

# ROLE OF META-HEURISTIC ALGORITHMS IN CLOUD COMPUTING

**Ritu Kapur**

*ME Student, Computer Science Department NITTTR, Chandigarh (India)*

## **ABSTRACT**

*Efficient Resource Scheduling in cloud has become a challenge for Cloud Service Providers. Resource Scheduling should be such that it minimizes the overall resource cost, minimizes the overall execution time, increases throughput, fails less frequently, balances the load on server and delivers the best possible QoS. The algorithms that help in achieving these aims while doing an efficient resource scheduling come under the category of meta-heuristic algorithms. Some of the examples are Genetic Algorithms, Ant Colony Optimization, and Particle Swarm Optimization etc.*

**Keywords:** *Cloud, Load Balancing, Meta-Heuristic Techniques, QOS, Self- Learning Algorithms*

## **I. INTRODUCTION**

Cloud Computing has emerged as a popular area of research and development in the past few years. Although no unique definition of Cloud Computing exists many standard organizations continue to give various definitions for it [1]. According to National Institute of Standards and Technology (NIST) definition in September 2011 in [2], Cloud computing is defined as “a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction”. Cloud computing as utility provides pay-per use and pay-as-you-go service [3] i.e. a client/user has to only pay for the amount of services used by him [4].

According to [5] and [6], International Standards Organization (ISO) on 24<sup>th</sup> October 2014 defined cloud computing as “an evolving paradigm” and identified and described its “key characteristics” including broad network access, measured service, multi-tenancy, on-demand self-service, rapid elasticity and scalability, and resource pooling. Cloud computing basically provides three service models: Software as a Service (SaaS), Platform as a Service (PaaS) and Infrastructure as a Service (IaaS) as shown in [2] and [7]. According to [8], [9] and [3], IaaS provides the access to various resources which provide processing, storage, networks, operating systems and various other capabilities of scaling up and down the infrastructure dynamically. The services or resources required by a particular user are listed in the Service Level Agreement (SLA) established as a contract between the cloud service provider and the cloud consumer respectively [10].

According to D. Linthium in [11], Amazon is providing lowest priced instances and thus most of the researchers in [14-17] have deployed Amazon Cloud in [12] and Azure Cloud in [13] for their simulations. CloudSim, being an open source tool can also be used for simulations as used in [18]. The detailed implementation and basic design of CloudSim is explained in [19].

The authors in [20-39] discuss the concept of various meta-heuristic algorithms like Genetic Algorithms (GAs), Ant Colony Optimization (ACO) algorithms, Particle Swarm Optimization (PSO) algorithms etc. Meta-Heuristic Algorithms [39] are inspired from the natural phenomena and may provide a sufficiently good solution to an optimization problem, especially with incomplete or imperfect information or limited computation capacity. The authors in [20-23, 37, 38] discuss various genetic algorithms. GA provide a solution to complex combinational optimization problems [26], which come under the category of NP-Hard problems and thus their solution is hard to find.

The authors in [24-32, 38] discuss various Ant Colony Optimization (ACO) algorithms. ACO uses a colony of artificial ants which behave as cooperative agents in a mathematical space where they are allowed to search and reinforce pathways (solutions) in order to find the optimal solution [32]. K. Li et al. [24] presents a Load Balanced version of basic ACO proposed by M. Dorigo et al. [28-29]. The authors in [23] also give a load balancing factor for the same. The authors in [31] describe an optimized version of basic PSO algorithm shown by [32]. PSO is similar to GA [32], but there is no direct operation on the individuals of the population. Instead, PSO relies on the social behavior of the particles. In every generation, each particle adjusts its path based on its best position (local best) and the best position among its neighbors.

Meta-Heuristic Algorithms [32] inspired by evolution programs, assume that the solutions to be passed on to the next generations are selected by a particular mechanism. Thus the meta-heuristic algorithms are also called evolutionary or self-learning.

## II. ARCHITECTURE

### 2.1 Cloud Computing Architecture [40]

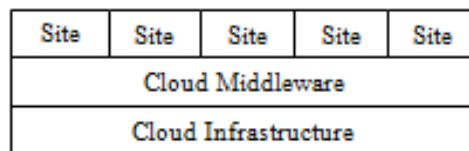


Fig. 1 Cloud Computing Architecture [40]

Cloud Service Providers or Site owners provide various services to the Cloud Consumers through the Middleware Layer from the actually executing Infrastructure Layer below which has the actual hardware resources like processors etc.

