

A Review of Methods for Resource Allocation and Operational Framework in Cloud Computing

Hadi Moei Emamqeyisi¹, Nasim Soltani², Masomeh Robati³, Mohamad Davarpanah⁴

Received (2017-05-21)

Accepted (2017-09-12)

Abstract — The issue of management and allocation of resources in cloud computing environments, according to the breadth of scale and modern technology implementation, is a complicated issue. Issues such as: the heterogeneity of resources, resource dependencies to each other, the dynamics of the environment, virtualization, workload diversity as well as a wide range of management objectives of cloud service providers to provide services in this environment. In this paper, first the description of a cloud computing environment and related issues have been reported. According to the performed studies, challenges such as: the absence of a comprehensive management of resources in the cloud environment, the method of predicting the resource allocation process, optimum resource allocation methods to reduce energy consumption and reducing the time to access resources and also implementation of dynamic resource allocation methods in the mobile cloud environments, have been addressed. Finally, with regard to the challenges, some recommendations to improve the process of allocation of resources in a cloud computing environment is has been proposed.

Keywords — cloud computing, resource allocation, resource management, virtualization, virtual machine migration

I. INTRODUCTION

With the rapid development of processing and storage technology, the Internet successfully converted the computational resources to the cheapest, most powerful and most achievable resources compared to the past. The process of this technology has achieved a model called cloud computing, in which resources (such as CPU, storage devices) are as public utilities that can be rented or released (redeemed) by users via the Internet in the form of a simple request [1, 2]. Two definitions are provided for cloud computing by the National Institute of Standards and Technology Cisco. Cloud Computing from the perspective of the National Institute of Standards and Technology is a model which by using its users will be able to receive configurable shared resources such as networks, provider, storage, application and services when they apply. These sources with minimal management effort or without any need for interaction with a service provider, promptly can be prepared and released [3]. But the Cisco's definition of cloud computing is: Resources and information technology services which in a abstract form and independent of infrastructure, provides a multi-tenant environment by demand and with any scale [4].

In recent years, cloud computing has emerged as a new form of utility based on a computational model for hosting and providing hardware and software as a type of services. That gives the impression to users that they are able to have unlimited access to computing resources and storage at any time and any place instantaneously [5]. Cloud computing has emerged as a new

1- Department of Electrical and Computer Engineering, Foolad Institute of Technology, Fooladshahr, Isfahan , Iran. (hadimoei@gmail.com)

2- Department of Software Engineering, Allame Naeini Higher Education Institute, Naein, Isfahan, Iran.

3- Department of Software Engineering Pooyesh Higher Education Institute, Qom, qom, Iran.

4- Department of Electrical and Computer Engineering, Foolad Institute of Technology, Fooladshahr, Isfahan , Iran.

model for the delivery of applications, platforms or computing resources (processing power, memory, bandwidth, etc.) To customers as “pay for use”. This model in the cloud computing is very economical because customers pay for their real consumption, without no additional cost. Also, this model is very flexible because it can be used depending on customer requirements. Due to these benefits, cloud computing is increasingly being utilized in various environments such as banking, e-commerce, Academic environments and etc [6, 7].

II. RESOURCE IN CLOUD COMPUTING

Now cloud computing is a business. So the price and quality of service is important in this area. Providers of cloud computing, deliver resources and software to clients in terms of service level agreement (SLA). SLA is a contract between the user and the cloud provider that contains the resource requirements of users, service time and cost constraints of service which has very important benefits for a business investor. Cloud service provider is willing to benefit from cloud computing service users also do not want to pay a lot of money. Cloud computing users, receive good quality services from service providers that service fee is based on the process of allocating resources in a particular service environment. Cloud service providers must allocate resources to clients specific to a particular method [8]. There are several resource allocation models used in the field of cloud computing. Each of these models uses a particular method and algorithms to achieve this goal. Consumers in cloud environments are not involved in significant investment in information technology's (IT) infrastructures and complex issues related to construction and maintenance of them. In this model, users regardless of knowledge about where services are hosted have access to needed services. Cloud computing based on “pay- use” model hosts practical, commercial and scientific programs. Data centers hosting for applications, consume a lot of energy. Maintenance of large data centers requires high energy consumption. It has been estimated that the cost of maintaining data centers increased the main cost of the original investment, in this case, the maintenance of data centers with this situation could be even adverse, and unfortunately it has been seen that

