



# A Ranking Chaos Algorithm for dual scheduling of cloud service and computing resource in private cloud



Yuanjun Laili<sup>a</sup>, Fei Tao<sup>a</sup>, Lin Zhang<sup>a,\*</sup>, Ying Cheng<sup>a</sup>, Yongliang Luo<sup>a</sup>, Bhaba R. Sarker<sup>b</sup>

<sup>a</sup> School of Automation Science and Electrical Engineering, Beihang University, Beijing 100191, China

<sup>b</sup> Department of Mechanical & Industrial Engineering, Louisiana State University, Baton Rouge, LA 70803, USA

## ARTICLE INFO

### Article history:

Received 20 September 2012

Received in revised form 22 January 2013

Accepted 15 February 2013

Available online 16 March 2013

### Keywords:

Private cloud

Service Composition Optimal Selection (SCOS)

Optimal Allocation of Computing Resources (OACR)

Dual-scheduling

Chaos optimization

## ABSTRACT

Private cloud as an important branch of cloud computing has brought significant benefit to many kinds of conglomerates in resource sharing. With central management of centre console, Service Composition Optimal Selection (SCOS) and Optimal Allocation of Computing Resources (OACR) are two critical steps for implementing high flexible and agile service provision and resource sharing among sub-enterprises and partner-enterprises under the key technologies of virtualization. However, two steps decision-making are inefficient and cumbersome. To overcome this deficiency, the idea of combining SCOS and OACR into one-time decision in one console is first presented in this paper, named Dual Scheduling of Cloud Services and Computing Resources (DS-CSCR). The mutual relations between the upper layer cloud services and the underlying infrastructures and their properties in the private cloud of conglomerate are deeply analyzed. For addressing large-scale DS-CSCR problem, a new Ranking Chaos Optimization (RCO) is proposed. With the consideration of large-scale irregular solution spaces, new adaptive chaos operator is designed to traverse wider spaces within a short time. Besides, dynamic heuristic and ranking selection are introduced to control the chaos evolution in the proposed algorithm. Theoretical analysis and simulations demonstrate that the new DS-CSCR outperforms the traditional two-level decision making with the improvements in both cloud service composition and computing resource allocation. In addition, RCO can remarkably give much prominent solutions with low time-consuming and high stability than a few typical intelligent algorithms for solving DS-CSCR in private cloud. With the new DS-CSCR and RCO, cloud services and computing infrastructures can then be quickly combined and shared with high efficient decision.

© 2013 Elsevier B.V. All rights reserved.

## 1. Introduction

Newly developing cloud computing [1,2] has brought about great benefits to both enterprises and individuals. With advanced technologies of virtualization and service, it incorporates various resources for user on-demand with open interfaces and transparent remote operations. While IBM, Google and Amazon are taking the lead in building general public cloud [3–5] under the modes of SaaS (Software as a Services), IaaS (Infrastructure as a Service) and PaaS (Platform as a Service) [6], many conglomerates have also obtained cost reduction and higher flexibility of resource sharing with the establishment of their own private cloud.

Private cloud of conglomerate usually consists of a set of virtualized distributed infrastructures and application services which are provided by couples of sub-enterprises and partner-enterprises [7,8], as shown in Fig. 1. For outside, such conglomerate

could be a large SaaS provider. For inside, it turns to a shared resource pool. In a fairly secure environment, all resources are under the ownership and control of a single administrative domain. On one hand, the virtualization of multiple distributed infrastructures can greatly improve the computing capability for the whole organization with lower-cost. On the other hand, upper layer application services, no matter provided to outside Internet or inside members, need no longer to be deployed on a fixed computing resource with specific maintenance. Services with central control become more flexible with dynamic allocation. Thus, private cloud in conglomerate also contains two aspects of significance, one is the integration and sharing of underlying distributed infrastructure, another is the flexible deployment and usage of upper layer application services.

Besides, with the development of cloud, the concept of “service” in traditional Service-Oriented Architecture (SOA) is extended from software application to generalized “cloud service” with the inclusion of both software applications and hardware equipments with good interoperability, self-organization and scalability [9]. The properties of cloud services have

\* Corresponding author. Tel.: +86 01082339211.  
E-mail address: [johnlin9999@163.com](mailto:johnlin9999@163.com) (L. Zhang).

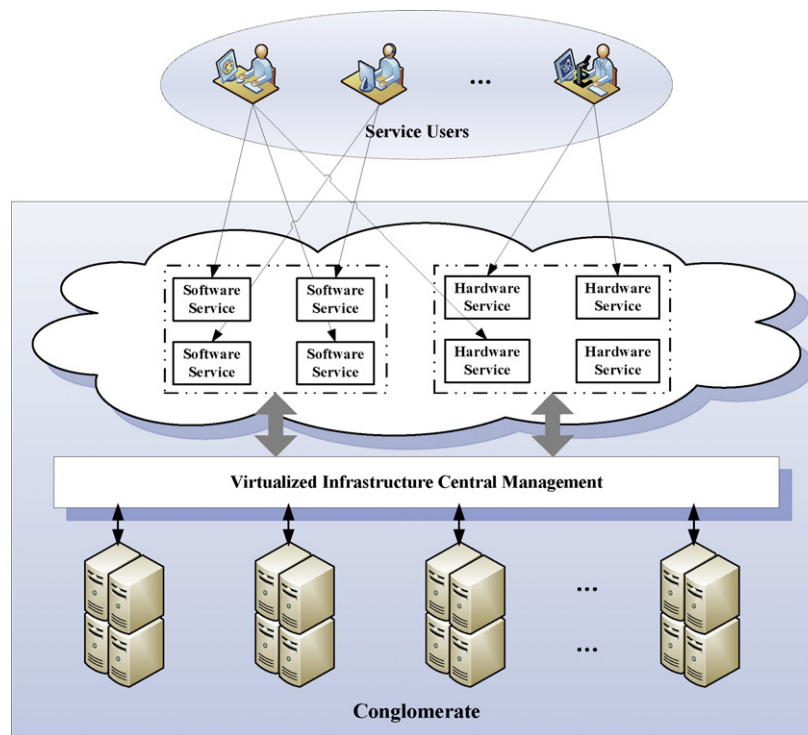


Fig. 1. Structure of private cloud and actors in conglomerate.

become more complex and most of them need higher computing ability to drive.

In such environment, when a composite project (which contains a set of tasks) is submitted, the console of conglomerate needs not only to aggregate suitable cloud services with different functionalities and generate service portfolio for user on-demand, but also choose available computing resources to support the running of cloud services. How to achieve high-quality and low-cost Services Composition Optimal Selection (SCOS) and Optimal Allocation of Computing Resources (OACR) simultaneously are critical for efficient project execution, green resource sharing and flexible service management.

At present, service composition and computing resource allocation in cloud have been studied preliminarily. Most researches are carried out according to the methodology of cluster computing, grid computing and high performance computing and consider the two problems independently. For one thing, computing availability and communication route of computing resources are analyzed. For another, QoS (Quality of Service) indexes and description languages are also discussed. In general public cloud, SCOS and OACR are performed in two steps and in the charge of different actors. Service providers are not infrastructure providers [2]. However, in private cloud of conglomerates with typical SaaS mode, they would provide suitable service portfolio and deploy corresponding services on their own infrastructure for customers on demand. The actors of SCOS and OACR turn out to be the same one.

With such two-step decision by a single administrator, the properties of upper layer selected cloud services in SCOS will limit the range of the underlying available computing resources for each service in OACR. Better portfolios of cloud services and computing resources are easily overlooked. Furthermore, as all knows, both SCOS and OACR are proved to be NP combinatorial optimization problem. Under the condition of large-scale cloud services and computing resources and complex relationship between them, addressing SCOS and OACR step by step with two different algorithms independently becomes very cumbersome and inefficient.

Therefore, we propose the idea of combining two stages decision-making into one and put forward the concept, Dual Scheduling of Cloud Service and Computing Resource (DS-CSCR), in private cloud of conglomerate. In the guidance of this idea, we analyze the complex features of hardware/software cloud service and computing resource in cloud computing in two levels and explore their mutual relations in-depth. Aiming at green efficient decision, the formulation of DS-CSCR with multi-objectives and multi-constraints is presented in this paper. Additionally, in order to achieve high efficient one-time decision in DS-CSCR, a new Ranking Chaos Optimization (RCO) is designed in this paper. Take the advantage of chaotic random ergodicity, this algorithm combines new adaptive chaos optimal strategy with ranking selection and dynamic heuristic mechanism to balance the exploration and exploitation in optimization. With adaptive control of chaotic sequence length, it is especially good at searching in large-scaled irregular solution space and shows remarkable performance for addressing DS-CSCR compared with other general intelligent algorithms.

The structure of this paper is as follows. Section 2 gives the related works of cloud computing, private cloud and the researches of SCOS and OACR in cloud. Section 3 introduces part machining project as an example to illustrate the existing problem of SCOS and OACR. Section 4 elaborates the main attribute indexes of tasks, cloud services, virtual machines and computing resources and establishes the formulation of DS-CSCR. Section 5 proposes a new RCO for DS-CSCR and describes its procedures in detail. The performance evaluation of our algorithm and comparisons with other intelligent algorithms in solving DS-CSCR by experiments are given in Section 6. Finally we conclude this paper with future works in Section 7.

## 2. Related works

Nowadays, the most commonly used and analyzed cloud computing platforms are “Google cloud computing” platform, Amazon “elastic cloud” platform and IBM “blue cloud” platform.

Private cloud with closed sharing are researched less and attracted criticism owing to the less hands-on management [3]. But it can notably reduce the cost of resources and improve the quality of services in large conglomerate. After years of development, large enterprises, academic institutions and new emerging internet service providers are building their own cloud platform, too, such as Eucalyptus, Red Hat's cloud, OpenNebula and so on. Though various platforms differ on their usage mode and openness, most of them share the same key technologies and target of resource sharing.

In cloud computing, two crucial optimization factors in determining resource sharing efficiency and platform application performance are SCOS and OACR exactly.

In recent years, researches on service composition are generally based on the environment of grid computing and other SOA mode [10]. These researches spread from service description language, service QoS indexes [11], reliability and trust evaluation [12], and optimal selection of services [13] and so on. Since cloud computing mode has been proposed, the concept and content of cloud service are broadened. According to the characteristics of cloud computing, semantic properties of cloud service are studied [14]. The classification, management, provision, storage and evaluation of cloud services are investigated widely. Pre-decision and online-decision of SCOS are also deliberated in different ways, such as [15]. Among these, QoS indexes of cloud services are discussed most widely. From the perspective of non-functional properties of cloud services, the existing indexes consider no more than cost, time and reliability factors. It is hard to describe various cloud services with different classification and attributes in a unified form. Thus the existing QoS indexes cannot satisfy all types of cloud services.

For computing resource allocation, traditional researches mostly focus on the modeling and evaluation of computing resources based on homogeneous/heterogeneous cluster systems or distributed grid computing systems [16]. User's demand for resources, the cost and computation and communication capabilities of resources are the major considerations among these studies. In cloud computing mode, virtualization is the main support of flexible resource sharing [17]. In this context, Endo et al. introduced the concept, classification of resource allocation in distributed cloud [18], Ma et al. [19] and Xiong et al. [20] investigated the management of cloud computing resources based on ontology and virtualization respectively. Zhang et al. [21] proposed a method for the deployment of upper layer software cloud services from virtual machines. Ghanbari et al. [22] have studied the feedback-based optimization problem including the allocation of resources especially in private cloud. Besides, considering the virtual division of computing resources and its influences on the quality of cloud services, researchers also built new models for computing resources from the rules, reliability and dynamic partition point of view, and so on, and presented various methods to solve OACR problem in cloud computing [23,24]. Most of these studies concentrated on the expansion of characteristics of computing resources based on traditional models and the algorithm designing for OACR in cloud computing. However, the mutual relations between cloud services and the underlying computing resources and the influence of virtualization on quality of cloud services, as two of the key factors in cloud computing, have not been studied.

