



**RESEARCH ARTICLE**

# Online Optimization for Scheduling Preemptable Tasks on IaaS Cloud Systems

S. Gandhimathi<sup>1</sup>, M. Madhushudhanan<sup>2</sup>, A. Srivishnu Paraneetharan<sup>3</sup>

<sup>1</sup>Department of Computer Science, PGP College of Arts and Science, Periyar University, India

<sup>2</sup>Department of Computer Science, PGP College of Arts and Science, Periyar University, India

<sup>3</sup>Department of Computer Science, PGP College of Arts and Science, Periyar University, India

<sup>1</sup> [gandhipgp@gmail.com](mailto:gandhipgp@gmail.com); <sup>2</sup> [mshudhanan@gmail.com](mailto:mshudhanan@gmail.com); <sup>3</sup> [vishnu\\_sri@yahoo.com](mailto:vishnu_sri@yahoo.com)

---

**Abstract**— *In Infrastructure-as-a-Service (IaaS) cloud computing, computational resources are provided to remote users in the form of leases. For a cloud user, he/she can request multiple cloud services simultaneously. In this case, parallel processing in the cloud system can improve the performance. When applying parallel processing in cloud computing, it is necessary to implement a mechanism to allocate resource and schedule the execution order of tasks. Furthermore, a resource optimization mechanism with preemptable task execution can increase the utilization of clouds. In this paper, we propose two online dynamic resource allocation algorithms for the IaaS cloud system with preemptable tasks. Our algorithms adjust the resource allocation dynamically based on the updated information of the actual task executions. And the experimental results show that our algorithms can significantly improve the performance in the situation where resource contention is fierce.*

---

## I. INTRODUCTION

In cloud computing, a cloud is a cluster of distributed computers providing on-demand computational resources or services to the remote users over a network. In an Infrastructure-as-a-Service (IaaS) cloud, resources or services are provided to users in the form of leases. The users can control the resources safely, the capacity used. Cloud computing is emerging with growing popularity and adoption. However, there is no data center that has unlimited capacity. Thus, in case of significant client demands, it may be necessary to overflow some workloads to another data center. These workload sharing can even occur between private and public clouds, or among private clouds or public clouds. The workload sharing is able to enlarge the resource pool and provide even more flexible and cheaper resources. To collaborate the execution across multiple clouds, the monitoring and management mechanism is a key component and requires the consideration of provisioning, scheduling, monitoring, and failure management. The two major contributions of this paper are:

- We present a resource optimization mechanism in heterogeneous IaaS federated multi-cloud systems, which enables preemptable task scheduling. This mechanism is suitable for the autonomic feature within clouds and the diversity feature of VMs.
- We propose two online dynamic algorithms for resource allocation and task scheduling. We consider the resource contention in the task scheduling.

## II. MODEL AND BACKGROUND

### 2.1 Cloud system

In this paper, we consider an infrastructure-as-a-service (IaaS) cloud system. In this kind of system, a number of data center participates in a federated approach. These data centers deliver basic on-demand storage and compute capacities over Internet. The provision of these computational resources is in the form of virtual machines (VMs) deployed in the data center. These resources within a data center form a cloud. Virtual machine is an abstract unit of storage and compute capacities provided in a cloud. Without loss of generality, we assume that VMs from different clouds are offered in different types, each of which has different characteristics. For example, they may have different numbers of CPUs, amounts of memory and network bandwidths. As well, the computational characteristics of different CPU may not be the same.

For a federated cloud system, a centralized management approach, in which a super node schedule tasks among multiple clouds, may be an easy way to address the scheduling issues in such system. Thus we propose a distributed resource allocation mechanism that can be used in either federated cloud system or the future cloud system with multiple providers.

As shown in Fig. 1, in our proposed cloud resource allocation mechanism, every data center has a manager server that knows the current statuses of VMs in its own cloud. And manager servers communicate with each other. Clients submit their tasks to the cloud where the dataset is stored. Once a cloud receives tasks, its manager server can communicate with manager servers of other clouds, and distribute its tasks across the whole cloud system by assigning them to other clouds or executing them by itself.

When distributing tasks in the cloud system, manager servers should be aware of the resource availabilities in other clouds, since there is not a centralized super node in the system. Therefore, we need the resource monitoring infrastructure in our resource allocation mechanism. In cloud systems, resource monitoring infrastructure involves both producers and consumers. Producers generate status of monitored resources. And consumers make use of the status information. Two basic messaging methods are used in the resource monitoring between consumers and producers: the pull mode and the push model. Consumers pull information from producers to inquire the status in the pull mode. In the push mode, when producers update any resource status, they push the information to the consumers. The advantage of the push mode is that the accuracy is higher when the threshold of a status update, i.e., trigger condition, is defined properly. And the advantage of the pull mode is that the transmission cost is less when the inquire interval is proper.

In our proposed cloud system resource allocation mechanism, we combine both communication modes in the resource monitoring infrastructure. In our proposed mechanism, when the manager server of cloud A assigns an application to another cloud B, the manager server of A is the consumer. And the manager server of B is the producer. Manager server of A needs to know the resource status from the manager server of B in two scenarios: (1) when the manager server of A is considering assigning tasks to cloud B, the current resource status of cloud B should be taken into consideration. (2) When there is a task is assigned to cloud B by manager server of A, and this task is finished, manager server of A should be informed.

We combine the pull and the push mode as the following:

- A consumer will pull information about the resource status from other clouds, when it is making scheduling decisions.
- After an application is assigned to another cloud, the consumer will no longer pull information regarding to this application.
- When the application is finished by the producer, the producer will push its information to the consumer. The producer will not push any information to the consumer before the application is finished.