this process is being continued without any limit. From the perspective of service providers, the method of maximizing the profits given the high cost of energy is a problem. Since the profit is pared to expenses, which its profit should be compensated through provision of services to customers. One way to reduce costs and increase profits is reducing energy consumption. The rising of cost of energy is a big and a potential threat to increase the cost of ownership.

A. Management and Resource Allocation

One of the most important tasks of cloud providers is managing and allocating resources. Cloud users are sending service requests to cloud from anywhere in the world. It should be noted that there are differences between cloud users and users of deployed services. Consumers could be an enterprise deploying web applications that offer different workflows according to the number of users who have access. Cloud service providers should reassure customers that their needs will be fully met. Until recently obtaining high performance was the only concern when allocating resources, regardless of increasing energy costs. In cloud computing, resource allocation (RA) is the process of allocating resources to required applications on the internet. Resource allocation, the allocation of services robs if not managed carefully. If the resource allocation is not managed carefully, it restricts services. Providing resources solves this problem by the possibility of managing resources by service providers for each of the models.

Resource allocation strategy (RAS) is about the incorporation of activities of cloud provider for exploitation and the allocation of scarce resources in the range of cloud environment to meet the needs of cloud applications. This strategy requires the type and amount of required resources by each application in order to complete the task of the user. The order and the time of resource allocation for an optimal RAS is an input. An optimal RAS should avoid the following criteria as follows [9].

- Resource Contention: Contention occurs when two applications try to achieve the same resources at the same time.
- Scarcity of resources: Scarcity of resources occurs when that there are limited resources and high demand for resources.

- **Separation of resources:** Separation of resources occurs when the resources are separated. There are sufficient resources but it can't be allocated to the required application due to segregation to small units.
- **Oversupply:** Oversupply occurs when the application reaches additional resources more than requested.
- **Supplying less than criterion:** It occurs when much lower resources are assigned to the application than it has been demanded.

Energy consumption in cloud environments is examined from two perspectives, the first view is the static power management, which is more related to hardware equipment and the second view is the dynamic energy management.

B. Variety of Resources

In this section of the main types of resources as a service cloud computing environments are explained.

- **Processing resources:** Processing resource are a set of physical machines (PM) that each includes one or more processors, memory, network interfaces, local input and output. All of these components are offered as computational capacity of a cloud computing environment. Usually PMs and VMs are implemented are produced virtually. VMs operate independently of others, and each may have different operating systems and applications. Most scholars limit PM and VM to increase the processing and memory capacity in the field of processing resources.
- **Network resources:** processing resources are organized in the data center shelves which usually are ready for hosting multitude groups of applications that need to be supplied. There is a problem that is associated with processing resources. Physical machines need Gigabit bandwidth and more to be able to respond all resource allocation requests that are sent by various applications. In addition to the applications, applications that include parallel processing require high bandwidth to transmit information. According to these problems, programs should be connected without creating

additional overload to the data centers and network protocols should maintain their effectiveness. Here, there are two important aspects:

The first aspect is related to the designed topology to the network which has a significant impact on the performance and fault tolerance in a cloud environment. Current network topologies are in the data center of hierarchical topologies, such as tree topology that was used in the initial mobile networks. The new topologies including fat trees and hyper-cubes or randomized small-world topologies have also emerged. In all these topologies the aim is to increase the number of input and output ports linearly with network bandwidth [10].

