

Workgroup: Mboned
Internet-Draft: draft-ietf-mboned-cbacc-04
Published: 7 March 2022
Intended Status: Standards Track
Expires: 8 September 2022
Authors: J. Holland

Akamai Technologies, Inc.

Circuit Breaker Assisted Congestion Control

Abstract

This document specifies Circuit Breaker Assisted Congestion Control (CBACC). CBACC enables fast-trip Circuit Breakers by publishing rate metadata about multicast channels from senders to intermediate network nodes or receivers. The circuit breaker behavior is defined as a supplement to receiver driven congestion control systems, to preserve network health if misbehaving or malicious receiver applications subscribe to a volume of traffic that exceeds capacity policies or capability for a network or receiving device.

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 <https://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 8 September 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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 Revised BSD License text as described in

Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

- [1. Introduction](#)
 - [1.1. Background and Terminology](#)
 - [1.2. Venues for Contribution and Discussion](#)
 - [1.3. Non-obvious doc choices](#)
- [2. Circuit Breaker Behavior](#)
 - [2.1. Functional Components](#)
 - [2.1.1. Bitrate Advertisement](#)
 - [2.1.2. Circuit Breaker Node](#)
 - [2.1.3. Communication Method](#)
 - [2.1.4. Measurement Function](#)
 - [2.1.5. Trigger Function](#)
 - [2.1.6. Reaction](#)
 - [2.1.7. Feedback Control Mechanism](#)
 - [2.2. States](#)
 - [2.2.1. Interface State](#)
 - [2.2.2. Flow State](#)
 - [2.3. Implementation Design Considerations](#)
 - [2.3.1. Oversubscription Thresholds](#)
 - [2.3.2. Fairness Functions](#)
- [3. YANG Module](#)
 - [3.1. Tree Diagram](#)
 - [3.2. Module](#)
- [4. IANA Considerations](#)
 - [4.1. YANG Module Names Registry](#)
 - [4.2. The XML Registry](#)
- [5. Security Considerations](#)
 - [5.1. Metadata Security](#)
 - [5.2. Denial of Service](#)
 - [5.2.1. State Overload](#)
- [6. Acknowledgements](#)
- [7. References](#)
 - [7.1. Normative References](#)
 - [7.2. Informative References](#)
- [Appendix A. Overjoining](#)
- [Author's Address](#)

1. Introduction

This document defines Circuit Breaker Assisted Congestion Control (CBACC). CBACC defines a Network Transport Circuit Breaker (CB), as described by [[RFC8084](#)].

The CB behavior defined in this document uses bit-rate metadata about multicast data streams coupled with policy, capacity, and load

information at a network location to prune multicast channels so that the network's aggregate capacity at that location is not exceeded by the subscribed channels.

To communicate the required metadata, this document defines a YANG [[RFC7950](#)] module that augments the DORMS [[I-D.draft-ietf-mboned-dorms](#)] YANG module. DORMS provides a mechanism for senders to publish metadata about the multicast streams they're sending through a RESTCONF service, so that receivers or forwarding nodes can discover and consume the metadata with a set of standard methods. The CBACC metadata MAY be communicated to receivers or forwarding nodes by some other method, but the definition of any alternative methods is out of scope for this document.

The CB behavior defined in this document matches the description provided in Section 3.2.3 of [[RFC8084](#)] of a unidirectional CB over a controlled path. The control messages from that description are composed of the messages containing the metadata required for operation of the CB.

CBACC is designed to supplement protocols that use multicast IP and rely on well-behaved receivers to achieve congestion control. Examples of congestion control systems fitting this description include [[PLM](#)], [[RLM](#)], [[RLC](#)], [[FLID-DL](#)], [[SMCC](#)], and WEBRC [[RFC3738](#)].

CBACC addresses a problem with "overjoining" by untrusted receivers.

