

AUTOMATIC SCALING OF INTERNET APPLICATIONS FOR CLOUD COMPUTING SERVICES

Kotta Pratap Reddy¹, Mr.Y.Durga Prasad², Dr.Y.Venkateswarulu³

¹MtechStudent ,CSE, Giet Engineering College, Rajahmundry, A,P, India

²Asst. Professor, Dept of CSE, Giet Engineering College, Rajahmundry, A,P, India

³Professor and HOD, Dept of CSE, Giet Engineering College, Rajahmundry, A,P, India

ABSTRACT

An automatic scaling property utilized by many Internet applications in the cloud service provider and get benefit from where their resource usage can be scaled up and down automatically. The present paper proposes a system that provides automatic scaling for Internet applications in the cloud environment. The present method encapsulate each application instance inside a virtual machine (VM) and use virtualization technology to provide fault isolation and this method is called Class Constrained Bin Packing (CCBP) problem where each server is a bin and each class represents an application. The class constraint reflects the practical limit on the number of applications a server can run simultaneously. The present paper develops an efficient semi-online color set algorithm that achieves good demand satisfaction ratio and saves energy by reducing the number of servers used when the load is low. Experiment results demonstrate that the proposed system can improve the throughput by 180% over an open source implementation of Amazon EC2 and restore the normal QoS five times as fast during flash crowds.

Index Terms—Cloud computing, virtualization, auto scaling, CCBP, green computing

INTRODUCTION

One of the most useful benefits of cloud computing service is the resource flexibility: a business customer can scale up and down its resource usage as needed without upfront capital investment or long term commitment. The Amazon EC2 service [1], for example, allows users to buy as many virtual machine (VM) instances as they want and operate them much like physical hardware. However, the users still need to decide how much resources are necessary and for how long. A user only needs to upload the application onto a single server in the cloud, and the cloud service will replicate the application onto more or fewer servers as its demand comes and goes.

Fig. 1 shows the typical architecture of data center servers for Internet applications. It consists of a load balancing switch, a set of application servers, and a set of backend storage servers. The front end switch is typically a Layer 7 switch [2] which parses application level information in Web requests and forwards them to the servers with the corresponding applications running. Each application can run on multiple server machines and the set of their running instances are often managed by some clustering software such as WebLogic [3]. Each server machine can host multiple applications. The applications store their

state information in the backend storage servers. The storage servers may also become overloaded, but the focus of this work is on the application tier. The Google AppEngine service, for example, requires that the applications be structured in such a two tier architecture and uses the BigTable as its scalable storage solution [4].

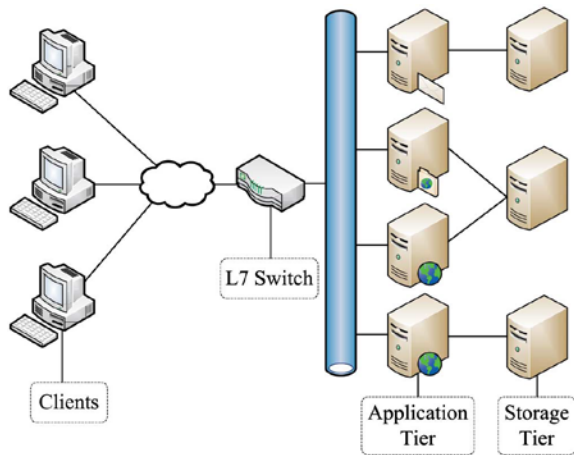


Figure 1: Internet applications two-tiered architecture

Even though the cloud computing model is sometimes advocated as providing infinite capacity on demand, the capacity of data centers in the real world is finite. The illusion of infinite capacity in the cloud is provided through statistical multiplexing. The present method defines the demand satisfaction ratio as the percentage of application demand that is satisfied successfully. The amount of computing capacity available to an application is limited by the placement of its running instances on the servers. Various studies have found that the cost of electricity is a major portion of the operation cost of large data centers. At the same time, the average server utilization in many Internet data centers is very low: real world estimates range from 5% to 20% [5], [6].

