

# Measuring and Testing the Scalability of Cloud-based Software Services

Amro Al-Said Ahmad  
School of Computing and Mathematics  
Keele University, UK  
a.m.k.al-said.ahmad@keele.ac.uk

Peter Andras  
School of Computing and Mathematics  
Keele University, UK  
p.andras@keele.ac.uk

**Abstract**—Performance and scalability testing and measurements of cloud-based software services are critically important in the context of rapid growth of cloud computing and supporting the delivery of these services. Cloud-based software services performance aspects are interrelated, both elasticity and efficiency are depending on the delivery of a sufficient level of scalability performance. In this work, we focused on testing and measuring the scalability of cloud-based software services in technical terms. This paper uses technical scalability metrics that address both volume and quality scaling, that inspired by earlier technical metrics of elasticity. We show how our technical scalability metrics can be integrated into an earlier utility oriented metric of scalability. We demonstrate the application of the metrics using a practical example and discuss the importance of them.

**Keywords**— Measurement, Performance, Testing, Scalability, Software-as-a-Service (SaaS), Metrics

## I. INTRODUCTION

In any software system, scalability and performance assessments provide an important basis for future optimizations, and for developing new opportunities aimed to maximize scalability and performance [1]. The performance assessment and testing of cloud-based software services is critically important in order to support the Service Level Agreement (SLA) compliant quality of delivery of these services, especially in the context of rapidly expanding the quantity of service delivery [2]. There are three cloud-specific performance aspects that are key determinants of service quality delivery in the context of variable service demand: scalability, elasticity and efficiency [3, 4].

Following [5] we adopt the following definitions of these three performance aspects. Scalability is the ability of the cloud layer to increase the capacity of the software service delivery by expanding the quantity of the software service that is provided. Elasticity is the level of autonomous adaptation provided by the cloud layer in response to variable demand for the software service. Efficiency is the measure of matching the quantity of software service available for delivery with the quantity of demand for the software service. These definitions focus on the technical side of cloud-based software services, however we note that alternative, utility oriented (i.e. economic cost/benefit focused), approaches are also used in the literature [6, 7].

Cloud-based applications should be scalable, and with auto-scaling and load-balancing features such applications should be able to deal with sudden workload by adding more of the application instance(s). Furthermore, as cloud-based applications been offered as Software as a Services (SaaS), and the use of multi-tenancy architectures [8], emphasizes the need for scalability that supports the availability and productivity of the services and on-demand resources.

Recently a series of papers has been published addressing the topic of measuring elasticity of cloud-based provision of software services [9, 10]. There have been also works on scalability of cloud-based software services from the utility perspective [6, 7, 9, 11]. However, relevant recent systematic reviews report only a very small number of works (mainly in the grey literature, e.g. project reports, MSc theses) which try to address the assessment of scalability of cloud-based software services from the technical perspective [5].

Measuring and testing scalability of cloud-based software services from a technical perspective are key for the assessment and testing of performance [1,12]. Both elasticity and efficiency performance depend on the delivery of a sufficient level of scalability performance. Understanding how components of the cloud-based software service system contribute to the scalability performance of the system helps in designing appropriate test scenarios and identifying options for changes and upgrades that can improve the scalability performance of the system.

Utility oriented assessment of scalability [6] i.e. measuring the scalability from an economic and cost perspective, is insufficient for the above purpose, since it measures scalability from a perspective that is abstract relative to the technical components and features of the system. Thus it becomes very difficult and possibly even practically impossible to associate specific technical components and features with specific impact on the utility oriented scalability performance. This is due to the potential multiple impacts of such technical components and features on utility features of the system that get integrated into the utility oriented scalability measurement of the system.

Here we follow ideas proposed in the context of measurements and metrics for cloud elasticity [13–15] to propose technical measurement and metrics for scalability of cloud-based software services. To sum up, the main contributions of this paper:

- The work uses metrics [16] that address both volume and quality scaling for evaluating cloud-based software services scalability performance. This work provides an extension to the previous work [16] which offers explanation in further detail, introduce the demand scenarios, and demonstrate a practical example of the metrics.
- The metrics can be useful in order to support effective measurement and testing of scalability performance of those services from technical perspective.
- The paper proposed how those technical metrics can be integrated with earlier utility oriented scalability metrics proposed by [11].

- We demonstrate the application of the metrics to a concrete cloud-based software service (OrangeHRM) run through the Amazon EC2 Cloud.
- We show how the metrics can be used to identify differences in the behaviour of the assessed system in the context of different usage scenarios.

The rest of the paper is structured as follows. First we review briefly the relevant recent literature. Next we present our approach to measure and quantify scalability of cloud-based software services and explain the metrics based on the measurement approach. Next we present an application example using two different usage scenarios to demonstrate the measurement approach and metrics. Next we discuss the implications and importance of the approach and metrics. Finally, the paper is closed by the conclusions section.

## II. RELATED WORKS

A recent review [17] on provisioning of cloud resources and related research challenges, identify, among others predictable performance and scalable resource management as promising challenges. Gao et al. [18] reviewed testing in relation with cloud-based software services. They highlight scalability and performance testing as key research directions. Other similar recent surveys [19, 20] focus primarily on cloud service elasticity. The systematic literature review by Lehrig et al. [5] provides very useful definitions of key cloud performance concepts such as capacity, scalability, elasticity and efficiency, which we adopt in this paper.

