

Original Research Article

AN IMPROVED RASA FOR LOAD-BALANCING IN CLOUD COMPUTING

ABSTRACT

Cloud is specifically known to have difficulty in managing resource usage during task scheduling, this is an innate from distributed computing and virtualization. Efficient procedures are therefore necessary to achieve maximum resource utilization at a minimized cost. This study implemented a load balancing scheme called Improved Resource Aware Scheduling Algorithm (I-RASA) for resource provisioning to cloud users on a pay-as-you-go basis using CloudSim 3.0.3 package. I-RASA was compared with recent load balancing algorithms and the results showed that it performed better than Max-min and RASA. However, it sometimes outperforms or on equal balance with Improved Max-Min when using makespan, flow time, throughput, and resource utilization as the performance metrics.

Keywords: Improved RASA, Load-balancing, Cloud Computing, Resource Utilization.

1. INTRODUCTION

Cloud Computing consists of a network of infrastructure which allows enterprises to achieve more efficient use of their IT hardware, software, and other services investments. These infrastructures are available to users as a utility on an on-demand basis and are charged proportionally to the amount of resources consumed by the users. Cloud is a heterogeneous pool of resources, and to enable one to have access to these resources; there must be a Service Level Agreement (SLA) from the Cloud Provider [1]. Efficient resource utilization can be optimized by the Cloud Service Providers (CSP) via a better scheduling algorithm in cloud computing [2]. Better scheduling could be achieved by breaking down the physical barrier that is essential to isolated system and automate it as a single entity [3].

In Cloud Computing, effective task scheduling needs an infrastructure from the cloud provider. The scheduling of task, jobs are mapped on available resources and submitted to a cloud environment in such a way that the total response time and the makespan are minimized [4]. Makespan is the maximum time taken by any of the computing resources assigned a set of tasks to complete execution. Efficient resource utilization can only be guaranteed when the workload was efficiently shared among the resources [5]. As the cloud computing environment continues to enjoy massive migration by enterprises, it would be important to improve the distribution of tasks among resources. Cloud service providers are greatly concerned about addressing the challenge of ensuring that Virtual Machines (VMs) are not overloaded or underloaded. This optimization problem has been gaining attention in recent times.

Muthusamy and Chandran [6] proposed an artificial bee foraging optimization for load balancing in the cloud. The study adopted a preemptive task scheduling approach to minimize the response and execution time. This showed a significant improvement in QoS metric in comparison to the Honey Bee Based (HBB) load balancing.

In a bid to enhance machine performance for balanced sharing of load, maximize virtual machine throughput and optimize the waiting time of tasks, Jena, Das, and Kabat [7] hybridized Modified Particle Swarm Optimization (MPSO) and improved Q-learning. The experimental analysis results showed that the hybrid QMPSO outperformed MPSO and Q-learning.

In [8], an adaptive Starvation Threshold Load Balancing (STLB) algorithm was proposed for load balancing. The objectives of the study were to minimize response time and migration cost, and maximize server utilization rate. The key feature used in the algorithm was not invoked until at least one of the VM is close to starvation; this lowers the number of migration. The proposed algorithm was compared to the HBB load balancing algorithm in terms of makespan, average response time, average idle time, and number of tasks migrated, the STLB showed significant improvement.

In [9], a hybrid algorithm that integrated the Elephant Herding Optimization (EHO) into the Grey Wolf Optimizer (GWO) was proposed and its performance was compared to some existing load balancing algorithms such as Constraint Measure (CMBLB), Fractional Dragonfly, EHO, and GWO load balancing algorithms using makespan and minimum load as metrics. The makespan of EHGWO algorithm was better than EHO and GWO, and it has the lowest minimum load value compared with CMBLB and Fractional Dragonfly algorithms.

Quadri and Ravi [10] developed a new algorithm to reduce the overall time taken to complete a set of assigned tasks and ensure a balanced distribution of tasks to computing resources. The two objectives (minimize makespan and maximize utility) were achieved by applying the Weighted Sum Method (WSM), which was one of the multi-criterion decision-making methods. The percentage deviation (Coefficient of Variation) of the resources with maximum and minimum execution time of all tasks from the mean were measured. The result showed that the proposed algorithm has lower makespan than

scheduling algorithms such as NHTBS, Opportunistic Load Balancing (OLB), MET and RASA. In terms of utilization only max-min scheduling technique could equal its performance.

According to [11], many algorithms have been used in the past; amongst them are Max-min, Min-min, and RASA, which are very popular during scheduling of tasks on resources. The study focused on max-min algorithm, based on completion time of processed tasks. This algorithm was improved upon using expected execution time [11]. The results showed that Improved Max-min (I-MM) performed better by 2% in completion time than RASA.

This research focused on improving RASA algorithm to give a better performance on load-balancing in cloud computing, and develop a model that can choose between Improved Max-min and I-RASA based on the cloudlet size at runtime.

2. META-TASK

CloudSim works on the Linux operating system, which has different schedulers in operation since the advent of Linux 2.4. These schedulers include Deadline scheduler, the Anticipatory scheduler, Complete Fair Queuing (CFQ) scheduler, and Noop scheduler. During the implementation of this research, the Anticipatory scheduler is no longer in operation; it was found that CFQ was a modified version of the Anticipatory scheduler. In other words, from Linux version 2.6, the Anticipatory scheduler was not in use, leaving only Deadline, CFQ, and Noop in Fedora Linux operating system, which was the environment used for this research. All the schedulers can perform merging and sorting process except the Noop scheduler, which does not permit sorting. Hence, it is referred to as the basic scheduler. The operating system used was tuned to CFQ because aside that tasks were put in batches, it also overcomes deceptive idleness [12]. Report also shows that the Apache web server achieved up to 71% more throughput from using this modified anticipatory scheduler [13].

