



PARADIGM SHIFT TO GREEN CLOUD COMPUTING

¹TAMANNA JENA, ²J.R. MOHANTY, ³RAMAKRISHNA SAHOO

¹School of Computer Engineering, KIIT University, Bhubaneswar, India

²School of Computer Engineering, KIIT University, Bhubaneswar, India

³Infosys Limited, Bhubaneswar, India

E-mail: ¹tamannasinghdeo@gmail.com, ²jmohantyfca@kiit.ac.in, ³ramakrishna.sahoo@gmail.com

ABSTRACT

Objective of this paper is to propose a multi-dimensional pricing model cloud providers can adopt to have long run sustainability in the market along with server/ job request consolidation leading paradigm shift to Green Cloud Computing. Energy spent on cooling cloud infrastructures is causing alarming high carbon emission. We present auto scaling of servers depending on the number of job-requests waiting in the queues, interactive transparent pricing strategy by using reward and penalty in accordance to timing of job request and impact of sharing resources among multiple job-requests to attain high server utilization. Job-requests are profiled in detail while registration with the cloud provider. End-user gets to decide for its budget for the service, extent of multitenancy, and reward for the job. Multiple attributes are identified in the profile which carries different weights towards service charge. Probability of returning customers is also predicted.

Keywords: *Green Cloud Computing Green, Pricing Model, Multi-tenancy, Virtualization, Queuing*

1. INTRODUCTION

Cloud Computing is a technology in which a user can access data, services, compute, store, retrieve huge amount of data for domestic as well as commercial purpose without owning data center, software, hardware, server etc. by only paying for its usage (computing facilities, memory, hardware, software and network capacity). The deployment ease and investment free, maintenance free, hassle free regarding installation or update application features make cloud computing an instant hit. Drastic demand of cloud computing has resulted huge datacenters and many more infrastructures. Large amount of energy is consumed in cooling, storage and network systems of the consolidated datacenters resulting high carbon emission. Energy spent for cooling these is resulting contradicting speculations on Cloud Computing being Green. Green cloud computing [1] is an adaptive energy efficient environmentally responsible use of resources for computing along with waste management. It is basically conscious cloud computing. Change of perception to green cloud computing will be highly beneficial in long run to human lives and environment. Server sprawl is a concern, where underutilized servers are not

justified with workload. Its adverse impact shouts environment sustainable approach, though a lot was discussed but very little has done towards it. The main reason behind server sprawl is the requirement of the customer/user that their application need to run in isolation. It is found that in most cases server utilization lies between 10 percentages to 50 percentages. Most cloud computing providers have similar basic strategies like: pay as you go, pay less when you use more, pay less when you reserve. Each cloud providers use non-standard terminology which make difficult for user to compare while deciding for a provider. Cloud federation has 4 deployment models which are classified: on the basis of isolation among users. They are as follows:

- Private Cloud: Private cloud is leased or owned by a single organization otherwise termed as enterprise or internal cloud. Security is its top priority mainly used for military, government, health etc. High speed and compliance are its major advantages. Drawback is capital intensive and capacity ceiling. Need a detail domain expertise and technical planning for getting the right functional setup. Recurrent

- costs like installation, hardware, software, licensing, maintenance is also quite high.
- **Public Cloud:** Public Cloud is a platform where numerous users can access the infrastructure generally maintained by a third party. Users have least visibility or control over the components of the cloud. Its hassle free for individuals and organizations in terms of installation, infrastructure, employing staff, maintenance etc., is done by cloud provider. Its scalable, affordable, energy efficient. Security is not its strength. Microsoft Azure, Google App Engine is some of the most successful public cloud providers.
 - **Hybrid Cloud:** Hybrid Cloud is deployment model which gives the benefit of both public and private clouds. Successful implementation of the concept is difficult. Sensitive data can use smaller private cloud features whereas remaining data can be moved to public cloud for cost effectiveness. Organizations using hybrid cloud need to closely track multiple security platforms and their inter dependability.

