# TCC : A Novel Cloud Service for Cloud-deployed Applications

**[1]Banoth Seetha Ramulu, [2] H. Balaji**

[1]Associate Professor, Department of CSE, Vardhaman College of Engineering, Shamshabad, Hyderabad, TS, India
[2]Associate Professor, Department of CSE, Sreenidhi Institute of Science and Technology, Ghatkesar, Hyderabad, TS, India

## ABSTRACT

Applications with a dynamic workload demand need access to a flexible infrastructure to meet performance guarantees and minimize resource costs. While cloud computing provides the elasticity to scale the infrastructure on demand, cloud service providers lack control and visibility of user space applications, making it difficult to accurately scale the infrastructure. Thus, the burden of scaling falls on the user. That is, the user must determine when to trigger scaling and how much to scale. Scaling becomes even more challenging when applications exhibit dynamic changes in their behavior. In this paper, we propose a new cloud service, Trusty Compute Cloud (TCC), which spontaneously scales the infrastructure to meet the user-specified performance requirements, even when multiple user requests execute concurrently.

**Keywords:** Cloud Computing, Trusty Compute Cloud (TCC). SLA.

## I. INTRODUCTION

### Motivation

With the coming of cloud processing, numerous application owners have begun moving their arrangements to the cloud. Cloud figuring offers numerous advantages over customary physical usage including lower foundation expenses and versatile asset portion. These advantages are particularly invaluable for applications with a dynamic workload request. Such applications can be sent in the cloud based on the present request, and the organization can be scaled powerfully because of changing workload request. Unfortunately, , it is difficult to completely understand the capability of cloud figuring. While Cloud Service Providers (CSPs, for example, Amazon [3], give clients simple access to cloud assets for their processing needs, they don't offer any assurances on the execution of a client's organization or any rules on how clients should set their asset portions. Thus, clients either fall back on inefficient practices, for example, overprovisioning or forsake the cloud by and large and return to top provisioned physical arrangements. This is clear by the poor cloud selection for execution touchy applications [14].Given these perceptions, we declare that giving execution assurances to cloud clients will significantly enhance cloud use and advance effective cloud usage.

### Problem Statement and Goal

Giving execution certifications to cloud clients is troublesome on the grounds that client arrangements are dark: CSPs can't control or access a client's workload or application. Further, CSPs won't not know the client application because of protection concerns. Given these confinements, the best that CSPs can do is to give straightforward, lead-based answers for overseeing client applications. These administer based arrangements enable the clients to indicate a few conditions on the checked measurements which, when met, will trigger a pre-characterized scaling activity. Indeed, even with the assistance of lead-based arrangements, notwithstanding, the weight still rests with the client. For instance, keeping in mind the end goal to utilize a CPU use based trigger for scaling, the client must

decide the CPU edges at which to trigger scale-up and scaledown, and the quantity of occurrences to scale. To exacerbate the situation, the ideal limit esteems ordinarily rely upon (dynamic) workload attributes, for example, entry rate and workload blend, in this way requiring steady tuning and manual intercession. Given these entanglements, it isn't amazing that client organizations are tormented with execution issues [6].

Note that clients can oversee and scale their applications in a cloud situation. Nonetheless, this is a testing undertaking since it: (I) requires master information about the progression of all the included programming, including the service necessities of the application at every level, and (ii) requires advanced displaying skill to decide when and how to resize the organization. These obstacles are not an issue for enormous organizations that have enough assets to utilize a group of specialists for managing these issues. Be that as it may, for little and medium undertakings (which contain the focused on client base for some CSPs [5]), and for the easygoing cloud client, these obstacles are non-inconsequential to overcome. Such clients would much rather get a cloud service that deals with their application.