In addition, SCOS and OACR are both combinatorial optimization problems. For this kind of problems, the most widely used algorithms are intelligent algorithms due to its NP complexity. It includes Genetic Algorithm (GA) [25], Particle Swarm Optimization (PSO) [26] and so on and has the virtues of brachylogy, universality and rapidity. According to different specific problems, abundant researches mainly focus on the balance of exploration

and exploitation in searching process based on evolutionary iteration of population and presented many kinds of improved hybrid intelligent algorithms such as [27]. Nevertheless, these improved hybrid intelligent algorithms are mostly problem-dependent with local convergence more or less. For addressing large-scaled DS-CSCR problem in private cloud of large conglomerate with irregular solution space efficiently, the design of high performance intelligent algorithm is imperative.

### 3. Motivation example

Currently, the concept of cloud is studied and applied in almost every field. Based on the technology of cloud computing, manufacturing equipments and simulation software as cloud services can be realized [28,29]. Various software and hardware can be dynamically shared for product customization of both inside or outside organizations without repeat-purchase. Under this background, we use "the design and NC (Numerical Control) machining of a complex surface part in conglomerate cloud" as a case to describe the whole process from tasks' submission to tasks' execution. As shown in Fig. 2, it can be divided into five sub-tasks: (1) technical and mathematical analysis, (2) CAD modeling and NC programming, (3) verification simulation and post-processing, (4) first NC machining and measuring, and (5) batch production.

During this process, task (1)–(3) can be implemented directly by manufacturing software cloud services, such as CATIA, MasterCAM or Pro/E, etc., and task 4 and 5 can be executed by manufacturing hardware cloud service with users' supervision and control, such as 3-axis, 4-axis or 5-axis linkage CNC (Computer Numerical Control) machines, etc. When user submitted the tasks of designing and machining a customized part, four steps are needed to be done by centre console: (1) Requirement analysis of tasks, (2) Services Composition Optimal Selection, (3) Optimal Allocation of Computing Resources, and (4) Execution.

The scheduling of computing resources totally depends on the corresponding upper layer selected cloud services. With the distributed characteristics of services and infrastructures, the available computing resources are reduced and the OACR are constrained by the upper layer decision. For example, for task (4), assume the suitable CNC hardware service Nos. 1 and 2 are provided in Location A and Location B respectively. CNC service No. 1 is with higher QoS than CNC service No. 2. But the idle computing resources in Location A are less than Location B. If CNC service No. 1 is selected for task 4 in step 3, the low computing ability of computing resource in Location A and the remote communication overhead of computing resource in Location B would both cause the low execution efficiency of CNC service No. 1. If we select CNC service No. 2, the better available adjacent computing resource would then improve the overall execution efficiency of task 4. However, the decision of SCOS in step 2 usually disregards the influence of the underlying support computing resources due to the traditional binding mode of service and infrastructure. The latter strategy of choosing MasterCAM service No. 2 is then overlooked. At this time, you might say, if SCOS and OACR are performed at the same time, then bad decision would not be happened.

Therefore, in order to reduce the time and improve the quality of decision, we merge SCOS and OACR into one dual-scheduling decision. With the purpose of efficient DS-CSCR decision, the following three issues are needed to be studied.

- (1) QoS indexes of software/hardware cloud services and computing resources respectively and the mutual relation between them;
- (2) The problem formulation of DS-CSCR with multi-objectives and multi-constraints in private cloud;

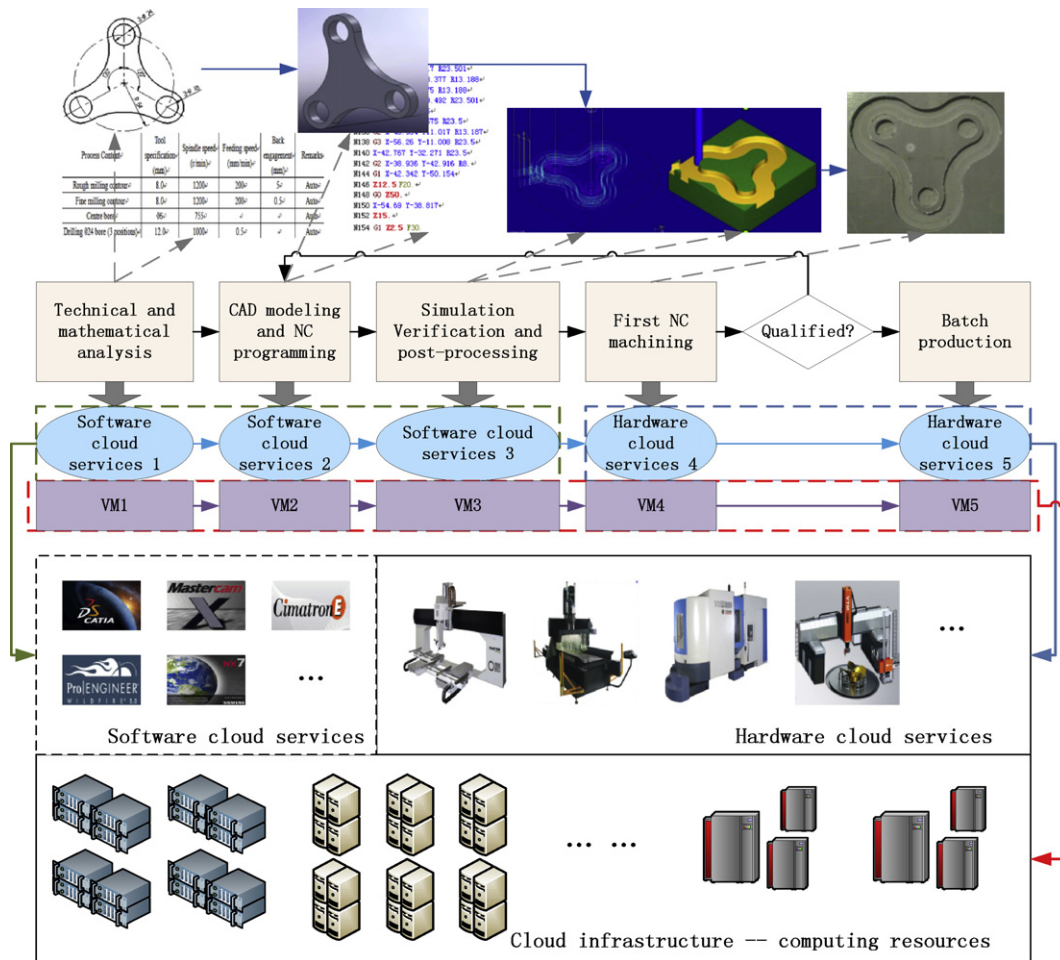


Fig. 2. The design and NC machining of a customized part in conglomerate cloud.

(3) The efficient scheduling algorithm for addressing large-scale DS-CSCR problem.

This paper will directly focus on these three issues.

#### 4. The modeling of DS-CSCR in private cloud

In conglomerate, services and the support infrastructures are provided by distributed sub-enterprises and controlled by central head. Traditionally, service provider usually deploy the service to a fixed computer, put service and computing resource together to ensure the quality of service. The support computing resources are always occupied by fixed service and needed specific maintenance. With new cloud mode, services can be encapsulated and registered to cloud and deployed to virtual machines dynamically. Through the collaborative development of upper layer applications and underlying resources, all of the resources can be shared flexibly on-demand with more energy-saving, higher redundancy and reliability.

Moreover, based on such a flexible environment, cloud services with the support of VMs contain not only software cloud services, but also hardware services with further expansion. For hardware cloud services, the computing resources are no longer support carriers, but controlling and monitoring facilities for these manufacturing equipments.

##### 4.1. The characteristics and QoS indexes of cloud services

From the perspective of QoS evaluation, only simplified quantitative cost, time and reliability cannot comprehensively summarize the characteristics of software/hardware cloud services

and their requirement for VMs' performance. With the consideration of the difference between software and hardware cloud services and their demands for VM configuration, this section gives new evaluation indexes for software/hardware cloud service and virtual resources respectively.

##### 4.1.1. The characteristics and QoS indexes of software cloud services

Software applications in cloud computing are running with the support of VMs. Each software service is deployed to a single VM and mapped to a corresponding computing resource. Thus the minimum requirements of VM which represents the required volume of services should be defined to facilitate the allocation of computing resources. Based on the functional description of services, we consider mainly the following non-functional factors of software cloud services in this paper.

- $s$  – service execution efficiency under the minimum required configuration of VM;
- $c$  – the rent cost of service;
- $r$  – trustiness of service, which is the ratio of the success execution time and the total execution time;
- $v$  – the minimum required speed of VM.

**Remarks:** The performance of the required VM is determined by many factors, such as the CPU and memory of the corresponding computing resources. In a computer, the speed of CPU is in proportion to the power supply voltage [30]. It is a constant value. The speed of VM can mainly be calculated by the number and speed of occupied CPUs. So that the minimum required speed of



VM is adopted here for evaluation. The higher the speed of VM is, the faster the service runs.

#### 4.1.2. The characteristics and QoS indexes of hardware cloud services

Unlike the software service, hardware service is energy-consuming and needs supervision or control during execution. Real-time supervision or control will produce large amount of communication and increase service execution time (i.e. the time-consumption of data transmission). Different hardware service needs different amount of supervision and control. For this reason, based on the above four factors of software service, two more factors need to be considered.

- $s$  – service execution efficiency under the minimum required configuration of VM;
- $c$  – the rent cost of service;
- $r$  – trustiness of service, which is the ratio of the success execution time and the total execution time;
- $v$  – the minimum required speed of VM;
- $e$  – the average energy-consumption of hardware service;
- $\zeta$  – the average control rate, which is the ratio of the amount of control commands and the amount of tasks;
- $\eta$  – the transmission rate between VM (computing resources) and hardware service.

*Remarks:* For hardware services, there are two conditions of control. One is inputting all control commands beforehand, and then executing tasks without interaction. Another is controlling during execution. Owing to the large amount of task in hardware service,  $\zeta$  in the first condition can usually be ignored (i.e.  $\zeta = 0$ ). We mainly focus on the second condition. Besides, if the hardware service needs no control or supervision any more, then  $\zeta = 0$ , too.

Usually, the transmission path of the control commands of software service is “user – VM”, while which of hardware service is “user – VM – hardware service”. Without the consideration of task interactions and energy-consumption of VMs, if the amount of submitted task is  $W$ , the total execution time  $T$ , the total cost  $C$  and the total energy-consumption  $E$  of the software and hardware service can be calculated as follows respectively.

For software services,

$$T = \frac{W}{s} \quad (1)$$

$$C = Tc = \frac{cW}{s} \quad (2)$$

For hardware services,

$$T = \frac{W}{s} + \frac{W\zeta}{\eta} = W \frac{\eta + s\zeta}{s\eta} \quad (3)$$

$$C = Tc = cW \frac{\eta + s\zeta}{s\eta} \quad (4)$$

$$E = Te = eW \frac{\eta + s\zeta}{s\eta} \quad (5)$$

#### 4.2. The characteristics and QoS indexes of VMs

VMs are the virtual division of the underlying computing resources. The performance of VM are mainly embodied in the running speed, transmission rate and energy consumption of the corresponding computing resources. It is still hard to locate one VM into multiple computers by existing technologies of virtualization. Hence, we assume each VM maps into only one

physical node. In accordance with the characteristics of cloud services, we primarily concentrate on five factors below.

- $p$  – the running speed of VM, which depends on the occupancy rate and the speed of CPUs;
- $q$  – the transmission rate of VM;
- $g$  – the average energy-consumption of VM;
- $f$  – the failure probability of VM;
- $u$  – the recovery time of VM when fails.

*Remarks:*  $q$  reflects the transmission rate between the occupied physical computing resources and the objects. If the transport object and the VM are in the same local network, then evaluate the transmission rate by local bandwidth. Else, the transmission rate is evaluated with the synthetic consideration of the transport object, the central console and the VM itself. Besides, the energy function of CPU per unit time can be represented as [29]:  $P_0 = AV^2 f + Z$ . Where  $A$  and  $Z$  are constant,  $V$  is the power supply voltage and  $f$  is the dominant frequency. Thus  $g$  is in proportion to  $p$ , too. In cloud platform, the way to handle the failures of physical nodes is usually dynamic migration of VMs. So,  $u$  is no longer the recovery time of the corresponding computing resource but the dynamic migration time. Computing resources with low reliability can easily cause dramatically increase of task execution time, cost and energy consumption.

