

SIP 프록시 서버를 위한 우선순위 스케줄링 기반의 과부하 제어 알고리즘

이장현*, 조인휘*

*한양대학교 전자컴퓨터통신공학과

e-mail: gms1987@hanyang.ac.kr, iwjoe@hanyang.ac.kr

An Overload Control Algorithm based on Priority Scheduling for SIP Proxy Server

Janghyun Lee *, Inwhae Joe *

*Dept of Electronics and Computer Engineering, Hanyang University

요 약

SIP (Session Initiation Protocol)는 사용자간의 멀티미디어 세션의 초기화, 변경 및 종료를 위한 응용 계층의 시그널링 프로토콜이다. SIP는 실시간 멀티미디어 서비스 통신에 많이 이용되기 때문에 주로 Unreliable Transport Protocol을 사용한다. 따라서 손실된 SIP 메시지들의 신뢰성을 보장하기 위하여 재전송 메커니즘을 제공한다. 하지만 이러한 재전송 메커니즘은 SIP 프록시 서버가 과부하 상황일 경우 신뢰성보다는 오버헤드가 증가되는 문제점을 가진다. 기존의 여러 SIP 과부하 제어 방법이 제안되었지만 네트워크 혼잡이 증가함에 따라 프록시 서버의 처리율이나 호 설정 시간의 지연 문제를 효율적으로 해결하지 못한다. 본 논문에서는 SIP 시그널링 네트워크에서 프록시 서버가 과부하 상황일 경우 낮은 호 설정 지연시간과 높은 처리율을 위해 우선순위 스케줄링 기반의 과부하 제어 알고리즘을 제안한다. 그리고 기존 과부하 제어에서 사용하는 알고리즘과 제안하는 알고리즘을 비교하여 보았다. 성능 평가 결과 부하에 따른 프록시 서버의 처리율과 호 설정 시간의 지연 측면에서 기존의 과부하 제어 알고리즘보다 향상됨을 보였다.

1. 서론

SIP (Session Initiation Protocol)는 사용자간의 멀티미디어 세션의 초기화, 변경 및 종료를 위한 응용 계층의 시그널링 프로토콜로 인터넷을 이용한 전화, 컨퍼런스, 인스턴트 메시징 등에 주로 널리 사용된다 [1]. SIP는 하위 계층의 전송프로토콜과 독립적이지만, 실시간 멀티미디어 서비스 통신에 많이 이용되기 때문에 주로 Unreliable Transport Protocol (UDP)을 사용한다 [5]. 따라서 손실된 SIP 메시지들의 신뢰성을 보장하기 위하여 RFC3261에서는 재전송 메커니즘을 제공한다 [3]. 하지만 이러한 재전송 메커니즘은 SIP 프록시 서버가 과부하 상황일 경우 신뢰성보다는 오버헤드가 증가되는 문제점을 가진다 [4]. 최근 인터넷 기반 멀티미디어 서비스의 급격한 증가로 SIP 프록시 서버의 과부하를 제어하는 알고리즘은 중요한 이슈로 연구되고 있다 [8].

현재 SIP 프로토콜은 과부하 제어를 위하여 503 (Service Unavailable) 응답 메시지를 사용한다. 기존의 여러 SIP 과부하 제어 방법이 제안되었지만 네트워크 부하에 따른 처리율 감소나 호 설정 시간의 지연 문제를 효율적으로 해결하지 못한다 [2]. 본 논문에서는 SIP 시그널링 네트워크에서 프록시 서버가 과부하 상황일 경우 낮은 호 설정 지연시간을 위하여 INVITE 메시지의 재전송이 적게 이루

어진 메시지에 가중치를 높이고, 프록시 서버의 높은 처리율을 위하여 INVITE 메시지보다 Non-INVITE 메시지에 가중치를 높이는 우선순위 스케줄링 기반의 과부하 제어 알고리즘을 제안한다. 본 논문의 구성은 다음과 같다. 2장에서는 관련 연구로서 기존의 호 설정을 위한 SIP의 기본적인 흐름과 503 응답메시지를 이용한 과부하 제어를 설명하고 3장에서는 본 문에서 제안하는 시스템의 설계 및 구현과 신호 처리 절차를 설명하고, 4장에서는 성능 평가를 통하여 논문에서 제안한 시스템에 대한 성능 분석을 한 후 5장에서 결론을 맺는다.