and how to resize the deployment. These hurdles are not a problem for big businesses that have enough resources to employ a team of experts for dealing with these issues. However, for small and medium businesses (which comprise the targeted customer base for many CSPs [5]), and for the casual cloud user, these hurdles are non-trivial to overcome. Such users would much rather contract a cloud service that manages their application.
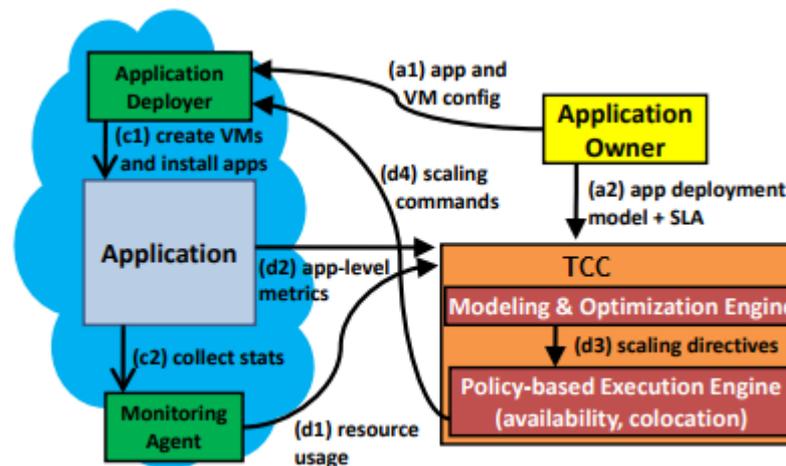


Figure 1. System architecture TCC

The objective of this paper is to give a cloud service to clients that naturally scale their murky cloud applications in light of changing workload and cloud conditions without falling back on intrusive and frequently infeasible (in a cloud domain) approaches, for example, disconnected profiling and benchmarking.

Existing Solutions

Many CSPs today, including Amazon [1] and RightScale, offer rule-based solutions (not necessarily for free) to users for dynamically managing their deployments. These solutions are typically meant to be used with CSP-provided monitoring solutions such as Amazon's Cloud Watch [2] and Rackspace's Cloud Monitoring. Such rule-based solutions are also offered by cloud software solutions such as OpenStack. Even with the help of such services, however, the user still requires expert knowledge about the application and the performance modeling expertise to convert the monitored information into scaling rules. This is a non-trivial task and requires extensive testing. There is also a lot of prior research work on dynamically scaling user applications in order to provide performance guarantees. Unfortunately, almost all of these works are infeasible for opaque cloud applications since they require access to the user deployment for instrumentation, benchmarking [9], or assume that expert application knowledge (such as per-tier service times) is available a priori [8], [11].

## II. OUR APPROACH

We propose a computerized cloud service, Trusty Compute Cloud (TCC), that proactively and progressively scales the application organization based on client determined execution prerequisites without requiring any extra instrumentation, benchmarking or master application learning. TCC use accessible asset level and application-level measurements to induce the fundamental framework parameters of the application(s) and decides the required scaling activities to meet the execution objectives in a financially savvy way. These scaling mandates would then be able to be passed on to the cloud figuring programming, for instance, OpenStack [7], to execute the scaling.

Note that the scaling mandates can likewise be passed on to an approach based execution motor to guarantee colocation and high-accessibility limitations if necessary. A nitty-gritty talk of our approach is displayed in Section 2. At the core of TCC lies the displaying and execution motor that disguises the observed insights and gathers the essential framework parameters. While this motor can utilize any dim box or discovery demonstrating approach, in this paper we utilize Kalman sifting to comprehend the framework parameters.

## 1. ARCHITECTURE

We now give a structural perspective of our TCC arrangement. Figure 1 demonstrates the proposed framework design for the TCC service condition. The Application is facilitated in the (blue) cloud delineated in the inside. The Application Owner (client or client) is in charge of: (errand (a1)) giving the underlying application organization demonstrate and the Virtual Machine (VM) designs to the cloud with the goal that the application can be propelled, and (undertaking (a2)) giving a similar application arrangement display and the execution SLA necessities to our TCC service. The sending model contains data on the quantity of VMs and their system associations, as a diagram or a setup record; structures, for example, OpenStack Heat or Weaver [15]

regularly utilize such arrangement demonstrate data as a feature of their information documents.