It is estimated that, by 2020 US organizations that move to the cloud could save up to 12.3 billion dollar in energy costs and the equivalent of 200 million barrels of oil and reduce carbon emission [6]. It is high time to find the correct trade in between commercialization of cloud computing and environmental sustainability. Cloud users can be broadly classified into individual or organizations (small, medium and large). Amazon Web Service (AWS) is (according to Gartner magic quadrant) much ahead of its competitors like Microsoft, Rackspace, and CenturyLink. Amazon Web Services has started free tier usage since 2010 having an upper cap of computing power and memory. Its business rules are: no penalty, no reward, and no connection charges [6]. When a job is requested with the cloud then the service charge is estimated in accordance with Service Level Agreement of the provider among user and provider. Business policy has come a long way from making high profits to sustain in market for longer. Till date cloud computing is exclusively business oriented. Some discussion has done on carbon emission and after effect on human lives. Cloud has distributed architecture where they have servers and data centers in few locations and centrally monitors its usage per user. It is more or less like having internet providers, customers having different requirements for different purposes having different budgets. Basically cloud computing have 3 different work models i.e.,

Platform as a Service (PaaS), Software as a Service (SaaS) and Infrastructure as a Service (IaaS). User can create its web services or reservation instances from any system having internet.

This paper is organized as follow. Section2 provides background on cloud computing, virtualization and its strength and weaknesses used in cloud computing. Section 3 analyzes our Dynamic Rewarding Model along with its algorithm. Section 4 describes the experimental evaluation and simulated results. We conclude in Section 5.

2. RELATED RESEARCH

Outsourcing business applications to Cloud will cut-down carbon footprint of organization and save power (90% for small businesses, 60-90 % medium and 30- 60% for large businesses) [24]. In near future all businesses will move to cloud sooner or

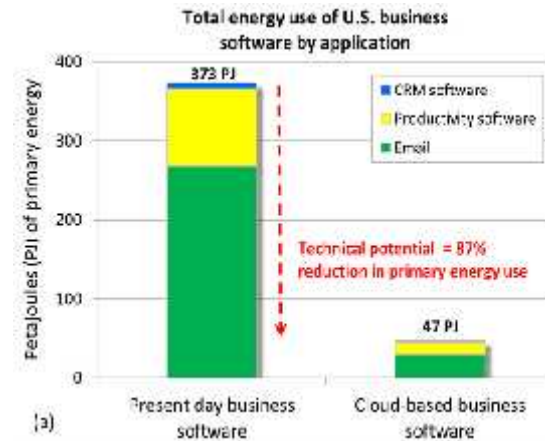


Fig.1: Energy Saving On Transfer To Cloud

later. Cloud provider need to come up with a dynamic pricing strategy which need to be profit oriented, environment friendly as well as sustainable in market for longer. Cloud user need to do their bit by deploying voluminous job request when workload on the server is less. CLEER model[3] gives a detail idea on energy saving and carbon footprint on business move to Cloud illustrated in Fig1, it has worked on datacenters, all type and size of servers across the world. It provides an open source code for researchers to know what and how much can be saved when a business moves to cloud. Greater importance should be given to user as they pay for the cloud services. Till date very limited is achieved towards dynamic negotiation of SLAs between user and mechanisms for automatic resource allocation to numerous competing job requests [2]. In current



scenario cloud provider have inflexible billing and restriction on user to swap one provider for another. There is a gap between market strategy of cloud provider and end user usage pattern. Environment sustainability is to be included in the state of art cloud computing along with market oriented strategies. Survey on pricing strategy study shows, quantity discount, initial setup charges and flexible penalty settings are missing in most of the cloud providers[23]. So the aim here is to address over provisioning of resources and gap in the pricing strategy. Server consolidation needs to be researched extensively in order to deploy smart energy efficient integrated cloud federation worldwide. Energy consumption of any job request is directly proportional to job profile [7]. Extent of virtualization [5] and multi-tenancy needed to be designed and implemented efficiently. Virtualization avoids starvation and provides illusion of owning the server. It is found that server utilization ranges from 10% to 50% in most cases. Utilization is driven by user demand. The work [8] proposed on power constrained performance optimization by heterogeneous server while load balancing. Its multi-objectives are optimal server speed determination, optimal workload distribution among the servers and minimization of response time. Some proposed analytically a finite multiserver queuing model, where applications are modeled as queues and service providers are modeled as virtual machines.

