

Stability Condition for SIP Retransmission Mechanism: Analysis and Performance Evaluation

Yang Hong, Changcheng Huang, James Yan

Dept. of Systems and Computer Engineering, Carleton University, Ottawa, Canada

E-mail: {yanghong, huang}@sce.carleton.ca, jim.yan@sympatico.ca

Abstract—SIP (Session Initiation Protocol) has been widely adopted as a signaling protocol to establish, modify and terminate multimedia sessions between end-users in the Internet. SIP introduces a retransmission mechanism to ensure the reliability of its real-time message delivery. However, retransmission makes server overload worse, as indicated by the recent server crashes in the real carrier networks. In this paper, we use a discrete time model to describe the queuing dynamics of an overloaded SIP server with the retransmission mechanism. We then derive a sufficient stability condition that a SIP server can handle the overload effectively under the retransmission mechanism. Discrete time model allows us to run fluid-based Matlab simulation directly to evaluate the overload performance. This approach is much simpler than event-driven simulation. Event-driven OPNET simulation was also conducted to observe the transient behaviour of an overloaded server in a SIP network. Our simulation results demonstrate that: (1) The sufficient stability bound is quite tight. The bound indicates that effective CPU utilization as low as 20% can still lead to an unstable system after a short period of demand burst or a temporary server slowdown. Resource overprovisioning is not a viable solution to the server crash problem; (2) By satisfying the stability condition, the initial queue size introduced by a transient overload can avoid a system crash. Such stability condition can help the operator to determine whether and when to activate overload control mechanism in case of heavy load.

Keywords—SIP, Retransmission, Overload, Stability Condition, CPU Utilization

1. INTRODUCTION

SIP (Session Initiation Protocol) [1] has been widely deployed for significantly growing session-oriented applications in the Internet, such as Voice-over-IP, instant messaging and video conference. As a signaling protocol, SIP is responsible for creating, modifying and terminating sessions in a mutual real-time communication [2]. 3GPP (3rd Generation Partnership Project) has adopted SIP as the basis of the IMS (IP Multimedia Subsystem) architecture [3, 4, 5].

Fig. 1 illustrates a simplified configuration of a SIP network which consists of two basic elements: UA (User Agent) and P-Server (Proxy Server) [1]. A UA may perform two roles: in the UAC (User Agent Client) role, the originating UA sends requests; in the UAS (User Agent Server) role, the terminating UA receives requests and sends responses. The task of a P-server is to receive SIP requests and forward them to the terminating UA (or to another P-server that is closer to the terminating UA). Each P-server is assigned to serve multiple individual UAs. UA and P-server cooperate to establish, modify and terminate sessions for multimedia communication.

SIP is designed to be an application layer protocol independent of the underlying transport mechanism which

may be TCP (Transmission Control Protocol) or UDP (User Datagram Protocol). SIP introduces a retransmission mechanism to maintain its reliability [5, 6]. The mechanism provides reliability by retransmitting lost SIP messages either end-to-end or hop-by-hop. A SIP source uses a delay to detect a message loss. It would produce one or more retransmissions if the corresponding reply message is not received in a predetermined time interval. If a retransmission triggered by a delay caused by the overload, it would introduce the overhead rather than reliability into the network. Such redundant retransmissions increase the memory and CPU loads for a SIP server, which may cause a system overload and deteriorate the signaling performance [6-18]. In an overload situation, the throughput drops down to a small fraction of the original processing capacity, thus poses a serious problem for a SIP network [11]. This kind of behaviour has happened in real carrier networks where large scale collapses of SIP servers have been observed in present of sudden heavy SIP traffic (e.g., emergency-induced call volume) [12].

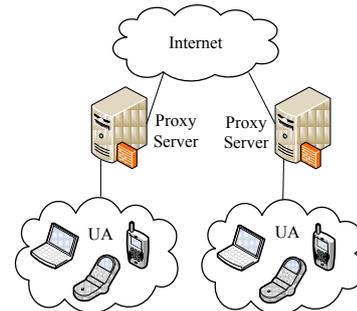


Fig. 1. Simplified configuration of a SIP network

A SIP server can be overloaded due to various reasons such as poor capacity planning, dependency failures, component failures, avalanche restart, flash crowds, denial of service attacks, etc., as indicated by RFC 5390 [20]. In general, a short period of demand burst or a server slowdown may bring a server overload and lead to server crash. The built-in SIP overload control mechanism has proven to be ineffective in practice, because it attempts to mitigate the overload by rejecting some calls, but the cost of rejecting a SIP session is comparable with the cost of serving a session.

SIP retransmission mechanism should be disabled for hop-by-hop transaction when running SIP over TCP to avoid redundant retransmissions at both SIP and TCP layer [1]. However, almost all real SIP networks run SIP over UDP mainly because the following reasons [8-18, 21-22]: (1) TCP is optimized for accurate delivery by sacrificing its timeliness which is a critical requirement for real-time application such as SIP; (2) SIP works at application layer while TCP works at transport layer. Even TCP can provide reliability at transport

layer, SIP messages can still be dropped or corrupted while being processed at application layer; (3) TCP keeps retransmitting outstanding packets until an ACK is received. Each retransmitted packet is pushed to the application layer to be a SIP message which costs extra CPU time and introduces more delay, therefore making CPU overload worse.

Recent collapse of SIP servers due to emergency-induced call volume or “American Idol” flash crowd in the real carrier networks has motivated some overload control solutions. For example, three window-based feedback algorithms were proposed to adjust the message sending rate of the upstream SIP servers based on the queue length [11]. Both centralized and distributed overload control mechanisms for SIP were investigated in [12]. Retry-after control, processor occupancy control, queue delay control and window based control were proposed to improve goodput and prevent overload collapse in [8]. A small buffer size has been proved to be a simple overload control mechanism at a cost of temporary call rejection rate hike in [23]. It has been revealed that the retransmission mechanism is a main factor to deteriorate the overload performance [11, 18]. Thus it is necessary to investigate the impact of the retransmission mechanism on the SIP overload. A demand burst or routine server maintenance such as database synchronization may accumulate the signaling messages to create a long queue. Excessive queuing delay, introduced by a long initial queue size, may continue to trigger the redundant retransmissions and crash the server after an overloaded server resumes its normal service with a low effective CPU utilization. It would be interesting to find a sufficient stability condition for the initial queue size, which indicates whether the SIP server can handle overload effectively. Such stability condition can help the SIP operator to decide whether and when to activate the overload control algorithm. It also can help researchers propose more effective solutions to avoid SIP overload collapse caused by the SIP retransmissions.