Let the task execution time in the corresponding VM without failure be  $t$ , the average task execution time of VM can be evaluated as:

$$\tilde{t} = t(1 - f) + (t + u)f = t + fu \quad (6)$$

Assume the set of the predecessor nodes of the task  $i$  to be  $\mathbf{L}_i$ , and the input communication amount from the predecessor node  $j$  to be  $U_{ij}$ , then the total communication time between the task and its predecessor nodes are:

$$U = \max_{j \in \mathbf{L}_i} \frac{U_{ij}}{q_j} \quad (7)$$

If the performance of VM can satisfy the minimum requirement of service, then the total execution time  $T$ , the total cost  $C$  and the total energy consumption  $E$  of the task can be calculated as follows.

(a) If the selected service is software cloud service, then

$$T = \frac{vW}{ps} + U + fu \quad (8)$$

$$C = Tc = \left( \frac{vW}{ps} + U + fu \right) c \quad (9)$$

$$E = Te = \left( \frac{vW}{ps} + U + fu \right) e \quad (10)$$

(b) If the selected service is hardware cloud service, then

$$T = \frac{vW(\eta + s\zeta)}{ps\zeta} + U + fu \quad (11)$$

$$C = Tc = \left( \frac{vW(\eta + s\zeta)}{ps\zeta} + U + fu \right) c \quad (12)$$

$$E = T(g + e) = \left( \frac{vW(\eta + s\zeta)}{ps\zeta} + U + fu \right) (g + e) \quad (13)$$

#### 4.3. Problem formulation of DS-CSCR in private cloud

According to the analysis of the characteristics and QoS indexes of cloud services and virtual resources, the abstract formal description of cloud services, VMs and computing resources are elaborated in this section.

**Definition 1.** The set of tasks in cloud computing environment can be presented as a directed acyclic graph (DAG)  $G = (N, W, U, H_t, H_c, H_e, H_r)$ . Where:

- The set  $\mathbf{N} = \{N_i | i = 1 : n\}$  represents tasks with serial numbers, where  $n$  is the total number of tasks.
- The set  $\mathbf{W} = \{W_i | i = 1 : n\}$  indicates the size of tasks.
- The set  $\mathbf{U} = \{U_{ij} | i = 1 : n, j = 1 : n\}$  represents the communication relationships among tasks, where  $U_{ij}$  reflects the communication from task  $N_i$  to task  $N_j$ . We should note that  $U_{ij} \neq U_{ji}$ . If there is no communication between the two tasks, then  $U_{ij} = 0$ .
- $\mathbf{H}_t = \{H_t(i) | i = 1 : n\}$ ,  $\mathbf{H}_c = \{H_c(i) | i = 1 : n\}$ ,  $\mathbf{H}_e = \{H_e(i) | i = 1 : n\}$  and  $\mathbf{H}_r = \{H_r(i) | i = 1 : n\}$  represent the lowest time, cost, energy and reliability requirements of tasks respectively.

Besides, let the predecessor tasks set of  $N_i$  be  $\mathbf{L}_i$ , and the successor tasks set be  $\mathbf{R}_i$ . The node with no predecessor task  $\mathbf{L}_i = \emptyset$  is named *source* node, and the node with no successor task  $\mathbf{R}_i = \emptyset$  is called *sink* node. All tasks strictly observe the tasks' priority rules, that is to say, a node can only be started after all output communication data of its predecessor tasks are obtained.

According to the QoS indexes analyzed in the previous sections, the general model of cloud computing can be defined as follow.

**Definition 2.** The software/hardware cloud services in cloud computing mode can be presented respectively as:

$S: \begin{cases} \text{software service : } S_1 = (s, c, r, v) \\ \text{hardware service : } S_2 = (s, c, r, v, e, \zeta) \end{cases}$   
 $\mathbf{S}_1 = \{s_1(i) | i = 1 : n_{s_1}\}$  represents the set of software cloud services, where the number of services is  $n_{s_1} = |\mathbf{S}_1|$ .  $\mathbf{S}_2 = \{s_2(i) | i = 1 : n_{s_2}\}$  represents the set of hardware cloud services, where the number of services is  $n_{s_2} = |\mathbf{S}_2|$ . Therefore the total number of cloud services is  $n_s = n_{s_1} + n_{s_2}$ . In the definition,  $s, c, r, v, e$ , and  $\zeta$  represents the execution efficiency, rent cost, reliability, the minimum required speed of VM, energy-consumption and the average control rate of cloud services respectively. All of these attributes stored according to the type and the serial number of services.

Because the performance of VM is decided by the corresponding computing resources, so this paper just define the formal description of computing resources as follow.

**Definition 3.** The computing resources in cloud computing mode can be presented as  $P = (x, \varphi, \phi, \sigma, f, \lambda)$ , where:

- $\mathbf{P} = \{P_{kl} | k = 1, 2, \dots, d, l = 1, 2, \dots, m_k\}$  indicates the computing resources with different groups and different serial number, where  $k$  is the group number of the whole set,  $l$  is the number of computing resources in each group and  $d$  is the number of groups.
- $\mathbf{x} = \{x_{kl} | k = 1, 2, \dots, d, l = 1, 2, \dots, m_k\}$  represents the speed of computing resources. It is related to the configuration characteristics of these computers.
- $\Psi = \{\varphi_{kl} | k = 1, 2, \dots, d, l = 1, 2, \dots, m_k\}$  means the bandwidths of computing resources in local networks, and  $\Phi = \{\phi_k | k = 1, 2, \dots, d\}$  be the bandwidths between the switches of various sub-infrastructure groups and cloud centre console.
- $\sigma = \{\sigma_{kl} | k = 1, 2, \dots, d, l = 1, 2, \dots, m_k\}$  represents the average energy-consumption per unit time of these computing resources. According to the analysis above,  $\sigma_{kl}$  is in proportional to  $x_{kl}$ .
- $\mathbf{f} = \{f_{kl} | k = 1, 2, \dots, d, l = 1, 2, \dots, m_k\}$  means the failure probability of each computing resource. This factor is changed after each time of task execution.

- $\mathbf{\lambda} = \{\lambda_{kl} | k = 1, 2, \dots, d, l = 1, 2, \dots, m_k\}$  represents the number of task loads in each computing resource at present. It changed during task execution. If multiple VMs map into one single computing resource, the running speed of the resource will be dramatically declined. For simplified the evaluation, we assume the VMs share the same computing resource with average division.

In the definition of computing resources, the failure recovery time is not defined. Because of the dynamic migration in cloud computing system, we assume the average dynamic migration time (i.e. the recovery time) as a constant  $u = \text{Const}$ .

For two tasks  $N_i$  and  $N_j$ , if the support VMs are  $v_i$  and  $v_j$ , and the allocated computing resources are  $P_{kl}$  and  $P_{k'l'}$ , the running speed of  $v_i$  and  $v_j$  can be expressed as  $p_i = x_{kl}/\lambda_{kl}$  and  $p_j = x_{k'l'}/\lambda_{k'l'}$ . If the selected computing resources are in the same group, i.e.  $k = k'$ , the transmission rate is  $q_{ij} = \min(\varphi_{kl}, \varphi_{k'l'})$ . If the allocated computing resources are distributed, the transmission rate can be represented as  $q_{ij} = \min(\phi_k, \phi_{k'})$ . In addition, the energy-consumption of the two VMs are  $g_i = \sigma_{kl}/\lambda_{kl}$  and  $g_j = \sigma_{k'l'}/\lambda_{k'l'}$ . And the rent cost, failure probability and recovery time of VMs are defined the same as the attributes of computing resources.

Corresponding to Fig. 3, the DS-CSCR model can be defined as a quadric-tuple  $M = (G, S, V, P)$ . Based on the above definitions, the decision of DS-CSCR can be made and evaluated with multi objectives of the lowest execution time, energy-consumption and cost and the highest reliability for tasks.

Take the serial tasks as a case, let the number of tasks be  $n$ , the type of the selected cloud service for each task  $N_i$  is  $y_i$ .  $y_i$  can be 1 or 2 which represents software and hardware cloud service respectively. So that the serial number of the selected service is  $S_{y_i}(i)$ . Assume the allocated computing resource for the support VM  $v_i$  of each task is  $P_{k_i l_i}$ . Then the overall optimal objectives and constraints can be calculated as follows.

$$\text{MAX Objective Function} = w_1 \prod_{i=1}^n R_i + \frac{w_2}{\sum_{i=1}^n T_i} + \frac{w_3}{\sum_{i=1}^n C_i} + \frac{w_4}{\sum_{i=1}^n E_i} \quad (14)$$

The variables in the objective function are calculated according to Table 1.

The main constraints of DS-CSCR are shown as following:

$$\forall i \in [1, n] \quad 0 < \rho_i < 1 \quad (15)$$

$$\forall k \in [1, g], l \in [1, m_k] \quad \sigma_{kl} \geq 0 \quad (16)$$

$$\forall i \in [1, n] \quad T_i < H_t(i), C_i < H_c(i), E_i < H_e(i), R_i < H_r(i) \quad (17)$$

The first constraint means that the occupancy rates of VMs in computing resources are no less than 0 and no more than 1, that is to say, one VM can only be allocated in one computing resource with full occupancy at most. The second constraint indicates that the load of computing resources must be no less than 0. When  $\sigma_{kl} = 0$ , the computing resource is idle. When  $0 < \sigma_{kl} < 1$ , the computing resource is not fully occupied, the running speed can be hold. However, when  $\sigma_{kl} \geq 1$ , the tasks need to be executed in queue, the running speed of computing resource will be dramatically decreased. The third constraint represents that each attributes of cloud services and computing resources must satisfy the lowest requirement of tasks.

## 5. Ranking Chaos Algorithm (RCO) for DS-CSCR in private cloud

From the above analysis it is clear that the model of DS-CSCR is more complex than the traditional SCOS and OACR. The upper layer cloud services and the underlying computing resources interact with

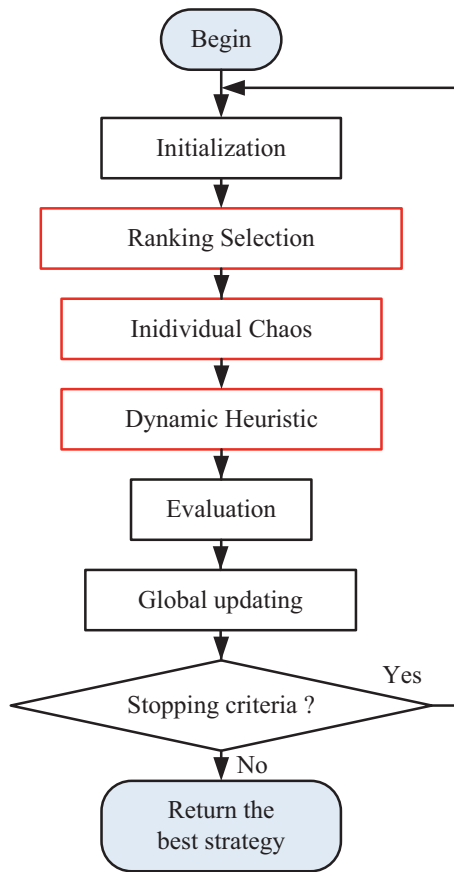


Fig. 3. The flowchart of RCO.

each other. Their complex attributes together directly determine the efficiency of task execution. In large-scale solution space, it is hard to find optimal solution of DS-CSCR by a deterministic algorithm. The general methods for solving these kinds of problems are searching for sub-optimal solutions by intelligent algorithms, such as GA, PSO and ACO and so on. ACO is designed particularly for path optimization. PSO is presented for continuous numerical optimization. GA is more universal but with serious local convergence. In the condition of complex mutual relations among the attributes of the problems with large-scaled irregular solution space, these typical algorithms are quite unsuitable.

Therefore, a new RCO is presented in this paper for DS-CSCR. The flowchart of this algorithm is shown in Fig. 3. It contains three main operators: ranking selection operator, adaptive chaos operator and dynamic heuristic operator. Their initialization (coding scheme), operators and evolutionary strategy for solving DS-CSCR are elaborated as follows.

