

# The Combined Method for Load Distribution in Cloud Computing

<sup>1</sup>Sakshi Singh, <sup>2</sup>Pradeep Tripathi

<sup>1</sup>Research Scholar, <sup>2</sup>Assistant Professor  
Department of Computer Science  
Vindhya Group of Institutions  
VITS, Engineering in Satna, India.

## Abstract-

Cloud registration enables the flow of information and provides users with valuable resources. Customers are only billed for the use of resources. Cloud computing is a technology that saves data and ensures that information remains accessible. When situations are transparent, the propensity for information hoarding increases rapidly. Stack adjustment serves as a test specifically designed for cloudy weather conditions. Load adjustment is a process that evenly distributes the dynamic workload across hubs in order to avoid overloading. It facilitates the process of legalising assets. Additionally, it enhances system performance. The bulk of the presently available calculations enable stack modification and enhanced asset utilisation. Cloud computing utilises memory, central processing unit (CPU), and system stacks. The load adjustment system identifies hubs that are carrying excessive load and redistributes the additional burden to hubs that are carrying less load. Load balancing distributes workloads around the cloud data centres of a shared system to ensure that none of them are overwhelmed or underutilised. This study presents a proposed approach that combines Honey Bee (HB) with Particle Swarm Optimisation (PSO) to achieve an acceptable response time and implement a load balancing strategy. The hybrid method was evaluated using the CloudSim simulator. The load balancing approaches of Honey Bee (HB) and Particle Swarm Optimisation (PSO) outperform the hybrid approach. The hybrid algorithm's enhanced responsiveness is seen in its accelerated response time. This study examines the response time, request processing, data centre utilisation, and cost of virtual machines via the use of a simulator.

**Keywords - CPU, CC, Load sharing, SaaS, PaaS, IaaS, PSO, HB.**

## I. INTRODUCTION

Cloud computing (CC) is yet another example of cutting-edge technology. It provides the customer access to their online assets and storage space. It provides every piece of information at a more affordable price. Customers using cloud computing have continuous web-based access to their stored assets. They are responsible for paying for just the portion of the assets they utilize. In cloud computing, the cloud provider outsources every asset to the company they serve as a client. Cloud computing now has a lot of problems that need to be fixed. Adjusting the stack is the primary challenge presented by cloud computing. The load adjustment process moves all the loads between the various hubs in the system. In addition, it assures that every registered asset is distributed effectively and reasonably. It helps mitigate framework bottlenecks that may occur due to load lopsidedness, and it does this by providing support. The customers report a high level of satisfaction as a result. The approach of load adjustment is only somewhat new, yet it offers great asset utilization and improved response time. [1] [2] [3] [4] Customers gain a great deal in a variety of ways by using cloud processing.

**A. Cloud computing may be broken down into its parts, which are as follows: [5] [6].**

- Provides services to customers based on their requests. Cloud registration provides clients with on-demand benefits.
- Users can access the administration whenever they need it.

- Capabilities to access a Broad Network Through the System Cloud computing capabilities may be accessed through the system.
- Access to each of the capabilities may be gained through various means.
- Instantaneous Elasticity: The number of assets may be increased at any time following the customer’s requirements.

**B. Obstacles in the Field of Cloud Computing**

Cloud computing presents several challenges, including the following:

1. Safekeeping
2. Skillful manipulation of the burden
3. Monitoring of the Execution
4. Discussions about Reliable and Robust Service options
5. Asset Scheduling
6. Administration of scale and quality of service
7. Requires an Internet connection with high bandwidth and a fast speed.

**II. CLOUD COMPUTING MODEL**

The Cloud Figuring Model is shown in Figure 1, which includes various organization types and cloud-based services.

A. The Services Provided by Cloud Computing Administration refer to the numerous uses made available by a network of servers in the cloud. A multitude of services is made available to customers via cloud computing. [7]

**1) Software as a Service (SaaS):** SaaS made it possible for the buyer to access all of the applications made available by the seller. Applications are now operating on a framework that is hosted in the cloud. Access to the programs is provided by various interfaces, such as an internet browser. The ability to take in new items is not within the shopper’s control. [8] [9]

Layer	Cloud computing Component
<p style="text-align: center;"><b>Five Characteristics</b></p>	<div style="display: flex; flex-wrap: wrap; justify-content: space-around;"> <div style="border: 1px solid black; padding: 5px; margin: 5px;">On – demand self service</div> <div style="border: 1px solid black; padding: 5px; margin: 5px;">Broad Network Access</div> <div style="border: 1px solid black; padding: 5px; margin: 5px;">Resource Pooling</div> <div style="border: 1px solid black; padding: 5px; margin: 5px;">Rapid elasticity</div> <div style="border: 1px solid black; padding: 5px; margin: 5px; width: 100%;">Measured Service</div> </div>
<p style="text-align: center;"><b>Three Delivery Model</b></p>	<div style="display: flex; justify-content: space-around; align-items: center;"> <div style="border: 1px solid black; padding: 5px; margin: 5px;">LAAS</div> <div style="border: 1px solid black; padding: 5px; margin: 5px;">PAAS</div> </div> <div style="text-align: center; margin: 5px;"> <div style="border: 1px solid black; padding: 5px; width: 100%;">SAAS</div> </div>
<p style="text-align: center;"><b>Four Deployment Model</b></p>	<div style="display: flex; flex-wrap: wrap; justify-content: space-around;"> <div style="border: 1px solid black; padding: 5px; margin: 5px;">Public</div> <div style="border: 1px solid black; padding: 5px; margin: 5px;">Private</div> <div style="border: 1px solid black; padding: 5px; margin: 5px;">Community</div> <div style="border: 1px solid black; padding: 5px; margin: 5px;">Hybrid</div> </div>

**Figure 1 Model of Cloud Computing**

Customers unable to design their software but need custom applications may also benefit from using a software as a service (SaaS) platform. The following are examples of the administration’s uses for computer programming:-

- Customer resource management (CRM)
- Video conferencing

- The administration of the benefits of IT
- Financial Reporting
- Research on the World Wide Web
- Web content administration

**2) Platform as a Service (PaaS):** PaaS is short for “platform as a service,” and it works by providing customers with all of the resources necessary to construct apps. It offers all conceivable administration that may be found online. The user is exempt from having to download and install the product. Customers upload all of the applications to the cloud-computing platform. Regarding the development of apps, users have access to various tools and programming languages; the client does not have control over the arrangement of the servers, the operating systems, or the capacity. The purchaser retains complete authority over any applications they submit. Drawbacks.