The second aspect is directly tied to resource management, that how to predict the latency and the network bandwidth in data centers according to the different types of traffic patterns. Traditionally, this problem should be resolved on issues related to network but in large data centers due to lack of necessary communication and lack of constant traffic pattern, the cost and difficulty of work have been increased. Due to this, it has been tried to make a progress to implement and deploy different and high quality services which their performance is separated from traffic policy. This work requires a high level of engineering to design traffic patterns. The natural extension of this approach is towards technologies that provide network resources virtually. Creation of virtual networks provides the opportunity to address these projects and networking protocols.

- **Storage Resources:** public cloud providers such as Amazon often offer a wide variety of storage services. Services include virtual drives, database services, and other storage facilities, which are at different levels of stability and reliability. A challenging issue for providing storage services is their performance in dynamic environments. Such that a service should be able to adapt to the requests optimally by increasing the number of users and increasing the workload or increasing the volume of data or vice versa by reducing these criteria can be optimally adapted to the requests. As in traditional and centralized systems stored data should

be compatible with transaction features including Atomicity, Consistency, Isolation, Durability, Distributed (ACID), these properties are essential in cloud computing systems. Fortunately, many applications implemented in a cloud environment has a great capacity for adaptation, and this allows designers to reconcile between consistency and efficiency of the implemented systems in a cloud environment. It is implemented by a wide range of data storage technologies with the capability to implement in different operational and non-operational conditions. Examples of these technologies include: text storage, vertical storage, graphics storage and storage using distributed keys. Recently the method of using distributed key for data storage has attracted much attention. This method also has the ability to support dynamic environments. In the storage method using distributed keys, insertion, retrieve and data modification in cloud computing environments are also possible[11].

Energy Resources: energy costs account for a significant share of the costs of their data centers. Energy costs include charges for consumed energy by systems and data center hardware, as well as the energy used to cool the data center. As noted above, a significant share of the costs accounts for cost of energy. For example, according to research by Mr. Hamilton, the energy consumption of systems in data centers is equal to 19 percent of overall costs as well as costs related to cooling data centers account for 23 percent overall costs [12]. In a data center energy is consumed by server equipment, network equipment, power distribution equipment, cooling and infrastructure related to the consumption support. Data centers typically provide one or more energy supplies. Recently, data centers have turned to the production of energy through renewable energy. As well as to reduce the costs, they have a special focus on methods of reducing energy consumption. In this area there are four main approaches:

- establishment and implementation of hardware components that have minimal energy consumption to improve energy

efficiency.

- Implementation of efficient resources allocation methods to reduce energy consumption.
- Implementation of useful applications in order to reduce energy consumption.
- Deploying cooling systems consistent with the geographic region and climate to reduce energy consumption.

III. THE OPERATIONAL FRAMEWORK FOR RESOURCES ALLOCATION

Currently, many academic and industrial researchers are addressing the issues related to the allocation of resources in cloud environments. Management and allocation of resources in cloud environments consists of many subjects, which can be addressed from different aspects. In the following an overall outline of the issues related to the allocation of resources in cloud environments is shown and then some items of works carried out in every area also will be expressed.

The framework shown in Figure 1, generally shows objectives and tasks done by different researchers in order to allocate resources in cloud environments. Each of these cases include algorithms and solutions to optimal provision of resources with respect to the objectives identified in the management of resources.

- Global Scheduling of Virtual Resource

Global Scheduling includes a general vision of allocation of physical and virtual resources in the cloud environment. In recent years, extensive proposals for determining a scheduling framework is provided. More offers are related to academic research[13]. The objectives of research conducted associated with the global scheduling of resources include the following items:

- Locating and Allocating Virtual Machines

when a user makes a request for one or more VM's resource management system that is controlled by the cloud service provider, adjusts the VM's time schedule as well as a method of use for the user, according to management objectives and the agreements. The first problem to locate VMs is non-dynamic and fixed placement of VMs on physical machines. To achieve a stable and reliable service in which the virtual machines could be placed on physical machines dynamically, as well as any VM in accordance with the specified service does not get out of their group, it is needed to resolve the issue of dynamic

