



Swarm Intelligence Optimization for Distributed Cloud Workloads

Aditi Namdeo

AI Researcher, Amazon, Seattle, USA

ABSTRACT: Today with Cloud Computing, however, workload management is more complicated as they are all heterogeneous (data centre, edge nodes and virtual resources). To minimize the execution time, energy consumption, service delay and resource imbalance, an efficient workload allocation is needed. Later, a research paper titled “Swarm Intelligence Optimization for Distributed Cloud Workloads.” claims that the framework based on the swarm intelligence for workload scheduling in distributed cloud system is missing. There are three components to it: Intelligent Task Profiling, Resource Monitoring and the Swarm based decision making and adaptive workload migration. The ants behavior, group decision making of birds and bee group to achieve near optimal resource task mapping, inspired optimization algorithm like Particle Swarm Optimization algorithm, Ant Colony Optimization algorithm and Artificial Bee Colony algorithm are the motivations of the proposed model. The features of the workloads (e.g., size of the tasks, priority, execution time and resource requirement) are the first ones collected through the framework. It then scans the cloud and finds out the available resources in the cloud based on the processing, memory, bandwidth, energy and load. To make the allocation decision, a swarm intelligence optimizer and a feedback module are used to ensure that the different scheduling policies are used with the various workload conditions. Analysis of the research results shows that it is possible to increase the scalability, fault tolerance, balancing and quality of service of the distributed system of clouds with the help of swarm optimization. Considering the highly unpredictable and dynamic nature of the cloud environment, the proposed framework is more flexible as compared to the static and heuristic scheduling.

KEYWORDS: Swarm Intelligence, Cloud Computing, Distributed Workloads, Workload Scheduling, Load Balancing, Resource Optimization, Particle Swarm Optimization.

I. INTRODUCTION

One of the most crucial elements of technologies behind the delivery of digital services is Cloud Computing specifically with regard to distributed and scalable computing resources, storage, applications and platforms. Cloud systems are deployed with many different industries, such as eCommerce, healthcare, education, banking, AI and IoT where flexibility in allocation of resources, fast response to requests and high availability are essential. More and more workloads are being distributed and the management of distributed cloud workloads has to be efficient—a vast research problem. The use of cloud workloads is typically unpredictable, varied and dynamic. These may be latency sensitive applications, machine learning applications, data analytics workloads, to batch jobs and real-time user requests. The workloads must also be assigned to the correct compute resource that has the necessary capabilities to meet the requirements of performance, reliability, energy-efficiency and cost [1].

Each of the distributed cloud environments can include one or more computing resources, such as one or more data centres, virtual machines, containers, edge nodes and storage systems and other resources. Unlike traditional centrally distributed cloud computing systems, distributed cloud infrastructures generally consist of numerous different resources with geographical distribution. This does contribute to scalability and fault tolerance, but adds to the complexities of workload scheduling and resource management. Some servers would not have been able to handle the distribution if some were overloaded, while other servers would not have been utilized. This imbalance can lead to service reaction time, more energy usage, SLA violations and ultimately, a poor user experience. To ensure that the distributed cloud system is efficient and sustainable, hence, optimisation of the workload is crucial [2].

Workload scheduling for cloud computing is about scheduling workloads to resources based on some objectives that are subject to other constraints. They may be for performance, efficiency, throughput or to balance server loads, quality of service, lower operational cost, or anything else. The clouds, however, are dynamic environments, not easy to accomplish all of the above at the same time all in the same cloud. Resources are continuously coming on and off-line as a result of failures, maintenance, migrations or simply because of fluctuation in demand. Workloads can also occur

randomly and have different sizes, priorities, deadlines and computation requirements. In many cases, however, the static scheduling techniques don't deal with the dynamicity of the cloud environment.

The various classical workload allocation techniques like FCFS (First Come First Served), Round Robin, Min-Min, Max-Min and other rules based allocation are adapted in cloud systems. Easy and simple to implement, but generally don't help to achieve best solutions in large scale and dynamic environments. They may ignore factors which influence performance, like the different conditions of the resources, delays in the network, dependency of tasks, energy use and priority of the workload. On top of this, many of the traditional approaches use hardcoded scheduling rules, and are inflexible to deal with dynamic changes to the cloud environment. This compels the researchers to embark on a study of how workloads can be optimally managed by smart routing when performing in cloud.

