

Internet Engineering Task Force
INTERNET-DRAFT
Intended Status: Standards Track
Expires: January 2, 2016

X. Wei
L.Zhu
Huawei Technologies
L.Deng
China Mobile
B.Briscoe
July 1, 2015

Tunnel Congestion Feedback
draft-wei-tsvwg-tunnel-congestion-feedback-04

Abstract

This document describes a mechanism to calculate congestion of a tunnel segment based on RFC 6040 recommendations, and a feedback protocol by which to send the measured congestion of the tunnel from egress to ingress. A basic model for measuring tunnel congestion and feedback is described, and a protocol for carrying the feedback data is outlined.

Status of this Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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."

The list of current Internet-Drafts can be accessed at
<http://www.ietf.org/lid-abstracts.html>

The list of Internet-Draft Shadow Directories can be accessed at
<http://www.ietf.org/shadow.html>

Copyright and License Notice

Copyright (c) 2015 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	3
3. Congestion Information Feedback Models	4
3.1 Direct Model	4
3.2 Centralized Model	4
4. Congestion Level Measurement	5
5. Congestion Information Delivery	7
5.1 IPFIX Extentions	7
5.1.1 ce-cePacketTotalCount	7
5.1.2 ect-nectPacketTotalCount	8
5.1.3 ce-nectPacketTotalCount	8
5.1.4 ce-ectPacketTotalCount	8
5.1.5 ect-ectPacketTotalCount	9
6. Congestion Management	9
7. Security	9
8. IANA Considerations	10
9. References	10
9.1 Normative References	10
9.2 Informative References	10
Authors' Addresses	11

1. Introduction

In IP network, persistent congestion (or named congestion collapse) would cause transport throughput to drop down, lead to waste of network resource, so appropriate congestion control mechanisms are critical to make sure the network not fall into persistent congestion state. Currently, transport protocols such as TCP, SCTP, DCCP, has their built-in congestion control mechanism, and even for certain single transport protocol like TCP there could be a couple of different congestion control mechanism to choose. All these congestion control mechanisms are implemented on host side, and there are reasons that only host side congestion control is not sufficient for the whole network to keep away from persistent congestion, e.g., (1) some protocol's congestion control scheme might has internal design flaws; (2) improper software implementation of protocol; (3) some transport protocols even don't provide congestion control at all.

In order to have a better control on network congestion status, it's necessary for the network side to do certain kind of traffic control. For example, ConEx [ConEx] provides a method for network operator to learn about traffic's congestion contribution information, and then congestion management action could be taken based on this information.

Tunnels are widely deployed in various networks including public Internet, datacenter network, and enterprise network etc, a tunnel consists of an ingress, an egress and a set of interior routers. For the tunnel scenario, a tunnel-based mechanism which is different from ConEx is introduced for network traffic control to keep network away from persistent congestion; in this case, tunnel ingress will implement congestion management function to control the traffic entering the tunnel.

In order to do congestion management at ingress, the ingress must first get the inner tunnel congestion level information. But the ingress cannot use the locally visible traffic rates, because it would require additional knowledge of downstream capacity and topology, as well as cross traffic that does not pass through this ingress.

This document provide a mechanism of feeding back inner tunnel congestion level to ingress, using this mechanism the egress could feed the tunnel congestion level information it collects back to ingress, after receiving the information ingress could do congestion management according to network management policy.

2. Conventions

In this model, after egress collects network congestion level information, it feeds back the information to controller instead of ingress, and then the controller makes congestion management decision and sends the decision to ingress.

4. Congestion Level Measurement

This section describes how to measure congestion level in tunnel.

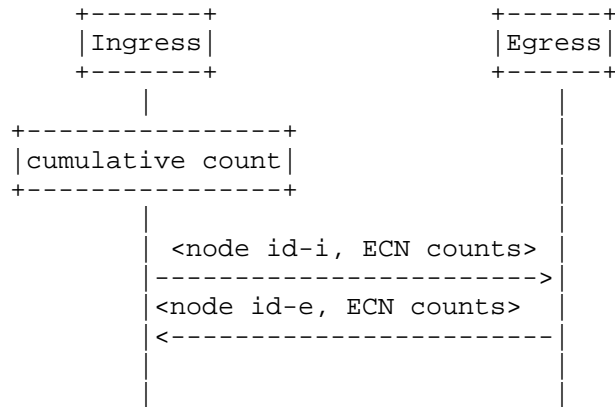
There may be different approaches of packet loss detection for different tunneling protocol scenarios, for instance, if there is a sequence field in tunneling protocol header, it will be easy for egress to detect packet loss through the gaps in sequence number space; another approach is to compare the number of packets entering ingress and the number of packets arriving at egress over the same span of packets. This document will focus on the latter one which is a more general approach.