In an overjoining condition, receivers (either malicious, misconfigured, or with implementation errors) subscribe to multicast channels but do not respond appropriately to congestion. When sufficient multicast traffic is available for subscription by such receivers, this can overload any network.

The overjoining problem is relevant to misbehaving receivers for both receiver-driven and feedback-driven congestion control strategies, as described in Section 4.1 of [[RFC8085](#)].

Overjoining attacks and the challenges they present are discussed in more detail in [Appendix A](#).

CBACC offers a solution for the recommendation in Section 4 of [[RFC8085](#)] that circuit breaker solutions be used even where congestion control is optional.

1.1. Background and Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in

BCP 14 [[RFC2119](#)] [[RFC8174](#)] when, and only when, they appear in all capitals, as shown here.

1.2. Venues for Contribution and Discussion

This document is in the Github repository at:

<https://github.com/GrumpyOldTroll/ietf-dorms-cluster>

Readers are welcome to open issues and send pull requests for this document.

Please note that contributions may be merged and substantially edited, and as a reminder, please carefully consider the Note Well before contributing: <https://datatracker.ietf.org/submit/note-well/>

Substantial discussion of this document should take place on the MBONED working group mailing list (mboned@ietf.org).

*Join: <https://www.ietf.org/mailman/listinfo/mboned>

*Search: <https://mailarchive.ietf.org/arch/browse/mboned/>

1.3. Non-obvious doc choices

*Since nothing is necessarily being actively measured by a network component at the ingress, referring to the bitrate advertisement as an "ingress meter" for this context was considered confusing by reviewers, so the section was renamed with just a note pointing to the link. Likewise the egress meter and "CB node".

*TBD: might need more and better examples explaining the point in [Section 2.1.5.1](#)? Some reason to believe it's not sufficiently clear...

*Another TBD: consider Dino's suggestion from 2020-04-09 to include an operational considerations section that addresses some possible optimizations for CB placement and configuration.

*TBD: add a section walking through the requirements in <https://datatracker.ietf.org/doc/html/rfc8084#section-4> and explaining how this matches.

*I'm unclear on whether <https://datatracker.ietf.org/doc/html/rfc8407#section-3.8.2> applies here, such that providing an augmentation inside the DORMS namespace causes an update to the DORMS document.

2. Circuit Breaker Behavior

2.1. Functional Components

This section maps the functional components described in Section 3.1 of [\[RFC8084\]](#) to the operational components of the CBACC CB defined by this document.

2.1.1. Bitrate Advertisement

The metadata provides an advertised maximum data bit-rate, namely the "max-speed" field in the YANG model in [Section 3](#). This is a self-report by the sender about the maximum amount of traffic a sender will send within any time interval given by the "data-rate-window" field, which is the measurement interval for the CB. This value refers to the total IP Payload data for all packets in the same (S,G), and its units are in kilobits per second.

The sender MUST NOT send more data for a data stream than the amount of data declared according to its advertised data rate within any measurement window, and it's RECOMMENDED for the sender to provide some margin to account for the possibility of burst forwarding after traffic encounters a non-empty queue, e.g. as sometimes observed with ACK compression (see [\[ZSC91\]](#) for a description of the phenomenon). If a CB node observes a higher data rate transmitted within any measurement window, it MAY circuit-break that flow immediately.

In the terminology of [\[RFC8084\]](#), the bitrate advertisement qualifies as an ingress meter.

2.1.2. Circuit Breaker Node

A circuit breaker node (CB node) is a location in a network where the constraints of the network and the observations about active traffic are compared to the bitrate advertisement in order to make the decision loop about when and whether to perform the circuit breaking behavior. In the terminology of [\[RFC8084\]](#), the CB node qualifies as an egress meter.

The CB node has access to several pieces of information that can be used as relevant egress metrics that may include:

1. Physical capacity limits on each interface.
2. Configured capacity limits for multicast traffic for each interface.
3. The observed received data rates of subscribed multicast channels with CBACC metadata.

