Adaptive SLA-based Elasticity Management Algorithms for a Virtualized IP Multimedia Subsystem

Hani Nemati^{*}, Arjun Singhvi[†], Nadjia Kara^{*}, May El Barachi[‡] ^{*}Ecole de Technologie Superieure, University of Quebec, Montreal, Canada [†]PES Institute of Technology, Bengaluru, Karnataka, India [‡]Zayed University, Abu Dhabi, United Arab Emirates

*Hani.Nemati.1@ens.etsmtl.ca, Nadjia.Kara@etsmtl.ca,[†]Arjunsinghvi29@gmail.com,[‡]May.ElBarachi@zu.ac.ae

Abstract—The IP Multimedia System (IMS) is an important reference service delivery platform for next generation networks and is considered as a de-facto standard for IP-based multimedia communication services. In its current design, the IMS faces important challenges in terms of scalability and elasticity, and lacks the ability to adaptively manage the network resources and dynamically dimension the network nodes based on load and demand. Network function virtualization and cloud computing are two important concepts that can be leveraged to address those challenges in IMS environments. In this work, we propose two adaptive SLA-based elasticity management algorithms for virtualized IMS environments. Our proposed algorithms use two SLA attributes (the call setup delay and user priority) to dynamically control the CPU resources allocated/de-allocated to virtualized IMS nodes. The aims of our proposed algorithms are: 1) to ensure efficient usage and sharing of CPU resources by various IMS components; 2) to reduce the overall power consumption in virtualized IMS platforms; and 3) to enhance the user experience when using IMS networks. We have tested the proposed algorithms by setting up a virtualized IMS environment using OpenIMS Core and Xen as the hypervisor. The results obtained show that our proposed algorithms meet the SLA constraints, even when subjected to dynamic load, thereby enhancing the overall QoS. We have also compared the proposed algorithms with Xen Server's existing CPU resource scaling governors and the results indicate that our algorithms work better when compared to the existing governors.

Keywords—IP Multimedia Subsystem; Elasticity management; Network Function Virtualization; Cloud Computing; Xen. I. INTRODUCTION

Next Generation Networks (NGNs) are IP-based networks that aim at offering their users ubiquitous access to a multitude of feature-rich, QoS-enabled, broadband multimedia services [1]. The IP Multimedia Subsystem (IMS) is the key component in the NGN architecture, which enables the realization of this vision [2]. The IMS consists of a horizontal control and service layer that is deployed on top of IP-based networks. This layer encompasses a set of common functions and service logics needed for the seamless provision of IP multimedia services to users. Among the key IMS nodes, we mention: The CSCFs (Call/Session Control Functions) responsible of session control and signaling operations; The HSS (Home Subscriber Server) acting as a central database containing users related information; and the ASs (Application Servers) hosting and executing value-added services offered to users.

Cloud computing is a promising paradigm that promotes a computing-as-a-service model, in which a dynamic pool of virtualized computational resources can be leased and released on demand [3]. Some of the key benefits of the cloud computing model include: resource efficiency through sharing; high scalability and elasticity through dynamic resource management and pooling; and ease of introduction of new applications and services through substrates reuse.

Leveraging cloud technologies (e.g. hardware virtualization, virtual switches, smart NICs, poll-mode Ethernet drivers, orchestration and management mechanisms) as key enabler, an emerging concept called Network Function Virtualization (NFV) is being contemplated in the telecom domain. The main goal behind the NFV concept, which is specified by the European Telecommunications Standard Institute (ETSI) [4], is to enable the consolidation and sharing of various software-based, virtualized, telecom networking resources, running on cloud infrastructures [5]. This network virtualization and telecom infrastructure cloudification and sharing is expected to play an important role in reducing the deployment and operation costs of future telecom infrastructures, while opening the door for innovation and performance enhancements in the current telecom networking architecture.

Despite its merits, the IMS still faces important challenges in terms of scalability and elasticity. In fact, due to its reliance on text-based, bandwidth-hungry, and delay-inducing protocols (such as SIP), and with the constant increase in users number and demands, IMS nodes can become quickly overloaded, as shown in [6]. Such overload impacts the QoS offered to users and results in the lack of ability to meet SLA requirements. In addition, the current IMS design lacks the ability of dynamically dimensioning network nodes based on load and demand, and does not implement adaptive resource management mechanisms.

