

Network Working Group
Internet-Draft
Intended status: Standards Track
Expires: April 27, 2015

Y. Nishida
GE Global Research
P. Natarajan
Cisco Systems
A. Caro
BBN Technologies
P. Amer
University of Delaware
K. Nielsen
Ericsson
October 24, 2014

Quick Failover Algorithm in SCTP
draft-ietf-tsvwg-sctp-failover-08.txt

Abstract

One of the major advantages of SCTP is that it supports multi-homed communication. A multi-homed SCTP end-point has the ability to withstand network failures by migrating the traffic from an inactive network to an active one. However, if the [\[RFC4960\]](#) specified failover operation is followed there can be a significant delay in the migration to the active destination addresses, thus severely reducing the effectiveness of SCTP multi-homed operation.

The memo complements [RFC4960](#) by the introduction of the Potentially Failed state and associated new Quick Failover operation to apply during network failure and specifies for SCTP senders to support this more performance optimal failover procedure as an add-on to the [\[RFC4960\]](#) failover operation. The memo in addition complements [\[RFC4960\]](#) by introduction of alternative switchover operation modes for the data transfer path management after a failover. These operation modes offer for more performance optimal operation in some network environments. From the perspective of this memo the implementation of the additional switchover operation modes is considered optional.

The procedures defined require only minimal modifications to the current specification. The procedures are sender-side only and do not impact the SCTP receiver.

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 April 27, 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	3
2.	Conventions and Terminology	4
3.	Issues with the SCTP Path Management	4
4.	SCTP with Potentially-Failed Destination State (SCTP-PF) . .	5
4.1.	SCTP-PF Concept	5
4.2.	SCTP-PF Algorithm Detail	6
4.3.	Optional Feature: Permanent Failover	9
5.	Socket API Considerations	10
5.1.	Support for the Potentially Failed Path State	11
5.2.	Peer Address Thresholds (SCTP_PEER_ADDR_THLDS) Socket Option	12
5.3.	Exposing the Potentially Failed Path State (SCTP_EXPOSE_POTENTIALLY_FAILED_STATE) Socket Option . .	13
6.	Security Considerations	13
7.	IANA Considerations	14
8.	Proposed Change of Status (to be Deleted before Publication)	14
9.	References	14
9.1.	Normative References	14
9.2.	Informative References	14

Appendix A . Discussions of Alternative Approaches	15
A.1 . Reduce Path.Max.Retrans (PMR)	15
A.2 . Adjust RTO related parameters	16
Appendix B . Discussions for Path Bouncing Effect	17
Authors' Addresses	17

1. Introduction

The Stream Control Transmission Protocol (SCTP) as specified in [\[RFC4960\]](#) supports multihoming at the transport layer -- an SCTP association can bind to multiple IP addresses at each endpoint. SCTP's multihoming features include failure detection and failover procedures to provide network interface redundancy and improved end-to-end fault tolerance.

In SCTP's current failure detection procedure, the sender must experience Path.Max.Retrans (PMR) number of consecutive failed retransmissions on a destination before detecting a path failure. The sender fails over to an alternate active destination only after failure detection. Until detecting the failover, the sender continues to transmit data on the failed path, which degrades the SCTP performance. Concurrent Multipath Transfer (CMT) [\[IYENGAR06\]](#) is an extension to SCTP and allows the sender to transmit data on multiple paths simultaneously. Research [\[NATARAJAN09\]](#) shows that the current failure detection procedure worsens CMT performance during failover and can be significantly improved by employing a better failover algorithm.

This document specifies an alternative failure detection procedure for SCTP that improves the SCTP performance during a failover.

