Modern Education
and Computer Science
PRESS

# A Novel Approach for Optimization Auto-Scaling in Cloud Computing Environment

**Khosro Mogouie**
Department of Computer Engineering, Mahallat Branch, Islamic Azad University, Mahallat, Iran
Email: khosro_mogouie@yahoo.com

**Mostafa Ghobaei Arani**
Department of Computer Engineering, Parand Branch, Islamic Azad University, Tehran, Iran
Email: mostafaghobaei@piau.ac.ir

**Mahboubeh Shamsi**
Department of Computer Engineering, Qom Branch, University Of Technology Qom, Iran
Email: shamsi@qut.ac.ir

*Abstract*—In recent years, applications of cloud services have been increasingly expanded. Cloud services, are distributed infrastructures which develop the communication and services. Auto scaling is one of the most important features of cloud services which dedicates and retakes the allocated dynamic resource in proportion to the volume of requests. Scaling tries to utilize maximum power of the available resources also to use idle resources, in order to maximize the efficiency or shut down unnecessary resources to reduce the cost of running requests. In this paper, we have suggested an approach based on learning automata auto- scaling, in order to manage and optimize factors like cost, rate of violations of user-level agreements (SLA Violation) as well as stability in the presence of traffic workload. Results of simulation show that proposed approach has been able to optimize cost and rate of SLA violation in order to manage their trade off. Also, it decreases number of operation needed for scaling to increase stability of system compared to the other approaches.

*Index Terms*—Cloud computing, Scalability, Auto-scaling, Learning automata.

## I. INTRODUCTION

Cloud computing is a new technology that has found its place in human life as a basic need, and its popularity increases among Internet users day by day. Cloud service users only pay the price based on the amount of resources they have used (Pay-per-use). On the other hand, high accuracy is required to establish trade-off between price and efficiency. Providing high-quality cloud services at the lowest possible price, meeting the requests and maintaining system stability are the ultimate desire. Scalability is one of the important challenges and characteristics of cloud computing technology that can provide this feature for cloud services users [1, 2, 3].

One of the management techniques for scaling in the cloud computing is using a threshold approach. Threshold is determined based on factors such as the strength and speed of processing, and storage capacity. When utilization of resource is greater than the value of threshold, requests will be referred to other sources and when use of resource is less than threshold, some of the unnecessary resources will be selected and disabled to decrease costs. Therefore, the amount of implementation activities on a resource has measured periodically then the scaling will done on its basis [4].

Scalability which means allocation or withdrawal resources in proportion of requests, tries to maximize efficiency of available resources and using inactive resources in order to increase productivity or disabling available resources according to low volume of requests in order to decrease cost of performing requests. Auto-scaling enables users to make larger or smaller infrastructures according to volume of activities, desirable efficiency and other dynamic behaviors [5]. Such automation increases advantages of cloud dynamic scalability substantially and can use more resources actively at the time of high workload also it can manage cost of using resources by disabling unnecessary resources at time of low requests. Efficiency index in cloud auto- scaling mechanisms is included the amount of using CPU, disk operations and bandwidth. Also we need accurate plan to have trade-off among these factors [6]. The Service quality in cloud services is an ability of dependent on resources and unlimited access based on the workload of the user at different times. Existence of these factors will affect the cost. Because service quality requires using resources with higher capacity, speed and power. On the other hand, utilization of the resources required to pay more. Therefore, we need a mechanism to manage the trade-off between these factors.

In this paper, in addition to use previous researches results in determining threshold [7-11] and volume of workload [12], we have tried to offer an auto-scaling

approach by learning automata [13], in order to management trade-off between rate of SLA violation and cost of scalability to preserve stability in system. Learning automata is able to select the best response among many received responses from environment. Using learning automata causes to have a better management in trade-off between cost and SLA violation also enable us to have simultaneously minimum cost, reduction in SLA violation and system stability (reduction in scalability operations).

The rest of paper is organized as follows: second part is dedicated to related works. In third part, we will explain our proposed approach and in forth part we will evaluate it. Finally the last part is dedicated to conclusion and future work.

## II. RELATED WORKS

Various researches have been accomplished in the field of auto scaling. The aim of some of them is to offer a scalability approach for a special application like server of web and others have done for optimizing the mechanism of scalability. In this part we are going to consider some of the previous researches about auto-scaling.

M. Humphrey et al [6], considered auto scaling in clouds in order to dedication fast resources with minimum cost for group of independent works. Results showed efficiency of the approach in fast resource dedication but it is not recommended for situation of unequal priority and if there is needed a special level of productivity.

