



## ASQL: A New Approach for Resource Auto-Scaling Using Q-Learning in Cloud Computing Environment

Bahar Asgari<sup>1</sup> and Mostafa Ghobaei Arani<sup>2</sup>

<sup>1</sup>Department of Computer Engineering, Mahallat Branch, Islamic Azad University, Mahallat, Iran

<sup>2</sup>Department of Computer Engineering, Parand Branch, Islamic Azad University, Tehran, Iran

E-mail: <sup>1</sup>[bahar\\_asgari88@yahoo.com](mailto:bahar_asgari88@yahoo.com), <sup>2</sup>[mostafaghobaei@piau.ac.ir](mailto:mostafaghobaei@piau.ac.ir)

### ABSTRACT

Cloud services have become more popular among users these days. Providing automatic resource for cloud services is one of the important challenges. In cloud computing environment, resource providers shall offer required resources to users automatically without any limitations. It means whenever a user needs more resources, the required resources should be dedicated to the users without any problems. On the other hand if resources are more than user's needs extra resources should be turn off temporarily and turn back on whenever they needed. In this article the approach is to represent enforcement learning-aware for auto-scaling resources according to Markov decision process (MDP). Results would show the rate of SLA (service level agreement) violation and stability that proposed better functions compared to the similar approaches.

Keywords: *Cloud Computing, Scalability, Auto-Scaling, Reinforcement Learning.*

### 1 INTRODUCTION

Cloud computing is the number of virtualized connected computers which offers single computational resource dynamically and to compute complex computation [1, 2]. In other word cloud computing states to both applicable programs offered as services on the Internet, hardware and software systems in data centers. By definition, in data center hardware and software are called "cloud". Scalability is one of the basic concepts in cloud computing which is important in using higher efficiency of cloud computing [4]. Scalability is referred to increase system functional power to have suitable response against increased work load of course by adding software and hardware resources [5]. Whereas applicable programs, especially application program on web, do not have regular work load patterns so that scalability functions (increase or decrease of scale) should have be done immediately with minimum human intervention to provide resources for applications as soon as possible. Resource scaling with minimum human intervention is called auto-scaling [6]. Various workloads are of the biggest challenges in

different times, so whenever provider wants to meet all the requirements in all times, it should reserve maximum needed resources previously for peak work load to support them. In this situation provider sometimes will be over-provisioning and it is going to be very costly for them (to buy maximum resources at peak times) which leads to lower profit. Therefore functional expenses will be reduced by turning off idle nodes on idle times, but it cannot decrease financial expenses related purchasing and hosting IT equipment's and their depreciation. If provider possesses only enough resources (average capacity) to support average number of requests, the providers may be utilized, but the provider might not have enough local resources to meet clients' request which leads to under-provisioning in some situations so provider has to reject new customers or cancel previous services operating on system. We should design a system which will be able to manage uncertainty and remove any problems in cloud environment. Also it should be able to impact parameters like expense, efficacy, SLA violation etc.

We offer auto-scaling according to reinforcement learning. Reinforcement learning (RL) is a kind of

decision making that determines a goal performing functional model, applies policy without previous information. RL has been performed successfully in extensive fields to support auto-control and dedicate resources [7-10] which works on the basic assumption of penalty and reward so the factors move toward operations which lead to highest profit. Major part of RL is on the basis of determination of optimal policies in Markov [11].

In this paper we want to propose auto-scaling approach using MDP to manage SLA violation and scaling expense and to preserve system stability. RL has the capacity to response suitably using environment experiences. RL leads to better management of compromise SLA violation and number of scales but it causes higher expenses.

The rest of this paper is organized as follows: we review related works about RL in second part; the proposed approach comes in third part in detail. The evaluation of proposed approach will be explained in fourth part. Finally conclusion and suggestion will be presented in fifth section.

## 2 RELATED WORKS

Various studies have been carried out about auto-scaling and its implementation. Current approaches have advantages and disadvantages. As the proposed approach in this paper is based upon RL, we review researches related to this technique in this section.