In a pull operation, the trigger manager server sends a task check inquire to manager servers of other clouds. Since different cloud providers may not be willing to share detailed information about their resource availability, we propose that the reply of a task check inquire should be as simple as possible. Therefore, in our proposed resource monitoring infrastructure, these target manager servers give only responses at the earliest available time of required resources, based on its current status of resources & no guarantee or reservation is made. Before target manager servers check their resource availability, they first check the required dataset locality. If the required dataset is not available in their data center, the estimated transferring time of the dataset from the trigger cloud will be included in the estimation of the earliest available time of required resources. Assuming the speed of transferring data between two data centers is  $S_c$ , and the size of the required dataset is  $MS$ , then the preparation overhead is  $MS/S_c$ . Therefore, when a target cloud already has the required dataset in its data center, it is more likely that it can respond sooner at the earliest available time of required resources, which may lead to an assignment to this target cloud. In a push operation, when B is the producer and A is the consumer, the manager server of B will inform the manager server of A the time when the application is finished.

When a client submits his/her workload, typically an application, to a cloud, the manager server first partitions the application into several tasks, as shown in Fig. 2. Then for each task, the manager server decides which cloud will execute this task based on the information from all other manager servers and the data dependencies among tasks. If the manager server assigns a task to its own cloud, it will store the task in a queue.

And when the resources and the data are ready, this task is executed. If the manager server of cloud A assigns a task to cloud B, the manager server of B first checks whether its resource availabilities can meet the requirement of this task. If so, the task will enter a queue waiting for execution. Otherwise, the manager server of B will reject the task. Before a task in the queue of a manager server is about to be executed, the manager server transfers a disk image to all the computing nodes that provide enough VMs for task execution.

Fig. 1. An example of our proposed cloud resource allocation mechanism. Heterogeneous VMs are provided by multiple clouds. And clouds are connected to the Internet via manager servers

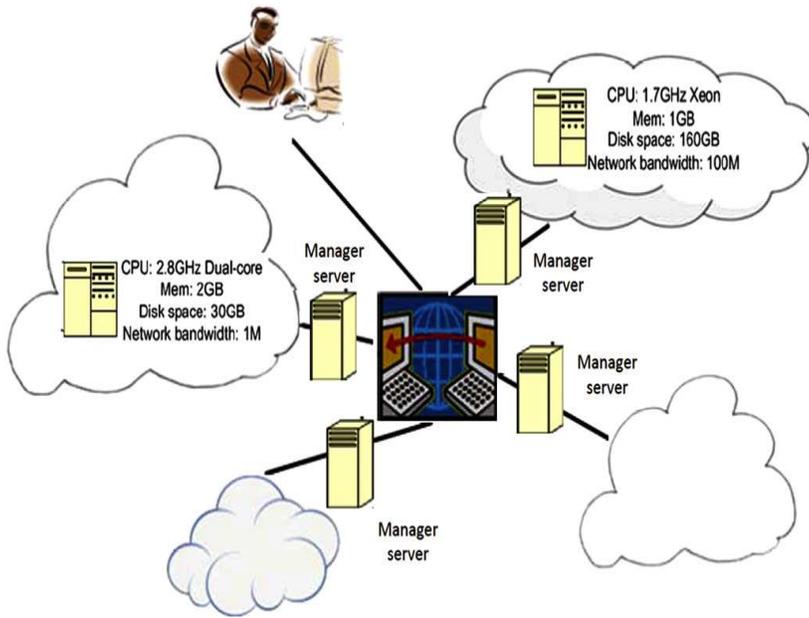
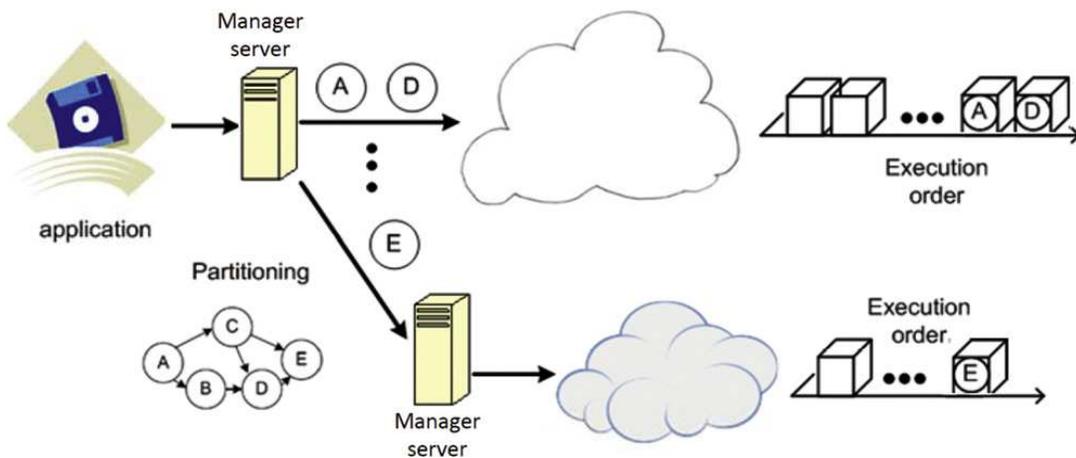


Fig. 2. When an application is submitted to the cloud system, it is partitioned, assigned, scheduled, and executed in the cloud system.



We assume that all required disk images are stored in the data center and can be transferred to any clouds as needed. We use the multicasting to transfer the image to all computing nodes within the data center. Assuming the size of this disk image is  $SI$ , we model the transfer time as  $SI/b$ , where  $b$  is the network bandwidth. When a VM finishes its part of the task, the disk image is discarded from computing nodes.

## 2.2 Resource allocation model

In cloud computing, there are two different modes of renting the computing capacities from a cloud provider.

- Advance Reservation (AR): Resources are reserved in advance. They should be available at a specific time.
- Best-effort: Resources are provisioned as soon as possible. Requests are placed in a queue.