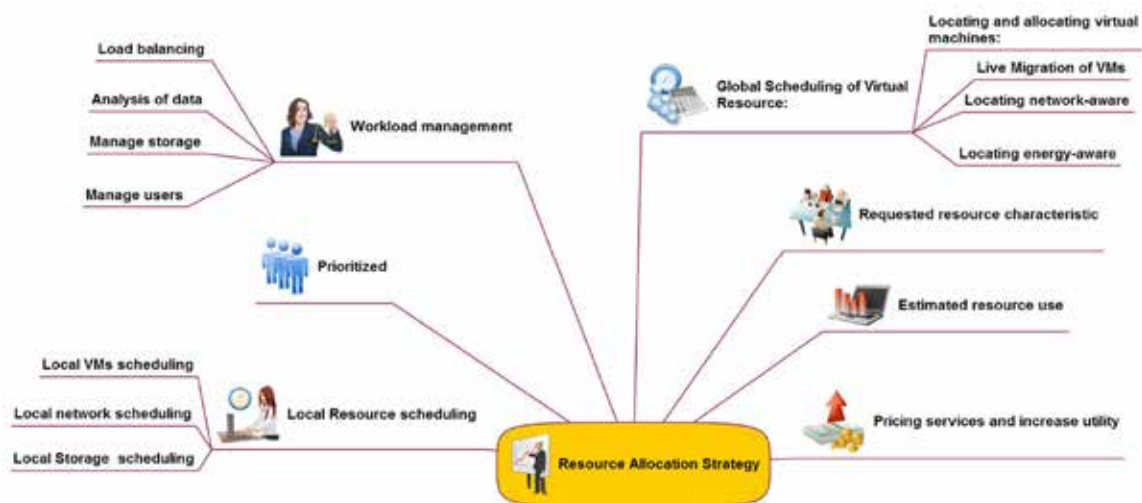


Figure 1: Operations related to the management of resources in a cloud environment

Locating within the group. Such issues are the kinds of Non-deterministic polynomial-time (NPhard) issues. Heuristic algorithms are used to solve these issues. Most of the proposed heuristic algorithms are based on greedy algorithms that use simple rules [14, 15].

- Live Migration of Virtual Machines

Following the initial Locating, location of VMs can change dynamically according to the hosting changes. From the perspective of global scheduling this is called live migration of virtual machines. According to this, a virtual machine running, can temporarily be stopped and be transferred from one physical machine to another physical machine and went back to work and continue to do their service. VM live migration is one of the major challenges[16]. Algorithms for improving Live Migration of VM are seeking to achieve benefits such as improved fault tolerance, reduce the cost of maintenance of physical machines, reduced response times, dynamic and arbitrary use of and physical machines in order to reduce energy costs and also achieve management objectives, and other benefits.

- Location-aware Locating of Virtual Machines

VMs communicate through physical machines network interface with other applications and other system components. Usually, PMs interface is considered as an unmanaged source that cannot

be to be available to the VMs as a guaranteed resource. So it appears that the nominal bandwidth of virtual machines is not provided, according to this problem VMs compete to acquire network resources that could be influential on the efficiency of VMs [17, 18]. In the proposed solutions it has been tried to minimize the competition between the VMs as much as possible.

- Energy location-aware Virtual Machines

As mentioned earlier, resource management systems utilize features such as migrating virtual machines to reduce their energy consumption. Many researchers have tried to do better and more efficient this task by proposed methods. Parallel to research, data centers are also developing technologies that support live VM migration. More focus on applied research is on this issue that virtual machines can be deployed on fewer physical machines. However, yet it is a challenge for cloud service providers[19].

- Requested Resource Characteristic

Cloud computing provide through Internet their services to others. Considering that the volume of applications on the network, dynamically varies according to different times, so it is needed to provide techniques for forecasting source demand requests. To this aim demand classification technique are used.

