DAG 分割模型下的云工作流调度策略

薛 凡

(黄淮学院 创新创业学院,河南 驻马店 463000)

摘 要:为了优化云工作流调度的经济代价和执行效率,提出一种基于有向无循环图(DAG)分割的工作流调度 算法 PBWS。以工作流调度效率与代价同步优化为目标,算法将调度求解过程划分为三个阶段进行:工作流 DAG 结构分割、分割结构调整及资源分配。工作流 DAG 结构分割阶段在确保任务间执行顺序依赖的同时求解 初始的任务分割图;分割结构调整阶段以降低执行跨度为目标,在不同分割间对任务进行重分配;资源分配阶段 旨在选择代价最高效的任务与资源映射关系,确保资源的总空闲时间最小。利用五种科学工作流 DAG 模型对 算法进行了仿真实验。结果表明。PBWS 算法仅以较小的执行跨度为开销,极大降低了工作流执行代价,实现 了调度效率与调度代价的同步优化,其综合性能是优于同类型算法的。

关键词:云计算;科学工作流;调度优化;DAG分割;执行跨度

中图分类号: TP391 文献标志码: A 文章编号: 1001-3695(2019)12-043-3725-04 doi;10.19734/j.issn.1001-3695.2018.04.0060

Cloud workflow scheduling strategy in DAG partition model

Xue Fan

(College of Innovation & Entrepreneurship, Huanghuai University, Zhumadian Henan 463000, China)

Abstract: For optimizing the economical cost and scheduling efficiency of cloud workflow scheduling, this paper proposed a workflow scheduling algorithm PBWS based on DAG (directed acyclic graph) partition. With the goal of optimizing synchronously the workflow scheduling efficiency and cost, this algorithm divided the scheduling solution into three stages: DAG partition of the workflow structure, partition structure adjustment and resource allocation. DAG partition of the workflow structure was to get the initial tasks partition graph when guaranteeing the execution order-dependency between tasks. The partition structure adjustment was to re-allocate tasks in different partitions with a goal of reducing execution makespan. The resource allocation was to determine the most cost-efficient matches between tasks and resources ensuring the minimization of the total idle time of resource. This paper constructed some simulation experiments for algorithms by the five types of scientific workflow DAG model. The experimental results show PBWS algorithm can greatly reduce the execution cost of workflow in terms of cost by a little of overhead on execution makespan and realize the synchronous optimization of the scheduling efficiency and the scheduling cost, whose overall performance performs better than the same type of algorithms.

Key words: cloud computing; scientific workflow; scheduling optimization; DAG partition; execution makespan

0 引言

云数据中心可以提供巨大的 IT 能力,服务于社交、邮件、 金融和工业等诸多领域。尤其在科学计算和工程控制应用中, 爆炸式、复杂大规模的数据增长对计算环境提出了更高的要 求。该类领域中的多数应用任务通常体现为工作流(work-flow)形式^[1],其结构可参见图 1 所示的五种常规科学领域的 工作流。工作流中,每个任务作为工作流整体的一部分,或者 作为初始数据的传送者,或者作为中间数据的中继者,进而形 成一种任务间的顺序依赖,协作完成一种应用任务。

工作流调度本质上是 NP 问题,多数工作通过设计单目标 函数进行优化。例如文献[2~4]设计调度算法以寻找满足预 先定义的工作流执行截止期限的最低价调度方案;文献[5~ 8]所设计的算法则是寻找满足预先定义的预算的最快调度方 案;相比较而言,文献[9~11]则是以降低执行时间为目标,而 未考虑费用约束问题。虽然以上方法引入了预算或截止期限 约束,但总体来说,算法优化均是单目标优化为限制的,未能很 好地实现符合云环境资源使用特征的工作流调度优化。

本文同步考虑执行跨度和执行代价的同步最小化,并设计

一种基于 DAG 分割的工作流调度算法。算法将工作流 DAG 结构划分为包括多个任务的多个群组(分割 partition),进而实现分割与云资源动态能力的优化匹配。算法通过三阶段的调度优化最终可以实现工作流调度效率与代价的均衡。



收稿日期: 2018-04-19; 修回日期: 2018-06-08

作者简介:薛凡(1978-),男,河南驻马店人,讲师,博士研究生,主要研究方向为智能计算、云计算(xuefan78@ sohu. com).