Scheduling algorithms are with many policies but could be subdivided into immediate and batch scheduling, preemptive and non-preemptive scheduling, static and dynamic scheduling, and so on [14], [15], [16]. In Immediate mode, tasks were scheduled using First Come First Serve (FCFS) in the computing environment, while in the batch mode, tasks were grouped into a batch; which means that, a set of meta-tasks would be assigned at a mapped out time depending on the scheduler's algorithm [15]. It is the algorithm that determines how the load balancing of the resource usage. In the next section, different load balancing algorithms were described based on policies earlier discussed.

3. OVERVIEW OF LOAD-BALANCING

Different techniques have been used to improve the performance and resource usage on task scheduling, quality of service, load balancing, and resource utilization. In [17], it was suggested that load balancing in cloud avoids overloading/underloading of virtual machines, which itself is an obstacle in cloud computing, thereby making it a requisite for researchers to develop a suitable load balancer for parallel and distributed cloud environments.

Load balancing is a method that enables frameworks on resources by giving a Maximum throughput with the least response time [7]. In cloud computing, the burden (overloading/underloading) could occur in the resources used in the datacentres. The altering of this burden was a way to divide the action between all servers (in the case of more than one datacentre) or resources (say, virtual machines), so that job could be sent and response got quickly while the stack was on modification. Various load balancing algorithms that give better throughput with quick response time in cloud condition do exist [18], but, each of them has favourable circumstance [19, 20, 21]. These include:

3.1. Static Algorithm (SA)

Static Algorithms, as shown in Fig. 1, are used where the load was of low variations. This algorithm needs a prior knowledge of server resources for the processors to perform better and this was determined at the beginning of the implementation [22]. The major limitation of Static Load Balancing Algorithm is that the load balancing tasks only work after being created.

```

Start
create different classes of job
assign next job to the belong class
Sort the class tasks based on weight of the execution time
for all classes
    for all task in selected class
        Assign task with min-weight to the next resource for processing
        remove the task from the selected class list
        next task
    end
    next class
end
End

```

Fig 1: Load Balancing Static Algorithm

3.2 Dynamic Algorithm (DA)

Dynamic Algorithm searches the lightest server resource, and gives it a preference for load balancing [22]. The current state of the machine was used to control the load. This was explained in Fig 2.

```

Start
create different classes of job , C
assign next job to the belong class
Sort the class tasks based on weight of the execution time
For all Resource, R
    for (i = 0; i < |C|; i++) // |C| is the total number of classes
        Assign task with min-weight in Ci on R
        remove the task from the selected class list
        Next class
    end
    Next Resource
end
End

```

Fig 2: Load Balancing Dynamic Algorithm

3.3 Round Robin Algorithm (RRA)

This algorithm assigns tasks to server resources. Its mode of operation uses the FCFS algorithm when the quantum time is high firstly, but once it reduces, tasks were assigned on a random basis in round- robin [23]. The round-robin method circularly assigned tasks without defining any priority. All the processes have different loading times. Some resources might be heavily loaded, while others remain under-utilized. The algorithm was shown in Fig 3.

```

Start
create different classes of job , C
assign next job to the belong class
Sort the class tasks based on weight of the execution time
For all Resource, R
    for (i = 0; i < |C|; i++) // |C| is the total number of classes
        if(task with min-weight in Ci < quantum_time)
            Assign the task in Ci on R
            remove the task from the selected class list
        else
            Process the task for based on quantum_time
            Update the task burst_time
            Relocate the task to the end of the class list
        end
        Next class
    end
    Next Resource
end
End

```

Fig 3: Round-robin Algorithm

3.4 Opportunistic Load Balancing (OLB) algorithm

This algorithm keeps each server resource busy, as shown in Fig 4, without considering machines' current workload. Irrespective of the current workload on each of the resources, OLB distributes all the unfinished tasks to them randomly [24].

```

Start
create different classes of job
assign next job to the belong class
Sort the class tasks based on weight of the execution time

 $N = \sum_{i=1}^{|C|} |T_{C_i}|$ 
Counter=0;
Do
    i = Random (1, |C|)
    for all tasks in selected class  $C_i$ 
        Assign task with min-weight to the next resource for processing
        remove the task from the selected class list
        Counter++;
    next task
end
While (Counter < N)
End

```

Fig 4: Opportunistic Load Balancing Algorithm

3.5 Minimum to Minimum (Min-Min) algorithm

The concept of Min-min algorithms in Load Balancing is to assign tasks with minimum completion time first for execution on resource with minimum execution time [25]. This procedure continues until all were mapped. This algorithm, as shown in Fig 5, seems to be the fastest in a situation where many smaller tasks are more than larger ones.

```

Start
Create different classes of job, C
Assign next job to the belong class
Sort the class tasks based on weight of the execution time
//R is the resources
for (i = 0; i < |C|; i++) // |C| is the total number of classes
    Assign task with min-weight in  $C_i$  on  $\min (R_{1_{\text{executiontime}}}, \dots, R_{|R|_{\text{executiontime}}})$  resource
    Remove the task from the selected class list
Next class
end
End