Event-driven simulation has been widely used for evaluating network performance. Its computation cost grows linearly with network sizes and message volumes [24]. When event-driven simulation is used to evaluate a SIP network, each outstanding SIP message requires a timer being maintained. When an overload happens, outstanding messages are built up, and the simulator needs to increase the number of timers dramatically in order to track message retransmissions. Tracking and manipulating these timers consume large amount memory and CPU time which make the simulation process extremely slow, thus in some cases, cause the simulator to crash and terminate simulation unexpectedly. In order to simplify the CPU-consuming timer-tracking process, fluid-based simulation tracks time slot instead of individual messages. Messages arriving within the same time slot are aggregated and processed together. This will greatly simplify the complexity caused by large number of messages and allow smooth scalability by choosing different granularities as required.

The contributions of this paper are: (1) Deriving a sufficient stability condition that a SIP server can handle the overload effectively under the retransmission mechanism; (2) Developing a discrete-time, fluid-based model to reduce significantly simulation effort; (3) Comparing the results of both fluid-based Matlab simulation and event-driven OPNET

simulation to demonstrated that the fluid-based simulation is relatively accurate and scalable for evaluating the performance of a SIP network; (4) Performing fluid-based Matlab simulation and event-driven OPNET simulation to verify that the stability bound is quite tight, or an initial queue size (created by a transient overload) which is 5% higher than the bound will bring a server crash. Slightly different initial queue sizes (the difference is less than 10% in this paper) create totally different dynamic behaviour patterns; (5) Simulating an application scenario to demonstrate that an effective CPU utilization as low as 20% cannot prevent a SIP server from overload during a short period of maintenance service and such overload continues to spread even after the normal service resumes.

The paper is organized as follows. Section 2 describes the SIP retransmission mechanism. Section 3 analyzes the queuing dynamics of an overloaded SIP server under retransmission mechanism. Section 4 derives a stability condition for SIP retransmission mechanism in the case of server overload. Section 5 evaluates the performance of an overloaded server. Some conclusions are made in Section 6.

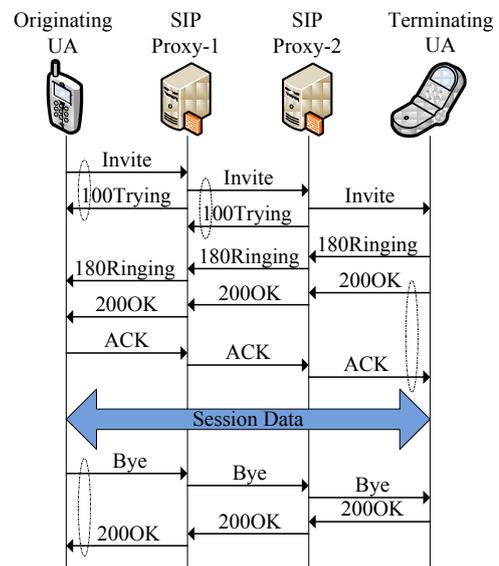


Fig. 2. A typical procedure of session establishment

2. SIP RETRANSMISSION MECHANISM OVERVIEW

SIP works in the application-layer for multimedia session establishment and tear-down. To briefly describe the basic SIP operation, we only consider originating UA, SIP P-server and terminating UA, as shown in Fig. 2. To set up a call, an originating UA sends an “Invite” request to a terminating UA via two P-servers. The P-server returns a provisional “100 (Trying)” response to confirm the receipt of the “Invite” request. The terminating UA returns an “180 (Ringing)” response after confirming that the parameters are appropriate. It also evicts a “200 (OK)” message to answer the call. The originating UA sends an “ACK” response to the terminating UA after receiving the “200 (OK)” message. Finally the call session is established and the multimedia communication is created between the originating UA and the terminating UA through the SIP session. The “Bye” request is generated to terminate the session thus cancel the communication.

SIP has two types of message retransmission: (a) a sender starts the first retransmission of the original message at T_1 seconds, the time interval doubling after every retransmission (exponential backoff), if the corresponding reply message is not received. The last retransmission is sent out at the maximum time interval $64 \times T_1$ seconds. Default value of T_1 is 0.5s, thus there is a maximum of 6 retransmissions. The hop-by-hop “Invite”-“Trying” transaction shown in Fig. 2 follows this rule [1]; (b) a sender starts the first retransmission of the original message at T_1 seconds, the time interval doubling after every retransmission but capping off at T_2 seconds, if the corresponding reply message is not received. The last retransmission is sent out at the maximum time interval $64 \times T_1$ seconds. Default value of T_2 is 4s, thus there is a maximum of 10 retransmissions. The end-to-end “OK”-“ACK” and “Bye”-“OK” transactions shown in Fig. 2 follows this rule [1].

3. QUEUING DYNAMICS OF SIP RETRANSMISSION MECHANISM

A real SIP network consists of a series of geographically distributed P-servers and a large amount of UAs. Each P-server is responsible for setting up a session call between two UAs. It forwards the requests and also generates a provisional response to confirm the receipt of every request from the upstream sender (an originating UA or a P-server) [1]. It provides a retransmission mechanism to guarantee a reliable delivery of a SIP message [5]. However, the arrival of too many SIP messages may cause an unnecessary queuing delay, stimulate redundant retransmissions and accelerate the overload, thus eventually bring down the entire network [11]. Therefore, it is necessary to describe the queuing dynamics of an overloaded SIP server (e.g., [16, 23]), before we derive a stability condition for SIP retransmission mechanism.