If the routers support ECN, after router's queue length is over a predefined threshold, the routers will mark ECN packets as CE packets or drop not-ECN packets with the probability proportional to queue length, if the queue overflows all packets will be dropped; if the routers don't support ECN, after router's queue length is over a predefined threshold, the routers will drop both ECN packets and not-ECN packets with the probability proportional to queue length. It's assumed all routers in the tunnel support ECN.

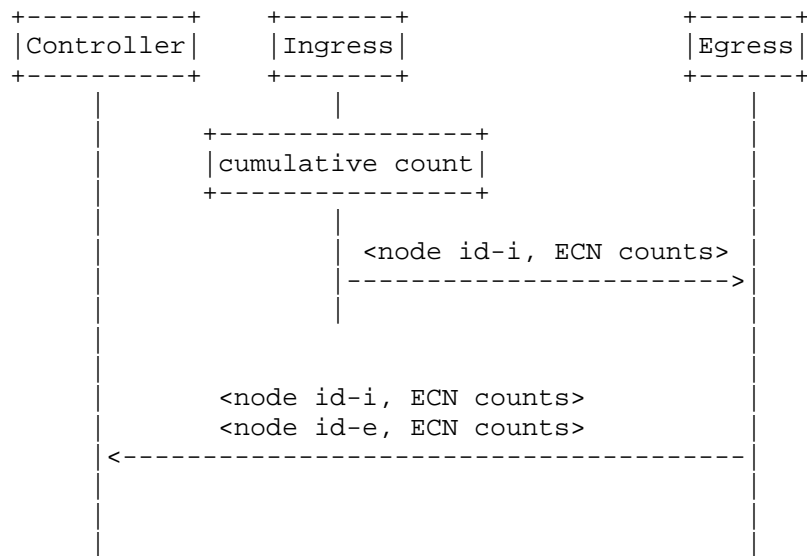
Faked ECT is used at ingress to defer packet loss to egress. The basic idea of faked ECT is that, when encapsulating packets, ingress first marks tunnel outer header according to RFC6040, and then remarks outer header of Not-ECT packet as ECT, there will be three kinds of combination of outer header ECN field and inner header ECN field: CE|CE, ECT|N-ECT, ECT|ECT (in the form of outer ECN| inner ECN).

In case all interior routers support ECN, the network congestion level could be indicated through the ratio of CE-marked packet and the ratio of packet drop, the relationship between these two kinds of indicator is complementary. If the congestion level in tunnel is not high enough, the packets would be marked as CE instead of being dropped, and then it is easy to calculate congestion level according to the ratio of CE-marked packets; if the congestion level is so high that ECT packet will be dropped, then the packet loss ratio could be calculated by comparing total packets entering ingress and total packets arriving at egress over the same span of packets, if packet loss is detected, it could be assumed that severe congestion has occurred in the tunnel, because loss is only ever a sign of serious congestion, so it doesn't need to measure loss ratio accurately.

The basic procedure of congestion level measurement is as follows:



(a) Direct model feedback procedure



(b) Centralized model feedback procedure

Ingress encapsulates packets and marks outer header according to faked ECT as described above. Ingress cumulatively counts packets for three types of ECN combination (CE|CE, ECT|N-ECT, ECT|ECT) and then the ingress regularly sends cumulative packet counts message of each type of ECN combination to the egress. When each message arrives, the

egress cumulatively counts packets coming from the ingress and adds its own packet counts of each type of ECN combination (CE|CE, ECT|N-ECT, CE|N-ECT, CE|ECT, ECT|ECT) to the message and either returns the whole message to the ingress, or to a central controller.

The counting of packets could be at the granularity of the all traffic from the ingress to the egress to learn about the overall congestion status of the path between the ingress and the egress; or at the granularity of individual customer's traffic or a specific set of flows to learn about their congestion contribution.

5. Congestion Information Delivery

As described above, the tunnel ingress needs to convey message of cumulative packet counts of each type of ECN combination to tunnel egress, and the tunnel egress also needs to feed the message of cumulative packet counts of each type of ECN combination to the ingress or central collector. This section describes how the messages could be conveyed.

The message could be along the same path with network data traffic, referred as in band signal; or go through a different path with network data traffic, referred as out of band signal. Because out of band scheme needs additional separate path which might limit its actual deployment, so the in band scheme will be discussed here.

Because the message is transmitted in band, so the message packet might get lost in case of network congestion. To cope with the situation that message packet gets lost, the packet counts values are sent as cumulative counters, so if a message is lost the next message will recover the missing information.

IPFIX [RFC7011] is selected as a choice of candidate protocol. IPFIX is preferred to use SCTP as transport, and because SCTP allows partially reliable delivery [RFC3758], which makes sure the feedback message will not be blocked to be sent in case of SCTP packets lost due to network congestion.