**3) Infrastructure as a Service (IaaS):** The client is not responsible for managing or controlling the core cloud infrastructure with this cloud computing service. Initially, as an administrator client who was ready to manage all of the operating systems, storage, and programs provided to them. On the sections of the systems administration that deal with clients, there is only little control. Putting away and handling limits are within the responsibility of Foundation Providers. Virtualization is used to assign resources and then gradually resizes them to build systems that precisely meet the requirements of individual customers. Customers are responsible for sending in the product stacks necessary to operate their administrations. The provider will, upon request, arrange to provide the requested advantages. Clients only make use of these administrative services. It can be used to avoid acquiring, housing and managing the necessary equipment and programming foundation pieces, and it grows quickly to accommodate demand.

## B. The Multiple Layers of Service

Every single administration is comprised of several different levels. Which oversight is provided by the customers: -

### Cloud Deployment Models:

1. A public cloud is a cloud foundation controlled by an organization and made accessible to the general public or to a sizeable section of the population of a certain industry. Cloud computing resources are considered open when they are made accessible to anybody and everyone without limitation.
2. A private cloud is a kind of cloud computing in which a single organization uses the underlying cloud infrastructure at any moment. A private cloud is exclusively administered by the organization that owns it or a third party completely separate from the organization. The general public is not yet ready to use the private cloud extensively.
3. Community Cloud: The fundamental infrastructure of the cloud is a resource that several organizations share. The community cloud may be beneficial to a particular network with challenges common to other networks, such as security needs, strategic concerns, and consistency issues. An independent third party or perhaps the associations themselves may be in charge of monitoring it.
4. Hybrid Cloud: A hybrid cloud is formed when at least two different kinds of clouds, either open, network, or private, are combined. This remaining component, made up of one-of-a-kind substances, is, nevertheless, kept together by institutionalised innovation, enabling information and application mobility. For instance, employing cloud blasting to shift the stack between mists is one example of how this may be done.

## III. VIRTUALISATION

The term “virtualization” refers to things that do not exist in the real world, yet virtualization provides an experience identical to the real. Virtualization refers to using a computer to make it seem like it is doing an entire system’s tasks. Because users can access the many apps and services provided by the cloud thanks to virtualization, this component of the cloud environment is considered its most important aspect. In the cloud environment, many different forms of virtualization may be used.

Two sorts of virtualization are:

1. Full virtualization
2. Paravirtualization

**1. Full Virtualisation:** When referring to full virtualization, it is important to understand that a whole computer is installed on another machine. That virtual computer provides access to the full capabilities of the primary machine. When the client's actual computer is not available, the offices employ the virtual machine instead.

**2. Para Virtualization:** Refers to a configuration that permits several operating systems to coexist on a single piece of hardware. In addition, it makes it possible to efficiently use the system's resources, such as the memory and the CPU.

#### IV. LOAD BALANCING

The term "load adjustment" refers to dividing a bigger processing load over a greater number of smaller preparatory hubs to increase the framework's overall performance. To maintain a uniform distribution of the dynamic local workload across all of the hubs in a distributed computing system, it is important to make appropriate adjustments to the condition stack. [10][11][12][13]

- Load adjustment helps properly identify registered assets, leading to increased user satisfaction and legal resource use. This is done via the suitable allocation of available resources. Appropriately changing the load may help limit the usage of assets, particularly when the consumption of such assets is considerable.
- Load adjusting is a strategy that has been useful to systems and assets since it has delivered a maximum throughput with the least reaction time. It also helps execute bomb over, preserves flexibility, and retains a strategic distance from bottlenecks. By distributing the burden among all of the servers that are part of the system, stack adjusting makes it possible to send and receive information instantly. Adjusting the load is the name given to this operation.

- When clouds are present in the sky, several different algorithms may be used to aid in assessing real rush hour congestion. A load balance is maintained among all accessible servers. Most of them can be connected to the cloud environment with the required confirmations. In distributed computing, batch mode heuristic planning calculations and online mode heuristic computations can be partitioned into groups for condition stack modifying calculations. The first calculations are known as Batch mode heuristic planning calculations (BMHA), and the subsequent calculations are known as online mode heuristic calculations. Only when jobs have a point of interaction with one another in the framework, do they get merged in BMHA. Following the establishment of the day and the age, the BMHA planning calculation will then get underway.

It's first come, first served for everything. Calculations based on the BMHA include but are not limited to, First-Come, First-Served scheduling calculations (FCFS), Round Robin booking calculations (RR), Min calculations, and Max-Min calculations. When the heuristic calculation for online mode booking is used, every task is planned at the time it is now contacting base in the framework. As a result of the cloud environment being a heterogeneous structure, the speed at which each processor functions may alter unexpectedly and without any effort being exerted. The heuristic calculations in the online mode are more suitable for the cloud environment and function more effectively there.

- When developing a calculation for adjusting the load that is on a heap, it is essential to measure the appropriate amount of load, conduct an inspection of the entire heap, ensure the safety of each system, and ensure the successful operation of the intended system, ensure the connection between all of the hubs, and determine the nature of the work that will be traded. The most significant aspect of this procedure is the selection of the hubs that will be included and the provision of a wide range of unique ones. The size of the system's heap is based on a combination of the stack on the CPU and the required quantity of RAM.

- One such illustration of load balancing in our day-to-day lives may be observed at several sites. In the absence of load adjustment, clients risk encountering various problems, such as deferrals, timeouts, and sluggish system answers.

1) The static approach: This method is often characterized by either the planning or the actual execution of the framework. It is also known as the “traditional” approach. The algorithms that are used to alter the static load distribute the movement among all the servers proportionately.

2) The dynamic method: This approach considered the system’s present state before making judgments on the stack modification. When working with Cloud frameworks like distributed computing, a dynamic approach is the best way to go about things.

The dynamic load adjustment algorithms are comprised of two distinct components altogether.

The first method is known as the cloud strategy, while the second is called the integrated approach and does not use the cloud. The following is a list of features that it possesses:

a) The Centralised Technique: In a system that uses the centralized method, there is only one central hub that is responsible for controlling and transmitting information across the whole system. There is not a single one of the other hubs that bears any responsibility for this matter.

b) The Cloud Method: Each hub is responsible for independently creating its own heap vector in the Cloud technique. At this time, Vector is gathering the heap data from several different hubs. Every choice is made at the neighborhood level with the assistance of the stack vectors that are close by. The cloud technique is best suited for cloud-based frameworks that are mostly cloud-based, such as distributed computing.