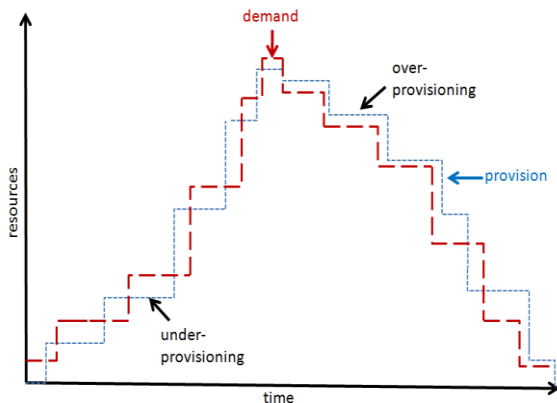


Fig. 1. Key concepts for measuring elasticity

There is a considerable number of recent papers that address the issue of measuring elasticity of cloud services in technical terms [4, 9, 15, 21–25]. Herbst et al. [4] define a useful set of key concepts that allow technical measurement of cloud service elasticity (see Fig. 1) such as the quantity and time extents for periods of time when the service provision is either below or above what is required by the service demand. They [4, 21] define as elasticity measures: the time shares and average time lengths in under-provisioned and over-provisioned states; the amounts of excess (over-provisioned) and lacking (under-provisioned) resources per time unit; the averages of the excess and lacking resources; and the jitter, which is the number of resource adaptations during a given time period of provisioning the service. The up-elasticity metric is defined as the reciprocal value of the product of the average under-provisioned time length and average lack of resource. The down-elasticity is defined similarly. Further elaboration on

these metrics is provided by [22], who introduced further components and ways of considering the above factors (e.g. scalability, functions of resource inaccuracy and reconfiguration time).

In terms of measuring and quantifying scalability we note the work of Hwang et al. [7, 11], which uses a utility oriented definition of scalability. Their production-driven scalability measure includes the consideration of a quality-of-service measure and the cost of service, in addition to a more technically oriented performance metric [7, 11]. While this approach is likely to be useful from the perspective of utility, because of its reliance of multiple facets of the system (including cost measures), it is unlikely to be able to provide sufficiently specific and useful information in terms of contribution of system components to system scalability in a technical sense. Thus the usefulness of the utility oriented scalability metric is limited in the context of testing and technical improvement of the cloud-based provision of the software service.

Attempts to provide a more technically oriented measure or metric for cloud-based software service scalability are also limited. For example, Herbst et al. [4] provide a technical scalability metric, however, this is a rather elasticity driven metric (sum of over- and under-provisioned resources over the total length of time of service provision). Jamal et al. [26] describe practical measurements of throughput in systems with and without multiple virtual machines, without clearly formulating a specific measure or metric of scalability. Similarly, Jayasinghe et al. [13, 14] provide a technical scalability measure practically the system scalability in terms of throughput and CPU utilization of a set of virtual machine system settings, but does not provide a generic metric or measure. Gao et al. [15] evaluates SaaS performance and scalability from the system capacity perspective, using the system load and capacity as measurements for scalability. Another recent work [27] focuses on building a model to help measuring and comparing different deployment configurations in terms of capacity, elasticity and costs.

## III. SCALABILITY PERFORMANCE MEASUREMENT

While this as noted in the Introduction, we define scalability as the ability of the cloud-based system to increase the capacity of the software service delivery by expanding the quantity of the software service that is provided when such increase is required by increased demand for the service [5]. We are not concerned with the short-term flexible provision of the resources, which we term elasticity of the service provision [21]. Our focus is whether the system can expand the quantity of the service when this expansion is required by demand over a sustained period of service provision.

In principle the increase of capacity could happen either by increasing the volume of service requests served by a single instance of the service provision software or by deploying multiple software instances, or by a combination of these two approaches. In general, we expect that if a service scales up ideally then the increase in demand for service should be matched by proportional increase in the provision of the service such that the quality of the service does not change. Here quality of the service may be seen for example in terms of average response time. This ideal scaling behaviour of the system should be valid over a sufficiently long time scale, i.e. short-term mismatches between provision and demand, which are the subject of

elasticity, are not relevant from the perspective of scalability. If the system does not scale according to the ideal manner, it recruits insufficient resources to deliver the increased volume of service without change in the quality of the service. In general, real systems are expected to operate below the level of the ideal scaling behaviour and the aim of measuring scalability is to quantify the extent to which the real system behaviour differs from the ideal behaviour.

To deliver the ideal scaling, we expect that the system increases the number of instances of the software proportionally with the increase in demand for software services, i.e. if the demand increases by 50% we would ideally expect the base number of software instances to increase by 50%. We expect also that the system maintains the quality of service in terms of maintaining the same average response time irrespective of the volume of service requests, i.e. an increase of 50% of demand we would ideally expect no increase in average response time. Formally, let us assume that  $D$  and  $D'$  are two service demand volumes,  $D' > D$ . Let  $I$  and  $I'$  be the corresponding number of software instances that are deployed to deliver the service, and let  $t_r$  and  $t'_r$  be the corresponding average response times. If the system scales ideally we expect that for any levels of service demand  $D$  and  $D'$ .

$$D' / D = I' / I \quad (1)$$

$$t_r = t'_r \quad (2)$$

Equation (1) expresses that the volume of software instances providing the service scales up with the demand for the service. Equation (2) expresses that the quality of service, in terms of average response time, remains unchanged for any level of service demand.

