

Received 15 October 2013; revised 18 May 2014; accepted 22 July 2014. Date of publication 17 August 2014;
date of current version 30 October 2014.

Digital Object Identifier 10.1109/TETC.2014.2348196

Evolutionary Scheduling of Dynamic Multitasking Workloads for Big-Data Analytics in Elastic Cloud

FAN ZHANG^{1,7}, (Senior Member, IEEE), JUNWEI CAO², (Senior Member, IEEE),
WEI TAN³, (Senior Member, IEEE), SAMEE U. KHAN⁴, (Senior Member, IEEE),
KEQIN LI⁵, (Senior Member, IEEE), AND ALBERT Y. ZOMAYA⁶, (Fellow, IEEE)

¹Kavli Institute for Astrophysics and Space Research, Massachusetts Institute of Technology, Cambridge, MA 02139 USA

²Research Institute of Information Technology, Tsinghua University, Beijing 100084, China

³IBM T. J. Watson Research Center, Yorktown Heights, NY 10598 USA

⁴North Dakota State University, Fargo, ND 58102 USA

⁵Department of Computer Science, Tsinghua University, Beijing 100084, China

⁶School of Information Technologies, The University of Sydney, Sydney, NSW 2006, Australia

⁷Shenzhen Institute of Advanced Technology, Chinese Academy of Science, Guangdong 51800, China

CORRESPONDING AUTHOR: F. ZHANG and J. CAO (zf19830622@gmail.com; jcao@mail.tsinghua.edu.cn)

This work was supported in part by Ministry of Science and Technology of China under National 973 Basic Research Program (grants No. 2013CB228206 and No. 2011CB302505), National Natural Science Foundation of China (grant No. 61472200 and No. 61233016) and National Science Foundation under grant CCF-1016966.

ABSTRACT Scheduling of dynamic and multitasking workloads for big-data analytics is a challenging issue, as it requires a significant amount of parameter sweeping and iterations. Therefore, real-time scheduling becomes essential to increase the throughput of many-task computing. The difficulty lies in obtaining a series of optimal yet responsive schedules. In dynamic scenarios, such as virtual clusters in cloud, scheduling must be processed fast enough to keep pace with the unpredictable fluctuations in the workloads to optimize the overall system performance. In this paper, ordinal optimization using rough models and fast simulation is introduced to obtain suboptimal solutions in a much shorter timeframe. While the scheduling solution for each period may not be the best, ordinal optimization can be processed fast in an iterative and evolutionary way to capture the details of big-data workload dynamism. Experimental results show that our evolutionary approach compared with existing methods, such as Monte Carlo and Blind Pick, can achieve higher overall average scheduling performance, such as throughput, in real-world applications with dynamic workloads. Furthermore, performance improvement is seen by implementing an optimal computing budget allocating method that smartly allocates computing cycles to the most promising schedules.

INDEX TERMS Big-data, cloud computing, evolutionary ordinal optimization, multitasking workload, virtual clusters.

I. INTRODUCTION

Large-scale business and scientific applications are usually composed of big-data, multitasking, time-variant, and fluctuating workloads [5], [35]. Cloud computing [2], with virtualization [3] as the key enabling technology, provides an elastic scaling-up and scaling-down provisioning mechanism. Agile and appropriate computational resource provisioning that keeps pace with the fluctuations of big-data multitasking workloads is very important in the scheduling paradigm [41].

For example, there are unprecedented amounts of requests from the customers of Amazon and eBay during the holiday seasons. Thereafter, the website traffic drops down dramatically. Over-provisioning of computing resource to constantly satisfy the requirements at the peak level leads to high and unnecessary cost, while under-provisioning leads to user churn [28].

In general, scheduling big-data multitasking workloads onto distributed computing resources is an NP-complete

problem [13]. The main challenge is to maintain a reasonable tradeoff between the scheduling overhead and scheduling accuracy.

Consider a continuous scheduling scenario with multiple scheduling periods, and within each period the scheduling scheme is being repeatedly updated. If the scheduling period is too long, while a seemingly more descent solution can be achieved through elaborate workload analysis, the overall system performance may degrade. This is due to the fact that the workloads might have changed considerably during such a long time. Therefore, it is necessary that scheduling solutions are provided in an evolving fashion, so that during each iteration, a “good-enough”, suboptimal, but fast-paced solution can be obtained. In the meanwhile, evolving iterations of optimization can adapt to the nature/details of system dynamism, such as dynamic workload fluctuation and resource provisioning in virtual clusters to achieve better performance.

In our previous work [40]–[43], for the purpose of fast scheduling, we have explored the possibility of applying simulation-based optimization methods. Ordinal Optimization (OO) [15] has been applied for suboptimal but fast searching. The OO is utilized for searching suboptimal solutions in much shorter time and with reduced runs and analysis of the workloads. Targeting at a sequence of suboptimal schedules instead of an optimal one results in much lower scheduling overhead by reducing the exhaustive searching time [40]. Experiments show that the methods capture workload characteristics and lead to an overall improved performance. However, the scheduling overhead of applying those methods still refrains its applicability in highly fluctuated workloads.

In this paper, we propose an iterative and evolutionary usage of ordinal optimization to further reduce the scheduling overhead. Our major contribution consists of the following:

- The iterative OO (iOO) method published in our earlier paper [40], is applied in the new and real multitasking-scheduling model in this paper. As expected, the iOO method shows up to 20% performance speedup than competing methods.
- In a series of rapidly fluctuating workload periods, we contribute a procedure of using short-phase scheduling for fine-grained workload analysis. In stable workload phases, long-term scheduling is accordingly contributed which intends to save time to analyze the workloads for accurate scheduling.
- A step further, as an extension of iOO, the evolutionary OO (eOO) analyzes the workload patterns among consecutive phases and adjusts the scheduling based on the patterns. We develop a series of algorithms to partition and merge workload for efficient scheduling.
- We use a dynamic scenario of scheduling multitask scientific workloads to a group of virtual clusters on Amazon EC2 cloud. The experiment results show that the eOO approach achieves up to 30% performance speedup regarding task throughput, compared with

Monte Carlo [23] and Blind Pick. As far as we know, this is the first time an OO is applied in an evolutionary way for dynamic optimization scenarios to meet special requirements of cloud workload scheduling.

The rest of the paper is organized as follows. Section II provides a brief description of the dynamic multitasking workload scheduling problem and existing methods. In Section III, we describe our proposed eOO technique that iteratively applies the OO for scheduling workloads on virtual clusters of a cloud. The experimental results, performance evaluation, and comparative analysis are presented in Section IV that also includes detailed information on the applications and system configurations. The related work is reviewed in Section V, and we provide concluding remarks in Section VI.

II. DYNAMIC WORKLOAD SCHEDULING

In this section, first, we introduce our dynamic multitasking-scheduling model. Thereafter, the necessity of using simulations to is presented. To follow, three existing approaches, Monte Carlo, Blind Pick, and iterative Ordinal Optimization, are introduced. To each readership, Table 1 summarizes the most frequently used symbols, notations, and definitions.

TABLE 1. Notations of workload scheduling.

Notations	Definition and Description
T_i	Throughput at the i -th scheduling period
S_i	The i -th scheduling period
$\Delta\delta_c(t)$	New generated workload for VC_c at time t
$\delta_c(t)$	Remaining workload for VC_c at time t
$\delta_c^t(t)$	Total workload for VC_c at time t
$CET_c(t_{i,m_i}, t_{i,m_{i+1}})$	Makespan in VC_c from time t_{i,m_i} to $t_{i,m_{i+1}}$
$EET(t_{i,m_i}, t_{i,m_{i+1}})$	Makespan from time t_{i,m_i} to $t_{i,m_{i+1}}$
c or C	Index of a VC or the number of VCs
i or I	Index of time or the last scheduling point
$\theta_c(t_i)$	Number of VMs in VC_c that are allocated at time t_i for the S_i stage
θ	Number of VMs available for scheduling
$p_c(t_i)$	Expected execution time in VC_c at time t_i
$\beta_c(t_i)$	Job processing rate in VC_c at time t_i
$r_c(t_i)$	Remaining execution time in VC_c at time t_i

A. DYNAMIC MULTITASKING WORKLOAD SCHEDULING MODEL

In this model, we define the *task class* as a set of tasks that are of the same type and can be executed concurrently. Suppose there are C task classes in all, and the index is denoted as c , $c \in [1, C]$. Tasks within one task class can be either interdependent, such as scientific workflow, or independent of each other, such as scientific simulations with multiple parameters. Tasks across different task classes are independent.

Virtual Machines (VMs) are the scheduling and computing units built on top of the physical machines. Given a fixed number of VMs, our scheduling method organizes them into C groups, each group serves for one task class. Each group

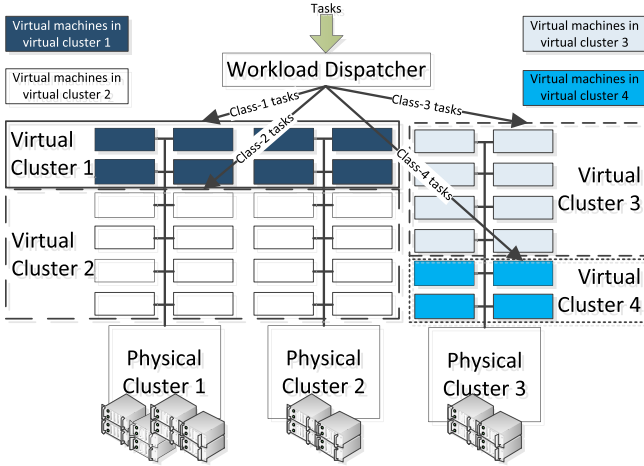


FIGURE 1. Virtual cluster resource allocation model for workload execution in a virtualized cloud platform. Three physical clusters with each cluster having twelve VMs are shown. The 36 VMs are topologically partitioned into four VCs. The workload dispatcher assigns the class- c task to the c -th VC to execute.

is called a Virtual Cluster (VC). All of the tasks in the c -th task class are assigned to the c -th VC. In this way, we must partition all of the VMs into C VCs to serve the C task classes.