### **B. Metrics for load balancing is as follows:**

1. Throughput: This measure is used to assess whether or not all of the tasks whose execution has been finished have been completed. 2. Productivity: This statistic evaluates how efficiently a process operates. The efficiency of the operation of any framework may be considerably enhanced by increasing the throughput.

2. Fault Tolerance: This implies being able to bounce back after experiencing a setback. The stack adjustment should use a fault-tolerant strategy that is at a satisfactory level.

3. Migration time: - It is a fair opportunity to migrate the professions or assets from one hub to numerous hubs. - It is the time between migration windows. - It is the period between migration windows. Migration takes place whenever there is a relocation of one or more of the network’s hubs. It should be constrained with a certain end goal, and that end objective is improving the framework’s performance.

4. Response Time is the time it takes for a specific calculation that modifies the load to reply to an assignment set in a framework. To improve how a framework is put into action, the range of this parameter needs to be constrained.

5. Scalability refers to the capability of a computing system to carry out load adjustment for any number of hubs included inside a framework. The concept of scalability refers to the same thing. This metric has to be strengthened for the framework to be regarded as acceptable.

### **C. Principles of the algorithm for load balancing**

Stack adjustment computations make use of a variety of different methodologies, including the following: [14] [15]

- Guidelines for information policy: It specified what data are necessary and how they should be acquired appropriately. In addition, it is further defined when these statistics are compiled.

- Resource type policy: This method describes the large variety of resources available throughout the heap adjustment process.

- Selection policy: This method is applied to locate the assignment that transfers from an overburdened hub to a free hub.

### **D. The primary purposes served by load-balancing algorithms**



1. Efficient use of resources: load balancing may improve system performance while reducing operational expenses.
2. For future stack adjustment calculations, it is important to ensure they are scalable and flexible. So the calculation needs to consider these kinds of things. Because of this, the calculation process must be adaptable and sensitive.
3. Priority: The prioritization of assets or employment must be completed. Therefore, greater necessity employments are exhibiting signals of improvement opportunities to carry them out.

## V. EXISTING LOAD BALANCING ALGORITHMS

Many different load-adjusting formulas may assist you in getting higher throughput and improve your response time while working in cloud conditions. Every one of the computations comes with its own set of benefits. [16] [17] [18]

1. Task Scheduling Based on LB: This calculation primarily consists of a two-level task planning component dependent on stack modification to satisfy customers' demanding requirements. It achieves a high use of assets. This computation brings about the desired effect of stack adjustment by mapping assignments to virtual machines and then determining whether or not all virtual machines own assets. The response time for the errand is being improved as a result. It equally provides improved asset utilization.

2. Opportunistic Load Balancing: OLB attempts to keep every hub busy and does not consider the current amount of work done on any individual computer. OLB assigns each errand submitted as a free request to the exhibit hub of helpful. The preferred perspective is very easy and can accomplish stack adjustments. Still, its flaw is that it does not account for every wish execution time of the task, which results in the overall completion time (Make range) being incredibly bad.

3. Round Robin: - In this computation, each operation is carried out in a separate thread on each processor. Each operation is passed on to the processor like a round robin in this arrangement. No difference in the work stack dissemination occurs across processors. Different methods do not need the same amount of time to handle employment. When web servers have HTTP requests similar to Cloud computing, RR computation is used. Some hubs may be very stacked at certain times, while others may remain idle. The time quantum is an extremely important component of the Round Robin Scheduling method. When a significant amount of time is available, the RR Scheduling Algorithm is identical to the FCFS scheduling. In addition, the Round Robin Scheduling algorithm is referred to as the Processor Sharing Algorithm when the time quantum available is insufficient.

4. Randomised: The nature of this computation is to be considered static. Within the context of this computation, a process may be managed by a particular hub  $n$  with a probability  $p$ . This computation will work wonderfully when all operations have been placed in the same order as before. The problem arises when the workloads have varying degrees of computational complexity. The deterministic method isn't being followed by this computation at all.

5. Min-Min The algorithm starts with an arrangement of every assignment that has not yet been allocated. The time needed to complete any errand may be determined using this basis. At that point after that, among these basic occurrences, the base worth is decided. At that point, design the errand on the machine to take the least amount of time. After that, the execution time for every other task is updated on that machine, and a similar technique is followed until all tasks have been assigned to the available resources. This calculation's primary concern is that it has a deficiency in starving.

6. The max-min algorithm The max-min calculation and the min-min calculation are similar. The following is the major contrast: In this computation, finding the least amount of time needed to complete tasks comes first. After that, the most significant value is determined, the amount of time required to complete all tasks using any resources. After the most extensive amount of time was spent determining it, the task was assigned to the chosen machine. [19] After that, the time it takes to complete each task on that computer is brought up to date. This is accomplished by adding the amount of time it takes to complete the assigned task to the total time it takes to complete each of the other tasks on that machine. When this occurs, any task that has been delegated is removed from the list of tasks that are being carried out by the system.

Honey bee foraging behavior is an example of a self-association algorithm inspired by the natural world. Bumble bees can achieve global load adjustment via local server operations. The framework's functionality

has been improved with the expansion of the diversity of framework descents. The primary problem is that the throughput does not increase proportionally with the size of the framework estimate. This calculation is the most suited when it is necessary to have a diverse population that composes the administration.

8. Active Clustering: This method of computation involves bringing hubs of the system that have the same composition together so that they may work together as a group. A system is rewired to adjust the load placed on the framework. This approach functions similarly to a self-total load adjusting method. By bringing together several administrations that perform similar functions, frameworks make it easier to use comparable task assignments. The performance of the framework was improved with the addition of enhancements. When each asset is effectively used, the operation's throughput may be increased.

9. Compare and Balance: - This computation produces a harmonious condition and manages stacks of uneven systems. Based on the probability number of virtual machines operating on the current host and overall cloud framework, the current host chooses a host randomly and considers their heap. If the heap of the current host is more than what they decided to have, it transfers the extra heap to that particular hub. At that moment, every host of the framework will carry out a strategy comparable to the one described. This computation for changing the heap is also planned to reduce the time spent relocating virtual machines. Shared capacity memory is implemented to reduce time spent relocating virtual machines.