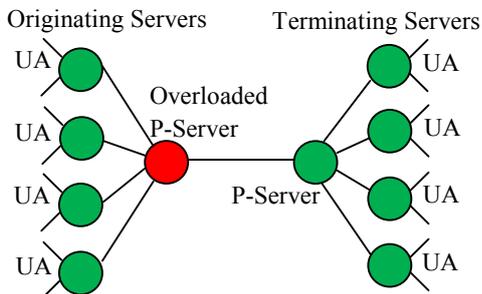


Fig. 3. SIP network topology with an overloaded P-server (which is marked in red color) and its multiple upstream originating servers

When overloads happen in the network, at any time, one of the servers will be the most overloaded one among all the overloaded servers. Without loss of generality, we consider a typical SIP network which consists of an overloaded P-server and its multiple upstream originating servers [11], as shown in Fig. 3. For a clear presentation, we use difference equations to describe the queuing dynamics of an overloaded P-server. In our discrete time model, we make the following assumptions according to SIP RFC [1]:

(a) We investigate the retransmissions which are mainly caused by long queuing delay of the overloaded server. Therefore, for the round trip response time between the overloaded server and its neighbouring server, the queuing and processing delays are dominant, while transmission and propagation delay are negligible [12]. This assumption is valid

because signaling messages are typically CPU capacity constrained rather than bandwidth constrained;

(b) Time is divided into discrete time slots. This makes it easy to describe how many retransmitted messages are triggered by a delay caused by the overload. The errors introduced by the discrete time slot can be made arbitrarily small by making the interval of a timeslot smaller and smaller. We use t and n to denote time and timeslot respectively;

(c) The SIP RFC [1] does not specify the queuing and scheduling discipline to be deployed by a SIP server. We assume that a SIP server maintains a First-In-First-Out (FIFO) queue for messages arriving at different time-slots. This FIFO queuing model reflects the common practice by most vendors today [11]. Within the same time slot, original request messages enter the tail of the queue prior to retransmitted request messages. Such enqueueing priority has negligible impact if the interval of the time slot is very small. There is no enqueueing difference for the messages arriving at different time slots;

(d) The time to process a response message or a timer timeout is typically much smaller than a request message [1]. We assumed that, within a time slot, the server has enough CPU capacity to process the incoming response messages, thus response messages will not be enqueued as long as they are treated with higher priority such as interrupt. They will not be dropped either when the queue for request messages are overflowed. The service capacity of the overloaded server includes the rate for processing response messages;

(e) In order to focus our analysis on the overloaded server, we assume multiple upstream originating servers and the downstream server of the overloaded P-server have sufficient capacity to process all requests, retransmissions, and response messages immediately without any delay;

(f) Practical buffer sizes vary with the actual service rates and system configuration plans. With the memory becoming cheaper and cheaper, typical buffer sizes are likely to become larger and larger. The buffer sizes for all servers are assumed to be large enough to hold all the incoming messages. Therefore there is no message loss at all servers.

(g) The hop-by-hop Invite-100Trying transaction is the major workload contributor due to its role for call setup and its hop-by-hop retransmission mechanism [1]. Given the proportionate nature and the general similarity of the retransmission mechanisms between the “Invite” and “non-Invite” messages in a typical session [1], we will focus on the hop-by-hop Invite-100Trying transaction and ignore other end-to-end transactions in this paper.

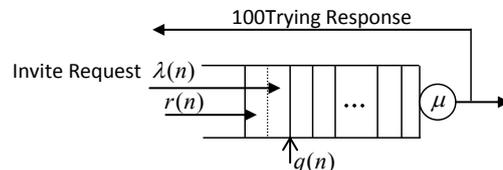


Fig. 4. Queuing dynamics of an overloaded SIP server ($\lambda(n)$ denotes aggregated original message arrivals, $r(n)$ denotes aggregated retransmitted message arrivals from multiple upstream servers, $q(n)$ denotes queue size, $\mu(n)$ denotes service rate)

Fig. 4 depicts the queuing dynamics of an overloaded SIP server in a SIP network (as shown in Fig. 3). The overloaded server receives the original Invite requests with an aggregate

rate $\lambda(n)$ at time slot n , where $\lambda(n)$ can be arbitrary. We can obtain the queue size $q(n+1)$ at next time slot $n+1$ based on the information at the current time slot n , i.e.,

$$q(n+1) = [q(n) + \lambda(n) + r(n) - \mu(n)]^+ \quad (1)$$

where $q(n)$ denotes the queue size; $r(n)$ denotes the aggregated retransmitted messages; $\mu(n)$ denotes an arbitrary service process for the request messages, which is equal to the server service capacity minus the service rate for the response messages. $\lambda(n)$ plus $r(n)$ give the total arrival messages at current time slot n . Adding $q(n)$ and deducting $\mu(n)$ would generate a new queue size $q(n+1)$ in the next time slot $n+1$, as described by Eq. (1). We use $[\]^+$ to indicate that the queue size at each time slot should be nonnegative.

If the server does not receive the corresponding response message for an original request message at a specific time-out, it would trigger retransmission. There are maximum 6 retransmissions for every original request message [1]. Thus we can obtain the total retransmitted messages $r(n)$ at current time slot n as

$$r(n) = \sum_{j=1}^6 r_j(n), \quad (2)$$

where $r_j(n)$ denotes the j^{th} -time retransmission for the original request messages arriving at time $n-T_j$, $T_j=(2^j-1)T_1$ and $1 \leq j \leq 6$.