The Application Deployer and Monitoring Agent are services that are ordinarily given by the CSP and are in this manner appeared as a major aspect of the cloud. The Application Deployer, for example, AWS Elastic Beanstalk [4] or OpenStack Heat, tweak the picture and VM for an organization and ties up the endpoints for the application amid establishment and setup (errand (c1)). When utilizing these Application Deployers, TCC can specifically acquire the application sending model from them without requiring the client to give these subtle elements. The Monitoring Agent, for example, AWS CloudWatch [2] or OpenStack Ceilometer, tracks and stores asset use measurements, for example, CPU and plate usage, of the VMs (assignment (c2)).

The TCC part gathers asset use insights from the Monitoring Agent (undertaking (d1)). It likewise oversees application-level figures, for example, asks for rate, from the application (undertaking (d2)). Application-level observing is given, however to a restricted degree, by some CSPs, for example, Amazon (when utilizing their heap balancer [4]). These insights are then encouraged to the Modeling and Optimization Engine, which models the hidden application based on the client gave organization demonstrate and the deliberate measurements, for example, CPU usage and demand rate. It likewise construes undetectable framework parameters, for example, per-level service necessities and foundation CPU usage, based on the model and estimated insights. Utilizing the model and the deliberate and construed parameters, the Modeling and Optimization Engine decides the scaling orders, for example, VM scale up/down, required for keeping up the client gave execution SLA (errand (d3)). These orders are passed on to the Policy-based Execution Engine that issues orders to the hidden cloud API (assignment (d4)), that thus plays out the scaling tasks. The Policy-based Execution Engine can likewise decide the arrangement of VMs on the genuine Physical Machines (PMs) based on accessibility, security, or colocation requirements..

## 2. MODELING

The displaying motor lies at the core of our TCC approach. We utilize a queueing-organize model to inexact our multi-level cloud application. Nonetheless, since we can't get to the client application to determine the parameters of our model, we utilize a Kalman separating method to derive these inconspicuous parameters. Further, by utilizing the present observing data by means of the checking specialist, we refine our model to adjust to any adjustments in the framework progressively. Critically, by utilizing the Kalman channel to use the genuine checked esteems, we limit our reliance on the estimated queueing model of the framework. We now portray our queueing model and Kalman sifting method, trailed by an investigation of our demonstrating motor, lastly, a clarification of how our displaying motor decides the required scaling activities for SLA consistence.
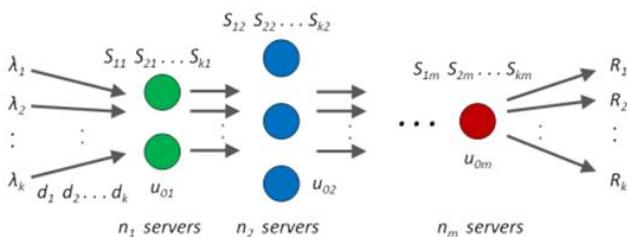


Fig. 2. Our queueing model for an $m$-tier system with $k$ request classes. The system parameters are: $\lambda_i$, arrival rate of class $i$; $R_i$, response time for class $i$; $d_i$, network latency for class $i$; $u_{0j}$, background utilization for tier $j$; $S_{ij}$: service time of class $i$ at tier $j$.

### 2.1 Queuing-network model

Figure 2 shows our queueing-network model for a generic m-tier system with each tier representing a collection of homogeneous servers. We assume that the load at each tier is distributed uniformly across all the servers in that tier. The system is driven by a workload consisting of k distinct request classes, with each class, say class i, characterized by its arrival rate, $\lambda i$ , and end-to-end response time, Ri . Let nj be the number of servers at tier j. With homogeneous servers and perfect load-balancing, the arrival rate of requests at any server in tier j is $\lambda ij := \lambda i/nj$ . Since servers at a tier are identical, for ease of analysis, we model each tier as a single representative server.

With some abuse of terminology, we refer to the representative server at tier j as tier j. Let uj $\in$ [0, 1) be

the utilization of tier j. The background utilization of tier j is denoted by u0j , and models the resource utilization due to other jobs (not related to our workload) running on that tier and the virtualization overhead due to multi-tenancy, if any. We believe that u0j can also account for resource interference in highly contended cloud environments; we will investigate models for interference ridden environments as part of future work. The end -to-end network latency for a class i request is denoted by di . Let Sij ($\geq$ 0) denote the average service time of a class i request at tier j.