Consider a virtualized cloud platform with $C = 4$ VCs as illustrated in Fig. 1. The workload dispatcher directs the tasks to the related VC for execution. A resource-reservation *schedule* specifies the sets of VMs to be provisioned at continuous time periods. For example, the i -th schedule $\theta(t_i)$ is represented by a set of VMs allocated in C clusters in a *schedule space* U . A C -dimensional vector represents this schedule:

$$\theta(t_i) = [\theta_1(t_i), \theta_2(t_i), \dots, \theta_c(t_i), \dots, \theta_C(t_i)], \quad (1)$$

where $\theta_c(t_i)$ is the number of VMs assigned in VC c , and θ is the total number of VMs that can be allocated. Therefore, we have $\sum_{c=1}^C \theta_c(t_i) = \theta$.

In Fig. 2, we show stacked workloads that are dispatched to the C VCs in a timeline. From t_{i-1} (or also named $t_{i-1,0}$), schedule $\theta(t_{i-1})$ is applied until t_i (or $t_{i,0}$), where a new schedule $\theta(t_i)$ is used. The new schedule $\theta(t_i)$ is generated during the previous simulation stage from t_{i-1} to t_i . This stage is named S_{i-1} as shown in the figure. Between t_i and t_{i+1} , new workloads (at time points $t_{i,1}, t_{i,2}, \dots$) arrive that are also shown in the figure. Therefore, the dynamic workload scheduling model is built on such a sequentially overlapped simulation-execution phases. In each phase, one schedule is applied and in the meanwhile, the workloads of the subsequent stage are analyzed to simulate and generate the schedule of the following stages.

We use $\Delta\delta_c(t)$ to denote the newly produced workload of the i -th task class at any time t . The $\Delta\delta_c(t)$ is an incremental value, which represents the accumulated unfinished workloads from the previous stages. As shown in Fig. 2, $\Delta\delta_c(t_{i+1})$ denotes the workload generated at t_{i+1} for the i -th VC.

The real workload at t is the sum of the remaining workload of the previous stages and the newly generated workload at time t . We use $\delta_c(t)$ and $\delta_c(t)$ to denote the real and remaining workload at time t for the c -th VC, which gives us:

$$\delta_c(t) = \delta_c(t) + \Delta\delta_c(t) \quad c \in [1, C].$$

The aforementioned results are $\delta(t) = [\delta_1(t), \delta_2(t), \dots, \delta_c(t), \dots, \delta_C(t)]^T$.

The tasks arriving at different VCs are time-variant. Our target is to look for a time series of schedules:

$$\theta(t) = [\theta(t_1), \theta(t_2), \dots, \theta(t_i), \dots, \theta(t_I)], \quad (2)$$

to maximize the throughput of the whole workload. The throughput is defined as the ratio of the total number of tasks finished and the time spent by the VCs to process the tasks. To calculate throughput, we need to find out the schedule time points, as well as the corresponding schedules for each time point.

From time point t_i to t_{i+1} , the phrasal optimal throughput given the workload and phase can be represented by the following:

$$\begin{aligned} T_i^* &= \max T(\theta(t_i)|\delta(t), S_{i-1}) \\ &= \max T(\theta(t_i)|\delta(t), t_i, t_{i-1})t \in [t_i, t_{i-1}), \theta(t_i) \in U \end{aligned}$$

We first introduce the concepts of *Class Execution Time (CET)* and *Effective Execution Time (EET)*. Suppose that $t_{i,0}$ and $t_{i,1}$ are the time points of two consecutive workload stages, as shown in Fig. 2. $CET_c(t_{i,0}, t_{i,1})$ is the runtime from the starting point $t_{i,0}$ to the finishing time when all the tasks in task class c are finished. We show $CET_c(t_{i,0}, t_{i,1})$, $c = \{1, 2, C\}$ in Fig. 2 by the three double end-arrow dashed lines. If $CET_c(t_{i,0}, t_{i,1})$ is larger than the scheduling period length $t_{i,1} - t_{i,0}$ which means all the tasks in this task class cannot be finished within $[t_{i,0}, t_{i,1}]$, then $t_{i,1} - t_{i,0}$ is used as CET and the rest tasks are rolled over to $\Delta\delta_c(t_{i,1})$. $EET(t_{i,0}, t_{i,1})$, on the other hand, means the longest execution time across all

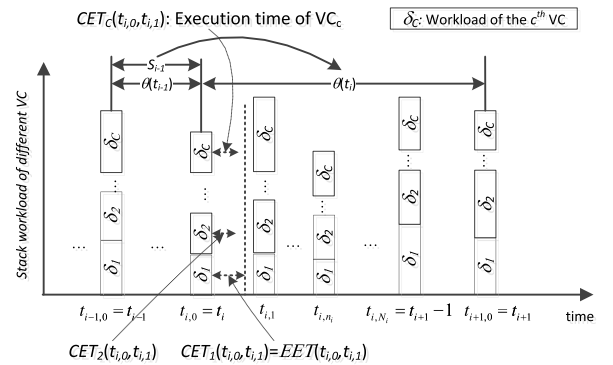


FIGURE 2. Demonstration of the stacked workload of different VC, given in a timeline. Each rectangular box represents the workload generated for the corresponding virtual cluster. Between schedule periods $[t_i, t_{i+1}]$, there are time points at which new workloads are generated (or arrive). The time points are represented as $t_{i,1}, t_{i,2}, \dots, t_{i,N_i}$.

of the task classes within $[t_{i,0}, t_{i,1}]$. It can be formally defined as follows.

$$EEF(t_{i,n_i}, t_{i,n_i}) = \begin{cases} \max CET_c(t_{i,n_i}, t_{i,n_i+1}) \\ \quad \forall CET_c(t_{i,n_i}, t_{i,n_i+1}) < t_{i,n_i+1} - t_{i,n_i} \\ t_{i,n_i+1} - t_{i,n_i} \\ \quad \exists CEF_c(t_{i,n_i}, t_{i,n_i+1}) \geq t_{i,n_i+1} - t_{i,n_i} \end{cases}$$

As shown in Fig. 2, $EET(t_{i,1}, t_{i,0})$ equals to $CET_1(t_{i,1}, t_{i,0})$. Therefore, the period throughput from t_i to t_{i+1} is calculated by:

$$T_i^* = \frac{\sum_{j=1}^{N_i-1} \sum_{c=1}^C (\delta_c(t_{j-1}) - \underline{\delta}_c(t_j))}{\sum_{j=0}^{N_i-1} EET(t_{i,j}, t_{i,j+1})}$$

Let $\delta_c(t_i)$ be the number of tasks in the c -thVC at time t_i , $p_c(t_i)$ be the expected execution time for each VM in the c -th VC at time t_i . Consequently, $\beta_c(t_i) = \theta_c(t_i)/p_c(t_i)$ represent the corresponding job processing rate, where $\theta_c(t_i)$ is the number of VMs in VC_c that are allocated at time t_i and $r_c(t_i) = \delta_c(t_i)/\beta_c(t_i)$ is the remaining execution time in the c -thVC at time t_i .

The task throughput is used as the performance metric for the schedules. Given a schedule $\theta(t) = [\theta(t_1), \theta(t_2), \dots, \theta(t_i), \dots, \theta(t_l)]$, the task throughput is calculated as the total number of finished tasks divided by the total task execution time (or the *makespan*). At different time periods, different schedules may be applied. All of the candidate schedules at each successive time period form a schedule space U . The *cardinality* of U is calculated by the following expression:

$$u = (\theta - 1)! / [(\theta - C)! (C - 1)!], \quad (3)$$

where θ is the total number of VMs used in C virtual clusters. The parameter u counts the number of ways to partition a set of θ VMs into C nonempty clusters.

For example, if we use $\theta = 20$ VMs in $C = 7$ VCs for seven task classes, then we need to assess $u = 27,132$ possible schedules to search for the best schedule at each time t_i . Each schedule $\theta(t_i)$ takes the form of seven dimensional tuple, having a sum equal to 20, such as [2, 3, 3, 5, 2, 3, 2] or [6, 2, 1, 3, 2, 4, 2]. The search space U is very large given the large number of VMs θ , which leads to the ineffectiveness of the scheduling method. Therefore, the schedule search space must be reduced significantly, especially in a dynamic environment.

B. SIMULATION-BASED OPTIMIZATION METHODS

In this section, we introduce a series of existing simulation-based optimization solutions for this problem, namely Monte Carlo Simulation, Blind Pick, and iOO.

Simulations are needed due to the fact that the expected execution time $p_c(t_i)$ is unknown in real-time. The value is stochastic and may be subjected to unknown runtime conditions, such as the contention of other virtual resources.