### 2.2 Algorithms

#### 2.2.1 Genetic Algorithm [20-23, 37]

Jingyi Ma [20] proposes a genetic algorithm named as a heuristic genetic load balancing algorithm (HGLBA) which helps in reducing the overall response time of the tasks through its load balancing phenomenon. The algorithm executes in the following manner:-

**A) Generate:** The first step is to generate the initial population,  $P(0)$ . This initial population is basically some initial solutions coded in the form of chromosomes. The maximum iteration named as *MaxSize* is initialized to zero i.e.  $k=0$ ;

**B) Evaluate:** In the evaluation step, the Fitness value of the Initial Population,  $P(k)$ , is evaluated on the basis of some function. e.g. the evaluation function in [44] is shown below in equation 1.

$$T_j = \sum_{i=1}^m (w_{ij} * Q_{ij}) \text{ for } \varphi=1, 2, \dots, n \quad (1)$$

where  $w_{ij}$  represents the weight listed to various QoS parameters (according to the user preferences),  $Q_{ij}$  represents the expected values of QoS parameters,  $T_j$  is the fitness function of the task  $j$ ,  $m$  represents the total no. of parameters and  $n$  represents the total no. of resource types.

**C) Improve:** In this improve step, various operations like mutation and crossover etc. are applied in order to generate the next generation,  $P(k+1)$ .

**D) Termination:** If the iteration  $k > \text{MaxSize}$ , we take the best solution and terminate the process otherwise we increment  $k$ , i.e.  $k=k+1$ , and continue the process from step B.

As this algorithm provides optimal allocation of resources, it serves the aim of the load balancing i.e. the minimization of overall execution time. The authors in [53] present a variation of genetic algorithm in an unreliable cloud scenario.

### 2.2.2 Ant Colony Optimization (ACO) Algorithm [24-31]

M. Dorigo et al. [28-29] present the basic ACO algorithm. K. Li et al. [24] embeds the load balancing parameter into it and thus presents a Load Balanced ACO i.e. LBACO algorithm. According to the authors in [23-28] ants have a general property of locating their path to food. In this phenomenon of locating their food, ants keep on depositing a chemical substance known as pheromone. This pheromone in turn helps new ants on their path to find the shortest route to the food. The path having largest amount of pheromone represents the shortest route. So this process thus leads the new ants to their location in minimum time through the shortest path.

The LBACO follows the following steps:

#### A) Initialization

The ants i.e. resources are randomly distributed onto various VMs initially. The pheromone is initialized as follows:

$$\tau_j(0) = pe\_num_j * pe\_mips_j + vm\_bw_j \quad (2)$$

where  $pe\_num$ ,  $pe\_mips_j$  and  $vm\_bw_j$  represent the no. of available VM processors, MIPS (Million Instructions Per Second) and the bandwidth capacity of the processor  $VM_j$  respectively.

#### B) Choosing the next task

Every  $k^{\text{th}}$  ant chooses the next task with the probability  $p_j^k(t)$  defined as:

$$p_j^k(t) = \frac{(\tau_j)^{\alpha} * (EV_j)^{\beta} * (LB_j)^{\gamma}}{\sum (\tau_j)^{\alpha} * (EV_j)^{\beta} * (LB_j)^{\gamma}} \text{ if } j \in 1 \dots n \quad (3)$$

else  $p_j^k(t) = 0$

where  $\tau_j(t)$  represents the pheromone value at time  $t$ ,  $EV_j$  represents the computing capacity of  $j^{\text{th}}$  VM defined as follows:

$$EV_j = pe\_num_j * pe\_mips_j + vm\_bw_j \quad (4)$$

$LB_j$  represents the Load balancing factor for the  $j^{\text{th}}$  VM represented as follows:

$$LB_j = 1 - \frac{res_j - lastAver\_res}{res_j + lastAver\_res} \quad (5)$$

where  $lastAver\_res$  represents the average execution time of all the VMs and  $res_j$  represents the current expected execution time, defined as follows:

$$res_j = \frac{total\_tasklength}{EV_j} + \frac{InputFilesize}{vm\_bw_j} \quad (6)$$

where total\_tasklength is the total length of all the tasks submitted to the VM<sub>j</sub> and InputFilesize is the length of the task before execution and  $\alpha, \beta, \gamma$  represent the parameters to control the relative weight of pheromone laid, the computing capacity and the load balancing factor of various VMS.