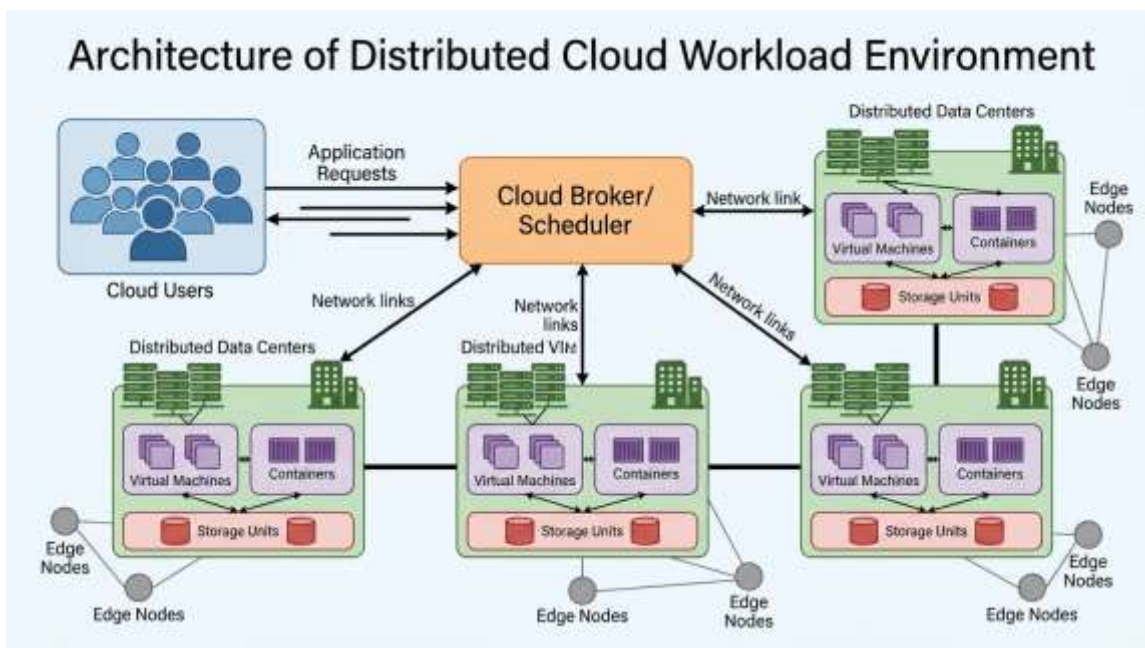


Figure 1: Architecture of Distributed Cloud Workload Environment

Inspired from the collective behaviour of the natural systems, Swarm Intelligence is a great optimisation technique. Ants, birds, bees, fish and termites are examples of simple agents, which in nature, solve complex problems by cooperating, self-organizing and decentralizing the decision making process. This is because ants can take the shortest route from their nest to the food, birds flocking can co-ordinate their movements and it is efficient for bees to seek out the food source. They are all natural phenomenon, which can be applied in a plethora of computational algorithms that can solve complex optimisation problems without the need of a central control system. Swarm Intelligence algorithms has been widely utilized for PSO, ACO, ABCO, FA, GWO and BA [5].

Swarm behaviour is not only shared between swarm and distributed cloud workload, but also their decisions are made in a distributed manner; this makes the concept of 'swarm-optimized' distributed cloud workload even more relevant. In a world of cloud computing there are plenty of scenarios and contingencies, and workloads must be broken down into a set of resources. Algorithms that use swarm intelligence are particularly suitable for such a task since they can explore large spaces of possible solutions, adapt to the changing environment and improve the solution step by step until near the optimum one is found. These models, known as swarm based models, perform several experiments with different task-resource mapping and then based on the frameworks results of the performance evaluation, take turn repeating the experiments again, in order to get the best mapping arrangement of the tasks with resources. These are also useful in the case of multiple objectives to be considered simultaneously develop the solution of multi-objective scheduling problems such as: execution time, energy consumption, cost and load balance [6].

This paper presents the framework for an effective workload allocation in distributed cloud environment and designing and analysis of the same using swarm intelligence. It is aimed to work for intelligent task scheduling, where workload can be profiled, it is also aimed to monitor the resource, resource can also be used for swarm based optimization;



Adaptive feedback control, etc. The first step is to give workload breakdown of the workloads entering into the system in terms of priority and estimated execution time (or deadline constraints) and also computation and task size. During the second phase, the cloud resources considered for performance monitoring include the processing power, memory of the cloud resources, as well as the bandwidth, workload, energy consumption and response time of cloud resources. Finally, the swarm intelligence algorithm is applied to make an optimized scheduling decision, based on the most suitable resources, according to the workload. Finally a feedback process compares the actual results with the amount allocated in previous process and alters the scheduling process to improve results of amount allocation in next process [7].