4. The observed received data rates of subscribed multicast channels without CBACC metadata.
5. The observed received data rates of competing non-multicast traffic.
6. The loss rate for subscribed multicast channels, when available. The loss rate is only sometimes observable at a CB node; for example, when using AMBI [[I-D.draft-ietf-mboned-ambi](#)], or when the data stream carries a protocol that is known to the CB node by some out of band means, and whose traffic can be monitored for loss. When available, the loss rates may be used.

Note that any on-path router can behave as a CB node, even though there may be other CB nodes downstream or upstream covering the same data streams. When viewing CB nodes as egress meters in the context of [[RFC8084](#)], it's important to recall there's not a single egress meter in the network, but rather an egress meter per CB node, representing potentially multiple overlaid circuit breakers that may redundantly cover parts of the same path, with potentially different constraints based on the network location where the egress meter operates. All of the CB nodes anywhere on a path constitute separate circuit breakers that may trip independently of other circuit breakers.

Also note that other kinds of components besides on-path routers forwarding the traffic can act as CB nodes, for example the operating system or browser on a device receiving the traffic, or the receiving application itself.

2.1.3. Communication Method

CBACC generally operates at a CB node, where metrics such as those described in [Section 2.1.2](#) are available through system calls, or by communication with various locally deployable system monitoring applications. However, the CBACC processing can equivalently occur on a separate device that can monitor statistics gathered at a CB node, as long as the necessary control functions to trigger the CB can be invoked.

The communication path defined in this document for the CB node to obtain the bitrate advertisement in [Section 2.1.1](#) is the use of DORMS [[I-D.draft-ietf-mboned-dorms](#)]. Other methods MAY be used as well or instead, but are out of scope for this document.

2.1.4. Measurement Function

The measurement function maintains a few values for each interface, computed from the metrics described in [Section 2.1.2](#) and [Section 2.1.1](#):

1. The aggregate advertised maximum bit-rate capacity consumed by CBACC data streams. This is the sum of the max-speed values in the CBACC metadata for all data streams subscribed through an interface
2. An oversubscription threshold for each interface. The oversubscription threshold will be determined differently for CB nodes in different contexts. In some network devices, it might be as simple as an administratively configured absolute value or proportion of an interface's capacity. For other situations, like a CB node operating in a context with loss visibility, it could be a dynamically changing value that grows when data streams are successfully subscribed and receiving data without loss, and shrinks as loss is observed across subscribed data streams. The oversubscription threshold calculation could also incorporate other information like out-of-band path capacity measurements with bandwidth detection techniques such as [\[PathChirp\]](#) or [\[CapProbe\]](#).

This document covers some non-normative examples of valid oversubscription threshold functions in [Section 2.3.1](#). In general, the oversubscription threshold is the primary parameter that different CBs in different contexts can tune to provide the safety guarantees necessary for their context.

2.1.5. Trigger Function

The trigger function fires when the aggregate advertised maximum bit-rate exceeds the oversubscription threshold for any interface.

When oversubscribed, the trigger function changes the states of subscribed channels to "blocked" until the aggregate subscribed bit-rate is below the oversubscription threshold again.

2.1.5.1. Fairness and Inter-flow Ordering

The trigger function orders the monitored flows according to a fairness function and a within-sender priority ordering (chosen by the sender as part of the CBACC metadata). When flows are blocked, they're blocked in order until the aggregate bitrate of the permitted flows do not exceed the oversubscription thresholds monitored by the CB node.

Flows from a single sender MUST be ordered according to their priority field from the CBACC metadata when compared with each other. This takes precedence over the fairness function ordering, since certain flows from the same sender may need strict priority over others.

For example, consider a sender using File Delivery over Unidirectional Transport (FLUTE, defined in [\[RFC6726\]](#)) that sends File Delivery Table (FDT) Instances (see section 3.2 of [\[RFC6726\]](#)) in one (S,G) and data for the various referenced files in other (S,G)s. In this case the data for the files will not be consumable without the (S,G) containing the FDT. Other transport protocols may similarly send control information (often with a lower bitrate) on one channel, and data information on another. In these cases, the sender may need to ensure that data channels are only available when the control channels are also available.