In this paper, we present a system that provides automatic scaling for Internet applications in the cloud environment. The present approach includes the following. The automatic scaling problem in the cloud environment, and model it as a modified Class Constrained Bin Packing (CCBP) problem where each server is a bin and each class represents an application. The present study develops an innovative auto scaling algorithm to solve the problem and presents a rigorous analysis on the quality of it with provable bounds. Experiments and simulations show that our algorithm is highly efficient and scalable which can achieve high demand satisfaction ratio, low placement change frequency, short request response time, and good energy saving.

The present paper builds a real cloud computing system which supports our auto scaling algorithm. The present method compares the performance of our proposed system with an open source implementation of the Amazon EC2 auto scaling system in a testbed of 30 Dell PowerEdge blade servers. Experiments show that the proposed system can restore the normal QoS five times as fast when a flash crowd happens.

The rest of the paper is organized as follows. Section 2 presents the architecture of the system and formulates the auto scaling problem. Section 3 describes the details of our proposed algorithm. Experiment results are presented in Sections 4. Section 7 gives the conclusions of the paper.

SYSTEM ARCHITECTURE

The architecture of our system is shown in Fig. 2. The proposed method encapsulates each application instance inside a virtual machine (VM). The use of VMs is necessary to provide isolation among untrusted users. Each server in the system

runs the Xen hypervisor which supports a privileged domain 0 and one or more domain U [7]. Each domain U encapsulates an application instance, which is connected to shared network storage (i.e., the storage tier). The multiplexing of VMs to PMs (Physical Machines) is managed using the Usher framework [8]. The main logic of our system is implemented as a set of plug-ins to Usher. Each node runs an Usher local node manager (LNM) on domain 0 which keeps track of this set of applications running on that node and the resource usage of each application.

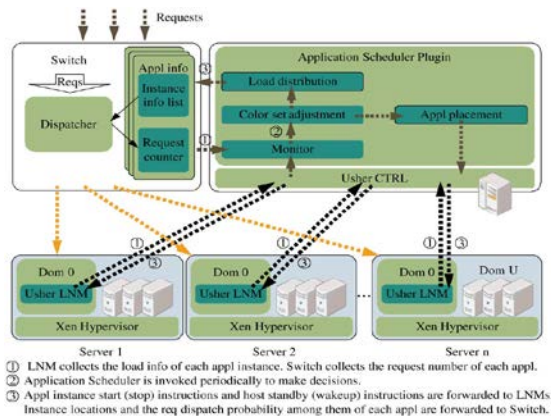


Figure 2: the system architecture.

The schedule procedure of our system can be described as follows.

1. The LNM at each node and the L7 switch collect the application placement, the resource usage of each instance, and the total request number of each application periodically.
2. The Application Scheduler is invoked periodically to make the following decisions:

- Application placement: for each application decide the set of servers its instances run on.

- load distribution: for each application, predict its future resource demands based on the request rate and past statistics, and then decide how to allocate its load among the set of running instances. The load of an Internet application is largely driven by the rate of user requests.

3. The decisions are forwarded to the LNM and the L7 switch for execution. The list of action items for each node includes:

- standby or wake up instructions

- application starts and stops

- the allocation of local resource among the applications

After that the Scheduler notifies the L7 switch of the new configuration including:

- the list of applications

- for each application, the location of its running instances and the probability of request distribution among them

The L7 switch then starts processing Web requests according to the new configuration. It may seem from the discussion above that the UsherCTRL is a central point of failure.

From above architecture notice that complicated applications can take along time (several minutes or much longer) to start and finish all the initializations. The present method takes the advantage of this feature to bypass the application start process by suspending a fully started and initialized application instance to the disk

PROPOSED METHOD

Automatic scaling: The auto scaling problem is defined as follows: Suppose, have a server set S on which need to run a set of applications (A). The CPU capacity of server s ($s \in S$) is C_s , the maximum number of application instances which can run on server simultaneously according to memory factor is M_s , and the CPU demand of application ($a \in A$) is C_a .