The primary goal of this scheme is to perform management of the cloud workload that will ensure balancing of the resources, minimize delay in cloud execution, optimize resource utilization with a significant boost in the overall quality of service. The framework can be especially beneficial in distributed cloud deployments, where workloads are constantly in flux and resources are distributed across various locations. The system can also be decentralized and adaptive, continuously optimizing its performance, based on the instantaneous change of the actual load, and applying for this swarm intelligence. This system can also scale up since the swarm based approach can be capable of performing a large number of tasks and resources without relying completely on the central scheduler [8].

The study has a great significance for the smart computing and sustainable cloud computing. As more and more services are deployed to the cloud, energy use/cost begins to be a concern for service providers. Otherwise this might result in increased use of energy and reduced productivity of the system. These issues could be eradicated by optimizing these resources with Swarm Intelligence approach and minimizing this unnecessary moving workload to other resources. In addition, improved scheduling can improve the performance of the application, minimize service failures and keep users happy. So, future work can be done by optimizing the cloud by distributing the workload of cloud [9] while using Swarm Intelligence.

The concept of this paper is the intelligent workloads scheduling and is one of the most crucial concepts in the upcoming generation of cloud computing. It summarizes and introduces the popular method of the workload distribution in the distributed system; It points out the limitations of traditional scheduling methods, and proposes a flexible scheduling method based on swarm intelligence. From the study results, it is concluded that the cloud workload management can be done in a efficient, scalable and autonomous way using SWO. It is combined with Natural Intelligence attribute with Algorithm of the Computational Resource Managing to enhance the efficiency of the Distributed cloud system, reliable and sustainable operation of cloud system in presented framework.

II. CURRENT OBSTACLES

The use of the swarm intelligence is highly promising to optimise the work load in the cloud, but still has some technical and operation points to consider before using swarm intelligence. Optimizing the placement of workloads is a tough problem in distributed cloud setups, given the size, dynamicity and diversity of distributed cloud. Details of some of the major problems related to swarm intelligence approach to efficient and reliable Workload Scheduling is discussed below.

2.1 Dynamic and Unpredictable Workload Patterns

The distributed cloud workload management critical factors are: workload arrival and resource demands. There are different types of workloads generated by Cloud Apps – real-time requests, data intensive workloads, batch and latency sensitive workloads. Unexpected surges may be caused by network congestion and application problems, client traffic and business requirements in these workloads. The swarm intelligence algorithms usually provide nearly optimal solutions, several iterations are required. However if the workload conditions change rapidly then this optimisation process may not be as effective as the scheduling decision may be stale before it can be fully exploited.

2.2 Resource Heterogeneity and Complexity

Distributed cloud infrastructures are the combination of a diverse set of resources—including virtual machines, containers, edge servers, storage and networking. They have various strengths, memory sizes, bandwidth and energy consumption and availability. What makes this task-to-resource mapping hard is the kind and extent of the aforementioned heterogeneity. There may be some resources that are beneficial for one resource but do not work for another. Therefore, the swarm intelligence models have to take several parameters of resources to be considered. If these differences are not accurately identified by the algorithm there is one or more of the following consequences: resources are poorly used, performance decreases and misallocation occurs.



2.3 Scalability of Optimization Algorithms

When attempting to make a scheduling decision, with a large number of cloud users, tasks and computing nodes, the search space will be huge. This could be a serious problem for the time that might take to compute the solution space, in swarm intelligence based methods: Artificial Bee Colony Optimization, Ant Colony Optimization or Particle Swarm Optimization for exploring the complex solution space in distributed cloud systems. The optimization algorithm may be a burden and consume too much resources, or slow down the scheduling process. Thus reaching a high level of quality optimization is difficult.

2.4 Balancing Multiple Objectives

Typically, you'll find several different goals that can often conflict with each other with respect to cloud workload scheduling. For instance, if the execution time decreases, it can be achieved by employing high performance servers which can lead to higher energy consumption and operating costs. Likewise, cost reductions may impact timeliness and/or service performance. In the implementing frameworks based on swarm intelligence, some of these factors, such as execution time, energy consumption, load balancing, cost and reliability and service level agreement (SLA) requirements should be taken into account. Of the most important problems is thinking of an appropriate fitness function that will consider all these target objectives, but not be too complex.