1 系统模型

工作流表示为一个有向无循环图 $G = \langle V, E \rangle$,其中 V 为任 务集,E 为边集。每条边连接两个任务,代表任务间的顺序约 束。给定资源集合 $R = R_1, R_2, \dots, R_m$,资源集以 Amazon EC2 的 实例类型进行配置,包括 m4. xlarge、c4. 2xlarge 和 r3. 4xlarge。 每个资源集 $R_i \in R$ 由同质资源构成,配置不同数量的处理器内 核、处理能力、内存和存储空间。基于资源能力对集合 R 进行 排序,若j > i,则属于集合 R_i 中的资源拥有低于集合 R_j 的资源 能力。仅当接收父任务产生的数据后,子任务 v_i 才可以开始 执行。在所有父任务中,最迟发送数据的任务称为最有影响父 任务 MIP。令 $c_{(i,j)}$ 为任务 $v_i = v_j$ 间的通信代价,工作流的完 成时间可定义为执行跨度 makespan。对于任务 $v_i \in V$ 的最早 开始时间 EST 和最早完成时间 EFT 定义为

$$EST(v_i) = \begin{cases} 0 & v_i \end{pmatrix} (1)$$
$$EFT(MIP(v_i)) + c_{MIP(v_i),i}] 否则$$

$$EFT(v_i) = EST(v_i) + w_i$$
(2)

其中: w_i 为任务 v_i 的执行时间。当调度至相同资源上的另一个任务的实际完成时间晚于 EST(v_i)时,任务的实际开始时间AST 和实际完成时间 AFT 不同于 EST 和 EFT。

工作流的调度跨度 makespan 主要受关键路径上任务的执 行时间的影响,该路径拥有最高的代价(通信和计算代价),并 结束于出口任务。任务 $v_i \in V$ 的计算代价(执行时间 w_i)取决 于调度任务的资源能力,资源能力越高,执行时间越短。令 w_i^i 为计算代价(时间),其中,*i*为任务 $v_i \in V$ 的索引,*j*为资源集 R_j $\in R$ 的索引。为了简化计算过程,为每个资源集合 $R_i \in R$ 引入 一个常量因子 $cv_i > 0$ 。给定任务 v_i ,每个资源集合常量 $cv_i(R_i \in R)$ 代表任务的近似相对执行时间。例如若 $cv_3 = 4$,则任务 的执行时间 $w_i^3 = (1/4)w_i^1$ 。令 t_i 为每小时利用资源集 $R_i \in R$ 中资源的代价,该值取决于资源能力,若 i < j,则 $t_i < t_i$ 。

工作流调度目标即为分配给定工作流的任务至 R 中资源 集 S,同时确保资源利用总代价和任务执行跨度最小化。

2 算法设计

2.1 算法概述

PBWS 算法的目标是分割工作流任务形成任务群组,进而 产生有效的资源与任务间的映射,分割目标是简化资源分配过 程。基于任务分割区域间的数据依赖性,决定分配至每个分割 的资源能力,并确保工作流执行跨度与代价最小化。PBWS 算 法由 DAG 分割、分割调整和资源分配三个步骤组成。DAG 分 割步骤的目标是划分给定的工作流 DAG 结构,使得到的 DAG 分割在确保任务间依赖的同时,简化资源分配过程,而在后续 步骤中,每个分割均视为单个节点对待。分割调整步骤中,为 了进一步降低执行跨度,某些任务需要重新分配并重新划分至 不同分割中。资源分配步骤的目标是实现任务与资源间代价 最有效的映射关系,确保资源总空闲时间最小化。

2.2 DAG 分割步骤

算法1是DAG分割步骤的伪代码。算法中,首先需要在考虑工作流的关键路径(CP)的执行时间情况下决定分配至每个分割的任务数量。由于关键路径上的任务执行时间之和(即关键路径长度,表示为*l*(CP))代表执行跨度 makespan 的下限(最优解),DAG分割需要确保属于相同分割中的任务的总执行时间小于或等于*l*(CP)(第8行的 while 循环)。由于处于关键路径上任务被调度至同一资源上,所以关键路径长度仅包括执行时间,而在任意其他分割(算法1中的*P'*)中的任务调度长度则包括执行时间和通信时间。分割*P*的长度(即代

价)定义为

$$l(P) = \sum_{i \in P} w_i^1 + \sum_{e(v_i, v_i) \in P} c_{(i,j)}$$
(3)