A lease of resource is implemented as a set of VMs. And the allocated resources of a lease can be described by a tuple  $(n, m, d, b)$ , where  $n$  is number of CPUs,  $m$  is memory in megabytes,  $d$  is disk space in megabytes, and  $b$  is the network bandwidth in megabytes per second. For the AR mode, the lease also includes the required start time and the required execution time. For the best-effort and the immediate modes, the lease has information about how long the execution lasts, but not the start time of execution. The best-effort mode is supported by most of the current cloud computing platform. The Haizea, which is a resource lease manager for OpenNebula, supports the AR mode. The “map” function of “map/reduce” data-intensive applications are usually independent. Therefore, it naturally fits in the best-effort mode. However, some large scale “reduce” processes of data intensive applications may need multiple reducers. For example, a simple “wordcount” application with tens of PBs of data may need a parallel “reduce” process, in which multiple reducers combine the results of multiple mappers in parallel. Assuming there are  $N$  reducers, in the first round of parallel “reduce”, each of  $N$  reducers counts  $1/N$  results from the mappers. Then  $N/2$  reducers receive results from the other  $N/2$  reducers, and counts  $2/N$  results from the last round of reducing. It repeats  $\log_2 N + 1$  rounds. Between two rounds, reducers need to communicate with others. Therefore, an AR mode is more suitable for these data-intensive applications.

When supporting the AR tasks, it may lead to a utilization problem, where the average task waiting time is long, and machine utilization rate is low. Combining AR and best-effort in a preemptable fashion can overcome these problems. In this paper, we assume that a few of applications submitted in the cloud system are in the AR mode, while the rest of the applications are in the best-effort mode. And the applications in AR mode have higher priorities, and are able to preempt the executions of the best-effort applications.

When an AR task A needs to preempt a best-effort task B, the VMs have to suspend task B and restore the current disk image of task B in a specific disk space before the manager server transfers the disk image of tasks A to the VMs. When the task A finishes, the VMs will resume the execution of task B. We assume that there is a specific disk space in every node for storing the disk image of suspended task.

There are two kinds of AR tasks: one requires a start time in future, which is referred to as “non-zero advance notice” AR task; the other one requires to be executed as soon as possible with higher priority than the best-effort task, which is referred to as “zero advance notice” AR task. For a “zero advance notice” AR task, it will start right after the manager server makes the scheduling decision and assign it a cloud. Since our scheduling algorithms, mentioned heuristic approaches, this waiting time is negligible, compared to the execution time of task running in the cloud system.

## 2.3 Local mapping and energy consumption

From the user’s point of view, the resources in the cloud system are leased to them in the term of VMs. Meanwhile, from the cloud administrator’s point of view, the resources in the cloud system are utilized in the term of servers. A server can provide the resources of multiple VMs, and can be utilized by several tasks at the same time. One important function of the manager server of each cloud is to schedule its tasks to its server, according to the number of required VMs. Assuming there are a set of tasks  $T$  to schedule on a server  $S$ , we define the remaining workload capacity of a server  $S$  is  $C(S)$ , and the number of required VM by task  $t_i$  is  $wl(t_i)$ . The server can execute all the tasks in  $T$  only if:

$$C(S) \geq \sum_{t_i \in T} (wl(t_i)).$$

We assume servers in the cloud system work in two different modes: the active mode and the idle mode. When the server is not executing any task, it is switched to the idle mode. When tasks arrive, the server is switched back to the active mode. The server consumes much less energy in the idle mode than that in the active mode.

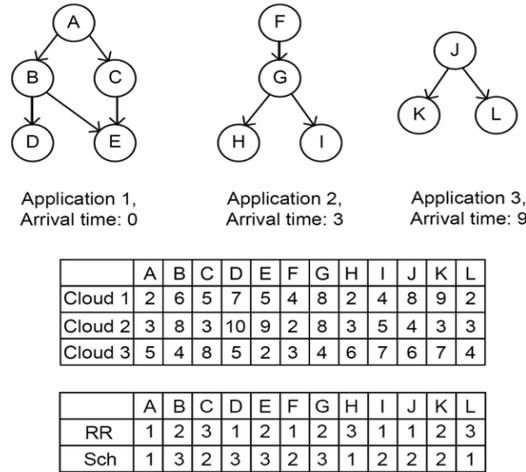
## 2.4 Application model

In this paper, we use the Directed Acyclic Graphs (DAG) to represent applications. A DAG  $T = (V, E)$  consists of a set of vertices  $V$ , each of which represents a task in the application, and a set of edges  $E$ , showing the dependences among tasks. The edge set  $E$  contains edges  $e_{ij}$  for each task  $v_i \in V$  that task  $v_j \in V$  depends on. The weight of a task represents the type of this task. Given an edge  $e_{ij}$ ,  $v_i$  is the immediate predecessor of  $v_j$ , and  $v_j$  is called the immediate successor of  $v_i$ . A task only starts after all its immediate predecessors finish. Tasks with no immediate predecessor are entry-node, and tasks without immediate successors are exit node.

Although the compute nodes from the same cloud may equip with different hardware, the manager server can treat its cloud as a homogeneous system by using the abstract compute capacity unit and the virtual machine. However, as we assumed, the VMs from different clouds may have different characteristics. So the whole cloud

system is a heterogeneous system. In order to describe the difference between VMs’ computational characteristics, we use an  $M \times N$  execution time matrix (ETM)  $E$  to indicate the execution time of  $M$  types of tasks running on  $N$  types of VMs. For example the entry  $e_{ij}$  in  $E$  indicate the required execution time of task type  $i$  when running on VM type  $j$ . We also assume that a task requires the same lease  $(n,m,d,b)$  no matter on which type of VM the task is about to run.

Fig. 3. (a) The DFG of three applications, (b) the execution time table, and (c) two different task assignments, where ‘RR’ is the round-robin approach, and ‘Sch’ is using the list scheduling.



### III. RESOURCE ALLOCATION AND TASK SCHEDULING ALGORITHM

Since the manager servers neither know when applications arrive, nor whether other manager servers receive applications, it is a dynamic scheduling problem. We propose two algorithms for the task scheduling: dynamic cloud list scheduling (DCLS) and dynamic cloud min–min scheduling (AMMS).