2.5 Security, Fault Tolerance, and Reliability

The vulnerabilities of the distributed cloud are the exposure to security threats, the failure of nodes, delay in communication and data privacy concerns. The workloads may be more susceptible if the security controls are not strong, when migrating to multiple nodes. Furthermore, communication of information is necessary between the agents/decision units in a permanent manner in order to obtain the function of swarm intelligence algorithms. Optimization process could make bad decisions if the resource monitoring information is inaccurate, late or corrupted. Therefore fault-tolerant and secure workload optimization mechanism plays a significant part in a practical world of deployments to cloud.

2.6 Integration with Existing Cloud Platforms

One of the biggest challenges is integrability of the swarm intelligence model(s) and integration of the models in existing platforms for cloud management. Lots of cloud-based systems have their own built-in schedulers, monitoring tools and virtualization technologies. If presented in an introduction, for instance the swarm-base optimization framework, it is critical that the framework can be harmonized to the systems and don't affect the normal function of the systems. Prior to its widespread introduction careful and careful design, testing and standardization is required.

III. SWARM INTELLIGENCE OPTIMIZATION FRAMEWORK FOR DISTRIBUTED CLOUD WORKLOADS

In order to overcome the workload scheduling, workload allocation, load balancing and quality of service problems in a distributed cloud environment, the Swarm Intelligence Optimization for Distributed Cloud Workloads is proposed. Although urethral sounds and vaginal sounds continue to be employed for treatment of constriction of the urethra and vagina (stricture).The urethral and vaginal sounds are still in use today to dilate strictures (narrowings) of the urethra and vagina. The size, priority, deadline, computational requirement and memory requirement of these workloads varies, as does also the cost of the communication. So now, there is a need for flexible, adaptable for effective distribution of tasks among resources. The proposed framework is based on the principles of Swarm Intelligence (SI) which attempts to emulate the groups of natural swarms such as flocks of birds, ant colonies, bee swarm, etc., to arrive at one which is very close to the optimum solution to the scheduling problem.

The framework consists of various layers that are interdependent with each other such as workload submission and profiling, resource monitoring, swarm based optimization, workload allocation, adaptive migration, performance evaluation and feedback control. Each layer has a specific function and messages are transmitted between the layers to aid in intelligently deciding what to do. To attain all this, the overall goal of the framework is to cut down on the execution time, reduce energy usage, avoid resource overload, increment the throughput of the framework system and follow service level agreement (SLA) conditions.



Figure 2: Swarm Intelligence Optimization Framework

3.1 Workload Submission and Profiling Layer

This “workload submission and profiling” layer is the first layer of this framework. Here all user, application and service requests are registered. These can be jobs that are Web service calls, data analysis jobs, machine learning jobs, database query jobs, multimedia analysis jobs or real-time workload jobs from an IoT (Internet of Things) application. All the workloads are analysed in the framework prior to the scheduling process, as workloads do not all have the same characteristics.

By using workload profiler, important properties related to a task (e.g. task length, input size, output size, CPU requirement, memory requirement, bandwidth requirement, task priority, task deadline) are identified and dependency with other tasks are identified. A lower response time means that less time is given to carry out a task, such as a latency sensitive task, which should be completed as quickly as possible or it could mean that more time is given for a task, such as a task for batch processing, which could be done later, and thus have a longer response time, for example. For example, a computation-intensive task could need more processing power while a task to handle large amounts of data could need more band width and storage. That's where this framework can come in and determine these attributes to decide not only random mapping of workloads, but incorporating the actual need of the workload.

Another feature of this layer is that the set of tasks is indicated as high priority, medium priority, low priority, compute intensive, memory intensive, delay sensitive workloads etc. This classification is appropriate for futher to make the allocation decision of swarm optimiser appropriately. This layer ends up producing a structured workload profile which will be passed to the optimization layer.

3.2 Resource Monitoring Layer

The next is the monitoring layer of resources as they are continuously monitored for “health” with distributed cloud resources. Cloud environments are dynamic, with the ability of virtual machines, containers, servers and edge nodes to have varying loads over time. Resources becoming overloaded and other resources may be under-used. Hence, it is highly important that a real time and accurate information on resource is available to make an efficient workload scheduling.

The information for the availability of the CPU, memory used, disk free space, whether any network bandwidth is being used, whether the load is being lifted or not, consumption of energies, length of the queue, processing speed, failure status of each resource etc is collected in this layer. It also keeps track of distributed nodes' geographical position, since geographical position can impact the communication delay and response time. For instance, if the task

gets assigned to a data center that is far away, then no matter how many processing from the data center, the latency could be high.