- Enda Barrett et al. [11] have been considered the parallel Q learning to reduce time of determination about optimal policies and online learning. Their proposed approach uses MDP along with RL.
- Fouad Bahrpeyma et al. [12] suggests RL-DRP approach. They use neural networks in their proposed mechanism. The approach enable cloud service providers to meet high volume of requests without wasting any time, valuable work and at the same time control resources optimally.
- Xavier Dutreilh et al. [13] have proposed using proper initialization in primary stages also increasing the rate of convergence in process of learning to solve problem. They have offered experiments results. Also they have introduced an efficient model to detect changes then completed learning process management based on that.
- Bauer et al. [14] proposed using RL to manage threshold orders. First controller applies these

orders to the goal program to reinforce its quality features. Second controller supervises the orders, adapts thresholds and changes conditions, also it deactivates unrelated orders.

- Jia Rao et al. [15] represent a RL aware virtualized machine configuration (VCONF). Central design of VCONF is prepared based on RL aware model to scale and adapt.
- Amoui et al. used RL successfully in management quality of web programs to optimize program's output [16]. Using simulation in initialization of learning functions is one of the interesting aspects of it.

Finally table 1 shows the comparison of above techniques.

Table 1: Comparison of technique

Reference	Auto-scaling technique	Advantages and disadvantages	contribution
Enda Barrett[11]	Parallel Q learning	Decreasing time of optimal policy determination and online learning Disadvantages: challenges in determination of initial policies	It uses inherent parallelism in distributed computing platforms like cloud
Fouad Bahrpeyma[12]	RL	Fast convergence process Higher productivity	It introduces a new decision making process to use predictably analysis of demand which considers parameters of offer and demand
Xavier	RL	Horizontal	Integration

Dutreilh[13]		scaling Increase in convergence rate in learning stages	in a real cloud controller and auto programming
Bahati[14]	RL	It limits situation as pair of operation-condition and provides the possibility of re-use of learned models in an order set for next stage Reinforcement of load based on effective limit	They proposed using RL to manage threshold orders. First controller applies order to the goal program to improve the features of quality
JiaRao[15]	VCONF	VCONF is good adaptation with online auto configuration policies with heterogeneous VMs VCONF is enable to guide initial setting without decreasing in function of VMs	Central design of VCONF using RL aware model works to scale and adapt
Amoui[16]	RL	Quality management in application of new web design to optimize program output	Using simulation for initialization of learning functions

### 3 PROPOSED APPROACH

Final goal is to make auto scaling system to have the ability of decreasing costs and increasing system stability; at the same time with SLA requirements and system efficiency It means to use an online policy to dedicate resources with scaling automatically. Proposed approach will be introduced according to RL and MDP. The offered MDP constitutes from 4 categories included conditions, operations, transmitted possibilities and rewards so that decision making about scale up/scale down will be accomplished based on it.

#### 3.1 Reinforcement Learning (RL)

RL [6, 11, 12] is a computational approach to understand automatic base learning to make the best decisions. It insists on learning via direct involvement of agent and environment. Decision maker refers to the agent who learns from experience and its best action is to perform at its maximum in any environment. An auto scaler is responsible for decisions about scaling without human involvement and its objective is to adapt resources dynamically to the applications according to input work load. It decides to allocate or deallocate resources to the applications based on work load. In any  $t$  time which  $t = 0, 1, 2, \dots$  time sequences are separated, agent shows condition of environment  $s_t \in S$  where  $S$  is all possible conditions and it selects  $a_t \in A(s_t)$  where  $A(s_t)$  is all variables in the condition of  $s_t$ , but in a determined time, sequence of these functions and agent will be the next reward  $r_{t+1}$  which finds itself in new condition of  $s_{t+1}$ . Agent will select from condition possibilities then operates the possible action. This will be agent's policy that shows  $\pi$  in which  $\pi(s, a)$  as  $a = a_t$  at the condition  $s = s_t$ .