At time $(n-T_j)$, the original request message arrivals were $\lambda(n-T_j)$ and the queue size was $q(n-T_j)$. Since the overloaded server can process $\sum_{k=1}^{T_j} \mu(n-T_j+k)$ messages during the T_j time slots, the remaining messages at current time slot n become

$$[\lambda(n-T_j) + q(n-T_j) - \sum_{k=1}^{T_j} \mu(n-T_j+k)]^+, \text{ which should be nonnegative. This may include both the original arrival messages at time } (n-T_j) \text{ and the queued messages right before the time slot } (n-T_j). \text{ However, only the remaining original arrival messages } \lambda(n-T_j) \text{ need to be retransmitted at current time } n, \text{ we use } \min\{\} \text{ function to get the } j^{\text{th}}\text{-time retransmitted messages at current time slot } n, \text{ i.e.,}$$

$$r_j(n) = \min\{[\lambda(n-T_j) + q(n-T_j) - \sum_{k=1}^{T_j} \mu(n-T_j+k)]^+, \lambda(n-T_j)\} \quad (3)$$

Eqs. (1) to (3) shows the dynamic behaviour of an overloaded SIP server. Due to its nonlinear characteristic, it may show complex, sometimes chaotic, patterns that bring a potential server collapse.

4. STABILITY CONDITION FOR SIP RETRANSMISSION MECHANISM

The retransmission can provide a reliable delivery of SIP messages. However, it also increases the queuing size and enhances the overload. It would be interesting to derive a stability condition that the server can handle overload effectively. The messages accumulated by a transient overload (e.g., a demand burst or a server slowdown) create an initial queue size when the server returns to its normal service state. Such initial queue size may bring a queuing delay long enough for the retransmissions of old remaining original messages in the queue as well as all the new incoming original messages. We would like to investigate whether the SIP server can serve

the original messages in the initial queue size and their retransmissions under a low effective CPU utilization.

Without loss of generality, we consider the ‘‘Invite-Trying’’ request-response pair with a deterministic arrival rate λ and a deterministic service rate μ ; there are i retransmissions for the new arrival original messages; the initial queue size is $q(0)$.

Theorem 1: If the initial queue size $q(0)$ created by a demand burst can satisfy a sufficient stability condition described by Eq. (4), then the SIP server is stable.

$$q(0) < \min\left\{ (2^{j+1}-1)\mu T_1, \frac{(2^{j+1}+3 \times 2^i - i - 4)\mu T_1 - ((i-1)2^i + 1)\lambda T_1}{i+1}, 1 \leq i \leq j \leq 6 \right\} \quad (4)$$

Proof:

To prevent messages from accumulating unlimitedly in SIP server, the total average incoming rate should be less than the service rate. Assume that there would be i retransmissions for an arbitrary original Invite request message, a conservative condition to maintain stability is

$$(i+1)\lambda / \mu \leq 1,$$

which is equivalent to

$$\mu \geq (i+1)\lambda, \quad (5)$$

i.e.,

$$i \leq (\mu - \lambda) / \lambda. \quad (6)$$

To achieve the above sufficient stability condition, we need to guarantee that the original messages from both the initial queue size and the new arrivals are not retransmitted more than j times, where we denote j as $j = \lfloor (\mu - \lambda) / \lambda \rfloor$.

Then we update the equivalent stability condition in Eq. (6) as $i \leq j = \lfloor (\mu - \lambda) / \lambda \rfloor$. (7)

To avoid $j+1$ retransmissions for the original messages in the initial queue size, we obtain a stability condition for the initial queue size as

$$q(0) / \mu < T_{j+1} = (2^{j+1} - 1)T_1,$$

which is equivalent to

$$q(0) < \mu T_{j+1}. \quad (8)$$

To avoid $(j+1)$ retransmissions for any newly arrival original messages, the queue size in any time should satisfy

$$q(t) < \mu T_{j+1}. \quad (9)$$

Eq. (4) can certainly satisfy Eq. (8). To show that Eq. (4) can satisfy the requirement of Eq. (9), we consider five cases in the following discussion.

1. We first consider the queue sizes at each specified retransmission times $T_i=(2^i-1)T_1$ using Eqs. (1)-(3) as follows,

$$\begin{aligned} q(T_1) &= q(0) - (\mu - \lambda)T_1 + [q(0) - \mu T_1]^+, \\ q(T_2) &= q(T_1) - 2(\mu - 2\lambda)T_1 + [q(0) - \mu T_2]^+, \\ &\vdots \\ q(T_i) &= q(T_{i-1}) - 2^{i-1}(\mu - i\lambda)T_1 + [q(0) - \mu T_i]^+, \\ &\vdots \end{aligned} \quad (10)$$

$$q(T_6) = q(T_5) - 32(\mu - 6\lambda)T_1 + [q(0) - \mu T_6]^+.$$

2. We next consider the queue sizes between any two neighbouring retransmission times T_{i-1} and T_i . Eqs. (1)-(3), (7) and (10) lead to

$$q(t) = q(T_{i-1}) - (\mu - i\lambda)(t - T_{i-1}) < q(T_{i-1}), \quad (11)$$

The inequality in Eq. (11) indicates that the queue size is decreasing continuously with a slope of $\mu - i\lambda$ during the time period. However, at time $t=T_i$, the i th retransmission for the remaining $[q(0) - \mu T_i]$ messages from the initial queue size $q(0)$ is triggered, resulting in a sudden increase in the queue size described by (10). Then when $0 < t \leq T_j$, the condition described by (9) becomes

$$q(T_i) < \mu T_{j+1} \quad 1 \leq i \leq j \leq 6. \quad (12)$$

Given the condition of Eq. (8), we assume the worst case with $q(0) - \mu T_i \geq 0$. Using recursive substitution for Eq. (10), we can obtain

$$q(T_i) = (i+1)q(0) - \sum_{k=1}^i 2^{k-1}(\mu - k\lambda)T_1 - \sum_{k=1}^i (2^k - 1)\mu T_1,$$

which can be reorganized as

$$q(T_i) = (i+1)q(0) - \sum_{k=1}^i 2^{k-1} \mu T_1 + \frac{d}{dx} \sum_{k=1}^i \lambda T_1 x^k \Big|_{x=2} - \sum_{k=1}^i 2^k \mu T_1 + i \mu T_1,$$

or

$$q(T_i) = (i+1)q(0) - \frac{1-2^i}{1-2} \mu T_1 + \frac{d}{dx} \left[\frac{x-x^{i+1}}{1-x} \lambda T_1 \right] \Big|_{x=2} - \frac{2-2^{i+1}}{1-2} \mu T_1 + i \mu T_1.$$