It will generate a table of resource status based on the resource status layer, and the status of each cloud node and the level of performance of cloud node. This is a table which is continually updated and handed over to an Optimization engine called the and swarm Optimization engine. An important part of this monitoring process is that it should be accurate such that it leads to good scheduling decisions. Because the system times aren't current and/or the resources information isn't accurate, the system may inefficiently allocate resources.

3.3 Swarm Intelligence Optimization Layer

The proposed framework has two layers to make sure that the swarm intelligence is built: A layer: Swarm intelligence optimization layer. Using this layer, swarm based algorithms are able to take correct decisions regarding workloads to cloud resources. The advantage of using swarm intelligence is that it is able to identify a number of potential scheduling solutions and refine these solutions by cooperation and feedback.

In this case the proposed algorithms are applied in their original and hybridized fashion such as particle swarm optimization (PSO), ant colony optimization (ACO) and artificial bee colony optimization (ABC). PSO is a process that is based on particles that represent the possible solutions for workload allocation. The particles move in the search space, and update their position based on their personal best and swarm best. ACO builds these paths by using the pheromone value (quality of already constructed allocation) and a scheduling for a given time is executed for the paths. This search process of Artificial Bee Colony Optimization utilize pseudo edition on the reward for searching a better task – resource combination to be employed as employed bee, onlooker bee or scout bee.

Initial population of potential schedules must be devised first as the initial step towards optimisation. Both these will be Task to Cloud Resource mappings. A fitness function is used for each schedule. The fitness function can consist of several performance metrics like makespan, response time, resource utilization, energy consumption, cost, load balancing and satisfying SLA (service level agreement). The lower the execution time and energy consumed and the higher the load balance and resource utilization will be the better your solution will be compared to others, and therefore, the higher the fitness value.

After that the optimizer carries out iterative search, to get new candidate solutions. As the algorithm goes on, poor solutions will be replaced by better solutions, until some stopping criterion is reached. Any of stopping at a number of iterations, stopping after a maximum amount of time or stopping when an acceptable fitness is reached can be used. Optimizer will be able to restore optimal workload allocation to workload allocation layer.

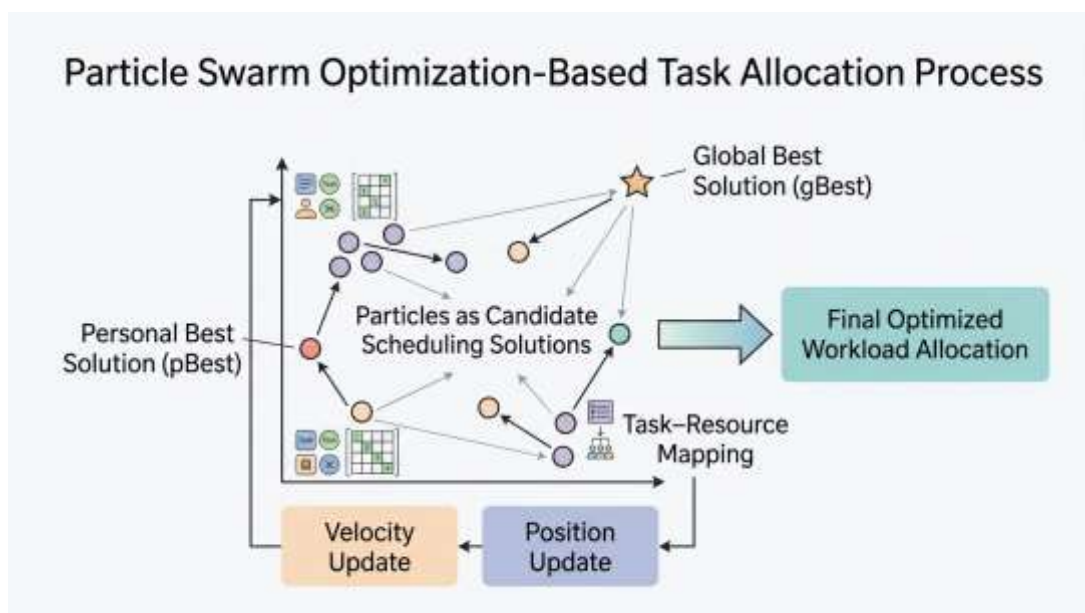


Figure 3: Particle Swarm Optimization-Based Task Allocation Process



3.4 Fitness Function Design

A major application of the framework used is the fitness function that will be used to steer the swarm optimizer in achieving good scheduling decisions. The one fitness factor can get better and one can get worse, depending upon the state of the fitness function. Therefore the fitness used in the proposed framework is a multi-objective fitness; it considers several indicators related to the performance of the clouds.

