

RTP Media Congestion Avoidance
Techniques
Internet-Draft
Intended status: Experimental
Expires: April 13, 2015

D. Hayes, Ed.
University of Oslo
S. Ferlin
Simula Research Laboratory
M. Welzl
University of Oslo
October 10, 2014

Shared Bottleneck Detection for Coupled Congestion Control for RTP
Media.

draft-hayes-rmcat-sbd-00

Abstract

This document describes a mechanism to detect whether end-to-end data flows share a common bottleneck. It relies on summary statistics that are calculated by a data receiver based on continuous measurements and regularly fed to a grouping algorithm that runs wherever the knowledge is needed. This mechanism complements the coupled congestion control mechanism in draft-welzl-rmcat-coupled-cc.

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 13, 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
1.1. The signals	3
1.1.1. Packet Loss	3
1.1.2. Packet Delay	3
1.1.3. Path Lag	4
2. Definitions	4
2.1. Parameter Values	5
3. Mechanism	5
3.1. Key metrics and their calculation	6
3.1.1. Mean delay	6
3.1.2. Skewness Estimate	7
3.1.3. Variance Estimate	7
3.1.4. Oscillation Estimate	8
3.1.5. Packet loss	8
3.2. Flow Grouping	8
3.2.1. Flow Grouping Algorithm	8
3.2.2. Using the flow group signal	9
4. Measuring OWD	10
4.1. Time stamp resolution	10
5. Acknowledgements	10
6. IANA Considerations	10
7. Security Considerations	10
8. References	11
8.1. Normative References	11
8.2. Informative References	11
Authors' Addresses	12

1. Introduction

In the Internet, it is not normally known if flows (e.g., TCP connections or UDP data streams) traverse the same bottlenecks. Even flows that have the same sender and receiver may take different paths and share a bottleneck or not. Flows that share a bottleneck link usually compete with one another for their share of the capacity. This competition has the potential to increase packet loss and delays. This is especially relevant for interactive applications that communicate simultaneously with multiple peers (such as multi-party video). For RTP media applications such as RTCWEB, [I-D.welzl-rmcat-coupled-cc] describes a scheme that combines the congestion controllers of flows in order to honor their priorities and avoid unnecessary packet loss as well as delay. This mechanism relies on some form of Shared Bottleneck Detection (SBD); here, a measurement-based SBD approach is described.

1.1. The signals

The current Internet is unable to explicitly inform endpoints as to which flows share bottlenecks, so endpoints need to infer this from packet loss and packet delay.

1.1.1. Packet Loss

Packet loss is often a relatively rare signal. Therefore, on its own it is of limited use for SBD, however, it is a valuable supplementary measure when it is more prevalent.

1.1.2. Packet Delay

End-to-end delay measurements include noise from every device along the path in addition to the delay perturbation at the bottleneck device. The noise is often significantly increased if the round-trip time is used. The cleanest signal is obtained by using One-Way-Delay (OWD).

Measuring absolute OWD is difficult since it requires both the sender and receiver clocks to be synchronised. However, since the statistics being collected are relative to the mean OWD, a relative OWD measurement is sufficient. Clock drift is not usually significant over the time intervals used by this SBD mechanism (see [RFC6817] A.2 for a discussion on clock drift and OWD measurements).

Each packet arriving at the bottleneck buffer may experience very different queue lengths, and therefore waiting times. A single OWD sample does therefore not characterize the actual OWD of a path well. However, multiple OWD measurements do reflect the distribution of

delays experienced at the bottleneck.

1.1.3. Path Lag

Flows that share a common bottleneck may traverse different paths, and these paths will often have different base delays. This makes it difficult to correlate changes in delay or loss. This technique uses the long term shape of the delay distribution as a base for comparison to counter this.

2. Definitions

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

Acronyms used in this document:

OWD -- One Way Delay

RTT -- Round Trip Time

SBD -- Shared Bottleneck Detection

Conventions used in this document:

T -- the base time interval over which measurements are made.

N -- the number of base time, T, intervals used in some calculations.

sum_T(...) -- summation of all the measurements of the variable in parentheses taken over the interval T

sum_N(...) -- summation of N terms of the variable in parentheses

sum_NT(...) -- summation of all measurements taken over the interval N*T

E_T(...) -- the expectation or mean of the measurements of the variable in parentheses over T

E_N(...) -- The expectation or mean of the last N values of the variable in parentheses

max_T(...) -- the maximum recorded measurement of the variable in parentheses taken over the interval T

p_l, p_f, p_pdf, p_s, p_d, p_v -- various thresholds used in the mechanism.

2.1. Parameter Values

Reference [Hayes-LCN14] uses T=350ms, N=50, p_l = 0.1, p_f = 0.2, p_pdf = 0.3, p_s = p_d = p_v = 0.2. These are values that seem to work well over a wide range of practical Internet conditions.

3. Mechanism

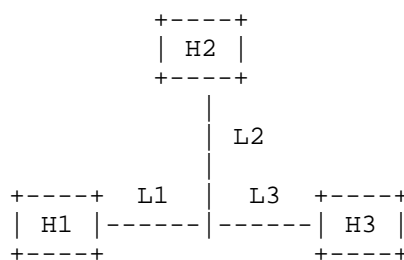
The mechanism described in this document is based on the observation that the delay measurements of flows that share a common bottleneck have similar shape characteristics. The shape of these characteristics are described using 3 key summary statistics:

variance (estimate PDV, see Section 3.1.3)

skewness (estimate skewest, see Section 3.1.2)

oscillation (estimate freqest, see Section 3.1.4)

Summary statistics help to address both the noise and the path lag problems by describing the general shape over a relatively long period of time. This is sufficient for their application in coupled congestion control for RTP Media. They can be signalled from a receiver, which measures the OWD and calculates the summary statistics, to a sender, which is the entity that is transmitting the media stream. An RTP Media device may be both a sender and a receiver. SBD can be performed at both the Sender and the Receiver.



A network with 3 hosts (H1, H2, H3) and 3 links (L1, L2, L3).

Figure 1

In Figure 1, there are two possible cases for shared bottleneck detection: a sender-based and a receiver-based case.

1. Sender-based: consider a situation where host H1 sends media streams to hosts H2 and H3, and L1 is a shared bottleneck. H2 and H3 measure the OWD and calculate summary statistics, which they send to H1 every T. H1, having this knowledge, can determine the shared bottleneck and accordingly control the send rates.
2. Receiver-based: consider that H2 is also sending media to H3, and L3 is a shared bottleneck. If H3 sends summary statistics to H1 and H2, neither H1 nor H2 alone obtain enough knowledge to detect this shared bottleneck; H3 can however determine it by combining the summary statistics related to H1 and H2, respectively. This case is applicable when send rates are controlled by the receiver; then, the signal from H3 to the senders contains the sending rate.

A discussion of the required signaling for the receiver-based case is beyond the scope of this document. For the sender-based case, the messages and their data format will be defined here in future versions of this document. We envision that an initialization message from the sender to the receiver could specify which key metrics are requested out of a possibly extensible set (losscnt, PDV, skewest, freqest). The grouping algorithm described in this document requires all four of these metrics, and receivers MUST be able to provide them, but future algorithms may be able to exploit other metrics (e.g. metrics based on explicit network signals). Moreover, the initialization message could specify T, N, and the necessary resolution and precision (number of bits per field).