```

Fig 5: Opportunistic Load Balancing Algorithm

3.6 Maximum to Minimum (Max-Min) Algorithm

Maximum-Min Load Balancing Algorithm is similar to Min-min load balancing algorithm, still the difference is that the task with maximum completion time was selected after searching and assigned to the machine (resource) with minimum execution time [26]. The execution time of all tasks were updated, and the assigned task was removed from the list.

```

Start
Create different classes of job, C
Assign next job to the belong class
Sort the class tasks based on Max-weight of the execution time
//R is the resources
for (i = 0; i < |C|; i++) // |C| is the total number of classes
    Assign task with Max-weight in  $C_i$  on  $\min (R_{1_{\text{executiontime}}}, \dots, R_{|R|_{\text{executiontime}}})$  resource
    Remove the task from the selected class list
Next class
end
End

```

Fig 6: Opportunistic Load Balancing Algorithm

3.7 Resource Aware Scheduling Algorithm (RASA)

This is a combination of both Max-min and Min-min algorithms. In RASA, the appraisal of the completion time for each task on available resources was calculated, after which the Max-min and Min-min algorithms were applied alternatively, as shown in Fig 6, thereby making use of the advantage of both algorithms and avoiding their drawbacks [4]. RASA executes small tasks to avoid delays in large ones. It also supports simultaneous executions of large and small tasks.

```

Start
Create different classes of job, C
Assign next job to the belong class
Sort the class tasks based on Max-weight, Min-Weight of the execution time interchangeably
//R is the resources
for (i = 0; i < |C|; i++) // |C| is the total number of classes
    Assign next task in Ci on min (R1completiontime, ..., R|R|completiontime) resource
    Remove the task from the selected class list
    Next class
end
End

```

Fig 7: Resource Aware Scheduling Algorithm

3.8 Improved Max-Min

The basic operations of improved Max-min (I-Max-Min) was to assign task with maximum execution time to a resource that has minimum completion time. The original Max-min assigned task with maximum completion time to resources with minimum execution time [27].

```

Start
Create different classes of job, C
Assign next job to the belong class
Sort the class tasks based on Max-weight of the execution time
//R is the resources
for (i = 0; i < |C|; i++) // |C| is the total number of classes
    Assign next task in Ci on min (R1completiontime, ..., R|R|completiontime) resource
    Remove the task from the selected class list
    Next class
end
End

```

Fig 8: Improved Max-Min Algorithm

4. RESOURCE PROVISIONING

The task of mapping resources to different entities in cloud on-demand that is pay-as-you-go basis is known as resource provisioning. Resources were allocated in cloud so that the processing elements (resources) are not overloaded and that none is undergoing wastage either. Resources mapping in cloud entities were done in two levels:

4.1 Host

Host, in cloud computing, can contain more instances of VM. The VM were mapped to a single host subject for availability and capabilities. The Host then assigns processing cores to VM based on the provisioning policy, and this defines the basis of allocating processing cores to VM. The allocation policy ensures that the critical characteristics of the Host and VM do not mismatch.

4.2 Cloudlet or task mapping onto VM

Cloudlets were executed on VM, and each requires a certain amount of processing power for their completion. VM provides this processing power to the task(s) mapped on it. These tasks were mapped on VM based on their configuration and availability.

4.3 Task scheduling policy

Tasks were scheduled after the resources had been allocated to the cloud entities. These activities allow multiprogramming capabilities in a cloud environment and were enabled in two modes: Space shared and Time shared policies.

In Space Shared policy at VM level, one task can be scheduled to a virtual machine at a time and when it was completed, another task is scheduled to the virtual machine. This policy behaves same as the First Come First Serve (FCFS) scheduling algorithm [28]. The algorithm of space shared policy is as follows:

- Step 1: Tasks are arranged in a queue.
- Step 2: First task is scheduled on the given virtual machine.
- Step 3: When first task is completed it assigns the next task from the queue.
- Step 4: If queue is empty it checks for new tasks. Step 5: Then repeat Step 1.
- Step 6: End.

Same algorithm is applicable for both Host level and VM level scheduling.

In Time Shared policy at Host level, virtual machines are scheduled on the CPU cores simultaneously amongst the VM, while at VM level, the scheduling policy schedules all the tasks on the VM at the same time, and this is done among all tasks. This algorithm is the same like the Round Robin (RR) scheduling algorithm [29]. The algorithm of time shared policy can be represented as follows:

- Step 1: All the tasks are arranged in a queue.
- Step 2: Then schedule the tasks simultaneously on the virtual machine.
- Step 3: When queue is empty it checks for new tasks.
- Step 4: If new task arrives it schedules similarly as in Step 2.
- Step 5: End.

The algorithm can be applied for both Host level and VM level of scheduling.

This study implements resource mapping at both host and VM levels via load balancing using CloudSim 3.0.3 package. Load balancing in cloud provides an efficient solution to various issues applicable to cloud computing set-ups and usage. However, this does not necessarily result in shorting makespan [11]. Hence, we proposed a new load balancer called I-RASA.