其中:等式右边第一项为分割 P 中的任务在最低价资源 w¹上的计算代价之和;第二项为分割 P 中任务间的通信代价之和。

DAG 分割步骤从工作流结构中最低层次的任务开始执行 (如图 1(b)中的工作流不同层次),该步骤添加任务至当前构 造的分割中,直到在不违背 *l*(CP)限制的情况下没有新任务添 加进来为止,即第 12 行。一旦属于相同层次的所有任务分配 至分割中,处于下一层次(高层次)的任务将继续进行分割,直 到所有任务被分配至不同分割中为止。

算法1 DAG 分割过程

1 procedure partitioning($G, l(CP)$)					
2	input: $G = \langle V, E \rangle$, $l(CP)$.				
3	output : $P = \langle p_1, p_2, \cdots \rangle$.				
4	\searrow (collection of partitions)				
5	$N \leftarrow \text{order tasks based on level}$				
6	while N is not empty do				
7	$P' \leftarrow \emptyset$				
8	while $l(P') < l(CP)$ do				
9	add first task in N to P'				
10	remove the task from ${\cal N}$				
11	end while				
12	if $l(P') > l(CP)$ then				
13	undo last step				
14	end if				
15	add P' to P				
16	end while				
17	return P				
18	end procedure				
	-				

2.3 分割调整步骤

为了确保 DAG 分割步骤中产生的任务分割在最终资源分 配过程中拥有最优的粒度,不同分割的任务需要重新分配调 整,寻找最优的分割主要涉及每个分割中的任务数量及其分配 的资源能力。由于任务分割与任务间的数据依赖存在相互关 系,不同分割间的任务重新分配时需要涉及 EST 和 EFT 时间 值的重新计算。对于每个任务,计算两个值分别为 sest 和 pest 。 任务 v_i 的 s_{est} 值表示 v_i 的子任务的平均 EST, 且子任务与 v_i 不 在同一分割中;任务 v_i 的 p_{est} 值表示与 v_i 同属一个分割中的任 务的平均 EST。基于这两个值,对于每个任务 $v_i \in V$,可以决定 是否通过将 v_i 重新分配至另一分割(拥有最小的EST(v_i))从 而调整 v_i 的分割以降低执行跨度。该步骤可以通过寻找 s_{est} ≤ pest的任务完成。这种分割调整对执行跨度和执行代价性能具 有重要影响,它可以降低由于处于不同分割的父任务的执行导 致其子任务的等待时间。如果重新分配的任务在其分割内拥 有子任务,则其子任务也需考虑重新分配至新的分割。仅当其 执行时间小于任务所在分割中任务的平均执行时间时,该任务 才考虑需要重新分配。通过该标准可以确保任务在重新分配 的分割中尽早执行。

将分割调整步骤得到的结果表示为有向无循环图 $G' = \langle V', E' \rangle$,其中,V'表示分割集合,E'表示分割间的数据依赖集合。任意两个分割 $v'_i, v'_j \in V'$ 通过一条以上的有向边连接,每条有向边 $e'(v'_i, v'_j) \in E'$ 表示任务 $v_p \in v'_j$ 和其在 v'_i 中的父任务;同时,边展示了属于 v'_j 的任务的最小 EST,即 e_{est} 。

2.4 资源分配步骤

算法2给出了资源分配步骤的伪代码。该步骤包括资源 集合识别和资源分配两个阶段。资源集合识别阶段中,需要决 定分配至每个分割的资源集合类型($i, R_i \in R$)。资源分配阶段 中,每个任务 v_j 分配至资源 $r \in R_i$,使得r属于与 v_j 所处分割相 同的资源集合类型;资源集合识别阶段的目标是寻找分配至分 割的资源集合类型,使得分割间的数据依赖关系达到最小化。 资源分配阶段的目标则是基于松驰参数β和所识别的资源集 合类型决定最快的调度方案。

资源集合识别阶段开始于将分割划分为不同的分割群组, 即第4行。一个分割群组包括一个群组领导 $P_i \in P$ 和与 P_i 连 接的父分割。拥有一个以上依赖关系的分割可属于多个群组 中,对于每个群组(行8中的 for 循环),从拥有出口分割的群 组开始作为一个领导群组(出口群组即包括出口任务的群 组),属于该群组的分割(除领导以外)开始识别其资源集合类 型(行9的 for 循环)。在每个群组中,每个分割的资源集合类 型 $(i, R_i \in R)$ 取决于群组领导的 e_{est} 值。为了决定资源集合类 型,每个分割需要识别其执行开始时间,即EST,将该时间记为 e,;此外,每个分割需要计算其任务的计算时间 cpi。获得这两 个值后,将每个分割的 $e_i + c_{Pi}$ 除以其群组领导的 e_{est} 值,即行 10。该除法结果表示为了降低执行跨度该分割要求执行时的 计算时间降低量,该降低量可以通过将与该除法结果最相近的 c。的资源集合类型至该分割来得到,即行11。如果任一分割 均包含单个分割,则该分割被分配至最低资源集合类型 R₁。 一旦确定了一个分割的资源集合类型,属于该分割中的每个任 务的 EST 和 EFT 即可重新计算。在该分割中的子任务的 EST 和 EFT 同时被重新计算,即行 12~14。分割识别其资源集合 类型后,如果该分割是一个群组领导,则它将告之其群组成员/ 分割开始执行资源集合识别阶段。当一个分割属于多个群组 时,该分割被分配至最高的资源集合类型。