To measure the values of  $I$  and  $t_r$  the system must perform the delivery of the service over some sustained time, such that short-term variations, due to elastic response of the system, do not influence the system measurements. In practice this means that the number of software instances and the average response time should be calculated by averaging over a number of measurements during the execution of a demand scenario (e.g. every second), and a number of repeated applications of the same demand scenario, i.e. a pattern of demand presentation, which may include variation in the demand.

Demand scenarios may follow certain patterns expected to test the scalability of the system in specific ways. Two kinds of demand patterns that appear as natural and typical choices are the steady increase followed by steady decrease of the demand with a set level of the peak, and the stepped increase and decrease, again with a set peak level of demand. These two demand scenarios are shown in Fig.2. Other demand scenarios may reflect prior knowledge about the system or its service market (e.g. the market may be characterized by spikes of demand or by seasonal variation of the demand). Any demand scenario has to be characterized by a summary measure of the demand level, which may be the peak level or the average or total demand level. This characteristic demand of a demand scenario is represented by  $D$ .

Naturally, real world cloud systems are unlikely to deliver the ideal scaling behaviour. The difference between

the ideal and the actual scaling behaviour of the system offers the possibility of defining technical scalability metrics for cloud-based software services.

In terms of provision of software instances for the delivery of the services the scaling is deficient if the number of instances is lower than the ideally expected number of software instances.

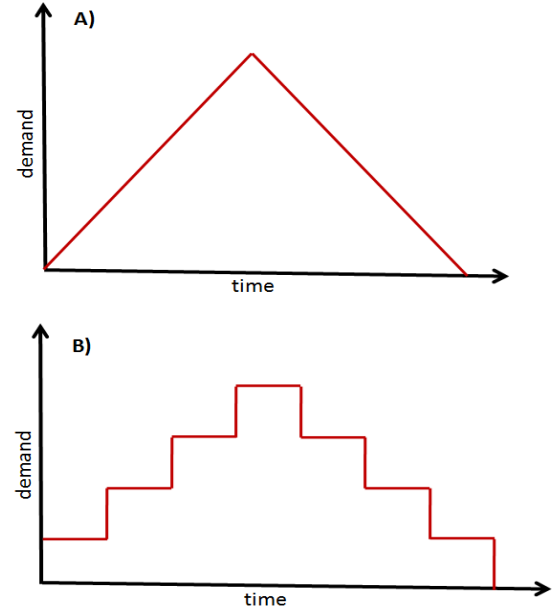


Fig. 2. Demand scenarios: A) steady rise and fall of demand; B) stepped rise and fall of demand.

To quantify the level of deficiency we pick a demand scenario and start with a low level of demand  $D_0$  and measure the corresponding volume of software instances  $I_0$ . Then measuring the number of software instances  $I_k$  corresponding to a number ( $n$ ) of demand levels  $D_k$  following the same demand scenario, we can calculate how close are the  $I_k$  values to the ideal  $I_k^*$  values ( $I_k < I_k^*$ ). Following the ideal scalability assumption of equation (1) we get for the ideal  $I_k^*$  values:

$$I_k^* = (D_k / D_0) \cdot I_0 \quad (3)$$

Considering the ratio between the area defined by the  $(D_k, I_k)$  values,  $k = 0, \dots, n$ , and the area defined by the  $(D_k, I_k^*)$  values we get a metric of service volume scalability of the system:

$$A^* = \sum_{k=1, \dots, n} (D_k - D_{k-1}) \cdot (I_k^* + I_{k-1}^*) / 2 \quad (4)$$

$$A = \sum_{k=1, \dots, n} (D_k - D_{k-1}) \cdot (I_k + I_{k-1}) / 2 \quad (5)$$

$$\eta_l = A / A^* \quad (6)$$

where  $A$  and  $A^*$  are the areas under the curves calculated for actual and ideal  $I$  values and  $\eta_l$  is the volume scalability performance metric of the system. If  $\eta_l$  is close to 1 the system is close to ideal volume scalability, if it is close to 0,

then the volume scalability of the system is much less than ideal.

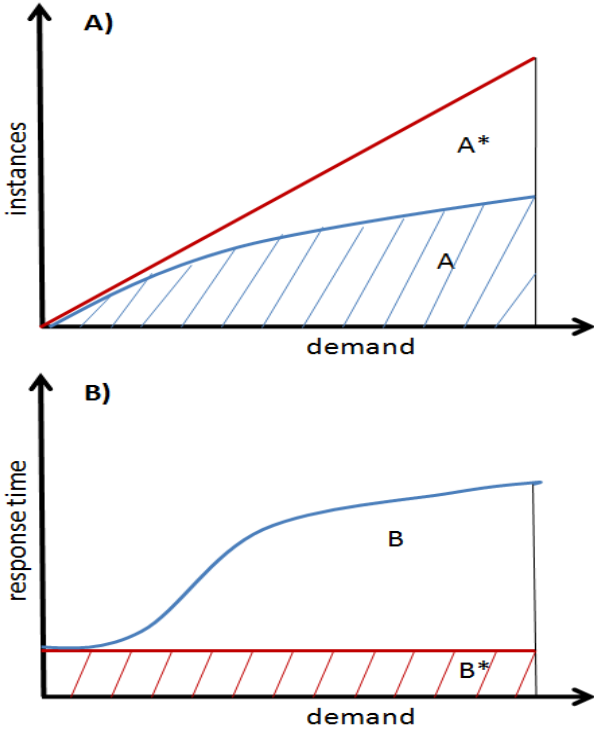
Similarly, we can define the quality scalability of the system by measuring the service average response times  $t_k$  corresponding to the demand levels  $D_k$ . We approximating the ideal average response time as  $t_0$ , following the ideal assumption of equation (2). The quality scalability of the system is less than ideal if the average response times for increasing demand levels increase, i.e.  $t_k > t_0$ . By considering the ratio between the areas defined by the  $(D_k, t_k)$  values,  $k = 0, \dots, n$ , and the area defined by the  $(D_k, t_0)$  values we get a ratio that defines a metric of service quality scalability for the system:

$$B^* = \sum_{k=1, \dots, n} (D_k - D_{k-1}) \cdot t_0 = (D_n - D_0) \cdot t_0 \quad (7)$$

$$B = \sum_{k=1, \dots, n} (D_k - D_{k-1}) \cdot (t_k + t_{k-1}) / 2 \quad (8)$$

$$\eta_t = B^* / B \quad (9)$$

where  $B$  and  $B^*$  are the areas under the curves calculated for actual and ideal  $t$  values and  $\eta_t$  is the quality scalability performance metric of the system. If  $\eta_t$  close to 1 the system is close to ideal quality scalability, if it is close to 0 the quality scalability of the system is much less than ideal.



**Fig. 3.** The calculation of the scalability performance metrics: A) the volume scalability metric is  $\eta_v$ , which is the ratio between the areas  $A$  and  $A^*$  – see equation (6); B) the quality scalability metric is  $\eta_t$ , which is the ratio between the areas  $B^*$  and  $B$  – see equation (9). The red lines indicate the ideal scaling behaviour and the blue curves show the actual scaling behaviour.

The calculation of the two scalability performance metrics is illustrated in Fig. 3. In Fig. 3A,  $A^*$  is the area under the red line showing the ideal expectation about the scaling behaviour (see equation (1)) and  $A$  is the shaded area

under the blue curve. The blue curve is under the ideal red line, indicating that the volume scaling is less efficient than the ideal scaling. In Fig. 3B,  $B^*$  is the shaded area under the red line indicating the expected ideal behaviour (see equation (2)) and  $B$  is the area under the blue curve. The blue curve is above the ideal red line, indicating that the quality scaling is less than ideal. We chose nonlinear curves for the examples of actual scaling behaviour to indicate that the practical scaling of the system is likely to respond in a nonlinear manner to changing demand.

These scalability metrics allow the effective measurement of technical scalability of cloud-based software services. These metrics do not depend on other utility considerations (e.g. price of service, non-technical quality aspects), which makes them appropriate for testing the technical scalability of the system. This makes possible the use of these metrics in scalability tests that aim to identify parts of the system that have significant impact on the technical scalability, and also the testing of the impact of any change made to the system on the technical scalability of the system.

Applying the scalability metrics to different demand patterns allows the testing and tuning of the system for particular usage scenarios and the understanding of how system performance can be expected to change as the pattern of demand varies. Such application of these metrics may highlight trade-offs between volume scaling and quality scaling of the system that characterize certain kinds of demand pattern variation (e.g. the impact of transition from low frequency peak demands to high frequency peak demands or to seasonal change of the demand). Understanding such trade-offs can help in tailoring the system to its expected or actual usage.

The scalability metrics can be integrated into the utility oriented scalability metric proposed by Hwang et al. [9], by considering a combination of our metrics as the performance and/or quality components of the utility oriented scalability metric. In [9] the utility oriented scalability of the system is defined as the ratio of two utility oriented productivity metric values associated with two different configurations of the system (i.e. one configuration is a scaled-up version of the other). The utility oriented productivity metric ( $P(\Lambda)$ ) is given as [9]:

$$P(\Lambda) = p(\Lambda) \cdot \omega(\Lambda) / c(\Lambda) \quad (10)$$

Where  $\Lambda$  is the system configuration,  $p(\Lambda)$  is the performance component of the metric,  $\omega(\Lambda)$  is the quality component of the metric and  $c(\Lambda)$  is the cost component of the metric. A natural way of integrating our technical metrics into this utility oriented framework is to use our volume and quality scaling metrics for the performance and quality components in (10) and thus re-define the productivity metric as

$$P(\Lambda) = \eta_v(\Lambda) \cdot \eta_t(\Lambda) / c(\Lambda) \quad (11)$$

by adopting  $p(\Lambda) = \eta_v(\Lambda)$  and  $\omega(\Lambda) = \eta_t(\Lambda)$ .

As important as measuring and testing scalability is, so is collecting the right measurements, in order to interpret those measurements by the right metrics. This will develop a consistent interpretation of the fine grained performance

measurement data through the lenses of externally relevant scalability performance metrics. This interpretation will allow understanding better the factors that influence performance metrics of the scalability of cloud-based systems and will support software engineers to fine-tune such systems to achieve better performance.

#### IV. APPLICATION EXAMPLE AND RESULT

To demonstrate the applicability of the scalability metrics we used the Amazon AWS cloud environment, and the OrangeHRM<sup>1</sup> open source human resource software system as the cloud-based software service. To measure the scalability we simulate the user demand scenarios using the Apache JMeter script<sup>2</sup>, and run through Redline13<sup>3</sup> services after connecting our Amazon account to the service. To provide the scaling of the service we relied on the Auto-Scaling and Load-Balancer services provided by the Amazon AWS cloud.