3.1. Key metrics and their calculation

Measurements are calculated over a base interval, T. T should be long enough to provide enough samples for a good estimate of skewness, but short enough so that a measure of the oscillation can be made from N of these estimates. Reference [Hayes-LCN14] uses T = 350ms and N = 50, which are values that seem to work well over a wide range of practical Internet conditions.

3.1.1. Mean delay

The mean delay is not a useful signal for comparisons, however, it is a base measure for the 3 summary statistics. The mean delay, $E_T(\text{OWD})$, is the average one way delay measured over T.

To facilitate the other calculations, the last N $E_T(\text{OWD})$ values will need to be stored in a cyclic buffer along with the moving average of $E_T(\text{OWD})$:

$$E_N(E_T(\text{OWD})) = \text{sum}_N(E_T(\text{OWD})) / N$$

3.1.2. Skewness Estimate

Skewness is difficult to calculate efficiently and accurately. Ideally it should be calculated over the entire measurement for the entire period ($N * T$), however this would require storing every delay measurement over the period. Instead, an estimate is made over T using the previous calculation of $E_T(\text{OWD})$. Comparisons are made using the mean of N skew estimates.

The skewness is estimated using two counters, counting the number of one way delay samples above and below the mean:

$$\text{skewest} = (\text{sum}_T(\text{OWD} < E_T(\text{OWD})) - \text{sum}_T(\text{OWD} > E_T(\text{OWD}))) / \text{num}(\text{OWD})$$

where

$$\text{if} (\text{OWD} < E_T(\text{OWD})) \text{ 1 else } 0$$

$$\text{if} (\text{OWD} > E_T(\text{OWD})) \text{ 1 else } 0$$

skewest is a number between -1 and 1

$$E_N(\text{skewest}) = \text{sum}_N(\text{skewest}) / N$$

For implementation ease, $E_T(\text{OWD})$ is the mean delay of the previous T interval. Care must be taken when implementing the comparisons to ensure that rounding does not bias skewest.

3.1.3. Variance Estimate

Packet Delay Variation (PDV) ([RFC5481] and [ITU-Y1540]) is used as an estimator of the variance of the delay signal. We define PDV as follows:

$$\text{PDV} = (\text{max}(\text{OWD}) - E_T(\text{OWD}))$$

$$E_N(\text{PDV}) = \text{sum}_N(\text{PDV}) / N$$

This modifies PDV as outlined in [RFC5481] to provide a summary statistic version that best aids the grouping decisions of the algorithm (see [Hayes-LCN14] section IVB).

3.1.4. Oscillation Estimate

An estimate of the low frequency oscillation of the delay signal is calculated by counting and normalising the significant mean, $E_T(OWD)$, crossings of $E_N(E_T(OWD))$:

$$freqest = \text{number_of_crossings} / N$$

Where

we define a significant mean crossing as a crossing that extends $p_v * E_N(PDV)$ from $E_N(E_T(OWD))$. In our experiments we have found that $p_v = 0.2$ is a good value.

$freqest$ is a number between 0 and 1. $freqest$ can be approximated incrementally as follows:

With each new calculation of $E_T(OWD)$ a decision is made as to whether this value of $E_T(OWD)$ significantly crosses the current long term mean, $E_N(E_T(OWD))$, with respect to the previous significant mean crossing.

A cyclic buffer, $last_N_crossings$, records a 1 if there is a significant mean crossing, otherwise a 0.

The counter, $number_of_crossings$, is incremented when there is a significant mean crossing and subtracted from when a non zero value is removed from the $last_N_crossings$.

This approximation of $freqest$ was not used in [Hayes-LCN14], which calculated $freqest$ every T using the current $E_N(E_T(OWD))$. Our tests show that this approximation of $freqest$ yields results that are almost identical to when the full calculation is performed every T .

3.1.5. Packet loss

The proportion of packets lost is used as a supplementary measure:

$$PL_{NT} = \text{sum}_{NT}(\text{lost packets}) / \text{sum}_{NT}(\text{total packets})$$

3.2. Flow Grouping

3.2.1. Flow Grouping Algorithm

The following grouping algorithm is RECOMMENDED for SBD in this context and is sufficient and efficient for small to moderate numbers of flows. For very large numbers of flows, hundreds, a more complex clustering algorithm may be substituted.

Flows determined to be experiencing congestion are successively divided into groups based on freqest, PDV, and skewest.

The first step is to determine which flows are experiencing congestion. This is important, since if a flow is not experiencing congestion its delay based metrics will not describe the bottleneck, but the "noise" from the rest of the path. Skewness, with proportion of packets loss as a supplementary measure, is used to do this:

1. Grouping will be performed on flows where:

$$E_N(\text{skewest}) < 0 \quad || \quad PL_NT > p_l.$$

These flows, flows experiencing congestion, are then progressively divided into groups based on the freqest, PDV, and skewest summary statistics. The process proceeds according to the following steps:

2. Group flows whose difference in sorted freqest is less than a threshold:

$$\text{diff}(\text{freqest}) < p_f$$

3. Group flows whose difference in sorted $E_N(\text{PDV})$ is less than a threshold:

$$\text{diff}(E_N(\text{PDV})) < (p_pdv * E_N(\text{PDV}))$$

4. Group flows whose difference in sorted $E_N(\text{skewest})$ or PL_NT is less than a threshold:

$$\text{if } PL_NT < p_l$$

$$\text{diff}(E_N(\text{skewness})) < p_s$$

otherwise

$$\text{diff}(PL_NT) < p_d$$

This procedure involves sorting the groups, according to the measure being used to divide them. It is simple to implement, and efficient for small numbers of flows, such as are expected in RTCWEB.

3.2.2. Using the flow group signal

A grouping decisions is made every T. Network conditions can cause bottlenecks to fluctuate. A coupled congestion controller MAY decide only to couple groups that remain stable, say grouped together 90% of the time, depending on its objectives. Recommendations concerning

this are beyond the scope of this draft and will be specific to the coupled congestion controllers objectives.

4. Measuring OWD

This section discusses the OWD measurements required for this algorithm to detect shared bottlenecks.

The SBD mechanism described in this draft relies on differences between OWD measurements to avoid the practical problems with measuring absolute OWD (see [Hayes-LCN14] section IIIC). Since all summary statistics are relative to the mean OWD and sender/receiver clock offsets are approximately constant over the measurement periods, the offset is subtracted out in the calculation.

4.1. Time stamp resolution

The SBD mechanism requires timing information precise enough to be able to make comparisons. As a rule of thumb, the time resolution should be less than one hundredth of a typical paths range of delays. In general, the lower the time resolution, the more care that needs to be taken to ensure rounding errors don't bias the skewness calculation.

Typical RTP media flows use sub-millisecond timers, which should be adequate in most situations.

5. Acknowledgements

This work was part-funded by the European Community under its Seventh Framework Programme through the Reducing Internet Transport Latency (RITE) project (ICT-317700). The views expressed are solely those of the authors.

6. IANA Considerations

This memo includes no request to IANA.

7. Security Considerations

The security considerations of RFC 3550 [RFC3550], RFC 4585 [RFC4585], and RFC 5124 [RFC5124] are expected to apply.

Non-authenticated RTCP packets carrying shared bottleneck indications

and summary statistics could attackers to alter the bottleneck sharing characteristics for private gain or disruption of other parties communication.