所有分割确定其资源集合类型后,资源分配阶段基于每个 分割所识别的资源集合类型开始向任务分配资源,即行17~ 20。该阶段中,基于松驰参数β的值确定每个任务的 AST 和 AFT。参数 $\beta(0 \le \beta \le 1)$ 用于控制执行任务的资源数量,增加 该参数值可以降低资源使用的数量,此时相比执行跨度将有利 于执行代价的优化;而降低该参数值将有利于优化执行跨度。 该阶段开始于计算每个任务 AST 的上限 UB 和下限 LB,两个 值代表每个任务执行时所允许的延时。每个任务的 AST 根据 参数 β 进行设置,即行 18。如果 β = 1,则任务 AST 设置为 UB, 且算法利用最少数量的资源;如果 $\beta = 0$,则任务AST设置为 LB,且算法利用最多数量的资源。任务 v_i 的 LB 为(MIP_i) + $c_{(MIPi,i)}$,表示无资源限制时执行任务 v_i 的最早时间。当属于相 同分割的任务分配至相同资源时,任务 AST 为其上限 UB,在 这种情况下,任务 v_i 的 UB 为 AFT(v_i) + c_{ii} 。其中, v_i 为插入 分割中的 v_i 间的中间任务, $c_{i,i}$ 为任务 v_i 与 v_i 间的通信时间。 每次迭代中,拥有最低 AST 的未分配任务考虑为待分配任务, 该任务将分配至与任务所处分割相匹配且能保证及时执行的 资源集合类型中的资源,如果不存在该资源,则启用资源集合 类型中的新实例。

```
算法2 资源分配过程
    procedure reallocation((G', P, \beta, R))
1
        input: G' = \langle V', E' \rangle, P partitions,
2
3
                 \beta, R = R_1, \cdots, R_m resources.
        C {\leftarrow} \text{group } P \text{ based on successor partition ID}
4
5
           \searrow each group C_i has one or more partitions
        order {\it C} based on the depandancy between the tasks
6
7
            \subseteq C_i, C_i (i < j), C_i depends on C_i data
8
        for each C_i \in C do
9
              for each P_i \in C_i do
                   h \leftarrow (e_t + c_{Pi}) / e_{est}
10
11
                   call AssignResource (h, P_i, R)
                   for each v_i \in P_i do
12
                      call UpdateESTEFT(v_i, h, G')
13
14
                   end for
15
                end for
16
          end for
          for each v_i \in V do
17
               AST(v_i) \leftarrow \beta \times (UB_i - LB_i) + LB_i
18
```

- 19 $\operatorname{AFT}(v_i) \leftarrow \operatorname{AST}(v_i) + w_i$
- 20 end for
- 21 end procudure

2.5 算例说明

以图 2 为例对 PBWS 算法的思想进行阐述。为了简化描 述,算例中未显示计算和通信代价。图 2 中,(b)是 DAG 分割 步骤得到的结果,每个任务被分配至一个分割中;(c)是分割 调整步骤得到的结果,该步骤中任务 t。满足拥有最小 EST 条 件,被重新分配至分割 p_2 ;(d)的资源分配步骤中,每个分割 (除根分割)标志其 $e_i + e_{Pi}$ 为H及其 e_{est} 值。在该算例中, $H_1 =$ $40, e_{est1} = 20, H_2 = 20, e_{est2} = 20$ 。由于 p_3 是分割 p_1 和 p_2 的子分 割,这两个分割通过其H值与ees值相除识别其需要的资源集 合类型,这表明 p_1 的资源集合类型为 $2(R_2), p_2$ 的资源集合类 型为 $1(R_1)$ 。这表明在 p_1 中的所有任务被分配至 R_2 中的资 源上, p_2 中的所有任务被分配至 R_1 中的资源上。一旦相同群 组中的分割(拥有相同的后代)识别了其资源集合类型,即可 获得每个分割中任务的执行次序和每个分割的 eart,然后拥有 已识别资源集合类型的领导分割告之其群组成员开始执行资 源集合识别阶段。该算例中, p_1 为分割 p_2 的父分割,然而 p_1 并未作改变,由于所识别的资源集合类型 $2(R_2)$ 高于 $1(R_1)$, 任务的执行次序取决于所在的层次,如 p1 中任务的执行次序 为 t_{10} 、 t_7 、 t_8 、 t_4 。每个分割中任务的AST取决于 β 值。换言之, 如果分配单个资源至每个分割($\beta = 1$),则 p_1 中任务的 AST 分 别为0、10、20、30。该算例中假设p1中每个任务的计算代价为 10。如果假设 $\beta = 0, p_1$ 中任务的 AST 则分别为 0、10、10、20。 此时,t7 和 t8 分配至不同的资源上。