We set-up and configured an EC2 instance to host the targeted application through the Amazon EC2 management console. Both Auto-Scaling and Load-Balancer services have been connected to the application instance, and we set up the CloudWatch service to monitor the scaling performance and parameters. The experimental data has been collected through both Redline13 and Amazon's CloudWatch services. In this study, the system average response time was measured as the average amount of time that the application takes to process a HTTP request after it has received one. The parameters of the Amazon EC2 virtual machines, and Auto-scaling policies that have been used for the experiments are given in Table I. The service requests consisted of a HTTP request to the main page of software with gaining login access by the following steps using the Apache JMeter:

- Path = /.
- Method = GET.
- Parameters = username, password and login button.

We used the Redline13 services by uploading the test script into our account; which allows us to easily deploy JMeter test plans inside our Amazon AWS domain and repeat the tests without the need to reset the test parameters again, this allows efficient extraction of the data.

TABLE I. AMAZON AWS EC2 VIRTUAL MACHINE PARAMETERS

Virtual Machine Parameters			
Instance type: t2.micro			
vCPUs	RAM (GiB)	CPU Credits/hr	Storage (GB)
1	1.0	6	10
Auto Scaling Policies			
Add Instance	When 80% <= CPUUtilization < +infinity		
Remove Instance	When 30% >= CPUUtilization > -infinity		

We used two demand scenarios. The first scenario follows the steady rise and fall of demand pattern shown in Fig. 2A. The second scenario consists of a series of stepwise increases and falls in demand, conceptually similar to the demand pattern shown in Fig. 2B. Examples of the two kinds of experimental demand patterns are shown in Fig. 4. We varied the volume of demand and experimented with four volume settings: 100, 200, 400 and 800 service requests in total.

We ran all experimental settings (i.e. demand pattern and demand volume combinations) 20 times, in total 160 experimental runs. We calculated the average number of simultaneously active software instances and the average response time for all service requests for each experimental run. We also calculated the averages and standard deviations of simultaneously active software instances and average response times over the 20 experimental runs. We note that the standard deviations are small relative to the averages over the 20 runs. The average number of software instances for both scenarios and for the four demand levels are shown in Fig. 5. The average response times for the two scenarios and four demand levels are shown in Fig. 6.

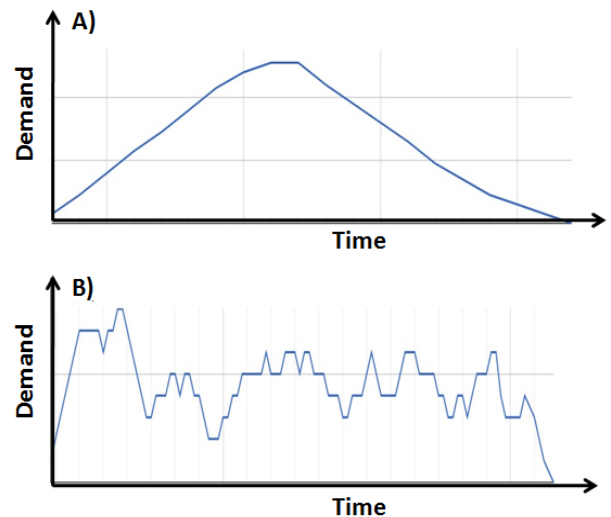


Fig. 4. Typical experimental demand patterns: A) steady rise and fall of demand; B) series of step-wise increases and decreases of demand

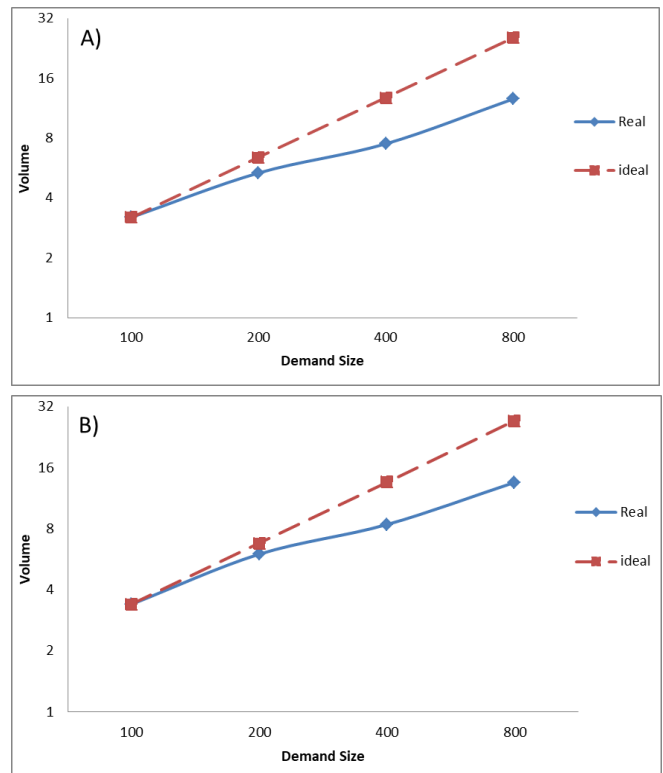


Fig. 5. The average number of software instances: A) steady rise and fall of demand; B) series of step-wise increases and decreases of demand.