Also the operation after a failover impacts the performance of the protocol. With [\[RFC4960\]](#) procedures, SCTP will, after a failover from the primary path, switch back to use the primary path for data transfer as soon as this path becomes available. From a performance perspective, as confirmed in research [\[CAR002\]](#), such a switchback of the data transmission path is not optimal in general. As an optional alternative to the switchback operation of [\[RFC4960\]](#), this document specifies for SCTP to support the Permanent Failover switchover procedures proposed by [\[CAR002\]](#). Additional discussions for alternative approach that does not require modifications to [\[RFC4960\]](#) and path bouncing effects that might be caused by frequent switchover are provided in Appendix.

2. Conventions and 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 [[RFC2119](#)].

3. Issues with the SCTP Path Management

This section describes issues in the current SCTP to be fixed by the approach described in this document.

SCTP can utilize multiple IP addresses for a single SCTP association. Each SCTP endpoint exchanges the list of its usable addresses during initial negotiation with its peer. Then the endpoints select one address from the peer's list and define this as the primary destination. During normal transmission, SCTP sends all user data to the primary destination. Also, it sends heartbeat packets to all idle destinations at a certain interval to check the reachability of the path. Idle destinations normally include all non-primary destinations.

If a sender has multiple active destination addresses, it can retransmit data to secondary destination address, when the transmission to the primary times out.

When a sender receives an acknowledgment for DATA or HEARTBEAT chunks sent to one of the destination addresses, it considers that destination to be active. If it fails to receive acknowledgments, the error count for the address is increased. If the error counter exceeds the protocol parameter 'Path.Max.Retrans', SCTP endpoint considers the address to be inactive.

The failover process of SCTP is initiated when the primary path becomes inactive (error counter for the primary path exceeds Path.Max.Retrans). If the primary path is marked inactive, SCTP chooses a new destination address from one of the active destinations and start using this address to send data to. If the primary path becomes active again, SCTP uses the primary destination for subsequent data transmissions and stop using non-primary one.

One issue with this failover process is that it usually takes significant amount of time before SCTP switches to the new destination. Let's say the primary path on a multi-homed host becomes unavailable and the RTO value for the primary path at that time is around 1 second, it usually takes over 60 seconds before SCTP starts to use the secondary path. This is because the recommended value for Path.Max.Retrans in the standard is 5, which requires 6 consecutive timeouts before failover takes place. Before SCTP

switches to the secondary address, SCTP keeps trying to send packets to the primary and only retransmitted packets are sent to the secondary and can thus be reached at the receiver. This slow failover process can cause significant performance degradation and will not be acceptable in some situations.

Another issue is that once the primary path is active again, the traffic is switched back. This is not optimal in some situations. This is further discussed in [Section 4.3](#).

[4.](#) SCTP with Potentially-Failed Destination State (SCTP-PF)

To address the issues described in [Section 3](#), this section updates SCTP path management scheme with the Potentially Failed state and associated Quick Failover operation. We use the term SCTP-PF to denote the resulting SCTP path management operation.

[4.1.](#) SCTP-PF Concept

SCTP-PF as defined stems from the following two observations about SCTP's failure detection procedure:

- o To minimize performance impact during failover, the sender should avoid transmitting data to the failed destination as early as possible. In the current SCTP path management scheme, the sender stops transmitting data to a destination only after the destination is marked Failed (inactive). Thus, a smaller PMR value is ideal so that the sender transitions a destination to the Failed (inactive) state quicker.
- o Smaller PMR values increase the chances of spurious failure detection where the sender incorrectly marks a destination as Failed (inactive) during periods of temporary congestion. As [\[RFC4960\]](#) recommends for a coupling of the PMR value and the AMR value such spurious failure detection risks to carry over to spurious association failure detection and closure. Larger PMR values are preferable to avoid spurious failure detection.

From the above observations it is clear that tuning the PMR value involves the following tradeoff -- a lower value improves performance but increases the chances of spurious failure detection, whereas a higher value degrades performance and reduces spurious failure detection in a wide range of path conditions. Thus, tuning the association's PMR value is an incomplete solution to address performance impact during failure.

