

Actions Addressing Identified Issues with the
Session Initiation Protocol's (SIP) Non-INVITE Transaction

Status of This Memo

This document specifies an Internet standards track protocol for the Internet community, and requests discussion and suggestions for improvements. Please refer to the current edition of the "Internet Official Protocol Standards" (STD 1) for the standardization state and status of this protocol. Distribution of this memo is unlimited.

Copyright Notice

Copyright (C) The Internet Society (2006).

Abstract

This document describes modifications to the Session Initiation Protocol (SIP) to address problems that have been identified with the SIP non-INVITE transaction. These modifications reduce the probability of messages losing the race condition inherent in the non-INVITE transaction and reduce useless network traffic. They also improve the robustness of SIP networks when elements stop responding. These changes update behavior defined in RFC 3261.

Table of Contents

1. Introduction	2
2. Improving the Situation When Responses Are Only Delayed	2
2.1. Action 1: Make the best use of provisional responses	2
2.2. Action 2: Remove the useless late-response storm	3
3. Improving the Situation When an Element Is Not Going to Respond	4
4. Normative Updates to RFC 3261	4
4.1. Action 1	4
4.2. Action 2	5
5. Security Considerations	5
6. Contributors	5
7. Normative References	6

1. Introduction

There are a number of unpleasant edge conditions created by the SIP non-INVITE transaction (NIT) model's fixed duration. The negative aspects of some of these are exacerbated by the effect that provisional responses have on the non-INVITE transaction state machines. These problems are documented in [3]. In summary:

A non-INVITE transaction must complete immediately or risk losing a race

Losing the race will cause the requester to stop sending traffic to the responder (the responder will be temporarily blacklisted)

Provisional responses can delay recovery from lost final responses

The 408 response is useless for the non-INVITE transaction

As non-INVITE transactions through N proxies time-out, there can be an $O(N^2)$ storm of the useless 408 responses

This document specifies updates to RFC 3261 [1] to improve the behavior of SIP elements when these edge conditions arise.

2. Improving the Situation When Responses Are Only Delayed

There are two goals to achieve when we constrain the problem to those cases where all elements are ultimately responsive and networks ultimately deliver messages:

- o Reduce the probability of losing the race, preferably to the point that it is negligible
- o Reduce or eliminate useless messaging

2.1. Action 1: Make the best use of provisional responses

- o Disallow non-100 provisionals to non-INVITE requests
- o Disallow 100 Trying to non-INVITE requests before Timer E reaches T2 (for UDP hops)
- o Allow 100 Trying after Timer E reaches T2 (for UDP hops)
- o Allow 100 Trying for hops over reliable transports

Since non-INVITE transactions must complete rapidly ([3]), any information beyond "I'm here" (which can be provided by a 100 Trying) can be just as usefully delayed to the final response. Sending non-100 provisionals wastes bandwidth.

As shown in [3], sending any provisional response inside a NIT before Timer E reaches T2 damages recovery from failure of an unreliable transport.

Without a provisional, a late final response is the same as no response at all and will likely result in blacklisting the late-responding element ([3]). If an element is delaying its final response at all, sending a 100 Trying after Timer E reaches T2 prevents this blacklisting without damaging recovery from unreliable transport failure.

Blacklisting on a late response occurs even over reliable transports. Thus, if an element processing a request received over a reliable transport is delaying its final response at all, sending a 100 Trying well in advance of the timeout will prevent blacklisting. Sending a 100 Trying immediately will not harm the transaction as it would over UDP, but a policy of always sending such a message results in unnecessary traffic. A policy of sending a 100 Trying after the period of time in which Timer E reaches T2 had this been a UDP hop is one reasonable compromise.

2.2. Action 2: Remove the useless late-response storm

- o Disallow 408 to non-INVITE requests
- o Absorb stray non-INVITE responses at proxies

A 408 to non-INVITE will always arrive too late to be useful ([3]), The client already has full knowledge of the timeout. The only information this message would convey is whether or not the server believed the transaction timed out. However, with the current design of the NIT, a client cannot do anything with this knowledge. Thus, the 408 is simply wasting network resources and contributes to the response bombardment illustrated in [3].

Late non-INVITE responses by definition arrive after the client transaction's Timer F has fired and the client transaction has entered the Terminated state. Thus, these responses cannot be distinguished from strays. Changing the protocol behavior to prohibit forwarding non-INVITE stray responses stops the late-response storm. It also improves the proxy's defenses against malicious users counting on the RFC 3261 requirement to forward such strays.