#### 3.1 Static resource allocation

When a manager server receives an application submission, it will first partition this application into tasks in the form of a DAG. Then a static resource allocation is generated offline. We proposed two greedy algorithms to generate the static allocation: the cloud list scheduling and the cloud min–min scheduling.

##### 3.1.1 Cloud list scheduling (CLS)

Our proposed CLS is similar to CPNT. Some definitions used in listing the task are provided as follow. The earliest start time (EST) and the latest start time (LST) of a task are shown as in the following equations.

The entry-tasks have EST equals to 0. And The LST of exit-tasks equal to their EST.

$$EST(v_i) = \max_{vm \in pred(v_i)} \{EST(vm) + AT(vm)\}$$

$$LST(v_i) = \min_{vm \in succv_i} \{LST(vm) - AT(v_i)\}$$

Because the cloud system concerned in this paper is heterogeneous, the execution times of a task on VMs of different clouds are not the same.  $AT(v_i)$  is the average execution time of task  $v_i$ . The critical node (CN) is a set of vertices in the DAG of which EST and LST are equal. Algorithm 1 shows a function forming a task list based on the priorities.

Once the list of tasks is formed, we can allocate resources to tasks in the order of this list. The task on the top of this list will be assigned to the cloud that can finish it at the earliest time. Note that the task being assigned at this moment will start execution only when all its predecessor tasks are finished and the cloud resources allocated to it are available. After assigned, this task is removed from the list. The procedure repeats until the list is empty. A static resource allocation is obtained after this assigning procedure that is shown in Algorithm 2.

Algorithm 1 Forming a task list based on the priorities

Require: A DAG, Average execution time AT of every task in the DAG.

Ensure: A list of tasks P based on priorities.

- 1: The EST of every tasks is calculated.
- 2: The LST of every tasks is calculated.
- 3: Empty list P and stack S, and pull all tasks in the list of task U.
- 4: Push the CN task into stack S in the decreasing order of their LST
- 5: while the stack S is not empty do
- 6: if top(S) has un-stacked immediate predecessors then
- 7: S ←the immediate predecessor with least LST
- 8: else
- 9: P ←top(S)
- 10: pop top(S)
- 11: end if
- 12: end while

Algorithm 2 The assigning procedure of CLS

Require: A priority-based list of tasks P, m different clouds, ETM matrix

Ensure: A static resource allocation generated by CLS

- 1: while The list P is not empty do
- 2: T = top(P)
- 3: Pull resource status information from all other manager servers
- 4: Get the earliest resource available time for T , with the consideration of the dataset transferring time response from all other manager servers
- 5: Find the cloud Cmin giving the earliest estimated finish time of T, assuming no other task preempts T
- 6: Assign task T to cloud Cmin
- 7: Remove T from P
- 8: end while

### 3.1.2 Cloud min–min scheduling (CMMS)

Min–min is another popular greedy algorithm. The original min–min algorithm does not consider the dependences among tasks. So in the dynamic min–min algorithm used in this paper, we need to update the mappable task set in every scheduling step to maintain the task dependences. Tasks in the mappable task set are the tasks whose predecessor tasks are all assigned. Algorithm 3 shows the pseudo codes of the DMMS algorithm.

### 3.1.3 Energy-aware local mapping

A manager server uses a slot table to record execution schedules of all resources, i.e., servers, in its cloud. When an AR task is assigned to a cloud, the manager server of this cloud will first check the resource availability in this cloud. Since AR tasks can preempt best-effort tasks, the only case where an AR task is rejected is that most of the resources are reserved by some other AR tasks at the required time, no enough resources left for this task. If the AR task is not rejected, which means there are enough resources for this task, a set of servers will be reserved by this task, using the algorithm shown in Algorithm 4. The time slots for transferring the disk image of the AR task and the task execution are reserved in the slot tables of those servers. The time slots for storing and reloading the disk image of the preempted task are also reserved if preemption happens.

When a best-effort task arrives, the manager server will put it in the execution queue. Every time when there are enough VMs for the task on the top of the queue, a set of servers are selected by the algorithm shown in Alg. 5. And the manager server also updates the time slot table of those servers

Algorithm 3 Cloud min–min scheduling (CMMS)

Require: A set of tasks, m different clouds, ETM matrix.

Ensure: A schedule generated by CMMS.

- 1: Form a mappable task set P.
- 2: while there are tasks not assigned do
- 3: Update mappable task set P.
- 4: for i: task  $v_i \in P$  do
- 5: Pull resource status information from all other manager servers.
- 6: Get the earliest resource available time, with the consideration of the dataset transferring time response from all other manager servers.
- 7: Find the cloud Cmin( $v_i$ ) giving the earliest finish time of  $v_i$ , assuming no other task preempts  $v_i$
- 8: end for
- 9: Find the task-cloud pair( $v_k, Cmin(v_k)$ ) with the earliest finish time in the pairs generated in for-loop.

10: Assign task  $vk$  to cloud  $Dmin(vk)$ .  
 11: Remove  $vk$  from  $P$ .  
 12: Update the mappable task set  $P$   
 13: end while

Algorithm 4 Energy-aware local mapping for AR tasks

Require: A set of AR tasks  $T$ , which require to start at the same time. A set of servers  $S$ .

