

SOC Working Group
Internet-Draft
Intended status: Standards Track
Expires: January 22, 2015

Eric Noel
AT&T Labs
Philip M Williams
BT Innovate & Design

July 22, 2014

**Session Initiation Protocol (SIP) Rate Control
draft-ietf-soc-overload-rate-control-09.txt**

Abstract

The prevalent use of Session Initiation Protocol (SIP) in Next Generation Networks necessitates that SIP networks provide adequate control mechanisms to maintain transaction throughput by preventing congestion collapse during traffic overloads. Already a loss-based solution to remedy known vulnerabilities of the SIP 503 (service unavailable) overload control mechanism has been proposed. This document proposes a rate-based control scheme to complement the loss-based control scheme, using the same signaling.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on January 22, 2015.

Copyright Notice

Copyright (c) 2014 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents

(<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the [Trust Legal Provisions](#) and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
2. Terminology	3
3. Rate-based algorithm scheme	4
3.1. Overview	4
3.2. Via header field parameters for overload control	4
3.3. Client and server rate-control algorithm selection	5
3.4. Server operation	5
3.5. Client operation	6
3.5.1. Default algorithm	6
3.5.2. Priority treatment	10
3.5.3. Optional enhancement: avoidance of resonance	11
4. Example	12
5. Syntax	14
6. Security Considerations	14
7. IANA Considerations	14
8. References	14
8.1. Normative References	14
8.2. Informative References	15
Appendix A. Contributors	16
Appendix B. Acknowledgments	16

1. Introduction

The use of SIP in large scale Next Generation Networks requires that SIP based networks provide adequate control mechanisms for handling traffic growth. In particular, SIP networks must be able to handle traffic overloads gracefully, maintaining transaction throughput by preventing congestion collapse.

A promising SIP based overload control solution has been proposed in [[draft-ietf-soc-overload-control-15](#)]. That solution provides a communication scheme for overload control algorithms. It also includes a default loss-based overload control algorithm that makes it possible for a set of clients to limit offered load towards an overloaded server.

However, such loss control algorithm is sensitive to variations in load so that any increase in load would be directly reflected by the clients in the offered load presented to the overloaded servers. More importantly, a loss-based control cannot guarantee clients to produce a bounded offered load from the clients towards an overloaded server and requires frequent updates which may have implications for stability.

In accordance with the framework defined in [[draft-ietf-soc-overload-control-15](#)], this document proposes an alternate overload control, the rate-based overload control algorithm. The rate-based control guarantees an upper bound on the rate, constant between server updates, of requests sent by clients towards an overloaded server. The tradeoff is in terms of algorithmic complexity, since the overloaded server is more likely to use a different target (maximum rate) for each client, than the loss-based approach.

The proposed rate-based overload control algorithm mitigates congestion in SIP networks while adhering to the overload signaling scheme in [[draft-ietf-soc-overload-control-15](#)] and presenting a rate based control as an optional alternative to the default loss-based control in [[draft-ietf-soc-overload-control-15](#)].

2. Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#) [[RFC2119](#)].

The normative statements in this specification as they apply to clients and servers assume that both the clients and servers support this specification. If, for instance, only a client supports this specification and not the server, then follows that the normative statements in this specification pertinent to the behavior of a server do not apply to the server that does not support this specification.

3. Rate-based algorithm scheme

3.1. Overview

The server is the one protected by the overload control algorithm defined here, and the client is the one that throttles traffic towards the server.

Following the procedures defined in [[draft-ietf-soc-overload-control-15](#)], the server and clients signal one another support for rate-based overload control.

Then periodically, the server relies on internal measurements (e.g. CPU utilization, queueing delay...) to evaluate its overload state and estimate a target SIP request rate in number of request per second (as opposed to target percent loss in the case of loss-based control).

When in overload, the server uses the Via header field oc parameters [[draft-ietf-soc-overload-control-15](#)] of SIP responses in order to inform the clients of its overload state and of the target SIP request rate.

Upon receiving the oc parameters with a target SIP request rate, each client throttles new SIP requests towards the overloaded server.