This new method introduces a new "Potentially-Failed" (PF) destination state in SCTP's path management procedure. The PF state

was originally proposed to improve CMT performance [[NATARAJAN09](#)]. The PF state is an intermediate state between Active and Failed states. SCTP's failure detection procedure is modified to include the PF state. The new failure detection algorithm assumes that loss detected by a timeout implies either severe congestion or failure en-route. After a number of consecutive timeouts on a path, the sender is unsure, and marks the corresponding destination as PF. A PF destination is not used for data transmission except in special cases (discussed below). The new failure detection algorithm requires only sender-side changes.

4.2. SCTP-PF Algorithm Detail

SCTP PF operation is specified as follows:

1. The sender maintains a new tunable parameter called Potentially-Failed.Max.Retrans (PFMR). The RECOMMENDED value of PFMR = 0 when Quick Failover is used. When PFMR is larger or equal to PMR, Quick Failover is turned off.
2. The error counter of an active destination address is incremented as specified in [[RFC4960](#)]. This means that the error counter of the destination address will be incremented each time the T3-rtx timer expires, or at times where a HEARTBEAT sent to an idle, active address is not acknowledged within an RT0. When the value in the destination address error counter exceeds PFMR, the endpoint MUST mark the destination transport address as PF.
3. The sender SHOULD avoid data transmission to PF destinations. When the destinations are all in PF state or some in PF state and some in inactive state, the sender MUST choose one destination in PF state and transmit data to this destination. The sender SHOULD choose the destination in PF state with the lowest error count (fewest consecutive timeouts) for data transmission and transmit data to this destination. When there are multiple PF destinations with same error count, the sender SHOULD let the choice among the multiple PF destination with equal error count be based on the [[RFC4960](#)], [section 6.4.1](#), principles of choosing most divergent source-destination pairs when executing (potentially consecutive) retransmission. This means that the sender SHOULD attempt to pick the most divergent source - destination pair from the last source - destination pair on which data were transmitted or retransmitted. Rules for picking the most divergent source-destination pair are an implementation decision and are not specified within this document. A sender may choose to deploy other strategies than the above when choosing among multiple PF destinations with

equal error count. In all cases the sender MUST NOT change the state of chosen destination and it MUST NOT clear the destination's error counter as a result of choosing the destination for data transmission.

4. Heartbeats SHOULD be sent to PF destination(s) once per RT0. This means the sender MUST ignore HB.interval for PF destinations. If an heartbeat is unanswered, the sender SHOULD increment the error counter and exponentially back off the RT0 value. If error counter is less than PMR, the sender SHOULD transmit another heartbeat immediately after T3-timer expiration. When data is transmitted to a PF destination, the transmission of heartbeats may be omitted as SACK or T3-rtx timer expiration can provide equivalent information. It is RECOMMENDED that heartbeats be sent to PF destinations regardless of whether the Path Heartbeat function ([Section 8.3 of \[RFC4960\]](#)) is enabled for the destination address or not.
5. When the sender receives an heartbeat ACK from a PF destination, the sender MUST clear the destination's error counter and transition the PF destination back to Active state. When the sender resumes data transmission on the destination it MUST do this following the prescriptions of [Section 7.2 of \[RFC4960\]](#).
6. Additional (PMR - PFMR) consecutive timeouts on a PF destination confirm the path failure, upon which the destination transitions to the Inactive state. As described in [\[RFC4960\]](#), the sender (i) SHOULD notify ULP about this state transition, and (ii) transmit heartbeats to the Inactive destination at a lower frequency as described in [Section 8.3 of \[RFC4960\]](#) (when this function is enabled for the destination address).
7. When all destinations are in inactive state (association dormant state) the sender MUST also choose one destination to transmit data to. The sender SHOULD choose the destination in inactive state with the lowest error count (fewest consecutive timeouts) for data transmission and transmit data to this destination. When there are multiple destinations with same error count in inactive state, the sender SHOULD attempt to pick the most divergent source - destination pair from the last source - destination pair on which data were transmitted or retransmitted following [\[RFC4960\]](#). Rules for picking the most divergent source-destination pair are an implementation decision and are not specified within this document. Therefore, a sender SHOULD allow for incrementing the destination error counters up to some reasonable limit larger than PMR+1, thus changing the prescriptions of [\[RFC4960\]](#), [section 8.3](#), in this respect. The exact limit to apply is not specified in this document but it is