Number of customers waiting in the queue depends on arrival rate, service rate, upper threshold/ lower threshold which trigger the sensor resulting addition/ deletion of server. Recursive method is used to achieve steady state [9] in which size of the queue remain steady. In [10], a single queue is used with dynamic upper threshold and lower threshold showing the number of customers waiting in the queue and its relation with rate of arrival. In [1], the minimum active physical machine is determined by monitoring gross resource weight ratio of VM, it includes VM migration as well in case of heavy load on active physical machine. Threshold of adding or deleting a physical machine depends on gross resource weight ratio of VM, calculated by considering CPU and memory usage. Important metrics of Cloud Computing performance are its scheduling technique and utilization rate of servers. Lot of research is done on schedulers. Assuming all jobs are arriving on same time, [11] have proposed Min-Max Dispersion Round Robin Algorithm. It is a high impact improvised version of RR algorithm where time-slot is taken as difference between max-

burst time and min-burst time of the job requests. Throughput is high with least number of context switches and less waiting times. Genetic Algorithm based task scheduler is also a smart option for executing job requests. Each chromosome is represented as strings of genes. Gene is represented as schedule or slots assigned to tasks initialized randomly. They have divided each job request into two or more tasks and compared the task/ server consolidation using FIFO scheduler, Delay Scheduler and GA based scheduler. The objective function of GA is latest completion time of all tasks [12]. Assumptions taken are computational capacity of each server and upcoming workload which is capacity sensitive is known prior to execution. In [4], Voltage of lightly or idle loaded processor is reduced using Dynamic Voltage Scaling Mechanism(DVSM) and load balancing in case of heavily loaded computer by migrating to lesser loaded system. Implementing DVFS needs some procedural hardware installment to automate the voltage scaling. It is suggested that sustainability of Cloud Provider depends on its QoS and pricing strategy [13]. Using Nash Equilibrium, proposed on economic model for cloud provider. Neither only Cloud provider nor end user can gain profit by changing only their strategy. Pricing strategy need to fluctuate with market all the time.

Determining the best price for complex load situation is quite crucial. Game theory is used to determine extent of inter organizational economics into cloud, i.e., what all and how much of services should be sent to cloud (both public and private) [14]. User having job requests of higher capacity than the processing capacity of server are less likely to have high monetary gain even after deploying service to cloud. Simulating real life scenarios in the commercial world put hurdles as simulators [15] are neither user friendly or efficient enough to consider realistic scenarios. Financial aspect of cloud computing is studied by analytically proposed a model M/M/m, where service charge of multiserver system, net business gain, profit maximization is dependent on factors like workload, satisfaction rate, rental cost, consumption of energy, SLA, and penalty(in case of not meeting SLA)[16]. Consumption of energy in memory can also be reduced by scheduling technique like memory contention, reducing data replication, data transfer etc. Various kinds of virtualizations are used in cloud computing: para virtualization (involves hyper call) is higher performed than full virtualization (involves system call). Virtualization where instead of real memory,



server, hardware etc. virtual entities are created and used which does the job without adding cost price. It is actually the powerful technique needed to be managed efficiently. To make the system the energy-efficient all layers need to be monitored spanning hardware, network, and server [17]. Traditionally pricing strategy was simple and straight forward, basically service are charged for deployment of service. Modern economy incorporates relational, temporal and behavioral matters. A SBIFT model (Scope, Base, Influential, Formula, and Temporal) is proposed while determining service charge of renowned telecom company [6]. Till date Cloud Computing is business benefit oriented. AWS (Amazon are modeled as queues and service providers are modeled as virtual machines. Number of customers waiting in the queue depends on arrival rate, service rate, upper threshold/ lower threshold which trigger the sensor resulting addition/ deletion of server. Recursive method is used to achieve steady state in which size of the queue remain steady. In a single queue is used with dynamic upper threshold and lower threshold showing the number of customers waiting in the queue and its relation with rate of arrival.

Amazon Web Services) has slashed its price more than 40 times by 2014(incorporated Moore's Law [18]) and trend is followed by other providers too. Leading computing service providers have also recently formed a global consortium known as The Green Grid was recently formed to encourage energy efficient datacenters and minimize their environmental impact [19]. Latency is the difference of time between arrival of the job request and complete execution of the job request. It is one of the vital factors while deciding a cloud provider. Nowadays users are more impatient than ever. Cloud computing providers use different types of pricing strategy. Pricing is broadly divided into 3 types:

- 1] On-demand prices (which is subscription less, costliest instances hourly rate for its usage)
- 2] Reserved prices (which is subscription based otherwise called reserved instance (RI) where user pay some upfront price for the usage and pay hourly for its usage)
- 3] Spot services prices (basically unused instances are auctioned by the cloud provider and can be used by the user who bids higher than the fluctuating price of the instance).