When comparing flows between senders, (S,G)s from the same sender with different priorities should be treated as aggregated (S,G)s with regard to their declared bitrate consumption, to ensure that if any flows from the same sender need to be pruned by the circuit-breaker, the least preferred priority flows from that sender are pruned first.

Between-sender flows and flows from the same sender with the same priority are ordered according to the fairness function. TBD: need to work thru details, this does not work as written. Sample fairness function would reward senders for splitting a flow in 2 (more total subscribers). Maybe should count offload instead? This has trouble from favoring padding in your flow, but is (i think?) dominated by subscriber count where that's known. The fairness function can be different for CBs in different contexts.

A CBACC CB implementation SHOULD provide mechanisms for administrative controls to configure explicit biases, as this may be necessary to support Service Level Agreements for specific events or providers, or to block or de-prioritize channels with historically known misbehavior.

Subject to the above constraints, where possible the default fairness behavior SHOULD favor streams with many receivers over streams with few receivers, and streams with a low bit-rate over streams with a high bit-rate. See [Section 2.3.2](#) for further considerations and examples.

2.1.6. Reaction

When the trigger function fires and a subscribed channel becomes blocked, the reaction depends on whether it's an upstream interface or a downstream interface.

If a channel is blocked on one or more downstream interfaces, it may still be unblocked on other downstream interfaces. When this is the case, traffic is simply not forwarded along blocked interfaces, even though clients might still be joined downstream of those interfaces.

When a channel is blocked on all downstream interfaces or when the upstream interface is oversubscribed, the channel is pruned so that data no longer arrives from the network on the upstream interface. The prune would be performed with a PIM prune (Section 3.5 of [\[RFC7761\]](#)), or a "leave" operation to be communicated via IGMP, MLD, or another multicast group signaling mechanism, according to the expected signaling within the network.

Once initially pruned, a flow SHOULD remain pruned for a minimum amount of time. The minimum hold-down duration SHOULD be no less than 2.5 minutes by default, even if available bitrate space clears up, to ensure downstream subscriptions will notice and respond. The hold-down duration SHOULD be extended from the minimum by a randomly chosen number of seconds uniformly distributed over a configurable desynchronization period, to avoid synchronized recovery of different circuit breakers along the path. The default length of the desynchronization period should be at least 30 seconds.

2.5 minutes is chosen to exceed the default maximum lifetime of 2 minutes that can occur if an IGMP responder suddenly stops operation, and ceases responding to IGMP queries with membership reports, and 30 seconds is chosen to allow for some flexibility in lost packets. The values MAY be administratively tuned as needed by network operators to meet performance goals specific to their networks or to the traffic they're forwarding.

When enough capacity is available for a circuit-broken stream to be unblocked and the circuit-breaker hold-down time is expired, flows SHOULD be unblocked according to the priority order until no more flows can be unblocked without exceeding the circuit breaker limits.

2.1.7. Feedback Control Mechanism

The bitrate advertisement metadata from [Section 2.1.1](#) should be refreshed as needed to maintain up to date values. When using DORMS and RESTCONF, the Subscription to YANG Notifications for Datastore Updates [\[RFC8641\]](#) is the preferred method to receive changes if available.

If datastore subscriptions are not supported by the client or server, the HTTP Cache Control headers provide valid refresh time properties from the server, and SHOULD be used if present. If No-Cache is used, the default refresh timing SHOULD be 30 seconds. A uniformly distributed random value between 0 and 10 seconds SHOULD be added to the Cache Control or the default refresh timing to avoid synchronization across multiple clients.

2.2. States

2.2.1. Interface State

A CB holds the following state for each interface, for both the inbound and outbound directions on that interface:

*aggregate bandwidth: The sum of the bandwidths of all non-circuit-broken CBACC flows that transit this interface in this direction.