The fitness function they are trying to minimize is supposed to minimize the overall execution time (also called as response delay), maximize the utilization of the resources, minimize the energy consumed, balance the workload and of course achieve service level agreements. Fitness value is computed from some weighted values for each of the goals each fitness function is based on, depending on the system requirements. But if the application in question is going to be deployed to a real-time cloud-based application, response time could be more important. Consumption of energy plays a larger role in making a Green cloud Computing. In a cost sensitive environment, operational cost is one cost, which can be taken as a factor.

3.5 Workload Allocation Layer

The workload allocation layer will allocate workload to selected cloud resources, when the swarm optimizer finds out the best solution of the scheduling. It's an implementation of the execution controller of the framework. Based on the optimized tasks-resource mapping, it sends tasks to virtual machines, containers, physical machines or edge nodes.

In allocation, the following factors are considered: Tasks need and Resources availability. Numerous queues and significant amounts of processors could be provided to more important jobs. If nodes have a large storage and/or network bandwidth, then they can be used to solve data-intensive tasks. Workloads that are latency sensitive can be placed on local edge nodes, to minimize latency. Smart allocation can be used to avoid resource bottlenecks and it can give some of the resources a performance boost.

Also the workload allocation layer keeps a record of task execution. The following information is stored in this record: When the task was created, what user is responsible for the task, when is the task intended to be completed, when did the task end and a list of the resources that are used to complete the task and the task status. This records are not utilized until later within the performance evaluation and feedback levels.

3.6 Adaptive Workload Migration Layer

The distributed Cloud environment is a flexible application. Acceptable resource could be used more than necessary or could not be available later in time. Hence, it is proposed to inject a new layer to the proposed framework: Adaptive workload migration layer. This layer triggers recognition and invocation of a task reallocation to another resource if the service level of the resource is degraded, there are resources overload, the node fails or there is chance of violating the deadline.

The migration of workloads can help increase the level of fault-tolerance and load-balancing, but it could be a communication overhead and execution delay, if not done as needed. So migration is only considered relative to where the migration is likely to have greater benefits than costs. Migration is decided upon looking at the current load, progress being made on the tasks and available time to finish the tasks, taking network delays and resources available at the target in consideration.

An example of such scenario might be if a VM is resource constrained, the framework knows how to scramble a high priority task off to another node where they might have more resource, or the deadline may be in jeopardy. Similarly, if one of the edge nodes fails, then workload can be routed to another edge node in the cloud. It's a migration technique, that adapts/adjusts the system so it makes it more reliable, stable and resilient to disrupt the service.

3.7 Performance Evaluation Layer

A measure of the effectiveness of workload scheduling and after task executions is performance evaluation layer. It is used to measure how well he/she does, compared to the standard. The sets of performance metrics are – makespan, response time, throughput, resource utilization, energy, migration cost, load balancing and SLA violation rate.

The total time to finish all the tasks assigned - Maximum time to complete the tasks of this system - Makespan. Response time: It is the time measured as elapsed to respond from the system to the user's request. Throughput: Number of tasks that can be completed in a unit of time. Resource utilization – This is the evaluation of CPU, memory, storage and bandwidth utilization-levels. Energy Consumption is a good indicator of the energy spent by the 'at

runtime' cloud resources. With load balance, a visualisation of the load evenly spread over the resources is given. SLA violation rate (SLA violation rate): Percentage rate of SLAs violation of the tasks performed.

This layer will give you information on how well the framework that underlies the swarm is working and how well the 'cloud' is performing during your interactions with the swarm. It also facilitates to compare the result of the proposed algorithm with the traditional scheduling algorithm such as round robin, FCFS, min-min and max-min scheduling.

3.8 Feedback and Learning Layer

Feedback and learning on the last level within the framework. The layer is given the responsibility of properly altering the framework thereafter depending on the result of the performance in the next schedules. It is periodically assessed whether the workload is efficient/effective. Therefore, in the future, an allocation strategy which has lower execution time and uses less resources naturally chooses to be the "high priority" to undergo allocation. Otherwise, the consumption or delay of energy is less significant – and thus not the optimal approach – when the set is "overloaded."