2.6 算法时间复杂度分析

PBWS 算法首先需要计算每个任务的 EST 和 EFT,时间复 杂度为 O(|V|);然后计算关键路径的时间复杂度为 O(|V|); DAG 分割步骤至多花费 $O(|V|^2)$ 可以收敛。分割调整步骤 中,对于每个任务,计算其子任务的平均 EFT 的时间复杂度为 $O(|V|^2)$;此外,重新分配任务的时间复杂度为 O(|V|(|P| - 1)),因此,分割调整步骤的时间复杂度总体为 $O(|V|^2 + |V|$ (|P| - 1))。资源分配步骤中,对分割进行排序的时间复杂度 $为 <math>O(|P|^2)$;将分割划分为群组的时间复杂度为 O(|P|),资源分配过程的时间复杂度为 <math>O(|P||R|),因此,资源分配步骤 的总体时间复杂度为 $O(|P|^2 + |P|(1 + |R|))$ 。在最差情况 下,DAG 分割数量等手任务数量,即一个任务为一个分割,因 此,PBWS 算法的最差时间复杂度为 $O(|V|^2 + |P|^2 + (|V||P)$ $|) + |P||R|) = O(|V|^2 + |V||R|)$ 。

3 仿真实验

本章通过仿真实验分析 PBWS 算法的性能,选择 IC-PCP^[4]和 HEFT^[11]两种算法作为基准算法,选择的云计算环境 仿真工具为 CloudSim,输入工作流由 Pegasus 工作流发生器得 到(https://confluence. pegasus. isi. edu/display/pegasus/WorkflowGenerator)。利用图 1 所示的五种现实世界中的科学工作 流进行实验测试,包括 CyberShake、Epigenomics、LIGO、Montage 和 SIPHT,为每种工作流配置 1 000 个任务。资源集合 $R = R_1$, R_2 , R_3 ,具体来说,最快资源 R_3 的价格是最慢资源 R_1 的三倍, 资源节点间的带宽设置为1 Gbps,且假设可用的资源数量是不 受限制的,β 的取值为(0,0.5,1)。

3.1 性能指标

实验主要观察算法的执行跨度和代价指标。代价包括任 务的执行代价、资源初始化代价和数据传输代价。显然,在调 度方案中增加资源数量对总体代价具有重要影响。因此,采用 标准化代价 N。比较该指标性能,考虑了任务执行代价和所使 用的资源数量。对于每个调度,将获得的代价除以在价格最低 的资源上的调度代价,该值代表相比最低价调度所获得调度方 案的代价程度;然后将该值乘以所获调度方案中所使用的资源 数量,即

$$N_c = \frac{c}{s} \times r \tag{4}$$

其中:c为资源上执行任务的数量;r为该调度中所使用的资源 数量;s。为任务调度至最低价资源上的代价。

3.2 实验结果

图 3 和 4 是三种算法在执行跨度和标准化代价方面的性能比较。可以看出,在代价上 PBWS 算法始终优于其他算法, 而 HEFT 算法在执行跨度上优于其他算法,这表明 PBWS 算法 为了降低代价牺牲了部分执行效率。具体地,本文根据五种工 作流的结构特征分成三个部分进行讨论。