considered reasonable to require for such to be an order of magnitude higher than the PMR value. A sender MAY choose to deploy other strategies than the above. For example, a sender could choose to prioritize the last active destination during dormant state. The strategy to prioritize the last active destination is optimal when some paths are permanently inactive, but suboptimal when paths' instability is transient. While the increment of the error counters above PMR+1 is a prerequisite for the error counter values to serve to guide the path selection in dormant state, then it is noted that by virtue of the introduction of the Potentially Failed state, one may deploy higher values of PMR without compromising the efficiency of the failover operation, and thus making the increase of path error counters above PMR+1 less critical as the dormant state will be less likely to happen. The downside of increasing the PMR value relative to the AMR value, however, is that the per destination address failure detection and notification of such to ULP thereby is weakened. In all cases the sender MUST NOT change the state of the chosen destination and it MUST NOT clear the destination's error counter as a result of choosing the destination for data transmission.

8. ACKs for chunks that have been transmitted to multiple destinations (i.e., a chunk which has been retransmitted to a different destination than the destination to which the chunk was first transmitted) SHOULD NOT clear the error count of an inactive destination and SHOULD NOT transition a PF destination back to Active state, since a sender cannot disambiguate whether the ACK was for the original transmission or the retransmission(s). The same ambiguity concerns the related congestion window growth. The bytes of a newly acknowledged chunk which has been transmitted to multiple destinations SHOULD be considered for contribution to the congestion window growth towards the destination where the chunk was last sent. The contribution of the acked bytes to the window growth is subject to the prescriptions described in [Section 7.2 of \[RFC4960\]](#) is fulfilled. A SCTP sender MAY apply a different approach for both the error count handling and the congestion control growth handling based on unequivocally information on which destination (including multiple destinations) the chunk reached. This document makes no reference to what such unequivocally information could consist of, neither how such unequivocally information could be obtained. The implementation of such an alternative approach is left to implementations.
9. ACKs for chunks which has been transmitted to one destination address only MUST clear the error counter of the destination address and MUST transition a PF destination back to Active

state. This situation can happen when new data is sent to a destination address in PF state. It can also happen in situations where the destination address is in PF state due to the occurrence of a spurious T3-rtx timer and ACKs start to arrive for data sent prior to occurrence of the spurious T3-rtx and data has not yet been retransmitted towards other destinations. This document does not specify special handling for detection of or reaction to spurious T3-rtx timeouts, e.g., for special operation vis-a-vis the congestion control handling or data retransmission operation towards a destination address which undergoes a transition from active to PF to active state due to a spurious T3-rtx timeout. But it is noted that this is an area which would benefit from additional attention, experimentation and specification for Single Homed SCTP as well as for Multi Homed SCTP protocol operation.

10. SCTP stack SHOULD provide the ULP with the means to expose the PF state of its destinations as well as the means to notify the state transitions from Active to PF, and vice-versa. When doing this, such SCTP stack MUST provide the ULP with the means to suppress exposure of PF state and association state transitions as well.

4.3. Optional Feature: Permanent Failover

In [[RFC4960](#)], an SCTP sender migrates the traffic back to the original primary destination once this destination becomes active again. As the CWND towards the original primary destination has to be rebuilt once data transfer resumes, the switch back to use the original primary path is not always optimal. Indeed [[CAR002](#)] shows that the switch back to the original primary may degrade SCTP performance compared to continuing data transmission on the same path, especially, but not only, in scenarios where this path's characteristics are better. In order to mitigate this performance degradation, Permanent Failover operation was proposed in [[CAR002](#)]. When SCTP changes the destination due to failover, Permanent Failover operation allows SCTP sender to continue data transmission on the new working path even if the old primary destination becomes active again. This is achieved by having SCTP perform a switch over of the primary path to the alternative working path rather than having SCTP switch back data transfer to the (previous) primary path.