Then we can obtain

$$q(T_i) = (i+1)q(0) + ((i-1)2^i + 1)\lambda T_1 - (3 \times 2^i - i - 3)\mu T_1. \quad (13)$$

Combining Eqs. (12) and (13), we can obtain the second condition in Eq. (4) as

$$q(0) < \frac{(2^{j+1} + 3 \times 2^j - i - 4)\mu T_1 - ((i-1)2^i + 1)\lambda T_1}{i+1} \quad 1 \leq i \leq j \leq 6 \quad (14)$$

3. We then consider the case that $T_j < t < T_{j+1}$. From Eqs. (1)-(3), (7), (9) and (12), we have

$$q(t) = q(T_j) - (\mu - (j+1)\lambda)(t - T_j) \leq q(T_j) < \mu T_{j+1}. \quad (15)$$

This means the queue size is non-increasing during the time period $T_j < t < T_{j+1}$.

4. Next, we consider the case that $t = T_{j+1}$. Since Eq. (8) indicates $[q(0) - \mu T_{j+1}]^+ = 0$, from Eqs. (1)-(3), (7) and (12), we can obtain

$$q(T_{j+1}) = q(T_j) - 2^j(\mu - (j+1)\lambda)T_1 + [q(0) - \mu T_{j+1}]^+ \leq q(T_j) < \mu T_{j+1}. \quad (16)$$

5. Finally, we consider the case that $t > T_{j+1}$. From Eqs. (1)-(3), (7) and (16), we have

$$q(t) = q(T_{j+1}) - (\mu - (j+1)\lambda)(t - T_{j+1}) \leq q(T_{j+1}) < \mu T_{j+1}. \quad (17)$$

Combining Eq. (8), (11), (12), (14), (15), (16) and (17), we can reach a sufficient stability condition for the initial queue size described by Eq. (4). \square

5. Performance Evaluation and Simulation

Since violating the sufficient stability condition does not always bring the instability to a SIP system, we would like to investigate how tight the sufficient stability bound for the retransmission mechanism is when a SIP overload happens. To achieve this goal, we will evaluate the performance of an overloaded SIP server by performing fluid-based Matlab simulation using the analytical model described by Eqs. (1) to

(3), where the time slot is 50ms. In the mean time, in order to validate the accuracy and scalability of fluid-based simulation, we also performed event-driven OPNET simulation in a real SIP network as depicted by Fig. 3. Four originating servers generated original request messages with equal rate, and then sent them to four terminating servers via two P-servers. In our OPNET simulation, messages were enqueued based on first-come-first-in principle. That is, Assumption (c) was unnecessary for OPNET simulation. The default timer for the first retransmission was $T_1 = 0.5$ s [1].

The retransmission messages triggered by the overload are redundant messages. Therefore, only the CPU consumed by the original messages can be regarded as effective use of resources. We define effective CPU utilization ρ as the ratio between the total mean arrival rate for the original messages and the mean service rate, i.e., $\rho = \lambda/\mu$.

To verify the sufficient stability condition for the initial queue size, we have considered two scenarios: (1) Arrival rate and service rate were deterministic and the overload was caused by a demand burst; (2) Arrival rate and service rate were Poisson distributed¹ and the overload was caused by a server slowdown.

5.1. Constant Arrival Rate and Service Rate

In this scenario, a demand burst overloaded the server and created an initial queue size at time $t=0$ s. This emulated a short surge of user demands. Normal original request messages arrived at the overloaded server with a constant rate $\lambda=200$ messages/sec. This emulated regular user demands. The overloaded server maintained a constant service rate $\mu=1000$ messages/sec. Thus the effective CPU utilization for regular user demands is $\rho = \lambda/\mu = 20\%$. The simulation time is 50s.

Eq. (7) gives $j = \lfloor (\mu - \lambda)/\lambda \rfloor = 4$. Then using Eq. (4), we can obtain the stability condition for the overloaded server as $q(0) < \min\{15500, 8200, 6167, 5700, 6220\} = 5700$ messages. We will consider two sub-scenarios with different initial queue sizes.

Sub-scenario (a)

In this sub-scenario, a demand burst created an initial queue size as $q(0) = 5500$ messages < 5700 messages, obeying the stability condition described by Eq. (4).

Figs. 5 and 6 show the dynamic behaviour of the overloaded SIP server using both Matlab simulation and OPNET simulation. One can observe that the curves obtained by Matlab simulation are very close to the curves obtained by OPNET simulation. The difference for instantaneous retransmission rate shown in Fig. 6 was caused by enqueueing priority within the same time slot (Assumption (c)). The similarity between Matlab simulation result and OPNET simulation result demonstrates that fluid-based simulation is a relatively accurate and cost-effective approach for performance evaluation of a SIP network, while it can simplify a CPU-consuming timer-tracking process by tracking single time slot instead of individual message timer.

¹ Currently there is no measurement result for the workload in the real SIP networks. Poisson distributed message arrival rate and service rate are widely adopted by most existing research work (e.g., [11]).

Fig. 5(b) shows that the queue size decreased linearly with 800 messages/sec at the beginning.

At time $t=T_1=0.5s$, the overloaded SIP server had processed 500 messages, the 1st-time retransmission for the residual 5000 original messages in the initial queue happened (as shown in Fig. 6). The new 100 original messages arriving between $t=0s$ and $t=T_1=0.5s$ joined the queue together with 5000 retransmitted SIP messages, so the queue size became 10,100 messages (as shown in Fig. 5(b)). The new arrival original messages at time $t=0s$ started to trigger the first-time retransmissions (as shown in Fig. 6). Similarly, due to the 2nd-time and 3rd-time retransmissions, the queue size increased dramatically at time $t=T_2=1.5s$ and $t=T_3=3.5s$ respectively.