8. References

8.1. Normative References

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

8.2. Informative References

- [Hayes-LCN14]
Hayes, D., Ferlin, S., and M. Welzl, "Practical Passive Shared Bottleneck Detection using Shape Summary Statistics", Proc. the IEEE Local Computer Networks (LCN) p150-158, September 2014, <http://heim.ifi.uio.no/davihay/hayes14_pract_passiv_shared_bottl_detec-abstract.html>.
- [I-D.welzl-rmcat-coupled-cc]
Welzl, M., Islam, S., and S. Gjessing, "Coupled congestion control for RTP media", draft-welzl-rmcat-coupled-cc-03 (work in progress), May 2014.
- [ITU-Y1540]
ITU-T, "Internet protocol data communication service - IP packet transfer and availability performance parameters", Series Y: Global Information Infrastructure, Internet Protocol Aspects and Next-Generation Networks , March 2011, <<http://www.itu.int/rec/T-REC-Y.1540-201103-I/en>>.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.
- [RFC5124] Ott, J. and E. Carrara, "Extended Secure RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/SAVPF)", RFC 5124, February 2008.

[RFC5481] Morton, A. and B. Claise, "Packet Delay Variation Applicability Statement", RFC 5481, March 2009.

[RFC6817] Shalunov, S., Hazel, G., Iyengar, J., and M. Kuehlewind, "Low Extra Delay Background Transport (LEDBAT)", RFC 6817, December 2012.

Authors' Addresses

David Hayes (editor)
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Phone: +47 2284 5566
Email: davihay@ifi.uio.no

Simone Ferlin
Simula Research Laboratory
P.O.Box 134
Lysaker, 1325
Norway

Phone: +47 4072 0702
Email: ferlin@simula.no

Michael Welzl
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Phone: +47 2285 2420
Email: michawe@ifi.uio.no

RTP Media Congestion Avoidance
Techniques
Internet-Draft
Intended status: Experimental
Expires: September 4, 2015

D. Hayes, Ed.
University of Oslo
S. Ferlin
Simula Research Laboratory
M. Welzl
University of Oslo
March 3, 2015

Shared Bottleneck Detection for Coupled Congestion Control for RTP
Media.
draft-hayes-rmcat-sbd-02

Abstract

This document describes a mechanism to detect whether end-to-end data flows share a common bottleneck. It relies on summary statistics that are calculated by a data receiver based on continuous measurements and regularly fed to a grouping algorithm that runs wherever the knowledge is needed. This mechanism complements the coupled congestion control mechanism in draft-welzl-rmcat-coupled-cc.

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 September 4, 2015.

Copyright 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
1.1.	The signals	3
1.1.1.	Packet Loss	3
1.1.2.	Packet Delay	3
1.1.3.	Path Lag	4
2.	Definitions	4
2.1.	Parameter Values	5
3.	Mechanism	6
3.1.	Key metrics and their calculation	7
3.1.1.	Mean delay	7
3.1.2.	Skewness Estimate	8
3.1.3.	Variance Estimate	9
3.1.4.	Oscillation Estimate	9
3.1.5.	Packet loss	10
3.2.	Flow Grouping	10
3.2.1.	Flow Grouping Algorithm	10
3.2.2.	Using the flow group signal	12
3.3.	Removing Noise from the Estimates	12
3.3.1.	Oscillation noise	12
3.3.2.	Clock drift	13
3.3.3.	Bias in the skewness measure	14
3.4.	Reducing lag and Improving Responsiveness	14
3.4.1.	Improving the response of the skewness estimate	15
3.4.2.	Improving the response of the variance estimate	15
4.	Measuring OWD	16
4.1.	Time stamp resolution	16
5.	Acknowledgements	16
6.	IANA Considerations	16
7.	Security Considerations	16
8.	Change history	17
9.	References	17
9.1.	Normative References	17
9.2.	Informative References	17
	Authors' Addresses	18

1. Introduction

In the Internet, it is not normally known if flows (e.g., TCP connections or UDP data streams) traverse the same bottlenecks. Even flows that have the same sender and receiver may take different paths and share a bottleneck or not. Flows that share a bottleneck link usually compete with one another for their share of the capacity. This competition has the potential to increase packet loss and delays. This is especially relevant for interactive applications that communicate simultaneously with multiple peers (such as multi-party video). For RTP media applications such as RTCWEB, [I-D.welzl-rmcat-coupled-cc] describes a scheme that combines the congestion controllers of flows in order to honor their priorities and avoid unnecessary packet loss as well as delay. This mechanism relies on some form of Shared Bottleneck Detection (SBD); here, a measurement-based SBD approach is described.

1.1. The signals

The current Internet is unable to explicitly inform endpoints as to which flows share bottlenecks, so endpoints need to infer this from whatever information is available to them. The mechanism described here currently utilises packet loss and packet delay, but is not restricted to these.

1.1.1. Packet Loss

Packet loss is often a relatively rare signal. Therefore, on its own it is of limited use for SBD, however, it is a valuable supplementary measure when it is more prevalent.

1.1.2. Packet Delay

End-to-end delay measurements include noise from every device along the path in addition to the delay perturbation at the bottleneck device. The noise is often significantly increased if the round-trip time is used. The cleanest signal is obtained by using One-Way-Delay (OWD).

Measuring absolute OWD is difficult since it requires both the sender and receiver clocks to be synchronised. However, since the statistics being collected are relative to the mean OWD, a relative OWD measurement is sufficient. Clock drift is not usually significant over the time intervals used by this SBD mechanism (see [RFC6817] A.2 for a discussion on clock drift and OWD measurements). However, in circumstances where it is significant, Section 3.3.2 outlines a way of adjusting the calculations to cater for it.

Each packet arriving at the bottleneck buffer may experience very different queue lengths, and therefore different waiting times. A single OWD sample does not, therefore, characterize the path well. However, multiple OWD measurements do reflect the distribution of delays experienced at the bottleneck.

1.1.3. Path Lag

Flows that share a common bottleneck may traverse different paths, and these paths will often have different base delays. This makes it difficult to correlate changes in delay or loss. This technique uses the long term shape of the delay distribution as a base for comparison to counter this.

2. Definitions

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

Acronyms used in this document:

OWD -- One Way Delay
PDV -- Packet Delay Variation
RTT -- Round Trip Time
SBD -- Shared Bottleneck Detection

Conventions used in this document:

T -- the base time interval over which measurements are made.
N -- the number of base time, T, intervals used in some calculations.
sum_T(...) -- summation of all the measurements of the variable in parentheses taken over the interval T
sum(...) -- summation of terms of the variable in parentheses
sum_N(...) -- summation of N terms of the variable in parentheses

sum_NT(...) -- summation of all measurements taken over the interval $N \cdot T$

E_T(...) -- the expectation or mean of the measurements of the variable in parentheses over T

E_N(...) -- The expectation or mean of the last N values of the variable in parentheses

E_M(...) -- The expectation or mean of the last M values of the variable in parentheses, where $M \leq N$.

max_T(...) -- the maximum recorded measurement of the variable in parentheses taken over the interval T

min_T(...) -- the minimum recorded measurement of the variable in parentheses taken over the interval T

num_T(...) -- the count of measurements of the variable in parentheses taken in the interval T

num_VM(...) -- the count of valid values of the variable in parentheses given M records

PC -- a boolean variable indicating the particular flow was identified as experiencing congestion in the previous interval T (i.e. Previously Congested)

CD_T -- an estimate of the effect of Clock Drift on the mean OWD per T

CD_Adj(...) -- Mean OWD adjusted for clock drift

p_l, p_f, p_pdv, c_s, c_h, p_s, p_d, p_v -- various thresholds used in the mechanism.