The manner of switch over operation that is most optimal in a given scenario depends on the relative quality of a set primary path versus the quality of alternative paths available as well as it depends on the extent to which it is desired for the mode of operation to enforce traffic distribution over a number of network paths. I.e., load distribution of traffic from multiple SCTP associations may be

sought to be enforced by distribution of the set primary paths with [\[RFC4960\]](#) switchback operation. However as [\[RFC4960\]](#) switchback behavior is suboptimal in certain situations, especially in scenarios where a number of equally good paths are available, it is recommended for SCTP to support also, as alternative behavior, the Permanent Failover switch over modes of operation.

The Permanent Failover operation requires only sender side changes. The details are:

1. The sender maintains a new tunable parameter, called Primary.Switchover.Max.Retrans (PSMR). The PSMR MUST be set greater or equal to the PFMR value. Implementations MUST reject any other values of PSMR.
2. When the path error counter on a set primary path exceeds PSMR, the SCTP implementation MUST autonomously select and set a new primary path.
3. The primary path selected by the SCTP implementation MUST be the path which at the given time would be chosen for data transfer. A previously failed primary path MAY come in use as data transfer path as per normal path selection when the present data transfer path fails.
4. The recommended value of PSMR is PFMR when Permanent Failover is used. This means that no forced switchback to a previously failed primary path is performed. An implementation of Permanent Failover MUST support the setting of PSMR = PFMR. An implementation of Permanent Failover MAY support setting of PSMR > PFMR.
5. It MUST be possible to disable the Permanent Failover and obtain the standard switchback operation of [\[RFC4960\]](#).

This specifications RECOMMENDS a default configuration that uses standard [RFC4960](#) switchback, i.e., switch back to the old primary destination once the destination becomes active again. However, to support optimal operation in a wider range of network scenarios, an implementation MAY implement Permanent Failover operation as detailed above and MAY enable it based on network configurations or users' requests.

5. Socket API Considerations

This section describes how the socket API defined in [\[RFC6458\]](#) is extended to provide a way for the application to control and observe the quick failover behavior.

Please note that this section is informational only.

A socket API implementation based on [\[RFC6458\]](#) is, by means of the existing `SCTP_PEER_ADDR_CHANGE` event, extended to provide the event notification when a peer address enters or leaves the potentially failed state as well as the socket API implementation is extended to expose the potentially failed state of a peer address in the existing `SCTP_GET_PEER_ADDR_INFO` structure.

Furthermore, two new read/write socket options for the level `IPPROTO_SCTP` and the name `SCTP_PEER_ADDR_THLDS` and `SCTP_EXPOSE_POTENTIALLY_FAILED_STATE` are defined as described below. The first socket option is used to control the values of the PFMR and PSMP parameters described in [Section 4](#). The second one controls the exposition of the potentially failed path state.

Support for the `SCTP_PEER_ADDR_THLDS` and `SCTP_EXPOSE_POTENTIALLY_FAILED_STATE` socket options need also to be added to the function `sctp_opt_info()`.

[5.1](#). Support for the Potentially Failed Path State

As defined in [\[RFC6458\]](#), the `SCTP_PEER_ADDR_CHANGE` event is provided if the status of a peer address changes. In addition to the state changes described in [\[RFC6458\]](#), this event is also provided, if a peer address enters or leaves the potentially failed state. The notification as defined in [\[RFC6458\]](#) uses the following structure:

```
struct sctp_paddr_change {
    uint16_t spc_type;
    uint16_t spc_flags;
    uint32_t spc_length;
    struct sockaddr_storage spc_aaddr;
    uint32_t spc_state;
    uint32_t spc_error;
    sctp_assoc_t spc_assoc_id;
}
```

[\[RFC6458\]](#) defines the constants `SCTP_ADDR_AVAILABLE`, `SCTP_ADDR_UNREACHABLE`, `SCTP_ADDR_REMOVED`, `SCTP_ADDR_ADDED`, and `SCTP_ADDR_MADE_PRIM` to be provided in the `spc_state` field. This document defines in addition to that the new constant `SCTP_ADDR_POTENTIALLY_FAILED`, which is reported if the affected address becomes potentially failed.