Menasce et al [14], presented a structure for scheduling tasks with deadlines in a heterogeneous environment. In this project, a working set (DAG), is received as input, and it is assumed that there are available a variety of services for any of sub-jobs. Decision making on the timing of the sub-jobs is done by considering its desired performance and access cost. Timing process is done by using genetic algorithm.

N.Chohan et al [15], offered sample stabilization mechanism for accelerating the implementation of tasks in the context of cloud, to reduce their expenses. Overall, this research [16], has been done with the aim of providing a cost-effective structure for cloud services (profitability for cloud service provider) within the Service Level Agreement (SLA). Note that from the perspective of the customer, cost comparison is done between costs of implementation in the context of cloud and the implementation cost in the normal situation.

Roy Nilabja et al [17], have proposed an algorithm which does not use reaction approach in order to auto-scaling web resources. Rather, they have suggested a predictor solution by applying the concepts proposed by Sharif et al. [18,19], which can be used by systems that show a mixed behavior (continuous dynamic or discrete) and it has a large set with controlling limits. The base of approach is controlling principles of volume in the future that called advanced optimization. Optimal situation is

calculated through repetition in a determined period of time considering current and future limits.

Trieu C.Chiue et al [20], offered an approach to auto-scale web applications in a virtual cloud environment. The most important issue in web applications is inability to design or even to predict the number of active users. The mechanism includes a first-end load balancer, and a sub-system (to keep recording) and an intelligent control system (scaling index). Load balancer is utilizing for tracking and balance for users' requests in order to access to web services on the web servers in cloud VMs. Scaling is done according to number of active users. Results show that algorithm is able to manage high traffic meeting factors of cost and good efficiency.

Jing Qi Yang et al [21], offered cloud architecture for the purpose of auto-scaling resources based on anticipated workload. This method consists of two sub division pre-scaled and real time scaling. Also, the linear regression model has been used to predict traffic scaling is classified in three ways: self-healing scaling, scaling of source level and scaling of virtual machine level. The main idea in self-healing scaling is that two VMs will be able to operate overlapping. The purpose of anticipating the volume of work is estimation the number of service requests in a period of time. Deviation management between estimating traffic and real amount of traffic is the main issue in this approach.

## III. PROPOSED APPROACH

An efficient Scalable mechanism is able to meet the desired quality of service, also optimal cost for users, On the other hand, load balancer using appropriate distribution in system samples will be able to reduce the frequency of the system needed to be scaled up as much as possible to preserve stability of system. As mentioned before, trade-off between cost and SLA violation makes difficult to achieve such mechanism because having high quality services is required more expenses. On the other hand, cost minimization could face us with SLA Violation. We use our learning automata in our proposed approach to manage and trade-off between cost and rate of SLA violation. Following in this section firstly we describe needed scalability framework to implement proposed approach then briefly introduce the automata and finally we present our proposed algorithm.

### A. Scalability Framework

Determining and using an efficient scaling approach requires to create and utilize a suitable framework included effective factors in scaling mechanism. Therefore we offer a framework then the proposed scaling approach, will be presented based on this framework. Fig.1 shows the main components of scalability framework include: Load Balancer, Scaling Manager, Cost Manager, Virtualization Manager and Server Cluster.
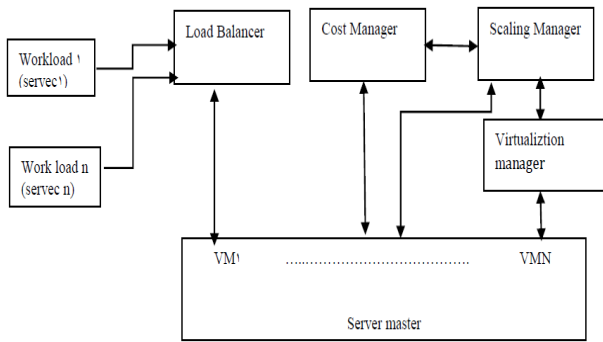
Fig.1. Scalability Framework

**Load balancer:** User requests, generated by the customer, are sent to the services through admission controller. First, admission controller filters out invalid requests then remained requests will be offered to load balancing. As there is more than one virtual machine on server cluster, quality of dynamic sending requests to the different VMs needs more precise management. So we use load balancer for this purpose. Load balancer collects information related to the various VMs conditions from Server cluster then sends requests to the suitable VM according to the information and based on load balancer strategy.