Based on this feedback some parameters like the particle velocity, amount of pheromone, probability of the search, the weights of the fitness function, etc. are updated in the swarm optimizer. Thus, the framework gets to be more clever with time. Feedback: dynamically change depending on the type of workload. The system can focus on response-time and loadbalancing for example during busy traffic speeds. Furthermore, it can be a way to focus attention on 'energy conservation' as with less traffic fewer machines will be in use.

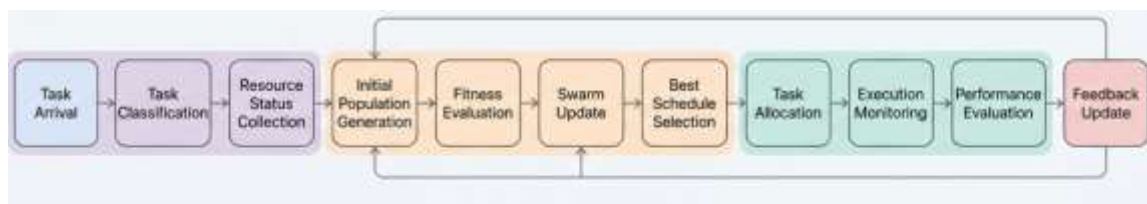


Figure 4: Workflow of Swarm-Based Workload Scheduling

IV. OPERATIONAL BENEFITS OF THE SWARM INTELLIGENCE OPTIMIZATION FRAMEWORK

The Swarm Intelligence Optimization Framework (SIOF) offers a variety of technical and operational benefits for distributed cloud workloads in today's cloud computing context. The stack of workload profiling and swarm optimization monitoring resources in real-time, dynamic migration and feedback control can be leveraged to make smart decisions on a distributed cloud system dealing with large, dynamic and heterogeneous workloads. Below is a summary of the benefits' highlights.

4.1 Improved Load Balancing

Proposed framework's one of the major advantages is the improvement of the load balancing of distributed cloud resources. Consuming more resources than required in some of the servers, idle in traditional scheduling. This imbalance can result in the increased execution delay and less overall system efficiency. Before initiating the assignment, the swarm intelligence optimizer evaluates constantly the resources and the amount of work required. This also provides an additional way to distribute a workload across VMs, containers, edge nodes and data centers. However, such a distribution decreases the amount of congestion, wait times during operations, bottlenecks and enhances cloud performance.

4.2 Reduced Execution Time and Response Delay

It is assumed that the tasks are duly allocated to the resources so that the time for execution of the tasks along with the delay time in receiving the actual response in the proposed framework. A workload profiling layer can assist the optimizer in matching a task to a cloud resource, as it includes information about the task, such as the amount of CPU, memory, bandwidth, priority and deadline it requires. They will be offloaded to the nearby nodes for latency sensitive tasks and computation intensive tasks are offloaded to high-capacity servers. This is a dynamic distribution that can help reduce the response time and hence the better QoS to end-user.

4.3 Enhanced Resource Utilization

To minimize the cost of operation and maximize the productivity of the system, it becomes extremely important to have an effective utilization of the cloud resources. The proposed framework has no stress or waste on the resources and is



able to do real time monitoring of the condition of the resources. The optimizer will figure out what tasks to execute and on which resources to best utilize the CPU, memory, storage and network bandwidth resources available. Optimized resource usage also leads to decrease in idle-time resources and opportunity of maximizing ROI with Cloud resources.

4.4 Energy Efficiency and Cost Reduction

Among the problems cloud data centers of this magnitude must contend with is energy usage. It is not advisable to change the workload allocation because it can result in setting too many servers to "up" or using too many of the power consuming resources. One of the optimizations objectives – Energy consumption – is added in the proposed one as a virtual concept. The framework might be implemented to optimise the assignment of the tasks to minimise task unnecessary migration and consequently minimising the power consumption and running cost. This makes the platform suitable to sustainable and green Cloud Computing.

4.5 Scalability and Adaptability

Distributed cloud environments are constantly expanding, experiencing more uses, apps and compute nodes. Proposed Swarm intelligence framework is scalable and the algorithms utilized by the swarm can explore huge space of solutions, adapting to change of workload situations. Another advantage of feedback layer is adaptability – decisions made and determined by the past performance. It's suitable for a dynamic cloud environment, fluctuating over time with regards to workloads and resources.

4.6 Improved Fault Tolerance and Reliability

The adaptive migration layer consists of maintaining the state of resources, handling node failures or deadline risky, thus increasing the fault tolerance. If it is not efficient or cannot fail over, then the framework can fail over to other nodes that are ready. This minimises disruption to the service and is improved reliability of the system. Hence, it can be concluded that the proposed framework can be able to attain best performance and also further strengthening the stability in distributed workload management in cloud too.