To simplify the problem described above, the present method assumes that the servers are homogeneous with uniform capacity. Then the auto scaling problem is similar to the Class Constrained Bin Packing (CCBP) problem when label each application as a class and treat the CPU demands of all classes as the items which need to be packed into bins. The only difference is that the CCBP problem does not have the “Minimize the placement change frequency” goal. Therefore, in order to solve our problem, the present method modified the CCBP model to support the “Minimize the placement change frequency” goal and provide a new enhanced semi online approximation algorithm to solve it. The proposed algorithm is given below.

Proposed Algorithm: the proposed algorithm belongs to the family of color set algorithms [13], but with significant modification to adapt to proposed problem. The proposed method labels each class of items with a color and organizes them into color sets as they arrive in the input sequence. The number of distinct colors in a color set is at most c (i.e., the maximum number of distinct classes in a bin). This ensures that items in a color set can always be packed into the same bin without violating the class constraint. The packing is still subject to the capacity constraint of the bin. All color sets contain exactly c colors except the last one which may contain fewer colors.

Items from different color sets are packed independently. A greedy algorithm is used to pack items within each color set: the items are packed into the current bin until the capacity is reached. Then the next bin is opened for packing. Thus each color set has at most one unfilled (i.e., non-full) bin. Note that a full bin may contain fewer than c colors. When a new item from a specific color set arrives, it is packed into the corresponding unfilled bin. If all bins of that color set are full, then a new bin is opened to accommodate the item.

Application Load Increase

The load increase of an application is modeled as the arrival of items with the corresponding color. A naive algorithm is to always pack the item into the unfilled bin if there is one. If the unfilled bin does not contain that color already, then a new color is added into the bin. This corresponds to the start of a new application instance which is an expensive operation. Instead, the proposed algorithm attempts to make room for the new item in a currently full bin by shifting some of its items into the unfilled bin. Let be the color of the new item and be any of the existing colors in the unfilled bin. The proposed method searches for a bin which contains items of both colors. Then the proposed method moves an item of color from bin to the unfilled bin. This makes room for an item in bin where we pack the new item.

The searching for bin process is

bin₁ contains colors c_1 and c_3

bin₂ contains colors c_2 and c_3

If proposed method can find such two bins, the following procedure follows:

- move an item of color c_2 from bin b_2 to the unfilled bin
- move an item of color c_3 from bin b_1 to bin b_2
- pack the item in bin b_1

This process is illustrated in Fig. 3

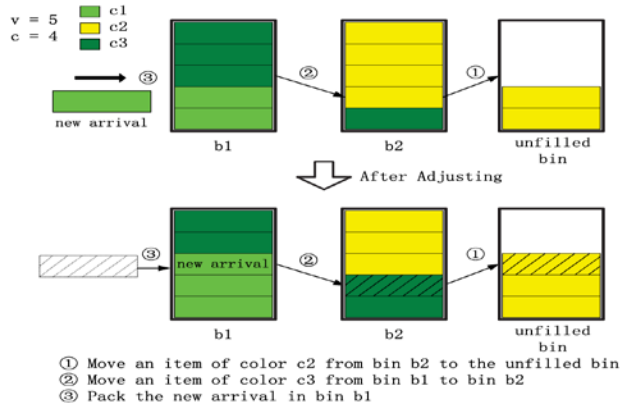


Figure 3: Schematic diagram of Arrival a new item

APPLICATION LOAD DECREASE

The load decrease of an application is modeled as the departure of previously packed items. Note that the departure event here is associated with a specific color, not with a specific item. The algorithm has the freedom to choose which item of that color to remove.

The proposed departure algorithm works as follows. If the color set does not have an unfilled bin, we can remove any item of that color and the resulting bin becomes the unfilled bin. Otherwise, if the unfilled bin contains the departing color, a corresponding item there can be removed directly. In all other cases, we need to remove an item from a currently full bin and then fill the hole with an item moved in from somewhere else. Let c_1 be the departing color and be any of the colors in the unfilled bin. The proposed method need to find a bin

which contains items of both colors. Let be such a bin. We remove the departing item from bin and then move in an item of color from the unfilled bin. The procedure is similar to the previous case for application load increase. Fig. 4 illustrates this process for a chain with three colors.