- Estimates of Resource Use

The process of resources management needs to estimate the resources required to achieve two

important goals: first to identify the different types of open source and second, identify the pattern of resource use and people who are affected by the workload of using resources. In most presented articles it is assumed that an immediate assessment and forecast of resources are easily accessible.

- Workload Management

In cloud computing environments, the use of applications from cloud resources is called workload. Requests for use of resources will be accepted or rejected. One important aspect of the workload management is load balancing that we will explain it in the following.

- Load Balancing

The level measurement usually to maintain a service, when one or more components fail, is mechanized. Components are monitored frequently, and when one component does not respond, load balancing pops up and does not send traffic on it. With appropriate surveys in place of consumption often it is possible to minimize the problems, which not only reduces the cost and provides green computing, but also holds down pressure on unique circuits that potentially extends their life.

In fact, the purpose of load balancing is finding efficient mapping of tasks on the processors in the system so that on each processor equal amount of tasks to be executed until the overall tasks reach to its minimum value.

The importance of load balancing: with the load balancing, it is possible to balance the load with dynamic transition of workload from one machine to a machine at remote node or a machine that is used less. This work maximizes the user satisfaction, minimizes response time, increases resource utilization, reduces the number of rejected tasks and increases system performance.

IV. SUMMARY AND CONCLUSION

Based on the analysis and studies, there are several benefits in resource allocation while cloud computing is being used regardless of the size of the organization and commercial markets. But there are limits, because cloud computing is an emerging technology. In cloud computing resources are provided as services to users, this resource allocation or services has some advantages and also disadvantages, for example It does not need hardware and software installation,

solving location restrictions, lack of user control over resources and more. In addition there are some challenges in the process of allocation of resources that could be the basis for future research. Some of the challenges are issues which are inherently complex and difficult, it is needed to discover new methods to resolve these kinds of problems and some of the challenges arise due to the management and design limitations or issues, the reason is due to the lack of a specific standard for cloud computing environments.

Given the extent of the issue, different people from different perspectives and with different goals focused on the issues related to the allocation of resources in cloud environments. In Table 1 it has been tried to present a general summary of the proposed methods along with the advantages and disadvantages of each one, and also the purpose of each method.

According to reviews and comparisons, one of the challenges in resource allocation process which requires specific attention and developing a standard method is achieving methods to predicting performance of applications and requested resources. Predicting and understanding the performance of applications such as exact run time, the time of requesting resources and the time of releasing resources as a challenge that must be specified at the time of formulation of SLA. For this purpose we need for an accurate and precise modeling and timing of applying requests, which due to the dynamics and complexity of the problem it is a difficult work. In addition there is no fixed control to predict the performance of applications in a cloud environment. In general, the prediction of performance in most of computer systems is a challenge. In cloud environments due to the presence of a further layer as virtual layer, the task of prediction is more difficult. As a result of prediction of demand for virtual resources is more difficult than physical resources. Therefore, the proposed method for this purpose, is using data mining algorithms to classify incoming requests. Due to the nature of classification algorithms and also the science of data mining, it is possible to provide a classification of received requests from previous operations of allocating resources in the cloud as well as reviewing the circumstances of each request, and accordingly provide the resources required for each one of the applications. It also requires the collection of data from history of demands. So after identifying

the different categories of requests, at the time of receipt of a request and then specification of category and pattern of request then resources related to this category of requests is offered as well as the release time and the type of use from these resources by request is also highlighted to perform management tasks.