10. A multiprocessing approach for LB that does not need locks: It offers a multiprocessing load adjusting arrangement that does not use shared memory and avoids the use of bolt as an alternative to conventional multiprocessing load adjusting arrangements that use shared memory and bolt to keep a client session active. The bit must be adjusted to do this. Running several load-adjusting forms inside a single load balancer allows this arrangement to improve the general performance of load balancers when used in conjunction with several courses.

11. The Optimisation of Ant Colonies: - Calculations based on ants are one approach to solving complex problems involving combinatorial improvement in several ways. A good example of this strategy is the voyaging salesman problem (TSP) and the quadratic task problem (QAP). These computations were given more life because of the impression of real insect communities. The behavior of subterranean insects is mostly motivated by the desire to survive in their environments. They do not consider the individual.

12. Fastest Response Time Comes First: The ability to do this calculation is not complicated. Each process is given a necessity in this scenario, after which it is permitted to execute. For this particular FCFS arrangement, identical need forms are scheduled to be used. The computation of the (SJF) is an exceptional example of universal need Scheduling calculation. The computation that is required for SJF is the inverse of the CPU burst that will follow. That is to say, if the CPU burst is allowed to continue for longer, the requirement will decrease. The SJF technique prioritizes completing the task that requires the least time to prepare for it. For the sake of this computation, shorter employments are carried out a short period after longer employments. When working with SJF, knowing or estimating the amount of time needed to handle each task is essential since this is the true difficulty with SJF.

13. Based Random Sampling: This computation depends on building a virtual chart with a network connecting all of the hubs of the framework, with each hub of the diagram being compared to the hub PC of the cloud framework. Incoming edges and active edges are the two types of edges that can be found between hubs. These edges are utilized to consider the load of a certain framework and apportion the hub's assets. [20] It is a fantastic method for adjusting the stack.

## VII. LITERATURE REVIEW

As a means of highlighting some of the current CC concerns and challenges, we explore what CC is and its many services in the first part. Second, based on the quality of service provided by CC, we identify several security threats. We highlight several remaining issues from the perspective of Cloud Computer Discovery and its long-term ramifications. This book briefly overviews the cloud platforms now available for research and development [21].

We propose using a method called "load allocation," which is conceptually very close to the load balancing function. This study's subject is liquid Galaxy, an open source project aiming to emulate Google Earth and other applications using several virtual computers [22].

The simulation results are compared to several previously suggested cloud load balancing approaches. Simulations show that jobs are dynamically divided across various available virtual machines of different configurations in different data centers such that relatively superior reaction times and makespan times may be obtained [23].

To get around the scheduling issue and increase throughput and resource usage without negatively impacting the CC platform's overall results, we used CS-SS load balancing and grasshopper optimization using MapReduce [24].

To overcome the problem associated with existing meta-heuristic methods, the proposed methodology investigated the MakeSpan parameters. Based on mutation, the Particle Swarm algorithm is used in the proposed approach to distributing work equally throughout data centers. Cloud computing's fitness function may be improved by reducing performance indicators like MakeSpan time and using an efficient load balancing strategy [25].

This article examines and evaluates existing cloud load balancing solutions. They're all compared in the system's state load-balancing algorithms to get the best result. This article examines the reliability, reaction time, adaptability, performance, resource utilization, and fault tolerance of various systems and services. System performance is enhanced by these adjustments [26].

This research examines the usefulness and limitations of benchmark load balancing techniques. Additionally, methods for opportunistic load balancing (OLB) and load balance min-min (LBMM) scheduling are offered. Analytical model findings and CloudSim simulation results are examined for validity [27].

As a result, HEC-Clustering Balance beats other techniques to load balancing. In two experiments, we reduce HEC server processing time by 19% and 73% [28].

These tactics are examined utilizing a cloud analyst simulator in this paper. Examine three approaches to load balancing (Round Robin, Throttled, and Active Monitoring). Algorithms used in cloud data centers include service broker techniques [29].

To begin, a mechanism for distributing load among sBSs is provided. Electrocardiogram-based encryption and the decryption key are utilized in an advanced encryption standard (AES) cryptographic technique for further security. If you combine load balancing with carbon offset (CO), you can save time and money. An analysis of the results shows that, in comparison to local execution, our technique saves between 68% and 72.4 percent of system utilization.

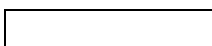
Cloud computing and load balancing enabled businesses to overlook network traffic and bad workload by distributing the load logically among virtual servers. VM load-balancing slopes were lowered as a result. Load-balancing features are encrypted from the bottom up using the Bat algorithm. A meta-heuristic algorithm is what we call our approach [31].

Using prototypes built using CloudSim 4.0 and Amazon Web Services, the researchers investigated three approaches to load balancing: First Come, First, Shortest Job First, and Least Connection First. Space-shared scheduling, rather than time-shared scheduling, performs better in simulations. First come, first served, and shortest jobs perform worse when the workload is lowered [32].

Researchers address the challenges of LB in clinical trials and the need for a new LB approach that uses functional tissue thermometry (FT). Because FT efficiency factors are not included in traditional LB algorithms, they are inadequate. The study's findings show FT efficiency measurements in LB algorithms are necessary, raising concerns about the cloud. [33] This paper suggests a brand-new FT-based LB algorithm.

Efficiencies, resource utilization, quality of service (QoS), and performance should all be considered while developing scheduling algorithms. Academics have developed several scheduling methods to deal with this issue. Job scheduling and resource usage in cloud computing is the focus of this research [34].

Experiments show that workload prediction lowers the number of SLO violations and migrations and enhances the data center load balancing performance. The RL-based VM migration technique outperforms the heuristic-based solution in heavily loaded systems [35].





Parameters	Table 1. Comparison of several LB approaches, depending on specified parameters					
	Round Robin [36]	Max-Min[36]	Throttled [36]	Shortest Job Scheduling [36]	Active Clustering [36]	Ant Colony Optimization [36]
Performance	Yes	Yes	Yes	No	No	Yes
Throughput	Yes	Yes	No	No	No	No
Overhead	Yes	Yes	No	No	Yes	No
Fault Tolerance	No	No	Yes	No	No	No
Migration	No	No	Yes	No	Yes	Yes
Response Time	Yes	Yes	Yes		No	No
Resource Utilization	Yes	Yes	Yes	Yes	Yes	Yes
Scalability	Yes	No	Yes	No	No	No
Power Saving	No	No	No	No	No	No

### VIII. RESEARCH GAP

- Several LB procedures must be carried out precisely, and particular algorithms must be designed. When designing these algorithms, it is important to consider many factors such as control rates, complex thresholds, fine-grained relocation costs, contact and data transmission lengths, and overhead.
- Three examples of the computational overheads associated with several fundamental processes are the migration of virtual machines and jobs and the monitoring of devices. It's important to keep them under control and organized as well.
- Workload forecasting algorithms must be improved to properly anticipate future overload or underload situations far ahead of the period in question.
- In general, load balancing algorithms aim to improve performance while simultaneously decreasing operating expenses. Consequently, finding an effective solution to the various competing goals is imperative.
- The existing approach to see whether an algorithm works in a "real world" cloud environment is to build it in a "real world" cloud environment.