a) CyberShake 工作流。对于 CyberShake, 当β值为0.5和 0时,即算法更倾向于最小化执行跨度时,PBWS和HEFT算法 拥有非常相近的执行跨度(图3(a)),且仅仅只带来较小的执 行代价的增加(图4(a))。具体地,由于降低β值可以导致更 多的资源使用,且 CyberShake 工作流(图1(d))不存在瓶颈任 务,PBWS所构造的多个分割能够并行执行从而降低了执行跨 度。尽管 IC-PCP 算法得到的执行跨度是整体最小的,但获得这 种更好的性能是以更多的资源使用和更高的代价为代价的(图 4(a))。此外,当β=0时 HEFT 算法拥有更高的代价,这是由于 HEFT 为了得到更高的执行跨度性能利用更高价的资源。

b) Epigenomics 和 LIGO 工作流。对于这两种工作流, PB-WS 算法获得与 IC-PCP 算法相近的执行跨度(图 3(b)(c)), 但 PBWS 拥有比 IC-PCP 更低的代价(图 4(b)(c))。对于标 准化代价,由于 IC-PCP 获得了与 β 无关的相同执行跨度,所以 期望得到的调度也拥有相同的代价。对于 PBWS,得到的执行 跨度取决于 β 值和工作流结构。降低 β 值可以增加使用资源 的数量,且可影响最终的执行跨度。然而只有在该值的降低使 得分配至每个分割中的资源数量增加才会拥有这种效果。两 种工作流在结构上拥有一致性,使得在满足数据依赖上可以使 分割的执行具有更好的连续性。因此在这种条件下,改变 β 值 并不会对 PBWS 的性能造成巨大影响。实验中设置 HEFT 使 用的资源数量与 PBWS 相同,对于 PBWS 算法,改变 β 值对利 用的资源数量带来了更小的变化,降低 β 值对于 HEFT 的执行 跨度影响很小。但总体而言,HEFT 作为一种贪婪算法,可以 获得更小的执行跨度,但比 PBWS 拥有更高的代价。

c) Montage 和 SIPHT 工作流。对于这两种工作流,结果表明 IC-PCP 算法在大多数情况下,其标准化代价均拥有最差的性能,这与工作流的结构相关。在 Montage 中,处于最后的层次中的任务形成了一种管道结构,这对 IC-PCP 具有巨大的负面影响,这主要源于此时最长关键路径的长度在决定算法效率上具有更重要的影响。在 SIPHT 中,由于出现了左侧路径,效果与 Montage 类似。

表1给出三种算法在两个性能指标上的具体值表现。其 结果再次证明 PBWS 算法得到的调度比较其他算法拥有更低

0.414 12 10 makespan/ $\times 10^2$ 6.4 6.3 0.2 0.1 8 6 4 2 0.5 0.5 β值 β值 (a)CyberShake工作流 (a)CyberShake工作流 0.4 24 22 20 18 16 14 12 makespan/ $\times 10^{2}$ 标准化代价 0.3 0.2 0.1 0.5 0.5 0 0 β值 *B*佰 (b)epigenomics工作流 (b)epigenomics工作流 16 14 12 0.8 makespan/ $\times 10^2$ 0.6 0.6 0.4 0.2 10 8 6 4 2 0.5 0 0.5 0 β值 β值 (c)LIGO工作流 (c)LIGO工作流 0.8 2 3.5 0.6 0.4 0.2 makespan/ × 2.3 1.5^{2} 0 0.5 0.5 0 0 β值 β值 (d)montage工作流 (d)montage工作流 10 makespan/ $\times 10^2$ 1.8 41.5 0.9 0.6 0.3 8 6 4 2 0.5 0.5 0 0 β值 β值 (e)SIPHT工作流 (e)SIPHT工作流 **PBWS** IC-PCP I HEFT B PBWS S IC-PCP I HEFT 图3 执行跨度 图4 标准化代价 表1 调度性能

的代价,且仅是以牺牲微小的执行跨度得到的。

工佐运米刊	算法	执行跨度		标准化代价	
工作机关型		$\beta = 0$	$\beta = 1$	$\beta = 0$	$\beta = 1$
	PBWS	194.88	1196.56	0.07	0.01
CyberShake	IC-PCP	107.22	213.90	0.14	0.18
	HEFT	185.22	580.19	0.35	0.06
	PBWS	23089.17	23186.92	0.07	0.07
epigenomics	IC-PCP	21919.12	21919.12	0.27	0.27
	HEFT	20583.27	20695.50	0.36	0.38
	PBWS	1394.06	1427.70	0.15	0.15
LIGO	IC-PCP	1389.97	1395.13	0.22	0.24
	HEFT	784.01	878.84	0.73	0.57
	PBWS	341.37	370.48	0.13	0.10
montage	IC-PCP	174.77	349.55	0.73	0.76
	HEFT	180.22	180.86	0.72	0.55
	PBWS	5261.95	7893.58	0.08	0.07
SIPHT	IC-PCP	3563.85	5261.95	1.04	1.09
	HEFT	2630.98	2630.98	1.66	0.82