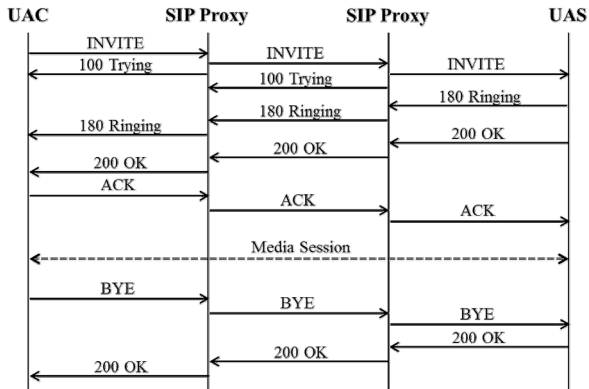
2. 관련연구

SIP의 구성요소는 기본적으로 SIP UAs (User Agents)와 SIP Server의 개체들로 구성된다. UAs는 UAC (User Agent Client)와 UAS (User Agent Server)로 구성되며 이들은 SIP 콜에 참여하는 클라이언트로 동작하는 요소들이다. SIP Server는 세션 라우팅을 위한 프록시 서버와 UA 등록을 위한 registration server로 구성된다. 그림 1은 SIP 호 세션을 설정하기 위한 SIP 메시지 흐름이다. UAC는 UAS와의 호 설정을 위해 Proxy Server에게 INVITE 요청 메시지를 전송한다. Proxy Server는 요청에 대한 응답으로 100 Trying 메시지로 응답함과 동시에 요

ACKNOWLEDGEMENT

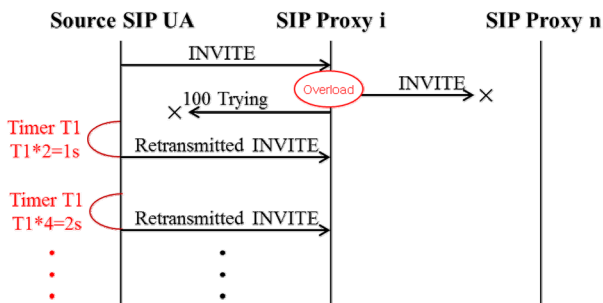
본 연구는 지식경제부 R&D 지원 프로그램의 일환으로, 한국산업기술평가관리원의 지원으로 수행되었음.

청에 대한 다음 홉을 결정하여 포워딩 한다. UAS는 INVITE에 대한 응답으로 180 Ringing 메시지와 200 OK 메시지를 SIP 프록시 서버를 통해 UAC에게 전달하고 UAC의 ACK 메시지에 의해 UAC와 UAS사이의 미디어 세션이 설정된다.