V. CONCLUSION

One of the important study papers is 'Swarm Intelligence Optimization for Distributed Cloud Workloads' as the distributed cloud workload system will be need of being intelligent and adaptive to the workload. Cloud infrastructures are becoming increasingly common, and heterogeneous/dynamic nature of cloud made the existing scheduling techniques ineffective to provide efficient resource allocation, load balancing and service quality. To overcome with such drawbacks, the swarm intelligence based approach towards collective optimization which is observed in nature – ants, birds, bees etc is proposed.

It includes workload profiling, real-time monitoring of resources, optimize workloads based on swarms, migrate workloads adaptively, performance assessment and feedback control. This hierarchical system allows processing of the received tasks, check the resources and optimally scheduling the task which further helps in the workload assignment process. The framework supports the algorithm of swarm intelligence (SSA) such as Particle Swarm Optimization (PSO), Ant Colony Optimization (ACO) and Artificial Bee Colony Optimization (ABC) to get close to optimal mapping of a task to resources.

We have several benefits to the proposed approach towards better load balancing, execution time, less energy consumption and utilization of resources, better scalability and better fault tolerance. Based on this framework, multiple metrics that include makespan, response time, cost, energy consumption and SLA can be used to more equitably and sustainably manage cloud workloads. Various adaptive migration capabilities, adaptation to node failures, workload fluctuations and feeding back resource conditions ensure system reliability.

In conclusion, the swarm intelligence is a potential approach to future distributed cloud optimisation. It helps to enable flexibility for self-governing decision making, dynamic adaptation and efficient resource utilization in complex clouds. Future work can be done to extend this architecture to accommodate Machine Learning, Deep Reinforcement Learning, edge-cloud coordination and security aware scheduling models. Their achievements will have the potential to take next generation cloud computing systems a step further towards more smarter, stronger and sustainable systems.



REFERENCES

- [1] S. Nabi, M. Ahmad, M. Ibrahim, and H. Hamam, "AdPSO: Adaptive PSO-based task scheduling approach for cloud computing," *Sensors*, vol. 22, no. 3, Art. no. 920, 2022.
- [2] S. A. Alsaidy, A. D. Abbood, and M. A. Sahib, "Heuristic initialization of PSO task scheduling algorithm in cloud computing," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 6, pp. 2370–2382, 2022.
- [3] K. Dubey and S. C. Sharma, "A novel multi-objective CR-PSO task scheduling algorithm with deadline constraint in cloud computing," *Sustainable Computing: Informatics and Systems*, vol. 32, Art. no. 100605, 2021.
- [4] E. H. Houssein, A. G. Gad, Y. M. Wazery, *et al.*, "Task scheduling in cloud computing based on meta-heuristics: Review, taxonomy, open challenges, and future trends," *Swarm and Evolutionary Computation*, vol. 62, Art. no. 100841, 2021.
- [5] H. Singh, S. Tyagi, P. Kumar, *et al.*, "Metaheuristics for scheduling of heterogeneous tasks in cloud computing environments: Analysis, performance evaluation, and future directions," *Simulation Modelling Practice and Theory*, vol. 111, Art. no. 102353, 2021.
- [6] R. Gong, D. Li, L. Hong, and N. Xie, "Task scheduling in cloud computing environment based on enhanced marine predator algorithm," *Cluster Computing*, vol. 27, no. 1, pp. 1–15, 2024.
- [7] Z. Zhang, M. Zhao, H. Wang, Z. Cui, and W. Zhang, "An efficient interval many-objective evolutionary algorithm for cloud task scheduling problem under uncertainty," *Information Sciences*, vol. 583, pp. 56–72, 2022.
- [8] Q. Hu, X. Wu, and S. Dong, "A two-stage multi-objective task scheduling framework based on invasive tumor growth optimization algorithm for cloud computing," *Journal of Grid Computing*, vol. 21, no. 2, Art. no. 31, 2023.
- [9] B. Pourghebleh, A. A. Anvigh, A. R. Ramtin, and B. Mohammadi, "The importance of nature-inspired meta-heuristic algorithms for solving virtual machine consolidation problem in cloud environments," *Cluster Computing*, vol. 24, no. 3, pp. 1–24, 2021.
- [10] A. R. Arunarani, D. Manjula, and V. Sugumaran, "Task scheduling techniques in cloud computing: A literature survey," *Future Generation Computer Systems*, vol. 91, pp. 407–415, 2019.