

De-Registration Based S-CSCF Load Balancing in IMS Core Network

Liang Xu¹, Changcheng Huang¹, James Yan¹ and Tadeusz Drwiega²

¹Department of Systems and Computer Engineering, Carleton University, Ottawa, Canada

²Nortel Networks, Ottawa, Canada

{liangxu, huang}@sce.carleton.ca, jim.yan@sympatico.ca, drwiega@nortel.com

Abstract—IP Multimedia Subsystem (IMS) provides a common control layer over which the services can be easily accessed by fixed and mobile users. In order to guarantee performance to both service providers and end service subscribers, careful management and maintenance within the control layer must be fulfilled. In this paper, we propose a solution that utilizes a controlling entity to track the utilization of hotspot Session Initiation Protocol (SIP) server (i.e., S-CSCF) and to initiate de-registration procedure so that a sufficient number of subscribers at the overloaded S-CSCF are re-associated with the other S-CSCF(s) for subsequent session services. The proposed solution redistributes and smoothes the future traffic load over individual S-CSCF in an automatic manner. The goodness of the proposal is assessed through simulations using OPNET.

Index Terms— IP Multimedia Subsystem, SIP Server, Load Balancing, Autonomic Management

I. INTRODUCTION

The IP Multimedia Subsystem (IMS), as the next-generation IP-based service framework, provides access to multimedia applications by wireless and wired terminals to achieve Fixed-Mobile Convergence (FMC). The framework was originally proposed by 3GPP to support GPRS network, and has been then actively updated and developed [1]. The key merit of the framework is its decoupling of the service and control logic. This decoupling enables service providers to develop and introduce new service features in a more plug-and-play manner without the need to consider underlying control logic. The framework also enables network operators to manage the control functions in a common control layer more consistently [2]. The control layer contains complete subscriber information and provides service session control and management. Because the control layer is critical in providing the right quality of experience to the end user, our research focus in this paper is to explore how the layer can improve service performance.

Figure 1 illustrates a simplified architecture of the IMS control layer. The main elements of the layer are: the Home Subscriber Servers (HSS) and different purposed Call/Session Control Function (CSCF) servers. The control layer connects to the Application Servers (AS) in the service layer and to the User Equipment (UE) through transport layer [3]. The Proxy-CSCF (P-CSCF) acts as the contact point with UE for core network service access; Interrogating-CSCF (I-CSCF)

provides routing for the signaling messages; and Serving-CSCF (S-CSCF) takes the full responsibility in UE registration, session control, and service routing with AS. Indeed, the S-CSCF unit has to virtually handle all major tasks in the control layer and it is where the majority of signaling traffic has to be processed.

In an IMS network, P-CSCF and S-CSCF servers are assigned to serve individual end users. Due to the random nature of traffic, it is inevitable to have traffic anomalies that can potentially overload certain elements. Overloads may cause severe performance degradation [4], and even activate the packet retransmission which can result in the unrecoverable server collapse [5]. We consider in this paper the S-CSCF as the most likely affected network element, and propose an autonomic load-balancing mechanism that redirects traffic among S-CSCFs to minimize the risk of congestion.

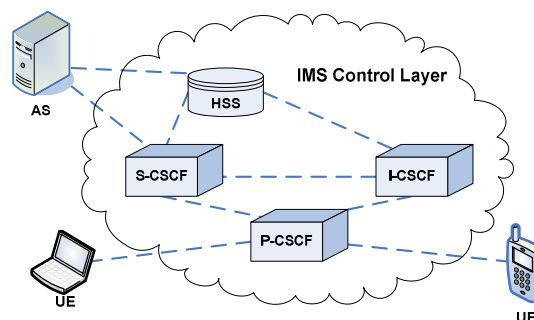


Fig.1. Simplified IMS control layer overview

Unlike traffic routing in usual packet-switched network, SIP-based session traffic [6] has to pass through different network elements in definite sequences to complete the specific tasks. Section II in this paper reviews the predominant interaction sequences in the IMS core network to aid in characterizing the traffic properties and to convince the need in proper load management among S-CSCFs. Section III presents the proposed S-CSCF load-balancing solution that utilizes SIP de-registration procedure to redirect subscriber associations and hence to facilitate load re-routing among the S-CSCFs. In Section IV, we introduce the simulation model in evaluating the goodness of the proposal, and give the primary results from simulations. Finally, we have conclusion and envisioned future work in Section V.