*bandwidth limit: The maximum aggregate CBACC advertised bandwidth allowed, not including circuit-broken flows.

When reducing the bandwidth limit due to congestion, the circuit breaker SHOULD NOT reduce the limit by more than half its value in 10 seconds, and SHOULD use a smoothing function to reduce the limit gradually over time.

It is RECOMMENDED that no more than half the capacity for a link be allocated to CBACC flows if the link might be shared with unicast traffic that is responsive to congestion.

2.2.2. Flow State

Data streams with CBACC metadata have a state for the upstream interface through which the stream is joined:

*'subscribed'

Indicates that the circuit breaker is subscribed upstream to the flow and forwarding packets through zero or more egress interfaces.

*'pruned'

Indicates that the flow has been circuit-broken. A request to unsubscribe from the flow has been sent upstream, e.g. a PIM prune (Section 3.5 of [[RFC7761](#)]) or a "leave" operation communicated via IGMP, MLD, or another group membership management mechanism.

Data streams also have a per-interface state for downstream interfaces with subscribers, where the data is being forwarded. It's one of:

`*'forwarding'`

Indicates that the flow is a non-circuit-broken flow in steady state, forwarding packets downstream.

`*'blocked'`

Indicates that data packets for this flow are NOT forwarded downstream via this interface.

2.3. Implementation Design Considerations

2.3.1. Oversubscription Thresholds

TBD.

2.3.2. Fairness Functions

As an example fairness function that makes good sense for a general case of unknown traffic:

Consider a network where the receiver count for multicast channels is known, for example via the experimental PIM extension for population count defined in [[RFC6807](#)].

A good fairness metric for a flow is max-bandwidth divided by receiver-count, with lower values of the fairness metric favored over higher values.

An overview of some other approaches to appropriate fairness metrics is given in Section 2.3 of [[RFC5166](#)].

3. YANG Module

3.1. Tree Diagram

The tree diagram below follows the notation defined in [[RFC8340](#)].

module: ietf-cbacc

augment /dorms:dorms/dorms:metadata/dorms:sender/dorms:group:

+-rw cbacc!

+-rw max-speed uint32

+-rw max-packet-size? uint16

+-rw data-rate-window? uint32

+-rw priority? uint16

3.2. Module

```

<CODE BEGINS> file ietf-cbacc@2022-03-07.yang
module ietf-cbacc {
  yang-version 1.1;

  namespace "urn:ietf:params:xml:ns:yang:ietf-cbacc";
  prefix "cbacc";

  import ietf-dorms {
    prefix "dorms";
    reference "I-D.jholland-mboned-dorms";
  }

  organization "IETF";

  contact
    "Author:   Jake Holland
              <mailto:jholland@akamai.com>
    ";

  description
    "Copyright (c) 2019 IETF Trust and the persons identified as
    authors of the code.  All rights reserved.

    Redistribution and use in source and binary forms, with or
    without modification, is permitted pursuant to, and subject to
    the license terms contained in, the Simplified BSD License set
    forth in Section 4.c of the IETF Trust's Legal Provisions
    Relating to IETF Documents
    (https://trustee.ietf.org/license-info).

    This version of this YANG module is part of
    draft-jholland-mboned-cbacc.  See the internet draft for full
    legal notices.

    The key words 'MUST', 'MUST NOT', 'REQUIRED', 'SHALL', 'SHALL
    NOT', 'SHOULD', 'SHOULD NOT', 'RECOMMENDED', 'NOT RECOMMENDED',
    'MAY', and 'OPTIONAL' in this document are to be interpreted as
    described in BCP 14 (RFC 2119) (RFC 8174) when, and only when,
    they appear in all capitals, as shown here.

    This module contains the definition for bandwidth consumption
    metadata for SSM channels, as an extension to DORMS
    (draft-ietf-mboned-dorms).";

  revision 2021-07-08 {
    description "Draft version, post-early-review.";
    reference
      "draft-ietf-mboned-cbacc";
  }
}