In spot services, when the market price exceeds the bid the service is pulled from the provider without any prior notification. It is actually very beneficial for small businesses or personal use where user is tolerant to interruption and instance hourly rates are cheapest. AWS provide Risk Analysis of Spot Services for spot service user as an added value service which helps in bidding. Cons of spot services are spanning frequent interruption, fluctuating prices, no guarantee of complete execution of job request. In [20] researched on the difference between spot services pricing, latency within and across markets of east and west of USA and its impact on arbitrage. Latency data is collected from CloudSleuth.com and concluded that even though technology is now highly integrated but geographical proximity still holds a bigger chunk of cost of leveraging applications leading arbitrage. So we have identified multiple factors which impact job request profile and service charge and came up with the pricing strategy taking Nash equilibrium into consideration how virtualization and different weightage to these factor can pull the business towards green computing.

3. MODEL DESCRIPTION

1. End-user or consumer- Cloud user can request from anywhere in the world having internet. It is categorized into 4 types:
 - a. Domestic/ Commercial,
 - b. Usage time(peak hour, less peak hour, off-peak hour)
 - c. Usage quantity(less, moderate, high)
 - d. Temporal factors(new user, old user, dedicate user)
2. Job request Profiler- Collects specific characteristics and choices of consumer and assign weightage to each factor while profiling each job request. Each job Service Scheduler- Each job requests are processed by scheduler depending on the choices at Job Request Profiler; it navigates to the respective queue.
3. Server- Homogeneous multiple physical machines are used as server, where each server creates multiple virtual machines to process multiple job requests in parallel.
4. Service Charges- Service is charged in accordance to SLA between user and cloud provider.
5. Reward- In our model, each job request is Job request Profiler- Collects specific characteristics and choices of consumer and assign weightage to each factor while

- profiling each job request. Each job request is registered and assigned a unique job identity, along with the estimated service charge, reward, congestion (if applicable).
6. rewarded. The amount of reward is dependent on multiple factor and each factor carries some weightage.
 7. Penalty (Congestion Cost) – Job request which places a voluminous request during the peak hour contributes towards congestion. This charge is called penalty or congestion cost.

We consider a finite buffer multiple queuing systems with queue dependent multi-heterogeneous virtual machines. In our paper system is modeled as M/M/3/K. Finite number of job requests are stored in the buffer called load balancer and after registration with Job Request Profiler it navigates into its respective queue followed by server for processing.

The load balancer will split the arrival stream into 3 sub streams such that:

$$= \lambda_1 + \lambda_2 + \lambda_3 \quad (1)$$

Requests in each queue are processed in FCFS manner considering their choice of server. Here auto scaling of server is used to avoid over provisioning of server. When the number of job requests waiting in the queue exceeds upper threshold, then a new server is added to process the job requests. Maximum 10 numbers of physical machines are considered in our simulation. Similarly when number of job requests waiting in the queue is lesser than the lower threshold, then number of active server is reduced by 1. Our scheduling algorithms are significantly different in their calculation of profitability. Here we are assuming removal of the active server will be done after completion of job requests. Migration of job requests (while processing) from one server to another server causes memory overhead and need lots of technical detailing without much value addition.

In our model given in Fig2 called Dynamic Rewarding Model, we have considered 3 queues (green queue, less-green queue and non green queue). Each incoming job request is assigned a reward. Identified attributes are assigned different weightage which sums the pricing rate. As the name suggests green queue does green computing by extensive multitancy and virtualization. 4

numbers of job requests are taken in one batch and processed in parallel and server capacity is shared with all allocated job requests. Job-requests in each batch are arranged in ascending order. Time quantum is calculated using Round Robin Algorithm using Min-Max Dispersion Round Robin Algorithm [11] for virtualization of the server. Difference between the minimum and maximum capacity of the job-requests in the batch is taken as the Time Quantum (TQ). Value of time quantum is decided iteratively.

Each time a job-request completes execution, new time quantum is decided iteratively depending on the difference between maximum and minimum of remaining job-length needed to be executed in the batch. In less-green queue, 2 number of job requests are processed in parallel in a batch. 2 virtual machines are generated by the physical machine and share the processing power of the server, similarly using Min-Max Dispersion Round Robin Algorithm. Non-green queue execute single job request sequentially. Rationing the demand of job-requests by introducing reward and penalty. Patient behavior of user can change the prospect of Cloud Computing. Rewarding off-peak hour usage, extensive virtualization, introducing a penalty for congestion (sudden voluminous capacity job requests at peak time) will reduce the carbon emission to a huge extent.