In this work, we show how the concepts of Network Function Virtualization and cloud computing could be used to improve the scalability and elasticity of IMS deployments. More specifically, we propose two adaptive elasticity management algorithms for the dynamic allocation of CPU resources, in virtualized IMS nodes, based on load and SLA parameters. In this work, we use the Call Setup Delay (CSD) and User Priority as SLA parameters. Furthermore, two existing CPU scaling algorithms [7] implemented in the Xen Server (the Performance and the On-demand governors) are compared to our proposed algorithms, for benchmarking and performance evaluation.

The next section presents related work in the area of dynamic resource management in virtualized IMS environments. The proposed dynamic elasticity management algorithms are presented in Section III. Section IV presents the test bed and an analysis of the performance evaluation results. We end the paper with our conclusions, and a highlight of future work.

II. RELATED WORK

In [8], the authors implement a virtualized IMS architecture with the aim of supporting dynamic resource allocation. They proposed a monitoring system to monitor the current CPU utilization and memory. In the case of overload, defined by CPU utilization threshold, the physical machine scheduler will allocate one or more virtual CPU to the VM. It worth mentioning that, the CPU utilization threshold in this work is statically preconfigured in a file, and not dynamically extracted as an SLA parameter. In [9], virtualized IMS components are installed on a set of host operator physical equipment to minimize user latencies while all capacity constraints are satisfied. The focus of the paper is on the P-CSCF, which is the first contact point with the IMS user terminal. By optimizing the location of these VMs into the suitable physical servers, the latency would be minimized. In [10] [11], authors use virtualization of IMS to dynamic reconfigure their network. In these papers, methods for migrating a session from one CSCF to another are developed. In case of failure or disaster, it is desirable to ensure reliability by continuing the operation of the CSCF on other servers. Thus, call session state migration in IMS is used to achieve the needed flexibility.

The concept of self-organizing IMS is used in [12] and [13]. Those works propose a centralized self-organizing node, which is a master node that maintains a database of the state information and capabilities for all nodes under its control. They also propose load balancing mechanisms for IMS networks.

III. PROPOSED ELASTICITY MANAGEMENT ALGORITHMS

In this section, we present two end-to-end adaptive SLAbased elasticity management mechanisms for delay sensitive applications hosted on the cloud. The main goal of these algorithms is to ensure the efficient usage of CPU resources while meeting the SLA requirements. In the proposed system, each physical CPU supports a number of working frequencies, which are termed as available CPU frequencies. The algorithms guarantee efficient usage of CPU resources by increasing/decreasing the CPU frequency to the nearest available CPU frequency and also allocating/de-allocating CPU cores based on the SLA attributes. The proposed algorithms have been tested on a distributed Virtualized IMS (VIMS) platform. The two SLA attributes taken into consideration in our proposed algorithms are: the Call Setup Delay (CSD) and the User Priority. The call setup delay, which consists of the time taken by the VIMS platform to successfully process a call initiation request, acts an indicator of the amount of load on the various VIMS components. The user priority allows preferential allocation of CPU resources to the various VIMS components when they claim the same resource, thereby enabling service differentiation based on users categories (e.g. regular user, business user, emergency user etc).

In our proposed solution, a monitoring system calculates the average CSD (Current-CSD) of the VIMS network periodically and an actuator system adjusts the CPU resources of the various VIMS components accordingly. The monitoring and the actuator system do not add any additional complexity to the VIMS network as they run independently. The monitoring system reports the Current-CSD periodically to the actuator system which compares it with maximum CSD mentioned in the SLA (SLA-CSD). The CPU frequency will be set to nearest available CPU frequency which is twice the current CPU frequency, whenever the Current-CSD exceeds the SLA-CSD. Also, a new core will be added to the VIMS component on which the algorithm is run, if the CPU frequency of current CPU is the maximum available CPU frequency. In contrast, the CPU frequency is decremented by one step whenever the Current-CSD is less than SLA-CSD. The actuator system also releases one CPU core if the CPU frequency is less than the minimum available CPU frequency. Increasing CPU resources should be fast, in order to prevent violation of the SLA and should be equal to the actual required CPU resources. In contrast, decreasing CPU resources should be gradual to prevent oscillations in the system [14]. This strategy of increasing/decreasing the CPU frequency is similar to the TCP window flow controller, which is an effective mechanism in use.