### 5.1. Initialization

Usually, initialization in intelligent algorithm is very important. It determines the initial location and the coding scheme of

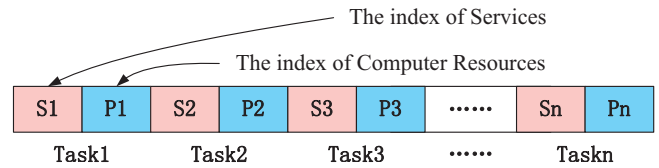


Fig. 4. The real number coding scheme for DS-CSCR problem.

population. The initial location ways of population include regular generation and random generation, and so on. For DS-CSCR, the solution space is quite complex, so that the random initialization scheme is selected in this paper.

Additionally, different coding style has different contribution to algorithm. Coding scheme in intelligent algorithm not only directly reflects the characteristics of the problems, but also affects the performance of the operators. Suitable coding scheme can even improve the searching capability of algorithms. In this paper, the real number coding scheme is adopted because of its characteristics of simplicity and intuitive.

Specifically, for the above mentioned DS-CSCR model, both service-genebit which represents the selected cloud services and resource-genebit which represents the allocated computing resources are needed to be set. One task corresponds to two genebits. Thus the real number coding is the most intuitive and space-saving scheme for DS-CSCR. When a set of tasks are submitted to cloud system, the system should choose suitable cloud services and computing resources with specific serial numbers at the same time. Assume the length of gene code be twice of the number of tasks, as shown in Fig. 4. Each two genebits represent the serial number of the selected service and the allocated computing resource for the corresponding task. It briefly demonstrates the relationship between cloud service and computing resource and makes the optimal process more convenient.

### 5.2. Ranking selection operator

In most chaos-based optimizations, chaotic operator is based on the individuals regardless of whether they are good or bad. In this case, the algorithm is easy to trap into bad conditions with large randomly searching range and extremely strong diversity. To obtain better seeds for chaotic random ergodicity, selection before it is needed.

The most commonly used selection operator in GA is roulette wheel selection. With high randomness, bad individuals may be selected more than good ones, higher diversity can be achieved in population. But high diversity has been implemented by chaos and what we need before chaos is just a set of good seeds. In this condition, roulette wheel selection becomes unsuitable. To make sure the high quality of good individuals (i.e. seeds), a dynamic ranking selection operator is designed in this section.

Normally, ranking selection means selection according to the descending sort of individual fitness values under a constant

**Table 1**  
The calculation of elements in the objective function.

Variables	Software services	Hardware services
$R_i$	$r_{s_1(i)}$	$r_{s_2(i)}$
$T_i$	$W_i \frac{v_{s_1(i)} \sigma_{k_i l_i}}{\lambda_{k_i l_i} s_{s_1(i)}} + \max_{j \in \text{pred}(i)} \frac{U_{ij}}{q_{ij}} + u f_{k_i l_i}$	$W_i \frac{v_{s_2(i)} \lambda_{k_i l_i} (\phi_{k_i} + s_{s_2(i)} \zeta_{s_2(i)})}{\lambda_{k_i l_i} s_{s_2(i)} \phi_{k_i}} + \max_{j \in \text{pred}(i)} \frac{U_{ij}}{q_{ij}} + u f_{k_i l_i}$
$C_i$	$T_i c_{s_1(i)}$	$T_i c_{s_2(i)}$
$E_i$	$T_i \frac{\sigma_{k_i l_i}}{\lambda_{k_i l_i}}$	$T_i \left( \frac{\sigma_{k_i l_i}}{\lambda_{k_i l_i}} + e_{s_2(i)} \right)$

proportion. That is to say, the numbers of individuals from best to worst are in arithmetic sequence. Here we adopt quick sort algorithm with the computation complexity  $O(n \log n)$ . Let  $\mathbf{I} = \{I_i | i = 1, 2, \dots, N\}$  be the population with  $N$  individuals, and  $I_i$  in the population be the  $i$ th individual. Assume the sorted population to be  $\mathbf{I}' = \{I'_i | i = 1, 2, \dots, N\}$  with the fitness value  $F'_N < F'_{N-1} < \dots < F'_1$ . Define  $P_{selection}$  to be the percentage of individuals to be selected on the whole. If  $P_{selection} = 1$ , then all individuals are selected at least once, if  $P_{selection} = 0.5$ , then only the first half individuals are selected, the other half individuals would not be selected any more. It represents the selection range in the sorted population. Thus the worst individual to be selected is the  $K$ th individual where  $K = NP_{selection}$ . Under the selection range, let the number of times that the best individual to be selected as  $\theta_1$  and the number of times that the worst individual to be selected as  $\theta_K$ . Then the difference between the numbers of two adjacent individuals can be calculated as follow.

$$\Delta\theta = \theta_{i-1} - \theta_i = \frac{\theta_1 - \theta_K}{K - 1} \quad (18)$$

$$\theta_i = \theta_1 - \Delta\theta(i - 1) = \theta_1 - (i - 1) \frac{\theta_K - \theta_1}{K - 1} \quad \text{where } 1 \leq i \leq K \quad (19)$$

It can be seen that  $\sum_{i=1}^K \theta_i = N$ . Therefore, we can deduce that,

$$\theta_1 + \theta_K = \frac{2N}{K} \quad (20)$$

Let  $\theta_K = 1$ , then

$$1 = \theta_K \leq \theta_1 \leq \frac{2N}{K} - 1 \quad (21)$$

To make the selection adaptively, a function for calculating  $\theta_1$  in the ranking selection is defined as follow.

$$\theta_1 = 1 + \left(\frac{2N}{K} - 2\right) \frac{F_a}{F_{best}} = 1 + \left(\frac{2N}{K} - 2\right) \frac{F_a}{F'_1} \quad (22)$$

$$\Delta\theta = \left(\frac{2N}{K} - 2\right) \frac{F_a}{F'_1} \frac{1}{(K - 1)} = \frac{2(N - 2K)F_a}{K(K - 1)F'_1} \quad (23)$$

where  $F_a$  represents of the average fitness value of the whole population. Thus the much closer  $F_{average}$  and  $F'$  are, the bigger  $\theta_1$  is, the bigger the number of times the better individuals to be selected. Otherwise, the number of times the worse individuals would be bigger and the selection of  $K$  individuals becomes more balance. The pseudo-code of this operator is shown below as Algorithm 1.

#### Algorithm 1. Ranking Selection Operator

##### Ranking\_Selection ( $\mathbf{I}$ )

Define the selection range according to  $P_{selection}$

Sort  $\mathbf{I}$  with quick sort algorithm and stored as  $\mathbf{I}'$

Calculate the number of times of  $I_1$  to be selected,

$$\theta_1 = 1 + (2N/K - 2)F_a/F'_1$$

Calculate  $\Delta\theta = 2(N - 2K)F_a/K(K - 1)F'_1$

Calculate  $\theta_2, \theta_3, \dots, \theta_K$  for other  $K - 1$  individuals

Select  $N$  individuals according to  $\theta_1, \theta_2, \dots, \theta_K$  and generate new  $\mathbf{I}$

#### 5.3. Individual chaos operator

Chaos is a universal non-linear phenomenon. It has the characteristics of strong randomness and internal regularity. With the generation of logistic chaos sequences, it can traverse almost

all states in a certain range without duplication and cause great changes in output with rich dynamism. Thus it can improve population diversity in many typical intelligent algorithms and help them to avoid local optimization. Nevertheless, it is non-directional and hard to control.

In general, the searching process of typical chaos-based optimization can be divided into two stages. In the first stage, a bunch of chaotic sequences with certain length are generated by logistic chaos generating function. Then one or more gene-bits of individuals are changed according to the chaotic sequences and a series of new individuals are generated. After the selection of good solution among these new individuals, the second stage will introduce a small disturbance to the local optimum individuals for further exploitation. The iteration will continue until the terminate standards are satisfied.

However, two problems come up to restrain the performance of chaos for large-scale problems with irregular solution spaces. First, small disturbance will not help to exploit in complex and irregular spaces. Besides, the length of chaotic sequence directly decides the time consumption and searching ability of the algorithm. For higher searching ability, the second problem is that fixed length of chaotic sequences may bring large time consumption in exploration. Thus, we design a new individual chaos operator in which the small disturbance is abandoned and adaptation of chaotic length is introduced for individuals with customization.

Specifically, the length of chaotic sequence for each individual is determined by its current evolutionary state. Let  $\mathbf{I} = \{I_i | i = 1, 2, \dots, N\}$  be the population with  $N$  individuals, and  $I_i$  be the  $i$ th individual. It includes its gene-bit values  $G_i = \{G_i(1), G_i(2), \dots, G_i(M)\}$  and fitness value  $F_i$ , where  $M$  represents the length of gene code (i.e. twice of the number of tasks). The specific pseudo-code is shown as Algorithm 2.

#### Algorithm 2. Individual Chaos Operator

##### Individual\_Chaos ( $\mathbf{I}$ )

For ( $i = 1 \rightarrow N$ )

$$L_{chaos} = A + B(F_{best} - F_i)/(F_{best} - F_{worst} + 1)$$

Generate  $X_1[L_{chaos}], X_2[L_{chaos}] \in [0, 1]$  by using Logistic chaos function

For ( $j = 1 \rightarrow L_{chaos}$ )

Map  $X_1(j)$  as genebit serial number  $k \in [1, M]$

If ( $k$  corresponds to service-bit)

Map  $X_2(j)$  as genebit value  $v \in [1, n_s]$

Else

Map  $X_2(j)$  as genebit value  $v \in [1, n_p]$

End if

Generate  $j$  new temporary individuals  $\{r(1), r(2), \dots, r(j)\}$  by replacing the value of  $G_i(k)$  with  $v$

Choose the best individual  $r_{best}$  from the temporary individuals

If ( $F_{r_{best}} > F_i$ )

$$I_i = r_{best}$$

Else

If ( $\exp((F_{r_{best}} - F_i)/t^0) > \gamma$ )

Replace  $I_i$  with  $r_{best}$

End if

End if

End for

$$t^0 = Dt^0$$

End for



Go in detail, the evolutionary state of the  $i$ th individual is defined as  $Q_i$ :

$$Q_i = \frac{F_{best} - F_i}{F_{best} - F_{worst}} \quad (24)$$

$$L_{chaos} = A + (B - A)Q \quad (25)$$

where  $A$  and  $B$  is the lower bound and upper bound of  $L_{chaos}$ , respectively.  $F_{best}$  and  $F_{worst}$  represent the serial numbers of the individual with the best and the worst fitness value. To be exact, the closer the average fitness value to the best fitness value in population, the better the evolutionary state is, and the shorter the length of chaotic sequence  $L_{chaos}$  is, so that the smaller the searching range is. Otherwise, the closer the average fitness value to the worst fitness value in population, the smaller the searching range is.

With the initialized definition of the length of chaotic sequences  $L_{chaos}$ , the operator generates two chaotic sequences  $X_1[L_{chaos}], X_2[L_{chaos}]$  for each individual  $I_i (i = 1, 2, \dots, N)$  by Logistic mapping chaotic function, as shown in Eq. (26).

$$z_{l+1} = \mu z_l (1 - z_l) \quad (26)$$

where  $\mu = 4$  according to general chaotic strategy. Then  $X_1$  and  $X_2$  are mapped to the serial number  $k$  and the value  $v$  of gene-bits respectively. If  $k \in [1, M]$  corresponds to service gene-bit, we should map  $X_2$  to relative service number and store it in  $v$ . Or we should map  $X_2$  to relative computing resource number and store it. In the pseudo-code,  $n_s$  and  $n_p$  represents the number of cloud services and the number of computing resources respectively. After the chaotic mapping step, new neighbor solutions  $\{r(1), r(2), \dots, r(j)\}$  can be generated by changing the value of  $G_i(j)$  into  $v$ . Further, choose the individual  $r_{best}$  with the best fitness value and accept it as new individual with probability  $P_{annealing} = \exp((F_{r_{best}} - F_i)/t^0)$ , where  $t^0$  is the annealing temperature and the initial value is 100. In the algorithm, the rate of  $t^0$  drop  $D$  is set to be 0.95 to gradually narrow down the accept probability. On the whole, in the individual chaos operator, searching is carried out with the adaptive changing of the length of chaotic sequences for each individual according to its evolutionary state  $Q_i$ . Chaos states can finally be controlled by population state.