N, M, and F -- number of values (calculated over T).

2.1. Parameter Values

Reference [Hayes-LCN14] uses $T=350\text{ms}$, $N=50$, $p_l = 0.1$. The other parameters have been tightened to reflect minor enhancements to the algorithm outlined in Section 3.3: $c_s = -0.01$, $p_f = p_s = p_d = 0.1$, $p_pdv = 0.2$, $p_v = 0.2$. $M=50$, $F=10$, and $c_h = 0.3$ are additional parameters defined in the document. These are values that seem to work well over a wide range of practical Internet conditions, but are the subject of ongoing tests.

3. Mechanism

The mechanism described in this document is based on the observation that the distribution of delay measurements of packets from flows that share a common bottleneck have similar shape characteristics. These shape characteristics are described using 3 key summary statistics:

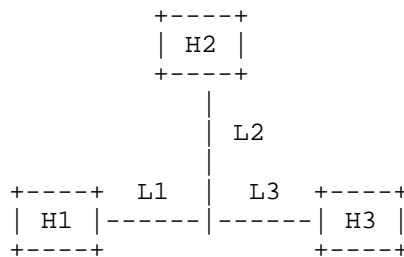
variance (estimate var_est, see Section 3.1.3)

skewness (estimate skew_est, see Section 3.1.2)

oscillation (estimate freq_est, see Section 3.1.4)

with packet loss (estimate pkt_loss, see Section 3.1.5) used as a supplementary statistic.

Summary statistics help to address both the noise and the path lag problems by describing the general shape over a relatively long period of time. This is sufficient for their application in coupled congestion control for RTP Media. They can be signalled from a receiver, which measures the OWD and calculates the summary statistics, to a sender, which is the entity that is transmitting the media stream. An RTP Media device may be both a sender and a receiver. SBD can be performed at either Sender or receiver or both.



A network with 3 hosts (H1, H2, H3) and 3 links (L1, L2, L3).

Figure 1

In Figure 1, there are two possible cases for shared bottleneck detection: a sender-based and a receiver-based case.

1. Sender-based: consider a situation where host H1 sends media streams to hosts H2 and H3, and L1 is a shared bottleneck. H2 and H3 measure the OWD and calculate summary statistics, which they send to H1 every T. H1, having this knowledge, can determine the shared bottleneck and accordingly control the send rates.

2. Receiver-based: consider that H2 is also sending media to H3, and L3 is a shared bottleneck. If H3 sends summary statistics to H1 and H2, neither H1 nor H2 alone obtain enough knowledge to detect this shared bottleneck; H3 can however determine it by combining the summary statistics related to H1 and H2, respectively. This case is applicable when send rates are controlled by the receiver; then, the signal from H3 to the senders contains the sending rate.

A discussion of the required signalling for the receiver-based case is beyond the scope of this document. For the sender-based case, the messages and their data format will be defined here in future versions of this document. We envision that an initialization message from the sender to the receiver could specify which key metrics are requested out of a possibly extensible set (`pkt_loss`, `var_est`, `skew_est`, `freq_est`). The grouping algorithm described in this document requires all four of these metrics, and receivers **MUST** be able to provide them, but future algorithms may be able to exploit other metrics (e.g. metrics based on explicit network signals). Moreover, the initialization message could specify `T`, `N`, and the necessary resolution and precision (number of bits per field).

3.1. Key metrics and their calculation

Measurements are calculated over a base interval, `T`. `T` should be long enough to provide enough samples for a good estimate of skewness, but short enough so that a measure of the oscillation can be made from `N` of these estimates. Reference [Hayes-LCN14] uses `T = 350ms` and `N=M=50`, which are values that seem to work well over a wide range of practical Internet conditions.

3.1.1. Mean delay

The mean delay is not a useful signal for comparisons between flows since flows may traverse quite different paths and clocks will not necessarily be synchronized. However, it is a base measure for the 3 summary statistics. The mean delay, `E_T(OWD)`, is the average one way delay measured over `T`.

To facilitate the other calculations, the last `N` `E_T(OWD)` values will need to be stored in a cyclic buffer along with the moving average of `E_T(OWD)`:

$$\text{mean_delay} = E_M(E_T(\text{OWD})) = \text{sum}_M(E_T(\text{OWD})) / M$$

where $M \leq N$. Generally $M=N$, setting M to be less than N allows the mechanism to be more responsive to changes, but potentially at the expense of a higher error rate (see Section 3.4 for a discussion on

improving the responsiveness of the mechanism.)

3.1.2. Skewness Estimate

Skewness is difficult to calculate efficiently and accurately. Ideally it should be calculated over the entire period ($M * T$) from the mean OWD over that period. However this would require storing every delay measurement over the period. Instead, an estimate is made over T using the previous calculation of `mean_delay`. Comparisons are made using the mean of M skew estimates (an alternative that removes bias in the mean is given in Section 3.3.3).

The skewness is estimated using two counters, counting the number of one way delay samples (OWD) above and below the mean:

$$\text{skew_est_T} = (\text{sum_T}(\text{OWD} < \text{mean_delay}) \\ - \text{sum_T}(\text{OWD} > \text{mean_delay})) / \text{num_T}(\text{OWD})$$

where

if (`OWD < mean_delay`) 1 else 0

if (`OWD > mean_delay`) 1 else 0

`skew_est_T` is a number between -1 and 1

$$\text{skew_est} = E_M(\text{skew_est_T}) = \text{sum_M}(\text{skew_est_T}) / M$$

For implementation ease, `mean_delay` does not include the mean of the current T interval.

Note: Care must be taken when implementing the comparisons to ensure that rounding does not bias `skew_est`. It is important that the mean is calculated with a higher precision than the samples.

3.1.3. Variance Estimate

Packet Delay Variation (PDV) ([RFC5481] and [ITU-Y1540]) is used as an estimator of the variance of the delay signal. We define PDV as follows:

$$\text{PDV} = \text{PDV_max} = \text{max_T(OWD)} - \text{E_T(OWD)}$$

$$\text{var_est} = \text{E_M(PDV)} = \text{sum_M(PDV)} / \text{M}$$

This modifies PDV as outlined in [RFC5481] to provide a summary statistic version that best aids the grouping decisions of the algorithm (see [Hayes-LCN14] section IVB).

The use of $\text{PDV} = \text{PDV_min} = \text{E_T(OWD)} - \text{min_T(OWD)}$ is currently being investigated as an alternative that is less sensitive to noise. The drawback of using PDV_min is that it does not distinguish between groups of flows with similar values of skew_est as well as PDV_max (see [Hayes-LCN14] section IVB).

3.1.4. Oscillation Estimate

An estimate of the low frequency oscillation of the delay signal is calculated by counting and normalising the significant mean, E_T(OWD) , crossings of mean_delay:

$$\text{freq_est} = \text{number_of_crossings} / \text{N}$$

Where

we define a significant mean crossing as a crossing that extends $p_v * \text{var_est}$ from mean_delay. In our experiments we have found that $p_v = 0.2$ is a good value.