Algorithm 1 Unary SLA adaptive [USA] Elasticity Management

1:	procedure USA
2:	$CoresVM \leftarrow$ Available cores for VM
3:	$i \leftarrow$ Index of last added core to CoresVM
4:	$F \leftarrow Available \ CPU_{(i)} \ frequencies$
5:	Current-CSD \leftarrow Average CSD of last 2s
6:	if $(Current-CSD > SLA-CSD)$ then
7:	$CurrentFreq_{(i)} \leftarrow 2 * CurrentFreq_{(i)}$
8:	if $(CurrentFreq_{(i)} > Max(F))$ then
9:	$CurrentFreq_{(i)} \leftarrow Max(F)$
10:	$CoresVM \leftarrow Add \ new \ core$
11:	end if
12:	else
13:	$CurrentFreq_{(i)} \leftarrow lower available frequency$
14:	if $(CurrentFreq_{(i)} < Min(F))$ then
15:	$CurrentFreq_{(i)} \leftarrow Min(F)$
16:	Release $CPU_{(i)}$ from $CoresVM$
17:	end if
18:	end if
19:	end procedure

A. The Unary SLA Adaptive (USA) Elasticity Algorithm

The USA elasticity algorithm depicted in Algorithm 1, adjusts the CPU resources of the various VIMS components based on one SLA attribute (namely, the Call Setup Delay). This algorithm can be implemented on a single or multiple VIMS components (e.g., P-CSCF, S-CSCF).

In the proposed algorithm, *CoresVM* contains the list of cores allocated to the VIMS component (line 2). Based on *Current-CSD*, the actuator increases (line 7) /decreases (line 13) the CPU frequency of the last added core.

If twice the *CurrentFreq* of the last added core is more than the maximum available frequency, then the *CurrentFreq* of the last added core is set to the maximum available frequency and a new core is added to *CoreVM* (line 8 to 11). In contrast, if the lowered *CurrentFreq* of the last added core is less than the minimum available frequency, then the *CurrentFreq* of this core is set to the lowest available frequency and the core is released (line 14 to 17).

B. The Binary SLA Adaptive (BSA) Elasticity Algorithm

The BSA elasticity algorithm adjusts the CPU resources of the various VIMS components based on two SLA attributes (the Call Setup Delay and the User Priority), thus introducing another dimension to the USA elasticity algorithm.

Algorithm 2 Binary SLA adaptive [BSA] Elasticity Management

ment
1: procedure BSA
2: $S \leftarrow$ Index of free available cores
 for j=Index of VIMSs base on priority do
4: $i \leftarrow Index \ of \ last \ added \ core \ to \ VM_{(i)}$
5: $F \leftarrow Available CPU_{(i)} frequencies$
6: $CurrentCSD_{(i)} \leftarrow Average CSD of last 2s$
7: if $(CurrentCSD_{(i)} > SLA-CSD_{(i)})$ then
8: $CurrentFreq_{(i)} \leftarrow 2 * CurrentFreq_{(i)}$
9: if $(CurrentFreq_{(i)} > Max(F))$ then
10: $CurrentFreq_{(i)} \leftarrow Max(F)$
11: $CoreVM_{(i)} \leftarrow Select \ core \ from \ S$
12: end if
13: else
14: $CurrentFreq_{(i)} \leftarrow lower available frequency$
15: if $(CurrentFreq_{(i)} < Min(F))$ then
16: $CurrentFreq_{(i)} \leftarrow Min(F)$
17: $S \leftarrow Release CPU_{(i)} from CoresVM_{(j)}$
18: end if
19: end if
20: end for
21: end procedure

In Algorithm 2, the actuator system will decide to allocate CPU resources from its shared CPU pool (line 2) based on the priority of the various VIMS components (line 3). The allocation/de-allocation strategy followed by the BSA algorithm is similar to the USA algorithm (line 4 to 19), but the order in which the adjustment of CPU resources occurs amongst the various VIMS components depends on the priority.