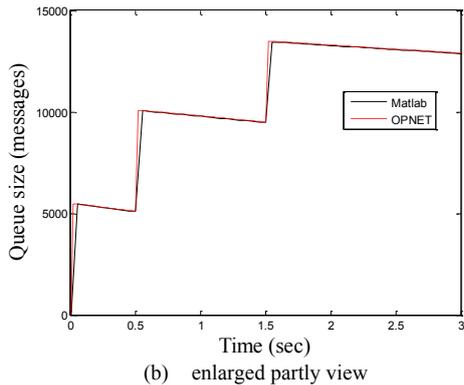
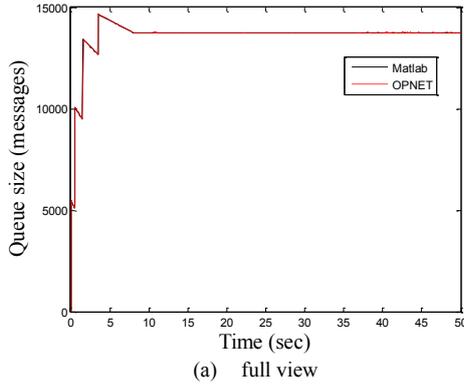


Fig. 5. Queue size q (messages) versus time for the overloaded server when the initial queue size obeys the stability condition

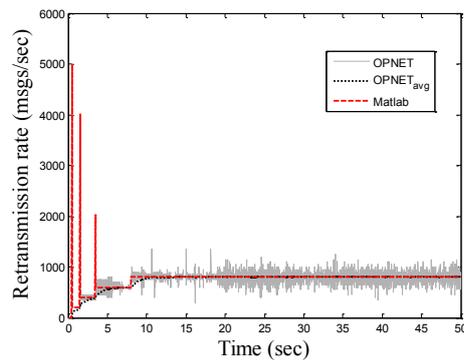


Fig. 6. Retransmission rate r and moving average retransmission rate r_{avg} (messages/sec) versus time for the overloaded server when the initial queue size obeys the stability condition

At time $t=8s$, the retransmission rate of new arrival original messages increased from 600 messages/sec to 800

messages/sec (as shown in Fig. 6), thus the total incoming traffic rate of both original messages and retransmitted messages was equal to the service rate $\mu=1000$ messages/sec (or $\rho'=5\lambda/\mu=1$). Between the time $t=3.5s$ and $t=8s$, 900 new incoming original messages and 2700 incoming retransmitted messages entered the overloaded SIP server, thus the queue size reached and stayed at a steady queue size as $14700+900+2700-4500=13800$ messages, well matching our theoretical analysis on the queuing dynamics in Section 3.

Sub-scenario (b)

In this sub-scenario, a demand burst created an initial queue size as $q(0)=6000$ messages > 5700 messages, violating the stability condition described by Eq. (4).

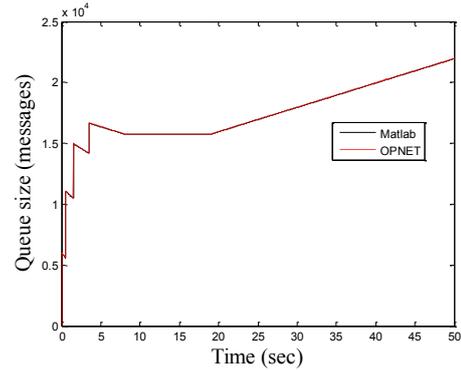


Fig. 7. Queue size q (messages) versus time for the overloaded server when the initial queue size violates the stability condition

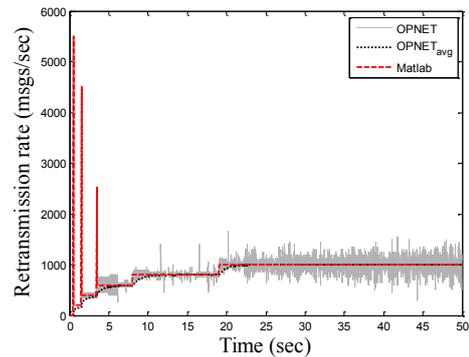


Fig. 8. Retransmission rate r and moving average retransmission rate r_{avg} (messages/sec) versus time for the overloaded server when the initial queue size violates the stability condition

Fig. 7 shows that the queue size decreased linearly except 4 spikes due to the dramatic retransmissions until the time $t=19s$. At time $t=19s$, the retransmission rate of new arrival original messages increased from 800 messages/sec to 1000 messages/sec (as shown in Fig. 8). The total incoming traffic rate of both original messages and retransmitted messages was larger than the service rate $\mu=1000$ messages/sec (or $\rho'=6\lambda/\mu=1.2>1$). Therefore, after the time $t=19s$, the queue size increased linearly and continuously with 200 messages/sec (as shown in Fig. 7), which would bring a SIP server crash eventually.

In summary, slightly different initial queue sizes due to the demand bursts (the difference is less than 10% in the two sub-scenarios) create totally different dynamic behaviour patterns. The slightly smaller initial queue size of 5500 messages

allows the server to handle the initial temporary overload effectively, while the slightly larger initial queue size of 6000 messages will result in infinitely increasing queue size, thus bring a SIP server to crash. This indicates that the sufficient stability bound is quite tight.

5.2. Poisson Distributed Arrival Rate and Service Rate

In this application scenario, the overloaded SIP server worked in one of the two states (i.e., normal service state and maintenance state) alternately. During the maintenance period, the overload may happen due to the server slow down. The mean service time at the normal service state was $m_1=600\text{sec}$; the mean service time at the maintenance state was $m_0=30\text{sec}$; all were exponential distributed. The mean service rate at the normal service state was $\mu_1=1000$ messages/sec; the mean service rate at maintenance state was $\mu_2=200$ messages/sec; the mean arrival rate of the SIP messages was $\lambda=200$ messages/sec; all were Poisson distributed. The simulation time is 2000s. The overall effective mean utilization was equal to $\rho = \lambda(m_1 + m_0)/(m_1\mu_1 + m_0\mu_0) \approx 0.2$. We will not show the OPNET simulation result because it was very close to the Matlab simulation result as the Subsection 5.1.