```

```

augment
  "/dorms:dorms/dorms:metadata/dorms:sender/dorms:group" {
    description "Definition of the manifest stream providing
      integrity info for the data stream";

  container cbacc {
    presence "CBACC-enabled flow";
    description
      "Information to enable fast-trip circuit breakers";
    leaf max-speed {
      type uint32;
      units "kilobits/second";
      mandatory true;
      description "Maximum bitrate for this stream, in Kilobits
        of IP packet data (including headers) of native
        multicast traffic per second";
    }
    leaf max-packet-size {
      type uint16;
      default 1400;
      description "Maximum IP payload size, in octets.";
    }
    leaf data-rate-window {
      type uint32;
      units "milliseconds";
      default 2000;
      description
        "Time window over which data rate is guaranteed,
          in milliseconds.";
      /* TBD: range limits? */
    }
    leaf priority {
      type uint16;
      default 256;
      description
        "The relative preference level for keeping this flow
          compared to other flows from this sender (higher
          value is more preferred to keep)";
    }
  }
}
}
}

```

<CODE ENDS>

4. IANA Considerations

4.1. YANG Module Names Registry

This document adds one YANG module to the "YANG Module Names" registry maintained at <https://www.iana.org/assignments/yang-parameters>. The following registrations are made, per the format in Section 14 of [\[RFC6020\]](#):

```
name:      ietf-cbacc
namespace: urn:ietf:params:xml:ns:yang:ietf-cbacc
prefix:    cbacc
reference: I-D.draft-ietf-mboned-cbacc
```

4.2. The XML Registry

This document adds the following registration to the "ns" subregistry of the "IETF XML Registry" defined in [\[RFC3688\]](#), referencing this document.

```
URI: urn:ietf:params:xml:ns:yang:ietf-cbacc
Registrant Contact: The IESG.
XML: N/A, the requested URI is an XML namespace.
```

5. Security Considerations

TBD: Yang Doctor review from Reshad said this should "mention the YANG data nodes". I think this means "do what <https://tools.ietf.org/html/rfc8407#section-3.7> says"?

5.1. Metadata Security

Be sure to authenticate the metadata. See DORMS security considerations, and don't accept unauthenticated metadata if using an alternative means.

5.2. Denial of Service

5.2.1. State Overload

Since CBACC flows require state, it may be possible for a set of receivers and/or senders, possibly acting in concert, to generate many flows in an attempt to overflow the circuit breakers' state tables.

It is permissible for a network node to behave as a CBACC circuit breaker for some CBACC flows while treating other CBACC flows as non-CBACC, as part of a load balancing strategy for the network as a whole, or simply as defense against this concern when the number of monitored flows exceeds some threshold.

The same techniques described in Section 3.1 of [RFC4609] can be used to help mitigate this attack, for much the same reasons. It is RECOMMENDED that network operators implement measures to mitigate such attacks.

6. Acknowledgements

Many thanks to Devin Anderson, Ben Kaduk, Cheng Jin, Scott Brown, Miroslav Ponec, Bob Briscoe, Lenny Giuliani, Christian Worm Mortensen, Dino Farinacci, and Reshad Rahman for their thoughtful comments and contributions.

7. References

7.1. Normative References

- [I-D.draft-ietf-mboned-ambi] Holland, J. and K. Rose, "Asymmetric Manifest Based Integrity", Work in Progress, Internet-Draft, draft-ietf-mboned-ambi-01, 31 October 2020, <<https://www.ietf.org/archive/id/draft-ietf-mboned-ambi-01.txt>>.
- [I-D.draft-ietf-mboned-dorms] Holland, J., "Discovery Of Restconf Metadata for Source-specific multicast", Work in Progress, Internet-Draft, draft-ietf-mboned-dorms-01, 31 October 2020, <<https://www.ietf.org/archive/id/draft-ietf-mboned-dorms-01.txt>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/

RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.