**Scaling manager:** scaling manager is the central component of the framework and its responsibility is monitoring and controlling the process of allocation resources to the requests in the clusters. The decision about time and quality of scaling is the responsibility of scaling manager. Scaling manager controls efficiency of cloud system at periodic intervals and determines strategy of scalability according to the software context. If the efficiency of system is less than the threshold value, license of releasing additional virtual machine will be issued in order to reduce the cost of scaling. Also, if the efficiency of system is more than threshold, management will issue the license of activation the virtual machine, in order to avoid additional SLA Violation.

**Virtualization manager:** resources on the clusters directly and operates the strategies offered by scalable management then adds or releases virtual resources.

**Costmanager:** cost manager is responsible to consider and select economical resources among SLA provider's resources. In fact whenever system needs to move into Scale up mode, cost manager operates optimal VM selection (simultaneous SLA and desired cost provider) according to learning automata technique. It means among the resources of SLA, instances with lowest cost will be selected during a certain number of repetitions and their probability will be raised. Also, if an instance doesn't have this situation, cost manager reduces its probability. At the end, Cost manager chooses the sample with highest probability and introduces it to scalable management for scaling up. Also at the time of scale down, above operations will be done to select the active instance with the highest cost to decrease the expenses of scalability, only difference is that firstly selection has been done among active instances and secondly at the end, the instance with highest penalty (lowest probability) will

be selected and introduced to the scalability management to make it idle (Scale down).

**Server Cluster:** as you can see in Fig.1, server cluster consists of a pool of VMs which have been used based on traffic and in order to offer services to the users to provide efficiency factors. In this paper, in order to ease the task, we have considered a separate virtual machine for any particular services.

### B. Learning Automata

Learning algorithms struggle to improve their performance and features toward a special aim by knowing related environmental situations. [22,23]. Learning automata is kind of such algorithms which try to change its conditions probabilities, based upon responses come from surrounding then it shows special reaction in any conditions. Learning automata [24] is an abstract thing with limited actions. A learning automata's performance is in a manner that it selects one of its actions among set of its actions and applies it to the environment. Then the mentioned action will be evaluated by a random environment and automata select its next action based on response of environment. The method which automata use it for selecting next action will be determined according to the used learning algorithm. Along the process automata learns how to select optimal action. Environment includes all internal and external conditions which affect automata. Generally it is possible to present environment as a set:

$E \equiv \{\alpha, \beta, c\}$ in which $\alpha$ is set of inputs, $\beta$ is set of outputs and $c$ is environment penalties. Fig. 2, shows relationship between learning automata and environment.
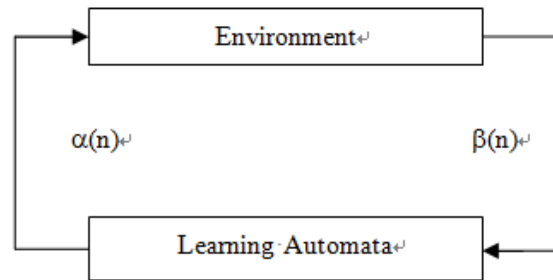


Fig.2. Relationship between Learning Automata and Environments

$\alpha$ is the number of actions automata is able to do, $\beta$ is output of environment which differs according to the model and kind of environment Various models defined for probable environments are divided to 3 groups: P-Model, Q-Model and S-Model. In Q-model, environment outputs are discrete values of zero and one and in S-model output is always a constant value between zero and one. In this paper, we use P-Model. In this model $\beta$ might be one of the two values of zero or one. Zero means desirable action and the probability of an action increases while probability of other actions decreases. One means undesirable action and probability of current action decreases and probability of other actions increase so that after receiving amount of one from environment, automata change the action. C is the set of probability of penalties for actions of automata which defines as follow:

$$c_i = \mathrm{Pr}\,ob[\beta(i) = 1 \mid \alpha(i) = \alpha_i], i = 1, 2, ...., r \qquad (1)$$

These values change over the time. Values of ci often are not clear and knowing them means thorough understanding of the potential environment that is not possible in most applications. Learning automata tries to know these values. In our sample environment will operate n times in any action:

- New input load (a) enters in environment.
- Includes reward and penalty comes from environment.
- Probability vector is as P={0.5, 0.5}.