#### 5.4. Dynamic heuristic operator

For further improving the searching direction in chaos optimization, dynamic heuristic is introduced in this algorithm after ranking selection and adaptive individual chaos. The principle of this operator is dynamically guiding the algorithm for local search with right direction by using some priori knowledge of the problem.

To be specific, for each individual, the operator randomly chooses a gene-bit, traverses part of the available values for the single gene-bit and dynamically calculates the heuristic of each value, then picks the most suitable value with the highest heuristic and generates new individual. It is quite like the mechanism of pheromone in ACO. Compared with the pheromone, dynamic heuristic here does not contain empirical information. It uses just the priori knowledge which is dynamically calculated according to the states or the gene-bit values of individuals in each generation. Define the traverse range for one gene-bit to be  $hn$ , where  $h \in [0, 1]$ , and  $n$  can be  $n_s$  or  $n_p$ . The specific pseudo-code is shown as follow.

#### Algorithm 3. Dynamic Heuristic Operator

*Dynamic\_Heuristic (I)*

For ( $i = 1 \rightarrow N$ )

```

     $r = I_i$ 
    Randomly choose a genebit  $k \in [1 : M]$ 
    If ( $p$  corresponds to service-bit)
        Randomly choose  $hn_s$  values from 1 to  $n_s$ 
        Choose the service  $s_i$  with the highest heuristic  $Y_s = \max_j y_s(j) (j \in [1, hn_s])$ 
         $G_r(k) = s_i$ 
    Else
        Randomly choose  $hn_p$  values from 1 to  $n_p$ 
        Choose the computing resource  $p_i$  with the highest heuristic  $Y_p = \max_j y_p(j) (j \in [1, hn_p])$ 
         $G_r(k) = p_i$ 
    End if
    Accept  $r$  with simulation annealing probability
End for

```

During the process,  $h$  can be set as 0.3. And  $p$  represents the randomly selected gene-bit for each individual. If  $p$  corresponds to service-genebit, search available services and calculate dynamic heuristic of each available service  $y_s(j) (j \in [1, hn_s])$  by service heuristic function. Then choose the service  $s_i$  with the highest heuristic  $Y_s$  to replace the original value of  $k$ th gene-bit. If  $k$  corresponds to computing resource gene-bit, search available computing resources and calculate dynamic heuristic of each computing resource  $y_p(j) (j \in [1, hn_p])$  by computing resource heuristic function. Then choose the computing resource  $p_i$  with the highest heuristic  $Y_p$  to replace the value of  $k$ th gene-bit. After these steps, a new individual  $r$  is generated for each individual. Then replace  $I_i$  with  $r$  in simulation annealing probability as well as in adaptive chaos operator  $P_{annealing} = \exp((F_r - F_i)/t^0)$ .

In this process, how to design service heuristic function and computing resource heuristic function is very important. Unsuitable heuristic function can cause wrong searching direction in algorithm and easily lead the algorithm to serious premature convergence. In this paper, the service and computing resource heuristic function are simply designed as follow.

$$y_s(j) = \begin{cases} \alpha_1 \frac{s_{s_j}}{v_{s_j}} + \alpha_2 \frac{1}{r_{s_j}} + \alpha_3 \frac{1}{e_{s_j}} + \alpha_4 \frac{1}{c_{s_j}} + \alpha_5 r_{s_j}, & \text{if } s_j \in S_1 \\ \alpha_1 \frac{s_{s_j}}{v_{s_j}} + \alpha_2 \frac{1}{c_{s_j}} + \alpha_3 r_{s_j}, & \text{if } s_j \in S_2 \end{cases} \quad (27)$$

$$y_p(j) = \alpha_1 \frac{x(j)}{\lambda(j)} + \alpha_2 \max(\varphi(j), \phi(j)) + \alpha_3 \frac{1}{\sigma(j)} + \alpha_4 \frac{1}{f(j)} \quad (28)$$

where  $\alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5$  represent the weights of services/ computing resources attributes respectively which corresponds to the weights setting in the objective function. Through the adjustment of weights, small range of local search in a single gene-bit could be guided in the algorithm according to the dynamic heuristics.

#### 5.5. The complexity of the proposed algorithm

Generally, the time complexity of the intelligent algorithms is dynamically varied with different problems. Let  $n$  be the scale of the population,  $m$  be the size of tasks,  $s$  be the number of the available cloud services for each task and  $p$  be the total scale of computing resources. The algorithms' complexities in each generation are shown in Table 2.

In GA, typical roulette wheel selection needs  $n$  times roulette operations to generating new population. Each roulette operation contains at least 1 and at most  $n$  times comparison according to the

**Table 2**

The complexity of the operators in GA and RCO.

Algorithms	The time complexities of operators			$n \rightarrow \infty$	$m \rightarrow \infty$	$s \rightarrow \infty$	$p \rightarrow \infty$
GA	Roulette wheel selection $O(n^2)$	Crossover $O(nm)$	Mutation $O(nm)$	$O(n^2)$	$O(m)$	$O(1)$	$O(1)$
RCO	Ranking selection $O(n \log n)$	Individual chaos $O(n)$	Dynamic heuristic $O(n'(\max(p, s)))$	$O(n \log n)$	$O(1)$	$O(s)$	$O(p)$

relative fitness values of individuals. Thus the average complexity of selection operator is  $O(n^2)$ . In RCO, the complexity of ranking individuals in selection is  $O(n \log n)$  (with quick sort method) and the selection step according to selective pressure needs at most  $n$  times. Thus the complexity of ranking selection is  $O(n \log n)$ .

Besides, crossover and mutation operation in GA are just executed once for each individual. The complexity are both at least  $O(n)$  and at most  $O(mn)$ . In RCO, chaotic sequences with constant length are generated for each individual. From the pseudo-code it can be seen that the complexity of individual chaos operator is  $O(nL_{chaos}) = O(n)$ . Because the adaptation of chaotic length is in a limited area, the complexity of chaos operator is also  $O(n)$ . It is lower than crossover operator. In addition, dynamic heuristic randomly chooses a gene-bit for each individual, traverse part of available value of this gene-bit with heuristics. If all of the selected gene-bits are service-bit, then the complexity is  $O(n_s)$ , else if all of the selected gene-bits are computing resource-bit, then the complexity is  $O(n_p)$ . Thus the average complexity of dynamic heuristic operator is  $O(n(s+p)/2) = O(n \max(p, s))$ .

In theory, if  $s \rightarrow \infty$  and  $p \rightarrow \infty$ , the complexity of RCO is a little higher than GA. But in the condition of  $n \rightarrow \infty$  and  $m \rightarrow \infty$ , the complexity of RCO is lower than GA.

## 6. Experiments and discussions

Based on the case “the design and NC (Numerical Control) machining process of a complex surface part” in Section 3, three typical DAG: two DAGs as shown in Fig. 5 [21] and the “j30” DAG of Resource-Constrained Project Scheduling Problem (RCPSP) in PSPLIB [31], are used as three task graphs in our experiments. In practical application of private cloud in manufacturing conglomerate or large-scale manufacturing service providers, a composite project contains multiple complex surface parts’ machining. Thus a composite project can be divided into far more than 5 tasks. Those tasks have several functional and non-functional requirements for cloud services. Some of them need hardware cloud services, some need software cloud services. In order to evaluate the performance of dual-scheduling optimization compared with the traditional two-level decision, we use basic real-coding GA uniformly to simulate the decision process in theory. At the OACR step, each gene-bit represents the selected computing resource number for the above selected service. The lengths of gene-bits at the two steps

are equal. Furthermore, At the SCOS step, we consider only the properties of cloud services and then set the objective function as following according to Eq. (14) and [32].

At the OACR step, we also use the objective function in Eq. (14) with the fixed properties of cloud services. In Eq. (18), let the weight to be  $w_1 = w_2 = w_3 = w_4 = 100n$ .

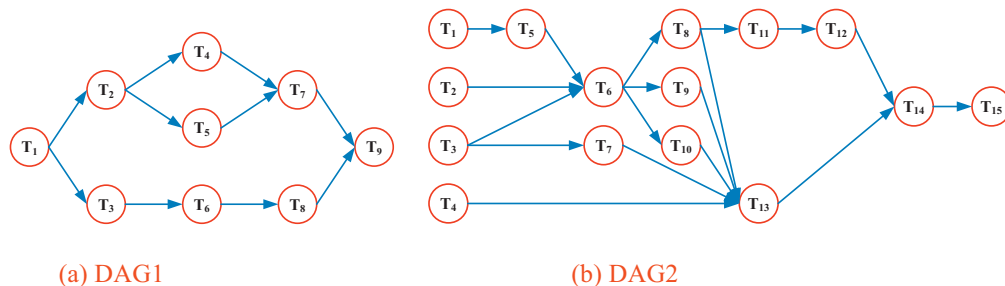
For simplifying the optimization process, we set that each task in a composite project has the same number of available cloud services. In the three cases, 3 composite scales of DS-CSCR are tested, as shown in Table 3. And in each scales, computing resources are equally divided into 5 distributed groups.

Assume the available number of cloud services for each task is  $s$  and the available number of computing resources is  $p$ , then the size of solution space is  $s^n p^n$ . From Scale 1 to Scale 9, it is range from  $10^9 \times 20^9$  to  $50^9 \times 100^9$ . Most deterministic algorithms cannot handle these situations due to composite exposition.

Moreover, because of the restriction of experimental environment, we set the ranges of properties of cloud services and computing resources as shown in Table 4.

For theoretical analysis, all the values are randomly generated with normalization and idealization and stored in a *txt* file. In order to distinguish the bandwidths inter-group and intra-group, the range of  $\varphi$  is set to be slightly larger than  $\phi$ . Initially, task load of all computing resources are 0.

Based on DS-CSCR with 9 scales, standard GA, chaos GA (CGA), typical chaos optimization (CO), chaos optimization with only individual chaos operator designed in this paper (RCO<sup>-2</sup>), chaos optimization with ranking selection and individual chaos operator (RCO<sup>-1</sup>) and chaos optimization with the addition of dynamic heuristics (RCO) are compared together for further testing the performance of the above designed algorithm. In the experiments, the classical roulette wheel selection operator, multiple-point crossover operator and single-point mutation operator are adopted in GA. And the crossover and mutation probabilities are set to be the typical values, i.e. 0.8 and 0.15, respectively. In chaos strategy of CGA and CO, the length of chaotic sequences is set as a constant 10. For a fairer comparison, in the new RCO, let  $A = 5$  and  $B = 15$  to make sure the same level of chaotic operation. Besides, the iterations of all experiments are set as 2000 uniformly and population sizes are all 20. Due to the randomness of intelligent algorithms, a total of 100 runs of each experiment

**Fig. 5.** Two typical DAG with 9 and 15 tasks respectively.

**Table 3**

The selected 4 composite scales of cloud services and computing resources.

	Scale 1	Scale 2	Scale 3	Scale 4	Scale 5	Scale 6	Scale 7	Scale 8	Scale 9
Number of tasks	9	9	9	15	15	15	30	30	30
Number of available cloud services	10	20	50	10	20	50	10	20	50
Number of available computing resources	20	50	100	20	50	100	20	50	100

**Table 4**

The property ranges of cloud services and computing resources.

	$s$	$c$	$r$	$\nu$	$e$	$\zeta$
Software service	[1, 10]	[1, 10]	(0, 1)	[1, 10]		
Hardware service	[1, 10]	[1, 10]	(0, 1)	[1, 10]	[1, 10]	[0, 1]
	$x$	$\varphi$	$\phi$	$\sigma$	$f$	$\lambda$
Computing resource	[1, 10]	[1, 10]	[1, 5]	[1, 10]	(0, 1)	0

are conducted and the average fitness value of the best solutions throughout the run is recorded.