Algorithm 1 Monte Carlo Simulation Method

Input:

$U = \{\theta = [\theta_0, \theta_1, \theta_2, \dots, \theta_C]\}$ //All of the possible schedules
 $\delta = \{[\delta(t_{i,0}), \delta(t_{i,1}), \dots, \delta(t_{i,0}), \dots, \delta(t_{i,N_i})]\}$ //Workloads
 α //The ratio of the search space U should be sampled

Output:

The best $\theta^*(t)$ for use at each schedule time t

Procedure:

1. Random sample U as per the ratio α ; //Generate search space
2. Estimate the overhead time of using Monte Carlo method: O_M ;
3. Initialize schedule period value $i = 1$;
4. **while** ($i \times O_M < t$) **do** //Not the last schedule period
5. $t_i = (i-1) \times O_M$, $t_{i+1} = i \times O_M$;
6. Load workload $\delta_i(t_i)$ between $[t_i, t_{i+1}]$; // $\delta_c(t_{i,0}), \delta_c(t_{i,1}), \dots, \delta_c(t_{i,N_i})$
7. **for each** θ in U **do** // Search Each Schedule
8. **for each** j in $[0, N_i]$ **do** //Each time point workload generated
9. **for each** VC c in $[1, C]$ **do** //Each virtual cluster
10. Simulate $p_c(t_i)$ N times to get an average $P_c(t_i)$;
11. $\delta_c(t_{i,j}) = \underline{\delta}_c(t_{i,j}) + \Delta \delta_c(t_{i,j})$; //Current tasks for VC_c
12. **if** ($\delta_c(t_{i,j}) - (t_{i,j+1} - t_{i,j}) \times P_c(t_i) > 0$) //Not all the tasks finished
13. $\underline{\delta}_c(t_{i,j+1}) = \delta_c(t_{i,j}) - (t_{i,j+1} - t_{i,j}) \times P_c(t_i)$; //Remaining tasks
14. $CET_c(t_{i,j}, t_{i,j+1}) = t_{i,j+1} - t_{i,j}$;
15. **else** //All the tasks finished
16. $\underline{\delta}_c(t_{i,j+1}) = 0$;
17. **end if**;
18. $CET_c(t_{i,j}, t_{i,j+1}) = (\delta_c(t_{i,j}) - (t_{i,j+1} - t_{i,j}) \times P_c(t_i)) / P_c(t_i)$;
19. **end for**;
20. $EET(t_{i,j}, t_{i,j+1}) = \max(CET_c(t_{i,j}, t_{i,j+1}))$;
21. $Makespan(\theta(t_i)) = EET(t_{i,0}, t_{i,1}) + \dots + EET(t_{i,N_i-1}, t_{i,N_i})$;
22. **end for**;
23. **if** ($Makespan(\theta(t_i)) < Makespan(\theta^*(t_i))$) //A better schedule
24. $\theta^*(t_i) = \theta(t_i)$;
25. $T^*(t_i) = (\delta_c(t_{i,0}) + \delta_c(t_{i,1}) + \dots + \delta_c(t_{i,N_i}) - \underline{\delta}_c(t_{i,1}) + \dots + \underline{\delta}_c(t_{i,N_i})) / \theta^*(t_i)$;
26. **end if**;
27. **end for**;
28. $i = i + 1$; //Next schedule period
29. **end while**

Moreover, the random runtime workload distribution in the c -thVC also has a huge impact on its value.

For a given number of cloud VM instances, configuration (2ECU in Amazon EC2, 1.7GB Memory), VM utilization value, and failure rate, we evaluate the execution time distribution of a particular task class with a number of tasks, and correspondingly create a table.

For example, the table tells that the execution time is a normal distribution ($X \sim N(20, 5^2)$) when running [100, 110] tasks of task class two with 10 VMs, each being small Amazon EC2 instance, failure rate being 0 and utilization being in a range of [70%, 80%]. Then, before each experiment runs, if we have profiled the average VM utilization is 75%, and there are 105 tasks and 10 VMs, then we must sample the normal distribution above to estimate $p_c(t_i)$, and apply this sample value to estimate the throughput of each schedule, and finally choose the best schedule for use. Algorithm 1 below introduces the process.

To yield such an estimation, however, multiple simulation runs must be conducted and an observed average is to be used for $p_c(t_i)$. This is a typical Monte Carlo simulation method that is known to be time consuming.

1) *Monte Carlo Method*: Suppose that the overhead of using the Monte Carlo method is denoted by O_M , which is the length of the scheduling period. Thereafter, we search through all of the time points between t_i and t_{i+1} , when new workloads are generated, to calculate the *CET* for each VC as shown in Line 14 and 18. This step is followed by calculating the *EET* during the scheduling period, as shown in Line 20. Following this, the Makespan and throughput are calculated.

2) *Blind Pick Method*: Instead of searching through the whole schedule space U as done in the Monte Carlo method, the Blind Pick (BP) method, randomly selects a portion of the schedules only within U for evaluation. The ratio to be selected is defined by a value α ($0 \leq \alpha \leq 1$). The algorithm of applying BP is exactly the same as Monte Carlo, except for the scheduling sample space U and scheduling overhead. The length of the scheduling period is $\alpha \times O_M$.

3) *Iterative Ordinal Optimization Method*: The OO, a sub-optimal low overhead scheduling method is thoroughly introduced in [15]. Different from the aforementioned methods, the OO uses a rough model (n repeated runs) to generate a rough order of all of the schedules in U , and uses the accurate model (N repeated runs, $n \ll N$) to evaluate the top-schedules in the rough order list. For the values of s , the readers are referred to Section III.B.

The OO method contains two steps. The rough evaluation step evaluates u schedules, each one being repeated n times, and the overhead becomes $u \times O \times n$. Suppose O denotes the time needed to evaluate one schedule, then the second stage will take $s \times O \times N$ amounts of time.

We have proven in our previous paper [40] that the OO simulation time is shorter than other two methods that enable us to develop an iterative OO (iOO) method. In other words, the iOO applies the OO method iteratively in a multi-stage scheduling manner, with each stage being a new OO routine.

The iOO approach exhibited superior performance in the research paper, especially in the highly fluctuated workload cases due to the reduced overhead, resulting in shorter scheduling periods and fine-grained scheduling. The length of the iOO scheduling period can be determined as $u \times O \times n + s \times O \times N$.

However, the iOO does not consider the characteristics of workloads, namely the similarity of the workloads in consequent stages of a multi-stage scheduling problem. In the next section, we introduce the evolutionary Ordinal Optimization (eOO) that extends the iOO method along those lines.

III. EVOLUTIONARY ORDINAL OPTIMIZATION

An intuitive interpretation of the eOO is the assumption that the simulation runs must be allocated based on the characteristics of the workloads. For example, fluctuating workload requires more intermediate scheduling time points. Because each of the short scheduling phases satisfies small periods of the workload, the eOO delivers better overall performance.

However, one disadvantage of the eOO method is that the short simulation phase leads to less simulation runs of

the $p_c(t_i)$, which on the contrary, might decrease the performance of each of the single schedule period. In this section, we demonstrate the methodologies and effectiveness of the eOO method, and propose a solution to partition and merge the scheduling periods based on the fluctuations within the workload.

A. EVOLUTIONARY ORDINAL OPTIMIZATION

Before we step into any further detail, we define the concept of *workload pattern*. At time t_0 in Fig. 3, one batch of workload for the seven VCs arrives. This workload pattern is very similar to the workload patterns at t_1 and t_2 . To give a formal definition, suppose that there are two batches of workloads at time t' , $t'' \in [t_i, t_{i+1}]$, i.e., $\delta(t') = [\delta_1(t'), \dots, \delta_c(t'), \dots, \delta_C(t')]$ and $\delta(t'') = [\delta_1(t''), \dots, \delta_c(t''), \dots, \delta_C(t'')]$, respectively, where $\delta_c(t')$ denotes the tasks for VC c at time t' . The similarity pattern between $\delta(t')$ and $\delta(t'')$ can be defined as:

$$\text{Sim}(\delta(t'), \delta(t'')) = \frac{\sum_{c=1}^C (\delta_c(t') \times \delta_c(t''))}{\sqrt{\sum_{c=1}^C \delta_c(t')^2 \times \sum_{c=1}^C \delta_c(t'')^2}}. \quad (4)$$

Intuitively, a high similarity pattern leads to the merger of two scheduling periods, while a low similarity pattern leads to a scheduling period partition. In Fig. 3, we would merge $[t_0, t_1]$, $[t_1, t_2]$, and $[t_2, t_3]$ as a single scheduling period, when supposing that t is the scheduling period of any method proposed in the previous section. However, such a merge does not apply to the period $[t_3, t_3 + t]$, where the two batches of workloads are dramatically different from each other. Therefore, we must divide the scheduling period adaptively into smaller periods for separate scheduling.

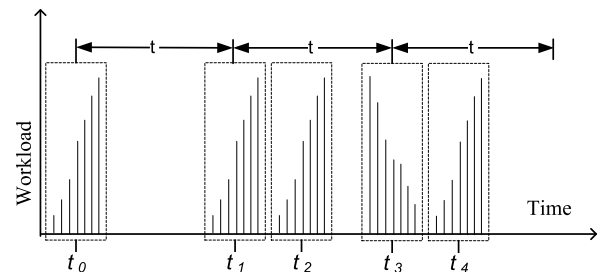


FIGURE 3. Demonstration of various workload characteristics. The dashed-line rectangular box with seven solid lines are the seven workloads for the $C = 7$ VCs. There are five batches of workloads arrive at time t_0, t_1, t_2, t_3 and t_4 for each of the virtual clusters.

The major tenets of the eOO lies in the automatic scheduling period partitioning and merging as introduced below.

(1) A pairwise comparison of the workload pattern must be made between each of the pairs of the workload in each scheduling period, such as $[t_i, t_{i+1}]$. If the maximum similarity across all the workload patterns is sufficiently small, then the interval must be sliced into two, or more, consecutive scheduling periods. Otherwise the period is just kept as it is.

(2) Consider two workload patterns, $\delta(t')$ and $\delta(t'')$ that are extracted from two consecutive scheduling periods $[t_i, t_{i+1}]$

Algorithm 2 Scheduling Period Partition

Input: All workloads in $[t_i, t_{i+1}]$ denoted by $\delta(t_i, t_{i+1}) = \{\delta(t_{i,0}), \delta(t_{i,1}), \dots, \delta(t_{i,N})\}$.

Output: Decision point (*DP*) set in $[t_i, t_{i+1}]$.

Procedure:

1. $DP \leftarrow \text{null}$;
2. **if** $((t_{i+1} - t_i) < w)$ //Interval is too small, cannot be partitioned anymore
3. Add t_i into DP ;
4. **return**;
5. **end if**;
6. **for each** workload $\delta(t_i)$ in $\delta(t_i, t_{i+1})$ **do**
7. **for each** workload $\delta(t_k) \neq \delta(t_i)$ in $\delta(t_i, t_{i+1})$ && $k > i$ **do**
8. **if** $(\text{Max}(\text{Sim}(\delta(t_i), \delta(t_k))) < \alpha)$ //small similarity
9. Partition($t_i, t_i + (t_{i+1} - t_i)/2$); //partition left half
10. Partition($t_i + (t_{i+1} - t_i)/2, t_{i+1}$); //partition right half
11. **end if**;
12. **end for**;
13. **end for**;
14. Sort and return DP .

and $[t_{i+1}, t_{i+2}]$, respectively, where $t' \in [t_i, t_{i+1}]$ and $t'' \in [t_{i+1}, t_{i+2}]$. We make similarity analysis among all of the possible pairs as mentioned above. If the least similarity value is sufficiently large, then the two scheduling periods must be merged. Such a merge operation leads to a sufficiently large simulation period of $[t_i, t_{i+2}]$ that promises a longer simulation period for the subsequent period.