According to automata operation, higher penalty in non-optimized sample, lower the chance of its selection so it can assure us of selection an optimum sample

If n repetition selects $a_i$ so that we have in (n +1) repetition:

Received a favorable response:

$$p_i(n+1) = p_i(n) + a[1 - p_i(n)] \qquad (2)$$

$$p_j(n+1) = (1-a)p_j(n) \quad \forall j \quad j \neq i \qquad (3)$$

Received an unfavorable response:

$$p_i(n+1) = (1-b)p_i(n) \qquad (4)$$

$$p_j(n+1) = \frac{b}{r-1} + (1-b)p_j(n) \quad \forall j \quad j \neq i \qquad (5)$$

### C. Proposed Algorithm

Proposed algorithm includes 3 sub-algorithms. First algorithm is responsible to manage scalability process. The process includes monitoring and measuring the system efficiency in order to perform scalability. Second algorithm is responsible to select optimum samples then introducing it to the scalable management at the time of activation a new sample (high scale). Finally, the third algorithm is responsible to select suitable sample at the time of deactivation operating samples (low scale) in order to saving cost of scalability.

Generally, the process works as an algorithm (Algorithm.1), measures efficiency of active virtual machines at the time of starting any period of time and periodically at the time of working then compares them with high and low thresholds. If the measured value is higher than threshold, algorithm calls scale up to increase VMs. Otherwise, if the measured value is lower than the threshold, algorithm will call scale down to decrease VMs. If there is none of these situations, cloud will stay in a stable condition.

---

**Algorithm 1: Scalability management**

Begin
While (the system is running and in the beginning of an interval)

    For (every $S_{i,j}$ in the cloud)

        Monitor $u_{i,j}(t)$

        If ($u_{i,j}(t) > u_{upper-threshold}$)
            Execute VM-level cost-aware scaling up

        If ($u_{i,j}(t) < u_{lower-threshold}$)
            Execute VM-level cost-aware scaling down
End

---

In the Scale up, firstly rate of SLA Violation is calculated because the amount of service requests has been increased more than the service capacity. Service capacity is specified according to the MIPS processor virtual machines of that service. Then requested amount of mentioned service– faced with lack of resources, will be calculated to meet the requests of service adding more virtual machines (Algorithm.2).

---

**Algorithm2: Cost-aware scaling up at *t* th interval**

Begin
Calculate SLA Violation
Calculate shortage Request
While (shortage Request is not empty)
    For (every VM in VM List)
        If (VM is suitable for shortage Request)
            Give reward to VM based on its cost
// cheapest VM has most reward
        else
            Give penalty to VM based on its cost
// expensive VM has most penalty
        Calculate chance(reward, penalty)
Cheapest VM← Select VM with biggest chance
// biggest chance is equal to cheapest VM
    Add Cheapest VM
End

---

When the capacity of cloud service is more than issued requests of users, scalable management will issue allowance for scale down to shut down additional VMs in order to save costs of scaling (Algorithm.3). Our criteria for determining the on removal of virtual machines is that servers still have the ability to respond to users requests after deleting the virtual machine.

---

**Algorithm 3: Cost-aware scaling down at t th interval**

Begin
Calculate extra Request
While (extra Request is not empty)

    For (every VM in $S_{i,j}$)

        If (VM can be removed)
            Give penalty to VM based on its cost
// expensive VM has most penalty
        Else
            Give reward to VM based on its cost
// cheapest VM has most reward
        Calculate chance(reward, penalty)
Expensive VM← Select VM with biggest chance
// biggest chance is equal to expensive VM
    Remove Expensive VM
End

---

## IV. PERFORMANCE EVALUATION

We are going to evaluate our proposed approach based on 3 criteria included cost, SLA violation and number of scaling. We have used Cloudsim [25] to simulate the proposed approach. We will use 4 kinds of VMs (Table 1) based on Amazon company VMs in simulation tests. Regarding variety in the existed services in cloud, we have performed four various kinds of services and we have not focused on a special service or program so that services are independent of program. The service combines all heterogeneous applications such as HPC, Web and more. The workload of the system has been modeled based on the normal distribution to make it closer to the real world.