3.2. Via header field parameters for overload control

The use of the via header oc parameter(s) inform clients of the desired maximum rate. They are defined in [[draft-ietf-soc-overload-control-15](#)] and summarized below:

oc: Used by clients in SIP requests to indicate [[draft-ietf-soc-overload-control-15](#)] support and by servers to indicate the load reduction amount.

oc-algo: Used by clients in SIP requests to advertise supported overload control algorithms and by servers to notify clients of the algorithm in effect. Supported values: loss (default), rate (optional).

oc-validity: Used by servers in SIP responses to indicate an interval of time (msec) that the load reduction should be in effect. A value of 0 is reserved for server to stop overload control. A non-zero value is required in conjunction for all other cases.

oc-seq: A sequence number associated with the "oc" parameter.

3.3. Client and server rate-control algorithm selection

Per [[draft-ietf-soc-overload-control-15](#)], new clients indicate supported overload control algorithms to servers by inserting oc and oc-algo, with the names of the supported algorithms, in Via header of SIP requests destined to servers. The inclusion by the client of the token "rate" indicates that the client supports a rate based algorithm. Conversely, servers notify clients of the selected overload control algorithm through the oc-algo parameter in the Via header of SIP responses to clients. The inclusion by the server of the token "rate" in the oc-algo parameter indicates that the rate based algorithm has been selected by the server.

Support of rate-based control MUST be indicated by clients setting oc-algo to the token "rate". Selection of rate-based control MUST be indicated by servers by setting oc-algo to the token "rate".

3.4. Server operation

The actual algorithm used by the server to determine its overload state and estimate a target SIP request rate is beyond the scope of this document.

However, the server MUST periodically evaluate its overload state and estimate a target SIP request rate beyond which it would become overloaded. The server must allocate a portion of the target SIP request rate to each of its client. The server may set the same rate for every client, or may set different rates for different clients.

The max rate determined by the server for a client applies to the entire stream of SIP requests, even though throttling may only affect a particular subset of the requests, since as per [draft-

ietf-soc-overload-control-15] and REQ 13 of [RFC 5390](#), request prioritization is the client responsibility.

When setting the maximum rate for a particular client, the server may need take into account the workload (e.g. cpu load per request) of the distribution of message types from that client. Furthermore, because the client may prioritize the specific types of messages it sends while under overload restriction, this distribution of message types may be different from (e.g., either higher or lower cpu load) the message distribution for that client under non-overload conditions

Note that the oc parameter for the rate algorithm is an upper bound (in messages per second) on the traffic sent by the client to the server. The client may send traffic at a rate significantly lower than the upper bound, for a variety of reasons

In other words, when multiple clients are being controlled by an overloaded server, at any given time some clients may receive requests at a rate below their target (maximum) SIP request rate while others above that target rate. But the resulting request rate presented to the overloaded server will converge towards the target SIP request rate.

Upon detection of overload, and the determination to invoke overload controls, the server MUST follow the specifications in [[draft-ietf-soc-overload-control-15](#)] to notify its clients of the allocated target SIP request rate.

The server MUST use [[draft-ietf-soc-overload-control-15](#)] oc parameter to send a target SIP request rate to each of its clients.

3.5. Client operation

3.5.1. Default algorithm

In determining whether or not to transmit a specific message, the client may use any algorithm that limits the message rate to 1/T messages per second. It may be strictly deterministic, or it may be

probabilistic. It may, or may not, have a tolerance factor, to allow for short bursts, as long as the long term rate remains below $1/T$. The algorithm may have provisions for prioritizing traffic in accordance with REQ 13 of [RFC5390](#).

If the algorithm requires other parameters (in addition to "T", which is $1/oc$), they may be set autonomously by the client, or they may be negotiated independently between client and server.

In either case, the coordination is out of scope for this document. The default algorithms presented here (one without provisions for prioritizing traffic, one with) are only examples. Other algorithms that forward messages in conformance with the upper bound of $1/T$ messages per second may be used

To throttle new SIP requests at the rate specified in the oc value sent by the server to its clients, the client MAY use the proposed default algorithm for rate-based control or any other equivalent algorithm.