In Algorithm 2 and Algorithm 3, we show the details of the scheduling period merge and partition based on the workload similarities. In general, for Algorithm 2, if the maximum similarity value between any two pairs of workloads in $[t_i, t_{i+1}]$ is less than α , then the interval will be partitioned into two parts. For the aforementioned case, we use a recursive algorithm, as shown in Line 9 and Line 10, to further analyze the two sub-intervals. The analysis stops when either the sub-interval is smaller than a threshold value w , or all of the workload patterns are similar in the interval that is being checked. The algorithm outputs all of the decision time points DP .

B. PARAMETER CONFIGURATION

Unknown parameters must be finalized before the experimental runs. These are the simulation runs N and n , size of selected set s_i in iOO, size of the selected set s_e using eOO, size of the selected set s_{BP} using BP, α and β in Algorithm 3 and Algorithm 4.

The BF Monte Carlo repeats itself N times to estimate the average value of each $p_c(t_i)$. As a common practice to achieve reliable simulation accuracy, Monte Carlo usually applies a large simulation runs, say $N = 1000$, for each single individual schedule for large-scale scientific workflow scheduling problems [42].

The rough model n of the iOO is also determined by the statistic theory. Based on the central tendency theory, it takes two orders of magnitudes more runs to improve one order of

Algorithm 3 Scheduling Period Merge

Input: Output of Algorithm 2: DP .

Output: New decision point set (NDP).

Procedure:

1. **for each** dp in DP **do** // Each existing decision point dp
2. Find dps associated with dp ; //decision point stage
3. **if** $(\text{Width}(dps) > W)$
4. Goto Step 1; //large decision stage, no more merge
5. **end if**;
6. **if** (dp is not the last stage)
7. Find dps after dp ; //next decision point
8. Find dps associated with dps ;
9. **for each** workload $\delta(t_i)$ in ds **do**
10. **for each** workload $\delta(t_k)$ in dps **do**
11. **if** $(\text{Min}(\text{Sim}(\delta(t_i), \delta(t_k))) > \beta)$ //large similarity
12. Remove dps from DP ;
13. Associate dp with $(dps + dps)$; //merge stages
14. Goto Step 4;
15. **else**
16. Goto Step 1;
17. **end if**;
18. **end for**;
19. **end for**;
20. **end if**;
21. **end for**;
22. Return ($NDP \leftarrow DP$).

magnitude of the estimation accuracy. When applying $n = 10$ runs in the rough model, which is two orders of magnitude less simulation runs than N , would lead to one order of higher estimation error. For example, suppose the estimation error of BF Monte Carlo stays in a range from 0.1 to 0.9, the estimation error of iOO would be in a range from 1 to 9. Details analysis of and a formal proof of the relationship between the simulation time and simulation accuracy can also be referred to the work [42].

Selection set size of each method, BP, iOO and eOO is different. The BP method picks a fixed ratio α that is normally larger than 80%. As we mentioned earlier, the iOO and eOO are two-stage scheduling algorithms. Their selection sets are determined by the time left for the second scheduling period. For example, if each scheduling period of iOO equals to $t = 100$ seconds and the rough evaluation of all schedules takes 10 seconds. Suppose a precise evaluation of one schedule takes 3 seconds, the maximum selection set size for iOO would be 30.

Around ten thousand groups of simulated data were generated randomly based on the distribution of the workloads. We plot the histogram of the statistic analysis of these workloads as shown in Fig. 4. Six consecutive ranges of pair-wise similarity degree values are used as the x-axis and the total number of workload pairs that shows the similarity degree is plotted as the y-axis. This method classifies the similarities among all the workloads in a few categories, and identifies quantitatively the similar value threshold we use to determine

whether one pair of workload being similar or not. We set the upper third quantile of the distribution as α and the lower third quantile as β . The results are reported in Fig. 4.

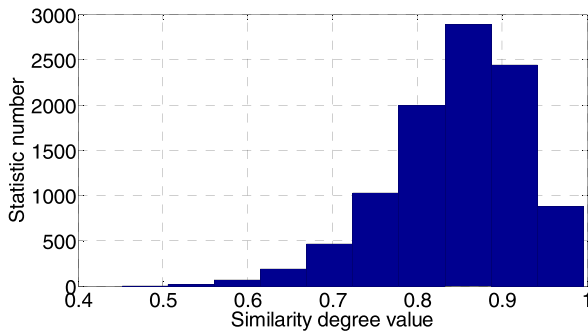


FIGURE 4. The cumulative number of workloads and the similarities between all of the workload pairs over the similarity degree values.

C. THE eOO FOR DYNAMIC WORKLOAD SCHEDULING

If the scheduling overhead is too high and the scheduling itself takes longer than the period of workload and resource changing, then the dynamic workload in the multitasking scheduling problem is difficult to capture. However, applying the OO in an evolutionary manner, leads to the full investigation of the workload patterns. For the highly dynamic and fluctuating workloads, the eOO adaptively partitions into small pieces, to generate the fine-grained schedules. Although the short-period scheduling means less accuracy, the gain in the adaptivity to such kinds of workloads, outweighs the minor inaccuracies in the solutions.

On the contrary, the eOO produces course-grained schedules for stable workloads. In this way, the merged stages lead to more simulation time for a subsequent stage that improves the accuracy. We illustrate the method graphically in Fig. 5.

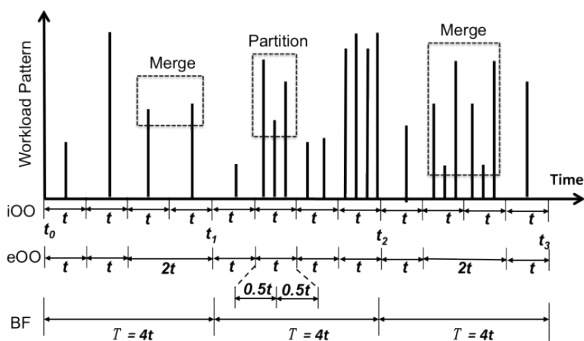


FIGURE 5. Illustration of OO adaption to dynamic workloads.

Let T be the time overhead of scheduling, using the BF or Monte Carlo, and let t be the time overhead of that of the iOO. For example, if at time t_0 , the BF is used for simulations, then it is not until t_1 that the BF can generate the optimal schedule for $[t_1, t_2]$. While the solution is optimal at time t_1 , no further and consecutive solutions can be generated again during t_1 and t_2 . As for the iOO at time t_1 , workload is used to generate an acceptable schedule at time $t_1 + t$, and then

$t_1 + 2t, \dots$. The aforementioned process is conducted repeatedly to capture a finer granularity variation of the workload to improve the overall performance.

As an improvement of iOO, the eOO takes the workload pattern into account. In $[t_1 + t, t_1 + 2t]$, workload patterns evolve significantly. In eOO, this interval is partitioned into two parts shown in the solid-rectangular box of Fig. 5. In $[t_0 + 2t, t_0 + 4t]$ or $[t_2, t_2 + 2t]$, the inter-workload patterns are similar; and as a result two intervals are merged into one.

The eOO scheduling is performed with a finer granularity, namely with $0.5t$ per scheduling period when partitioned. In case of a merge, the eOO scheduling is performed with a coarse granularity period, namely with $2t$ per scheduling period. This adaptability makes the scheduling method more accuracy with better performance.

D. OPTIMAL COMPUTING BUDGET ALLOCATION BASED SCHEDULING

All of the scheduling methods, namely Monte Carlo, Blind Pick, iOO, eOO can be further improved by combining them with the Optimal Computing Budget Allocation (OCBA) [7].

The OCBA is a simulation-based optimization approach that allocates computing resources effectively among all of the candidate schedules. Simply speaking, rather than equally allocate simulation time among all of the schedules, the OCBA allocates more simulation time to the schedules that show better overall performances in their early simulation stages. Such an exercise removes the potentially poor schedules at an early stage, and gives more observations to the promising schedules. In our experimental evaluations, we apply the OCBA method on each of the schedule stage derived by the eOO, and compare the performances of all of the methods.

IV. EVALUATIONS, RESULTS, AND DISCUSSIONS

In this section, detailed experimental design and results are presented for the virtual cluster allocation in multitasking workload scheduling applications.

A. EXPERIMENTAL SETTINGS

The cloud experiments are carried out using 16, 32, 64, and up to 128 standard VM instances on the Amazon EC2 platform. Each instance is by default, a small one with 1.7 GB of memory, one EC2 Compute Unit (one virtual core with one EC2 Compute Unit), 160 GB of local instance storage, a 32-bit platform on deployed Linux Base with EBS boot, and a 32-bit architecture with Amazon EC2 Tools shown in Fig. 6.

The benchmarking applications used are similar to [8]. The authors report a benchmark application suite includes bit-reversals, merge-sort, matrix multiplication, a trace-driven simulation, partitioning meshes, k-nearest neighbor classification, and AdaBoost algorithm. These seven applications correspond to our seven tasks classes, which need seven virtual clusters to host and execute.

Up to 10,000 tasks were concurrently executed. Each task corresponding to one benchmark application has a variable amount of input dataset size. The total amount of data used in the 128-VM scenarios approximately exceeds 1.6 TB.

We define four typical workload types by using, say, LFLD to represent Low Frequency Low Difference (LFLD) workload. This means the workload comes slowly while most of the workload patterns are alike. Similarly, the three other workload types are LFHD, HFLD and HFHD. In the HFHD case, the average arrival rate of the workload comes up to 500 tasks/second. The high velocity and volume of data used justifies a testing-scale big data analytic in these experiments.

Even though the data amount tested here is not extremely large, we argue that the method can still be extrapolated to a larger scale of hundred TB or even PB of data due to the loosely-coupled multitasking property.