Table 1. VM's Information

| VM type | MIPS (CPU) | Core | RAM (MB) | Price (cent) |
|---|---|---|---|---|
| Micro | 500 | 1 | 633 | 0.026 |
| Small | 1,000 | 1 | 1,700 | 0.070 |
| Extra Large | 2,000 | 1 | 3,750 | 0.280 |
| High-CPU Medium | 2,500 | 1 | 850 | 0.560 |

Scaling is operated in a 24-hour period (288 * 5-minute intervals) on four different services. We considered low threshold value of 0.2 and high threshold value as 0.8 in the simulation which have been used for comparing efficiency of active VMs in system. Also, we use the probability $P_i$ of an initial value of 0.5 and the probability $P_j$ initial value 0.5 in order to determine the level of chance, and penalty for samples of cloud. We have tried in selection periods of time to prevent from operational overhead and scarcity so that we have considered an average period of time.

We have compared our proposed approach to the two approaches, cost-aware scaling [26] and random scaling virtual machines in the cloud computing environment to evaluate our approach. Cost-aware policy is a simple and non-learning method which decreasing cost of services is its scaling priority. Also in the random policy, there is no management on virtual machines expenses. In other words, in the Scale down and Scale up operations, adding or decreasing VMs is done randomly. We evaluate our proposed approach based on 3 criteria included cost, the number of scaling (adding or deleting VMs) and rate of SLA violation.

**Cost:** The estimated cost of the cloud service is based on the amount of working hours. The users (cloud user), depending on the speed, power and capacity of demanded resource (CPU, memory, or disk...) and also the period of acquisition (minimum of acquisition time is one hour) of resource should pay the cost. Naturally, our costs will be lower when we use the resource with lower speed and capacity, in shorter time. Perhaps this attitude optimizes your costs, but it should be noted that other factors of quality of service, will be affected. Therefore, in order to achieve high quality service, we have to bear increasing costs.

**The number of scaling (system stability):** One of the issues that impact largely the dynamic scaling is the number of removed or added virtual machines. This firstly causes to accelerate responding in the computing environment and on the other hand the number of these processes plays an important role in estimation of provider's expenses. So, lower number of scaling means lower cost and higher speed in responses. Finally we will have a stable system with minimum cost.

**SLA Violation:** SLA Violation occurs when a provider fails to meet predefined criteria (Service Level Objectives (SLO)) in the SLA for the users. The number of all missed deadlines, lack of MIPS guarantee agreement, not guaranteed bandwidth agreement, the number of requests rejected due to lack of supply at peak times and so on, are examples of violation of the SLA.

We consider three scenarios to evaluate the proposed approach. In the first scenario we investigate three scaling policies in a cloud environment based on the assumptions (Table1) in terms of costs. In the second scenario we compare our scaling mechanism in terms of system stability (frequency of scaling) in a simulated environment with the characteristics mentioned above, to the Cost-aware and random methods and finally in the last scenario we compare all three methods in equal surrounding and situations in term of SLA violation.

### A. First Scenario

In the first scenario, we discuss cost of proposed approach against two other approaches. Cost is one important factor for user. The user is trying to run his/her request with the lowest possible cost. In Fig. 3, results of simulation have been offered according to the cost of scaling for all three compared approaches for 4 existed services in 24 hours. As you can see random approach has the worst result in scaling 4 mentioned services regarding to its greedy nature of performance as this approach selects its choice on the basis of current needs without regarding cost parameters whenever it needs to perform the scaling. Cost- aware non-automata approach has better results than the random approach in response to all requests. However it is obvious cost-aware approach based on learning automata can decrease costs of scaling substantially. Our proposed approach can be more successful in decreasing scaling expenses in considered services using reward and penalty technique.
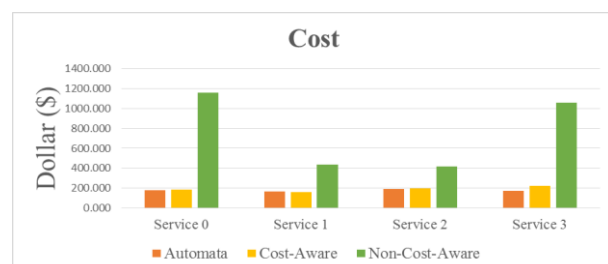


Fig.3. Comparing Cost of Different Approaches

In Fig.4, we have compared general cost of applying three approaches. According to the results of the comparison, our proposed approach has been able to

improve scaling costs up to 85% compared to random approach and 20 percent compared to cost-aware approach,
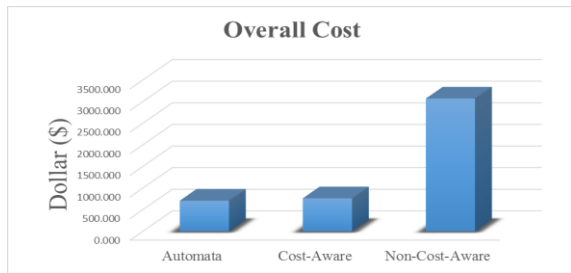