So MDP can be shown in four categories included conditions, operations, transmitted possibilities and rewards:

- **S:** Environmental state space
- **A:** total action space
- **P(.|s, a)** defines distribution of governed possibilities on transmitted conditions
 
$$s_{t+1} \sim p(. | s_t, a_t)$$
- **Q(.|s, a)** defines distribution of governed possibilities on received reward.

$$R ( s_t , a_t ) \sim q ( . | s_t , a_t )$$

The objective of learning process inside learning Q is to achieve the optimal policy which reflects by Q amount in general reward and continues by operating in current situation. The amount of Q will be calculated by equation 1 which includes discounted reward (decreased reward) and shows RL process policy.

$$Q ( s_t , a_t ) \leftarrow Q ( s_t , a_t ) + \alpha ( r_{t+1} + \gamma \max_a Q ( s_t , a ) - Q ( s_t , a_t ) )$$

Where  $r_{t+1}$  is medium received reward of selecting  $a_t$  in  $s_t$  condition.  $\alpha$  is learning rate and  $\gamma$  is discount coefficient (Reduction). The overall process of RL has been shown in Algorithm 1:

**Algorithm1:** Reinforcement Learning Algorithm (Q- learning)

1. Initialize Q(s,a) arbitrarily
2. Repeat ( for each episode)
3. Initialize s
4. Repeat
5. Choose a from s using policy derived from Q (ε- greedy)
6. Take action a and observe r,s'
7.  $Q ( s_t , a_t ) \rightarrow Q ( s_t , a_t ) + \alpha ( r_{t+1} + \gamma \max_a Q ( s_t , a ) - Q ( s_t , a_t ) )$   
 $s \leftarrow s'$ ;
8. Until s is terminal

### 3.2 Proposed Algorithm

The proposed algorithm has been offered based on RL, that is defined according to Markov process for auto scaling a MDP. Upper and lower threshold have been defined too and cloud service operation will be monitored after introduction proposed MDP.

Configuration of proposed MDP has been considered as follows:

S: the space of condition : Full utilization, Under utilization, Normal utilization.

A: the space of operation: Scale up, Scale down, No- op.

P(.|s, a) defines distributionpossibility governed on transmitted conditions.

Q(.|s, a) defines distribution possibility governed on received reward.

As we know MIPS means the number of instructions per second. There has been introduced two variables for proposed approach included available MIPS and Requested MIPS, both are variables of service inputs. Q Updating will be done using local regerssion according to history of instructions. Division of two amounts shows the amount of utillization (equation 2) and comparison of upper and lower thresholds determine space of condition.

$$\text{Utilization} = \text{Available MIPS} / \text{Requested MIPS}$$

Equation 3 shows the full utilization condition.

$$\begin{aligned} & (\text{Requested MIPS} / \text{Available MIPS}) > (\text{High-Threshold}) \Rightarrow \\ & (\text{Full-Utilization}) \Rightarrow (\text{Under-Provisioning}) \end{aligned}$$

Equation 4 shows the under utilization condition.

$$\begin{aligned} & (\text{Requested MIPS} / \text{Available MIPS}) < (\text{Low- Threshold}) \Rightarrow \\ & (\text{Under- utilization}) \Rightarrow (\text{Over- Provisioning}) \end{aligned}$$

Otherwise the condition will be normal.

After defining full utilization, under utilization and normal conditions, and operating equation 1(Q) (s, a), SLA violation amount will be acquired by Requested MIPS and available MIPS then decision will be made according to above functions to do current action, it means to increase or decrease virtual machine or no operation. table 2 represents process of decision making and figure 1 shows a diagram included provider condition changes regarded to productivity parameter.

Table 2: Decision making by Markov scaling

	Utilization>High-Threshold	Low-Threshold<Utilization < High - Threshold	Utilization <Low-Threshold
State(t)	Full-Utilization (Under-Provisioning)	Normal-Utilization (Normal-Provisioning)	Under-Utilization (Over-Provisioning)
Next-Action(t+1)	Scale_Up	No-op	Scale_Down

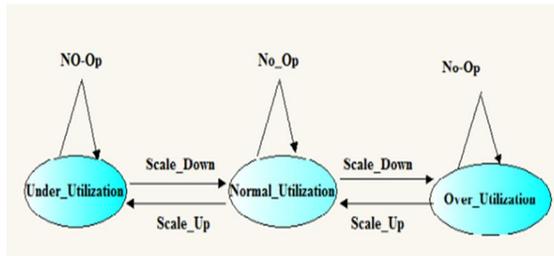


Fig. 1. Condition of Provider Changes Regarded to Productivity Parameter

Proposed algorithm introduced in this paper is represented as semi code offered in algorithm 2 according to Markov and decision making in table 2.