When sending message from ingress to egress, the ingress acts as IPFIX exporter and egress acts as IPFIX collector; when sending message from egress to ingress or controller, the egress acts as IPFIX exporter and ingress or controller acts as IPFIX collector.

5.1 IPFIX Extensions

5.1.1 ce-cePacketTotalCount

Description: The total number of incoming packets with CE|CE ECN

marking combination for this Flow at the Observation Point since the Metering Process (re-)initialization for this Observation Point.

Abstract Data Type: unsigned64

Data Type Semantics: totalCounter

ElementId: TBD1

Statuses: current

Units: packets

5.1.2 ect-nectPacketTotalCount

Description: The total number of incoming packets with ECT|N-ECT ECN marking combination for this Flow at the Observation Point since the Metering Process (re-)initialization for this Observation Point.

Abstract Data Type: unsigned64

Data Type Semantics: totalCounter

ElementId: TBD2

Statuses: current

Units: packets

5.1.3 ce-nectPacketTotalCount

Description: The total number of incoming packets with CE|N-ECT ECN marking combination for this Flow at the Observation Point since the Metering Process (re-)initialization for this Observation Point.

Abstract Data Type: unsigned64

Data Type Semantics: totalCounter

ElementId: TBD3

Statuses: current

Units: packets

5.1.4 ce-ectPacketTotalCount

Description: The total number of incoming packets with CE|ECT ECN

marking combination for this Flow at the Observation Point since the Metering Process (re-)initialization for this Observation Point.

Abstract Data Type: unsigned64

Data Type Semantics: totalCounter

ElementId: TBD4

Statuses: current

Units: packets

5.1.5 ect-ectPacketTotalCount

Description: The total number of incoming packets with ECT|ECT ECN marking combination for this Flow at the Observation Point since the Metering Process (re-)initialization for this Observation Point.

Abstract Data Type: unsigned64

Data Type Semantics: totalCounter

ElementId: TBD5

Statuses: current

Units: packets

6. Congestion Management

After tunnel ingress (or controller) receives congestion level information, then congestion management actions could be taken based on the information, e.g. if the congestion level is higher than a predefined threshold, then action could be taken to reduce the congestion level.

Congestion management action must be delayed by more than a worst-case global RTT, otherwise tunnel traffic management will not give normal e2e congestion control enough time to do its job, and the system could go unstable. The detailed description of congestion management is out of scope of this document, as examples, congestion management such as circuit breaker [CB] and congestion policing [CP] could be applied.

7. Security

This document describes the tunnel congestion calculation and

feedback. For feeding back congestion, security mechanisms of IPFIX are expected to be sufficient. No additional security concerns are expected.

8. IANA Considerations

This document defines a set of new IPFIX Information Elements (IE). New registry for these IE identifiers is needed.

TBD1~TBD5.

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, <<http://www.rfc-editor.org/info/rfc2119>>.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001, <<http://www.rfc-editor.org/info/rfc3168>>.
- [RFC3758] Stewart, R., Ramalho, M., Xie, Q., Tuexen, M., and P. Conrad, "Stream Control Transmission Protocol (SCTP) Partial Reliability Extension", RFC 3758, May 2004, <<http://www.rfc-editor.org/info/rfc3758>>.
- [RFC7011] Claise, B., Ed., Trammell, B., Ed., and P. Aitken, "Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of Flow Information", STD 77, RFC 7011, September 2013, <<http://www.rfc-editor.org/info/rfc7011>>.

9.2 Informative References

- [CONEX] Matt Mathis, Bob Briscoe. "Congestion Exposure (ConEx) Concepts, Abstract Mechanism and Requirements", draft-ietf-conex-abstract-mech-13, October 24, 2014
- [CB] G. Fairhurst. "Network Transport Circuit Breakers", draft-ietf-tsvwg-circuit-breaker-01, April 02, 2015
- [CP] Bob Briscoe, Murari Sridharan. "Network Performance Isolation in Data Centres using Congestion Policing", draft-briscoe-

conex-data-centre-02, February 14, 2014

Authors' Addresses

Xinpeng Wei
Beiqing Rd. Z-park No.156, Haidian District,
Beijing, 100095, P. R. China
E-mail: weixinpeng@huawei.com

Zhu Lei
Beiqing Rd. Z-park No.156, Haidian District,
Beijing, 100095, P. R. China
E-mail: lei.zhu@huawei.com

Lingli Deng
Beijing, 100095, P. R. China
E-mail: denglingli@gmail.com

Bob Briscoe
B54/77, Adastral Park
Martlesham Heath
Ipswich IP5 3RE
UK