The default Leaky Bucket algorithm presented here is based on [ITU-T Rec. I.371] [Appendix A.2](#). The algorithm makes it possible for clients to deliver SIP requests at a rate specified in the oc value with tolerance parameter TAU (preferably configurable).

Conceptually, the Leaky Bucket algorithm can be viewed as a finite capacity bucket whose real-valued content drains out at a continuous rate of 1 unit of content per time unit and whose content increases by the increment T for each forwarded SIP request. T is computed as the inverse of the rate specified in the oc value, namely $T = 1 / oc\text{-value}$.

Note that when the oc -value is 0 with a non-zero oc -validity, then the client should reject 100% of SIP requests destined to the overload server. However, when the oc -validity value is 0, the client should immediately stop throttling.

If, at a new SIP request arrival, the content of the bucket is less than or equal to the limit value TAU , then the SIP request is forwarded to the server; otherwise, the SIP request is rejected.

Note that the capacity of the bucket (the upper bound of the counter) is $(T + \text{TAU})$.

The tolerance parameter TAU determines how close the long-term admitted rate is to an ideal control that would admit all SIP requests for arrival rates less than $1/T$ and then admit SIP requests precisely at the rate of $1/T$ for arrival rates above $1/T$. In particular at mean arrival rates close to $1/T$, it determines the tolerance to deviation of the inter-arrival time from T (the larger TAU the more tolerance to deviations from the inter-departure interval T).

This deviation from the inter-departure interval influences the admitted rate burstiness, or the number of consecutive SIP requests forwarded to the server (burst size proportional to TAU over the difference between $1/T$ and the arrival rate).

Servers with a very large number of clients, each with a relatively small arrival rate, will generally benefit from a smaller value for TAU in order to limit queuing (and hence response times) at the server when subjected to a sudden surge of traffic from all clients. Conversely, a server with a relatively small number of clients, each with proportionally larger arrival rate, will benefit from a larger value of TAU.

Once the control has been activated, at the arrival time of the k -th new SIP request, $ta(k)$, the content of the bucket is provisionally updated to the value

$$X' = X - (ta(k) - LCT)$$

where X is the value of the leaky bucket counter after arrival of the last forwarded SIP request, and LCT is the time at which the last SIP request was forwarded.

If X' is less than or equal to the limit value TAU, then the new SIP request is forwarded and the leaky bucket counter X is set to X' (or to 0 if X' is negative) plus the increment T , and LCT is set to the current time $ta(k)$. If X' is greater than the limit value TAU, then the new SIP request is rejected and the values of X and LCT are unchanged.

When the first response from the server has been received indicating control activation ($oc_validity > 0$), LCT is set to the time of activation, and the leaky bucket counter is initialized to the parameter TAU_0 (preferably configurable) which is 0 or larger but less than or equal to TAU.

TAU can assume any positive real number value and is not necessarily bounded by T.

$TAU = 4 * T$ is a reasonable compromise between burst size and throttled rate adaptation at low offered rate.

Note that specification of a value for TAU, and any communication or coordination between servers, is beyond the scope of this document.

A reference algorithm is shown below.

No priority case:

```
// T: inter-transmission interval, set to 1 / TargetRate
// TAU: tolerance parameter
// ta: arrival time of the most recent arrival
// LCT: arrival time of last SIP request that was sent to the server
//      (initialized to the first arrival time)
// X: current value of the leaky bucket counter (initialized to
//    TAU0)

// After most recent arrival, calculate auxiliary variable Xp
Xp = X - (ta - LCT);

if (Xp <= TAU) {
    // Transmit SIP request
    // Update X and LCT
    X = max (0, Xp) + T;
    LCT = ta;
} else {
    // Reject SIP request
    // Do not update X and LCT
}
```

3.5.2. Priority treatment

As with the loss-based algorithm of [[draft-ietf-soc-overload-control-15](#)], a client implementing the rate-based algorithm also prioritizes messages into two or more categories of requests: Requests candidate for reduction, requests not subject to reduction (except under extenuating circumstances when there aren't any messages in the first category that can be reduced).

