M/M/c 排
 队系统的性能指标[J].工程数学学报,2012,29
 (2):179-191
- [10] 朱翼隽, 徐剑, 周宗好. 多重工作休假的 M/M/c 排队 系统[J]. 江苏大学学报, 2012, 33(3): 369-372
- [11] Schwartz C, Pries R, Tran-Gia P. A queuing analysis of an energy-saving mechanism in data centers [C] // Proceedings of IEEE International Conference on Information Networking, Bali, Indonesia, 2012: 70-75
- [12] 孙健,廖丹,李可,等.基于排队论的异构数据中心 性能及能源管理策略[J].电子科技大学学报,2018, 47(2):161-168
- [13] Dhiman G, Kumar V. Seagull optimization algorithm: theory and its applications for large-scale industrial engineering problems [J]. Knowledge-Based Systems, 2019, 165: 169-196

Research on semi-sleep mode based task offloading strategy and system performance in cloud computing

Song Jia, Yu Jing, Jin Shunfu

(School of Information Science and Engineering, Yanshan University, Qinhuangdao 066004)

Abstract

Task offloading in cloud computing is a mechanism for migrating intensive mobile applications to cloud environment. By considering the service quality of cloud users and the energy consumption level of cloud system, a periodic synchronous semi-sleep mode is introduced in the cloud and a task offloading strategy is proposed. According to the service rates of virtual machine under different loads, a multi-server queueing model with synchronous multiple working vacation is constructed. Using the quasi-birth-and-death process and matrix-geometric solution method, the average response time of cloud users and the energy consumption of cloud system are derived. System experiments are carried out to evaluate the system performance and verify the validity of the offloading strategy. With the weight factors of different performance measures, a system cost function is constructed. By improving a seagull optimization algorithm, the proposed task offloading strategy in cloud computing is optimized with the minimum system cost.

Key words: cloud computing, task offloading, semi-sleep, working vacation queue, intelligent optimization