5. PROPOSED ALGORITHM

The focused of this study is RASA, to derive an improved Resource Aware Scheduling Algorithm (I-RASA) as a load balancer for both small and large distributed system. The proposed algorithm, shown in Fig 8, calculates the expected completion time of the submitted tasks on each resource. The max-min algorithm is applied on the tasks length for the resources used first while, the min-min algorithm is also then applied for same length of resources used. The max-min algorithm is then applied on the tasks remaining with the overall minimum expected execution time assigned to the resource that had the minimum overall completion time. After scheduling the task is removed from meta-tasks and all calculated times are updated and the processing is repeated until all submitted tasks are executed. The algorithm shown is focused on minimizing the total makespan which is the total complete time in both small and large distributed system. The new algorithm will be compared with the last two discussed algorithms in section II of Table 5 in this paper.

```

Start
Create different classes of Metatask, C
Assign next Metatask to its class //This is done as request (Metatask) comes in

//After all the Metatasks has been assigned to each class then
for (i = 0; i < |C|; i++) // |C| is the total number of classes
    First arrange Metatasks with Max-weight of execution time based on number of resources
    Arrange the next Metatasks with Min-weight of execution time still based on number of resources
    Arrange the remaining Metatasks based on Max-weight of execution time
Next class
end
//R is the resources
for (i = 0; i < |C|; i++) // |C| is the total number of classes
    Assign next task in Ci on  $\min (R_{1_{completiontime}}, \dots, R_{|R|_{completiontime}})$  resource
    Remove the task from the selected class list
Next class
end
End

```

Fig 9: Pseudo-code for I-RASA

Fig 9. shows the flowchart processing. A procedure is created which performs the work of sorting based on execution time of each cloudlets (tasks) and completion time of each virtual machine installed on the host.

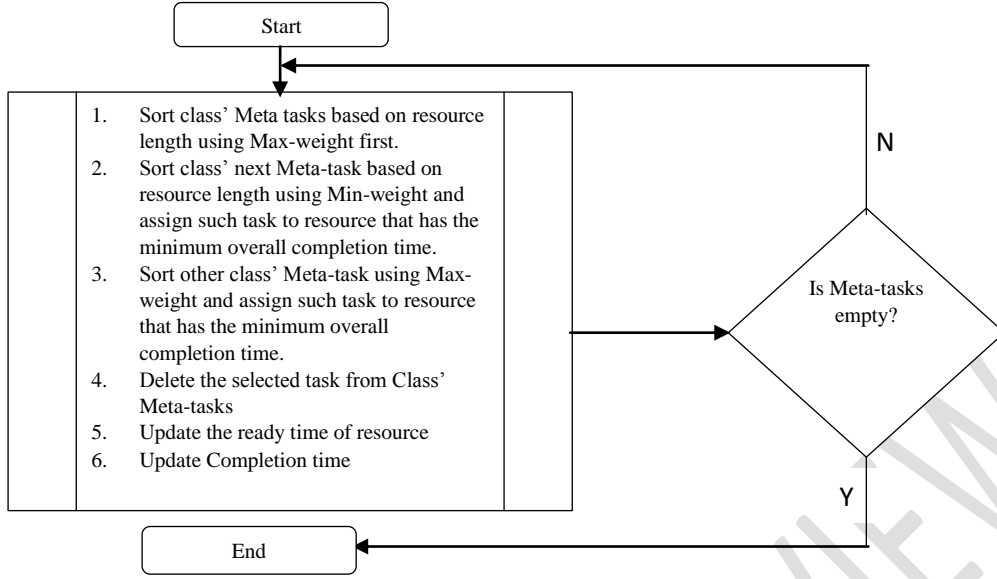


Fig 10: Flowchart for I-RASA

5.1 Simulation configuration and tools

The implementation was done using CloudSim 3.0.3 as a framework in the scalable simulator environment. The simulation was performed on core i3 processor with 6GB Ram and 500GB disk space, performed on Window 7 Ultimate Edition operating system. The experiment was written in java language on Jcreator IDE. The study considered two Datacenters, VM, host and cloudlet components from CloudSim for execution analysis on I-RASA, RASA and I-Max-Min algorithms. The simulation process comprises of different requests by the user to be processed. The total number of cloudlets or request is 50 and it ranges from 300 MB to 23000 MB. The range value was adopted from [30]. Table 1 shows the configuration of the datacenters which consist of four (4) virtual machines with the same configurations.

Table 1: VM Configuration

Parameter	Value
Size (MB)	10000
Ram (MB)	512
Processing Speed (MIPS)	1000
Bandwidth (MBBS)	1000
VM Name	"Xen"
PesNumber	1

6. PERFORMANCE METRIC USED FOR THE EXPERIMENT

The results of the performance of the algorithms I-RASA, I-Max-Min and RASA were evaluated using the following metrics tested in CloudSim toolkit.

6.1 Makespan

Makespan is the finishing time of the last task and in scheduling of task, one of the optimization criteria is minimization of makespan as most of the users desire fastest execution of their application. Equation 1 shows its mathematical representation.

$$makespan = \max_{i \in tasks} \{F_i\} \quad (1)$$

where F_i denotes the finishing time of the last task.

6.2 Economic cost

It indicates the total amount the user needs to pay to service provider for resource utilization. The mathematical representation is shown in equation 2.

$$\text{Economic Cost} = \sum_{i \in \text{resources}} \{C_i * T_i\} \quad (2)$$

where C_i denotes the cost of resources i per unit time and T_i denotes the time for which resource i utilized. From the CloudSim 3.0.3 version a pre-defined cost value for resources used is as follows:

the cost of using processing in VM resource	is	3.0
the cost of using memory in VM resource	is	0.05
the cost of using storage in VM resource	is	0.001

The cost in total is 3.051

6.3 Flow time

This indicates the total sum of finishing times of all the tasks. And to minimize this, tasks should be executed in ascending order of their processing time.

$$\text{Flowtime} = \sum_{i \in \text{tasks}} F_i$$

where F_i denotes the finishing time of task i .

6.4 Resource Utilization

This helps in understanding how busy the resources are. This is very important as service providers would like to earn maximum profit by renting limited number of resources. Equation 3 shows how it can be implemented mathematically.

$$\text{Average resource utilization} = \frac{\sum_{i=1}^n \text{Time taken by resource } i \text{ to finish all jobs}}{\text{Makespan} * n} \quad (3)$$

where n is the number of resources.