- [RFC7950] Bjorklund, M., Ed., "The YANG 1.1 Data Modeling Language", RFC 7950, DOI 10.17487/RFC7950, August 2016, <<https://www.rfc-editor.org/info/rfc7950>>.
- [RFC8084] Fairhurst, G., "Network Transport Circuit Breakers", BCP 208, RFC 8084, DOI 10.17487/RFC8084, March 2017, <<https://www.rfc-editor.org/info/rfc8084>>.
- [RFC8085] Eggert, L., Fairhurst, G., and G. Shepherd, "UDP Usage Guidelines", BCP 145, RFC 8085, DOI 10.17487/RFC8085, March 2017, <<https://www.rfc-editor.org/info/rfc8085>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8340] Bjorklund, M. and L. Berger, Ed., "YANG Tree Diagrams", BCP 215, RFC 8340, DOI 10.17487/RFC8340, March 2018, <<https://www.rfc-editor.org/info/rfc8340>>.

7.2. Informative References

- [CapProbe] Kapoor, R., Chen, L., Lao, L., Gerla, M., and M.Y. Sanadidi, "CapProbe: A Simple and Accurate Capacity Estimation Technique", September 2004, <<https://dl.acm.org/doi/pdf/10.1145/1015467.1015476>>.
- [FLID-DL] Byers, J.W., Horn, G., Luby, M., Mitzenmacher, M., Shaver, W., and IEEE, "FLID-DL: congestion control for layered multicast", DOI 10.1109/JSAC.2002.803998, n.d., <<https://ieeexplore.ieee.org/document/1038584>>.
- [PathChirp] Ribeiro, V.J., Riedi, R.H., Baraniuk, R.G., Navratil, J., Cottrell, L., Department of Electrical and Computer Engineering Rice University, and SLAC/SCS-Network Monitoring, Stanford University, "pathChirp: Efficient Available Bandwidth Estimation for Network Paths", 2003.
- [PLM] Biersack, Institut EURECOM, A.Legout, E.W., "PLM: Fast Convergence for Cumulative Layered Multicast Transmission Schemes", 1999, <<http://www.eurecom.fr/en/publication/340/download/ce-legoar-000601.pdf>>.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688, DOI 10.17487/RFC3688, January 2004, <<https://www.rfc-editor.org/info/rfc3688>>.

[RFC3738]

Luby, M. and V. Goyal, "Wave and Equation Based Rate Control (WEBRC) Building Block", RFC 3738, DOI 10.17487/RFC3738, April 2004, <<https://www.rfc-editor.org/info/rfc3738>>.

[RFC4609]

Savola, P., Lehtonen, R., and D. Meyer, "Protocol Independent Multicast - Sparse Mode (PIM-SM) Multicast Routing Security Issues and Enhancements", RFC 4609, DOI 10.17487/RFC4609, October 2006, <<https://www.rfc-editor.org/info/rfc4609>>.

[RFC5166]

Floyd, S., Ed., "Metrics for the Evaluation of Congestion Control Mechanisms", RFC 5166, DOI 10.17487/RFC5166, March 2008, <<https://www.rfc-editor.org/info/rfc5166>>.

[RFC6020]

Bjorklund, M., Ed., "YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF)", RFC 6020, DOI 10.17487/RFC6020, October 2010, <<https://www.rfc-editor.org/info/rfc6020>>.

[RFC6726]

Paila, T., Walsh, R., Luby, M., Roca, V., and R. Lehtonen, "FLUTE - File Delivery over Unidirectional Transport", RFC 6726, DOI 10.17487/RFC6726, November 2012, <<https://www.rfc-editor.org/info/rfc6726>>.

[RFC6807]

Farinacci, D., Shepherd, G., Venaas, S., and Y. Cai, "Population Count Extensions to Protocol Independent Multicast (PIM)", RFC 6807, DOI 10.17487/RFC6807, December 2012, <<https://www.rfc-editor.org/info/rfc6807>>.