Fig.4. Comparing overall cost of different approaches

### B. Second Scenario

In the second scenario, we evaluate the number of scaling for proposed approach against the other two approaches described above. Previously stated VMs deletion or addition is done based upon violation of threshold values. The number of scaling which indicates the system is balanced or not, also can impose operational overhead and cost to system. So that suitable management for the criteria and choosing suitable threshold values can help the scaling approach to access highest response speed and lowest cost resulted decreasing rate of SLA violation. In Fig.5, the results of simulation have been represented for three approaches compared to 4 per day based on the number of scaling for each approach in the cloud. According to the results, the number of scalability is changing continually in random approach in response to the various requests therefore system doesn't have suitable stability. On the other hand these values are approximately equal in our approach and Cost-aware approach so we can say that system stability has been partly preserved.
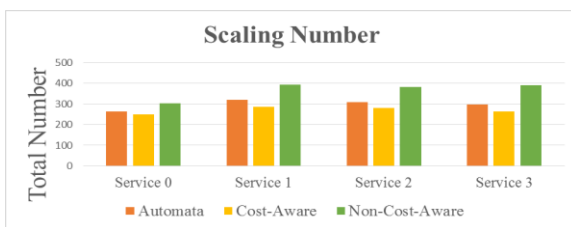


Fig.5. Comparing of the scaling number different approaches

Fig.6 shows general results of simulation the number of scaling in three scaling approaches. Accordingly, our proposed approach and cost-aware method have almost the same performance and better than random approach. So the system is more stable than random sampling.
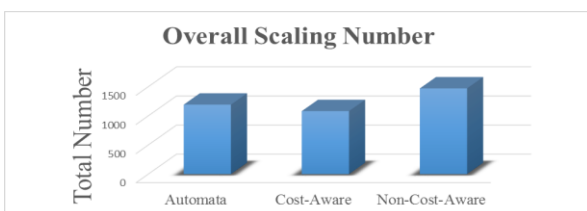


Fig.6. Comparing of Overall Scaling Number of Different Approaches

### C. Third Scenario

In the third scenario, we try to the evaluate criteria of SLA violation rate for the proposed approach compared to the other two approaches described above. The rate of SLA violation will be the least, when qualitative attributes will be provided such as //resource availability, high throughput and low respond time. In fact, when we will be able to keep the quality of the service at the optimal rate, the rate of SLA Violation for proposed approach will remain low; otherwise the amount will be increased. Increase of SLA Violation, indicates that quality service desired by user has not been met. In Fig.7, the results of the comparing values of SLA Violation have been shown for three compared approaches for 4 services in the cloud during 24 hours. It is obvious on the results of simulation the worst performance in SLA violation is in cost- aware approach. To reduce SLA violation, we should use instances with higher speed and power, and we should pay more for such services. Cost-aware approach has been encountered problems in SLA violation because lowest cost is its preference in selecting an operational sample to respond any kind of requests. Random approach in compare to other approaches has been able to improve its SLA violation because of paying high costs. Therefore none of two approaches (Cost-aware and random) have been able to manage trade-off between cost and rate of SLA violation. While our proposed approach used learning automata has been able to manage the trade-off. It means responding the requests of user with the lowest rate of SLA violation and the minimum possible cost.
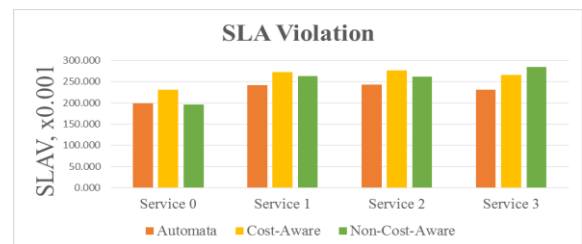


Fig.7. Comparing Rates SLA Violation of Scaling Approaches

In the Fig.8, the overall results of the simulation of three scaling approaches have been shown. Accordingly, in general our proposed method, compared to other methods, has been able to optimize rate of SLA violation for all services requested by user during the simulation. The optimal results of our approach are outcome of successful management the trade-off between cost and rate of SLA Violation.