### C. Pheromone Updating

$$\tau_j(t+1) = (1-\rho) * \tau_j(t) + \Delta\tau_j \quad (7)$$

where  $\rho$  represents the pheromone decay coefficient,

$$\Delta\tau_j = \frac{1}{T_{ik}} \quad (8)$$

where  $T_{ik}$  represents the shortest path searched by  $k^{\text{th}}$  ant in the  $i^{\text{th}}$  iteration and once an optimal solution is found, the global pheromone updating is done as follows:

$$\Delta\tau_j = \frac{D}{T_{op}} \quad (9)$$

where  $D$  represents the encouragement coefficient and  $T_{op}$  represents the found optimal solution.

### D. Terminating Step

If all the ants end their search i.e. the solution converges to an optimal value, current optimal solution is saved and the value of makespan variable,  $N_c$ , is incremented. If this makespan value proves to be the best value, the process ends else the ants continue to search for a better optimal solution.

Authors in [30-31] give a variation of ACO algorithm in which the workload is distributed uniformly among all the VMs.

## 2.2.3 Particle Swarm Optimization (PSO) Algorithm [32-36]

Particle Swarm Optimization (PSO) [35] is a population based stochastic optimization technique developed by Dr. Eberhart and Dr. Kennedy in 1995 inspired by social behavior of bird flocking or fish schooling. PSO is similar to GA in the sense that the system is initialized with a population of random solutions and searches for optimal solution by updating generations. However, as compared to GA, PSO does not perform any crossover and mutation operations. In PSO, the potential solutions, called particles fly through the problem space by following the current optimum particles. Each particle keeps track of its coordinates in the problem space, which are associated with the best solution (fitness), it has achieved so far. The fitness value stored is called pbest. Another best value found by the PSO optimizer by examining the particle and its neighbors is lbest. The global best value found by examining the global population is gbest.

The PSO works by varying the velocity per unit time, i.e. acceleration, of each particle towards its pbest and lbest locations. The change known as acceleration is brought about by generating random numbers by the random number generation method. The two basic equations used in PSO for calculation purpose are:

$$V_{i+1} = w * V_i + c_1 r_1 (X_{lb} - X_i) + c_2 r_2 (X_{gb} - X_i) \quad (10)$$

$$X_{i+1} = X_i + V_{i+1} \quad (11)$$

where  $V_{i+1}$  is new calculated velocity,  $w$  is inertia weight,  $V_i$  is current velocity of particle,  $c_1$  is first acceleration coefficient,  $r_1$  is first random variable,  $X_{lb}$  is lbest,  $X_i$  is current position,  $c_2$  is second acceleration coefficient,  $r_2$  is second random variable,  $X_{gb}$  is gbest,  $X_{i+1}$  is next calculated position value of variable.

The algorithm proposed by [31] has the following steps:

- A) **Initialize:** The population of particles, i.e. tasks, is initialized by initializing candidate solution, Inertia weight, pbest, lbest, gbest and other PSO variables.
- B) **Generate:** Particle mapping are generated with machines using smallest position value heuristic.

- C) Evaluate:** Next, the fitness parameter is evaluated by using the make span parameter. If the current candidate solution has best fitness value, the new mapping is updated with its value; else the previous candidate solution is retained.
- D) Update:** On the basis of the previous evaluate step, pbest and gbest value is evaluated
- E) Calculate:** New velocity value and new particle position is calculated based on the equation 10 and 11.
- F) Termination:** If the maximum no. of iterations are reached, candidate mapping's expected time to complete matrix is calculated and the control is passed onto the load balancing component else the process iterates to step B.
- G) Load Balancing:** Next we check if load balancing is required or not. If it is required, an appropriate max-max or min-min algorithm concept is chosen for shuffling the task mappings. Tasks are shuffled and reordered and mapping is reevaluated and the load balancing is rechecked. If the new mapping is found to be load balanced the process is terminated and the final task machine is ready for execution. If the load is not found to be balanced, the tasks are reshuffled and the process continues till a perfectly load balanced mapping is found.

The simulations by the authors in [31] prove that the proposed algorithm outperforms the basic algorithm.