总体来说,在不同类型的科学工作流结构下,PBWS的综合性能是较优化的。相比同类型算法,PBWS为了进行执行效率与代价的同步优化,将优化目标分两阶段进行,前一阶段通过初始任务DAG分割和分割结构调整后的任务重分配来提高执行效率,降低了执行跨度;后一阶段则在资源映射时优化了执行代价,使得资源空闲时间更少,资源利用更充分。虽然算法在单项指标上无法实现最优,但在多指标的综合性能上可以更加均衡,也更符合云环境的实际调度环境及工作流调度优化。

EEG signals [J]. Computers in Biology & Medicine, 2018, 100 (9):270-278.

- [7] Khan H, Marcuse L, Fields M, et al. Focal onset seizure prediction using convolutional networks[J]. IEEE Trans on Biomedical Engineering, 2018, PP(99);1-9.
- [8] Zhang Jian, Zou Junzhong, Wang Min, et al. Automatic detection of interictal epileptiform discharges based on time-series sequence merging method[J]. Neurocomputing, 2013, 110(8):35-43.
- [9] Wei Zuochen, Zou Junzhong, Zhang Jian. Automatic recognition of chewing noises in epileptic EEG based on period segmentation [J]. Neurocomputing, 2016, 190(5):107-116.
- [10] 李同庆, 邹俊忠, 张见, 等. 基于周期分割的睡眠自动分期研究
 [J]. 计算机工程与应用, 2019, 55(9):94-99. (Li Tongqing, Zhou Junzhong, Zhang Jian, et al. The Research of automatic staging of sleep based on period segmentation[J]. Computer Engineering and Applications, 2019, 55(9):94-99.)
- [11] Zhang Chong, Tan K C, Li Haizhou, et al. A cost-sensitive deep belief network for imbalanced classification [J]. IEEE Trans on Neural Networks & Learning Systems, 2018, PP(99):1-14.
- [12] Ma Li, Fan Suohai. CURE-SMOTE algorithm and hybrid algorithm for feature selection and parameter optimization based on random forests[J]. BMC Bioinformatics, 2017, 18(1):169.
- [13] Domingos P. MetaCost: a general method for making classifiers costsensitive [C]//Proc of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM Press, 1999: 155-164.
- [14] Su Chong, Ju Shenggen, Liu Yiguang, et al. Improving random forest and rotation forest for highly imbalanced datasets[J]. Intelligent Data Analysis, 2015, 19(6):1409-1432.
- [15] Olejarczyk E, Jozwik A, Zmyslowski W, et al. Automatic detection and analysis of the EEG sharp wave-slow wave patterns evoked by fluorinated inhalation anesthetics [J]. Clinical Neurophysiology Official Journal of the International Federation of Clinical Neurophysiology,2012,123(8):1512-1522.
- [16] Bingham E, Mannila H. Random projection in dimensionality reduction: applications to image and text data[C]//Proc of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data

(上接第页)

4 结束语

为了实现云计算中工作流的调度优化,本文提出了一种基于 DAG 分割的工作流调度算法。算法通过任务分割、分割调整及资源分配三个阶段的工作流调度模式,实现了任务执行跨度与执行代价的均衡调度。实验结果证明了所提算法能够以牺牲较小的执行效率开销得到更低的执行代价,其性能优于同类型算法。下一步的研究工作将集中于:将云资源方执行任务的能耗考虑到优化目标中,同时给出工作流执行的截止期限与用户费用约束,建立多约束条件下的多目标优化调度模型,并在新的调度模型下设计相应的 DAG 分割算法,实现多目标调度优化。

参考文献:

- [1] Pietri I, Malawski M, Juve G, et al. Energy-constrainted provisioning for scientific workflow ensembles [C]//Proc of International Conference on Cloud and Green Computing. Washington DC: IEEE Computer Society, 2013:34-41.
- [2] Mao Ming, Humphrey M. Auto-scaling to minimize cost and meet application deadlines in cloud workflows [C]//Proc of International Conference for High Performance Computing, Networking, Storage and Analysis. New York: ACM Press, 2014; articleNo 49.
- [3] Rodriguez M A, Buyya R. Deadline based resource provisioningand scheduling algorithm for scientific workflows on clouds [J]. IEEE Trans on Cloud Computing, 2014, 2(2):222-235.