Table 1: Review on Resources Allocation Methods

Researcher	Method	Energy Consumption	Time	Dynamic	cost	Fault Tolerance	Load Balancing
R Madhumathi, A. S. Balagopalan	Bin-Packing [20] (2015)	x	✓	x	x	x	✓
Ronak Patel	TARA [21] (2013)	x	✓	x	✓	x	✓
Aman Kumar	LSTR [22] (2012)	✓	✓	x	x	x	✓
Anand Prabu P, Sairamprabhus	Nephele Framework [23] (2014)	x	✓	✓	x	x	✓
Aman Kumar	EARA [24] (2013)	✓	✓	x	✓	x	✓
S Gopal Kirshna Shyam	Resource Allocation Using Agents [25] (2015)	✓	✓	x	✓	✓	x
Dilip Kumar, Bibhudatta Sahoo	MAXMAX [26] (2014)	✓	x	x	x	x	✓
Siva Theja Maguluri, Lei Ying	traffic optimal [27-29] (2012-214)	x	✓	x	x	x	✓
W. ju-Hu	Bee Colony Algorithm and evolutionary algorithms [30-32] (2014-2015)	x	✓	✓	x	x	✓

REFERENCES

- [1] Q. Zhang, L. Cheng, and R. Boutaba, "Cloud computing: state-of-the-art and research challenges," *Journal of internet services and applications*, vol. 1, pp. 7-18, 2010.
- [2] S. Taherian Dehkordi and V. Khatibi Bardsiri, "Optimization Task Scheduling Algorithm in Cloud Computing," *Journal of Advances in Computer Engineering and Technology*, vol. 1, pp. 17-22, 2015.
- [3] Z. Chen and J. Yoon, "IT auditing to assure a secure cloud computing," in *Services (SERVICES-1)*, 2010 6th

World Congress on, 2010, pp. 253-259.

- [4] K. Bakshi, "Cisco cloud computing-Data center strategy, architecture, and solutions," *CISCO White Paper*. Retrieved October, vol. 13, p. 2010, 2009.

- [5] S. Ray and A. De Sarkar, "Execution analysis of load balancing algorithms in cloud computing environment," *International Journal on Cloud Computing: Services and Architecture (IJCCSA)*, vol. 2, pp. 1-13, 2012.

- [6] V. Yarmolenko, R. Sakellariou, D. Ouelhadj, and J. M. Garibaldi, "SLA based job scheduling: A case study on policies for negotiation with resources," in *Proceedings of e-Science All Hands Meeting (AHM2005)*, 2005, pp. 20-22.

- [7] R. Buyya, C. S. Yeo, S. Venugopal, J. Broberg, and I. Brandic, "Cloud computing and emerging IT platforms:

Vision, hype, and reality for delivering computing as the 5th utility,” *Future Generation computer systems*, vol. 25, pp. 599-616, 2009.

[8] B. B. Nasim Soltani, Behzad Soleimani Neysiani, “Job Scheduling based on Single and Multi Objective Meta-Heuristic Algorithms in Cloud Computing: A Survey,” *Conference: International Conference on Information Technology, Communications and Telecommunications (IRICT)*, , vol. 2, March 2016.

[9] V. Vinothina and R. Sridaran, “A survey on resource allocation strategies in cloud computing,” *International Journal of Advanced Computer Science & Applications*, vol. 1, pp. 97-104, 2012.

[10] M. Al-Fares, A. Loukissas, and A. Vahdat, “A scalable, commodity data center network architecture,” in *ACM SIGCOMM Computer Communication Review*, 2008, pp. 63-74.

[11] I. Robinson, J. Webber, and E. Eifrem, *Graph databases: new opportunities for connected data: “O’Reilly Media, Inc.”*, 2015.

[12] J. Hamilton, “Cooperative expendable micro-slice servers (CEMS): low cost, low power servers for internet-scale services,” in *Conference on Innovative Data Systems Research (CIDR’09)(January 2009)*, 2009.

[13] S. Abirami and S. Ramanathan, “Linear scheduling strategy for resource allocation in cloud environment,” *International Journal on Cloud Computing: Services and Architecture (IJCCSA)*, vol. 2, pp. 9-17, 2012.

[14] R. Madhumathi, R. Radhakrishnan, and A. Balagopalan, “Dynamic resource allocation in cloud using bin-packing technique,” in *Advanced Computing and Communication Systems, 2015 International Conference on*, 2015, pp. 1-4.