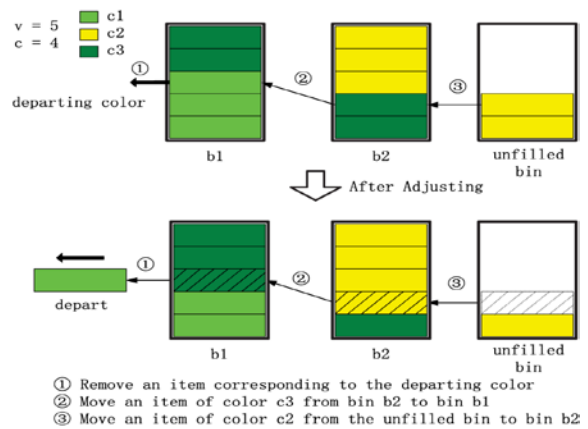


Figure 4: Schematic diagram of depicting an existing item

Analysis of the Approximation Ratio

The quality of a polynomial time algorithm A is measured by its approximation ratio $R(A)$ to the optimal algorithm OPT:

$$R(A) = \lim_{n \rightarrow \infty} \sup_{OPT(\sigma)=n} \frac{A(\sigma)}{OPT(\sigma)}$$

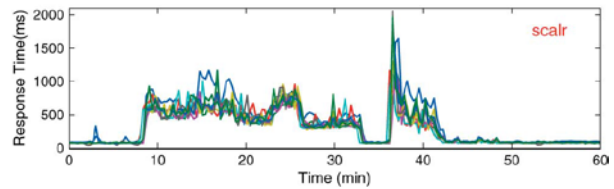
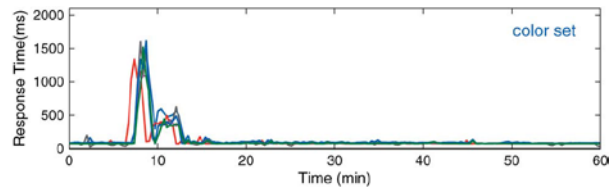
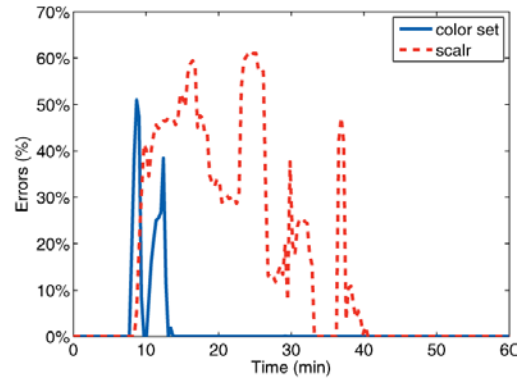
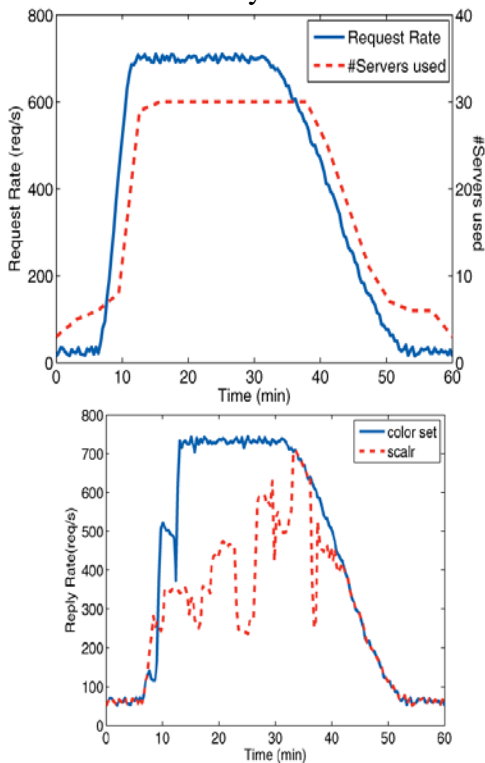
Where $A(\sigma)$ is the list bins used under the algorithm and the optimal algorithm, respectively [9, 10].