II. IMS SESSION PROCEDURES AND TRAFFIC PROPERTIES

The IMS core network supports the provision of services

This work was supported in part by NSERC Grant CRDPJ 354729-07 and NSERC Grant RGP 261469-2003.

through a complete set of signaling procedures [7]. The main procedures included are UE registration and de-registration, P-CSCF discovery, S-CSCF assignment, session establishment and termination, QoS negotiation and home network directing. In this paper, we consider session establishment, UE registration and de-registration as the predominant signaling traffic, where UE registration and de-registration allows the subscribers access IMS services while session establishment realizes the access.

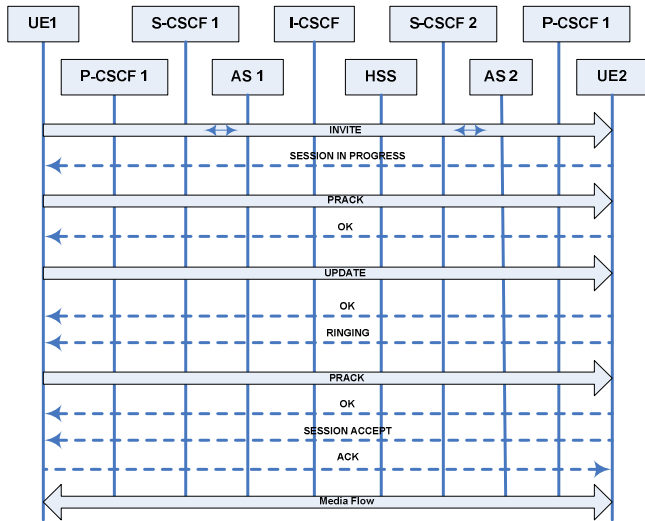


Fig.2. IMS Session establishment procedure

Figure 2 presents a basic session establishment procedure between two UEs, which can also be treated as analogous to VoIP call setup or session-based mode Instant Messaging (IM) setup [2] between two end subscribers. For the sake of simplicity, the service session is assumed to be established between UEs registered and visiting the same network. Figure 3 presents the UE registration procedure through which the UE requests for authorization to use IMS services. For management purpose, UE is required to perform re-registration periodically, which is virtually the same as registration procedure in case that network mandates the security challenge for each registration. Figure 4 presents the UE de-registration with two cases included. The first case is an UE initiated de-registration and the second case is deregistration initiated by an administrative network. More detailed definition of the procedures can be found in [3, 7].

It is important to note that once an UE completes its registration, it is bonded to a specific S-CSCF for the rest of session services unless another registration procedure re-introduces it to another S-CSCF. Hence, suppose that all the registered UEs are associated to a certain S-CSCF, the rest of S-CSCFs in the network would never be visited by any traffic. It is also noteworthy of having AS in our flow diagram, which is based on the assumption that the S-CSCFs always perform certain service control upon receiving the initial validated SIP request and triggers at least one AS.