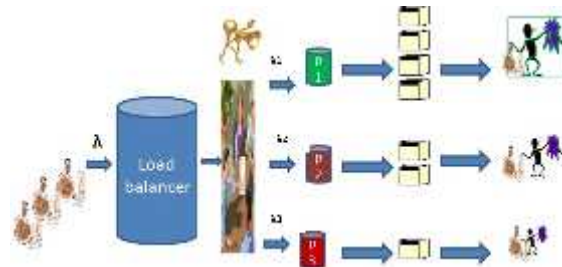


Fig.2: Dynamic Rewarding Model

3.1 DYNAMIC REWARD ALGORITHM

- Step1: Job requests arriving in Poisson distribution are registered with load balancer; profile of job request is done by profile scheduler. Jobs requests are created in accordance to Poisson distribution, all the incoming jobs are stored in form of tuples (jobid, userid, MIPS, com/ dom, usertype, sertype, timing of jobreq), Number of ActSer= 3, MaxActSer= 10
- Step2: Depending upon user choices of server and job profile, service charge and reward is estimated in accordance to pricing policy.



- Step3: If server selects green, job-requests goes to green queue. Job-requests are used to form a batch of 4, host machine creates 4 VMs and jobs requests are executed using Min-Max Dispersion RR algorithm
- Step4: If server selects less green, job-requests goes to less green queue. Job-requests are used to form a batch of 2, host machine creates 2 VMs and job requests are executed using Min-Max Dispersion RR algorithm
- Step5: If server choose is non green, job-requests are executed sequentially in non-green queue on FCFS basis
- Step6: If sum of number of job-requests waiting in any two queue is greater than upper-threshold, then number of active server is raised by 1. ActSer= ActSer + 1
- Step7: If the sum of number of job-request waiting in any two queues is less than lower-threshold, then the number of active server is reduced by 1. ActSer= ActSer- 1.
- Step8: Calculate total revenue, revenue across green server, less-green server and non-green server, reward for each job request, penalty (if applied)

$$\text{TotalReward} = \sum_{i=1}^3 H_i * B_i * D_i * E_i * V \quad (3)$$

H = Hour of job request, peak-hour job-request(h1= 0.2), less-peak-hour job- request(h2= 0.5), off-peak-hour(h3= 0.7); B= Extent of multi tenancy, non-green server(b1=0.2), less green server(b2= 0.5), green server(b3= 0.7), D = Volume of job requests, when low(d1= 0.2), medium(d2= 0.5), high(d3= 0.7) E = type of user, when new(e1= 0.2), old(e2= 0.5), frequent(e3= 0.7) and V = size of job requests in MIPS.

4 EXPERIMENTAL EVALUATION

Our model called Dynamic Rewarding consists of 3parts: Job request profile characteristics, CPU utilization [21] and reward/ penalty associated with each job request. Experimental methods, settings and its generation type are described in detail. Arrival of job requests is in accordance to Poisson distribution where lambda ranges from 4 to 20. Service rate of each server depends on the processing speed of the server and job length of job requests. Results are then presented based on consumption of energy and latency time of job request and waiting time in the queue. Energy consumption is being linear with CPU utilization [22]. We are taking average processor utilization as the metrics for energy used. In our simulation we

have considered homogeneous multiple servers for processing job requests, in all situations minimum 3 servers will remain active whereas additional servers will dynamically become active on trigger of sensor when number of job requests waiting in the queue crosses its upper threshold. Similarly number of server gets reduced when number of job requests waiting in the queue is lesser than lower threshold so that a better trade in can be established between optimal usage of server and frequent addition and removal of server. On completion of assigned job request next job which is on top of the list is picked by the processor. The utilization of the virtual CPU of a virtual machine, Vcpu can be calculated as:

$$Vcpu_i = \frac{VMmips_i}{HOSTmips_i} \quad (2)$$

In each experiment, 200-250 jobs are created from 50 users.

Each job requests capacity ranges are from 5 to 30 MIPS which are generated randomly. MIPS (Million instructions per second) are used to represent the serving speed of host machine, virtual machine and job length of each job request. The number of job requests poured in the buffer, waiting in the respective queues, completed is read every 10 seconds. Total latency time, waiting time, execution time is measured on completion of execution of job request.