IV. PREFORMANCE EVALUATION

In this section, we evaluate the performance of our proposed algorithms and compare them with the existing Xen Server's CPU scaling mechanisms namely: The Ondemand governor that increases the CPU frequency for active threads while other threads run at the lowest frequency; and the Performance governor that keeps the CPU frequency at the highest level. The test bed setup is first described, and then the performance measurements are analyzed.

A. Test Bed Setup

To setup the test environment of the proposed algorithms, a number of open source software were used. The IMS core was emulated using the OpenIMS Core [15]- an open source implementation of the IMS' core components developed by FOKUS. A modular approach was followed in which the IMS virtualized core components were installed on separate VMs. All VMs were running the Ubuntu 12.04 LTS operating system. The hardware and software specifications of the testbed are described in Table I below.

TABLE I: Test Bed Specification

Parameter	Specification/Value				
System Name	Xen Hypervisor 1	Xen Hypervisor 2	Xen		
CBU	Intel(R) Core(TM) i7-4770	Intel(R) Core(TM) i7	AMD Phenom(TM) 9600B		
cru	CPU@3.40 GHz	CPU X980@3.33 GHz	Quad-Core Processor		
Supported CPU Core Frequency (GHz)	3.401, 3.400, 3.200, 3.000,	3.334, 3.333, 3.200, 3.067,			
	2.800, 2.700, 2.500, 2.300,	2.933, 2.800, 2.667, 2.533,			
	2.100, 1.900, 1.700, 1.500,	2.400, 2,267, 2.133, 2.000,	-		
	1.400, 1.200, 1.000, 0.800	1.867, 1.733, 1.600			
Hypervisor	Xen Server 6.2	Xen Server 6.2	Xen Server 6.2		
Memory (MB)	24527	24567	4079		
Logical CPUs (Count)	8	12	4		
Virtual Machines (Deployed)	P-CSCF 1 P-CSCF 2	HSS I-CSCF	Manager IMS Bench SIPp		

As depicted in Fig. 1, the test bed consists of three physical servers connected by a switch. All the three physical servers use Xen Server 6.2 as their operating system [16]. and are



Fig. 1: The virtualized IMS test bed topology



Fig. 2: Comparison of running USA algorithm, on P-CSCF, S-CSCF and both, in terms of CSD

described as follows:

Xen Hypervisor 1: This server consists of three VMs: One VM is configured as an HSS; another is configured as an I-CSCF; and the last one is configured as a S-CSCF.

Xen Hypervisor 2 : This server consists of two VMs, each of them configured as a P-CSCF.

Xen : This server consists of two VMs, each configured with IMS Bench SIPp [17] - an open source implementation of a test system that generates SIP traffic and is capable of handling a large number of users. The IMS Bench SIPp consists of two components: the logical SIPp component and the Manager. These two components were configured on two different VMs. Also, our monitoring and actuator system is in the same VM which is configured with SIPp.

B. The USA Algorithm Performance Evaluation

Fig. 2 shows the CSD of using the USA algorithm on virtual P-CSCF (USA-P), virtual S-CSCF (USA-S) and on both virtual P-CSCF and S-CSCF (USA-PS). To study the effect of using our algorithm on P-CSCF, in USA-P we let S-CSCF work with 4 cores in the highest available frequency; and also for evaluating USA-S, P-CSCF works with 4 cores in the highest available frequency. The incoming call rate is Poisson and SLA-CSD is 8ms. During 0s to 70s we increase the call rate from 100 CPS to 300 CPS to warm-up the network. Our experimental results show that after 70s, USA-PS has almost the highest CSD and USA-P has the lowest CSD. The reason being, S-CSCF in IMS network needs more



Fig. 3: Comparison of running USA algorithm, Ondemand and Performance, in terms of CSD

CPU resources and in USA-P, S-CSCF is working with 4 cores (maximum number of available cores) in the highest CPU frequency. On the other hand, in USA-PS, both S-CSCF and P-CSCF work with the required CPU resources, which are adjusted based on the load. Thus, the CSD of both P-CSCF and S-CSCF will be more because the two components do not work with 4 cores all the time. Also, from the graph we can infer that the algorithm, when run in three different scenarios, does not violate the SLA-CSD of 8ms.