6.5 Throughput

This can be defined, as shown in equation 4, as the total number of jobs completing execution per unit time.

$$\text{Average throughput} = \text{Number of Tasks} / \text{makespan} \quad (4)$$

7. PERFORMANCE EVALUATION

Table 2 shows the output of 20 processed cloudlets using I-Max-Min algorithm. The Gantt chart for the flow of process is represented in Fig10. From Fig 10, the makespan of the last finished task is 59.

Table 2: 20 processed cloudlets by four virtual machines using Improved Max-min

Cloudlet ID	VM ID	Time	Cloudlet Length	Start Time	Finish Time
3	3	19	19259	0	19
2	2	20	20332	0	20
1	1	21	20581	0	21
0	0	23	22726	0	23
6	1	16	16377	21	37
4	3	19	19132	19	38
7	0	16	16010	23	39
5	2	19	18811	20	39
11	2	10	9895	39	49
10	0	11	10502	39	50
9	3	12	11666	38	50
8	1	14	13765	37	51
15	1	4	3960	51	55
13	0	7	7060	50	57
14	3	6	6455	50	56
12	2	8	7725	49	57
16	1	4	3843	55	59
19	2	2	2103	57	59
17	0	4	3535	57	61

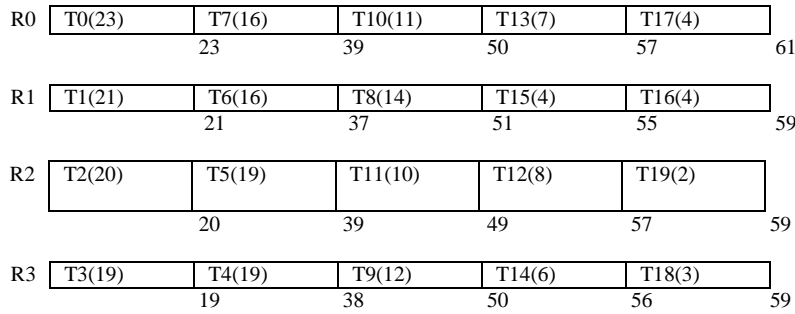


Fig 11: Gant-Chart of I-Max-Min

Table 3 shows the output of 20 processed cloudlets using Resource Aware Scheduling Algorithm. The Gantt chart for the flow of process in this algorithm is shown in Fig11. The makespan from Fig11 is 65. This is higher compare with the makespan in improved max-min.

Table 3: 20 processed cloudlets by four virtual machines using RASA

Cloudlet ID	VM ID	Time	Cloudlet Length	Start Time	Finish Time
1	1	2	2103	0	2
3	3	3	3421	0	3
5	3	4	3535	3	7
2	2	21	20581	0	21
4	1	20	20332	2	23
0	0	23	22726	0	23
7	2	4	3843	21	25
6	3	19	19259	7	26
9	0	4	3960	23	27
11	3	6	6455	26	32
13	3	7	7060	32	39
8	1	19	19132	23	42
12	0	16	16377	27	43
10	2	19	18811	25	44
15	1	8	7725	42	49
17	2	10	9895	44	54
14	3	16	16010	39	55
16	0	14	13765	43	57
18	1	12	11666	49	61
19	2	11	10502	54	65

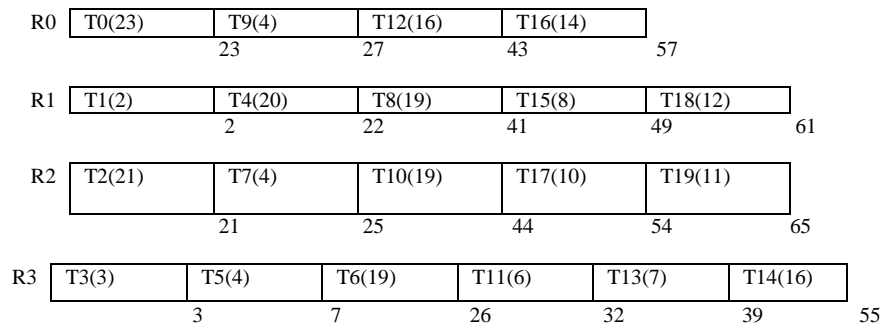


Fig 12: Gant-Chart of RASA

Table 4 shows the output of 20 processed cloudlets using our proposed algorithm called Improved Resource Aware Scheduling Algorithm. The Gantt chart for the flow of process in this algorithm is shown in Fig12. The makespan is 55. This is low when compared with the makespan in the two previous algorithms.

Table 4: 20 processed cloudlets by four virtual machines using I-RASA

Cloudlet ID	VM ID	Time	Cloudlet Length	Start Time	Finish Time
6	3	19	19259	0	19
4	2	20	20332	0	20
2	1	21	20581	0	21
1	3	2	2103	19	21
0	0	23	22726	0	23
3	2	3	3421	20	23
5	1	4	3535	21	25
7	3	4	3843	21	25
12	1	16	16377	25	41
14	3	16	16010	25	41
8	0	19	19132	23	42
10	2	19	18811	23	42
19	0	11	10502	42	53
17	2	10	9895	42	52
18	3	12	11666	41	53
16	1	14	13765	41	55
9	1	4	3960	55	59
11	3	6	6455	53	59
13	2	7	7060	52	59
15	0	8	7725	53	61