Ensure: A local mapping

1: for  $t_i \in T$  do  
 2: Calculate  $wlm(t_i)$   
 3: if  $wl(t_i) - wlm(t_i) < \sum_{s_i \in idle(C(s_i))}$  then  
 4: Schedule  $wl(t_i) - wlm(t_i)$  to the idle servers  
 5: else  
 6: First schedule a part of  $wl(t_i) - wlm(t_i)$  to the idle servers  
 7: Schedule the rest of  $wl(t_i) - wlm(t_i)$  to the active servers, preempting the best-effort tasks  
 8: end if  
 9: end for  
 10: Sort tasks in  $T$  in the descending order of marginal workload, form list  $L_d$ .  
 11: Sort tasks in  $T$  in the ascending order of marginal workload, form list  $L_a$   
 12: while  $T$  is not empty do  
 13:  $ta = top(L_d)$   
 14: if there exists a server  $j: C(j) = wlm(ta)$  then  
 15: Schedule the  $wlm(ta)$  to server  $j$   
 16: end if  
 17:  $sa = max_{s_i \in S}(C(s_i))$ .  
 18: Schedule  $ta$  to  $sa$ , delete  $ta$  from  $T$ ,  $L_d$ , and  $L_a$   
 19: for  $k: tk \in L_a$  do  
 20: if  $C(sa) > 0$  and  $C(sa) \geq wlm(tk)$  then  
 21: Schedule  $tk$  to  $sa$ , delete  $tk$  from  $T$ ,  $L_d$ , and  $L_a$   
 22: else  
 23: Break  
 24: end if  
 25: end for  
 26: end while

The objectives of Algorithms 4 and 5 are to minimize the number of active servers as well as the total energy consumption of the cloud. When every active server is fully utilized, the required number of active servers is minimized. When task  $t_i$  is assigned to cloud  $j$ , we define the marginal workload of this task as:

$$wlm(t_i) = wl(t_i) \bmod C(S_j) \quad (4)$$

where  $S_j$  represents the kind server in cloud  $j$ , and  $C(S_j)$  is the workload capacity of server  $S_j$ . To find the optimal local mapping,

Algorithm 5 Energy-aware local mapping for best-effort task

Require: A set of best-effort tasks  $T$ , which can start at the same time. A set of servers  $S$

Ensure: A local mapping

1: for  $t_i \in T$  do  
 2: Calculate  $wlm(t_i)$ .  
 3: Schedule  $wl(t_i) - wlm(t_i)$  to the idle servers.  
 4: end for  
 5: Form a set of active servers  $S_g$  that  $C(s_i) > 0, \forall s_i \in S_g$ .  
 6: Sort tasks in  $T$  in the descending order of marginal workload, form list  $L_d$   
 7: Sort tasks in  $T$  in the ascending order of marginal workload, form list  $L_a$   
 8: while  $T$  is not empty do  
 9:  $ta = top(L_d)$   
 10: if there exists a server  $j$  in  $S_g : C(j) = wlm(ta)$  then  
 11: Schedule the  $wlm(ta)$  to server  $j$   
 12: end if  
 13:  $sa = max_{s_i \in S_g}(C(s_i))$   
 14: if  $C(sa) < wlm(ta)$  then  
 15:  $sa = anyidleserver$   
 16: end if

```

17: Schedule ta to sa, delete ta from T , Ld, and La
18: for k: tk ∈ La do
19: if C(sa) > 0 and C(sa) ≥ wlm(tk) then
20: Schedule tk to sa, delete tk from T , Ld, and La
21: else
22: Break
23: end if
24: end for
25: end while
    
```

We group all the tasks that can be executed simultaneously, and sort them in the descending order of their marginal workloads. For each of the large marginal workload task, we try to find some small marginal workload tasks to fill the gap and schedule them on a server.

### 3.1.4 Feedback information

In the two static scheduling algorithms presented above, the objective functions when making decision about assigning a certain task is the earliest estimated finish time of this task. The estimated finish time of task *i* running on cloud *j*,  $\tau_{i,j}$ , is as below:

$$\tau_{i,j} = ERAT_{i,j} + SI/b + ETM_{i,j}. \quad (5)$$

SI is the size of this disk image, *b* is the network bandwidth.  $ERAT_{i,j}$  is the earliest resource available time based the information from the pull operation. It is also based on the current task queue of cloud *j* and the schedule of execution order. But the estimated finish time from (5) may not be accurate. For example, as shown in Fig. 5(a), we assume there are three clouds in the system. The manager server of cloud A needs to assign a best-effort task *i* to a cloud. According to Eq. (5), cloud C has the smallest  $\tau$ . So manager server A transfers task *i* to cloud C. Then manager server of cloud B needs to assign an AR task *j* to a cloud. Task *j* needs to reserve the resource at 8. Cloud C has the smallest  $\tau$  again. manager server B transfers task *j* to cloud C. Since task *j* needs to start when *i* is not done, task *j* preempts task *i* at time 8, as shown in Fig. 6. In this case, the actual finish time of task *i* is not the same as expected.

In order to reduce the impacts of this kind of delays, we use a feedback factor in computing the estimated finish time. As discussed previously in this paper, we assume once a task is done, the cloud will push the resource status information to the original cloud. Again, using our example in Fig. 5, when task *i* is done at time  $Tact\_fin(=14)$ , manager server C informs manager server A that task *i* is done. With this information, the manager server A can compute the actual execution time  $1\tau_{i,j}$  of task *i* on cloud *j*:

$$\Delta\tau_{i,j} = Tact\_fin - ERAT_{i,j}.$$

And the feedback factor  $fd_j$  of cloud *j* is :

$$fd_j = \alpha \times \frac{\Delta\tau_{i,j} - SI/b - ETM_{i,j}}{SI/b + ETM_{i,j}}$$

$\alpha$  is a constant between 0 and 1. So a feedback estimated earliest finish time  $\tau_{fdi,j}$  of task *i* running on cloud *j* is as follows:

$$\tau_{fdi,j} = ERAT_{i,j} + (1 + fd_j) \times (SI/b + ETM_{i,j}).$$

In our proposed dynamic cloud list scheduling (DCLS) and dynamic cloud min–min scheduling (DCMMS), every manager server stores feedback factors of all clouds. Once a manager server is informed that a task originally from it is done, it will update the value of the feedback factor of the task-executing cloud. For instance, in the previous example, when cloud C finishes task *i* and informs that to the manager server of cloud A, this manager server will update its copy of feedback factor of cloud C. When the next task *k* is considered for assignment, the  $\tau_{fdk,C}$  is computed with the new feedback factor and used as objective function.