Name	TestID	Instance	AMI ID	Root	Type	Status	Security Gr	Key Pair	Net	Monitoring	Virtualizat
empty	i-1e9693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-1c18693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-1a9693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-18693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-169693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-148693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-128693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-108693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-ee9693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-ec9693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-ec9693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-e9693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-e9693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-e49693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	
empty	i-e28693	ami-8c1fcec4	abs	m1	small	running	default	fzhang_sshk	basic	paravirtual	

FIGURE 6. Snapshot of multitasking applications on Amazon EC2.

The hierarchical structure of our experiments is shown in Fig. 7. Each server notation is an Amazon EC2 instance. There are 7 VCs denoted by the dashed-rectangular boxes. There is one level-2 seed server in each of the VC to ensure at least one VM is running to process the requests from level-1 seed server. It is also used to balance the workloads into the internal VMs. This is a dedicated VM that resides in each of the VC.

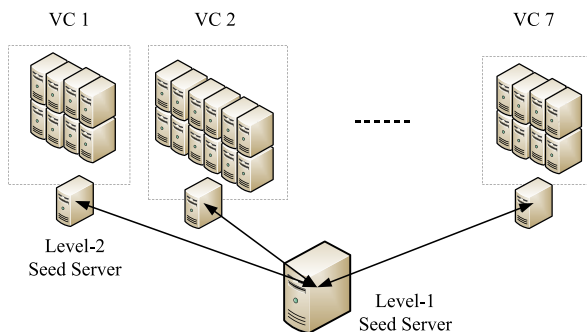


FIGURE 7. Architecture of the organized instances on Amazon EC2.

B. APPLICATION SCENARIOS

Two major application scenarios are depicted in Fig. 8 below. We describe both of them in detail.

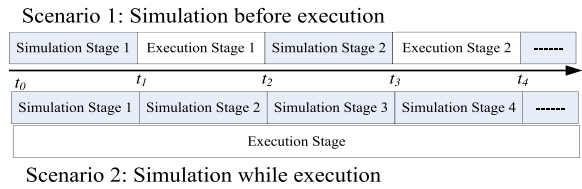


FIGURE 8. Two application scenarios for the benchmark.

Scenario 1: Simulation before execution. In the upper part of Fig. 8, one simulation stage is conducted before each experimental run. Before execution stage 1, there is sufficient time for a simulation to find a good allocation schedule from t_0 to t_1 , and the schedule is applied at t_1 . No simulation is carried out during execution stage 1 until t_2 . This scenario is commonly seen in the job scheduling of traditional parallel and distributed computing systems and spot instances in Amazon EC2. For example, in using spot instances, there is sufficient time for simulation before the price reaches an acceptable level, and execution stage can follow the simulation.

Scenario 2: Simulation while execution. In the lower part of Fig. 8, simulation and execution stages overlap. At time t_0 , a random schedule is used for the time being, while the workload information between $[t_1, t_2]$ is used to generate a better quality schedule to be used at time t_1 . Similar scheduling is conducted repeatedly until the end of experiment. Most of the real-time multitasking scheduling applications are carried out based on Scenario 2.

C. COMPARATIVE STUDIES ON SCENARIO 1

We compare the case studies of Scenario 1, on the above-mentioned four types of workloads. The results are reported in Fig. 9 (a)-(d).

It can be observed that the workload variation has a very limited impact on the three methods. The BF Monte Carlo is always the best while the Blind Pick is always the worst among all. The OO method lies in between. In Table 2, we

TABLE 2. Performance comparison of scenario 1.

Cluster Size	Methods	LFLD	LFHD	HFLD	HFHD
16	Monte Carlo	6.19%	5.50%	1.19%	0.80%
32		3.18%	3.32%	2.53%	6.48%
64		1.29%	3.51%	4.01%	2.96%
128		1.07%	1.60%	0.04%	0.13%
16	Blind Pick	-1.03%	-3.67%	-2.97%	-4.73%
32		-2.27%	-0.83%	-3.78%	-5.04%
64		-2.58%	-4.96%	-4.78%	-2.31%
128		-1.46%	-2.07%	-0.68%	-2.27%

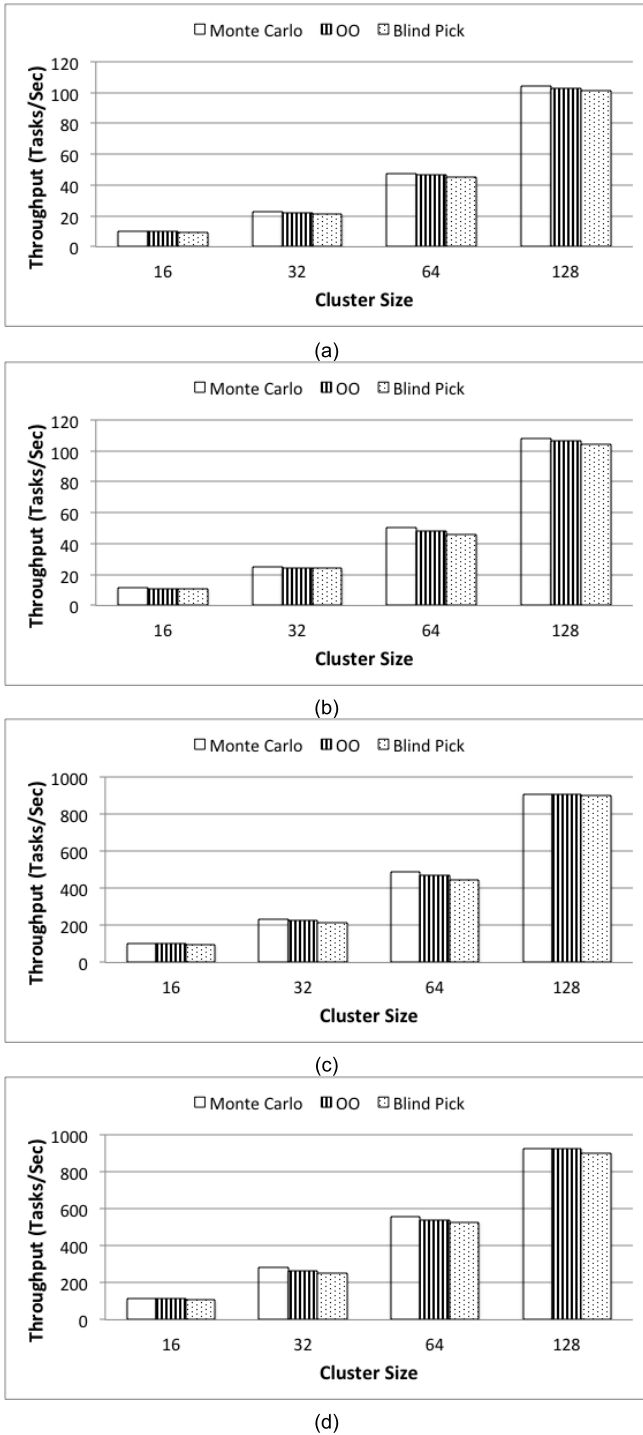


FIGURE 9. Comparative studies of application scenario 1. (a) LFLD. (b) LFHD. (c) HFLD. (d) HFHD.

demonstrate the performance difference of using OO compared with the other two methods. Positive values mean the method outperforms OO.

The aforementioned behavior is straightforward. Given sufficient simulation time, the BF Monte Carlo simulates longer than any other methods to evaluate $p_c(t_i)$, which

undoubtedly achieves the best performance because of its precise evaluation. Random selection of the Blind Pick method leads to the worst performance. This Method takes the least of the simulation time; however, the performance is mediocre at best. The OO method cannot take advantage of the long simulation time allowed and has improvement over Blind Pick, but not as good as BF Monte Carlo.

D. COMPARATIVE STUDIES ON SCENARIO 2

The simulation time budget in Scenario 2 is limited. Consequently, continuous simulation and scheduling are performed simultaneously.

In this scenario, the BF Monte Carlo method with extended simulation is no longer the best, with lower throughput compared with the iOO and eOO methods. From Fig. 10(a)-(d), we compare the throughput of the methods using a variable size of cluster. The corresponding performance comparison using eOO compared with other methods are summarized in Table 3.

TABLE 3. Performance comparison of scenario 2.

Cluster Size	Methods	LFLD	LFHD	HFLD	HFHD
16	Monte Carlo	-5.88%	-12.12%	-17.92%	-7.12%
32		-8.51%	-10.43%	-17.71%	-10.31%
64		-9.56%	-12.31%	-22.26%	-9.63%
128		-16.00%	-13.05%	-18.68%	-29.06%
16	Blind Pick	0.98%	-18.18%	-14.41%	-7.98%
32		-3.40%	-10.43%	-12.31%	-20.13%
64		-5.98%	-19.13%	-15.09%	-16.50%
128		-13.85%	-19.48%	-18.08%	-32.62%
16	iOO	-2.94%	-8.33%	-7.29%	-5.56%
32		-2.13%	-3.24%	-15.66%	-3.18%
64		-2.79%	-5.82%	-9.48%	-7.06%
128		-2.24%	-11.35%	-6.06%	-14.05%

The performance of all the four methods is similar when a small number of VMs is used. For example, the four methods perform similarly in the case of 16 VMs. However, the difference grows substantially as the number of VMs increase. In the case of 128 VMs, the eOO performs 32.62% and 29.06% better than the BP and the BF.

The relative performance of the BF Monte Carlo and BP varies given different types of workloads. In the LFHD and HFHD cases, the BP performs better; while in the LFLD and HFLD ones, the BF Monte Carlo slightly outperforms the BP. This is caused by the low overhead of BP. The BF Monte Carlo with long simulation time and high overhead, is inapplicable in Scenario 2.