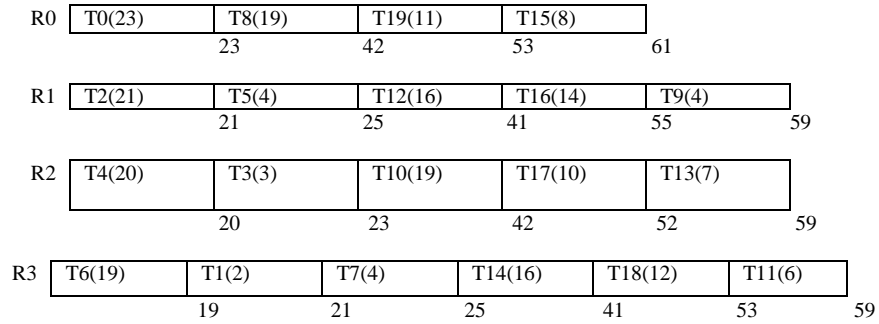


Fig 13: Gant-Chart of I-RASA

Following the report from Table 5 and the graph shown in Fig13, our proposed algorithm, I-RASA out-perform improved max-min and resource aware scheduling algorithms in terms of low makespan, high throughput and on the same merge with the improved max-min in the resource utilization. Although the flow time of RASA is lower compared with our proposed algorithm and that of improved max-min, however, the flow time of I-RASA still outperform improved max-min algorithm. Fig 14 shows the load balancing on the four (4) virtual machines used with respect to the completion time of each.

Table 5: Comparison table of algorithms

Metrics/Algorithms	RASA	Improved Max-Min	Proposed (I-RASA)
Makespan	65	59	53
Economic Cost	726.14	726.14	726.14
Flowtime	696	899	794
Throughput (%)	31	34	38
Resource Utilization (%)	92	100	100

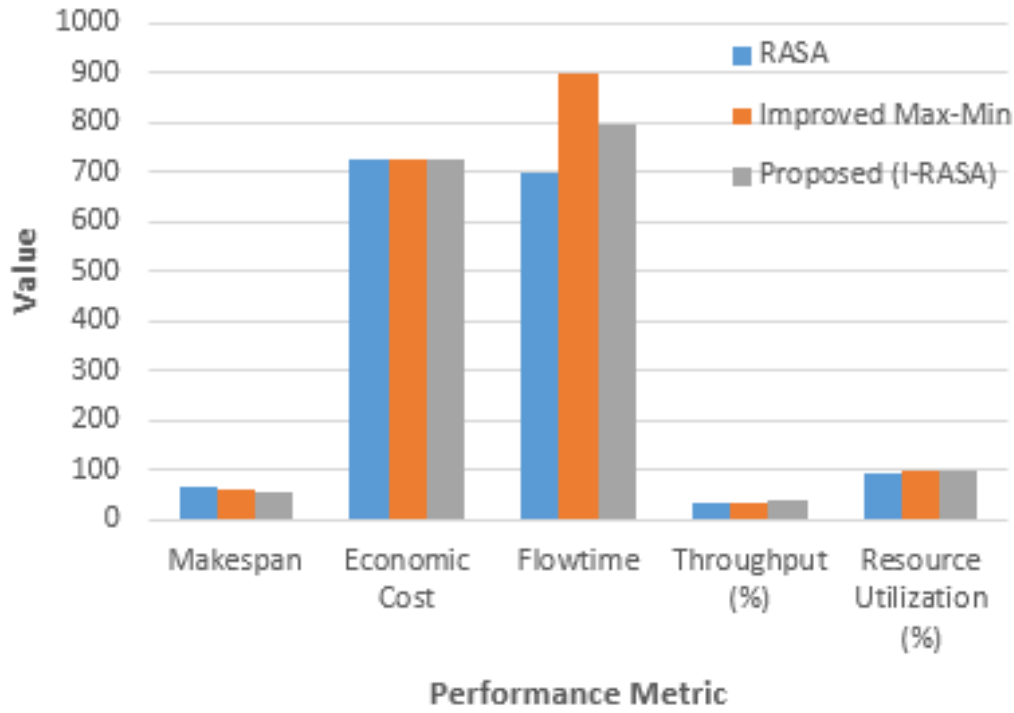


Fig14: Performance Metric Chart for the simulation

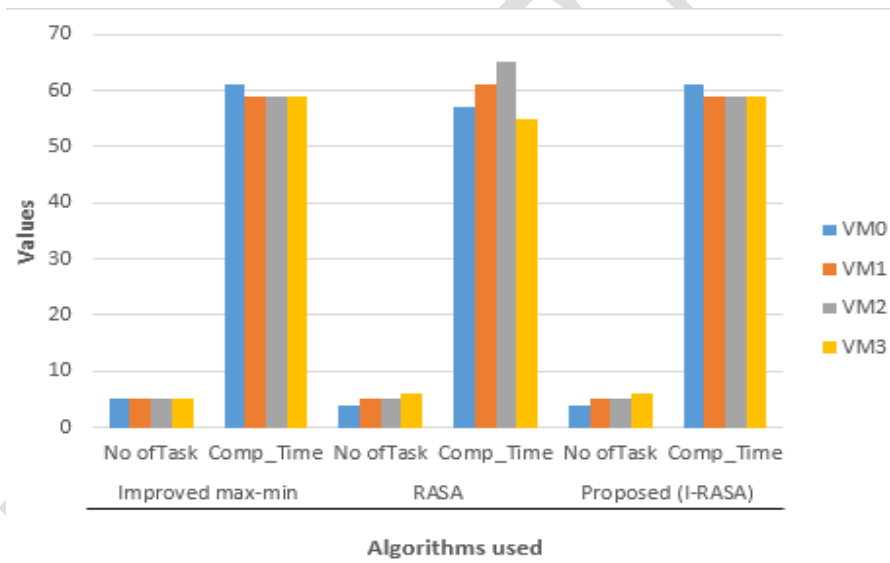


Fig15: Load-balancing Chart for the simulation

The performance evaluation of the proposed I-RASA was shown in Table 6 where the makespan of 20 processed cloudlets were taken alongside with two other algorithms for comparison. It was shown that from ten (10) different iterations, the makespan of the proposed algorithm is still the lowest, followed by improved max-min and lastly RASA. Fig 15 shows the bar-chart of the makespan difference for the ten iterations.