Assuming we have Poisson arrivals and a processor-sharing policy at each server, the stationary distribution of the queueing network is known to have a product-form, for any general distribution of service time at servers. Under the product-form assumption, we have the following analytical results from queueing theory:

$$u_j = u_{0j} + \sum_i \lambda_{ij} S_{ij}, \quad \forall j \tag{1}$$

$$R_i = d_i + \sum_j \frac{S_{ij}}{1 - u_j}, \quad \forall i \tag{2}$$

While uj , Ri and $\lambda i$ , $\forall i$, j, can be monitored easily and are thus observable, the parameters Sij , u0j , and di are non-trivial to measure and are thus unobservable. While existing work on auto-scaling typically obtains these values by directly accessing or modifying the application software (for example, by parsing the log files at each tier), our proposed applicationagnostic cloud service cannot encroach the user's application space. Instead, we employ a parameter estimation technique, Kalman filtering, to derive estimates for the unobservable parameters.

It is important to note that while the product-form is shown to be a reasonable assumption for tiered web services [7], we only use it as an approximation for our complex system. Also we approximate a multicore server as a single core with scaled-up capacity. Since we are interested in horizontal scaling, we do not need to explicitly model the scaling of service time

with cores as in prior work [10]. By employing the Kalman filter to leverage the actual monitored values, we minimize our dependence on these approximations.

## 2.2 Scaling directives

The estimated values of the system state are used to compute the required scaling actions for TCC. Specifically, given the response time SLA, we use Eqns. (1) and (2) to determine the minimum number of servers in each tier, $n_j$, $\forall j \in \{1, 2, \ldots, m\}$, required to ensure SLA compliance. In particular, recalling that $\lambda_{ij} = \lambda_i/n_j$, and substituting $u_j$ in Eqn. (2) using Eqn. (1), we have:

$$R_i = d_i + \sum_j \frac{S_{ij}}{1 - (u_{0j} + \sum_k \frac{\lambda_k}{n_j} S_{kj})}, \quad \forall i$$

where the summation over k represents multiple request classes. Eqn. (3) can now be solved for $n_j$ given the target response time(s). As a simple example, assume that we are concerned about the response time of one class, say class 1, and we are only concerned about scaling one tier, say tier 1. Then, for a given response time SLA for class 1, RSLA, we can determine the number of tier 1 VMs needed, $n_1$, as follows:

$$R_{SLA} > R_1 = d_1 + \sum_j \frac{S_{1j}}{1 - (u_{0j} + \sum_k \frac{\lambda_k}{n_j} S_{kj})}$$

$$\Rightarrow n_1 = \left\lceil \frac{\sum_k \lambda_k \cdot S_{k1}}{1 - u_{01} - \frac{S_{11}}{R_{SLA} - d_1 - \sum_{j>1} \frac{S_{1j}}{1 - (u_{0j} + \sum_k \frac{\lambda_k}{n_j} S_{kj})}}} \right\rceil \quad (4)$$

## 2.3 Rule-based approaches

Auto-scaling and load-balancing features are now being offered by almost every major CSP including Amazon Web Services (AWS) [1], and Google Cloud Platform [13]. However, to the best of our knowledge, all the existing CSP-offered autoscaling solutions are rule-based and typically require the user to specify the threshold values on the resource usage (e.g., CPU, memory, storage) for triggering scaling actions. While rule-based solutions are suitable for the cloud environment where the user application cannot be accessed, they ultimately place the burden of the auto-scaling logic on the user. Further, such rule-based approaches have to be tuned to the specific demand pattern and workload for best results, as demonstrated by the THRES policy. By contrast, TCC does not require the user to specify scaling rules. TCC automatically determines the required scaling actions and executes them in a timely manner to ensure SLA compliance. The authors in [12] use fuzzy logic to deduce threshold values for rulebased triggers. While this approach only uses online profiling, it does not leverage a queueing-theoretic system model to improve accuracy and convergence