Fig.8. Comparing Overall SLA Violation of Different Approaches

## V. Conclusion and Future Work

Cloud services are distributed infrastructure which develops services and communications. Scaling as one of the most important features of cloud computing, tries to allocate and pay back resources based on the requirements. Generally, we expect an efficient scaling mechanism will guarantee the quality of cloud services by the minimum cost. Important factors that have examined in the proposed approach are: the rate of SLA Violation, the cost of scaling and frequency of scaling. Using optimized virtual resources causes to decrease cost of a cloud service. On the other hand, efficiency improvement will be provided by using appropriate resources fixed by requests which cause to have higher costs. The frequency of scaling is a criterion to measure stability of system at the time or responding to the requests specifically at the peak of traffic. Quality assurance of a cloud service is dependent to meet the user service level agreement (SLA). In conclusion, we have trade-off between performance measures in our approach. Based on the evaluation results, the proposed scaling approach which has been offered based on learning automata can manage the trade-off between factors of SLA Violation and costs. The main objective of our Scalable oriented approach is management between two key benchmark including rate of SLA Violation, and cost. It is proposed as future work, to consider automatic scaling based on the automata can with respect to the deadlines. It should be noted that the deadline is applicable only for works in which the performance means efficiency not their catastrophic results. You can also use automatically scaling based on automata for Data intensive applications.

### References

[1] Anandhi, R. and K. Chitra. "A challenge in improving the consistency of transactions in cloud databases-scalability." International Journal of Computer Applications 52, no. 2 (2012): 12-14.

[2] Afife Fereydooni, Mostafa Ghobaei Arani and Mahboubeh Shamsi "EDLT: An Extended DLT to Enhance Load Balancing in Cloud Computing." International Journal of Computer Applications 108(7):6-11, December 2014.

[3] Md. Imran Alam, Manjusha Pandey, Siddharth S Rautaray,"A Comprehensive Survey on Cloud Computing", IJITCS, vol.7, no.2, pp.68-79, 2015. DOI: 10.5815/ijitcs.2015.02.09.

[4] Vaquero, Luis M., Luis Rodero-Merino and Rajkumar Buyya. "Dynamically scaling applications in the cloud." ACM SIGCOMM Computer Communication Review 41, no. 1 (2011): 45-52.

[5] Mao, Ming, Jie Li and Marty Humphrey. "Cloud auto-scaling with deadline and budget constraints." In Grid Computing (GRID), 2010 11th IEEE/ACM International Conference on, pp. 41-48. IEEE, 2010.

[6] Mao.Ming and Marty Humphrey. "Auto-scaling to minimize cost and meet application deadlines in cloud workflows." In Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis, p. 49. ACM, 2011.

[7] Monireh Fallah, Mostafa Ghobaei Arani and Mehrdad Maeen. "NASLA: Novel Auto Scaling Approach based on Learning Automata for Web Application in Cloud Computing Environment." International Journal of Computer Applications 113(2):18-23, March 2015.

[8] Monireh Fallah, Mostafa Ghobaei Arani "ASTAW: Auto-Scaling Threshold-based Approach for Web Application in Cloud Computing Environment." International Journal of u- and e- Service, Science and Technology (IJUNESST),Vol.8, No.3, pp.221-230, 2015.

[9] Hasan Masum Z., Edgar Magana, Alexander Clemm, Lew Tucker and Sree Lakshmi D. Gudreddi. "Integrated and autonomic cloud resource scaling." InNetwork Operations and Management Symposium (NOMS)، 2012 IEEE، pp. 1327-1334.

[10] Han Rui, Li Guo, Moustafa M. Ghanem and YikeGuo. "Lightweight resource scaling for cloud applications." In Cluster، Cloud and Grid Computing (CCGrid), 2012 12th IEEE/ACM International Symposium on, pp. 644-651.

[11] Maurer Michael, Ivona Brandic and RizosSakellariou. "Enacting SLA's in clouds using rules." In Euro-Par 2011 Parallel Processing, pp. 455-466. Springer Berlin Heidelberg, 2011.

[12] Dutreilh Xavier, Nicolas Rivierre, Aurélien Moreau، Jacques Malenfant and Isis Truck. "From data center resource allocation to control theory and back." In Cloud Computing (CLOUD) 2010 IEEE 3rd International Conference on، pp. 410-417.

[13] Narendra, Kumpati S. and Mandayam AL Thathachar. "Learning automata: an introduction." Courier Corporation, 2012.

[14] D. Menasc and E. Casalicchio. 2004. "A Framework for Resource Allocation in Grid Computing." In Proc. of the 12th Annual International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunications Systems, pp. 259-267, 2004.