#### 6.1. Performance of DS-CSCR compared with traditional two-level scheduling (TL-S)

Let TL-S to be the abbreviation of traditional Two-Level Scheduling, we compared it with new DS-CSCR in the above 9 scales of solution space. Fig. 6 shows the testing results from the perspectives of time consumption and solution quality respectively.

Firstly, we define the *decrease-rate* to be  $\tau_d = (T_{TL-S} - T_{CS/CR-DS})/T_{TL-S}$  in Fig. 6(a). As we have analyzed previously, the time consumption of SCOS and OACR in traditional TL-S are reduced by about 35–40 percent in DS-CSCR. Although the length of individual and the size of solution space are only half that of DS-CSCR. Traditional TL-S takes almost twice the time of DS-CSCR. For each task graph, as the numbers of cloud services and computing resources are enlarged, the *decrease-rate* increases gradually. Thus it can be seen, with the same algorithm (no matter deterministic algorithm or intelligent algorithm), TL-S is more and more time-consuming with the increase of solution space while DS-CSCR always maintains a relatively low level of time consumption. It proved that, with the same algorithm, no matter using deterministic or intelligent, two level decision is cumbersome.

Secondly, from the angle of solution quality in Fig. 6(b), we define the *growth-rate* to be  $\tau_g = (F_{CS/CR-DS} - F_{TL-S})/F_{TL-S}$ , where  $F_{CS/CR-DS}$  and  $F_{TL-S}$  represent the average result of the best fitness value in experiments. It increases along with the expansion of solution spaces in each kinds of task graph. For all of scales, the total level of quality in TL-S is improved by about 14–19 percent in DS-CSCR. In theory, the service properties are static in the second step of TL-S. With the splitting of SCOS and OACR under unified console, the mutual relations between cloud service and the underlying computing resources are ignored. This tells us the conclusion that in private cloud, the underlying support infrastructure must be considered in the process of SCOS. With fast development of dynamic network, service with dynamic deployment are more and more common. SCOS with the consideration of QoS only are impractical for many of the advanced system in large SaaS mode.

#### 6.2. Searching capability of RCO for solving DS-CSCR

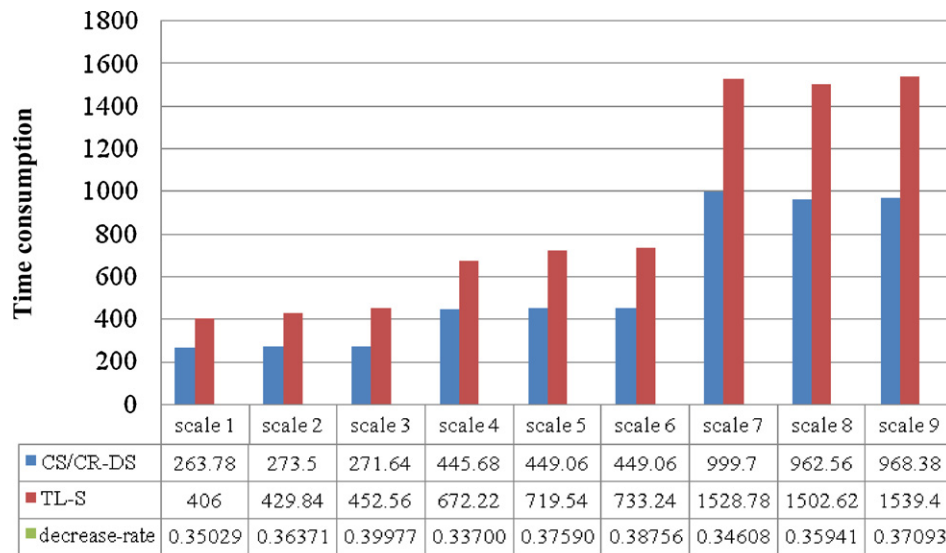
For addressing DS-CSCR more efficiently, we designed RCO especially aiming at the situation of large-scale solution space.

Fig. 7 recorded the average fitness value of the best solution during 2000 generations in 100 runs for 9 scales of DS-CSCR (i.e. the average evolutionary trend of the 6 algorithms in 100 runs). Fig. 8 shows the fitness value of the best solution, the worst solution and the average result in 100 runs for 9 scales of DS-CSCR. Note that the fitness value is the assessment value of each individual according to the objective function. So from the perspective of searching capability, the sort of the six algorithms from bad to good is:  $GA < CGA < CO < RCO^{-2} < RCO^{-1} < RCO$ . The step-by-step improvement from the design of individual chaos operator to the introduction of ranking selection and dynamic heuristic operators can be clearly observed.

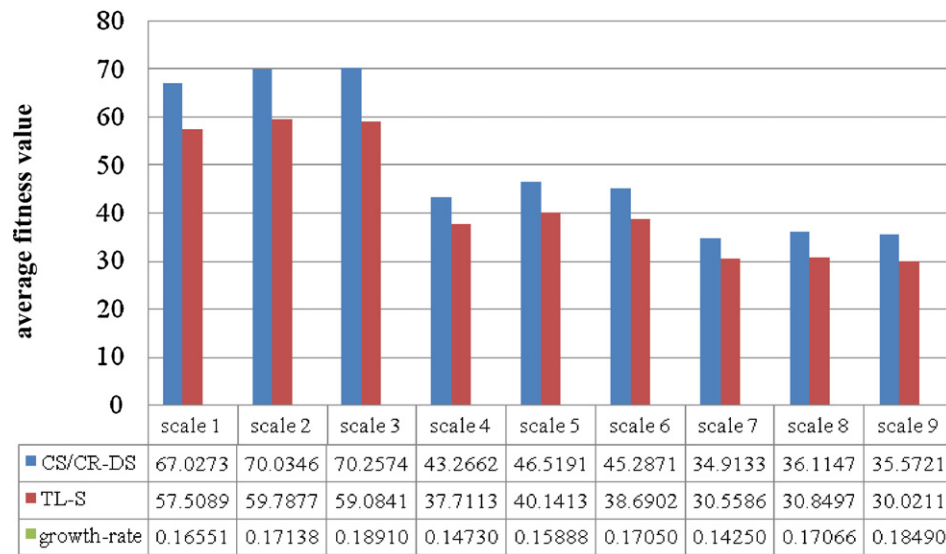
On the basis of GA, the average fitness value of the best solutions of CGA is about 30% higher than GA. At this moment, the average best fitness value of CO with single chaos optimal operator is about 1.5 times higher than GA. From here we can come to the conclusion that the basic operators of GA constrained the searching ability of chaos optimal operator in CGA to some degree. Simple chaos optimization can get much better solution than the traditional GA and improved CGA. Furthermore, the adaptive strategy adapts chaotic sequences according to the state of the whole population. When the population is in a good state, the adaptive strategy will reduce the chaotic sequences, so as to reduce the complexity of the algorithm. Compared with CO, the average best fitness value of  $RCO^{-2}$  in the 9 scales of DS-CSCR has been raised by about 3%. Afterwards, ranking selection was put in the front of  $RCO^{-2}$ . With the collaboration of selection and chaos, the average best fitness value of  $RCO^{-1}$  is improved again. Hence, it can be learned that the operation and collaboration of individual chaos operator and the “the survival of the fittest” ranking selection strategy can not only reduce the complexity of algorithm, but also improve the searching capability remarkably. Because the effect of mutation is similar to chaos operator, it may conclude that the crossover operator in GA mainly restrained the capability of chaos strategy in CGA. Based on the improved  $RCO^{-1}$ , for guiding chaos optimization further, dynamic heuristic operator was introduced at last. From Figs. 7 and 8 we can see that the new RCO performs better than  $RCO^{-1}$  with the guidance of heuristics. On the whole, the average best fitness value of RCO in 100 runs is about 2 times higher than GA. The overall improvements are extremely considerable.

#### 6.3. Time consumption and stability of RCO for solving DS-CSCR

Next, based on the above mentioned 9 scales with 3 kinds of task graphs (Table 3), the time efficiency and stability of the 6 algorithms are discussed below. Note that the time



(a) The average solution of DS-CSCR and TL-S based on GA in 9 scales



(b) The average solution of DS-CSCR and TL-S based on GA in 9 scales

**Fig. 6.** Comparison of DS-CSCR and TL-S based on GA.

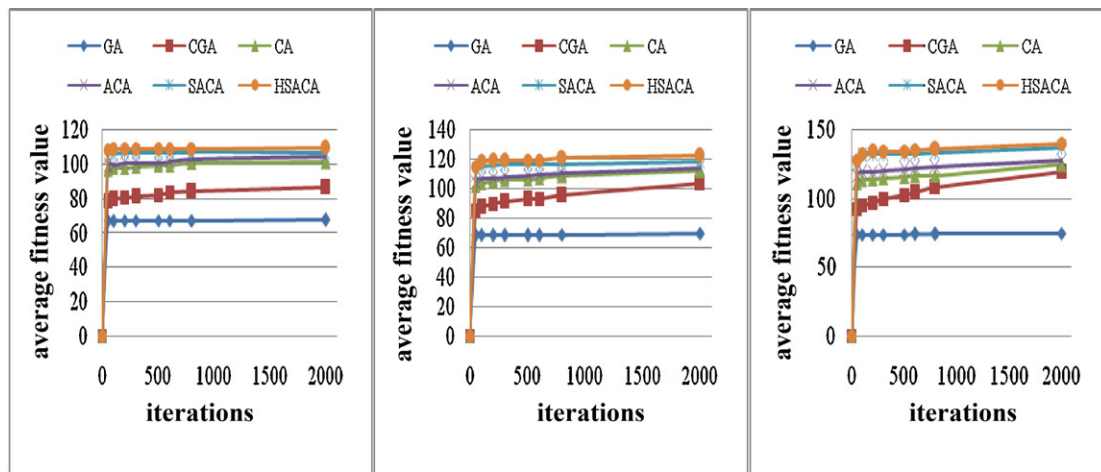
consumption are tested in milliseconds (ms) and the stability is measured by the standard deviation of the average fitness values in 100 runs.

Fig. 9(a) shows the average time-consumption of the 6 algorithms in 9 scales with 100 runs. The step-by-step improvement from CO to RCO compared with GA and CGA, the variation trends of time in all scales are the same. In CGA, there are four operators (selection, crossover, mutation and chaos), with lower searching capability, its time-consumption is the highest in these 6 algorithms. After wiping out the three operators of GA, the times of CO are just lower than CGA. It is clear that the most time-consuming operator in CGA is chaos operator. Only narrowing down the chaotic traverse range can reduce the total execution time of algorithm. Along with the decrease of chaotic sequences, the searching ability of algorithm will be reduced, too. Therefore, in order to reduce the time complexity of algorithm with the maintaining of the searching ability, individual chaos operator

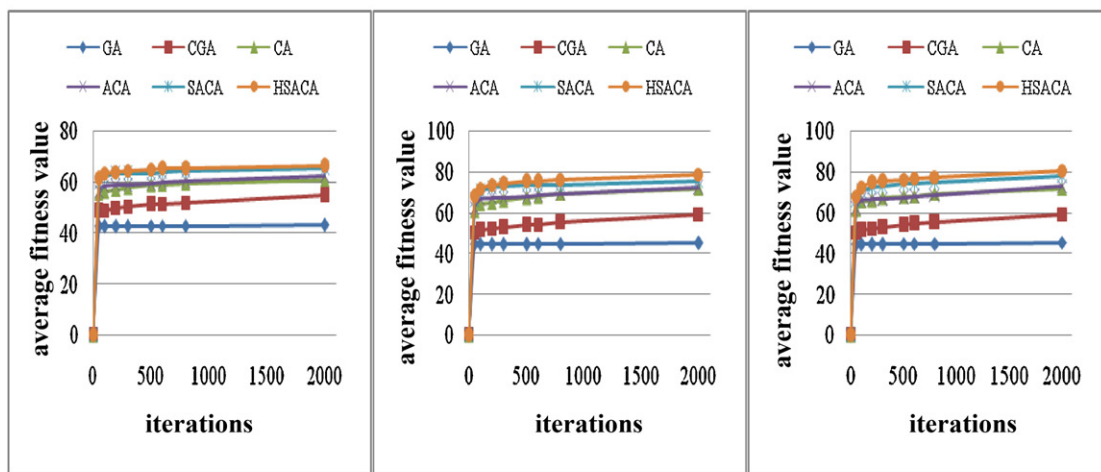
customized for individuals is designed in this paper. Experiments in  $RCO^{-2}$  show that the time-consuming is effectively reduced by about 20% based on CO with the improvement of searching ability. Especially in scale 7, 8 and 9 with very large solution spaces, time-consuming of chaotic operations are sharply reduced.