## III. CONCLUSIONS

In this paper, we show the plan and execution of another cloud service, Trusty Compute Cloud (TCC), that consequently scales client applications in a savvy way to give execution ensures. Since CSPs don't have finish control and permeability of a client's cloud organization, we outlined TCC to be application-sceptic. Specifically, not at all like a significant portion of the current auto-scaling research, TCC does not require any disconnected profiling or benchmarking of the application nor does it require a profound comprehension of the application elements. Instead, TCC utilizes a Kalman separating procedure in a blend with a queueing theoretic model to proactively decide the correct scaling activities for an application conveyed in the cloud utilizing effortlessly accessible measurements, for example, use and demand rate. We executed TCC as a service on OpenStack and exhibited its capacity to guarantee application SLA consistence by powerfully scaling virtual occasions and hypervisors. Our test comes to feature the vigour of TCC to changes in the requested design and to changes in the workload blend. As a feature of future work, we will explore coordinating vertical scaling choices with TCC to give all the more capable scaling alternatives. We will likewise more altogether investigate the mix of TCC with other autoscaling strategies, including prescient models.

## IV. REFERENCE

[1]. Amazon Auto Scaling. http://aws.amazon.com/autoscaling.

[2]. Amazon CloudWatch. http://aws.amazon.com/cloudwatch.

[3]. Amazon EC2. http://aws.amazon.com/ec2.

[4]. Elastic Beanstalk. http://aws.amazon.com/elasticbeanstalk.

[5]. Gartner's Advice for CSPs Becoming Cloud Service Providers. https://www.gartner.com/doc/2155315, 2012.

[6]. Sean Kenneth Barker and Prashant Shenoy. Empirical Evaluation of Latency-sensitive Application Performance in the Cloud. In Proceedings of the 1st Annual Conference on Multimedia Systems, pages 35–46, Phoenix, AZ, USA, 2010.

[7]. P. Dube, H. Yu, L. Zhang, and J.E. Moreira. Performance evaluation of a commercial application, trade, in scale-out environments. In Proceedings of the 15th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems, pages 252–259, 2007.

[8]. A. Gandhi, Y. Chen, D. Gmach, M. Arlitt, and M. Marwah. Minimizing Data Center SLA Violations and Power Consumption via Hybrid Resource Provisioning. In Proceedings of the 2011 International Green Computing Conference, pages 49–56, Orlando, FL, USA, 2011.

[9]. A. Gandhi, M. Harchol-Balter, R. Raghunathan, and M. Kozuch. AutoScale: Dynamic, Robust Capacity Management for Multi-Tier Data Centers. Transactions on Computer Systems, 30, 2012.

[10]. Anshul Gandhi, Parijat Dube, Alexei Karve, Andrzej Kochut, and Li Zhang. Modeling the Impact of Workload on Cloud Resource Scaling. In Proceedings of the 26th International Symposium on Computer Architecture and High Performance Computing, SBAC-PAD '14, Paris, France, 2014.

[11]. Hamoun Ghanbari, Bradley Simmons, Marin Litoiu, Cornel Barna, and Gabriel Iszlai. Optimal Autoscaling in a IaaS Cloud. In Proceedings of the 9th International Conference on Autonomic Computing, pages 173–178, San Jose, CA, USA, 2012.

[12]. D. Gmach, S. Krompass, A. Scholz, M. Wimmer, and A. Kemper. Adaptive quality of service management for enterprise services. ACM Transactions on the Web, 2:1–46, 2008.

[13]. Google Cloud Platform. Auto Scaling on the Google Cloud Platform. http://cloud.google.com/resources/articles/autoscaling-on-the-google-cloud-platform.

[14]. Internap. Internap Public Cloud Survey Reveals Performance as Top Challenge for Cloud-Wise Organizations. http://www.internap.com/press-release/internap-publiccloud-survey-reveals-performance-top-challenge-cloud-wiseorganizations.

[15]. M. Kalantar, E. Snible, F. Rosenberg, T. Roth, T. Eilam, A. Oliveira, M. Elder, and J. Doran. Weaver: Language and Runtime for Software Defined Environments. IBM Journal of Research and Development (Accepted for publication), 2014.