---

Algorithm2: Reinforcement Learning (Q-Learning)

---

1. Initialize  $Q(s, a)=0$ ,  $s=0$ ,  $a=0$ , high Range  $Q=0.8$ , low Range  $Q=0.2$ .
  2. Observe the Available MIPS and Requested MIPS.
  3. Observe the current state  $S$ .
  4. If (requested MIPS/available MIPS) > high Range  $Q$ , state  $[0]= 0$ ; /\*Full-Utilization state\*/
  5. Else if (requested MIPS/available MIPS) <Low Range  $Q$ , state $[1]= 1$ ; /\* Under-Utilization state \*/
  6. Else state  $[2]= 2$ ; /\* Normal-Utilization state\*/
  7. Loop
  8. Select action, choose for state ,based one of the action selection policy Utilization
  9. Take action, observe  $r$ , as well as the new state,  $s'$ .
  10. Update  $Q$  value for the state using the Regression and observed  $r$  and the maximum reward possible for the next state.
  11. 
$$Q(s_t, a_t) \rightarrow Q(s_t, a_t) + \alpha (r_{t+1} + \gamma \max_a Q(s_t, a) - Q(s_t, a_t))$$
  - Set the state  $s$  to the new state  $s'$ ,  $s \leftarrow s'$
  12. Until  $s$  is terminal
- 

#### 4 PERFORMANCE EVALUATION

There has been used of Cloudsim [17] simulator for simulation. Four kinds of virtual machine corresponded to Amazon EC2 [18] have been performed which their specifications offered in table 3. There have been used four kinds of services

regarded to the variety of available services in cloud and we have not focused on type of service or special program so that used services are independent to programs. These services are combination of all heterogeneous programs like HPC, Web and so on. Also work load has been modeled according to normal distribution to be closer to real world. Scaling will be done in 24 hour period and in 5 minutes intervals (288 five minutes), Low-Threshold is considered 0.2 and High-Threshold is considered 0.8 Standard deviation is 3000 MIPS and Diff Range is 0.4 There has been considered a function for initialization cost. As cost function is computed by hour and we have 5 minutes intervals so that we have to multiple overall costs by 300/3600.

Table 3: Specification of virtual machine

Type Of Virtual Machine	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

Algorithm works by updating  $Q$ . we have done  $Q$  updating and obtaining utilization by local regression. Updating  $Q$  will be accomplished according to instruction history; it means next amount will be determined according to prediction of previous amount. Predicted amount should be multiplied by 0.7, because error possibility has been considered as 30 percent. Regression function helps us to scale VMs in the way that decrease failed case along with minimum cost.

The amount of Available MIPS and Requested MIPS is calculated in the main function of proposed approach. Also the amount of SLA violation is calculated using their difference. Requested MIPS is divided to Available MIPS to introduce two dimensional array  $r[rstate]$  [0],  $r[rstate][1]$  and  $r[rsatate]$  [2] which determines underutilization, full utilization and normal functions. Then the amount of  $r$  will be updated and the amount of equation  $Q(s, a)$  will be calculated.

Overall MIPS amount is obtained by division of violation rate to the total MIPS then optimal current action will be selected using **Current Action()** and **Select action()** according to the amount of utilization. Then decision for scale down and scale up or null action will be made.