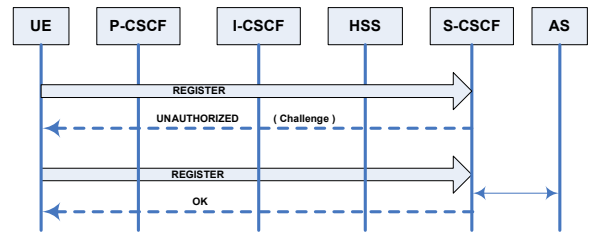


Fig.3. IMS UE registration procedure

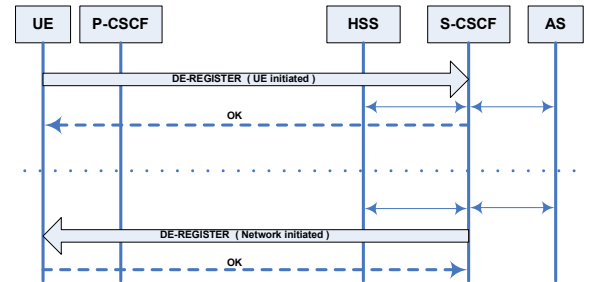


Fig.4. IMS UE de-registration procedure. The upper de-registration case is UE initiated, and the bottom case is network initiated.

Based on the procedures illustrated and the assumptions made above, the visiting frequency of different service session interaction flows over each core network elements are counted and listed in Table 1. The combined traffic loads over each network element are derived in Table 2 given the assumption that at a specific time period the SIP messages for session establishment and UE registration/de-registration overwhelm the messages for other session requests. Our assumption on traffic load composition is based on the fact that service accessing (e.g., VoIP and IM) represents the majority of the load, and there are constant subscriber registrations and de-registrations.

Type of service session	P-CSCF	I-CSCF	S-CSCF	HSS
	(visit per session)			
Session Establishment (SE)	24	6	26	1
UE Registration /Re-registration (UR)	4	6	5	2
UE initiated De-registration (UD)	2	0	3	1
Network initiated De-registration (ND)	2	0	3	1

Table.1. Network elements visit frequency per service session

Traffic composite SE / UR / UD,ND	P-CSCF	I-CSCF	S-CSCF	HSS
	(load percentage)			
0.9 / 0.05 / 0.05	0.417	0.108	0.453	0.022
0.8 / 0.1 / 0.1	0.414	0.113	0.451	0.022
0.7 / 0.2 / 0.1	0.406	0.123	0.444	0.027
0.7 / 0.1 / 0.2	0.411	0.112	0.451	0.026
0.6 / 0.2 / 0.2	0.402	0.124	0.443	0.031

Table.2. Network elements load percentage under different service traffic combinations based on Table.1

Through the preliminary quantitative analysis from Table 2, the first observation we can make is that P-CSCF and S-CSCF take the majority of traffic load in the network (over 80%) and

are at highest risk to become congested network element. Since the S-CSCF is a stateful proxy server involved in service sessions, it can potentially experience memory leaks and eventually server failure due to hanging state machine. Also, the S-CSCF has to possibly interact with more than one AS to satisfy the service control, and hence can induce itself more traffic load. Multiple S-CSCF (and P-CSCF) units are deployed in networks to serve different sets of subscribers and hence distribute the load, but this can not ultimately prevent individual units from overloading or even server failure due to the dynamics in subscriber and service behavior. The second observation from both tables is that the increased percentage in registration/de-registration traffic can only slightly affect the overall load distribution over different network elements. Hence, the changes in registration and de-registration rate may affect the consequent total network traffic load, but would not change significantly the general load distribution.

These observations are the basis for us to address the S-CSCF load-balancing problem in the rest of this paper.

III. DE-REGISTRATION BASED S-CSCF LOAD-BALANCING

As reviewed and discussed in the previous section, the individual S-CSCF are potentially the most affected network unit by traffic anomalies. Also, the traffic routing in the IMS core network is fundamentally based on S-CSCF and UE association. For this reason, to address load-balancing is to address S-CSCF/UE re-association in its root.

Procedures which can be used in S-CSCF/UE re-associating are registration, re-registration and de-registration (followed by registration) as presented previously. Registration associates new UE with S-CSCF, but the existing S-CSCF/UE associations are not changed. This means that by relying only on the registration procedure, the consequent traffic re-routing is limited and small-scaled in the short-run. Re-registration is primarily for the purpose to refresh the existing associations of UE; but since the S-CSCF load balancing is an administrative responsibility of network operators, a network-initiated de-registration procedure is more suitable to be used in our proposal to selectively break existing associations at overloaded S-CSCF and to force the re-registration.