<sup>1</sup> <https://www.orangehrm.com/>

<sup>2</sup> <http://jmeter.apache.org/>

<sup>3</sup> <https://www.redline13.com>

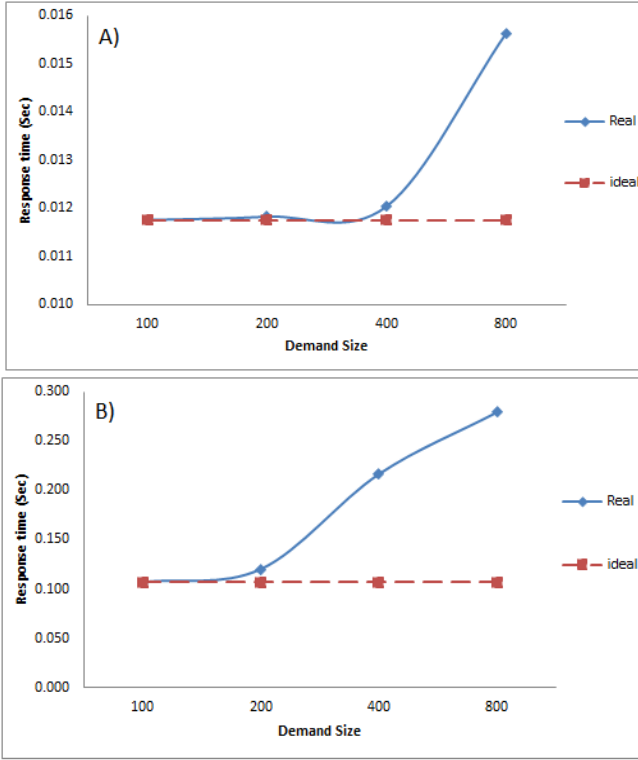


Fig. 6. The average response times: A) steady rise and fall of demand; B) series of step-wise increases and decreases of demand.

We note that the application performs similarly in term of volume (instances) scaling, while the observed average response time values for the stepped rise and fall of demand scenario are shown in Fig. 6b, starting from demand size of 200 the average response time increases significantly. In contrast, average response time values for the first scenario which shown in Fig. 6a, have increased gradually from demand size of 400 with less variation between values of average response times.

We calculated the scalability metrics  $\eta_i$  and  $\eta_t$  for the two demand scenarios that we considered. The values of the scalability metrics are shown in Table II. The calculated metrics show that in terms of volume scalability the two scenarios are similar, the scaling being slightly better in the context of the scenario with step-wise increase and decrease of demand. In terms of quality scalability, the system scales much better in the context of the first scenario, steady rise and fall of demand, than in the case of the second scenario with step-wise increase and decrease of demand.

TABLE II. SCALABILITY METRICS

Scenario	Metric	
	$\eta_i$	$\eta_t$
Steady rise and fall	0.5687	0.9041
Step-wise increase and decrease	0.5882	0.5201

The values of the metrics indicate that in the context of variable demand, which is likely to be the more realistic demand scenario for many cloud-based software services, the quality scaling performance drops considerably in comparison with the simpler demand scenario, while the volume scaling performance is retained (and even slightly improved).

## V. DISCUSSION AND LIMITATIONS

The proposed scalability metrics address both volume and quality scaling of cloud-based software services, and provide a practical measure of these features of such systems. The works do not yet integrate aspects of non-technical features [11] and also are distinct from elasticity oriented metrics [4]. This is important in order to support effective measurement and testing of scalability performance of the system.

Having an effective measure of the volume and quality scalability of the system allows exploring the contribution of various system components to the scalability performance of the system. For example, using mutation testing [28] we can test the impact of small changes to particular components on the scalability performance. Alternatively, by instrumenting the whole code of the system [29] and then measuring its scalability through a range of demand scenarios we can identify the components of the system at various resolutions (e.g. units, classes, functions, methods) that contribute critically to variations in scalability performance. Such identification of scalability-critical components can drive the design of scalability tests, system revision and upgrade focused on improvement of scalability, or development of fine-grained monitoring of system scalability performance.

In this paper the quality scaling is considered through measurement of average response time of the system. Other aspects of quality scaling could be also used to define further similar but functionally distinct quality scaling metrics. For example, system throughput (i.e. the rate of successful delivery of service provision in response to service demand), or slowdown, or recovery rate [11] can be used for alternative quality scaling metrics. Expanding the range of quality scaling metrics provides a multi-factor view of quality scaling supporting the identification and definition of trade-off options in the context of quality-of-service offerings in terms of service scaling. The equations of the quality metric can be amended based on the nature of the quality factor that could replace or combine with the current quality scaling feature.

The authors also note that over-provision of cloud service instances that exceed the ideal scaling behavior is as much of an issue as under-provision, which been taken into account in future research. On the other hand, the volume metric can be considered for extension to a larger volume.

Here we used two demand scenarios to demonstrate the effect of demands patterns on the scaling metrics. In principle, various demand scenarios may be used to fine-tune the cloud-based software service to fit particular demand scenario expectations. Similarly, considering a set of demand scenarios can also be used to identify changes in such scenarios that trigger interventions in terms of software upgrade or maintenance or direct investment of software engineering resources in development of focused upgrades for the system. Demand scenarios combined with multiple versions of quality scaling metrics can also be sued to determine reasonable quality-of-service expectations and likely variations of such expectations depending on changes in demand scenarios. We note the review [30] which concerns the study of the current practice of cloud service performance evaluation from system modelling perspective. It can be useful to adopt another demand scenario that already been used in the field, in order to track the impact of such scenarios.