The advantage of the eOO over the iOO resides in its adaptability in mutating the scheduling periods. In Fig. 10(b) and Fig. 10(d), where the workloads exhibit high difference

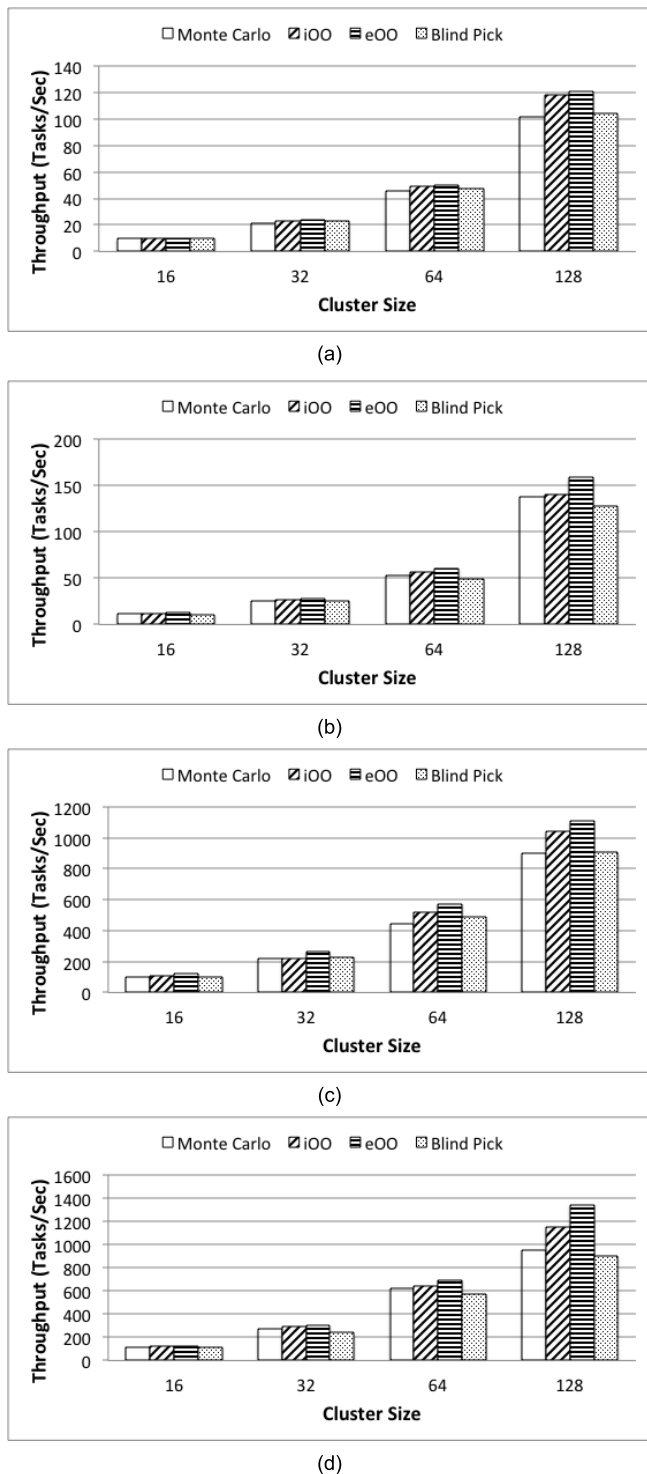


FIGURE 10. Comparative studies of application scenario 2. (a) LFLD. (b) LFHD. (c) HFLD. (d) HFHD.

over time, the eOO performs much better with its capability to partition one period into two or merge two neighboring periods into one. This is a workload pattern-aware method; therefore, it performs much better in highly dynamic workload cases.

E. SIMULATED COMPARATIVE STUDIES WITH OCBA

As mentioned previously (Section III.D), the OCBA attempts to judicially allocate the computing budget among all the multiple schedules, which breaks the equal budget allocation assumption implemented in the iOO and eOO methods. In simple words, the more a schedule shows better performance in a set of simulation runs, the more follow-up simulation runs the schedule tends to receive. By employing such a methodology, the potentially better schedules are evaluated more frequently and accurately, which translates into better system performance. Without explicitly mentioned, the OCBA used in this section is referred to using the eOO method by replacing its equal schedule simulation allocation time with the OCBA algorithm.

Unlike other methods, the OCBA calculates the simulation time for each schedule and adjusts it in runtime. In a real cloud-computing environment, this solution faces a difficulty by not being able to effectively predict the simulation cycles. Instead, we use simulations only to prove the effectiveness of this method. In our future works, we will explore by studying more scenarios to justify the applicability of this method in real cloud computing environments. The comparative results are shown in Fig. 11 by simulating a cluster of using 32 VMs. The performance metric is the makespan which aggregates all of the effective compute resource usage time, denoted by the *EET* (as detailed in Section II.A).

In Table 4, similar performance comparison of using OCBA compared with other methods is given. OCBA improves up to 29.69% performance than other methods. This improvement is observed mainly in HFHD workload. This validates our assumption again that, our methods are appealing to highly fluctuating workloads. We also demonstrate the effectiveness of using the OCBA method over all other methodologies.

TABLE 4. Performance comparison of using OCBA.

	LFLD	LFHD	HFLD	HFHD
Monte Carlo	12.64%	18.33%	16.54%	29.68%
Blind Pick	14.28%	16.53%	15.79%	29.65%
iOO	8.52%	11.11%	10.65%	23.73%
eOO	2.27%	3.63%	2.01%	7.14%

Overall, there is a 2% to 7% speedup using the eOO with OCBA compared to the eOO method only. Compared to other methods, such as the Monte Carlo and the Blind Pick in the HFHD workload scenario, the speedup can be as much as 29%. In other less dynamic workload scenarios, the speedup value varied between 12.6% and 18.3%.

V. RELATED WORK

We have witnessed an escalating interest in resource allocation for multi-task scheduling [21], [32]. Many classical optimization methods, such as opportunistic load balance, minimum execution time, and minimum completion time, are

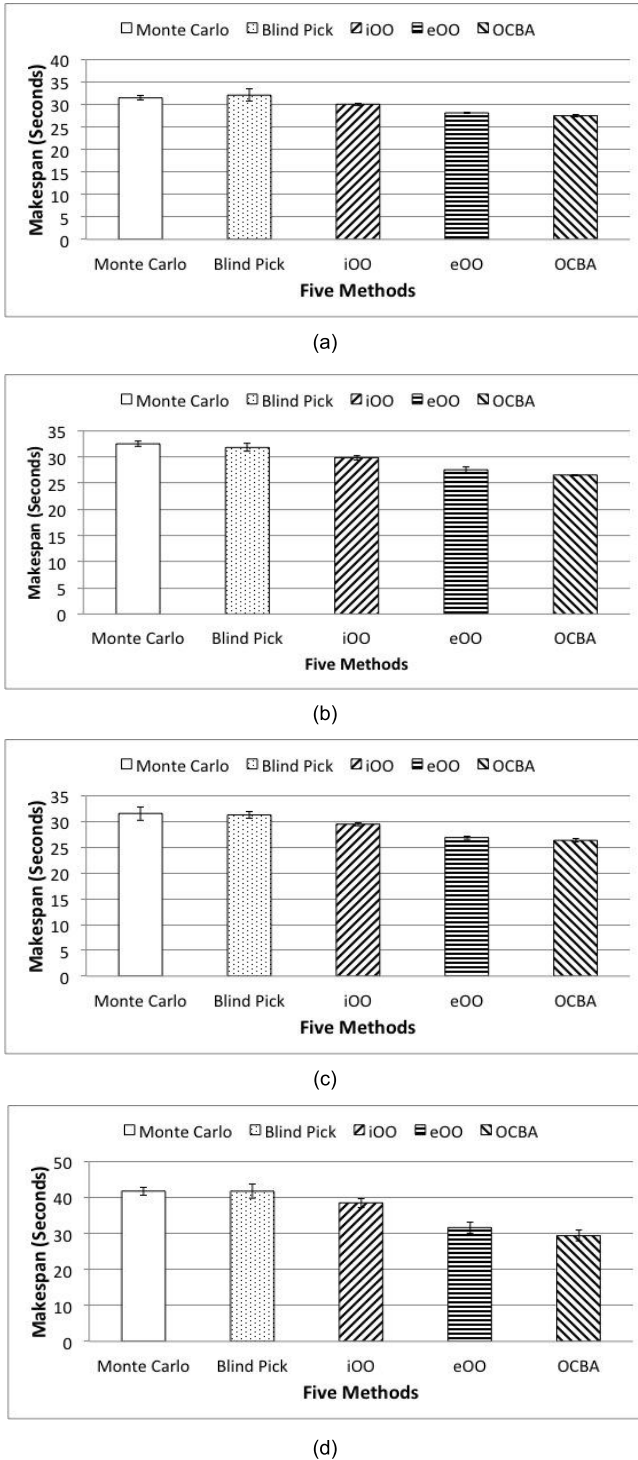


FIGURE 11. Comparative studies of application Scenario 2 with OCBA. (a) LFLD. (b) LFHD. (c) HFLD. (d) HFHD.

described in [12]. Suffrage, min-min, max-min, and auction based optimization procedures are discussed in [6] and [25]. Yu and Buyya [39] proposed economy-based methods to handle large-scale grid workflow scheduling under deadline constraints, budget allocation, and Quality of Service (QoS).

Benoit *et al.* [4] designed resource-aware allocation strategies for divisible loads. Li and Buyya [18] proposed a model-driven simulation grid scheduling strategies. Zomaya *et al.* [20], [33] have proposed a hybrid schedule and a cooperative game framework. There are also studies on making various tradeoffs under various constraints and objectives for the workflow scheduling. Wiczorek [37] analyzed five facets that may have a major impact on the selection of an appropriate scheduling strategy, and proposed taxonomies for the multi-objective workflow scheduling. Prodan *et al.* [29] proposed a novel dynamic constraint algorithm, which outperforms many existing methods, such as SOLOS and BDLS, to optimize bi-criteria problems.

Duan *et al.* [11] suggested a low complexity game-theoretic optimization method. Dogan *et al.* [10] developed a matching and scheduling algorithm for both the execution time and the failure probability. Moretti [24] suggested the All-Pairs to improve usability, performance, and efficiency of a campus grid. Smith *et al.* [30] proposed a robust static resource allocation for distributed computing systems operating under imposed QoS constraints.

Ozisikyilmaz *et al.* [26] suggested an efficient machine learning method for system space exploration. Similar to the work described in [19], our work is also carried out by using the search size based reduction. In [28], [33], and [34], petri-net and data-driven based methods are shown to compose services via the mediator, which motivates us to design the architectural cloud-scheduling platform for multitasking workloads.