[15] W. Ju-Hua, “Research of Resource Allocation in Cloud Computing Based on Improved Dual Bee Colony Algorithm,” *International Journal of Grid and Distributed Computing*, vol. 8, pp. 117-126, 2015.

[16] Z. Xiao, W. Song, and Q. Chen, “Dynamic resource allocation using virtual machines for cloud computing environment,” *IEEE transactions on parallel and distributed systems*, vol. 24, pp. 1107-1117, 2013.

[17] S. Parida and S. C. Nayak, “Emperical Resource Allocation Using Dynamic Distributed Allocation Policy in Cloud Computing.”

[18] G. K. Shyam and S. S. Manvi, “Resource allocation in cloud computing using agents,” in *Advance Computing Conference (IACC), 2015 IEEE International*, 2015, pp. 458-463.

[19] G. Portaluri, S. Giordano, D. Kliazovich, and B. Dorronsoro, “A power efficient genetic algorithm for resource allocation in cloud computing data centers,” in *Cloud Networking (CloudNet), 2014 IEEE 3rd International Conference on*, 2014, pp. 58-63.

[20] R. Madhumathi, R. Radhakrishnan, and A. S. Balagopalan, “Dynamic Resource Allocation in Cloud Using Bin-Packing Technique,” *2015 International Conference on Advanced Computing and Communication Systems*, 2015.

[21] Ronak Patel and S. Patel, “Survey on Resource

Allocation Strategies in Cloud Computing,” *International Journal of Engineering Research & Technology (IJERT)*, vol. 2, 2013.

[22] Abirami S and S. Ramanathan, “Linear Scheduling Strategy for Resource Allocation in Cloud Environment,” *International Journal on Cloud Computing: Services and Architecture(IJCCSA)*, vol. 12, 2012.

[23] Anand Prabu P, Dhanasekar P, and S. S. G, “Dynamic Resource Allocation Using Nephel Framework in Cloud,” *International Journal of Innovative Research in Computer and Communication Engineering*, vol. 2, 2014.

[24] Aman Kumar and E. S, “An Efficient Framework for Resource Allocation in Cloud Computing,” *4th ICCCNT 2013*, 2013.

[25] Gopal Kirshna Shyam, SunilKumar, and S. Manvi, “Resource Allocation in Cloud Computing Using Agents,” *2015 IEEE International Advance Computing Conference (IACC)*, 2015.

[26] Siva Theja Maguluri , R. Srikant , and L. Ying, “Heavy traffic optimal resource allocation algorithms for cloud computing clusters,” *journal homepage: www.elsevier.com/locate/peva*, 2014.

[27] Dilip Kumar and B. Sahoo, “Energy Efficient Heuristic Resource Allocation for Cloud Computing,” *Computer Science & Engineering, National Institute of Technology*, 2014.

[28] Mansoor Alicherry and T. V. Lakshman, “Net work Aware Resource Allocation in Distributed Clouds,” *Proceedings IEEE INFOCOM*, 2012.

[29] Mina Sedaghat, Francisco Hernández-Rodriguez, and E. Elmroth, “Autonomic Resource Allocation for Cloud Data Centers: A Peer to Peer Approach,” *2014 IEEE International Conference on Cloud and Autonomic Computing*, 2014.

[30] Giuseppe Portaluri, Stefano Giordano, Dzmitry Kliazovich, and B. e. Dorronsoro. (2014, A Power Efficient Genetic Algorithm for Resource Allocation in Cloud Computing Data Centers. *2014 IEEE 3rd International Conference on Cloud Networking (CloudNet)*.

[31] W. ju-Hu, “Research of Resource Allocation in Cloud Computing Based on Improved Dual Bee Colony Algorithm,” *International Journal of Grid Distribution Computing*, 2015.

[32] Wanneng Shu, W. Wang, and Y. Wang, “A novel energy-efficient resource allocation algorithm based on immune clonal optimization for green cloud computing,” *Shu et al. EURASIP Journal on Wireless Communications and Networking 2014*, 2014.