Freq_est is a number between 0 and 1. Freq_est can be approximated incrementally as follows:

With each new calculation of E_T(OWD) a decision is made as to whether this value of E_T(OWD) significantly crosses the current long term mean, mean_delay, with respect to the previous significant mean crossing.

A cyclic buffer, last_N_crossings, records a 1 if there is a significant mean crossing, otherwise a 0.

The counter, `number_of_crossings`, is incremented when there is a significant mean crossing and subtracted from when a non-zero value is removed from the `last_N_crossings`.

This approximation of `freq_est` was not used in [Hayes-LCN14], which calculated `freq_est` every `T` using the current `E_N(E_T(OWD))`. Our tests show that this approximation of `freq_est` yields results that are almost identical to when the full calculation is performed every `T`.

3.1.5. Packet loss

The proportion of packets lost is used as a supplementary measure:

$$\text{pkt_loss} = \text{sum_NT}(\text{lost packets}) / \text{sum_NT}(\text{total packets})$$

Note: When `pkt_loss` is small it is very variable, however, when `pkt_loss` is high it becomes a stable measure for making grouping decisions.

3.2. Flow Grouping

3.2.1. Flow Grouping Algorithm

The following grouping algorithm is RECOMMENDED for SBD in the RMCAT context and is sufficient and efficient for small to moderate numbers of flows. For very large numbers of flows (e.g. hundreds), a more complex clustering algorithm may be substituted.

Since no single metric is precise enough to group flows (due to noise), the algorithm uses multiple metrics. Each metric offers a different "view" of the bottleneck link characteristics, and used together they enable a more precise grouping of flows than would otherwise be possible.

Flows determined to be experiencing congestion are successively divided into groups based on `freq_est`, `var_est`, and `skew_est`.

The first step is to determine which flows are experiencing congestion. This is important, since if a flow is not experiencing congestion its delay based metrics will not describe the bottleneck, but the "noise" from the rest of the path. Skewness, with proportion of packets loss as a supplementary measure, is used to do this:

1. Grouping will be performed on flows where:

```
skew_est < c_s
|| ( skew_est < c_h && PC )
|| pkt_loss > p_l
```

The parameter `c_s` controls how sensitive the mechanism is in detecting congestion. `C_s = 0.0` was used in [Hayes-LCN14]. A value of `c_s = 0.05` is a little more sensitive, and `c_s = -0.05` is a little less sensitive. `C_h` controls the hysteresis on flows that were grouped as experiencing congestion last time.

These flows, flows experiencing congestion, are then progressively divided into groups based on the `freq_est`, `PDV`, and `skew_est` summary statistics. The process proceeds according to the following steps:

2. Group flows whose difference in sorted `freq_est` is less than a threshold:

```
diff(freq_est) < p_f
```

3. Group flows whose difference in sorted `E_N(PDV)` (highest to lowest) is less than a threshold:

```
diff(var_est) < (p_pdv * var_est)
```

The threshold, `(p_pdv * var_est)`, is with respect to the highest value in the difference.

4. Group flows whose difference in sorted `skew_est` or `pkt_loss` is less than a threshold:

```
if pkt_loss < p_l
```

```
diff(skew_est) < p_s
```

```
otherwise
```

```
diff(pkt_loss) < (p_d * pkt_loss)
```

The threshold, `(p_d * pkt_loss)`, is with respect to the highest value in the difference.

This procedure involves sorting estimates from highest to lowest. It is simple to implement, and efficient for small numbers of flows, such as are expected in RTCWEB.

3.2.2. Using the flow group signal

A grouping decisions is made every T from the second T, though they will not attain their full design accuracy until after the N'th T interval.

Network conditions, and even the congestion controllers, can cause bottlenecks to fluctuate. A coupled congestion controller MAY decide only to couple groups that remain stable, say grouped together 90% of the time, depending on its objectives. Recommendations concerning this are beyond the scope of this draft and will be specific to the coupled congestion controllers objectives.

3.3. Removing Noise from the Estimates

The following describe small changes to the calculation of the key metrics that help remove noise from them. Currently these "tweaks" are described separately to keep the main description succinct. In future revisions of the draft these enhancements may replace the original key metric calculations.

3.3.1. Oscillation noise

When a path has no congestion, the PDV will be very small and the recorded significant mean crossings will be the result of path noise. Thus up to N-1 meaningless mean crossings can be a source of error at the point a link becomes a bottleneck and flows traversing it begin to be grouped.

To remove this source of noise from `freq_est`:

1. Set the current PDV to `PDV = NaN` (a value representing an invalid record, ie Not a Number) for flows that are deemed to not be experiencing congestion by the first `skew_est` based grouping test (see Section 3.2.1).
2. Then `var_est = sum_M(PDV != NaN) / num_VM(PDV)`
3. For `freq_est`, only record a significant mean crossing if flow is experiencing congestion.

These three changes will remove the non-congestion noise from `freq_est`.

3.3.2. Clock drift

Generally sender and receiver clock drift will be too small to cause significant errors in the estimators. `Skew_est` is most sensitive to this type of noise. In circumstances where clock drift is high, making $M < N$ can reduce this error.

A better method is to estimate the effect the clock drift is having on the $E_N(E_T(OWD))$, and then adjust `mean_delay` accordingly. A simple method of doing this follows:

First divide the N $E_T(OWD)$ values into two halves ($N/2$ in each) -- old and new.

Calculate a mean of the old half:

$$\text{Older_mean} = E_{\text{old}}(E_T(OWD)) / N/2$$

Calculate a mean of the new (most recent) half:

$$\text{Newer_mean} = E_{\text{new}}(E_T(OWD)) / N/2$$

A linear estimate of the Clock Drift per T estimates is:

$$\text{CD}_T = (\text{Newer_mean} - \text{Older_mean}) / N/2$$

An adjusted mean estimate then is:

$$\text{mean_delay} = \text{CD_Adj}(E_M(E_T(OWD))) = E_M(E_T(OWD)) + \text{CD}_T * M/2$$

`CD_Adj` can be thought of as a prediction of what the long term mean will be in the current measurement period T . It is used as the basis for `skew_est` and `freq_est`.

3.3.3. Bias in the skewness measure

If successive calculations of `skew_est` are made with very different numbers of samples (`num_T(OWD)`), the simple calculation of `E_M(skew_est)` used for grouping decisions will be biased by the intervals that have few samples. This bias can be corrected if necessary as follows.

```
skew_base_T = (sum_T(OWD < mean_delay
               - sum_T(OWD > mean_delay)

skew_est = sum_MT(skew_base_T)/num_MT(OWD)
```

This calculation requires slightly more state, since an implementation will need to maintain two cyclic buffers storing `skew_base_T` and `num_T(OWD)` respectively to manage the rolling summations (note only one cyclic buffer is needed for the calculation of `skew_est` outlined previously).

3.4. Reducing lag and Improving Responsiveness

Measurement based shared bottleneck detection makes decisions in the present based on what has been measured in the past. This means that there is always a lag in responding to changing conditions. This mechanism is based on summary statistics taken over ($N \cdot T$) seconds. This mechanism can be made more responsive to changing conditions by:

1. Reducing N and/or M -- but at the expense of less accurate metrics, and/or
2. Exploiting the fact that more recent measurements are more valuable than older measurements and weighting them accordingly.