Table 6: Makespan of 20 processed cloudlets on 10 different iterations

Makespan			
Iteration	Improved Max-Min	RASA	Proposed(I-RASA)
1	57	63	51
2	53	59	48
3	50	54	44
4	45	50	41
5	55	59	47
6	58	65	52
7	72	78	61
8	55	62	50
9	59	65	53
10	61	66	54

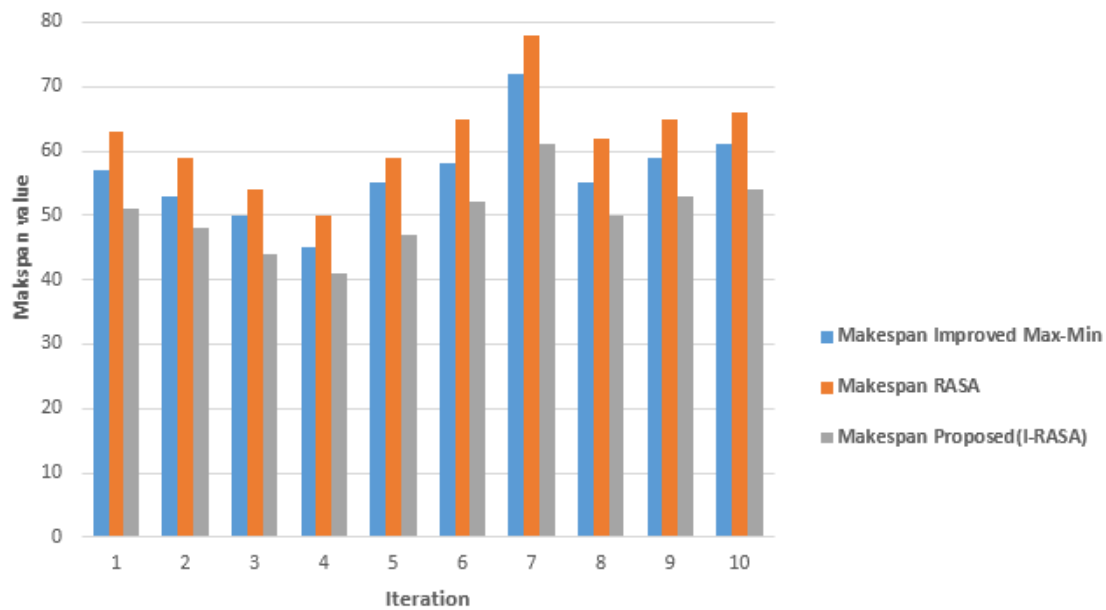


Fig16: Makespan chart 20 processed cloudlets on 10 different iterations

8. CONCLUSION

Cloud controls the lifting of computing-intensive jobs in cloud computing thereby placing enormous amounts of data on the platform especially in mobile cloud computing. Both the data processing and data storage are done in the cloud external of the mobile devices. As mobile applications increase, the leverage on the computing power of the cloud tends to increase as well, it therefore becomes imperative to efficiently manage computing resources for these applications for improving the performance. One of the criterions to improve the performance is by achieving a minimized makespan during load balancing. The study proposed a new load balancing technique, which is an improvement over RASA. The makespan is reduced compared with RASA and Improved Max-Min. The performance analysis from the presented results proved the proposed approach is efficient in optimizing task scheduling and load balancing.

COMPETING INTERESTS DISCLAIMER:

Authors have declared that no competing interests exist. The products used for this research are commonly and predominantly use products in our area of research and country. There is absolutely no conflict of interest between the authors and producers of the products because we do not intend to use these products as an avenue for any litigation but for the advancement of knowledge. Also, the research was not funded by the producing company rather it was funded by personal efforts of the authors.

REFERENCE

- [1] Proshikshya M., Prasant K., Tanmaya S., Amlan D., 2019. Task scheduling algorithm based on multi criteria decision making method for cloud computing environment: TSABMCDMCCE, Open Comput. Sci., 9, p.279–291.
- [2] Nayak S., Parida S., Tripathy C., 2018. Modeling of Task Scheduling Algorithm Using Petri-Net in Cloud Computing, Progress in Advanced Computing and Intelligent Engineering. Advances in Intelligent Systems and Computing, Springer, Singapore, 563, p.633–643.
- [3] Sumanpreet K., Navtej S., 2017. Review on Dynamic Resource Allocation Based on Lease Types in Cloud Environment, International journal of computers & technology, vol 16, p. 7581-7585.