Proposed approach which is learner automata aware will be compared to cost aware auto scaling approach which is a simple, automata approach by parameters like cost, SLA violation, initialization cost and number of scaling. There has been defined three scenario for evaluation proposed approach, in table 4:

Table 4: Evaluation scenarios

Scenario	goal
First scenario	Evaluation SLA violation
Second scenario	EvaluationCost
Third scenario	Evaluation number of scaling

#### 4.1 First scenario

Evaluation of SLA violation has been considered in first scenario compared to two other approaches. SLA violation will be happened when provider cannot provide predefined measures in SLA for users. Some examples of SLA violation is the number of lost deadlines, lack of warranty on agreed MIPS, lack of warranty on agreed bandwidth, number of rejected requests because of not having enough resources at the peak times.

Increasing rate of SLA violation causes lower quality in providing services for user. If Requested MIPS is not match with available MIPS, SLA violation will happen. Figure 1 represents results of comparison of SLA violation in three compared approaches for 4 services. As you can see SLA violation rate in proposed approach is less than the others.

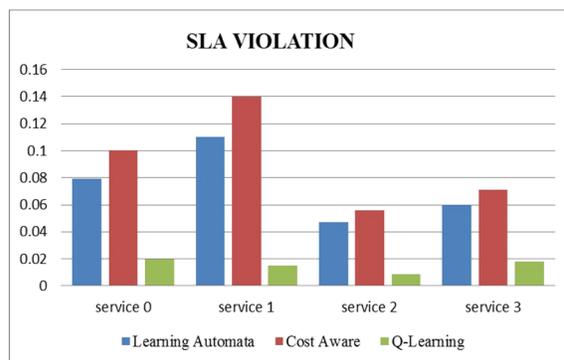


Fig. 1. Comparison of SLA violation in services

Figure 2 shows the comparison of overall SLA violation for services included cost aware, learning automata and proposed approach. As you can see results of proposed approach simulation compared to learning automata and cost aware approach has

lower rate of SLA violation at the time of simulation so that using Q learning technique in auto scaling leads to reduce SLA violation. So whenever SLA is important for auto scaling, we can use proposed approach.

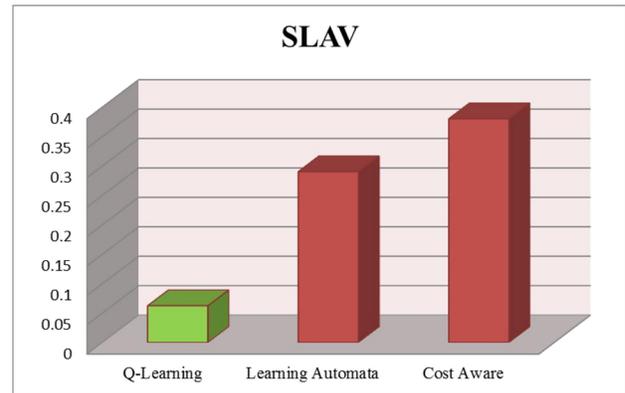


Fig. 2. Comparison of overall SLA violation in three approaches

#### 4.2 Second Scenario

We address evaluation of cost measure and comparing it with other approaches. Service cost will be calculated according to hours of utility. It means user pays the cost according to speed, power and capacity of requested resource (CPU, Memory, and disk and ...) also time of using resource. Naturally cost will be low when we use resource with lower speed and capacity in lower intervals. It can decrease cost but affects other quality factors. So to have a high quality service we have to increase cost. Cost is one of the most important factors for users. It means that user always struggles to accomplish the request with minimum cost. Parameter of cost introduces in three categories: Initialization cost, Runtime cost and Total cost. Initialization cost is initial cost for setting up VMs. Runtime cost equals to cost according to utility per hour which will be paid for VMs operation. Total cost calculates by equation 5:

(5)

$$\text{Total Cost} = \text{Initialization Cost} + \text{Runtime Cost}$$

Runtime cost of VM increases nearly 20 percent in simulation of vertical scale up. Cost measure in simulation is calculated according to addition of VM initialization cost and VM runtime cost.

Figure 3 shows VM initialization cost. It specifies that a Q aware approach has high initialization cost while an automata aware approach will save initialization cost substantially.