The limitations of the results presented here stem from the limited nature of the experimental investigation. We used only one cloud platform (Amazon AWS) and only one cloud-based software service (OrangeHRM) to demonstrate the application and usefulness of the proposed scalability metrics. Naturally, expanding the experiments to cover multiple cloud platforms and multiple cloud-based software services would provide a fuller picture of the application of the proposed metrics. We also used only two demand scenarios, while a wider range of these would offer a deeper understanding of how the proposed metrics vary depending on demand scenarios. Finally, we used one particular setting of the cloud service (i.e. virtual machine specification), one load generator and one auto-scaler to implement the demand scenarios and the scaling of the investigated cloud-based software service. Alternative load generators and auto-scalers might have an impact on the values of the calculated metrics due to their implementation details, although in principle we would not expect major impact of these on the reported results.

## VI. CONCLUSIONS

In this paper we present two scalability metrics for cloud-based software services. One of these addresses the volume scalability of the service, while the other the quality scalability of the service. The metrics are based on simple principles of proportional scaling of the service volume and constant provision of the service quality, and are defined using the differences between the real and ideal scaling curves for both the volume and quality scalability. The proposed metrics can be used alone or integrated into utility oriented metrics of cloud-based service scalability [11].

The proposed metrics are demonstrated using a cloud-based software service run on the Amazon AWS cloud platform and considering two demand scenarios. Our results show that the proposed metrics quantify explicitly the technical scalability performance of the system and also that they allow the clear assessment of the impact of demand scenarios on the cloud-based software service.

We believe that the proposed technical scalability metrics can be used to perform and design scalability testing of cloud-based software systems with the aim to identify system components that critically contribute to the technical scalability performance. Furthermore, the proposed metrics can be extended, by considering alternative service quality features, and combined with a range of demand scenarios to support the fine-tuning of the system, the identification of quality-of-service trade-offs, and estimation of realistic scalability performance expectations about the system depending on demand scenarios.

Future work will include the consideration of other cloud platforms (e.g. Microsoft Azure, Google Cloud, and IBM), demand workload generators and auto-scalers, and other cloud-based software services, so we get a wider range of measurements of the proposed metrics, extending the practical validity of the work. We also aim to consider further demand patterns (e.g. variable width sudden peaks in demand, seasonal demand) to show how these impact on the scalability performance of cloud-based software services.

## ACKNOWLEDGMENT

This research is supported by a PhD scholarship from Philadelphia University – Jordan for Amro Al-Said Ahmad.

We thank our colleagues, Fiona Polack and Pearl Brereton, for useful comments and suggestions.

## REFERENCES

- [1] H. H. Liu, *Software Performance and Scalability: A Quantitative Approach*. Hoboken, N.J.: Wiley Publishing, 2009.
- [2] T. Atmaca, T. Begin, A. Brandwajn, and H. Castel-Taleb, "Performance Evaluation of Cloud Computing Centers with General Arrivals and Service," *IEEE Trans. Parallel Distrib. Syst.*, vol. 27, no. 8, pp. 2341–2348, 2016.
- [3] M. Becker, S. Lehrig, and S. Becker, "Systematically Deriving Quality Metrics for Cloud Computing Systems," in *Proceedings of the 6th ACM/SPEC International Conference on Performance Engineering*, 2015, pp. 169–174.
- [4] N. R. Herbst, S. Kounev, and R. Reussner, "Elasticity in Cloud Computing: What It Is, and What It Is Not," in *Proceedings of the 10th International Conference on Autonomic Computing (ICAC'13)*, 2013, pp. 23–27.
- [5] S. Lehrig, H. Eikerling, and S. Becker, "Scalability, elasticity, and efficiency in cloud computing: A systematic literature review of definitions and metrics," in *2015 11th International ACM SIGSOFT Conference on Quality of Software Architectures (QoSA)*, 2015, pp. 83–92.
- [6] R. Buyya, R. Ranjan, and R. N. Calheiros, "InterCloud: Utility-Oriented Federation of Cloud Computing Environments for Scaling of Application Services," in *Algorithms and Architectures for Parallel Processing*, 2010, pp. 13–31.
- [7] K. Hwang, Y. Shi, and X. Bai, "Scale-Out vs. Scale-Up Techniques for Cloud Performance and Productivity," in *2014 IEEE 6th International Conference on Cloud Computing Technology and Science*, 2014, pp. 763–768.
- [8] H. AlJahdali, A. Albatli, P. Garraghan, P. Townend, L. Lau, and J. Xu, "Multi-tenancy in Cloud Computing," in *2014 IEEE 8th International Symposium on Service Oriented System Engineering*, 2014, pp. 344–351.
- [9] S. Islam, K. Lee, A. Fekete, and A. Liu, "How a Consumer Can Measure Elasticity for Cloud Platforms," in *Proceedings of the 3rd ACM/SPEC International Conference on Performance Engineering*, 2012, pp. 85–96.
- [10] U. Sharma, P. Shenoy, S. Sahu, and A. Shaikh, "A Cost-Aware Elasticity Provisioning System for the Cloud," in *Proceedings of the 2011 31st International Conference on Distributed Computing Systems*, 2011, pp. 559–570.
- [11] K. Hwang, X. Bai, Y. Shi, M. Li, W. G. Chen, and Y. Wu, "Cloud Performance Modeling with Benchmark Evaluation of Elastic Scaling Strategies," *IEEE Trans. Parallel Distrib. Syst.*, vol. 27, no. 1, pp. 130–143, 2016.
- [12] K. Blokland, J. Mengerink, and M. Pol, *Testing Cloud Services: How to Test SaaS, PaaS & IaaS*. Rocky Nook, 2013.
- [13] D. Jayasinghe, S. Malkowski, J. Li, Q. Wang, Z. Wang, and C. Pu, "Variations in Performance and Scalability: An Experimental Study in IaaS Clouds Using Multi-Tier Workloads," *IEEE Trans. Serv. Comput.*, vol. 7, no. 2, pp. 293–306, 2014.
- [14] D. Jayasinghe, S. Malkowski, Q. Wang, J. Li, P. Xiong, and C. Pu, "Variations in Performance and Scalability when Migrating n-Tier Applications to Different Clouds," in *IEEE 4th International Conference on Cloud Computing*, 2011.
- [15] J. Gao, P. Pattabhiraman, X. Bai, and W. T. Tsai, "SaaS performance and scalability evaluation in clouds," in *Proceedings of 2011 IEEE 6th International Symposium on Service Oriented System (SOSE)*, 2011, pp. 61–71.
- [16] A. A.-S. Ahmad and P. Andras, "Measuring the Scalability of Cloud-based Software Services," in *2018 IEEE World Congress on Services (SERVICES)*, 2018.
- [17] B. Jennings and R. Stadler, "Resource Management in Clouds: Survey and Research Challenges," *J. Netw. Syst. Manag.*, vol. 23, no. 3, pp. 567–619, Jul. 2015.
- [18] J. Gao, X. Bai, W. T. Tsai, and T. Uehara, "SaaS Testing on Clouds - Issues, Challenges and Needs," in *2013 IEEE Seventh International Symposium on Service-Oriented System Engineering*, 2013, pp. 409–415.
- [19] E. F. Coutinho, F. R. de C. Sousa, P. A. L. Rego, D. G. Gomes, and J. N. de Souza, "Elasticity in cloud computing: a survey," *Ann. des Télécommunications*, vol. 70, no. 7–8, pp. 289–309, 2015.