Fig 3 shows the number of job requests waiting to be executed and job requests executed in all three queues. The x-axis shows the cumulative number job requests waiting in green queue(in green *), less green queue(in blue *) and non green queue(in red *). Similarly executed green job are plotted in green triangle, executed less green job requests in blue and executed non-green job-requests in red triangle. Number of job requests waiting is much higher in non-green queue (expressed in red color) whereas number of job requests waiting in less-green queue (expressed in blue) lies between non-green and green queue. Number of job requests waiting in green queue is least. Green server executes job requests in batch. Number of job requests executed completely is most in green and least in non-green server. Each attribute is assigned with some weightS (ranging from 0.1 to 0.9). In our simulation 10 host machines and 22 virtual machines are considered. Each host supports up to 4 VMs. By promoting virtualization and multitenancy, higher utilization of servers can be achieved. Fig 4 shows the efficiency of servers. When servers are active means they are using power (higher granularity of

modes of power consumption can even show better results in terms of batch processing). Green server executes job-requests in a batch of 4, less-green server executes job requests in a batch of 2 and non-green server executes job requests sequentially. More MIPS is executed by green server whereas least in non-green server (expressed in red color) and less-green server (expressed in blue) efficiency lies between green and non-green. In order to attain green technology, processor utilization rate should be higher than 50%. In terms of energy efficiency, the more VMs can be placed on the host to get higher utilization rate of processor and lower carbon footprint.

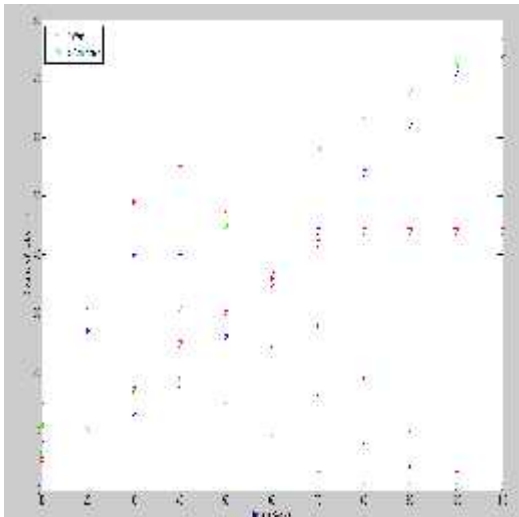


Fig.3: Number Of Job Requests Waiting For Execution And Executed

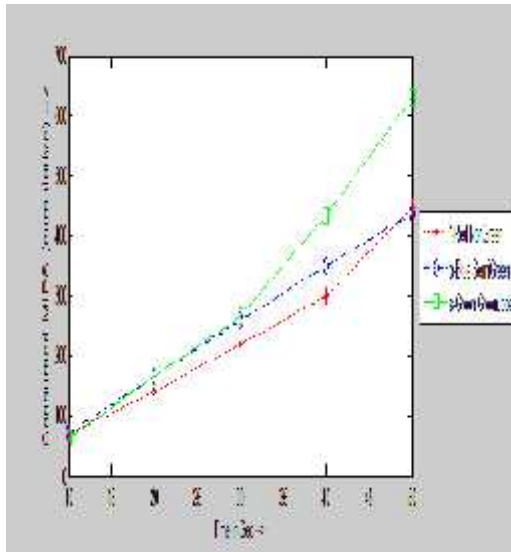


Fig. 4: MIPS-Consumed Per 10 Seconds

Regarding arrival time of job request, time is divided into 3 slots, examples of peak hour refers job requests placed between 9am to 5pm, less peak hour refers 5 pm to 9 pm and off-peak hours is from 9 pm to 9 am.

The transparency of the billing rate, the impact of the attributes on calculation of reward and congestion are discussed in the following

Reward calculation procedure is shown in Table 1.

Profile of few job requests is shown across row.

Job request having serial number 3 is a domestic request where user is categorized into dedicated. Server type: 1 for non-green, 2 for less green and 3 for green. Reward calculated in case of job request 1 is $0.2*0.2*0.5*0.5*V$, reward for job request 2 is $0.5*0.5*0.2*0.2$ and $0.7*0.2*0.5*0.7*V$ for job request 4.