Mining. New York: ACM Press, 2001:245-250.

- [17] Dasgupta S, Gupta A. An elementary proof of the Johnson-Lindenstrauss Lemma [J]. Random Structures & Algorithms, 1988, 36 (2):49-62.
- [18] Li Gen, Gu Yuantao. Restricted isometry property of Gaussian random projection for finite set of subspaces[J]. IEEE Trans on Signal Processing, 2018, PP(66):1705-1720.
- [19] Hoyos-Idrobo A, Varoquaux G, Thirion B. Fast brain decoding with random sampling and random projections [C]//Proc of International Workshop on Pattern Recognition in Neuroimaging. 2016:1-4.
- [20] Rodríguez J J, Kuncheva L I, Alonso C J. Rotation forest: a new classifier ensemble method [J]. IEEE Trans on Pattern Analysis and Machine Intelligence, 2006, 28(10):1619-1630.
- [21] 陆慧娟,刘亚卿,孟亚琼,等. 面向基因数据分类的核主成分分析 旋转森林算法[J]. 计算机科学与探索,2017,11(10):1570-1578. (Lu Huijuan, Liu Yaqing, Meng Yaqiong, et al. Classifier algorithm of genetic data based on kernel principal component analysis and rotation forest[J]. Journal of Frontiers of Computer Science and Technology,2017,11(10):1570-1578.)
- [22] Cieslak D A, Chawla N V. Learning decision trees for unbalanced data [C]//Proc of European Conference on Machine Learning and Knowledge Discovery in Databases. Berlin: Springer-Verlag,2008:241-256.
- [23] Shoeb A. Application of machine learning to epileptic seizure onset detection and treatment [D]. Cambridge: Massachusetts Institute of Technology, 2009.
- [24] 王凤琴,卢官明,柯亨进,等. 基于跨层全连接神经网络的癫痫发 作期识别[J]. 计算机应用研究,2019,36(7):2098-2103. (Wang Fengqin, Lu Guanming, Ke Hengjin, *et al.* Epileptic EEG identification with cross layer fully connected neural network [J]. Application Research of Computers,2019,36(7):2098-2103.)
- [25] Supratak A, Li Ling, Guo Yike. Feature extraction with stacked autoencoders for epileptic seizure detection [C]//Proc of the 36th Annual International Conference of IEEE Engineering in Medicine and Biology Society. Piscataway, NJ: IEEE Press, 2014:4184-4187.
- [26] Ahammad N, Fathima T, Joseph P. Detection of epileptic seizure event and onset using EEG[J]. Biomed Research International, 2014,2014(1):450573.
- [4] Abrishami S, Naghibzadeh M, Epema D H J. Deadline-constrained workflow scheduling algorithms for Infrastructure as a service clouds[J].
 Future Generation Computer Systems, 2013, 29(1);158-169.
- [5] Wu Qishi, Lin Xiangyu, Yu Dantong, *et al.* End-to-end delay minimization for scientific workflows in clouds under budget constraint
 [J]. IEEE Trans on Cloud Computing, 2015, 3(2):169-181.
- [6] Zeng Lingfang, Veeravalli B, Li Xiaorong. ScaleStar: budget conscious scheduling precedence-constrained many-task workflow applications in cloud [C]//Proc of the 26th IEEE International Conference on Advanced Information Networking and Applications. Washington DC: IEEE Computer Society, 2012:534-541.
- [7] Arabnejad H, Barbosa J G. A budget constrained scheduling algorithm for workflow applications [J]. Journal of Grid Computing, 2014,12(4):665-679.
- [8] Zheng Wei, Sakellariou R. Budget-deadline constrained workflow planning for admission control [J]. Journal of Grid Computing, 2013,11(4):633-651.
- [9] Lee Y C, Zomaya A. Stretch out and compact: workflow scheduling with resource abundance[C]//Proc of the 13th IEEE/ACM International Symposium on Cluster, Cloud, and Grid Computing. Washington DC: IEEE Computer Society, 2013:219-226.
- [10] Sakellariou Red, Zhao Hao. A hybrid heuristic for DAG scheduling on heterogeneous systems [C]//Proc of the 18th International Parallel and Distributed Processing Symposium. Washington DC: IEEE Computer Society, 2014:111-114.
- [11] Topcuoglu H, Hariri S, Wu Minyou. Performance-effective and lowcomplexity task scheduling for heterogeneous computing [J]. IEEE Trans on Parallel & Distributed Systems, 2012, 13(3):260-274.