We also compared the performance of the proposed algorithm with Xen Server's existing governors, namely, Ondemand and Performance. Both of the governors use 4 cores all the time. In case of the Performance governor, all the CPU cores work with the highest available frequency. The Ondemand governor, on the other hand, increases the CPU frequency based on demand [7]. CSD of different governors in Xen Server, along with the USA algorithm is depicted in Fig. 3. It is easily noticeable that the CSD at its maximum when the Ondemand governor is used. The CSD in case of Ondemand violates the agreed upon SLA-CSD of 8ms during time interval of 70s to 260s (Refer to Table II). Another inference that can be drawn is that the Performance governor outperforms the USA algorithm, when the comparison is made only in terms of CSD. The reason being, the Performance governor always works with 4 powered-on cores whereas the USA algorithm acquires and releases CPU cores based on the requirements.

TABLE II: Average CSD for USA, Ondemand, and Performance Algorithms

Time(s) Algorithm	0-50	50-70	70-120	120-260	260-360
USA	5.348	6.455	5.186	5.925	5.464
Ondemand	6.685	7.611	8.482	8.052	6.519
Performance	3.222	4.500	3.291	5.784	3.344

Fig. 4 illustrates two important results. The first being, that the USA algorithm uses a lower number of CPU cores in comparison to the existing Xen Server's governors. In case of the USA algorithm, the free cores can be allocated to other VMs whereas the other governors have zero free cores. The second being that the USA algorithm is load-sensitive. As it can be seen, when the incoming rate is 400 CPS (highload), the USA algorithm uses more resources and when the incoming rate is 25 CPS (low-load), it releases the unnecessary



Fig. 4: Comparison of running USA algorithm on P-CSCF, S-CSCF and both, in terms of number of free cores



Fig. 5: Comparison of running USA, Ondemand, and Performance algorithms, in terms of CPU utilization

CPU cores and works with a lower number of CPU cores.

Fig. 5 shows that despite the fact that the USA algorithm uses a lower number of CPU cores as compared to Performance governor, it uses the available cores efficiently and its CPU utilization is the same as that of Performance governor. Thus, we can conclude that the USA algorithm meets the SLA constraint despite using a smaller number of CPU cores and it ensures that the CPU power consumption is less than the existing Xen Server's governors.



(a) CSD of running BSA algorithm on (b) Number of powered on cores of P-CSCF1 and P-CSCF2 P-CSCF1 and P-CSCF2

Fig. 6: Illustration of the role of Priority in the BSA algorithm.

C. The BSA Algorithm Performance Evaluation

Continuing with the same test-bed that was used to test the performance of the USA algorithm, the concept of priority in the BSA algorithm is illustrated. P-CSCF2 has higher priority as compared to P-CSCSF1. Both the P-CSCFs share the CPU



Fig. 7: Comparison of running BSA, Ondemand, and Performance algorithms, for 4 different P-CSCFs, in terms of CSD



Fig. 8: Comparison of running BSA algorithm and Xen Server's existing governors, in terms of the sum of powered on cores of the 4 P-CSCFs

resources. Fig. 6a shows the CSD of the P-CSCFs, wherein, after 60s, the CSD exceeds the 8ms, which is the SLA-CSD. In this situation, both the P-CSCFs are competing for the last core, but P-CSCF2 gets it, as it has higher priority, which is indicated in Fig. 6b. This demonstrates that the CPU resources are allocated not only based on load but also based on priority.

The existing test-bed was modified to accommodate 2 more P-CSCFs. As shown in Fig. 7, the average CSD of 4 P-CSCFs in case of Ondemand violates the agreed upon SLA CSD of 8ms during 120s to 260s, but it never goes more than 8ms for BSA and Performance. Fig. 8 shows that the BSA algorithm allocates and de-allocates the CPU resources when the load is high (400 cps) and when the load is low (25 cps) respectively. It also indicates that our algorithm could reduce the power consumption by turning off unnecessary cores.