For different reasons (e.g., traffic management, service specification), a home network administrative function may determine a need to clear a subscriber's registration by initiating a de-registration procedure. De-registration can be treated as a mechanism by the network operator to ensure stable network operation and carrier-grade service, and this function initiator may reside in various elements depending on the exact reason for initiating the de-registration [3]. In our proposal, an AS is deployed in the application layer to communicate with S-CSCFs from the core network, and acting as the load-balancing decision maker. Its responsibility includes periodically collecting status information (e.g., server utilization) from the S-CSCFs, evaluating the S-CSCF load condition based on the information collected and initiating de-registration requests to the specific S-CSCF. Optionally, it can also communicate with HSS to zoom in on subscriber registration patterns and to initiate de-registration requests more intelligently.

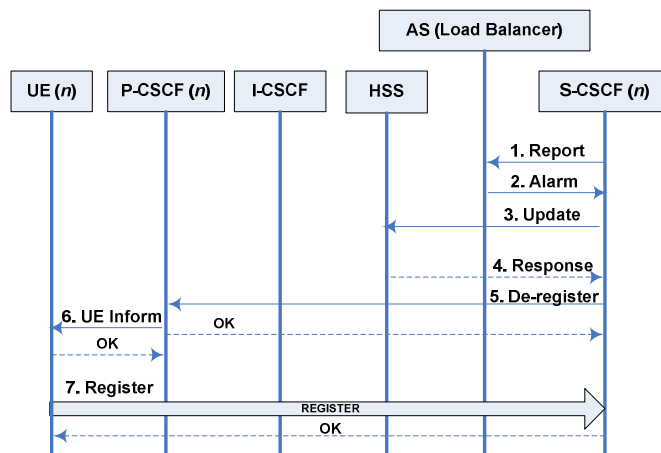


Fig.5. IMS UE registration procedure

Figure 5 presents the sequence diagram of the proposed de-registration based S-CSCF load balancing. The diagram emphasizes on the interactions involving the load-balancer, and the registration procedure after de-registration is again shown in brevity as it was shown in detail earlier. It should be noted that the de-registrations are only applied to the stand-by UEs, and not to the UEs actively involved in any on-going sessions. It is up to the individual S-CSCFs, once receiving any decision from the load-balancer, to decide which specific set of UEs should be de-registered. The detailed S-CSCF load balancing procedure is as follows:

1. Individual S-CSCF reports to the load-balancer the current server status (e.g., utilization) periodically.
2. The load-balancer evaluates the aggregated information from S-CSCFs and decides if any particular S-CSCF has, or is about to encounter traffic anomaly. The affected S-CSCF is alarmed and the request of de-registration is initiated.
3. The alarmed S-CSCF decides to which subscribers the de-registration request are sent, and informs the HSS.
4. Optionally, the HSS can aid S-CSCF in the decision process and suggest particular subscribers to be redirected.
5. The alarmed S-CSCF de-associates with the specific subscribers and informs the subscriber associated P-CSCFs.
6. Individual P-CSCF de-associates with the subscriber and informs the subscriber.
7. The de-registered subscriber re-initiates registration through the usual procedure. At this point, the S-CSCF selection policy at I-CSCF is open.

To effectively drive the traffic load, several factors have been further considered. First of all, the reporting cycle between S-CSCF and the load-balancer must be reasonably determined to address the tradeoff between response latency and messaging efficiency. It is desired to not have the individual S-CSCF too frequently reporting the current status to the load-balancer in order to avoid the extra traffic burden over the network; while a minimum information exchanging should be guaranteed so that the load-balancer can respond to the potential overloading in a timely fashion. Since the traffic load over individual S-CSCF is largely determined by the load of incoming requests, we associate the request arrival rate with the reporting cycle in our solution. Specifically, the S-CSCF reports the current status after every certain amount of request

arrivals. Therefore, the higher the arrival rate (i.e., the higher the presenting traffic load), the sooner the S-CSCF would report the current status. This approach enables the S-CSCF reaction being more adaptive to the network condition. The actual request arrival threshold can be further adjusted depending on the scale of the network.