Although more recent measurements are more valuable, older measurements are still needed to gain an accurate estimate of the distribution descriptor we are measuring. Unfortunately, the simple exponentially weighted moving average weights drop off too quickly for our requirements and have an infinite tail. A simple linearly declining weighted moving average also does not provide enough weight to the most recent measurements. We propose a piecewise linear distribution of weights, such that the first section (samples 1:F) is flat as in a simple moving average, and the second section (samples F+1:M) is linearly declining weights to the end of the averaging window. We choose integer weights, which allows incremental calculation without introducing rounding errors.

3.4.1. Improving the response of the skewness estimate

The weighted moving average for `skew_est`, based on `skew_est` in Section 3.3.3, can be calculated as follows:

$$\begin{aligned} \text{skew_est} = & ((M-F+1)*\text{sum}(\text{skew_base_T}(1:F)) \\ & + \text{sum}([(M-F):1].*\text{skew_base_T}(F+1:M))) \\ & / ((M-F+1)*\text{sum}(\text{numsampT}(1:F)) \\ & + \text{sum}([(M-F):1].*\text{numsampT}(F+1:M))) \end{aligned}$$

where `numsampT` is an array of the number of OWD samples in each `T` (ie `num_T(OWD)`), and `numsampT(1)` is the most recent; `skew_base_T(1)` is the most recent calculation of `skew_base_T`; `1:F` refers to the integer values 1 through to `F`, and `[(M-F):1]` refers to an array of the integer values `(M-F)` declining through to 1; and `.*` is the array scalar dot product operator.

3.4.2. Improving the response of the variance estimate

The weighted moving average for `var_est` can be calculated as follows:

$$\begin{aligned} \text{var_est} = & ((M-F+1)*\text{sum}(\text{PDV}(1:F)) + \text{sum}([(M-F):1].*\text{PDV}(F+1:M))) \\ & / (F*(M-F+1) + \text{sum}([(M-F):1])) \end{aligned}$$

where `1:F` refers to the integer values 1 through to `F`, and `[(M-F):1]` refers to an array of the integer values `(M-F)` declining through to 1; and `.*` is the array scalar dot product operator. When removing oscillation noise (see Section 3.3.1) this calculation must be adjusted to allow for invalid PDV records.

4. Measuring OWD

This section discusses the OWD measurements required for this algorithm to detect shared bottlenecks.

The SBD mechanism described in this draft relies on differences between OWD measurements to avoid the practical problems with measuring absolute OWD (see [Hayes-LCN14] section IIIC). Since all summary statistics are relative to the mean OWD and sender/receiver clock offsets should be approximately constant over the measurement periods, the offset is subtracted out in the calculation.

4.1. Time stamp resolution

The SBD mechanism requires timing information precise enough to be able to make comparisons. As a rule of thumb, the time resolution should be less than one hundredth of a typical path's range of delays. In general, the lower the time resolution, the more care that needs to be taken to ensure rounding errors do not bias the skewness calculation.

Typical RTP media flows use sub-millisecond timers, which should be adequate in most situations.

5. Acknowledgements

This work was part-funded by the European Community under its Seventh Framework Programme through the Reducing Internet Transport Latency (RITE) project (ICT-317700). The views expressed are solely those of the authors.

6. IANA Considerations

This memo includes no request to IANA.

7. Security Considerations

The security considerations of RFC 3550 [RFC3550], RFC 4585 [RFC4585], and RFC 5124 [RFC5124] are expected to apply.

Non-authenticated RTCP packets carrying shared bottleneck indications and summary statistics could allow attackers to alter the bottleneck sharing characteristics for private gain or disruption of other parties communication.

8. Change history

Changes made to this document:

- 01->02 : New section describing improvements to the key metric calculations that help to remove noise, bias, and reduce lag. Some revisions to the notation to make it clearer. Some tightening of the thresholds.
- 00->01 : Revisions to terminology for clarity

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.

9.2. Informative References

- [Hayes-LCN14]
Hayes, D., Ferlin, S., and M. Welzl, "Practical Passive Shared Bottleneck Detection using Shape Summary Statistics", Proc. the IEEE Local Computer Networks (LCN) p150-158, September 2014, <http://heim.ifi.uio.no/davihay/hayes14_pract_passiv_shared_bottl_detec-abstract.html>.
- [I-D.welzl-rmcat-coupled-cc]
Welzl, M., Islam, S., and S. Gjessing, "Coupled congestion control for RTP media", draft-welzl-rmcat-coupled-cc-04 (work in progress), October 2014.
- [ITU-Y1540]
ITU-T, "Internet Protocol Data Communication Service - IP Packet Transfer and Availability Performance Parameters", Series Y: Global Information Infrastructure, Internet Protocol Aspects and Next-Generation Networks , March 2011, <<http://www.itu.int/rec/T-REC-Y.1540-201103-I/en>>.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control

Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.

[RFC5124] Ott, J. and E. Carrara, "Extended Secure RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/SAVPF)", RFC 5124, February 2008.

[RFC5481] Morton, A. and B. Claise, "Packet Delay Variation Applicability Statement", RFC 5481, March 2009.

[RFC6817] Shalunov, S., Hazel, G., Iyengar, J., and M. Kuehlewind, "Low Extra Delay Background Transport (LEDBAT)", RFC 6817, December 2012.

Authors' Addresses

David Hayes (editor)
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Phone: +47 2284 5566
Email: davihay@ifi.uio.no

Simone Ferlin
Simula Research Laboratory
P.O.Box 134
Lysaker, 1325
Norway

Phone: +47 4072 0702
Email: ferlin@simula.no

Michael Welzl
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Phone: +47 2285 2420
Email: michawe@ifi.uio.no

RMCAT WG
Internet-Draft
Intended status: Informational
Expires: February 8, 2015

M. Zanaty
Cisco
V. Singh
Aalto University
S. Nandakumar
Cisco
Z. Sarker
Ericsson AB
August 7, 2014

RTP Application Interaction with Congestion Control
draft-ietf-rmcat-app-interaction-00

Abstract

Interactive real-time media applications that use the Real-time Transport Protocol (RTP) over the User Datagram Protocol (UDP) must use congestion control techniques above the UDP layer since it provides none. This memo describes the interactions and conceptual interfaces necessary between the application components that relate to congestion control, including the RTP layer, the higher-level media codec control layer, and the lower-level transport interface, as well as components dedicated to congestion control functions.

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 February 8, 2015.

Copyright Notice

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

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
2. Key Words for Requirements	3
3. Conceptual Model	4
4. Implementation Model	5
5. Interfaces and Interactions	6
5.1. Config - Codec Interactions	6
5.2. Config - RTP/RTCP Interactions	6
5.3. Codec - RTP Interactions	6
5.4. Codec - CC Interactions	7
5.5. RTP - CC Interactions	9
5.6. CC - UDP Interactions	9
5.7. CC - Shared State Interactions	10
6. Acknowledgements	10
7. IANA Considerations	10
8. Security Considerations	10
9. References	11
9.1. Normative References	11
9.2. Informative References	12
Authors' Addresses	13

1. Introduction

Interactive real-time media applications most commonly use RTP [RFC3550] over UDP [RFC0768]. Since UDP provides no form of congestion control, which is essential for any application deployed on the Internet, these RTP applications have historically implemented one of the following options at the application layer to address their congestion control requirements.