Table1: Reward Business Rule:

Sl No	D/C	N/O/D	Serv	Time	Usage	
1	1	1	1	2	2	4
2	2	2	2	1	1	8
3	1	3	3	3	1	10
4	2	3	1	2	3	10
5	2	1	3	2	2	10

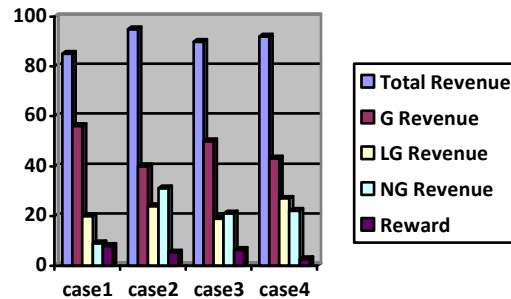


Fig.5: Simulation Result

We have used different pricing rate for domestic as well as commercial user. Depending upon the usage volume, it's subdivided into 3 slots i.e., low, moderate and high. Domestic user gets higher benefit for using less whereas commercial users get better benefit when using more. Higher utilization rate of processor decreases carbon emission and waiting time of job-requests too. In Fig 5, simulation result is plotted in the form of Candle Chart. It is found that when more job request is opting for green server then capital invested



towards reward is higher than otherwise. Revenue invested towards reward is ranges from 0.1% to 7%.

$$Prob_Of_Retn_A_User_t - r_i \cdot Reward_{t-1} \quad (4)$$

Here r_i is a constant whose value lies between 1 to 2.

Total Revenue= Total Revenue Collected towards service charge- Revenue spent towards Reward+ Congestion Cost (5)

Although reward is negative revenue for vendor but it can potentially increases survivability of the vendor in market for longer than its competitors. Higher reward increases the probability of returning a customer to the provider by 15% than with no reward. Limitation of our work is that data locality is not considered. Using more energy efficient hardware, energy efficient scheduling, multitenancy, techniques like Dynamic Voltage Scaling, energy efficient software, maximization of automation etc. altogether can impact huge positive change towards green cloud computing.

5 CONCLUSION

Long term survivability of cloud provider lies in strategic exploration of its pricing model, efficient utilization of resources, optimization of power consumption and transparency with the end user. Cloud federation need to incorporate renewable source of energy, extensive virtualization, energy efficient hardware, maximization of automation, efficient software, energy efficient disk type of memory, responsible usage of internet traffic etc., altogether to achieve Green Cloud Computing and minimize carbon foot print in true sense. Incorporating Nash equilibrium both at pricing model and end user demand can together attain environment sustainable green cloud computing, the most powerful parallel computation platform which integrates finance, retail, production, advertising, logistics etc. Simulation of the concept can be done at a larger volume to get better clarity on its pros and cons. Data locality needed to be considered for better insight. Virtualization needed to be incorporated efficiently so optimization of resources is achieved with lowered carbon emission.

REFERENCES:

[1] P. Mell and T. Grance, "The NIST Definition Of Cloud computing", *National Institute of Standards and Technology*, 2009, 1, pp 26–36.

- [2] R. Buyya, C. Yeo and S. Venugopal, "Market-Oriented Cloud Computing: Vision, hype, and Reality for delivering it services as computing Utilities", *Proceedings of the 10th IEEE International Conference on High Performance Computing and Communications, 2008*, pp 5-13
- [3] F. Masanet, A. Shehabi, J. Liang, L. Ramakrishnan, X. Ma, V. Hendrix, B. Walker, P. Mantha, "CLEER Model", *crd.lbl.gov/news-and-publications/news/2013/studymovingcomputer-services-to-the-cloud-promises-significant-energysavings*, 2013, Lawrence Berkeley National Laboratory.
- [4] Liang-The Lee, Kang-Yuan Liu, hu Yang, Chia Ying Tseng, "A Dynamic Resource Management with Energy Saving Mechanism For Supporting Cloud Computing", *International Journal of Grid Distributed Computing*, 2013, Vol 6, No.1, 2013, pp 67-76
- [5] W. Vogels, "Beyond Server Consolidation", *ACM Applicative*, 2008, Volume 6, pp 20–26.
- [6] E. Iveroth, A. Westelius, C.-J. Petri, N.-G. Olve, M. Cster and F. Nilsson, "How to differentiate by price: Proposal for a five-dimensional mode", *Future Generation Computer Systems*, 2012, 2, pp 16–26.
- [7] R. B. Mark Aggar, Francois Ajenstat, "Cloud Computing and Sustainability: The Environmental Benefits of Moving to the Cloud", *Accenture Microsoft Report*, 2010, <http://www.wspenvironmental.com/media/docs/newsroom/cloudcomputingandSustainability>.
- [8] Junwei Cao, Keqin Li, Ivan Stojmenovic, "Optimal Power Allocation and Load Distribution for Multiple Heterogeneous Multicore Server Processors across Clouds and Data Centers", *IEEE Transactions on Computers*, 2014, 63, pp 45-58
- [9] V. Goswami, S.S. Patra, G.B. Mund, "Performance Analysis of Cloud with Queue-Dependent Virtual Machines", *Proceedings of the 1st International Conference on Recent Advances in Information Technology*, Mar 2012, pp. 357-362
- [10] A.M.D Aljohani, D.R.W. Holton, I. Awan, "Modeling and performance analysis of Scalable Web Servers Deployed on the Cloud", *National Institute of Standards and Technology*, 2009, pp 238 - 242
- [11] Sanjaya Kumar Panda, Sourav Kumar Bhoi, "An Effective Round Robin Algorithm using Min-Max Dispersion Measure", *International Journal on Computer Science*