The `SCTP_GET_PEER_ADDR_INFO` socket option defined in [\[RFC6458\]](#) can be used to query the state of a peer address. It uses the following structure:


```
struct sctp_paddrinfo {
    sctp_assoc_t spinfo_assoc_id;
    struct sockaddr_storage spinfo_address;
    int32_t spinfo_state;
    uint32_t spinfo_cwnd;
    uint32_t spinfo_srtt;
    uint32_t spinfo_rto;
    uint32_t spinfo_mtu;
};
```

[RFC6458] defines the constants `SCTP_UNCONFIRMED`, `SCTP_ACTIVE`, and `SCTP_INACTIVE` to be provided in the `spinfo_state` field. This document defines in addition to that the new constant `SCTP_POTENTIALLY_FAILED`, which is reported if the peer address is potentially failed.

5.2. Peer Address Thresholds (`SCTP_PEER_ADDR_THLDS`) Socket Option

Applications can control the quick failover behavior by getting or setting the number of consecutive timeouts before a peer address is considered potentially failed or unreachable and before the primary path is changed automatically. This socket option uses the level `IPPROTO_SCTP` and the name `SCTP_PEER_ADDR_THLDS`.

The following structure is used to access and modify the thresholds:

```
struct sctp_paddrthlds {
    sctp_assoc_t spt_assoc_id;
    struct sockaddr_storage spt_address;
    uint16_t spt_pathmaxrxt;
    uint16_t spt_pathpfthld;
    uint16_t spt_pathcpthld;
};
```

`spt_assoc_id`: This parameter is ignored for one-to-one style sockets. For one-to-many style sockets the application may fill in an association identifier or `SCTP_FUTURE_ASSOC`. It is an error to use `SCTP_{CURRENT|ALL}_ASSOC` in `spt_assoc_id`.

`spt_address`: This specifies which peer address is of interest. If a wildcard address is provided, this socket option applies to all current and future peer addresses.

`spt_pathmaxrxt`: Each peer address of interest is considered unreachable, if its path error counter exceeds `spt_pathmaxrxt`.

spt_pathpfthld: Each peer address of interest is considered potentially failed, if its path error counter exceeds spt_pathpfthld.

spt_pathcpthld: Each peer address of interest is not considered the primary remote address anymore, if its path error counter exceeds spt_pathcpthld. Using a value of 0xffff disables the selection of a new primary peer address. If an implementation does not support the automatic selection of a new primary address, it should indicate an error with errno set to EINVAL if a value different from 0xffff is used in spt_pathcpthld. Setting of spt_pathcpthld < spt_pathpfthld should be rejected with errno set to EINVAL. An implementation MAY support only setting of spt_pathcpthld = spt_pathpfthld and spt_pathcpthld = 0xffff. In this case it shall reject setting of other values with errno set to EINVAL.

5.3. Exposing the Potentially Failed Path State

(SCTP_EXPOSE_POTENTIALLY_FAILED_STATE) Socket Option

Applications can control the exposure of the potentially failed path state in the SCTP_PEER_ADDR_CHANGE event and the SCTP_GET_PEER_ADDR_INFO as described in [Section 5.1](#). The default value is implementation specific.

This socket option uses the level IPPROTO_SCTP and the name SCTP_EXPOSE_POTENTIALLY_FAILED_STATE.