### III. CONCLUSION

The paper starts by introducing the concept of cloud computing and the importance of efficient resource scheduling in it. Then the importance of Meta-Heuristic algorithms like GA, ACO, LBACO, PSO etc. is explained. Meta-Heuristic algorithms help in balancing the workload. Thus a more balanced workload helps in reducing the load on server and thus decreasing the possibility of server failure and delivering a better QoS.

### REFERENCES

- [1] D. Linthicum (2014), "It's pointless to define cloud computing", <http://www.infoworld.com>
- [2] "Cloud Computing Bible", by Barrie Sosinsky, Chapter 2
- [3] A. Aggarwal and S. Jain, "Efficient Optimal Algorithm and Task Scheduling in Cloud Computing Environment", International Journal of Computer Trends and Technology (IJCTT), vol. 9, pp. 344-349, Mar 2014
- [4] D. Jensen, C. D. Locke, and H. Toluca (1985), A time-driven scheduling model for real-time systems, In IEEE Real-Time Systems Symposium.
- [5] J. Borne (20<sup>th</sup> October 2014), "ISO publishes new cloud computing standards and definitions", [www.cloudcomputing-news.net](http://www.cloudcomputing-news.net)
- [6] T. Roughton (17<sup>th</sup> October 2014), "Cloud Computing terms defined in new ISO standard", [www.out-law.com](http://www.out-law.com)
- [7] "Cloud Computing, A Practical Approach" by Anthony T. Velte, Toby J. Velte and Robert Elsenpeter, 2010 edition, Chapter 2.
- [8] "Real time task scheduling", last chapter, Version 2, IIT Kharagpur.
- [9] Casati and M. Shan (2001)., Definition, execution, analysis and optimization of composite e-service, IEEE Data Engineering.

- [10] R Buyya, C. S. Yeo and Venugopal, "Market oriented cloud computing: Vision, hype, and reality for delivering IT services as computing utilities" in *Proc. of the 10th IEEE International Conference on the High Performance Computing and Communications*, 2008.
- [11] D. Linthicum, "Faster and Cheaper: The cloud is becoming harder to resist", <http://www.infoworld.com>
- [12] Amazon elastic compute cloud [Online]. Available: <http://aws.amazon.com/ec2>
- [13] Microsoft windows azure. [Online]. Available: <http://code.google.com/windowsazure/>
- [14] X. Nan, Y. He, and L. Guan, "Optimization of Workload Scheduling for Multimedia Cloud Computing," in *Proc. IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1–4, 2013.
- [15] X. Nan, Y. He, and L. Guan, "Towards Optimal Resource Allocation For Differentiated Multimedia Services in Cloud Computing Environment," in *Proc. IEEE International Conference on Acoustic, Speech and Signal Processing (ICASSP)*, pp. 1–5, 2014.
- [16] X. Nan, Y. He, and L. Guan, "Optimal resource allocation for multimedia cloud based on queuing model," in *Proc. IEEE International Workshop on Multimedia Signal Processing (MMSP)*, pp. 1–6, 2011.
- [17] X. Nan, Y. He, and L. Guan, "Optimal resource allocation for multimedia cloud in Priority Service Scheme," in *Proc. IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1–4, 2012.
- [18] Aggarwal and S. Jain, "Efficient Optimal Algorithm and Task Scheduling in Cloud Computing Environment", *International Journal of Computer Trends and Technology (IJCTT)*, vol. 9, pp. 344-349, Mar 2014.
- [19] R. Raman, R. Calheiros, R.N.(2009), "Modeling and Simulation of Scalable Cloud Environment and the CloudSim Toolkit: Challenges and Opportunities", IEEE publication , pp1-11, 2009.
- [20] J. Ma, "A Novel Heuristic Genetic Load Balancing Algorithm in Grid Computing", in *Proc. of IEEE International Conference on Intelligent Human-Machine Systems and Cybernetics*, pp. 166-169, 2010.
- [21] J. Alcaraz and C. Maroto, "A robust genetic algorithm for resource allocation in project scheduling", *Annals of Operations Research* 102 (2001) pp. 83–109.
- [22] S. Hartmann, "A competitive genetic algorithm for the resource-constrained project scheduling", *Naval Research Logistics* 456 (1998) pp. 733–750.
- [23] U. Kohlmorgen, H. Schmeck and K. Haase, "Experiences with fine-grained parallel genetic algorithms", *Annals of Operations Research* 90 (1999) pp. 203–219.
- [24] K. Li, G. Xu, G. Zhao, Y. Dong and D. Wang, "Cloud Task scheduling on Load Balancing Ant Colony Optimization", in the *Proc. of IEEE Sixth Annual ChinaGrid Conference*, 2011, pp. 3-9
- [25] Merkle, M. Middendorf and H. Schmeck, "Ant Colony Optimization for Resource-Constrained Project Scheduling", in the *Proc. of IEEE transactions on Evolutionary Computation*, Vol. 6, NO. 4, August 2002, pp. 333-346
- [26] B. Tierney, W. Johnston, J. Lee and M. Thompson, "A data intensive distributed computing architecture for grid applications", *Future Generation Computing System* 2000, 16(5): pp. 473-481, 2000.
- [27] G. Guo-Ning and H. Ting-Lei, "Genetic Simulated Annealing Algorithm for Task Scheduling based on Cloud Computing Environment", in *Proc. of International Conference on Intelligent Computing and Integrated Systems*, pp. 60-63, 2010.
- [28] M. Dorigo and C. Blum, "Ant colony optimization theory: A survey in Theoretical Computer Science", 344 (2–3) (2005), DOI: 10.1016/j.tcs.2005.05.020, pp.243–278, 2005.