- and Engineering (IJCSSE), 2012, 4, pp 45–53.
- [12] G. W. Yujia Ge, “GA- Based Task Scheduler for the Cloud Computing”, *International Conference on Web Information Systems and Mining*, , 2010, pp 181–186.
- [13] Ranjan Pal, Pan Hui, “Economic Models for Cloud Service Markets Pricing and Capacity Planning”, *Distributed Computing and Networking (ICDCN)*, 2012, pp 113-124.
- [14] Jorn Kunsemoller, Holger Karl “A Game-Theoretical Approach to the Benefits of Cloud Computing”, *Economics of Grids, Clouds, Systems and Services Lecture Notes in Computer Science* ,Volume 7150, 2012, pp 148-160.
- [15] Georgia Sakellari, George Loukas, “A survey of Mathematical models, simulation Approaches and test beds used for research in cloud computing Market oriented Cloud Computing: Vision, hype, and reality for delivering it services as computing Utilities”, *Simulation Modelling Practice and Theory*, 2013, 39, pp 92–103.
- [16] Junwei Cao, Kai Hwang, Kequin Li, Albert Y. Zomaya, “Optimal Multiserver Configuration for Profit Maximization in Cloud Computing”, *IEEE Transaction on Parallel and Distributed Systems*, Volume 6, 2013, pp1087-1096
- [17] Andreas Berl, Erol Gelenbe, Marco Di Girolamo, Giovanni Giuliani, Hermann De Meer, Minh Quan Dang, Kostas Pentikousis, “Energy-Efficient Cloud Computing”, *The Computer Journal Advance Access*, Published by Oxford University Press on Behalf of the British Computer Society European Management Journal, 2009. pp 1045-1051
- [18] N.Clayton, “Meet the rain makers”, *National Institute of Standards and Technology*, 2013.
- [19] Jayant Baliga, Robert Ayre, Kerry Hinton, Rodney S. Tucker, “Green cloud computing: Balancing energy in Processing, storage and transport”, *Proceeding IEEE*, 2011, pp 149–167.
- [20] Hsing Kenneth Cheng, Zhi Li, Andy Naranjo, “Cloud Computing Spot Pricing Dynamics Latency and Limits to Arbitrage, April 2013,
- [21] Cha Tung Yang, Jung Chun Liu, Kuan Lung. Huang and Fuu Cheng Jiang, “A method for managing green power of a virtual machine clustering cloud”, *Future Generation Computer Systems*, 2014, Volume 37, pp 26–36
- [22] R.Yamini, “Power Management in Cloud Computing Using Green Algorithm”, *IEEE International Conference On Advances In Engineering, Science and Management*, March 30, 31, 2012, pp 128 133
- [23] J. Huang, ” Pricing Strategy for Cloud Computing Services”, In the proceedings of The Pacific Asia Conference on Information Systems (PACIS) 2013, pp 279
- [24] Saurabh Kumar Garg, Rajkumar Buyya, ” Green Cloud computing and Environmental Sustainability” *Harnessing Green IT: Principles and Practices*, S. M. a. G. G. (eds), Ed. UK: Wiley Press, 2012, pp 315-340