Secondly, proper mechanism should be applied in the load-balancer to evaluate the aggregated server utilization information and to decide whether to initiate the de-registration procedure in any specific S-CSCF. Intuitively, S-CSCF is supposed to react once its utilization reaches a critical level. And it is the objective of the load-balancer in our solution to analyze the server information collectively and to adaptively define the critical level for each individual S-CSCF. In the current proposal, the load-balancer actively maintains the averaged utilization of all S-CSCFs. For any newly reported server utilization, the value is to be compared to the current averaged utilization. If an S-CSCF has utilization certain percentage higher than the averaged value, it is deemed to be the affected S-CSCF. In other words, we are more emphasizing on the short term traffic changes in terms of the increments in server utilization. In this way, the traffic anomaly can be early detected and addressed before the server utilization hits the critical level. It is noteworthy here that as the averaged utilization going up, the percentage increment in determining the procedure triggering can be decreased. This is to increase the anomaly detection sensitivity in a more saturated network.

Lastly and most importantly, the actual subscribers who are to be re-associated must be well budgeted. On one hand, the total number of affected subscribers should be as small as possible. And on the other hand, the affected subscribers should not be the ones with any ongoing service sessions (or with as less as possible). This is primarily concerning the service experiences of the end users. In fact, it largely relies on the information exchange among the load-balancer, the S-CSCF and the HSS to address the issue. Ideally, if the load-balancer could have access to the HSS database, or specific subscriber behaviors such as the group activity pattern and demographic distribution can be derived. These can aid the load-balancer to initiate more intelligent requests to the individual S-CSCF in subscriber resettling. For simulations in the following section, we assume that the load-balancer acknowledges the subscriber activity and assume that the UEs involved in S-CSCF re-association are with no ongoing session (i.e., the stand-by UEs). In case for any triggering that the total number of re-associating requested UEs exceeds the total number of stand-by UEs, only the stand-by UEs will be re-associated.

In short, the combined strategies as discussed above are to enable the solution to monitor and respond to the network conditions in a more autonomic manner to achieve server load-balancing with also early avoidance of individual overloading.

IV. NETWORK SIMULATION

Most of current researches on IMS derive the network condition and evaluate engineering protocols quantitatively [8-10], and lack actual available testing networks. We further evaluate our proposal by detailed network simulations.

The IMS network largely depends on the performance of SIP

servers to minimize the latency in session signaling. The test bench we deployed in OPNET is a simplified IMS core network that mainly addresses on the S-CSCF SIP server. The test bench follows strictly the session interactions to allow messages being exchanged in the network to represent the IMS traffic properties to a best extent.