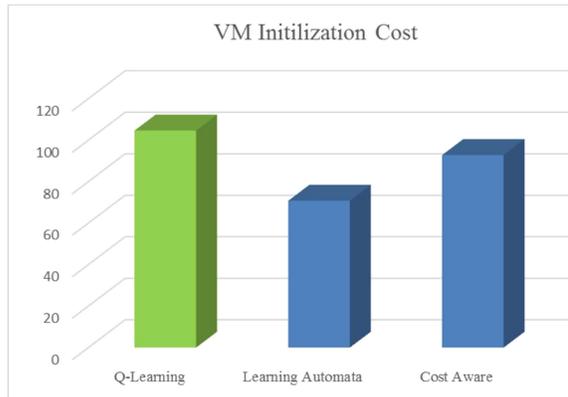


Fig. 3. Comparison of initialization cost in three approaches

Figure 4 represents VM runtime cost in 3 services in 24 hours. Results of simulation shows proposed approach has lower runtime cost.

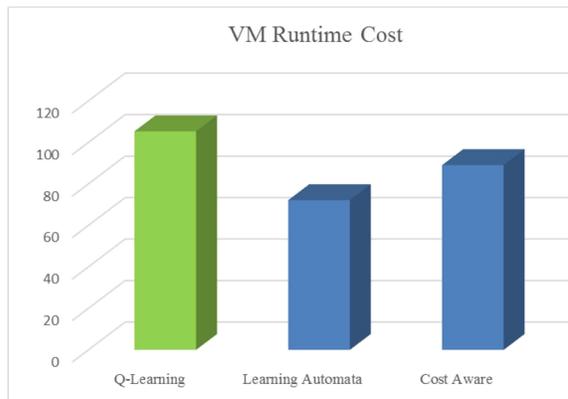


Fig. 4. Comparison of VM runtime cost in three approaches

Figure 5 represents results of simulation according to total cost of scaling for 3 compared approaches in a 24 hour period.

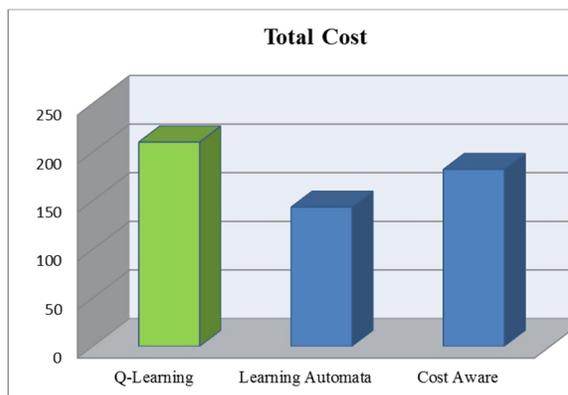


Fig. 5. Comparison of total cost in three approaches

Results of simulation according total cost in 24 hours for four services have been represented in Figure 6.

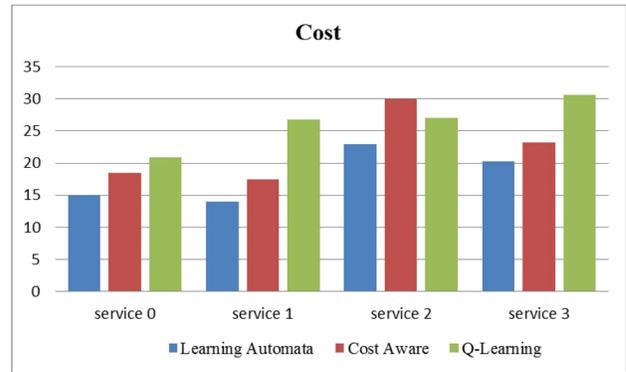


Fig. 6. Comparison of overall the total cost in three approaches

As it is obvious in Figures 5 and 6, scaling aware of learning automata has lower cost compared to proposed approach and cost aware approach. Proposed approach in this paper has the highest total cost in comparison.

Q aware approach has high initialization and runtime costs. Finally total cost which is addition of two mentioned costs shows that the approach will not be proper approach compared to learning automata and cost aware approaches whenever the cost measure is considered.