[15] N. Chohan, C. Castillo, M. Spreitzer, M. Steinder, A. Tantawi and C. Krintz. 2010. See Spot Run: Using Spot Instances for MapReduce Workflows. In 2nd USENIX Workshop on Hot Topics in Cloud Computing, HotCloud2010. Boston, MA. June 2010.

[16] Y. Yazir, C. Matthews, R. Farahbod, S. Neville, etc. 2010. "Dynamic Resource Allocation in Computing Clouds using Distributed Multiple Criteria Decision Analysis." In 3rd International Conference on Cloud Computing, Miami, Florida, USA, 2010.

[17] Ro.Nilabja, AbhishekDubey, and AniruddhaGokhale. "Efficient autoscaling in the cloud using predictive models for workload forecasting." In Cloud Computing (CLOUD), 2011 IEEE International Conference on, pp. 500-507. IEEE, 2011.

[18] S. Abdelwahed, J. Bai, R. Su, and N. Kandasamy, "On the application of predictive control techniques for adaptive performance management of computing systems", Network and Service Management, IEEE Transactions on, vol. 6, no. 4, pp. 212 –225, dec. 2009.

[19] S. Abdelwahed, N. Kandasamy, and S. Neema, "A controlbased framework for self-managing distributed computing systems," in WOSS '04: Proceedings of the 1st ACM SIGSOFT workshop on Self-managed systems. New York, NY, USA: ACM, 2004, pp. 3–7.

[20] Chieu, Trieu C., Ajay Mohindra, Alexei A. Karve, and Alla Segal. "Dynamic scaling of web applications in a virtualized cloud computing environment." In e-Business Engineering, 2009. ICEBE'09. IEEE International Conference on, pp. 281-286. IEEE, 2009.

[21] Yang, Jingqi, Chuanchang Liu, Yanlei Shang, Bo Cheng, Zexiang Mao, Chunhong Liu, LishaNiu, and Junliang Chen. "A cost-aware auto-scaling approach using the workload prediction in service clouds." Information Systems Frontiers 16, no. 1 (2014): 7-18.

[22] K. Narendra and M. A. L. Thathachar, "Learning Automata: An Introduction", Prentice Hall, Englewood Cliffs, New Jersey, 1989.

[23] K. Najim and A. S. Poznyak, "Learning Automata: Theory and Application", Tarrytown, NY: Elsevier Science Ltd., 1994.

[24] Behnaz Seyed Taheri, Mostafa Ghobaei Arani and Mehrdad Maeen. "ACCFLA: Access Control in Cloud Federation using Learning Automata." International Journal of Computer Applications 107(6):30-40, December 2014.

[25] Calheiros, R. N., Ranjan, R., Beloglazov, A., De Rose, C. A., & Buyya, R. (2011). "CloudSim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms." *Software: Practice and Experience*, *41*(1), 23-50.

[26] Yang, Jingqi, Chuanchang Liu, Yanlei Shang, Bo Cheng, ZexiangMao, Chunhong Liu, LishaNiu, and Junliang Chen. "A cost-aware auto-scaling approach using the workload prediction in service clouds." Information Systems Frontiers 16, no. 1 (2014): 7-18.

**Mostafa Ghobaei Arani** received the B.S.C degree in Software Engineering from IAU Kashan, Iran in 2009, and M.S.C degree from Azad University of Tehran, Iran in 2011, respectively. He's Currently a PhD Candidate in Islamic Azad University, Science and Research Branch, Tehran, Iran. His research interests include Grid Computing, Cloud Computing, Pervasive Computing, Distributed Systems and Software Development.

**Mahboubeh Shamsi** is Associate Prof. Qom University of Technology. She's received the B.S.C degree in Mathematics from Isfahan University, Iran in 2003, and received the M.S.C degree from Azad University of Isfahan, Iran in 2006, and also received the PhD degree in Software Engineering from UTM in 2011. Her research interests include Image Processing, Design and Implementation of Persian/Arabic OCR, Design and Implementation of Biometric Authentication System, Design and Test of Reliable Software, Databases, Software Engineering, UML, ERD, DFD, ERP.

## Authors' Profiles

**Khosro mogouie** received the B.S.C degree in Software Engineering from University Azad Arak, Iran in 2009, and M.S.C degree from Azad University of mahallat, Iran in 2014, respectively. Her research interests include Cloud Computing, Distributed Systems and Cloud Scalability technology.