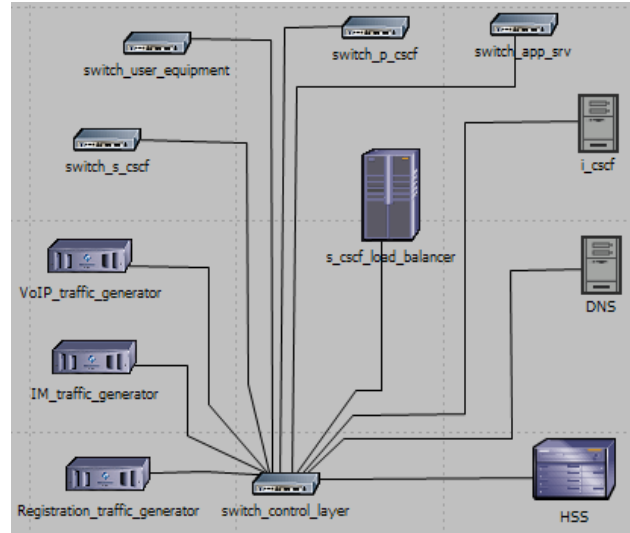


Fig. 6. IMS core network deployment in OPNET

Figure 6 illustrates the test bed developed using OPNET. The network topology includes the network elements we analyze in this paper. For brevity, the individual SIP servers and UEs are represented by single units in the figure. In actual testing, there are 10 individual P-CSCF and 10 S-CSCF deployed in total. Message exchanges among the individual elements are facilitated by the switching units prefixed with *switch_*. The network traffic is generated from the traffic generators shown in the figure with *_traffic_generator* as suffix. The primary testing parameters are listed in Table 3. We assume the arrival rate of session service requests follow a Poisson process, and the service response time for each request at S-CSCF unit is based on the exponential distribution. To avoid complicating the message flows, in the current configuration the session origination and termination are served by the same network operator.

Parameter Description	Basic Setting
VoIP traffic inter-arrival time (S-S#2)	Exponential (10ms)
IM traffic inter-arrival time (S-S#2)	Exponential (50ms)
Registration/Re-registration traffic (S-S#2)	Exponential (100ms)
Average SIP request / response message size	900 bytes / 500 bytes
HSS latency	50ms
S-CSCF service time	Exponential (0.5ms)

Table.3. Primary parameter setting (S-S#2, single network operator).

To create the traffic anomaly, an extra burst of VoIP session requests is generated at an arrival interval of exponentially 5ms for a duration of 5 minutes targeting one S-CSCF. It is expected to see initially the affected S-CSCF to experience higher traffic load and therefore higher server utilization, and we are interested to see how effective the proposal is to address overloading by enforcing the subscriber re-association. The server utilization reported by S-CSCF is calculated as:

$$u(n) = \alpha \frac{\text{total_busy_time}(n)}{\text{total_simulated_time}(n)} + (1 - \alpha) \frac{\text{total_busy_time}(n-1)}{\text{total_simulated_time}(n-1)} \quad (1)$$

where the function weight α is set as 0.7, and n indicates the every 10-second observation instance. For every 10,000 request arrivals, S-CSCF will perform the utilization calculation and send the result to load-balancer.

Figure 7 and 8 presents the server utilization conditions of two selected S-CSCFs without and with the load-balancing. Recordings 1, 2 and 3 (as also indicated in the figures) represent the instant utilization of overloaded S-CSCF, instant utilization of one arbitrary normal S-CSCF and the calculated utilization of overloaded S-CSCF, respectively. R4 in Figure 8 indicates the time instances at which the overloaded S-CSCF reports the current calculated utilization (i.e., Recording 3).

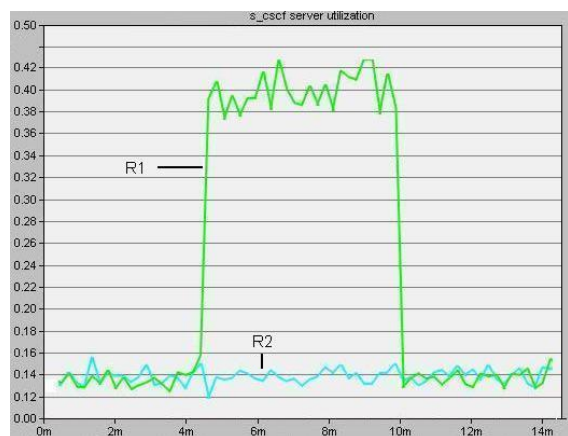


Fig.7. S-CSCF utilization recorded for 15 minutes without load-balancing.

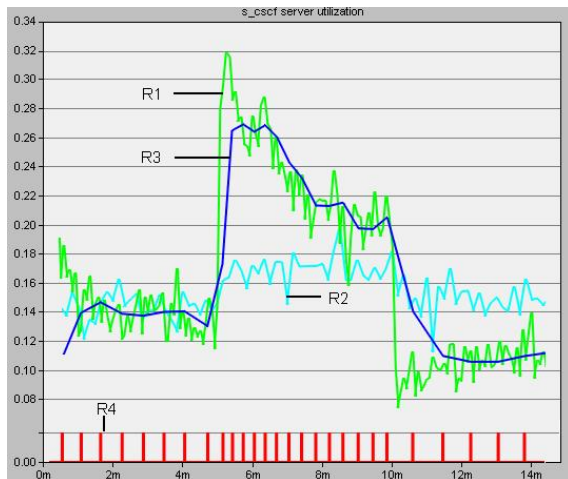


Fig.8. S-CSCF utilization recorded for 15 minutes with load-balancing, where the red vertical bars at the bottom indicates the reporting interval of the overloaded S-CSCF.