A few control-based strategies are also proposed [9], [27], [45] to manipulate the VM count for delivering optimal performance. Other related work, such as [22] and [44], focus on finding optimal price strategy for renting spot instances. The similarity of these works lies to relying on their individual statistic models. Our work, however, is a simulation-based strategy, which is applied to multi-stage scheduling problems.

Based on the theory of the OO, we proposed the eOO that has advantage in handling large-scale search space to solve the multitask scheduling problem for dynamic workloads. Ever since the introduction of OO in [14], one can search for a small subset of solutions that are sufficiently good and computationally tractable. Along the OO line, many OO based heuristic methods have been proposed [17], [36]. It quickly narrows down the solution to a subset of “good enough” solutions with manageable overhead. Typical and original applications of the OO include automated manufacturing [31], communications [38], power systems [1], and distributed computing systems [43]. Different selection rules of OO are compared in [16] to discuss the relative performance. Conclusion were made that no selection rule is absolutely better than the others under all circumstances. To the best of our knowledge, this paper reports the first successful attempt to apply the OO in an iterative and evolutionary fashion for the dynamic optimization scenarios to meet special requirements of the cloud workload scheduling.

VI. CONCLUSIONS

A low-overhead method for dynamic multitasking workload scheduling in elastic virtual clusters is introduced in this work. The advantage of the proposed low-overhead scheduling scheme is summarized below.

1) Modeling of multitasking workload scheduling in wide-area elastic virtualized clusters. Cloud computing allocates virtual cluster resources on demand. This model serves the cloud computing environments and proposes simulation-based method for optimization.

2) The OO methodology is enhanced to deal with dynamically evolving workloads. iOO, with lower scheduling overhead and eOO, with adaptivity in terms of variable scheduling interval to turbulent workloads, are introduced, respectively. These two enhancements to OO are agile enough to capture system dynamism and conduct real-time scheduling. The two methods improve the overall performance in fluctuating workload scenarios evidenced by the experimental studies.

3) Proposed methods are validated by real benchmark applications with large search space, VM resource constraints, and uncertainties. Very few previous works utilizes simulation-based method for dynamic multitasking workload scheduling. Our approach is simulation based which takes consideration of these factors and is more realistic to be applied in real-world cloud systems.

Ongoing work includes the establishment of new cloud infrastructure for enabling real-time large-scale data analysis using a variety of more benchmark applications with varying workloads. We also intend to release the simulation-based optimization approaches as profiling tools to enhance performance of scheduling multitasking workload on Amazon EC2 platform.

Acknowledgment

This work was supported in part by the Ministry of Science and Technology of China through the National 973 Basic Research Program under Grant 2013CB228206 and Grant 2011CB302505, in part by the National Natural Science Foundation of China under Grant 61472200 and Grant 61233016, and in part by the National Science Foundation under Grant CCF-1016966.

REFERENCES

- [1] E. H. Allen, "Stochastic unit commitment in a deregulated electric utility industry," Ph.D. dissertation, Dept. Elect. Eng. Comput. Sci., MIT, Cambridge, MA, USA, 1998.
- [2] M. Armbrust *et al.*, "Above the clouds: A Berkeley view of cloud computing," Univ. California, Berkeley, Berkeley, CA, USA, Tech. Rep. UCB/EECS-2009-28, 2009.
- [3] P. Barham *et al.*, "Xen and the art of virtualization," in *Proc. 19th ACM Symp. Operating Syst. Principles*, Bolton Landing, NY, USA, 2003, pp. 164–177.
- [4] A. Benoit, L. Marchal, J.-F. Pineau, Y. Robert, and F. Vivien, "Resource-aware allocation strategies for divisible loads on large-scale systems," in *Proc. IEEE Int. Parallel Distrib. Process. Symp. (IPDPS)*, Rome, Italy, May 2009, pp. 1–4.
- [5] J. Cao, S. A. Jarvis, S. Saini, and G. R. Nudd, "Gridflow: Workflow management for grid computing," in *Proc. 3rd IEEE/ACM Int. Symp. Cluster Comput. Grid*, Tokyo, Japan, May 2003, pp. 198–205.
- [6] H. Casanova, A. Legrand, D. Zagorodnov, and F. Berman, "Heuristics for scheduling parameter sweep applications in grid environments," in *Proc. 9th Heterogenous Comput. Workshop (HCW)*, Cancún, Mexico, May 2000, pp. 349–363.
- [7] C.-H. Chen, "An effective approach to smartly allocate computing budget for discrete event simulation," in *Proc. 34th IEEE Conf. Decision Control*, Dec. 1995, pp. 2598–2605.
- [8] S. Chen, L. Xiao, and X. Zhang, "Adaptive and virtual reconfigurations for effective dynamic job scheduling in cluster systems," in *Proc. 22nd Int. Conf. Distrib. Comput. Syst. (ICDCS)*, Vienna, Austria, 2002, pp. 35–42.
- [9] A. Demberel, J. Chase, and S. Babu, "Reflective control for an elastic cloud application: An automated experiment workbench," in *Proc. USENIX Workshop Hot Topics Cloud Comput. (HotCloud)*, San Diego, CA, USA, 2009, p. 8.
- [10] A. Doğan and F. Özgüner, "Biobjective scheduling algorithms for execution time–reliability trade-off in heterogeneous computing systems," *Comput. J.*, vol. 48, no. 3, pp. 300–314, 2005.
- [11] R. Duan, R. Prodan, and T. Fahringer, "Performance and cost optimization for multiple large-scale grid workflow applications," in *Proc. IEEE/ACM Int. Conf. Supercomput. (SC)*, Reno, NV, USA, 2007, p. 12.
- [12] R. F. Freund *et al.*, "Scheduling resources in multi-user, heterogeneous, computing environments with smartnet," in *Proc. 7th Heterogenous Comput. Workshop (HCW)*, Washington, DC, USA, Mar. 1998, pp. 184–199.
- [13] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. San Francisco, CA, USA: Freeman, 1979.
- [14] Y.-C. Ho, R. Sreenivas, and P. Vakkili, "Ordinal optimization of discrete event dynamic systems," *J. Discrete Event Dyn. Syst.*, vol. 2, no. 2, pp. 61–88, 1992.
- [15] Y.-C. Ho, Q.-C. Zhao, and Q.-S. Jia, *Ordinal Optimization, Soft Optimization for Hard Problems*. New York, NY, USA: Springer-Verlag, 2007.
- [16] Q.-S. Jia, Y.-C. Ho, and Q.-C. Zhao, "Comparison of selection rules for ordinal optimization," *Math. Comput. Model.*, vol. 43, nos. 9–10, pp. 1150–1171, 2006.
- [17] D. Li, L. H. Lee, and Y.-C. Ho, "Constraint ordinal optimization," *Inf. Sci.*, vol. 148, nos. 1–4, pp. 201–220, 2002.
- [18] H. Li and R. Buyya, "Model-driven simulation of grid scheduling strategies," in *Proc. 3rd IEEE Int. Conf. e-Sci. Grid Comput.*, Dec. 2007, pp. 287–294.
- [19] D. Lu, H. Sheng, and P. Dinda, "Size-based scheduling policies with inaccurate scheduling information," in *Proc. IEEE Comput. Soc. 12th Annu. Int. Symp. Model., Anal., Simul. Comput. Telecommun. Syst.*, Oct. 2004, pp. 31–38.
- [20] K. Lu and A. Y. Zomaya, "A hybrid schedule for job scheduling and load balancing in heterogeneous computational grids," in *Proc. 6th IEEE Int. Parallel Distrib. Process. Symp.*, Hagenberg, Austria, Jul. 2007, pp. 121–128.
- [21] M. Maheswaran, S. Ali, H. J. Siegel, D. Haensgen, and R. F. Freund, "Dynamic mapping of a class of independent tasks onto heterogeneous computing systems," *J. Parallel Distrib. Comput.*, vol. 59, no. 2, pp. 107–131, Feb. 1999.
- [22] M. Mattess, C. Vecchiola, and R. Buyya, "Managing peak loads by leasing cloud infrastructure services from a spot market," in *Proc. 12th IEEE Int. Conf. High Perform. Comput. Commun. (HPC)*, Melbourne, Australia, Sep. 2010, pp. 180–188.
- [23] N. Metropolis and S. Ulam, "The Monte Carlo method," *J. Amer. Statist. Assoc.*, vol. 44, no. 247, pp. 335–341, 1949.
- [24] C. Moretti, H. Bui, K. Hollingsworth, B. Rich, P. Flynn, and D. Thain, "All-pairs: An abstraction for data-intensive computing on campus grids," *IEEE Trans. Parallel Distrib. Syst.*, vol. 21, no. 1, pp. 33–46, Jan. 2010.
- [25] A. Mutz and R. Wolski, "Efficient auction-based grid reservations using dynamic programming," in *Proc. IEEE/ACM Int. Conf. Supercomput.*, Austin, TX, USA, Nov. 2007, p. 16.
- [26] B. Ozisikyilmaz, G. Memik, and A. Choudhary, "Efficient system design space exploration using machine learning techniques," in *Proc. 45th ACM/IEEE Design Autom. Conf.*, Jun. 2008, pp. 966–969.
- [27] P. Padala *et al.*, "Adaptive control of virtualized resources in utility computing environments," in *Proc. 2nd ACM SIGOPS/EuroSys Eur. Conf. Comput. Syst. (EuroSys)*, Lisbon, Portugal, 2007, pp. 289–302.
- [28] S. Pal, S. K. Das, and M. Chatterjee, "User-satisfaction based differentiated services for wireless data networks," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Seoul, Korea, May 2005, pp. 1174–1178.
- [29] R. Prodan and M. Wiczorek, "Bi-criteria scheduling of scientific grid workflows," *IEEE Trans. Autom. Sci. Eng.*, vol. 7, no. 2, pp. 364–376, Apr. 2010.
- [30] J. Smith, H. J. Siegel, and A. A. Maciejewski, "A stochastic model for robust resource allocation in heterogeneous parallel and distributed computing systems," in *Proc. IEEE Int. Parallel Distrib. Process. Symp. (IPDPS)*, Miami, FL, USA, Apr. 2008, pp. 1–5.
- [31] C. Song, X. Guan, Q. Zhao, and Y.-C. Ho, "Machine learning approach for determining feasible plans of a remanufacturing system," *IEEE Trans. Autom. Sci. Eng.*, vol. 2, no. 3, pp. 262–275, Jul. 2005.