### 4.3 Third scenario

In third scenario we address comparing number of scales with the other two approaches. The number of elimination or adding VMs is one of the important factors in dynamic scaling. It affects speed of response in computing environment. Also it can cause to operational overload and imposes cost to the system. Proper management of the measure helps us to achieve minimum cost, increase rate of response, consequently reduction the rate of SLA violation. Overall scaling function calculates total number of scales. The number of scaling functions for four services has been shown in Figure 7. As you can see the number of scaling functions for proposed approach will not change dramatically so that system will have a proper stability.

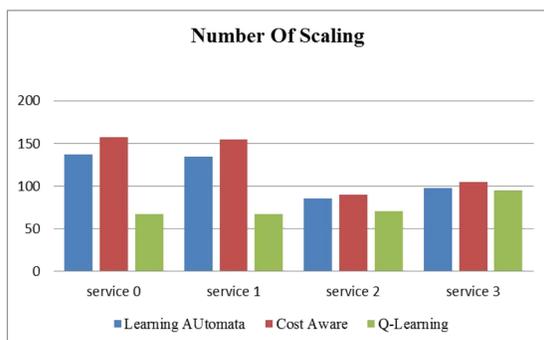


Fig. 7. Comparison the number of scaling in three approaches

As it is represented in Figure 8, the number of scaling functions has been decreased in proposed approach compared to two other approaches according to results of simulation. Reduction helps to optimize SLA violation rate, lower cost and higher system stability.

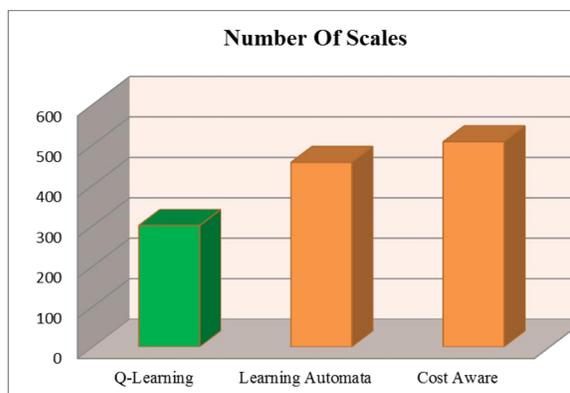


Fig. 8. Comparison of overall the number of scaling in three approaches

## 5 CONCLUSION AND FUTURE WORK

Cloud services are distributed infrastructures which extend space of communication and service. The resource providing has been very important because of daily grow of cloud services and scaling issue has been welcomed as one of the most important features of cloud computation. In this paper we have represented an approach based upon reinforcement learning also have addressed Markov model. There are 3 important factors in proposed approaches including SLA violation rate, scaling cost and number of scales. Regarding cost measure, Q aware approach is not proper approach compared to the automata and cost aware approaches. But proposed approach reduces number of scales which leads to optimize rate of SLA violation and system stability. Also proposed approach decreases SLA violation and optimizing SLA leads to increase

cost. As a result it makes difficult having minimum cost. On the other hand focusing on the minimum cost leads to SLA violation. So we can observe substantial reduction in SLA violation and higher system stability by using Q-learning technique in auto scaling.

Therefore it is possible to continue studies about auto-scaling regarded other effective factors and other approaches for example the condition space will be changed according to productivity or we can represent a novice approach in auto-scaling using parallel Q learning and combination of parallel factor and new condition. Also we can apply RL to predicate load in web aware software's. Also it is possible to merge RL and machine learning. Overload in proposed approach should be consider carefully too.