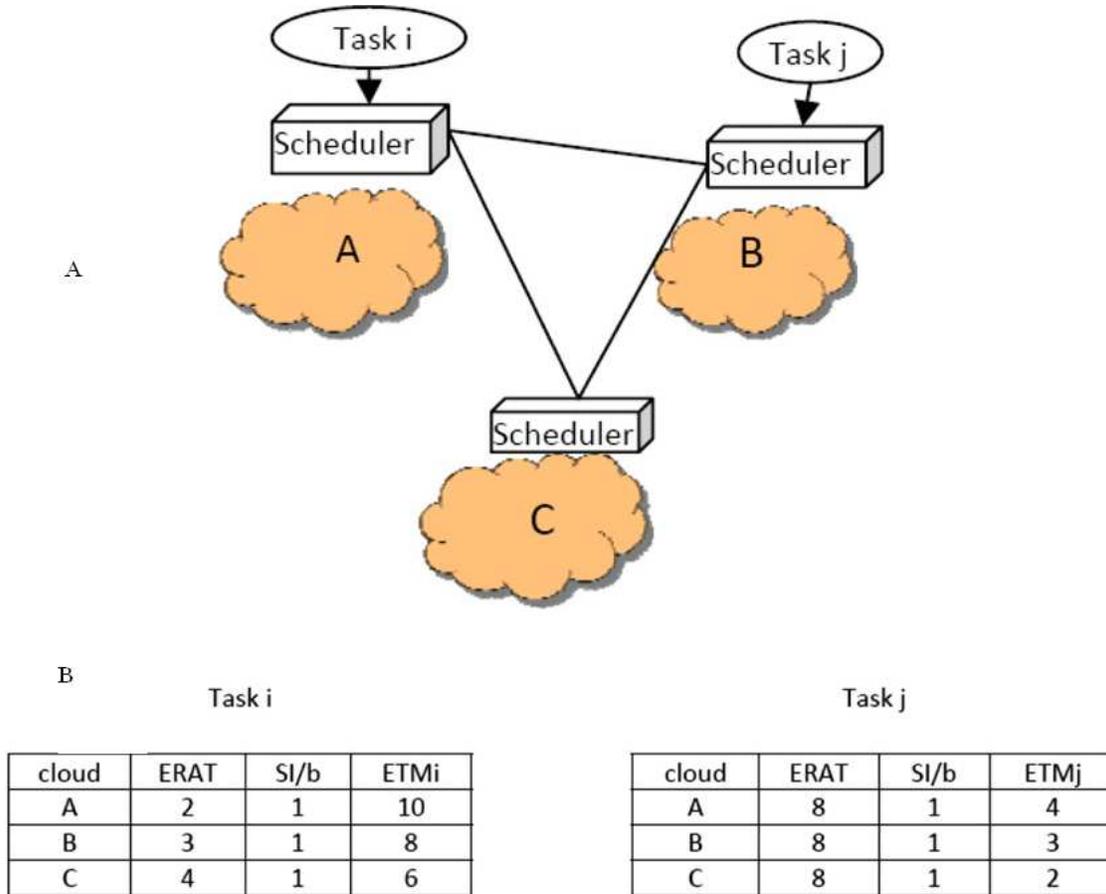


Fig. 5. Example of resource contention. (a) Two tasks are submitted to a heterogeneous clouds system. (b) The earliest resource available times (ERAT), the image transferring time (SI/b), and the execution time (EMT) of two tasks on different clouds

#### IV. EXPERIMENTAL RESULTS

##### 4.1 Experiment setup

We evaluate the performance of our dynamic algorithms through our own written simulation environment that acts like the IaaS cloud system. We stimulate workloads with job traces from the Parallel Workloads Archive. We select three different job traces: LLNL-Thunder, LLNL-Atlas, and LLNL-uBGL. For each job tracer, we extract four values: the job ID, the job start time, the job end time, and the node list. However, job traces from the Parallel Workloads Archive do not include information about data dependencies. To simulate data dependencies, we first sort jobs by their start time. Then we group up to 64 adjacent jobs as one application, represented by a randomly generated DAG. Table 1 shows how we translate those values from job traces to the parameter we use in our application model. Note that we map the earliest job start time in an application as the arrival time of this application, since there is no record about job arrival time in these job traces. There are three data center in our simulation: (1) 1024 node cluster, with 4 Intel IA-64 1.4 GHz Itanium processors, 8 GB memory, and 185 GB disk space per node; (2) 1152 node cluster, with 8 AMD Opteron 2.4 GHz processors, 16 GB memory, and 185 GB disk space per node; (3) 2048 processors BlueGene/L system with 512 MB memory, 80 GB memory. We select these three data center configuration based on the clusters where LLNL-Thunder, LLNL-Atlas, and LLNL-uBGL job traces were obtained. Based on the information in [24], we compare the computational power of these three data center in Table 2. With the normalized performance per core, we can get the execution time of all tasks on three different

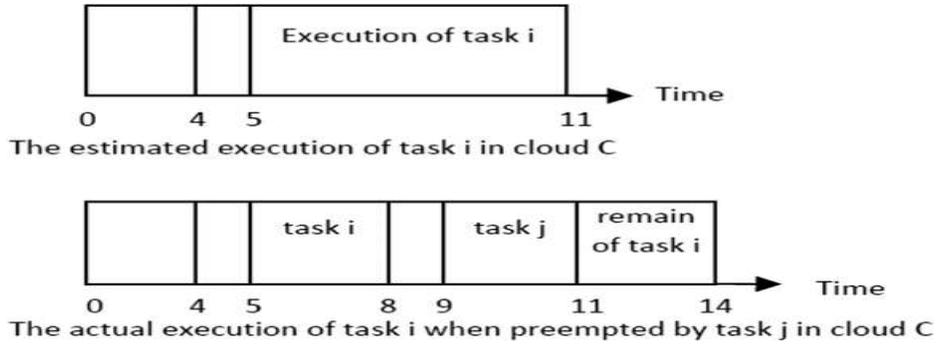


Fig. 6. The estimated and the actual execution order of the cloud C.