(그림 2) 세션 설정을 위한 절차

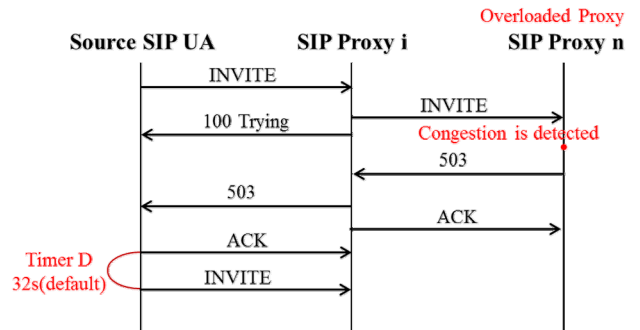
Unreliable Transport Protocol (UDP)에서 손실된 SIP 메시지들의 신뢰성을 보장하기 위하여 RFC3261에서는 재전송 메커니즘을 제공한다. INVITE 메시지가 손실 될 경우 UAC는 timer T1에 의하여 재전송 한다. Timer T1은 기본적인 RTT (round-trip time)값으로 500ms의 초기 값을 가지며 자신의 INVITE 메시지에 대한 100 Trying을 수신 할 때까지 $2 \times T1$, $4 \times T1$, $8 \times T1$, $16 \times T1$, $32 \times T1$, $64 \times T1$ 간격으로 재전송하게 된다. SIP 시그널링 네트워크에서 프록시 서버가 과부하일 경우 수신한 INVITE 메시지에 대한 100Trying 응답 메시지를 전송해주지 못한다. 따라서 UAC는 100Trying 메시지를 수신하지 못하여 자신이 보낸 INVITE 메시지가 손실되었다고 생각하여 Timer T1에 의한 재전송을 수행한다. Timer T1에 의한 재전송된 INVITE 메시지는 과부하 된 프록시 서버의 오버헤드를 증가시켜 신뢰성보다는 성능을 감소시킨다 [6].



(그림 2) INVITE 재전송으로 인한 과부하 문제

SIP 프로토콜은 SIP 프록시 서버가 일시적인 과부하 때문에 세션을 위한 요청 메시지를 포워딩할 수 없을때 과부하 제어를 위하여 503 (Service Unavailable) 응답 메시지를 사용한다. 과부하 된 프록시 서버는 업스트림 서버에

게 503 메시지를 전송한다. 503 메시지를 수신한 UAC는 TimerD를 수행한다. TimerD는 재시도 가능 타이머(Retry After timer) 로서 UDP의 경우 기본적으로 32초이다. 따라서 32초 동안은 새로운 INVITE 메시지에 대해 과부하된 프록시 서버에게 전달하지 않으므로 일시적으로 과부하를 해결 할 수 있다.

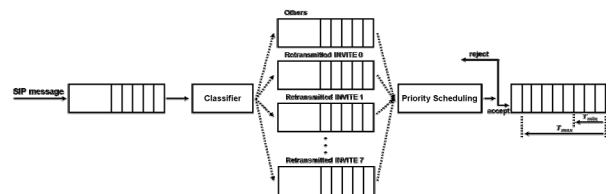


(그림 3) 503을 이용한 과부하 제어방법

[7]에서는 SIP 프록시 서버의 큐의 길이가 high watermark를 초과하면 과부하 상태로 판단하여 low watermark보다 작아 질 때까지 모든 SIP 메시지들을 503 응답 메시지를 통하여 거부하는 간단한 과부하 제어 알고리즘이 제안되었다. [9]에서는 SIP 프록시 서버의 큐의 길이에 비례하여 확률적으로 503 응답 메시지를 전송하여 기존 [9]의 알고리즘을 개선하였다. 하지만 모든 SIP 메시지에 대해서 503 응답 메시지를 랜덤으로 전송하게 되므로 프록시 서버의 처리율이나 호 설정 지연 측면에서 효율이 떨어진다.

3. 제안하는 알고리즘

SIP 시그널링 네트워크에서 프록시 서버는 SIP 메시지들을 응답하고 전달하는 역할을 한다. 호 설정을 위해서는 이 모든 메시지들이 성공적으로 전달되어야 한다. 하지만 네트워크 혼잡이 증가함에 따라 과부하 된 프록시 서버는 SIP 메시지를 정상적으로 응답하거나 전달하지 못하게 된다. 본 논문에서는 프록시 서버의 과부하 제어를 위하여 네트워크 혼잡시에 우선순위 스케줄링을 통하여 새로운 콜 요청을 시도하는 메시지를 503 응답 메시지로 거부한다. 그림 4는 SIP 시그널링 네트워크에서 과부하 제어를 위한 프록시 서버의 큐 구조이다. 프록시 서버에서 SIP 메시지를 수신하면 메시지들을 우선순위 스케줄링에 따라 분류한 뒤 실행 큐에서 실행하게 된다.