1. For media with relatively low packet rates and bit rates, such as many speech codecs, some applications use a simple form of congestion control that stops transmission permanently or temporarily after observing significant packet loss over a significant period of time, similar to the RTP circuit breakers [I-D.ietf-avtcore-rtp-circuit-breakers].

2. Some applications have no explicit congestion control, despite the clear requirements in RTP and its profiles AVP [RFC3551] and AVPF [RFC4585], under the expectation that users will terminate media flows that are significantly impaired by congestion (in essence, human circuit breakers).
3. For media with substantially higher packet rates and bit rates, such as many video codecs, various non-standard congestion control techniques are often used to adapt transmission rate based on receiver feedback.
4. Some experimental applications use standardized techniques such as TCP-Friendly Rate Control (TFRC) [RFC5348]. However, for various reasons, these have not been widely deployed.

The RTP Media Congestion Avoidance Techniques (RMCAT) working group was chartered to standardize appropriate and effective congestion control for RTP applications. It is expected such applications will migrate from the above historical solutions to the RMCAT solution(s).

The RMCAT requirements [I-D.ietf-rmcat-cc-requirements] include low delay, reasonably high throughput, fast reaction to capacity changes including routing or interface changes, stability without over-reaction or oscillation, fair bandwidth sharing with other instances of itself and TCP flows, sharing information across multiple flows when possible [I-D.welzl-rmcat-coupled-cc], and performing as well or better in networks which support Active Queue Management (AQM), Explicit Congestion Notification (ECN), or Differentiated Services Code Points (DSCP).

In order to meet these requirements, interactions are necessary between the application's congestion controller, the RTP layer, media codecs, other components, and the OS UDP stack. This memo discusses these interactions, presents a conceptual model of the required interfaces based on a simplified application decomposition, and proposes specific information exchange across these interfaces along with corresponding component behavior.

Note that RTP can also operate over other transports with integrated congestion control such as TCP [RFC5681] and DCCP [RFC4340], but that is beyond the scope of RMCAT and this memo.

2. Key Words for Requirements

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. Conceptual Model

It is useful to decompose an RTP application into several components to facilitate understanding and discussion of where congestion control functions operate, and how they interface with the other components. The conceptual model in Figure 1 consists of the following components.

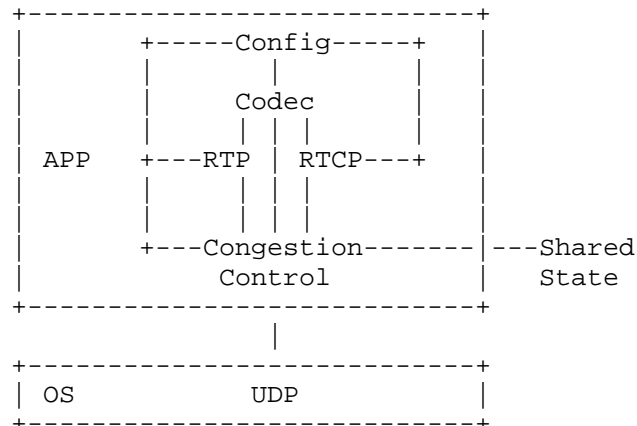


Figure 1

- o APP: Application containing one or more RTP streams and the corresponding media codecs and congestion controllers. For example, a WebRTC browser.
- o Config: Configuration specified by the application that provides the media and transport parameters, RTP and RTCP parameters and extensions, and congestion control parameters. For example, a WebRTC Javascript application may use the 'constraints' API to affect the media configuration, and SDP applications may negotiate the media and transport parameters with the remote peer. This determines the initial static configuration negotiated in session establishment. The dynamic state may differ due to congestion or other factors, but still must conform to limits established in the config.
- o Codec: Media encoder/decoder or other source/sink for the RTP payload. The codec may be, for example, a simple monaural audio format, a complex scalable video codec with several dependent layers, or a source/sink with no live encoding/decoding such as a mixer which selectively switches and forwards streams rather than mixes media.

- o RTP: Standard RTP stack functions, including media packetization / depacketization and header processing, but excluding existing extensions and possible new extensions specific to congestion control (CC) such as absolute timestamps or relative transmission time offsets in RTP header extensions. RTCP: Standard RTCP functions, including sender reports, receiver reports, extended reports, circuit breakers [I-D.ietf-avtcore-rtp-circuit-breakers], feedback messages such as NACK [RFC4585] and codec control messages such as TMMBR [RFC5104], but excluding existing extensions and possible new extensions specific to congestion control (CC) such as REMB [I-D.alvestrand-rmcat-remb] (for receiver-side CC), ACK (for sender-side CC), absolute and/or relative timestamps (for sender-side or receiver-side CC), etc.
- o Congestion Control: All functions directly responsible for congestion control, including possible new RTP/RTCP extensions, send rate computation (for sender-side CC), receive rate computation (for receiver-side CC), other statistics, and control of the UDP sockets including packet scheduling for traffic shaping / pacing.
- o Shared State: Storage and exchange of congestion control state for multiple flows within the application and beyond it.
- o OS: Operating System containing the UDP socket interface and other network functions such as ECN, DSCP, physical interface events, interface-level traffic shaping and packet scheduling, etc.

4. Implementation Model

There are advantages and drawbacks to implementing congestion control in the application layer. It avoids OS dependencies and allows for rapid experimentation, evolution and optimization for each application. However, it also puts the burden on all applications, which raises the risks of improper or divergent implementations. One motivation of this memo is to mitigate such risks by giving proper guidance on how the application components relating to congestion control should interact.

Another drawback of congestion control in the application layer is that any decomposition, including the one presented in Figure 1, is purely conceptual and illustrative, since implementations have differing designs and decompositions. Conversely, this can be viewed as an advantage to distribute congestion control functions wherever expedient without rigid interfaces. For example, they may be distributed within the RTP/RTCP stack itself, so the separate components in Figure 1 are combined into a single RTP+RTCP+CC component as shown in Figure 2.

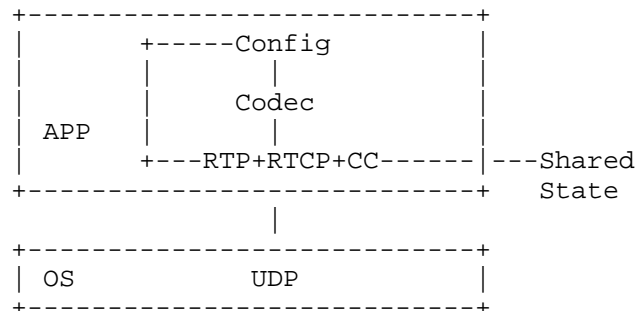


Figure 2

The conceptual model in Figure 1 will be used throughout this memo to establish clearer boundaries between functions. But actual implementations may be closer to the looser model in [Singh12].

5. Interfaces and Interactions

5.1. Config - Codec Interactions

The primary interactions between the config and the codec that are relevant to congestion control are the multiplexing of media streams [I-D.ietf-mmusic-sdp-bundle-negotiation] and RTP/RTCP [RFC5761] on the same UDP port.

The config also establishes limits for the codec such as maximum bit rate and other codec-specific parameters. For example, a video codec config often sets limits on maximum resolution and frame rate as well as bit rate.

5.2. Config - RTP/RTCP Interactions

The config establishes the negotiated RTP and RTCP attributes and extensions such as Extended Reports (XR), reduced size [RFC5506], codec control [RFC5104], transmission time [RFC5450], etc.