Accordingly, the proposed Leaky bucket implementation is modified to support priority using two thresholds for SIP requests in the set of requests candidate for reduction. With two priorities, the proposed Leaky bucket requires two thresholds $TAU1 < TAU2$:

- . All new requests would be admitted when the leaky bucket counter is at or below $TAU1$,
- . Only higher priority requests would be admitted when the leaky bucket counter is between $TAU1$ and $TAU2$,
- . All requests would be rejected when the bucket counter is above $TAU2$.

This can be generalized to n priorities using n thresholds for $n > 2$ in the obvious way.

With a priority scheme that relies on two tolerance parameters ($TAU2$ influences the priority traffic, $TAU1$ influences the non-priority traffic), always set $TAU1 \leq TAU2$ (TAU is replaced by $TAU1$ and $TAU2$). Setting both tolerance parameters to the same value is equivalent to having no priority. $TAU1$ influences the admitted rate the same way as TAU does when no priority is set. And the larger the difference between $TAU1$ and $TAU2$, the closer the control is to strict priority queueing.

$TAU1$ and $TAU2$ can assume any positive real number value and is not necessarily bounded by T .

Reasonable values for $TAU0$, $TAU1$ & $TAU2$ are: $TAU0 = 0$, $TAU1 = 1/2 * TAU2$ and $TAU2 = 10 * T$.

Note that specification of a value for $TAU1$ and $TAU2$, and any communication or coordination between servers, is beyond the scope of this document.

A reference algorithm is shown below.

Priority case:

```
// T: inter-transmission interval, set to 1 / TargetRate
// TAU1: tolerance parameter of no priority SIP requests
// TAU2: tolerance parameter of priority SIP requests
// ta: arrival time of the most recent arrival
// LCT: arrival time of last SIP request that was sent to the server
//      (initialized to the first arrival time)
// X: current value of the leaky bucket counter (initialized to
//    TAU0)

// After most recent arrival, calculate auxiliary variable Xp
Xp = X - (ta - LCT);

if (AnyRequestReceived && Xp <= TAU1) || (PriorityRequestReceived &&
Xp <= TAU2 && Xp > TAU1) {
    // Transmit SIP request
    // Update X and LCT
    X = max (0, Xp) + T;
    LCT = ta;
} else {
    // Reject SIP request
    // Do not update X and LCT
}
```

3.5.3. Optional enhancement: avoidance of resonance

As the number of client sources of traffic increases and the throughput of the server decreases, the maximum rate admitted by each client needs to decrease, and therefore the value of T becomes larger. Under some circumstances, e.g. if the traffic arises very quickly simultaneously at many sources, the occupancies of each bucket can become synchronized, resulting in the admissions from each source being close in time and batched or very 'peaky' arrivals at the server, which not only gives rise to control instability, but also very poor delays and even lost messages. An appropriate term for this is 'resonance' [[Erramilli](#)].

If the network topology is such that this can occur, then a simple way to avoid this is to randomize the bucket occupancy at two

appropriate points: At the activation of control, and whenever the bucket empties, as follows.

After updating the value of the leaky bucket to X' , generate a value u as follows:

if $X' > 0$, then $u=0$

else if $X' \leq 0$ then uniformly distributed between $-1/2$ and $+1/2$

Then (only) if the arrival is admitted, increase the bucket by an amount $T + uT$, which will therefore be just T if the bucket hadn't emptied, or lie between $T/2$ and $3T/2$ if it had.

This randomization should also be done when control is activated, i.e. instead of simply initializing the leaky bucket counter to TAU_0 , initialize it to $\text{TAU}_0 + uT$, where u is uniformly distributed as above. Since activation would have been a result of response to a request sent by the client, the second term in this expression can be interpreted as being the bucket increment following that admission.