Figs. 9 to 12 show the dynamic behaviour of the overloaded SIP server under two different service states.

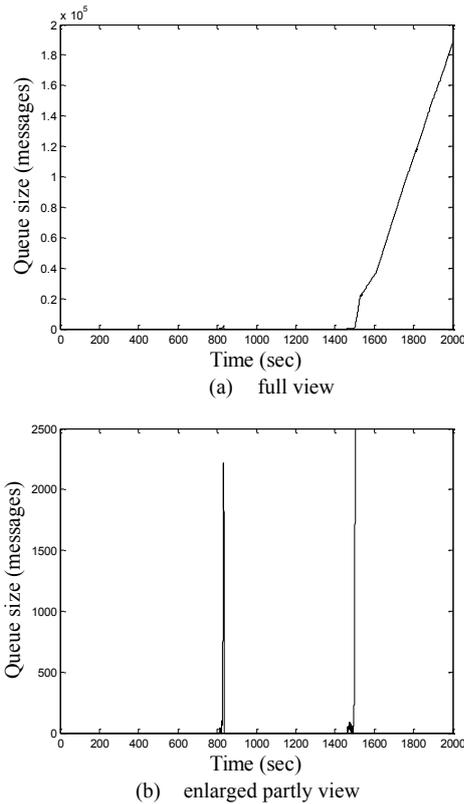


Fig. 9. Queue size q (messages) versus time for the overloaded server which performed normal service and maintenance service alternately

Between the time $t=812\text{s}$ and 833s , SIP server had short period of maintenance, service rate decreased (as shown in Fig. 12). The messages started to accumulate and the queue size increased to reach a peak around 2200 messages at time

$t=833\text{s}$ (as shown in Fig. 9(b)). When a mean service rate was 200 messages/sec, the queue size larger than 100 messages brought a queuing delay longer than 0.5s, and started to stimulate the retransmission (as shown in Fig. 11). After the server resumed normal service at time $t=833\text{s}$, the initial queue size was less than 5700 messages (the stability condition described by Eq. (4)). The server could process these accumulated messages in time, so the queue size decreased until the buffer was empty at time $t\approx 838\text{s}$ (as shown in Fig. 9(b)).

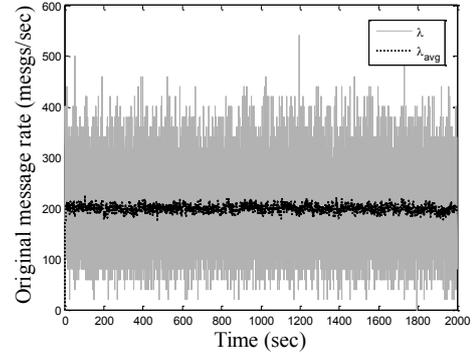


Fig. 10. Original message arrival rate λ and moving average original message arrival rate λ_{avg} (messages/sec) versus time for the overloaded server which performed normal service and maintenance service alternately

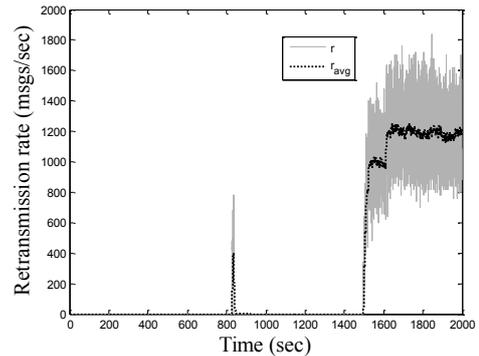


Fig. 11. Retransmission rate r and moving average retransmission rate r_{avg} (messages/sec) versus time for the overloaded server which performed normal service and maintenance service alternately

However, maintenance with a relatively long period (or a equivalent large demand burst) happened at time $t=1462\text{s}$, the queue size increased continuously and triggered more than 5 retransmissions that made the total arrival message arrival rate exceed the normal service rate (as shown in Fig. 11). After the server entered the normal service state at time $t=1527\text{s}$, the initial queue size was larger than 5700 messages. Since the stability condition for the initial queue size was violated, the SIP server cannot handle the overload effectively. The queue size tended to infinity (as shown in Fig. 9(a)), thus eventually crashed the server.

In summary, although the effective mean utilization is as low as 20%, if the accumulated messages in the SIP server during the short maintenance period violate the stability condition for the initial queue size, the server cannot mitigate the overload effectively after it resumed its normal service. Goodput collapse persists and the server would crash eventually, well matching our theoretical analysis.

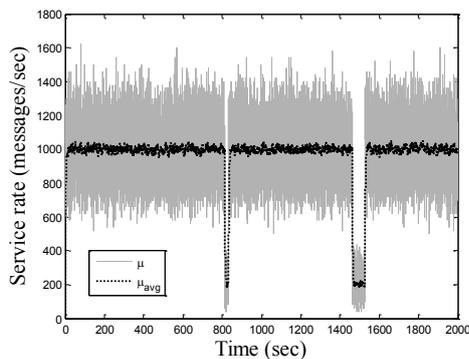


Fig. 12. Service rate μ and moving average service rate μ_{avg} (messages/sec) versus time for the overloaded server which performed normal service and maintenance service alternately

6. CONCLUSIONS

We have investigated the SIP retransmission mechanism in case of the overload. We have derived a sufficient stability condition that SIP server can handle the overload effectively under the retransmission mechanism. To prevent the system from crashing, the initial queue size caused by a transient overload should satisfy the stability condition. Such stability condition can help the SIP operator to trigger the overload control algorithm ahead of time to avoid the SIP server collapse.