## 6 REFERENCES

- [1] Buyya, R., Broberg, J., & Goscinski, A. M. (Eds.). (2010). Cloud computing: principles and paradigms (Vol. 87). John Wiley & Sons.
- [2] Buyya, R., Yeo, C. S., Venugopal, S., Broberg, J., & Brandic, I. (2009). Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation computer systems*, 25(6), 599-616.
- [3] Yang, J., Liu, C., Shang, Y., Cheng, B., Mao, Z., Liu, Chen, J. (2014). A cost-aware auto-scaling approach using the workload prediction in service clouds. *Information Systems Frontiers*, 16(1), 7-18.
- [4] Roy, N., Dubey, A., & Gokhale, A. (2011, July). Efficient autoscaling in the cloud using predictive models for workload forecasting. In *Cloud Computing (CLOUD)*, 2011 IEEE International Conference on (pp. 500-507). IEEE.
- [5] Kupferman, J., Silverman, J., Jara, P., & Browne, J. (2009). Scaling into the cloud. *CS270-advanced operating systems*.
- [6] Lorido-Botran, T., Miguel-Alonso, J., & Lozano, J. A. (2014). A review of auto-scaling techniques for elastic applications in cloud environments. *Journal of Grid Computing*, 12(4), 559-592.
- [7] Hu, R., Jiang, J., Liu, G., & Wang, L. (2014). Efficient Resources Provisioning Based on Load Forecasting in Cloud. *The Scientific World Journal*, 2014.
- [8] Chevaleyre, Y., Dunne, P. E., Endriss, U., Lang, J., Lemaitre, M., Maudet, Sousa, P. (2006). Issues in multiagent resource allocation. *Informatica*, 30(1).

- [9] Jacyno, M., Bullock, S., Payne, T., & Luck, M. (2007, May). Understanding decentralised control of resource allocation in a minimal multi-agent system. In Proceedings of the 6th international joint conference on Autonomous agents and multiagent systems (p. 208). ACM.
- [10] Scalas, E., Gallegati, M., Guerci, E., Mas, D., & Tedeschi, A. (2006). Growth and allocation of resources in economics: The agent-based approach. *Physica A: Statistical Mechanics and its Applications*, 370(1), 86-90.
- [11] Barrett, E., Howley, E., & Duggan, J. (2013). Applying reinforcement learning towards automating resource allocation and application scalability in the cloud. *Concurrency and Computation: Practice and Experience*, 25(12), 1656-1674.
- [12] Bahrpeyma, F., Haghghi, H., & Zakerolhosseini, A. (2015). An adaptive RL based approach for dynamic resource provisioning in Cloud virtualized data centers. *Computing*, 1-26.
- [13] Dutreilh, X., Kirgizov, S., Melekhova, O., Malenfant, J., Rivierre, N., & Truck, I. (2011). Using reinforcement learning for autonomic resource allocation in clouds: Towards a fully automated workflow. In ICAS 2011, the Seventh International Conference on Autonomic and Autonomous Systems (pp. 67-74).
- [14] Bahati, R. M., & Bauer, M. (2010, April). Towards adaptive policy-based management. In *Network Operations and Management Symposium (NOMS), 2010 IEEE* (pp. 511-518). IEEE.
- [15] Rao, J., Bu, X., Xu, C. Z., Wang, L., & Yin, G. (2009, June). VCONF: a reinforcement learning approach to virtual machines auto-configuration. In Proceedings of the 6th international conference on Autonomic computing (pp. 137-146). ACM.
- [16] Amoui, M., Salehie, M., Mirarab, S., & Tahvildari, L. (2008, March). Adaptive action selection in autonomic software using reinforcement learning. In *Autonomic and Autonomous Systems, 2008. ICAS 2008. Fourth International Conference on* (pp. 175-181). IEEE.
- [17] Rodrigo N. Calheiros, Rajiv Ranjan, Anton Beloglazov, César A. F. De Rose, and Rajkumar Buyya, *CloudSim: A Toolkit for Modeling and Simulation of Cloud Computing Environments and Evaluation of Resource Provisioning Algorithms*, *Software: Practice and Experience (SPE)*, Volume 41, Number 1, Pages: 23-50, 2011.
- [18] Amazon EC2 instance types, <http://aws.amazon.com/ec2/>.
- [19] Khosro Mogouie, Mostafa Ghobaei Arani, Mahboubeh Shamsi, "A Novel Approach for Optimization Auto-Scaling in Cloud Computing Environment", *IJMECS*, vol.7, no.8, (2015), pp.9-16.