As shown in Figure 8, the server utilization starts to increase at the targeted S-CSCF (traced by R1 and R3), and the dramatic increase in request arrival causes the target to report current condition more frequently (traced by R4). In the shown simulation, the load-balancer treats any S-CSCF with utilization 20% higher than the averaged value as affected S-CSCF, and the affected S-CSCF is then instructed to de-register 5% of the current associated UEs for each triggering. The de-registered UEs are then initiating the registration procedure immediately. It should be noticed of the under-utilizing from R1 (by comparing to R2), this is because

the S-CSCF accumulatively de-registered 22.7% of the original total associated stand-by UEs which results in less consequent requests in short run.

V. CONCLUSION AND FUTURE WORK

In this paper, we addressed the use of de-registration procedure based load-balancing mechanism in requesting re-association between subscriber and S-CSCF pair, and therefore facilitate in re-directing consequent traffic from the over-utilized S-CSCF to the others. The simulation results suggest that the proposal can effectively drive the traffic to avoid further overloading at potentially over-utilized S-CSCF while not inducing severe extra load.

Since our research did not explore in detail the mechanisms and strategies in the interaction among load-balancer, HSS and S-CSCF, further studies can be emphasized on the issue to have a load-balancer being more intelligent and adaptive in initiating UE re-association as discussed in this paper. Also, I-CSCF and load-balancer can collaborate in S-CSCF selection to improve the resulting traffic balancing.

VI. ACKNOWLEDGEMENT

We would like to thank Nortel Networks for their generous help and invaluable advice through this research.

REFERENCES

- [1] 3GPP TS 23.002, "Network Architecture," Release 7, v7.1.0, Mar. 2003.
- [2] G. Camarillo, Miguel-Angel, and Garcia-Martin, *The 3G IP Multimedia Subsystem (IMS): Merging the Internet and the Cellular Worlds*, John Wiley and Sons, August 2004.
- [3] 3GPP TS 23.228, "IP Multimedia Subsystem (IMS), Stage 2," Release 7, v7.6.0, Dec. 2006.
- [4] E.C. Noel and C.R. Johnson, "Initial Simulation Results That Analyze SIP Based VoIP Networks Under Overload," *20th International Teletraffic Congress*, Ottawa, Canada, June 2007.
- [5] V. Planat and N. Kara, "SIP Signaling Retransmission Analysis over 3G network," MoMM2006, Yogyakarta, Indonesia, December 2006.
- [6] J. Rosenberg, H. Schulzrinne, G. Camarillo, A. Johnson, J. Peterson, R. Sparks, M. Handley, and E. Schooler, "The session initiation protocol (SIP) RFC 3261," IETF, 2001.
- [7] 3GPP TS 29.228, "IP Multimedia Subsystem Cx and Dx Interfaces; Signaling Flows and Message Contents," Release 7, v7.6.0, Dec. 2006.
- [8] A. Kist and R. Harris, "SIP signaling delay in 3GPP," *In Proceedings of Sixth International Symposium on Communications Interworking of IFIP*, Perth, Australia, October 13-16 2002.
- [9] V. K. Gurbani, L. Jagadeesan, and V. B. Mendiratta, "Characterizing session initiation protocol (SIP) network performance and reliability," *International Service Availability Symposium*, April 2005.
- [10] V.S.Abhayawardhana and R.Babbage, "A Traffic Model for the IP Multimedia Subsystem," *Vehicular Technology Conference*, April 2007.