The following structure is used to control the exposition of the potentially failed path state:

```
struct sctp_assoc_value {
    sctp_assoc_t assoc_id;
    uint32_t assoc_value;
};
```

assoc_id: This parameter is ignored for one-to-one style sockets. For one-to-many style sockets the application may fill in an association identifier or SCTP_FUTURE_ASSOC. It is an error to use SCTP_{CURRENT|ALL}_ASSOC in assoc_id.

assoc_value: The potentially failed path state is exposed if and only if this parameter is non-zero.

6. Security Considerations

Security considerations for the use of SCTP and its APIs are discussed in [\[RFC4960\]](#) and [\[RFC6458\]](#). There are no new security considerations introduced in this document.

7. IANA Considerations

This document does not create any new registries or modify the rules for any existing registries managed by IANA.

8. Proposed Change of Status (to be Deleted before Publication)

Initially this work looked to entail some changes of the Congestion Control (CC) operation of SCTP and for this reason the work was proposed as Experimental. These intended changes of the CC operation have since been judged to be irrelevant and are no longer part of the specification. As the specification entails no other potential harmful features, consensus exists in the wg to bring the work forward as PS.

Initially concerns have been expressed about the possibility for the mechanism to introduce path bouncing with potential harmful network impacts. These concerns are believed to be unfounded. This issue is addressed in [Appendix B](#).

It is noted that the feature specified by this document is implemented by multiple SCTP SW implementations and furthermore that various variants of the solution have been deployed in Telco signaling environments for several years with good results.

9. References

9.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC4960] Stewart, R., "Stream Control Transmission Protocol", [RFC 4960](#), September 2007.

9.2. Informative References

- [CAR002] Caro Jr., A., Iyengar, J., Amer, P., Heinz, G., and R. Stewart, "A Two-level Threshold Recovery Mechanism for SCTP", Tech report, CIS Dept, University of Delaware , 7 2002.
- [CAR004] Caro Jr., A., Amer, P., and R. Stewart, "End-to-End Failover Thresholds for Transport Layer Multihoming", MILCOM 2004 , 11 2004.

- [CAR005] Caro Jr., A., "End-to-End Fault Tolerance using Transport Layer Multihoming", Ph.D Thesis, University of Delaware , 1 2005.
- [FALLON08] Fallon, S., Jacob, P., Qiao, Y., Murphy, L., Fallon, E., and A. Hanley, "SCTP Switchover Performance Issues in WLAN Environments", IEEE CCNC 2008, 1 2008.
- [GRINNEMO04] Grinnemo, K-J. and A. Brunstrom, "Performance of SCTP-controlled failovers in M3UA-based SIGTRAN networks", Advanced Simulation Technologies Conference , 4 2004.
- [IYENGAR06] Iyengar, J., Amer, P., and R. Stewart, "Concurrent Multipath Transfer using SCTP Multihoming over Independent End-to-end Paths.", IEEE/ACM Trans on Networking 14(5), 10 2006.
- [JUNGMAIER02] Jungmaier, A., Rathgeb, E., and M. Tuexen, "On the use of SCTP in failover scenarios", World Multiconference on Systemics, Cybernetics and Informatics , 7 2002.
- [NATARAJAN09] Natarajan, P., Ekiz, N., Amer, P., and R. Stewart, "Concurrent Multipath Transfer during Path Failure", Computer Communications , 5 2009.
- [RFC6458] Stewart, R., Tuexen, M., Poon, K., Lei, P., and V. Yasevich, "Sockets API Extensions for the Stream Control Transmission Protocol (SCTP)", [RFC 6458](#), December 2011.

[Appendix A](#). Discussions of Alternative Approaches

This section lists alternative approaches for the issues described in this document. Although these approaches do not require to update [RFC4960](#), we do not recommend them from the reasons described below.