Moreover, the introduction of ranking selection not only improved the searching capability of  $RCO^{-2}$ , but also reduced the time. The reason is that, based on ranking selection, the difference between the best fitness value and the average fitness value in the population is shortened, the population can always be adapted to a better state with "the survival of the fittest" strategy, then the chaotic sequences are shortened accordingly. With shorter chaotic sequences, the population can be guided to better areas based on fitter individuals and then find better solutions more quickly. In terms of the time measuring, the prominent performance of the collaborative operation of ranking selection and individual chaos operator

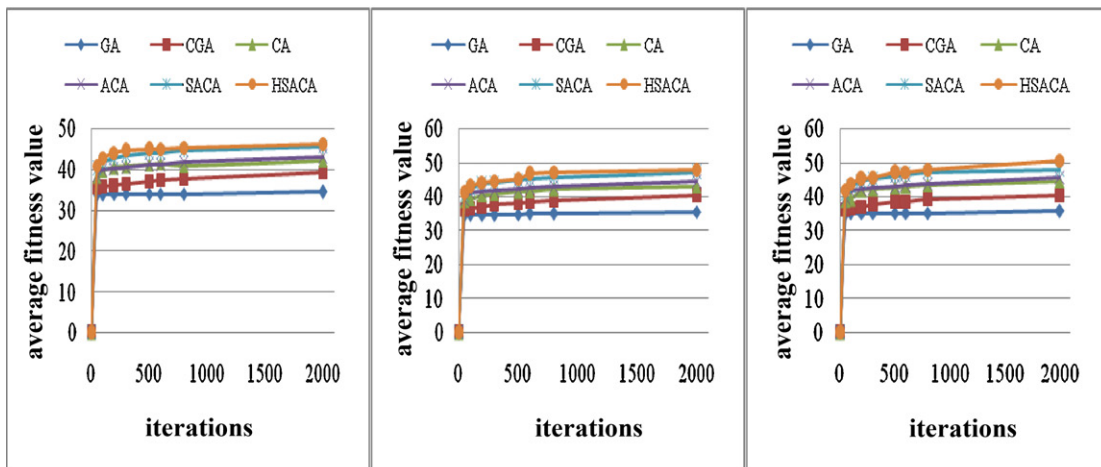




(a) DAG1 with 9 tasks in the scale 1, 2 and 3

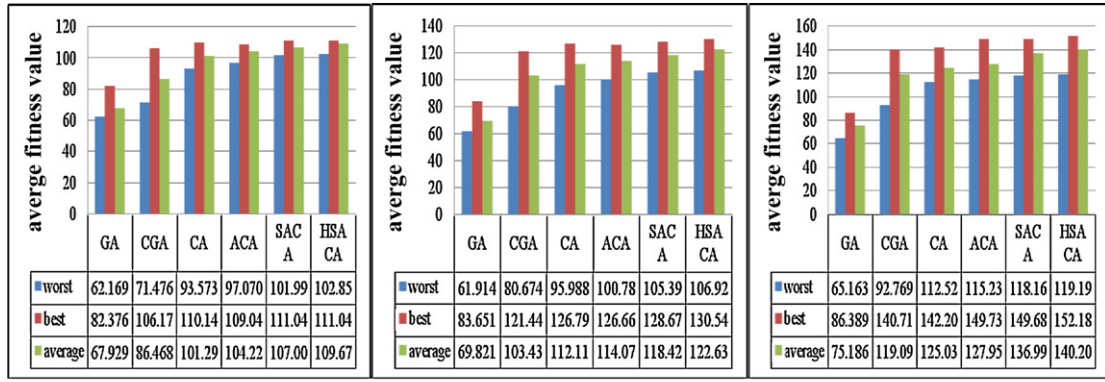


(b) DAG2 with 15 tasks in the scale 4, 5 and 6

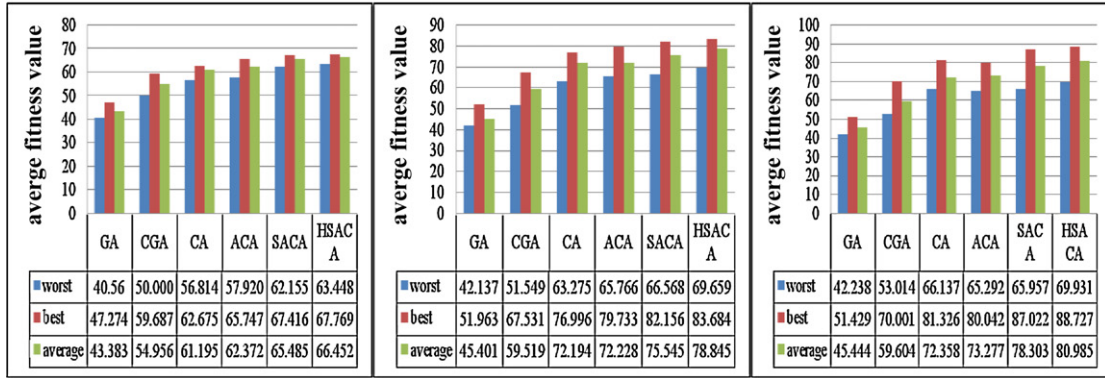


(c) DAG3 with 30 tasks in the scale 7, 8 and 9

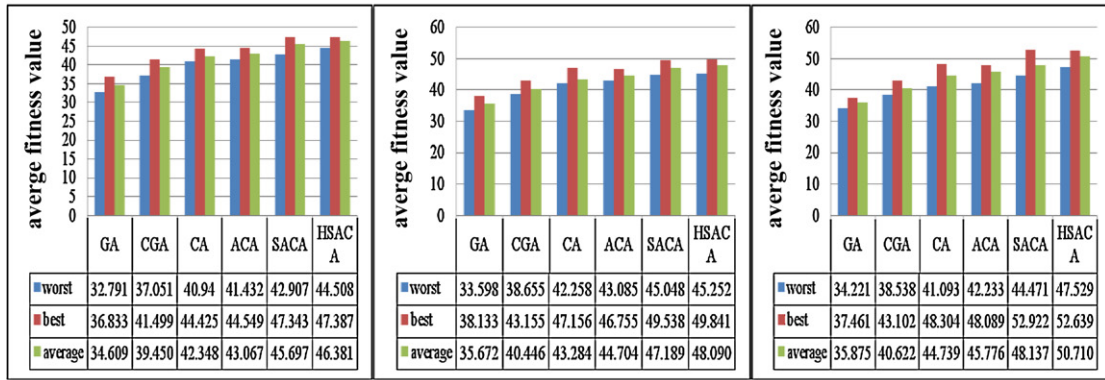
Fig. 7. The average evolutionary trend of the 6 algorithms in 100 runs for 9 scales of DS-CSCR.



(a) DAG1 with 9 tasks in the scale 1, 2 and 3



(b) DAG2 with 15 tasks in the scale 4, 5 and 6



(c) DAG3 with 30 tasks in the scale 7, 8 and 9

Fig. 8. The statistical results of the 6 algorithms in 100 runs for 9 scales of DS-CSCR.

has been verified again as  $RCO^{-1}$ . At the next step, the introduction of dynamic heuristic operator increase the time slightly based on  $RCO^{-1}$ , but the new complete RCO is much faster than  $RCO^{-2}$ , CGA and CO as a whole.

From the perspective of stability, as shown in Fig. 9(b), the six algorithms in the 9 problem scales changed irregularly. But from the 9 scales of tests, we can obtain the sort of stability of the six algorithms from bad to good is:  $CGA < CO < RCO^{-2} < RCO^{-1} < RCO < GA$ . Traditional GA is the most stable while the stability of CGA is the worst. With the adaptive improvement,  $RCO^{-2}$  is more stable than CO. That is because in large-scale solution space, chaotic sequences are generated based on no matter good or bad individuals, the

population is easy to be lead to bad areas during searching and the states of population in each generations are not stable any more. After the introduction of ranking selection, the stability of the algorithm has greatly improved. Each time of selection in iteration maintained the population state and reduced the chaotic sequences, so that the population can always be evolved based on fitter individuals with higher stability. Besides, the design of dynamic heuristic operator with the priori knowledge of DS-CSCR can always guide the population into better areas during evolution and then improve the stability further.

Thus it can be seen that the new designed RCO possesses plenty of advantages in searching capability, time-consumption and

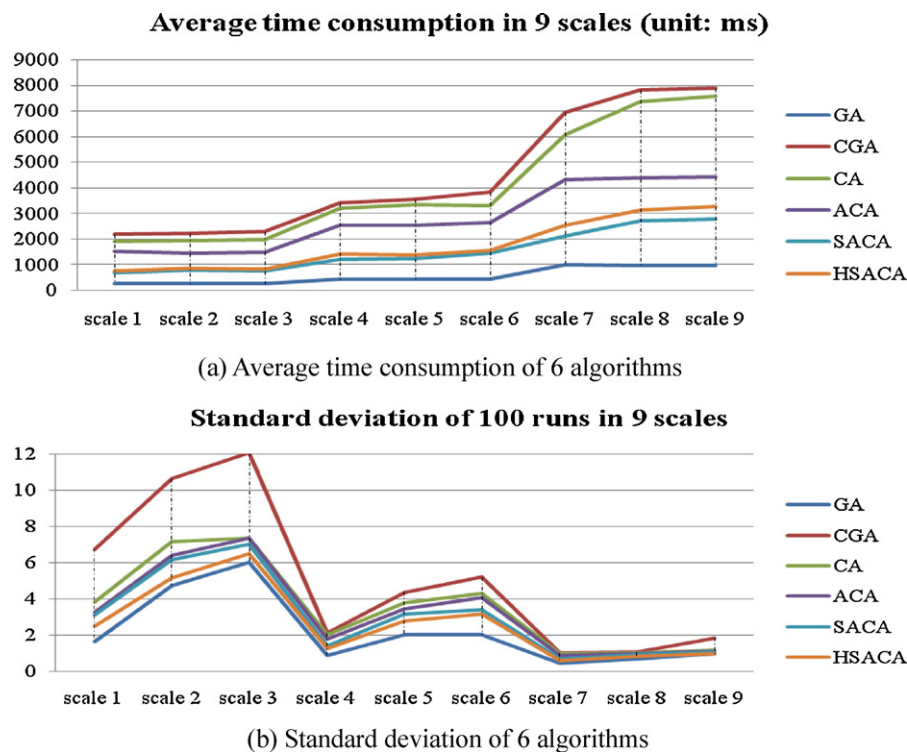


Fig. 9. The average time-consumption and standard deviation of the 6 algorithms in 100 runs.

stability for addressing DS-CSCR no matter with large or small scales solution spaces in private cloud.

## 7. Conclusions

Service Composition Optimal Selection (SCOS) and Optimal Allocation of Computing Resource (OACR) are both very critical in cloud system. Current works found that the two steps decision of SCOS and OACR in private cloud are quite cumbersome and the mutual relations between cloud services and underlying computing resources are always ignored. Thus this paper deeply analyzed the characteristics of these two problems and their interactions. Based on this, the idea of one-time decision of SCOS and OACR was presented accordingly. To sum up, the primary works and contribution of this paper can be concluded as follows.

- (1). New DS-CSCR model was presented in private cloud for high efficient one-time decision. Properties of software/hardware cloud services, VMs and computing resources are deeply analyzed. The formulation of DS-CSCR was clarified according to the aim of high efficient and low cost resource sharing.
- (2). For addressing the complex dual scheduling problem (DS-CSCR), a new intelligent algorithm – RCO was presented. Individual chaos operator was designed as the backbone operator of the algorithm. Then a new adaptive ranking selection was introduced for control the state of population in iteration. Moreover, dynamic heuristics were also defined and introduced to guide the chaos optimization. RCO with these three operators showed remarkable performances in terms of searching ability, time complexity and stability in solving the DS-CSCR problem in such private cloud compared with other algorithms.