(그림 4) SIP 프록시 서버의 큐 구조

$$I_i = T1 \times (2^i - 1)$$

프록시 서버의 과부하로 인해 INVITE 메시지가 큐에서 유실 될 확률을 b라고 가정하였을 때, i번째 재전송된 INVITE메시지가 큐에서 처리 되었을 확률은 다음과 같다.

$$p_i = (1 - b) \times b^{i-1}$$

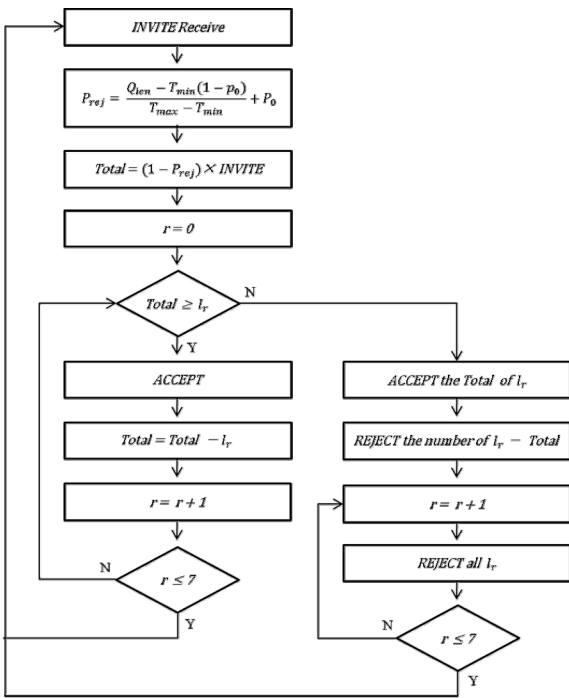
따라서 i번째 재전송된 INVITE 메시지로부터의 호 설정 지연 시간은 다음과 같다.

$$D = \sum_{i=1}^7 I_i \times P_i$$

본 논문에서는 호 설정 지연시간과 비례하는 재전송된 INVITE 메시지들에 대하여 우선순위 스케줄링을 한다.

$$\text{INVITE Priority} \propto \frac{1}{T1 \times (2^i - 1)}$$

따라서 네트워크 혼잡에 따라 증가하는 호 설정 지연시간을 줄일 수 있다.



(그림 5) INVITE 메시지 우선순위 스케줄링 흐름도

과부하 된 프록시 서버에서 SIP 메시지들을 수신 한 경우 INVITE 메시지는 Non-INVITE 메시지들에 비해 더 큰 처리 시간을 가지게 되므로 Non-INVITE 메시지에 대해 높은 우선순위 처리를 위해 가중치를 높인다. 그리고 INVITE 메시지들은 그림 5와 같이 작은 호 설정 지연 시간에 따라 재전송된 INVITE 메시지들을 분류하여 처리한다. 결국 실행되는 큐에서는 현재 네트워크 혼잡과 INVITE 대기시간, 처리 시간을 고려하여 처리하거나 503 응답 메시지를 통해 거절 한다. INVITE 메시지는 그림 6과 같다. 재전송 된 INVITE 메시지의 호 설정 지연 시간은 재전송 타이머와 비례하므로 INVITE 메시지의 CSeq를 이용하여 파악한다.

```
INVITE sip:bob@biloxi.com SIP/2.0
Via: SIP/2.0/UDP pc33.atlanta.com;branch=z9hG4bK776asdhds
Max-Forwards: 70
To: Bob <sip:bob@biloxi.com>
From: Alice <sip:alice@atlanta.com>;tag=1928301774
Call-ID: a84b4c76e66710@pc33.atlanta.com
CSeq: 314159 INVITE
Contact: <sip:alice@pc33.atlanta.com>
Content-Type: application/sdp
Content-Length: 142
```