[RFC7761]

Fenner, B., Handley, M., Holbrook, H., Kouvelas, I., Parekh, R., Zhang, Z., and L. Zheng, "Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol Specification (Revised)", STD 83, RFC 7761, DOI 10.17487/RFC7761, March 2016, <<https://www.rfc-editor.org/info/rfc7761>>.

[RFC8641]

Clemm, A. and E. Voit, "Subscription to YANG Notifications for Datastore Updates", RFC 8641, DOI 10.17487/RFC8641, September 2019, <<https://www.rfc-editor.org/info/rfc8641>>.

[RLC]

Rizzo, L., Vicisano, L., and J. Crowcroft, "The RLC multicast congestion control algorithm", 1999, <<http://www.iet.unipi.it/~a007834/rlc99.ps.gz>>.

[RLM]

McCanne, S., Jacobson, V., Vetterli, M., University of California, Berkeley, and Lawrence Berkeley National Laboratory, "Receiver-driven Layered Multicast", 1995,

<<http://www1.cs.columbia.edu/~danr/courses/6761/Fall00/week9/layering.pdf>>.

[SMCC] Kwon, G., Byers, J.W., and Computer Science Department, Boston University, "Smooth Multirate Multicast Congestion Control", 2002, <<http://www.cs.bu.edu/techreports/pdf/2002-025-smcc.pdf>>.

[ZSC91] Zhang, L., Shenker, S., and D.D. Clark, "Observations and Dynamics of a Congestion Control Algorithm: The Effects of Two-Way Traffic", Proc. ACM SIGCOMM, ACM Computer Communications Review (CCR), Vol 21, No 4, pp.133-147. , 1991.

Appendix A. Overjoining

[RFC8085] describes several remedies for unicast congestion control under UDP, even though UDP does not itself provide congestion control. In general, any network node under congestion could in theory collect evidence that a unicast flow's sending rate is not responding to congestion, and would then be justified in circuit-breaking it.

With multicast IP, the situation is different, especially in the presence of malicious receivers. A well-behaved sender using a receiver-controlled congestion scheme such as WEBRC does not reduce its send rate in response to congestion, instead relying on receivers to leave the appropriate multicast groups.

This leads to a situation where, when a network accepts inter-domain multicast traffic, as long as there are senders somewhere in the world with aggregate bandwidth that exceeds a network's capacity, receivers in that network can join the flows and overflow the network capacity. A receiver controlled by an attacker could do this at the IGMP/MLD level without running the application layer protocol that participates in the receiver-controlled congestion control.

A network might be able to detect and defend against the most naive version of such an attack by blocking end users that try to join too many flows at once. However, an attacker can achieve the same effect by joining a few high-bandwidth flows, if those exist anywhere, and an attacker that controls a few machines in a network can coordinate the receivers so they join disjoint sets of non-responsive sending flows.

This scenario will produce congestion in a middle node in the network that can't be easily detected at the edge where the IGMP/MLD join is accepted. Thus, an attacker with a small set of machines in a target network can always trip a circuit breaker if present, or can induce excessive congestion among the bandwidth allocated to

multicast. This problem gets worse as more multicast flows become available.

Although the same can apply to non-responsive unicast traffic, network operators can assume that non-responsive sending flows are in violation of congestion control best practices, and can therefore cut off flows associated with the misbehaving senders. By contrast, non-responsive multicast senders are likely to be well-behaved participants in receiver-controlled congestion control schemes.

However, receiver controlled congestion control schemes also show the most promise for efficient massive scale content distribution via multicast, provided network health can be ensured. Therefore, mechanisms to mitigate overjoining attacks while still permitting receiver-controlled congestion control are necessary.

Author's Address

Jake Holland
Akamai Technologies, Inc.
150 Broadway
Cambridge, MA 02144,
United States of America

Email: jakeholland.net@gmail.com