Future work includes intensive study on the operation mechanism of private cloud and the complex mutual relationships among tasks, cloud services and computing resources more deeply. The QoS properties for manufacturing capabilities and human

resources are needed to be studied, too. In addition, RCO presented in this paper still has some disadvantages. The design of heuristic function for specific problems in the dynamic heuristic operator is complex and hard though. And its convergence has not been proved. As a new improved intelligent algorithm, its effectiveness in various other complex combinatorial optimization problems remains to be further explored and validated.

## Acknowledgements

This work is partly supported by the Fundamental Research Funds for the Central Universities, the NSFC projects (Nos. 61074144 and 51005012), and National Hig-Tech. R&D (863) Program (No. 2011AA040501) in China.

## References

- [1] G. Boss, P. Malladi, D. Quan, L. Legregni, H. Hall, Cloud computing, IBM White Paper, 2007, [http://download.boulder.ibm.com/ibmdl/pub/software/dw/wes/hipods/Cloud\\_computing\\_wp\\_final\\_8Oct.pdf](http://download.boulder.ibm.com/ibmdl/pub/software/dw/wes/hipods/Cloud_computing_wp_final_8Oct.pdf).
- [2] M. Armbrust, A. Fox, R. Griffith, A.D. Joseph, R.H. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica, M. Zaharia, Above the Clouds: a Berkeley View of Cloud Computing, University of California, Berkeley, 2009.
- [3] T.Z. Xia, Z. Li, N.H. Yu, Research on cloud computing based on deep analysis to typical platforms, Lecture Notes in Computer Science 5931 (2009) 601–608.
- [4] X. Xu, From cloud computing to cloud manufacturing, Robotics and Computer-Integrated Manufacturing 28 (1) (2012) 75–86.
- [5] D. Wu, L. Thames, D. Rosen, D. Schaefer, Towards a cloud-based design and manufacturing paradigm: looking backward, looking forward, in: Proceedings of the ASME 2012 International Design Engineering Technical Conference & Computers and Information in Engineering Conference, Chicago, USA, 2012.
- [6] L.M. Vaquero, L. Roderio-Merino, J. Caceres, M. Lindner, A break in the clouds: towards a cloud definition, ACM SIGCOMM Computer Communication Review 39 (1) (2009) 50–55.
- [7] B.H. Li, L. Zhang, S.L. Wang, F. Tao, J.W. Cao, X.D. Jiang, X. Song, D. Chai, Cloud manufacturing: a new service-oriented networked manufacturing model, Computer Integrated Manufacturing Systems 16 (1) (2010) 1–16.
- [8] J.M. Nick, D. Cohen, B.S. Kaliski, Key enabling technologies for virtual private clouds, Handbook of Cloud Computing, Springer, vol. 1, 2010, pp. 47–63.
- [9] Gathering clouds of XaaS1, [https://www.ibm.com/developerworks/mydeveloperworks/blogs/sbose/entry/gathering\\_clouds\\_of\\_xaas](https://www.ibm.com/developerworks/mydeveloperworks/blogs/sbose/entry/gathering_clouds_of_xaas).



- [10] W. Tan, Y.S. Fan, M.C. Zhou, Data-driven service composition in enterprise SOA solution: a petri net approach, *IEEE Transactions on Automation Science and Engineering* 7 (3) (2010) 686–694.
- [11] F. Tao, Y.F. Hu, D. Zhao, Z.D. Zhou, H.J. Zhang, Z.Z. Lei, Study on manufacturing grid resource service QoS modeling and evaluation, *International Journal of Advanced Manufacturing Technology* 41 (9–10) (2009) 1034–1042.
- [12] F. Tao, Y.F. Hu, Z.D. Zhou, Application and modeling of resource service trust-QoS evaluation in manufacturing grid system, *International Journal of Production Research* 47 (6) (2009) 1521–1550.
- [13] F. Tao, D. Zhao, Y.F. Hu, Z.D. Zhou, Correlation-aware resource service composition and optimal-selection in manufacturing grid, *European Journal of Operational Research* 201 (1) (2010) 129–143.
- [14] K. Fujii, T. Suda, Semantics-based dynamic service composition, *IEEE Journal on Selected Areas in Communications* 23 (12) (2005) 2361–2372.
- [15] A.J. Ferrer, F. Hernandez, J. Tordsson, E. Elmroth, A. Ali-Eldin, C. Zsigri, R. Sirvent, J. Guitart, R.M. Djemame, W. Ziegler, T. Dimitrakos, S.K. Nair, G. Kousiouris, K. Konstanteli, T. Varvarigou, B. Hudzia, A. Kipp, S. Wesner, M. Corrales, N. Forgo, T. Sharif, C. Sheridan, OPTIMIS: a holistic approach to cloud service provisioning, *Future Generation Computer Systems* 28 (1) (2012) 66–77.
- [16] M. Mika, G. Waligora, J. Weglarz, Modeling and solving grid resource allocation problem with network resources for workflow applications, *Journal of Scheduling* 14 (3) (2011) 291–306.
- [17] J. Tordsson, R.S. Montero, R. Moreno-Vozmediano, I.M. Liorente, Cloud brokering mechanisms for optimized placement of virtual machines across multiple providers, *Future Generation Computer Systems* 28 (2) (2012) 358–367.
- [18] P.T. Endo, A.V.D. Palhares, N.N. Pereira, G.E. Goncalves, Resource allocation for distributed cloud: concepts and research challenges, *IEEE Network* 25 (4) (2011) 42–46.
- [19] Y.B. Ma, S.H. Jang, J.S. Lee, QoS and ontology-based resource management in cloud computing environment, *Information: An International Interdisciplinary Journal* 14 (11) (2011) 3707–3715.
- [20] P.C. Xiong, Y. Chi, S.H. Zhu, H.J. Moon, C. Pu, H. Hacigumus, Intelligent management of virtualized resources for database systems in cloud environment, in: 27th IEEE International Conference on Data Engineering, 2011.
- [21] Y.H. Zhang, Y.H. Li, W.M. Zheng, Automatic software deployment using user-level virtualization for cloud-computing, *Future Generation Computer Systems* (2011).
- [22] H. Ghanbari, B. Simmons, M. Litoiu, G. Iszlai, Feedback-based optimization of a private cloud, *Future Generation Computer Systems* 28 (1) (2012) 104–111.
- [23] Y.J. Laili, F. Tao, L. Zhang, B.R. Sarker, A study of optimal allocation of computing resources in cloud manufacturing systems, *International Journal of Advanced Manufacturing Technology* (2012), <http://dx.doi.org/10.1007/s00170-012-9393-0>.
- [24] A. Nathani, S. Chaudhary, G. Somani, Policy based resource allocation in IaaS cloud, *Future Generation Computer Systems* 28 (1) (2012) 94–103.
- [25] Y. Ma, C.W. Zhang, Quick convergence of genetic algorithm for QoS-driven web service selection, *Computer Networks* 52 (5) (2008) 1093–1104.
- [26] P.Y. Yin, J.Y. Wang, Optimal multiple-objective resource allocation using hybrid particle swarm optimization and adaptive resource bounds technique, *Journal of Computational and Applied Mathematics* 216 (1) (2008) 73–86.
- [27] H. Wada, J. Suzuki, Y. Yamano, K. Oba, Evolutionary deployment optimization for service-oriented clouds, *Software: Practice & Experience* 41 (5) (2011) 469–493.
- [28] F. Tao, L. Zhang, V.C. Venkatesh, Y.L. Luo, Y. Cheng, Cloud manufacturing: a computing and service-oriented manufacturing model, *Proceedings of the Institution of Mechanical Engineers, Part B: Journal of Engineering Manufacture* 225 (10) (2011) 1969–1976.
- [29] D. Schaefer, L. Thames, R.D. Wellman, D. Wu, Distributed collaborative design and manufacture in the cloud – motivation, infrastructure and education, in: *ASEE 2012 Annual Conference and Exposition*, Texas, 2012.
- [30] A. Beloglazov, J. Abawajy, R. Buyya, Energy-aware resource allocation heuristics for efficient management of data centers for cloud computing, *Future Generation Computer Systems* 28 (5) (2012) 755–768.
- [31] <http://129.187.106.231/psplib/>.
- [32] F. Tao, D.M. Zhao, Y.F. Hu, Z.D. Zhou, Resource service composition and its optimal-selection based on swarm optimization in manufacturing grid system, *IEEE Transactions on Industrial Informatics* 4 (4) (2008) 315–327.



**Yuanjun Laili** received the MS Degree and is studying for a Ph.D. Degree in the school of automation science and electrical engineering at Beihang University. Her main research interests are in the areas of service-oriented manufacturing, operations research, intelligent optimization and high performance computing. She has been involved in the project of cloud manufacturing (CMfg) and national projects of virtualization, and her main emphasis is on the decision-making and scheduling of distributed resource services in the cloud.



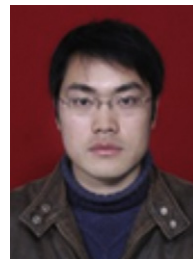
**Fei Tao** is currently an associate professor at School of Automation Science and Electrical Engineering in Beihang University (Beijing University of Aeronautics and Astronautics) since April 2009. He obtained his Ph.D from Wuhan University of Technology in 2008. From Sep. 2007 to Mar. 2009, he worked as a research scholar and postdoctoral researcher at University of Michigan-Dearborn, USA. His research interests include service-oriented manufacturing system such as cloud manufacturing and manufacturing grid, manufacturing service management, intelligent optimization theory and algorithm. He is the author of 2 monographs and over 60 journal and conference articles of these subjects. Dr. Tao was nominated and elected to be a research affiliate of CIRP (The International Academy for Production Engineering) in 2009. He is currently the editor of *International Journal of Service and Computing-oriented Manufacturing (IJSCOM)*, and the editorial board member of *International Journal of Modeling, Simulation, and Scientific Computing*.



**Lin Zhang** received the B.S. degree in 1986 from the Department of Computer and System Science at Nankai University, China. He received the M.S. degree and the Ph.D. degree in 1989 and 1992 from the Department of Automation at Tsinghua University, China, where he worked as an associate professor from 1994. He served as the director of CIMS Office, National 863 Program, China Ministry of Science and Technology, from December 1997 to August 2001. From 2002 to 2005 he worked at the US Naval Postgraduate School as a senior research associate of the US National Research Council. Now he is a full professor in Beihang University. He is an associate Editor-in-Chief of *International Journal of Modeling, Simulation, and Scientific Computing*. His research interests include integrated manufacturing systems, system modeling and simulation, and software engineering. Prof. Zhang is an IEEE senior member and a director of board of SCS.



**Ying Cheng** has a BS Degree in the school of mechanical and electronic engineering at Wuhan University of Technology and is studying for a Ph.D. Degree in the school of automation science and electrical engineering at Beihang University. Her main research interests are in the areas of service-oriented manufacturing and resource services management. She is involved in the projects of cloud manufacturing (CMfg) and the national projects of virtualization, and her main emphasis is on the utility model and equilibrium of resource service in the cloud.



**Yongliang Luo** has a MS Degree in the college of information and electrical engineering at Shandong University of Science and Technology and is studying for a Ph.D. Degree in the school of automation science and electrical engineering at Beihang University. His main research interests are in the areas of software engineering, the application of e-business, data mining and service-oriented manufacturing. He has been involved in the project of cloud manufacturing (CMfg) and is mainly focuses on the researches of manufacturing capability and the construction of cloud platform in CMfg.



**Bhabha R. Sarker** is the Elton G. Yates Distinguished Professor at the Department of Mechanical and Industrial Engineering, Louisiana State University at Baton Rouge, LA. He obtained his PhD degree from Texas A&M University. He is a professional engineer and won the 2006 Dr. David F. Baker Distinguished Research Award from the Institute of Industrial Engineers (IIE) for outstanding research. His teaching and research interests are in the areas of production and manufacturing systems, supply chain management, lean production systems, logistics and distribution systems, and operations research. He is a member ASEE, DSI, INFORMS, POMS, and a Fellow of IIE.