3. Improving the Situation When an Element Is Not Going to Respond

When we expand the scope of the problem to also deal with element or network failure, we have more goals to achieve:

- o Identifying when an element is non-responsive
- o Minimizing or eliminating falsely identifying responsive elements as non-responsive
- o Avoiding non-responsive elements with future requests

Action 1 helps with the first two goals, dramatically improving an element's ability to distinguish between failure and delayed response from the next downstream element. Some response, either provisional or final, will almost certainly be received before the transaction times out. So, an element can more safely assume that no response at all indicates that the peer is not available and follow the existing requirements in [1] and [2] for that case.

Achieving the third goal requires more aggressive changes to the protocol. As noted in [3], future non-INVITE transactions are likely to fail again unless the implementation takes steps beyond what is defined in [1] and [2] to remember non-responsive destinations between transactions. Standardizing these extra steps is left to future work.

4. Normative Updates to RFC 3261

4.1. Action 1

An SIP element **MUST NOT** send any provisional response with a Status-Code other than 100 to a non-INVITE request.

An SIP element **MUST NOT** respond to a non-INVITE request with a Status-Code of 100 over any unreliable transport, such as UDP, before the amount of time it takes a client transaction's Timer E to be reset to T2.

An SIP element **MAY** respond to a non-INVITE request with a Status-Code of 100 over a reliable transport at any time.

Without regard to transport, an SIP element **MUST** respond to a non-INVITE request with a Status-Code of 100 if it has not otherwise responded after the amount of time it takes a client transaction's Timer E to be reset to T2.

4.2. Action 2

A transaction-stateful SIP element MUST NOT send a response with Status-Code of 408 to a non-INVITE request. As a consequence, an element that cannot respond before the transaction expires will not send a final response at all.

A transaction-stateful SIP proxy MUST NOT send any response to a non-INVITE request unless it has a matching server transaction that is not in the Terminated state. As a consequence, this proxy will not forward any "late" non-INVITE responses.

5. Security Considerations

This document makes a number of small changes to the core SIP specification [1] to improve the robustness of SIP non-INVITE transactions. Many of these actions also prevent flooding and denial-of-service attacks.

One change prohibits proxies and user agents from sending 408 responses to non-INVITE transactions. Without this change, proxies automatically generate a storm of useless responses as described in [3]. An attacker could capitalize on this by enticing user agents to send non-INVITE requests to a black hole (through social engineering or DNS poisoning) or by selectively dropping responses.

Another change prohibits proxies from forwarding late responses. Without this change, an attacker could easily forge messages that appear to be late responses. All proxies compliant with RFC 3261 are required to forward these responses, wasting bandwidth and CPU and potentially overwhelming target user agents (especially those with low-speed connections).

The remainder of these changes do not affect the security of the SIP protocol.

6. Contributors

Rohan Mahy provided the Security Considerations section.

7. Normative References

- [1] 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.
- [2] Rosenberg, J. and H. Schulzrinne, "Session Initiation Protocol (SIP): Locating SIP Servers", RFC 3263, June 2002.
- [3] Sparks, R., "Problems Identified Associated with the Session Initiation Protocol's (SIP) Non-INVITE Transaction", RFC 4321, January 2006.

Author's Address

Robert J. Sparks
Estacado Systems
17210 Campbell Road
Suite 250
Dallas, TX 75252-4203

EMail: rjsparks@estacado.net

Full Copyright Statement

Copyright (C) The Internet Society (2006).

This document is subject to the rights, licenses and restrictions contained in BCP 78, and except as set forth therein, the authors retain all their rights.

This document and the information contained herein are provided on an "AS IS" basis and THE CONTRIBUTOR, THE ORGANIZATION HE/SHE REPRESENTS OR IS SPONSORED BY (IF ANY), THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Intellectual Property

The IETF takes no position regarding the validity or scope of any Intellectual Property Rights or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; nor does it represent that it has made any independent effort to identify any such rights. Information on the procedures with respect to rights in RFC documents can be found in BCP 78 and BCP 79.

Copies of IPR disclosures made to the IETF Secretariat and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the IETF on-line IPR repository at <http://www.ietf.org/ipr>.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights that may cover technology that may be required to implement this standard. Please address the information to the IETF at ietf-ipr@ietf.org.

Acknowledgement

Funding for the RFC Editor function is provided by the IETF Administrative Support Activity (IASA).