This method has the following characteristics:

- . If TAU_0 is chosen to be equal to TAU and all sources were to activate control at the same time due to an extremely high request rate, then the time until the first request admitted by each client would be uniformly distributed over $[0, T]$;
- . The maximum occupancy is $\text{TAU} + (3/2)T$, rather than $\text{TAU} + T$ without randomization;
- . For the special case of 'classic gapping' where $\text{TAU}=0$, then the minimum time between admissions is uniformly distributed over $[T/2, 3T/2]$, and the mean time between admissions is the same, i.e. $T+1/R$ where R is the request arrival rate;
- . At high load randomization rarely occurs, so there is no loss of precision of the admitted rate, even though the randomized 'phasing' of the buckets remains.

4. Example

Adapting [[draft-ietf-soc-overload-control-15](#)] example in [section 6.2](#) where client P1 sends requests to a downstream server P2:

```
INVITE sips:user@example.com SIP/2.0

Via: SIP/2.0/TLS p1.example.net;

    branch=z9hG4bK2d4790.1;received=192.0.2.111;

    oc;oc-algo="loss,rate"

...

SIP/2.0 100 Trying

Via: SIP/2.0/TLS p1.example.net;

    branch=z9hG4bK2d4790.1;received=192.0.2.111;

    oc=0;oc-algo="rate";oc-validity=0;

    oc-seq=1282321615.781

...
```

In the messages above, the first line is sent by P1 to P2. This line is a SIP request; because P1 supports overload control, it inserts the "oc" parameter in the topmost Via header that it created. P1 supports two overload control algorithms: loss and rate.

The second line --- a SIP response --- shows the top most Via header amended by P2 according to this specification and sent to P1. Because P2 also supports overload control, it chooses the "rate" based scheme and sends that back to P1 in the "oc-algo" parameter. It uses oc-validity=0 to indicate no overload control.

At some later time, P2 starts to experience overload. It sends the following SIP message indicating P1 should send SIP requests at a rate no greater than or equal to 150 SIP requests per seconds and for a duration of 1,000 msec.

```
SIP/2.0 180 Ringing

Via: SIP/2.0/TLS p1.example.net;

    branch=z9hG4bK2d4790.1;received=192.0.2.111;

    oc=150;oc-algo="rate";oc-validity=1000;

    oc-seq=1282321615.782
```

...

5. Syntax

This specification extends the existing definition of the Via header field parameters of [[RFC3261](#)] as follows:

algo-value /= "rate"

6. Security Considerations

Refer to [[draft-ietf-soc-overload-control-15](#)] Security Considerations section.

7. IANA Considerations

Header Field Parameter Name Predefined Values Reference

Via	oc-algo	Yes	RFCABCD RFCOPRQ
-----	---------	-----	-----------------

RFCABCD RFCOPRQ [NOTE TO RFC-EDITOR: Please replace with final RFC number of [draft-ietf-soc-overload-rate-control](#) & [draft-ietf-soc-overload-control](#)]

8. References

8.1. Normative References

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

[RFC3261] Rosenberg, J., Schulzrinne, H., Camarillo, G., Johnston, A., Peterson, J., Sparks, R., Handley, M., and E. Schooler, "SIP: Session Initiation Protocol", [RFC 3261](#), June 2002.

[[draft-ietf-soc-overload-control-15](#)]

Gurbani, V., Hilt, V., Schulzrinne, H., "Session Initiation Protocol (SIP) Overload Control", [draft-ietf-soc-overload-control-15](#).

8.2. Informative References

[ITU-T Rec. I.371]

"Traffic control and congestion control in B-ISDN", ITU-T Recommendation I.371.

[Erramilli]

A. Erramilli and L. J. Forys, "Traffic Synchronization Effects In Teletraffic Systems", ITC-13, 1991.

Appendix A. **Contributors**

Significant contributions to this document were made by Janet Gunn.

Appendix B. **Acknowledgments**

Many thanks for comments and feedback on this document to: Keith Drage, Vijay Gurbany, Volker Hilt, Christer Holmberg, Winston Hong, James Yu.

This document was prepared using 2-Word-v2.0.template.dot.

Authors' Addresses

Eric Noel
AT&T Labs
200s Laurel Avenue
Middletown, NJ, 07747
USA

Philip M Williams
BT Innovate & Design
Ipswich, IP5 3RE
UK