### IX PROPOSED METHOD

The queue load on accessible cloud computing is generally balanced using several approaches, and numerous of these methods exist. Several of which are discussed in this article; are as follows:

**A. In a Circular Motion** The round-robin algorithm is the simplest and easiest to implement of all the algorithms used in cloud computing for load balancing. It is also one of the most often used algorithms. as for numerous variations of the round robin algorithm that are accessible, each of which focuses on the measurement of a particular parameter, and here is the terminology that is common with this method [37]:

- Burst Time, often known as BT, is the needed amount of time to complete the request.
- Time Quantum (TQ) refers to the amount of time computed based on a request and allotted to access the VM.

**B. Reduced in Volume** This technique takes advantage of the system that is available in the VM status list (BUSY / AVAILABLE.) when the load balancer receives the request from the user. After displaying the available virtual machines (VM), it checks to identify the available VM and then assigns the resource to the free VM. Consequently, implementing these algorithms does not involve a particularly high degree of complexity, yet their performance is significantly improved [38].

**C. ESCE (Equally Spread Current Execution Load)** [39] This method operates constantly depending on the queue and distributes the work by passing it to a different virtual machine at random. After that, it distributes the work after determining the magnitude of the incoming load, and then it distributes the work to a virtual machine with light work. The approach utilized in this experimental paper will be taken from one of the algorithms described before (ESCE) and implemented in the CloudAnalyst simulator.

**D. Honey Bee Algorithm** [40]: This algorithm is based on how honey bees search for and collect food; their behavior inspired it. Forager bees are the kind of bees that go out into the world in search of new food sources. When they locate a new food source, forager bees return to the hive and perform a waggle dance to let the other bees know about it. The performance of this dance indicates the quality or amount of food and the distance from the hive to which it is located. The scout bees then follow the foragers to the food site and begin to harvest it when they arrive there. They then go back to the hive and perform a waggle dance, which estimates how much food is still available and, consequently, leads to either further consumption of the food source or its abandonment. The Honey Bees Algorithm in its most basic form.

- 1 for  $i=1, \dots, ns$  in which case:
  - i. scout[i]=Initialise scout()
  - ii. lower patch[i]=Initialise flower patch 1 for  $i=1, \dots, ns$  in which case: (scout[i])
    - 2 do until stopping condition=TRUE
    - I Recruitment ()
  - iii. for  $I = 1, \dots, nb$ 
    - flower patch[i] is equivalent to the local search phrase flower patch[i]
    - flower patch[i] is synonymous with "site abandonment" (flower patch[i]).
    - Neighborhood shrinking=flower patch(i), also known as flower patch(i)
  - iv. for  $I = nb, \dots, ns$ 
    - flower patch[i] equals Global search(flower patch[i])

As was previously noted, the Bees Method is an optimization algorithm that gets its name from the natural foraging activity of honey bees. The algorithm aims to identify the pseudo-code that corresponds to the best possible solution in its most basic form. The algorithm calls for several parameters to be set, including the initial size of patches (ngh), which includes the site and its surrounding area, the number of scout bees (n), the number of sites selected out of the n visited sites (m), the number of best sites out of the m selected sites (e), the number of bees recruited for the best e sites (nep), the number of bees recruited for the other (m-e) selected sites (ns). The first step of the method involves scattering the n scout bees around the search space randomly. In the second stage, the fatness of the locations that the scout bees have investigated is assessed.

- v. Start the population out with a selection of random solutions.
  - vi. Determine the overall health of the population.
- While (the halting requirement has not been satisfied), fresh populations are being formed.
- viii. Pick locations for the search in the neighborhood.
  - ix. Recruit bees for the specified places (more bees for the better sites) and assess their fitness for the job.
  - x. Choose the healthy and strongest bee from each patch.
  - xi. Assign the remaining bees to seek in a random pattern and assess their fullness level.
  - xii. End While.

In the fourth stage, "selected bees" are picked from the bees with the greatest fitnesses. The places these bees have visited are then chosen for the neighborhood search. After that, in stages 5 and 6, the algorithm performs searches in the vicinity of the sites that have been picked, allocating a greater number of bees to look for the finest e sites nearby.

### E. PSO algorithm [41]

1. Start

2. Randomly determine the starting positions and velocities of n different particles.
3. Determine the fitness of each particle and record the result as pbest .
  - 3.1 Use the pbest with the highest score as the gbest.
  - 3.2 While (if not meet the requirement)
    - 3.2.1 Bring the inertia weight W up to date (t)
    - 3.2.2 Make sure that the velocities and positions of each particle are up to date.
    - 3.2.3 Determine its level of usefulness.
    - 3.2.4 If position fitness is higher than pbest, then pbest should be updated.
  - 3.3 Conclude the while
4. Finish up while
5. Return gbest

## F. Proposed Hybrid algorithm

1. Start
2. Figure out which of the Vms is loaded and underloaded.
3. Determine which virtual machines are underloaded and which are overcrowded; reschedule cloudlets using the PSO algorithm
4. Determine the starting positions and velocities of all n particles.
5. Determine the value of fitness for each particle.
6. If the overall value of the workout is more beneficial than the prior pbest.
7. Use the current value of fitness as the new pbest value.
8. Return to steps 5 and 6 for each particle.
9. Classify the most desirable fragments as gbest.
10. Determine the velocities of all the particles and bring their locations up to date.
11. Exit if this is the maximum iteration.
12. Else, Repeat steps 5 & 6
13. End

## X. RESULT

In this scenario, for our software using eclipse OXYGEN.1 with Java programming language, and for hardware used in the form of Intel® Core™ i7 Processor with 8GB RAM, we assume the use of the large-scale internet that is on the application, such as Youtube, Facebook, Instagram, and other similar websites. In this particular scenario, we will suppose that more than 500 million people all over the globe utilize the program, which is only carried out in a single data center.