- [20] Y. Hu, B. Deng, F. Peng, B. Hong, Y. Zhang, and D. Wang, "A survey on evaluating elasticity of cloud computing platform," in 2016 World Automation Congress (WAC), 2016, pp. 1–4.
- [21] N. R. Herbst, S. Kounev, A. Weber, and H. Groenda, "BUNGEE: An Elasticity Benchmark for Self-Adaptive IaaS Cloud Environments," in 2015 IEEE/ACM 10th International Symposium on Software Engineering for Adaptive and Self-Managing Systems, 2015, pp. 46–56.
- [22] A. Bauer, N. Herbst, and S. Kounev, "Design and Evaluation of a Proactive, Application-Aware Auto-Scaler: Tutorial Paper," in Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering, 2017, pp. 425–428.
- [23] M. Beltran, "Defining an Elasticity Metric for Cloud Computing Environments," in Proceedings of the 9th EAI International Conference on Performance Evaluation Methodologies and Tools, 2016, pp. 172–179.
- [24] J. Kuhlenskamp, M. Klems, and O. Röss, "Benchmarking Scalability and Elasticity of Distributed Database Systems," *Proc. VLDB Endow.*, vol. 7, no. 12, pp. 1219–1230, Aug. 2014.
- [25] A. Ilyushkin, A. Ali-Eldin, N. Herbst, A. V Papadopoulos, B. Ghit, D. Epema, and A. Iosup, "An Experimental Performance Evaluation of Autoscaling Policies for Complex Workflows," in Proceedings of the 8th ACM/SPEC on International Conference on Performance Engineering, 2017, pp. 75–86.
- [26] M. H. Jamal, A. Qadeer, W. Mahmood, A. Waheed, and J. J. Ding, "Virtual Machine Scalability on Multi-Core Processors Based Servers for Cloud Computing Workloads," in 2009 IEEE International Conference on Networking, Architecture, and Storage, 2009, pp. 90–97.
- [27] S. Lehrig, R. Sanders, G. Brataas, M. Cecowski, S. Ivanšek, and J. Polutnik, "CloudStore — towards scalability, elasticity, and efficiency benchmarking and analysis in Cloud computing," *Futur. Gener. Comput. Syst.*, vol. 78, pp. 115–126, 2018.
- [28] I. Saleh and K. Nagi, "HadoopMutator: A Cloud-Based Mutation Testing Framework," in *Software Reuse for Dynamic Systems in the Cloud and Beyond*, 2014, pp. 172–187.
- [29] H. Jayathilaka, C. Krintz, and R. Wolski, "Performance Monitoring and Root Cause Analysis for Cloud-hosted Web Applications," in Proceedings of the 26th International Conference on World Wide Web, 2017, pp. 469–478.
- [30] Q. Duan, "Cloud service performance evaluation : status , challenges , and opportunities – a survey from the system modeling perspective," *Digit. Commun. Networks*, vol. 3, no. 2, pp. 101–111, 2017.