We have performed simulation using both fluid-based simulation approach and event-driven simulation approach to evaluate the performance of an overloaded SIP server. Our study indicated that the behaviour of the SIP server is highly sensitive to the temporary overload due to the demand burst or the server slow down. The sufficient stability bound for the initial queue size caused by the overload is quite tight. Effective resource utilization as low as 20% cannot prevent an overloaded server from crash, if an initial queue created by a short-term overload (due to a demand burst or a temporary server slowdown) exceeds the sufficient stability bound slightly.

Event-driven simulation, adopted by most existing literature on SIP study, requires a series of retransmission timers to track outstanding messages, thus makes the experiment computationally expensive. As the network size increases to a large scale, the number of timers may build up to consume excessive memory and CPU time, thus crashes the simulator eventually. On the contrary, fluid-based simulation tracks time on a slot-by-slot basis. Events happening within the same time slot will be aggregated and processed together. Individual timers do not need to be tracked anymore. Thus fluid-based approach is much simpler than event-driven approach. The similarity between fluid-based Matlab simulation result and event-driven OPNET simulation result demonstrate that fluid-based simulation can be a relatively accurate and cost-effective approach for evaluating the performance of a SIP network.

ACKNOWLEDGMENT

We appreciate the financial support from the NSERC grant #CRDPJ 354729-07 and the OCE grant #CA-ST-150764-8. This work is also supported in part by Nortel Networks.

REFERENCES

- [1] J. Rosenberg et al., "SIP: Session Initiation Protocol," RFC 3261, IETF, June 2002.
- [2] J. Rosenberg and H. Schulzrinne, "SIP: Locating SIP Servers," RFC 3263, IETF, June 2002.
- [3] 3GPP TS 24.228 v5.f.0 (2006-10), "Signaling flows for the IP Multimedia call control based on SIP and SDP; Stage 3 (Release 5)," October 2006.
- [4] 3GPP TS 24.229 v8.5.1 (2008-09), "IP Multimedia call control protocol based on SIP and SDP; Stage 3 (Release 8)," September 2008.
- [5] J. Rosenberg and H. Schulzrinne, "Reliability of provisional responses in the Session Initiation Protocol (SIP)," RFC 3262, IETF, June 2002.
- [6] M. Govind, S. Sundaragopalan, Binu K S, and Subir Saha, "Retransmission in SIP over UDP - Traffic Engineering Issues," Proceedings of International Conference on Communication and Broadband Networking, Bangalore, May 2003.
- [7] E. Noel and C.R. Johnson, "Initial simulation results that analyze SIP based VoIP networks under overload," *Proceedings of 20th International Teletraffic Congress*, 2007, pp. 54-64.
- [8] E. Noel and C.R. Johnson, "Novel Overload Controls for SIP Networks," *Proceedings of 21st International Teletraffic Congress*, 2009.
- [9] R.P. Ejzak, C.K. Florkey, and R.W. Hemmeter, "Network Overload and Congestion: A comparison of ISUP and SIP," *Bell Labs Technical Journal*, 9(3), 2004, pp. 173-182.
- [10] M. Ohta, "Overload Control in a SIP Signaling Network," *Proceeding of World Academy of Science, Engineering and Technology*, Vienna, Austria, March 2006, pp. 205-210.
- [11] V. Hilt and I. Widjaja, "Controlling Overload in Networks of SIP Servers," *Proceedings of IEEE ICNP*, Orlando, Florida, October 2008, pp. 83-93.
- [12] C. Shen, H. Schulzrinne, and E. Nahum, "SIP Server Overload Control: Design and Evaluation," *Proceedings of IPTComm*, Heidelberg, Germany, July 2008.
- [13] A. Abdelal and W. Matragi, "Signal-Based Overload Control for SIP Servers," *Proceedings of IEEE CCNC*, Las Vegas, NV, January 2010.
- [14] "SIP Express Router" <http://www.iptel.org/ser/>.
- [15] T. Warabino, Y. Kishi, and H. Yokota, "Session Control Cooperating Core and Overlay Networks for "Minimum Core" Architecture," *Proceedings of IEEE Globecom*, Honolulu, Hawaii, December 2009.
- [16] Y. Hong, C. Huang, and J. Yan, "Analysis of SIP Retransmission Probability Using a Markov-Modulated Poisson Process Model," *Proceedings of IEEE/IFIP Network Operations and Management Symposium*, Osaka, Japan, April 2010.
- [17] E.M. Nahum, J. Tracey, and C.P. Wright, "Evaluating SIP server performance," *Proceedings of international conference on Measurement and modeling of computer systems (ACM SIGMETRICS)*, San Diego, CA, US, 2007, pp. 349-350.
- [18] J. Sun, R.X. Tian, J.F. Hu, and B. Yang, "Rate-based SIP Flow Management for SLA Satisfaction," *Proceedings of 11th International Symposium on Integrated Network Management (IFIP/IEEE IM)*, New York, USA, June 2009, pp. 125-128.
- [19] V. Hilt and H. Schulzrinne, "Session Initiation Protocol (SIP) Overload Control," *IETF Internet-Draft*, draft-hilt-sipping-overload-07, October 2009.
- [20] J. Rosenberg, "Requirements for Management of Overload in the Session Initiation Protocol," *IETF RFC 5390*, December 2008.
- [21] W. R. Stevens, *TCP/IP Illustrated*, Volume 1, Addison-Wesley, Boston, 1994.
- [22] Y. Hong, O. W. W. Yang, and C. C. Huang, "Self-Tuning PI TCP Flow Controller for AQM Routers With Interval Gain and Phase Margin Assignment," *Proceedings of IEEE Globecom*, Dallas, TX, U.S.A., November 2004, pp. 1324-1328.
- [23] Y. Hong, C. Huang, and J. Yan, J., "Modeling and Simulation of SIP Tandem Server with Finite Buffer," To appear in *ACM Transactions on Modeling and Computer Simulation*, April 2011.
- [24] Y. Liu, F. L. Presti, V. Misra, D. F. Towsley, and Y. Gu, "Scalable fluid models and simulations for large-scale IP networks," *ACM Transactions on Modeling and Computer Simulation*, 14 (3), pp. 305-324.