**Table 1**  
The mapping of job traces to applications.

Parameter in our model	Values in job traces
Task id	Job ID
Application arrival time	Min(job start time)
Task execution time	Job end time–job start time
# of CPU required by a task	Length(node list) * cpu per node

**Table 2**  
Comparison of three data center. The job trace LLNL-uBGL was obtained from a small uBGL, which has the same single core performance as the one shown in this table.

Data center	Peak performance (TFLOP/s)	Number of CPUs	Normalized performance per core
Thunder	23	4 096	1
Altas	44.2	9216	0.85
uBGL(big)	229.4	81920	0.50

**Table 3**  
Feedback improvements in different cases.

Arrival gap reduces times	DLS	FDLS( $\alpha = 1$ )	Feedback improve. (%)	DMMS	FDMMS ( $\alpha = 1$ )	Feedback improve. (%)
1	237.82	253.59	-6.63	206.31	223.47	-8.32
20	309.35	286.55	7.37	262.66	255.44	2.75
40	445.74	397.15	10.9	385.48	336.52	12.7
60	525.32	420.83	19.89	448.04	343.60	23.31
80	729.56	537.28	26.36	648.37	440.05	32.13
100	981.41	680.22	30.69	844.33	504.66	40.23

**Table 4**

Average application execution time with various percentages of AR applications in the loose situation ( $\alpha = 0.8$ ).

	0%	20%	50%	80%	100%
FCFS	1	1	1	1	1
DCLS	0.81	0.75	0.61	0.55	0.49
DCMMS	0.77	0.56	0.52	0.46	0.44

In order to find out the relationship between resource data centers. Among these applications, 20% applications are in the AR modes, while the rest are in the best-effort modes. We assume the bandwidth between two data centers are 1 Gbps, the bandwidth of nodes inside the data center are 4 GBps, and the size of every dataset is 1 TB. We run these three jobs trace separately in our simulation. We set the arrival of applications in two different ways. In the first way, we use the earliest start time of an application in the original job trace as the arrival time of this application. We also set the required start time of an AR application as a random start time no later than 30 min after it arrives. In most of the cast, applications do not need to contend resources in this setting. We call this a loose situation. In the other way, we set the arrival time of applications close to each other. In this setting, we reduce the arrival time gap between two adjacent application by 100 time. It means that applications usually need to wait for resources in cloud. We call this a tight situation. In both these two setting, we tunes the constant  $\alpha$  to show how the dynamic procedure impacts the average application execution time. We define the execution time of an application as the time elapses from the application is submitted to the application is finished.

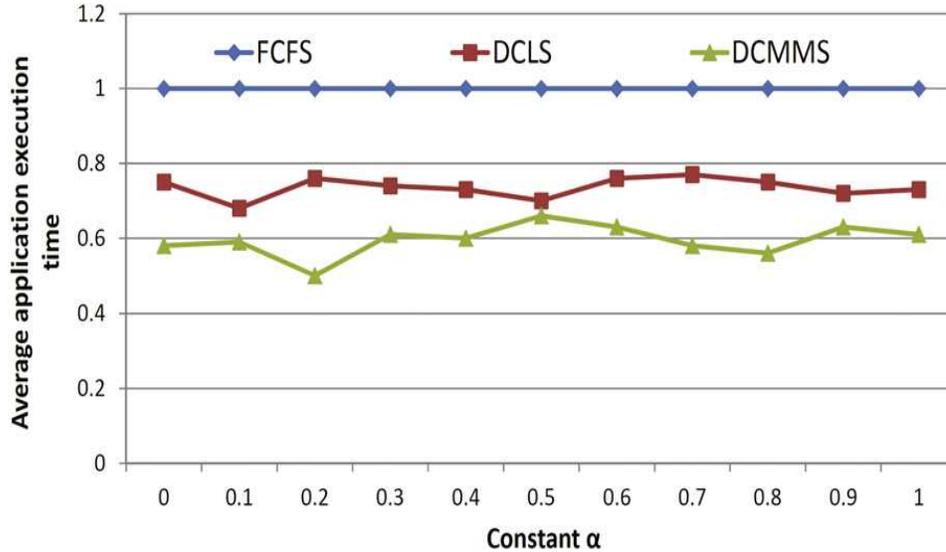


Fig. 7. Average application execution time in the loose situation.

**4.2 Result**

Fig. 7 shows the average application execution time in the loose situation. We compare our two dynamic algorithms with the First-Come-First-Serve (FCFS) algorithm. We find out that the DCMMS algorithm has the shorter average execution time. And the dynamic procedure with updated information does not impact the application execution time significantly. The reason the dynamic procedure do not has a significant impact on the application execution time is that the resource contention is not significant in the loose situation. Most of the resource contentions occur when an AR application preempts a best-effort application. So the estimated finish time of an application is usually close to the actual finish time, which limits the effect of the dynamic procedure. And the manager server does not call the dynamic procedure in most of the cases.

Fig. 8 shows that DCMMS still outperforms DCLS and FCFS. And the dynamic procedure with updated information works more significantly in the tight situation than it does in the loose situation. Because the resource contentions are fiercer in tight situation, the actual finish time of a task is often later than estimated finish time. And the best-effort task is more likely preempted some AR tasks. The dynamic procedure can avoid

tasks gathering in some fast clouds. We believe that the dynamic procedure works even better in a homogeneous cloud system, in which every task runs faster in some kinds of VMs than in some other kinds.