[A.1](#). Reduce Path.Max.Retrans (PMR)

Smaller values for Path.Max.Retrans shorten the failover duration. In fact, this is recommended in some research results [[JUNGMAIER02](#)] [[GRINNEMO04](#)] [[FALLON08](#)]. For example, if when Path.Max.Retrans=0, SCTP switches to another destination on a single timeout. This smaller value for Path.Max.Retrans can results in spurious failover, which might be a problem.

Unlike SCTP-PF, the interval for heartbeat packets is governed by 'HB.interval' even during failover process. 'HB.interval' is usually set in the order of seconds (recommended value is 30 seconds). When the primary path becomes inactive, the next HB can be transmitted only seconds later. Meanwhile, the primary path may have recovered. In such situations, post failover, an endpoint is forced to wait on the order of seconds before the endpoint can resume transmission on the primary path. However, using smaller value for 'HB.interval' might help this situation, but it will be the waste of bandwidth in most cases.

In addition, smaller Path.Max.Retrans values also affect 'Association.Max.Retrans' values. When the SCTP association's error count (sum of error counts on all ACTIVE paths) exceeds Association.Max.Retrans threshold, the SCTP sender considers the peer endpoint unreachable and terminates the association. Therefore, [Section 8.2 in \[RFC4960\]](#) recommends that Association.Max.Retrans value should not be larger than the summation of the Path.Max.Retrans of each of the destination addresses, else the SCTP sender considers its peer reachable even when all destinations are INACTIVE. To avoid such inconsistent behavior an SCTP implementation SHOULD reduce Association.Max.Retrans accordingly whenever it reduces Path.Max.Retrans. However, smaller Association.Max.Retrans value increases chances of association termination during minor congestion events.

[A.2.](#) Adjust RTO related parameters

As several research results indicate, we can also shorten the duration of failover process by adjusting RTO related parameters [[JUNGMAIER02](#)] [[FALLON08](#)]. During failover process, RTO keeps being doubled. However, if we can choose smaller value for RTO.max, we can stop the exponential growth of RTO at some point. Also, choosing smaller values for RTO.initial or RTO.min can contribute to keep RTO value small.

Similar to reducing Path.Max.Retrans, the advantage of this approach is that it requires no modification to the current specification, although it needs to ignore several recommendations described in the [Section 15 of \[RFC4960\]](#). However, this approach requires to have enough knowledge about the network characteristics between end points. Otherwise, it can introduce adverse side-effects such as spurious timeouts.

Appendix B. Discussions for Path Bouncing Effect

The methods described in the document can accelerate the failover process. Hence, they might introduce the path bouncing effect where the sender keeps changing the data transmission path frequently. This sounds harmful to the data transfer, however several research results indicate that there is no serious problem with SCTP in terms of path bouncing effect [[CAR004](#)] [[CAR005](#)].

There are two main reasons for this. First, SCTP is basically designed for multipath communication, which means SCTP maintains all path related parameters (CWND, ssthresh, RTT, error count, etc) per each destination address. These parameters cannot be affected by path bouncing. In addition, when SCTP migrates the data transfer to another path, it starts with the minimal or the initial CWND. Hence, there is little chance for packet reordering or duplicating.

Second, even if all communication paths between the end-nodes share the same bottleneck, the quick failover results in a behavior already allowed by [[RFC4960](#)].

Authors' Addresses

Yoshifumi Nishida
GE Global Research
2623 Camino Ramon
San Ramon, CA 94583
USA

Email: nishida@wide.ad.jp

Preethi Natarajan
Cisco Systems
510 McCarthy Blvd
Milpitas, CA 95035
USA

Email: prenatar@cisco.com

Armando Caro
BBN Technologies
10 Moulton St.
Cambridge, MA 02138
USA

Email: acar@bbn.com

Paul D. Amer
University of Delaware
Computer Science Department - 434 Smith Hall
Newark, DE 19716-2586
USA

Email: amer@udel.edu

Karen E. E. Nielsen
Ericsson
Kistavaegen 25
Stockholm 164 80
Sweden

Email: karen.nielsen@tieto.com