(그림 6) INVITE 메시지

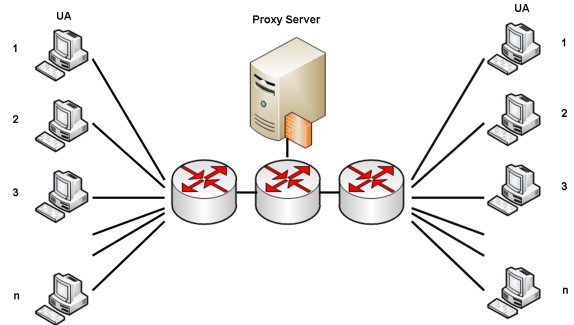
네트워크 혼잡이 증가함에 따라 프록시 서버가 INVITE 메시지를 처리할 수 없는 횟수가 비례하게 된다. 따라서 INVITE 메시지의 재전송도 이루어 질 것이다. 재전송횟수에 따라 INVITE 메시지의 대기 시간은 다음과 같이 증가하게 된다.

4. 성능평가

제안하는 알고리즘의 우수성을 판단하기 위해 기존의 [7], [9] 에서 제안된 과부하 제어 알고리즘에서 프록시 서버의 처리율과 호 설정 시간의 지연 측면에서 비교 분석하였다.

기존방식1: 프록시 서버의 큐를 이용한 과부하 제어 방식
 기존방식2: 프록시 서버의 큐를 이용한 확률적인 과부하 제어 방식

제안방식: 우선순위 스케줄링 기반의 과부하 제어 방식



(그림 7) INVITE 우선순위 스케줄링 흐름도

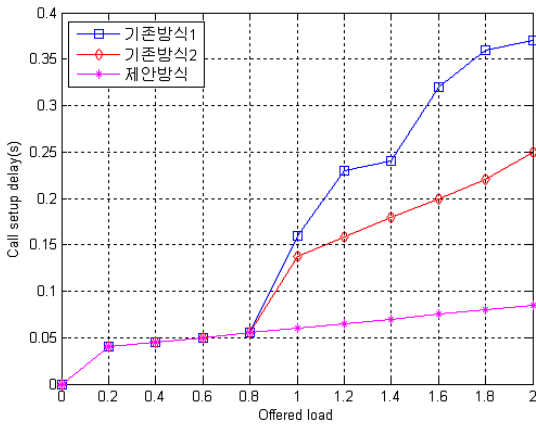
성능 평가를 위한 네트워크 모델은 그림 7과 같이 프록시 서버가 bottleneck 상황에서 고려하였다. SIP 메시지의 평균 도착률이 λ이고 프록시 서버의 처리율이 μ_{sip}일 때 프록시 서버의 인가된 부하(Offered load)는 다음과 같다.

$$\rho = \frac{\lambda}{\mu_{sip}}$$

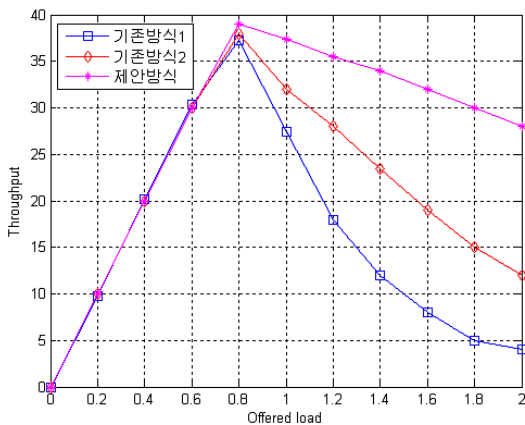
여기서 프록시 서버의 처리율은 정해진 상수 값 이므로 여기서는 n을 증가시키면서 프록시의 부하를 구할 수 있다. 그러므로 콜이 성공 했을 경우 프록시 서버에서 수락한 콜에 대한 호 설정 지연 시간은 다음과 같다.