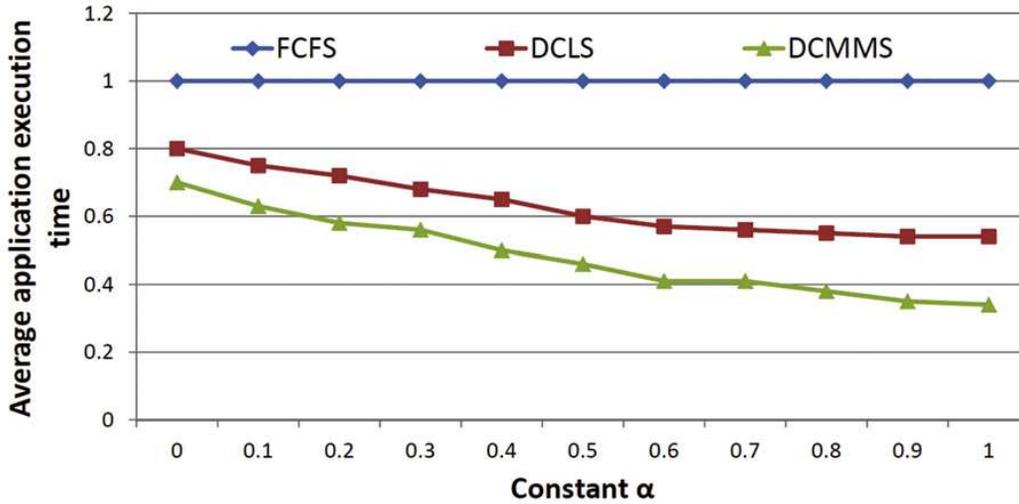


Fig. 8. Average application execution time in the tight situation

**Table 5**

Average application execution time with various percentages of AR applications in the tight situation ( $\alpha = 0.8$ ).

	0%	20%	50%	80%	100%
FCFS	1	1	1	1	1
DCLS	0.63	0.55	0.49	0.43	0.38
DCMMS	0.51	0.38	0.32	0.30	0.27

In order to find out the relationship between resource contention and feedback improvement, we increase the resource contention by reducing the arrival time gap between two adjacent applications. We reduce this arrival time gap by 20, 40, 60, 80, and 100 times, respectively. In the setting with original arrival time gap, an application usually come after the former application is done. Resource contention is light. And when arrival time gaps are reduced by 100 times, it means during the execution of an application, there may be multiple new applications arriving. Resource contention is heavy in this case. As shown in Table 3, the improvement caused by feedback procedure increases as the resource contention become heavier.

We also test our proposed algorithms in setups with various percentages of AR applications, as shown in Tables 4 and 5. The values in the first row represent how many applications are set as the AR applications. The values in the second, the third, and the fourth row are the average application execution time, normalized by the corresponding execution time with the FCFS algorithm. From these two tables, we can observe that higher percentage of AR applications leads to a better improvement of the DLS and the DCMMS algorithm, compared to the FCFS algorithm, in both the loose situation and the tight situation. The reason is that more AR applications cause longer delays of the best-effort applications. By using the feedback information, our DLS and DCMMS can reduce workload unbalance, which is the major drawback of the FCFS algorithm.

Furthermore, we compare the energy consumption of three algorithms, shown in Figs. 9 and 10. Both DCLS and DCMMS can reduce energy consumption compared to the FCFS algorithm. In addition, our energy-aware local mapping further reduce the energy consumption significantly, in all three algorithms.

In the future work, we will evaluate our proposed mechanism in existing simulators, so that results can be reproduced easier by other researchers. In addition, we will investigate the implementation of our design in the real-world cloud computing platform. A reasonable way to achieve this goal is to combine our design with the Hadoop platform. The multi-cloud scheduling mechanism and algorithms in our design can be used on the top of the Hadoop platform, distributing applications in the federated multicloud platform. When a give task is assigned to a cloud, the Hadoop will be used to distribute tasks to multiple nodes. And our proposed energy-aware local mapping design can be implemented in the Hadoop Distributed File System, which enables the ‘‘rack awareness’’ feature for data locality inside the data center.

## V. CONCLUSION

The cloud computing is emerging with rapidly growing customer demands. In case of significant client demands, it may be necessary to share workloads among multiple data centers, or even multiple cloud providers. The workload sharing is able to enlarge the resource pool and provide even more flexible and cheaper resources. In this paper, we present a resource optimization mechanism for preemptable applications in federated heterogeneous cloud systems. We also propose two novel online dynamic scheduling algorithms, DCLS and DCMMS, for this resource allocation mechanism. Experimental results show that the DCMMS outperforms DCLS and FCFS. And the dynamic procedure with updated information provides significant improvement in the fierce resource contention situation. The energy-aware local mapping in our dynamic scheduling algorithms can significantly reduce the energy consumptions in the federated cloud system.

## REFERENCES

- [1] Amazon AWS, <http://aws.amazon.com/>.
- [2] GoGrid, <http://www.gogrid.com/>.
- [3] RackSpace, <http://www.rackspacecloud.com/>.
- [4] Microsoft cloud, <http://www.microsoft.com/en-us/cloud/>.
- [5] IBM cloud, <http://www.ibm.com/ibm/cloud/>.
- [6] Google apps, <http://www.google.com/apps/intl/en/business/index.html>.
- [7] Eucalyptus, <http://www.eucalyptus.com/>.
- [8] RESERVOIR, [www.reservoir-fp7.eu](http://www.reservoir-fp7.eu).
- [9] E. Walker, J. Gardner, V. Litvin, E. Turner, Dynamic virtual clusters in a grid site manager, in: Proceedings of Challenges of Large Applications in Distributed Environments, Paris, France, 2006, pp. 95–103.
- [10] X. Wang, J. Zhang, H. Liao, L. Zha, Dynamic split model of resource utilization in map reduce, in: International Workshop on Data Intensive Computing in the Clouds, 2011, pp. 1–10.
- [11] M. Aron, P. Druschel, W. Zwaenepoel, Cluster reserves: a mechanism for resource management in cluster-based network servers, in: Proceedings of the ACM Sigmetrics, Santa Clara, California, USA, 2000.
- [12] A.I. Avetisyan, R. Campbell, M.T. Gupta, I. Heath, S.Y. Ko, G.R. Ganger, et al., Open Cirrus a global cloud computing testbed, *IEEE Computer* 43 (4) (2010) 35–43.
- [13] D. Cearley, G. Phifer, Case studies in cloud computing, <http://www.gartner>.