Figure 2: Scenario of Cloud Computing

Figure 2 depicted the several network scenarios that needed to be carried out in light of the fact that this experiment only used a single data centre and received requests from many User Bases. In this experiment, the method used is called Throttled on the load balancer.

In this case, we will suppose that the clock was accurate to a two-hour accuracy at the time. Evening usage of the app is common among users. Additionally, it is assumed that the user loads a new request at regular intervals of five minutes.

Table 2: Variables Relating to the User Population

User Base	Region	Time Zone	Peak Hours (GMT)	Peak Hours (Local Time)	Online Users Simultaneously During Peak Hour	Online Users Simultaneously During Off-Peak Hours
UB1	0	GMT - 6.00	14:00-16:00	8:00-10:00 pm	500,000	50,000
UB2	1	GMT - 4.00	16:00-18:00	8:00-10:00 pm	200,000	20,000
UB3	2	GMT + 1.00	21:00-23:00	8:00-10:00 pm	400,000	40,000
UB4	3	GMT + 6.00	02:00-04:00	8:00-10:00 pm	250,000	25,000
UB5	4	GMT + 2.00	22:00-24:00	8:00-10:00 pm	150,000	15,000
UB6	5	GMT + 10.00	10:00-12:00	8:00-10:00 pm	180,000	18,000

Table 3: Taking on the Cost Assumption

Cost per 1 Gb of Data Transfer (from/to Internet)	\$0.11
Cost VM per hour (1024Mb and 100MIPS)	\$0.11
Memory Cost	\$0.05
Storage Cost	\$0.1

Table 2 presents the information on each user’s base and the time the active user uses the cloud server simultaneously.

As seen in Table 3, for this experiment’s sake, we will assume a strategy that considers the prices currently being charged by cloud service providers like Amazon EC2. It can affect the performance produced by the virtual machine and the amount of pleasure the customer will experience because we make adjustments according to the real-world pricing criteria in our simulation.

Table 4 demonstrates the assumptions made at the level of the virtual machine specification and the physical machine that will be utilized in this experiment. These assumptions will be fed into the cloud computing simulator engine.

Table 4: A list of the data center’s physical parameters

Parameter of Data Center	The Value Used in The Simulation
VM Image Size	100,000
VM Bandwidth	1,000
VM Memory	1024
Architecture of Data Center	X86
OS of Data Center	Linux
VMM of Data Center	Xen
Number of Machines of Data Center	20
Memory per Machine of Data Center	4096
Storage per Machine of Data Center	100,000
Available per BW of Data Center	10,000
Number of Processors per Machine of Data Center	4
Processors Speed of Data Center	100
VM Policy of Data Center	1,000
Grouping Factors on Request	100
Executable Instruction Length	250

Table 5 presents the assumed value of the matrix’s delay (milliseconds)

Region/Region	0	1	2	3	4	5
0	35	110	160	260	260	110
1	110	35	260	510	360	210
2	160	260	35	160	160	210
3	260	510	160	35	510	510
4	260	360	160	510	35	510
5	110	210	210	510	510	35

The latency of the internet is shown in Table 5, which is one of two categories in the matrix that describes the properties of the currently operating Internet network. It will impact the simulation because it is calibrated to the real state.

Table 6 presents the assumed value of the matrix's delay (Mbps)

Region/Region	0	1	2	3	4	5
0	3000	1500	1500	1000	1000	1500
1	1500	700	1000	1000	1500	1000
2	1500	1500	2500	1000	1500	1000
3	1000	1000	1500	2000	1000	1500
4	1500	1000	1000	1500	600	1000
5	1000	1500	1000	1000	1500	2500

Table 6 presents the bandwidth parameters as one of the two groups of matrix parameters found on the current operational Internet network. It will impact the simulation because it is calibrated to the real state.

Table 7 presents an overview of the overall response times.

	Avg (ms)	Min (ms)	Max (ms)
<b>Overall response time</b>	534750.97	564.98	1062618.97
<b>Data Center processing time</b>	534515.92	500.24	1062540.24

Table 7 summarizes the average reaction time and the time required to analyze the data for the simulated output utilized in this experiment. This is reflected in the algorithm's simulation level efficiency in the load balancer to equalize user requests to the Data Center.

Table 8 displays the average response time for each location.

Userbase	Avg (ms)	Min (ms)	Max (ms)
UB1	532309.45	564.98	1062618.97
UB2	533961.01	34633.81	1028050.50
UB3	540151.29	45705.30	1037462.50
UB4	536031.90	49007.53	1043035.70
UB5	526644.95	69799.72	1011642.72
UB6	534195.94	37186.35	1025327.45

Table 8 summarizes the average reaction time based on each section of simulation output utilized in this experiment. Additionally, we can see the outcomes of the algorithmic efficiency level simulation carried out by the load balancer. This simulation was carried out to equalize the requests sent from users to the Data Center.

Table 9 displays the processing timeframes for data center requests

Data Center	Avg (ms)	Min (ms)	Max (ms)
DC1	524515.92	500.24	1062540.24

How the data center deals with requests from users is laid out in table 9, and it is based on the typical amount of service time. There also is the range of minimal service time up to the greatest service time, from which the average is determined.

Table 10 shows the total amount of money spent on the virtual machine (\$)

Data Center	VM Cost \$	Data Transfer Cost \$	Total \$
DC1	0.88	60.51	61.39

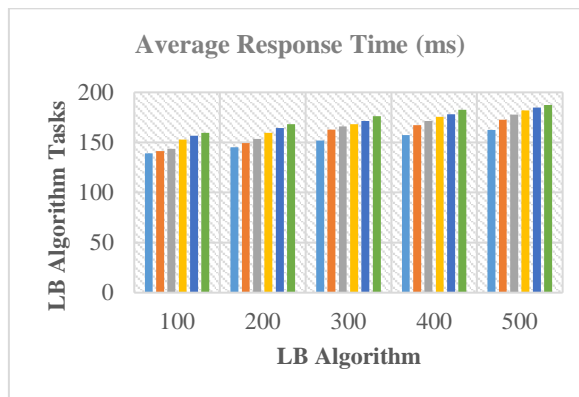


Table 10 presents the overall expenses incurred by a single data center, which are connected to the total amount of expenditures that are based on the costs associated with virtual machines and the costs associated with data transmission.

Table 11 presents the findings of an analysis that compared the performance of Honey Bee (HB) and Particle Swarm Optimization (PSO) algorithms to determine the average reaction time for a scenario involving fifty virtual machines, each of which had a varied amount of tasks to be evaluated.