- [32] D. P. Spooner, J. Cao, S. A. Jarvis, L. He, and G. R. Nudd, "Performance-aware workflow management for grid computing," *Computer J.*, vol. 48, no. 3, pp. 347–357, 2005.
- [33] R. Subrata, A. Y. Zomaya, and B. Landfeldt, "A cooperative game framework for QoS guided job allocation schemes in grids," *IEEE Trans. Comput.*, vol. 57, no. 10, pp. 1413–1422, Oct. 2008.
- [34] W. Tan, Y. Fan, M. Zhou, and Z. Tian, "Data-driven service composition in enterprise SOA solutions: A Petri net approach," *IEEE Trans. Autom. Sci. Eng.*, vol. 7, no. 3, pp. 686–694, Jul. 2010.
- [35] W. Tan, Y. Fan, and M. Zhou, "A Petri net-based method for compatibility analysis and composition of web services in business process execution language," *IEEE Trans. Autom. Sci. Eng.*, vol. 6, no. 1, pp. 94–106, Jan. 2009.
- [36] S. Teng, L. H. Lee, and E. P. Chew, "Multi-objective ordinal optimization for simulation optimization problems," *Automatica*, vol. 43, no. 11, pp. 1884–1895, 2007.
- [37] M. Wiecek, R. Prodan, and A. Hoheisel, "Taxonomies of the multi-criteria gridworkflow scheduling problem," CoreGRID, Barcelona, Spain, Tech. Rep. TR6-0106, 2007.
- [38] J. E. Wieselthier, C. M. Barnhart, and A. Ephremides, "Standard clock simulation and ordinal optimization applied to admission control in integrated communication networks," *Discrete Event Dyn. Syst.*, vol. 5, nos. 2–3, pp. 243–280, 1995.
- [39] J. Yu and R. Buyya, "Scheduling scientific workflow applications with deadline and budget constraints using genetic algorithms," *Sci. Programming*, vol. 14, nos. 3–4, pp. 217–230, Dec. 2006.
- [40] F. Zhang, J. Cao, K. Hwang, and C. Wu, "Ordinal optimized scheduling of scientific workflows in elastic compute clouds," in *Proc. 3rd IEEE Int. Conf. Cloud Comput. Technol. Sci.*, Athens, Greece, Dec. 2011, pp. 9–17.
- [41] F. Zhang, J. Cao, H. Cai, and C. Wu, "Provisioning virtual resources adaptively in elastic compute cloud platforms," *Int. J. Web Services Res.*, vol. 8, no. 3, pp. 54–69, 2011.
- [42] F. Zhang, J. Cao, K. Li, S. U. Khan, and K. Hwang, "Multi-objective scheduling of many tasks in cloud platforms," *Future Generat. Comput. Syst.*, vol. 37, pp. 309–320, Jul. 2014.
- [43] F. Zhang, J. Cao, L. Liu, and C. Wu, "Performance improvement of distributed systems by autotuning of the configuration parameters," *Tsinghua Sci. Technol.*, vol. 16, no. 4, pp. 440–448, 2011.
- [44] Q. Zhang, Q. Zhu, and R. Boutaba, "Dynamic resource allocation for spot markets in cloud computing environments," in *Proc. 4th IEEE/ACM Int. Conf. Utility Cloud Comput. (UCC)*, Melbourne, Australia, Dec. 2011, pp. 178–185.
- [45] H. Zhao, M. Pany, X. Liu, X. Liy, and Y. Fangy, "Optimal resource rental planning for elastic applications in cloud market," in *Proc. 28th IEEE Int. Parallel Distrib. Process. Symp. (IPDPS)*, Phoenix, AZ, USA, May 2012, pp. 808–819.



FAN ZHANG (S'08–M'12–SM'13) is currently a Post-Doctoral Associate with the Kavli Institute for Astrophysics and Space Research, Massachusetts Institute of Technology, Cambridge, MA, USA. He has been a Visiting Associate Professor with the Shenzhen Institute of Advanced Technology, Chinese Academy of Science, Shenzhen, China, since 2014. He received the Ph.D. degree from the Department of Control Science and Engineering, Tsinghua University,

Beijing, China, in 2012, the B.S. degree in computer science from the Hubei University of Technology, Wuhan, China, and the M.S. degree in control science and engineering from the Huazhong University of Science and Technology, Wuhan. From 2011 to 2013, he was a Research Scientist with the Cloud Computing Laboratory, Carnegie Mellon University, Pittsburgh, PA, USA. He was a recipient of an Honorarium Research Funding Award from the University of Chicago and the Argonne National Laboratory (2013), the Meritorious Service Award from the *IEEE Transactions on Service Computing* (2013), and two IBM Ph.D. Fellowship Awards (2010 and 2011). His research interests include big-data scientific computing applications, simulation-based optimization approaches, cloud computing, and novel programming models for streaming data applications on elastic cloud platforms.



JUNWEI CAO (M'99–SM'05) received the Ph.D. degree in computer science from the University of Warwick, Coventry, U.K., in 2001, and the bachelor's and master's degrees in control theories and engineering from Tsinghua University, Beijing, China, in 1996 and 1998, respectively. He is currently a Professor and the Deputy Director of the Research Institute of Information Technology, Tsinghua University. He is also the Director of the Open Platform and Technology Division,

Tsinghua National Laboratory for Information Science and Technology, Beijing. He is an Adjunct Faculty of Electrical and Computer Engineering with North Dakota State University, Fargo, ND, USA. Before joining Tsinghua University in 2006, he was a Research Scientist with the Laser Interferometer Gravitational-Wave Observatory Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA, and NEC Laboratories Europe, Heidelberg, Germany, for about five years. He has authored over 160 papers and cited by international scholars for over 7 000 times according to Google Scholar. He has authored or edited six books or conference proceedings. He is an Associate Editor of the *IEEE TRANSACTIONS ON CLOUD COMPUTING* and an Editor of *Journal of Computational Science*. He is the General Co-Chair of the International Conference on Networking and Distributed Computing. His research is focused on distributed computing and networking, and energy/power applications. He is a Senior Member of the IEEE Computer Society, and a member of the Association for Computing Machinery and the Civil Contractors Federation.



WEI TAN (M'12–SM'13) received the B.S. and Ph.D. degrees from the Department of Automation, Tsinghua University, Beijing, China, in 2002 and 2008, respectively. He is currently a Research Staff Member with the IBM T. J. Watson Research Center, Yorktown Heights, NY, USA. From 2008 to 2010, he was a researcher with the Computation Institute, University of Chicago, Chicago, IL, USA, and the Argonne National Laboratory, Lemont, IL, USA. At that time, he was the technical lead of the caBIG workflow system. His research interests include

NoSQL, big data, cloud computing, service-oriented architecture, business and scientific workflows, and Petri nets. He has authored over 50 journal and conference papers, and a monograph entitled *Business and Scientific Workflows: A Web Service-Oriented Approach* (272 pages, Wiley-IEEE Press). He is an Associate Editor of *IEEE Transactions on Automation Science and Engineering*. He served on the Program Committee of many conferences and co-chaired several workshops. He was a recipient of the Best Paper Award from the IEEE International Conference on Services Computing (2011), the Pacesetter Award from the Argonne National Laboratory (2010), and the caBIG Teamwork Award from the National Institute of Health (2008). He is a member of the Association for Computing Machinery.



SAMEE U. KHAN (S'02–M'07–SM'12) is currently an Associate Professor with North Dakota State University, Fargo, ND, USA. He received the Ph.D. degree from the University of Texas at Arlington, Arlington, TX, USA, in 2007. His research interests include optimization, robustness, and security of: cloud, grid, cluster, and big data computing, social networks, wired and wireless networks, power systems, smart grids, and optical networks. His work has appeared in over 225 publications with two receiving best paper awards. He is a fellow of the Institution of Engineering and Technology and the British Computer Society.

with two receiving best paper awards. He is a fellow of the Institution of Engineering and Technology and the British Computer Society.



actions on Cloud Computing and the *IEEE Transactions on Computers*.

KEQIN LI (M'90–SM'96) is currently a SUNY Distinguished Professor of Computer Science and an Intellectual Ventures Endowed Visiting Chair Professor with Tsinghua University, Beijing, China. His research interests are mainly in design and analysis of algorithms, parallel and distributed computing, and computer networking. He has over 300 research publications, and has received several best paper awards for his research work. He is currently on the Editorial Boards of *IEEE Transactions on Cloud Computing* and the *IEEE Transactions on Computers*.



the *IEEE Transactions on Computers* and Springer's *Scalable Computing*, and serves as an Associate Editor of 22 leading journals, such as, the *ACM Computing Surveys* and *Journal of Parallel and Distributed Computing*.

ALBERT Y. ZOMAYA (M'90–SM'97–F'04) is currently the Chair Professor of High Performance Computing and Networking with the School of Information Technologies, University of Sydney, Sydney, NSW, Australia, where he is also the Director of the Centre for Distributed and High Performance Computing, which was established in 2009. He has authored over 500 scientific papers and articles, and has authored, co-authored, or edited over 20 books. He is an Editor-in-Chief of the *IEEE Transactions on Computers* and Springer's *Scalable Computing*, and serves as an Associate Editor of 22 leading journals, such as, the *ACM Computing Surveys* and *Journal of Parallel and Distributed Computing*. Prof. Zomaya was a recipient of the IEEE Technical Committee on Parallel Processing Outstanding Service Award (2011), the IEEE Technical Committee on Scalable Computing Medal for Excellence in Scalable Computing (2011), and the IEEE Computer Society Technical Achievement Award (2014). He is a Chartered Engineer, and a fellow of the American Association for the Advancement of Science and the Institution of Engineering and Technology (U.K.). His research interests are in the areas of parallel and distributed computing and complex systems.