- [29] M. Dorigo and L.M. Gambardella, "Ant colony system: A cooperative learning approach to the traveling salesman problem", in *IEEE Transactions on Evolutionary Computation* (1997), pp. 53–66, 1997.
- [30] Bagherzadeh, J., MadadyarAdeh, M., "An Improved Ant Algorithm for Grid Scheduling Problem" in 14th International CSI Computer Conference (CSICC), pp.323-328, 2009.
- [31] Lorpunmanee, S., Sap, M.N, Abdul Hanan Abdullah, A.H., "An Ant Colony Optimization for Dynamic Job Scheduling in Grid Environment" in *Proceedings of World Academy of Science, English and Technology Volume 23 august 2007*, ISSN 1307-6884, 2007.
- [32] Anupinder Singh and Mala Kalra, "PSO Based Load Balancing Algorithm for Cloud Environment", M.E. diss., NITTTR, Chandigarh, 2014.
- [33] M. Clerc, "The Swarm and the Queen: Towards a Deterministic and Adaptive Particle Swarm Optimization", in *IEEE Proc. of the 1999 Congress on Evolutionary Computation, Washington, DC, vol. 3*, pp. 19-24,1999.
- [34] Y. Shi and R. C. Eberhart, "Parameter selection in particle swarm optimization", in *7<sup>th</sup> International Conference on Evolutionary Programming, San Diego, California, USA*, pp. 591-600, Springer 1998.
- [35] R. C. Eberhart and Y. Shi, "Particle Swarm Optimization: Developments, Applications and Resources", in *IEEE Proc. of the 2001Congress on Evolutionary Computation, Seoul, vol.1*, pp-81-86, 2001.
- [36] J. Kenndy and R. Eberhart, "Particle Swarm Optimization", in *IEEE Proc. of International Conference on Neural Networks, vol. 4*, pp. 1942-1948,1999.
- [37] Lovejit Singh and Sarbjeet Singh, "A Genetic Algorithm for Scheduling Workflow Applications in Unreliable Cloud Environment", *Second International Conference on Security in Computer Networks and Distributed Systems (SNDS-2014)*, Springer Berlin Heidelberg, pp. 139-150, 2014.
- [38] Ritu Kapur and Maitreyee Dutta, "Review of various Load Balancing and Green Computing Techniques in Cloud", in the *Proc. of 7<sup>th</sup> International Conference on Advances in Computing Sciences, Software Solutions, E-Learning, Information and Communication Technology (ACSEICT-2015)*, *Journal of Basic and Applied Engineering Research (JBAER)*, Print ISSN: 2350-0077, Online ISSN: 2350-0255;Vol 2, No. 2; pp. 122-127, 2015.
- [39] C.W. Tsai and J.J. Rodrigues, "Metaheuristic Scheduling for Cloud: A survey," *IEEE Systems Journal*, vol. PP, no. pp, pp. 279-291, 16 May 2013.
- [40] Verma, G. Dasgupta, T. K. Nayak, P. De, and R. Kothari, "Server workload analysis for power minimization using consolidation," in *USENIX'09*. Berkeley, CA, USA: USENIX Association, pp. 28–28, 2009.