Table 11: Average Response Time Results

Average Response Time (ms)						
LB Algorithm Tasks	Hybrid	PSO	Honey Bee	Round Robin	Throttled	ESCE (Equally Spread Current Execution Load)
100	138.95	141.35	143.72	152.94	156.82	159.6
200	145.24	149.34	153.62	159.45	164.32	168.3
300	151.82	162.78	165.82	168.21	171.28	176.3
400	157.22	167.29	171.52	175.56	178.24	182.5
500	162.47	172.51	177.85	181.97	184.69	187.3



Graph 3 Displays a Comparison of the Average Response Times of Different Load Balancing Algorithms

Honey Bee (HB), Round Robin, Throttled, ESCE (Equally Spread Current Execution Load), and Particle Swarm Optimization (PSO) has a lower average response time than the hybrid algorithm of HB and PSO, which shows that the hybrid algorithm is superior in terms of quality and effectiveness. This is because the hybrid algorithm is very straightforward and does not require any additional computation. We set up a data center with 50 virtual machines, and the number of tested cloudlets was equal to twice as many VMs. This is because the evaluated innovation allows the cloudlets to be rescheduled and carried out on particular VMs. PSO was utilized to randomly update the inertia weight for an improved performance after the fitness value was updated from the first round of scheduling on VM fitness based. This was done to improve overall system functionality. In this particular setting, hybridization works better than the conventional single population-based scheduling, and no single node is accountable for the entire scheduling decision. Since the goal is to enhance response time with a proportionate increase in the number of jobs, the burden should be distributed fairly across all of the virtual machines that are accessible.

## XI. CONCLUSIONS

DDistributed computing primarily handles the management of software, data retrieval, and storage capabilities, without requiring the end-user to have knowledge of the physical location and structure of the system providing these services. Stack adjustment is a significant consideration in the realm of distributed storage. It enhances the efficient utilisation of resources, hence enhancing the performance of the framework. By using recent calculations, it is possible to ensure proper stack modifications and enhance systems via more efficient booking and asset allocation procedures. This article introduces the notion of cloud computing, as well as the process of stack modification. The key concern in this context is the calculation of stack adjustment. In the realm of distributed computing, several computations have previously been enumerated, including various characteristics such as flexibility, enhanced asset utilisation, and improved reaction time. The scheduler finds the parameter ratio with the highest maximum value by taking into account the utilisation of RAM, Bandwidth, and CPU resources throughout the interval processing time. Instead of focusing on the CPU load, this strategy may determine which parameter has to be altered to contribute to the attained high utilisation. Among the several tactics considered, the recommended strategy resulted in a reduced workload in the specific location where it was applied, when used in a general context. The demonstration indirectly confirmed the validity of the performance analysis approach, suggesting that it may even surpass another method in terms of effectiveness.

## REFERENCES:

1. Mall, Peter and Grance, Tim, "The NIST definition of cloud computing", National Institute of Standards and Technology, 2009, vol53, pages50, Mell2009
2. <http://www.ancoris.com/solutions/cloudcomputing.html>, "Cloud Computing |Google Cloud - Acores."
3. <http://www.personal.kent.edu/~rmuhamma/E-Systems/Myos/prioritySchedule.htm>, "Priority Scheduling -Operating Systems Notes."
4. S. S. Moharana, R. D. Ramesh, D. Powar, "Analysis of load balancers in cloud computing" International Journal of Computer Science and Engineering, 2013, vol 2, pages 101-108
5. <http://blog.nextright.com/?cat=6> "Cloud Computing « Nexright Blog".
6. A. Sidhu, S. Kinger, "Analysis of load balancing techniques in cloud computing", International Journal Of Computers & Technology 4 (2) (2013) pages737-741.
7. <http://www.qualitytesting.info/group/cloudcomputing/forum/topics/software-as-a-s>, "Software as a Service (SAAS) - Quality Testing"
8. <http://letslearncloud.wordpress.com/> "Cloud Computing | Learn Cloud and its tips."
9. P. Gupta, M. Samvatsar, and U. Singh, "Cloud computing through dynamic resource allocation scheme," in Proceedings of the International Conference on Electronics, Communication and Aerospace Technology, ICECA 2017, 2017, vol. 2017-Janua, doi: 10.1109/ICECA.2017.8212723.
10. Chaudhari, Anand and Kapadia, Anushka, "Load Balancing Algorithm for Azure Virtualization with Specialized VM", 2013, algorithms, vol 1, pages 2, Chaudhari
11. Nayandeep Sran, Navdeep Kaur, "Comparative Analysis of Existing Load Balancing Techniques in Cloud Computing", vol 2, jan 2013
12. P. S. Chouhan, M. Samvatsar, and U. Singh, "Energetic SSource allotment scheme for cloud computing using threshold-based," in Proceedings of the International Conference on Electronics, Communication and Aerospace Technology, ICECA 2017, 2017, vol. 2017-Janua, doi: 10.1109/ICECA.2017.8212744.
13. Chaczko, Zenon and Mahadevan, Venkatesh and Aslanzadeh, Shahrzad and Mcdermid, Christopher, "Availability and load balancing in cloud computing", International Conference on Computer and Software Modeling, Singapore, chaczko2011availability
14. <http://www.writeaprogram.info/c/os-programs/priority-scheduling/> "Priority Scheduling Algorithm Example – Write a Program."
15. R. Alonso-Calvo, J. Crespo, M. Garcia-Remesal, A. Anguita, V. Maojo, "On distributing load in cloud computing: A real application for very-large image datasets", Procedia Computer Science 1 (1) (2010) pages 2669-2677