EXPERIMENTS

We in experiments.

To evaluate the effectiveness of our system, the Web applications used in the experiments are Apache servers serving CPU intensive PHP scripts. Each application instance is encapsulated in a

separate VM. The servers are connected over a Gigabit Ethernet. The client machines run `httperf` to invoke the PHP scripts on the Apache servers. This allows us to subject the applications to different degrees of CPU load by adjusting the client request rates. The proposed method considers a server as “full” when its capacity reaches 80%. This leaves some room for additional load increase. To save time on the experiments, the present approach configures the Application Scheduler with an aggressive two minutes interval between invocations. This allows us to complete the experiments in a timely manner.



The auto scaling capability of our algorithm with nine applications and 30 Dell PowerEdge servers with Intel E5620 CPU and 24 GB of RAM. The servers run Xen-4.0 and Linux 2.6.18. The results are shown in Figs.5, 6, 7 and 8. To increase the load of one application dramatically to emulate a “flash crowd” event while keeping the load of the other applications steady.

Fig. 5 shows the request rate of the flash crowd application and the number of active servers (i.e., APMs) used by all applications over the course of the experiment. Initially, the load in the system is low and only a small number of servers are used. When the flash crowd happens, the present approach detects the skyrocketing request rate quickly and scales up the server resources decisively. The above figures show that it uses up all 30 servers during the peak demand. Then the present method



reduces the request rate of the application gradually to emulate that the flash crowd is over. The algorithm scales down the server resources accordingly to conserve energy.

CONCLUSIONS:

In this paper proposed the design and implementation of a system that can scale up and down the number of application instances automatically based on demand. In this paper developed a color set algorithm to decide the application placement and the load distribution. The proposed system achieves high satisfaction ratio of application demand even when the load is very high. It saves energy by reducing the number of running instances when the load is low.

REFERENCES:

- [1] Amazon Elastic Compute Cloud (Amazon EC2). <http://aws.amazon.com/ec2/>. Accessed on May 10, 2012.
- [2] A. Cohen, S. Rangarajan, and H. Slye, "On the performance of tcp splicing for url-aware redirection," in Proc. 2nd Conf. USENIX Symp. Internet Technol. Syst., 1999, p. 11.
- [3] Oracle WebLogic Suite. <http://www.oracle.com/us/products/middleware/cloud-app-foundation/weblogic/overview/index.html>. Accessed on May 10, 2012.
- [4] Google App Engine. <http://code.google.com/appengine/>. Accessed on May 10, 2012.
- [5] M. Armbrust et al., "Above the clouds: A Berkeley view of cloud computing," EECS

Depart., Univ. California, Berkeley, CA, Tech. Rep. UCB/EECS-2009-28, Feb. 2009.

[6] L. Siegele, "Let it rise: A special report on corporate IT," in *The Economist*, London, U.K.: London Economist Newspaper, Oct. 2008, vol. 389, pp. 3-16.

[7] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, and A. Warfield, "Xen and the art of virtualization," in Proc. ACM Symp. Oper. Syst. Princ. (SOSP'03), Oct. 2003, pp. 164-177.

[8] M. McNett, D. Gupta, A. Vahdat, and G. M. Voelker, "Usher: A nextensible framework for managing clusters of virtual machines," in Proc. Large Install. Syst. Admin. Conf. (LISA'07), Nov. 2007, pp. 1-15.

[9] M. R. Garey and D. S. Johnson, "A 71/60 theorem for bin packing," *J. Complexity*, vol. 1, pp. 65-106, 1985.

[10] Scalr: The Auto Scaling Open Source Amazon EC2 Effort. <https://www>

[scalr.net/](https://www.scalr.net/). Accessed on May 10, 2012.