V. CONCLUSION

In this paper, we proposed two adaptive SLA-based elasticity management algorithms (the USA and the BSA algorithms) for virtualized IMS environments. The two SLA attributes used were the call setup delay and the user priority. By setting up a virtual IMS environment, we tested the performance of our proposed algorithms. Using the test-bed, we obtain results which show that the USA and the BSA algorithms never violate the SLA requirements. The comparisons made with Xen Server's existing CPU resource scaling governors indicate that our proposed algorithms outperform them. Both algorithms are efficient and they do not contribute to the load of the VIMS as they can be run on an independent monitoring system. Thus, the proposed algorithms achieve what they set out to do: Allocate CPU resources efficiently; reduce the power consumption by the CPU resources; allow sharing of CPU resources among the various IMS components; and deliver an overall enhanced end user experience.

As future work, we intend to do further analysis of the algorithms performance when introducing equal and unequal link delays to the test-bed, in order to mimic a more realistic setup. We also plan to introduce other parameters and dimensions to our proposed algorithms, such as the throughput and the geographic location of the users, which can have an influence on the virtual nodes dimensioning and their placement.

REFERENCES

- G. Camarillo and M. Garcia-Martin, The 3G IP Multimedia Subsystem (IMS): Merging the Internet and the Cellular Worlds, Second edition. England: John Wiley and Sons Ltd., May 2006.
- [2] The Third Generation Partnership Project, [online] available at: http://www.3gpp.org/
- [3] A. Mishra et al., guest editors, Cloud Computing: Networking and Communications Challenges, IEEE Communications Magazine, Vol.50, No.9, September 2012, pp.24-25.
- [4] The European Telecommunications Standard Institute (ETSI), [online] available at: http://www.etsi.org/
- [5] White paper on Network Function Virtualization, [online] available at: http://portal.etsi.org/NFV/NFV_White_Paper.pdf
- [6] Homayouni, M.; Nemati, H.; Azhari, V.; Akbari, A., "Controlling Overload in SIP Proxies: An Adaptive Window Based Approach Using No Explicit Feedback," Global Telecommunications Conference (GLOBECOM 2010), 2010 IEEE, 6-10 Dec. 2010
- [7] V. Pallipadi, A. Starikovskiy,"The Ondemand Governor," Proceedings of the Linux Symposium, 2006, Ottawa, Canada
- [8] F. Lu, H. Pan, X. Lei, X. Liao, and H. Jin, A Virtualization-Based Cloud Infrastructure for IMS Core Network, 2013 IEEE 5th Int. Conf. Cloud Comput. Technol. Sci., pp. 2532, Dec. 2013.
- [9] O. C. Imen Limam Bedhiaf, A Virtual Machine Location Problem for IP Multimedia Subsystem, Cyber Journals Multidiscip. Journals Sci. Technol. J. Sel. Areas Telecommun., vol. 3, no. 5, 2013.
- [10] S. Komorita, M. Ito, Y. Kitatsuji, and H. Yokota, Dynamic IMS Reconfiguration using Session Migration for Power Saving, in AICT 2013: The Ninth Advanced International Conference on Telecommunications, 2013, no. c, pp. 191196.
- [11] Y. K. and M. O. Takashi Shimizu, Tetsuya Nakamura, Michio Kiuchi, Atsushi Iwata, An experimental evaluation of dynamic virtualized networking resource control on an evolved mobile core network, in Humanitarian Technology Conference (R10-HTC), 2013 IEEE Region 10, 2013, pp. 270275.
- [12] A. Dutta, C. Makaya, S. Das, S. Komorita, T. Chiba, H. Yokota, and H. Schulzrinne, Self organizing IP Multimedia Subsystem, 2009 IEEE Int. Conf. Internet Multimed. Serv. Archit. Appl., pp. 16, Dec. 2009.
- [13] M. T. Alam, Cost Analysis for Self-Organizing IP Multimedia Subsystem, Res. India Publ., vol. 4, no. 3, pp. 291297, 2011.
- [14] Ali-Eldin, A.; Tordsson, J.; Elmroth, E., "An adaptive hybrid elasticity controller for cloud infrastructures," Network Operations and Management Symposium (NOMS), 2012 IEEE, pp.204,212, 16-20 April 2012
- [15] F. Fokus, Open IMS Core, [online] available at: http://www. openimscore. org.
- [16] "Xen Server, [online] available at: http://www.citrix.com/products/xenserver/overview.html."
- [17] [17]. "SIPp, [online] available at: http://sipp.sourceforge.net/index.html"