16. S.-S.Wang, K.-Q. Yan, S.-S.Wang, C.-W. Chen, "A three-phases scheduling in a hierarchical cloud computing network", in: Communications and Mobile Computing (CMC), 2011 Third International Conference on, IEEE, 2011, pp. 114–117.
17. O. Elzeki, M. Reshad, M. Elsouid, "Improved max-min algorithm in cloud computing, International Journal of Computer Applications" vol 50 (12) (2012) pages 22–27..
18. U. Bhoi, P. N. Ramanuj, "Enhanced max-min task scheduling algorithm in cloud computing", International Journal of Application or Innovation in Engineering and Management (IJAIEM), ISSN,2013,pages 2319--4847
19. P. Gupta, M. Samvatsar, and U. Singh, "Cloud computing through dynamic resource allocation scheme," in Proceedings of the International Conference on Electronics, Communication and Aerospace Technology, ICECA 2017, 2017, vol. 2017-Janua, doi: 10.1109/ICECA.2017.8212723.
20. P. S. Chouhan, M. Samvatsar, and U. Singh, "Energetic SSource allotment scheme for cloud computing using threshold-based," in Proceedings of the International Conference on Electronics, Communication and Aerospace Technology, ICECA 2017, 2017, vol. 2017-Janua, doi: 10.1109/ICECA.2017.8212744.
21. A. Jangra and H. Dubran, "Assessment of Load Balancing Techniques in Cloud Computing," 2021 6th International Conference on Signal Processing, Computing and Control (ISPCC), 2021, pp. 795-798, doi: 10.1109/ISPCC53510.2021.9609377.
22. V. Sivaraj, A. Kangaiammal and A. S. Kashyap, "Enhancing Fault Tolerance using Load Allocation Technique during Virtualization in Cloud Computing," 2021 7th International Conference on Advanced Computing and Communication Systems (ICACCS), 2021, pp. 1798-1801, doi: 10.1109/ICACCS51430.2021.9441779.
23. S. Swarnakar, R. Kumar, S. Krishn and C. Banerjee, "Improved Dynamic Load Balancing Approach in Cloud Computing," 2020 IEEE 1st International Conference for Convergence in Engineering (ICCE), 2020, pp. 195-199, doi: 10.1109/ICCE50343.2020.9290602.
24. M. Jeyakarthic and N. Subalakshmi, "Client Side-Server Side Load Balancing with Grasshopper optimization Mapreduce Enhancing Accuracy in Cloud Environment," 2020 Fourth International Conference on Inventive Systems and Control (ICISC), 2020, pp. 391-395, doi: 10.1109/ICISC47916.2020.9171130.
25. R. Agarwal, N. Baghel and M. A. Khan, "Load Balancing in Cloud Computing using Mutation Based Particle Swarm Optimization," 2020 International Conference on Contemporary Computing and Applications (IC3A), 2020, pp. 191-195, doi: 10.1109/IC3A48958.2020.233295.
26. F. Jamal and T. Siddiqui, "Comparative Analysis of Load Balancing Techniques in Cloud Computing, Based on LB Metrics," 2021 5th International Conference on Information Systems and Computer Networks (ISCON), 2021, pp. 1-5, doi: 10.1109/ISCON52037.2021.9702508.
27. S. K. Ojha, H. Rai and A. Nazarov, "Optimal Load Balancing In Three Level Cloud Computing Using Osmotic Hybrid And Firefly Algorithm," 2020 International Conference Engineering and Telecommunication (En&T), 2020, pp. 1-5, doi: 10.1109/EnT50437.2020.9431250.
28. C. S. M. Babou et al., "Hierarchical Load Balancing and Clustering Technique for Home Edge Computing," in IEEE Access, vol. 8, pp. 127593-127607, 2020, doi: 10.1109/ACCESS.2020.3007944.
29. A. I. El Karadawy, A. A. Mawgoud and H. M. Rady, "An Empirical Analysis on Load Balancing and Service Broker Techniques using Cloud Analyst Simulator," 2020 International Conference on Innovative Trends in Communication and Computer Engineering (ITCE), 2020, pp. 27-32, doi: 10.1109/ITCE48509.2020.9047753.
30. W. -Z. Zhang et al., "Secure and Optimized Load Balancing for Multitier IoT and Edge-Cloud Computing Systems," in IEEE Internet of Things Journal, vol. 8, no. 10, pp. 8119-8132, 15 May15, 2021, doi: 10.1109/JIOT.2020.3042433.
31. A. Hamidi, M. K. Goal and R. Astya, "Load Balancing in Cloud Computing Using Meta-Heuristic Algorithm: A Review," 2022 9th International Conference on Computing for Sustainable Global Development (INDIACom), 2022, pp. 639-643, doi: 10.23919/INDIACom54597.2022.9763131.
32. M. Kushwaha, B. L. Raina and S. Narayan Singh, "ImplementationAnalysis of Load Balancing Procedures for Cloud Computing Domain," 2021 International Conference on Computing,

- Communication, and Intelligent Systems (ICCCIS), 2021, pp. 287-292, doi: 10.1109/ICCCIS51004.2021.9397069.
33. M. A. Shahid, N. Islam, M. M. Alam, M. M. Su'ud and S. Musa, "A Comprehensive Study of Load Balancing Approaches in the Cloud Computing Environment and a Novel Fault Tolerance Approach," in IEEE Access, vol. 8, pp. 130500-130526, 2020, doi: 10.1109/ACCESS.2020.3009184.
  34. K. Pradeep and D. Pravakar, "Exploration on Task Scheduling using Optimization Algorithm in Cloud computing," 2022 6th International Conference on Trends in Electronics and Informatics (ICOEI), 2022, pp. 874-877, doi: 10.1109/ICOEI53556.2022.9777120.
  35. R. K. Ramesh, H. Wang, H. Shen and Z. Fan, "Machine Learning for Load Balancing in Cloud Datacenters," 2021 IEEE/ACM 21st International Symposium on Cluster, Cloud and Internet Computing (CCGrid), 2021, pp. 186-195, doi: 10.1109/CCGrid51090.2021.00028.
  36. N. Kumar and N. Mishra, "Load balancing techniques: Need, objectives and major challenges in cloud Computing- a systematic review," Int. J. Comput. Appl., vol. 131, no. 18, pp. 11-19, Dec. 2015
  37. R. Kaur and P. Luthra (2014). Load Balancing in Cloud Computing. Recent Trends in Information, Telecommunication and Computing, ITC, Association of Computer Electronics and Electrical Engineers.
  38. Klaithem Al Nuaimi, N. M.-J. (2012). A Survey of Load Balancing in Cloud Computing: Challenges and Algorithms.
  39. Younis, H. J. (2015). Efficient Load Balancing Algorithm in Cloud Computing.
  40. Hashem, W., Nashaat, H., & Rizk, R. (2017). Honey Bee Based Load Balancing in Cloud Computing. KSII Transactions on Internet and Information Systems (TIIS), 11(12), 5694-5711.
  41. Gudise, V. G., and Venayagamoorthy, G. K. (2003), Comparison of particle swarm optimization and backpropagation as training algorithms for neural networks, in Swarm Intelligence Symposium, Proceedings of the IEEE, pp. 110-117.