$$\sum_{r=0}^7 \frac{(1 - P_{rsj}) \times I_r}{Total} \times T1 \times (2^n - 1)$$

그림 8은 프록시 서버의 부하에 따른 호 설정 지연 시간을 나타내고 있다. 부하가 1이하의 경우 프록시 서버로 수신한 SIP 메시지들보다 SIP 프록시 서버의 처리 할 수 있는 능력이 크므로 모든 SIP 메시지가 받아들여진다. 따라서 부하가 1이하의 경우는 기존 모든 알고리즘의 호 설정 지연 시간은 비슷하다. 하지만 Load 가 1이상 이 될 경우 기존 알고리즘의 경우 재전송된 INVITE 메시지에 대해 랜덤으로 거부하므로 제안한 INVITE 대기시간을 고려한 알고리즘 보다 호 설정 지연 시간이 크다.



(그림 8) 프록시 서버의 부하에 따른 호 설정 지연



(그림 9) 프록시 서버의 부하에 따른 처리율

그림 9는 프록시 서버의 부하에 따른 처리율을 나타낸다. 프록시 서버가 모든 메시지를 처리 할 수 없을 때 INVITE 메시지보다 Non-INVITE 메시지에 가중치를 두어 다음과 같이 처리율을 향상 시킬 수 있다.

5. 결과

SIP 시그널링 네트워크에서 기존의 과부하 제어 알고리즘을 분석한 결과 이전 알고리즘에 비해 네트워크 부하에 따른 프록시 서버의 처리율과 호 설정 시간의 지연 문제를 줄일 수 있었다. 본 논문에서는 네트워크 혼잡상태에서 호 설정 지연을 개선하기 위해 호 설정 지연이 적은 INVITE 메시지에 대해 우선순위를 높였다. 또한 프록시 서버의 처리율을 높이기 위해 INVITE 메시지보다 Non-INVITE 메시지에 대해 우선순위를 높여 스케줄링 하여 기존 알고리즘에 비해 성능을 향상 시켰다.

참고문헌

- [1] H. Schulzrinne and J. Rosenberg: "The Session Initiation Protocol: Internet-Centric Signaling", IEEE Communication Magazine, vol.38, 10, pp-134-141, Oct.(2000)
- [2] T.Eyers and H. Schulzrinne: "Predicting Internet Telephony Call Setup Delay" Proc. 1st IP-Telephony Wksp.,Jan.(2000)
- [3] Rosenberg, J.Schulzrinne, H.Camarillo, G.Johnston, A.Peterson, J.Sparks, R.Handley, M.,and E.Schooler, "SIP: Session Initiation Protocol"[S], RFC 3261, June 2002.
- [4] M. Govind, S. Sundaragopalan, Binu K S, and Subir Saha: "Retransmission in SIP over UDP - Traffic Engineering Issues", Proc. of International Conference on Communication and Broadband Networking, Bangalore, May(2003)
- [5] G. Camarillo, R. Kantola and H. Schulzrinne: "Evaluation of Transport Protocols for the Session Initiation Protocol", IEEE Network, Vol.17, 5,pp.40-46, Sep.(2003)
- [6] 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, pp.173-182(2004)
- [7] Masataka Ohta. Overload Control in a SIP Signaling Network. PROCEEDINGS OF WORLD ACADEMY OF SCIENCE, ENGINEERING AND TECHNOLOGY VOLUME 12 MARCH 2006
- [8] J. Rosenberg, "Requirements for Management of Overload in the Session Initiation Protocol"[S], RFC5390, December 2008
- [9] J. Yang, F. Huang, and S. Gou, An Optimized Algorithm for Overload Control of SIP signaling Network, 5th International Conference on Wireless Communications, Networking and Mobile Computing (WiCom), 2009. pp. 1 - 4.