- [4] Saeed P., Reza E., 2009. RASA: A New Grid Task Scheduling Algorithm", International Journal of Digital Content Technology and its applications, Vol. 3, p. 91-99
- [5] Neelima, P., and Reddy, A.R.M. An efficient load balancing system using adaptive dragonfly algorithm in cloud computing. Cluster Computing, Vol. 23, Issue. 1, pp. 2891-2899, 2020.
- [6] Muthusamy, G, Chandran, S. R. Task scheduling using artificial bee foraging optimization for load balancing in cloud data centers. Comput Appl Eng Educ. Vol 28, 769– 778, 2020.
- [7] Jena U. K., Das, P. K. and Kabat M. R. Hybridization of meta-heuristic algorithm for load balancing in cloud computing environment. Journal of King Saud University – Computer and Information Sciences. <https://doi.org/10.1016/j.jksuci.2020.01.012>
- [8] Semmoud A, Hakem M, Benmammar B, and Charr J - C. Load balancing in cloud computing environments based on adaptive starvation threshold. Concurrency and Computation: Practice and Experience, Vol. 32 No. 11, pp. 259-277, 2020.
- [9] Arora, P. and Dixit, A. An elephant herd grey wolf optimization (EHGWO) algorithm for load balancing in cloud. International Journal of Pervasive Computing and Communications, Vol. 16 No. 3, pp. 259-277, 2020.
- [10] Abdulquadi, O. S. and Ravi G. Dual objective task scheduling algorithm in cloud environment. International Journal in Advanced Trends in Computer Science and Engineering, Vol. 9 No. 3, pp. 2527-2534, 2020.
- [11] Elzeki O., Reshad M., Elsoud M., 2012. Improved Max-Min Algorithm in Cloud Computing. International Journal of Computer Applications, vol 50, No 12, p.22-27.
- [12] Morton, A. (2003). IO scheduler benchmarking. linux-kernel (Mailing list). Archived from the original on 2 June 2007. Retrieved 23rd May 2007.
- [13] Iyer, S. and Druschel, P. (2001). Anticipatory scheduling: A disk scheduling framework to overcome deceptive idleness in synchronous I/O. 18th ACM Symposium on Operating Systems Principles. Retrieved 20th April, 2010.
- [14] George D., Amalarethnam and Muthalakshmi P., "An Overview of the scheduling policies and algorithms in Grid Computing ", International Journal of Research and Reviews in Computer Science, Vol. 2, No. 2, pp. 280-294, 2011.
- [15] FatosXhafa, Ajith A., "Computational models and heuristics methods for grid scheduling problems", Future Generation Computer systems, Vol. 26, pp. 608-621, 2010.
- [16] Casavant T. and Kuhl J., "A Taxonomy of scheduling in General purpose distributed computing systems", IEEE Trans on Software Engineering, Vol. 14, No. 2, pp. 141-154, 1988.
- [17] Pawan K., Rakesh R., 2019. Issues and Challenges of Load Balancing Techniques in Cloud Computing: A Survey. ACM Computing Surveys (CSUR) Volume 51, Issue 6, February 2019.
- [18] Udayraj P., Hemant G., 2019. Review of Load Balancing Technique in Cloud Computing, IJRAR- International Journal of Research and Analytical Reviews, Vol 6, issue 2, p.826-833.
- [19] Wang S., Yan K., Chen C., 2011. A three-phases scheduling in a hierarchical cloud computing network, in: Communications and Mobile Computing (CMC), 2011 Third International Conference on IEEE, p. 114–117.
- [20] Neetesh K., Deo P., 2019. A Green SLA Constrained Scheduling Algorithm for Parallel/Scientific Applications in Heterogeneous Cluster Systems," ELSEVIER, Sustainable Computing: Informatics and Systems 22: p. 107-119.
- [21] Bhoi U., Ramanuj P., 2013. Enhanced max-min task scheduling algorithm in cloud computing, International Journal of Application or Innovation in Engineering and Management (IIAEM), p. 2319—4847.
- [22] Venubabu K., 2012. Dynamic Load Balancing for the cloud, International Journal of Computer Science and Electrical Engineering, 2012.
- [23] Danuta S., Ignacio C., Deepak M., Barry O., 2015. On Energy- and Cooling-Aware Data Centre Workload Management. IEEE, pp. 1111-1114.
- [24] Che-Lun H., Hsiao-hsi W., Yu-Chen H., 2012. Efficient Load Balancing Algorithm for Cloud Computing Network. IEEE Vol. 9, p. 70-78.
- [25] Zhi Z., Fangming L., Ruolan Z., Jiangchuan L., Hong X, Hai J., 2015. Carbon-aware Online Control of Geo-distributed Cloud Services. IEEE, pp. 1-14.
- [26] Mao Y., Chen X., Li X., 2014. Max-min task scheduling algorithm for load balance in cloud computing. Proceedings of International Conference on Computer Science and Information Technology; Springer 2014.
- [27] Li X., Mao Y., Xiao X., Zhuang Y., 2014. An improved max-min task-scheduling algorithm for elastic cloud. Computer, Consumer and Control (IS3C), 2014 International Symposium on; 2014: IEEE.
- [28] Buyya R., Ranjan R., Calheiros, R., 2009. Modeling and Simulation of Scalable Cloud Computing Environments and the CloudSim Toolkit: Challenges and Opportunities. In *International Conference on High Performance Computing and Simulation (HPCS)*, 2009.
- [29] Singh A., Goyal P., and Batra S. (2010). Optimized Round Robin Scheduling Algorithm for CPU Scheduling. International Journal on Computer Science and Engineering, vol. 02, no. 07, 2010, p. 2383-2385.
- [30] Neha G., Parminder S., 2014. Load Balancing Using Genetic Algorithm in Mobile Cloud Computing. *International Journal of Innovations in Engineering and Technology (IJET)*, vol 1, Issue 4, 2014.