5.3. Codec - RTP Interactions

Packetization of codec frames into RTP packets can be an important interaction. Some network interfaces may benefit from small packet sizes well below the MTU, while others may benefit from large packets approaching the MTU. Equalizing packet sizes of a frame may also be beneficial in some cases, rather than a combination of large and small packets. For example, in some FEC schemes, the FEC bandwidth overhead depends on the largest source packet size. Equalizing the

source packet sizes can yield lower overhead than a combination of large and small packets.

5.4. Codec - CC Interactions

Allowed Rate (from CC to Codec): The max transmit rate allowed over the next time interval. The time interval may be specified or may use a default, for example, one second. The rate may be specified in bytes or packets or both. The rate must never exceed permanent limits established in session signaling such as the SDP bandwidth attribute [RFC4566] nor temporary limits in RTCP such as TMMBR [RFC5104] or REMB [I-D.alvestrand-rmcat-remb]. This is the most important interface among all components, and is always required in any RMCAT solution. In the simplest possible solution, it may be the only CC interface required.

Media Elasticity (from Codec to CC): Many live media encoders are highly elastic, often able to achieve any target bit rate within a wide range, by adapting the media quality. For example, a video encoder may support any bit rate within a range of a few tens or hundreds of kbps up to several Mbps, with rate changes registering as fast as the next video frame, although there may be limitations in the frequency of changes. Other encoders may be less elastic, supporting a narrower rate range, coarser granularity of rate steps, slower reaction to rate changes, etc. Other media, particularly some audio codecs, may be fully inelastic with a single fixed rate. CC can beneficially use codec elasticity, if provided, to plan Allowed Rate changes, especially when there are multiple flows sharing CC state and bandwidth.

Startup Ramp (from Codec to CC, and from CC to Codec): Startup is an important moment in a conversation. Rapid rate adaptation during startup is therefore important. The codec should minimize its startup media rate as much as possible without adversely impacting the user experience, and support a strategy for rapid rate ramp. The CC should allow the highest startup media rate as possible without adversely impacting network conditions, and also support rapid rate ramp until stabilizing on the available bandwidth. Startup can be viewed as a negotiation between the codec and the CC. The codec requests a startup rate and ramp, and the CC responds with the allowable parameters which may be lower/slower. The RMCAT requirements also include the possibility of bandwidth history to further accelerate or even eliminate startup ramp time. While this is highly desirable from an application viewpoint, it may be less acceptable to network operators, since it is in essence a gamble on current congestion state matching historical state, with the potential for significant congestion contribution if the gamble was

wrong. Note that startup can often commence before user interaction or conversation to reduce the chance of clipped media.

Delay Tolerance (from Codec to CC): An ideal CC will always minimize delay and target zero. However, real solutions often need a real non-zero delay tolerance. The codec should provide an absolute delay tolerance, perhaps expressed as an impairment factor to mix with other metrics.

Loss Tolerance (from Codec to CC): An ideal CC will always minimize packet loss and target zero. However, real solutions often need a real non-zero loss tolerance. The codec should provide an absolute loss tolerance, perhaps expressed as an impairment factor to mix with other metrics. Note this is unrecoverable post-repair loss after retransmission or forward error correction.

Throughput Sensitivity (from Codec to CC): An ideal CC will always maximize throughput. However, real solutions often need a trade-off between throughput and other metrics such as delay or loss. The codec should provide throughput sensitivity, perhaps expressed as an impairment factor (for low throughputs) to mix with other metrics.

Rate Stability (from Codec to CC): The CC algorithm must strike a balance between rate stability and fast reaction to changes in available bandwidth. The codec should provide its preference for rate stability versus fast and frequent reaction to rate changes, perhaps expressed as an impairment factor (for high rate variance over short timescales) to mix with other metrics.

Forward Error Correction (FEC): Simple FEC schemes like XOR Parity codes [RFC5109] may not handle consecutive or burst loss well. More complex FEC schemes like Reed-Solomon [RFC6865] or Raptor [RFC6330] codes are more effective at handling bursty loss. The sensitivity to packet loss therefore depends on the media (source) encoding as well as the FEC (channel) encoding, and this sensitivity may differ for different loss patterns like random, periodic, or consecutive loss. Expressing this sensitivity to the congestion controller may help it choose the right balance between optimizing for throughput versus low loss.

Probing for Available Bandwidth: FEC can also be used to probe for additional available bandwidth, if the application desires a higher target rate than the current rate. FEC is preferable to synthetic probes since any contribution to congestion by the FEC probe will not impact the post-repair loss rate of the source media flow while synthetic probes may adversely affect the loss rate [Nagy14]. Note that any use of FEC or retransmission must ensure that the total flow

of all packets including FEC, retransmission and original media never exceeds the Allowed Rate.

5.5. RTP - CC Interactions

RTP Circuit Breakers: The intent behind RTP circuit breakers [I-D.ietf-avtcore-rtp-circuit-breakers] is to provide a kill switch of last resort, not true congestion control. The breakers should never trip when an effective congestion control is operating. This may impose some boundaries on RMCAT solutions to ensure the congestion control never approaches situations which may trigger the breakers.

RTCP Feedback: The primary method of communicating CC information is RTCP.

RTP Header Extensions: While RTCP is likely to be the primary carrier of CC feedback, the RMCAT requirements also include the possibility of using RTP header extensions in bidirectional flows for CC feedback. Transmission time [RFC5450], or possibly absolute time, also use header extensions, as would any per packet priority markings expected to survive across different networks and administrative domains.

5.6. CC - UDP Interactions

Pacing / Shaping: Simple pacing / shaping strategies delay the transmission of packets to equalize inter-packet time intervals, assuming the bottleneck is most sensitive to packet rate. More complex pacing strategies may go beyond simple even distribution of transmission times. For example, Sprout [Winstein13] attempts to predict the optimal transmission time (and rate) with the highest probability of success for each packet based on channel statistics. Pacing may be always on, or adaptively enabled / disabled based on congestion state to minimize delay. Pacing may be performed within the CC for a single flow or across multiple flows. It may also be performed across all or selective traffic over the network interface if the OS supports interface-level traffic shaping.

Detection of Transport Capabilities: The CC can query the OS for useful transport capabilities such as ECN, DSCP, traffic shaping, etc. This may also aid upper layers in making better decisions such as whether or not to multiplex media streams. For example, if audio can be given differentiated network treatment from video when using separate ports.

ECN: If the OS and transport path support ECN, the CC can react faster than a loss-based CC and more reliably to congestion onset and abatement.

DSCP: If the OS and transport path support DSCP, the CC can map per-packet priority from RTP header extensions to DSCP (and layer 2 QoS if available) for better network handling under congestion.

AQM: If AQM is present in the bottleneck, and working effectively, there should be little or no excess delay observed when varying the transmission rate. The loss of such delay signals may hinder the performance of congestion control algorithms that are highly dependent on delay variation for adapting transmission rate. If the application has knowledge of the presence of AQM, through any means which are beyond the scope of this memo, it should communicate this to the CC. The CC may use this to alter its signal collection and rate adaptation strategies. The CC must not rely solely on this as a reliable indicator. It must continue to monitor statistics to validate this application hint, and react appropriately if the statistics suggest different network behavior.

5.7. CC - Shared State Interactions

Multiple Flows: Sharing state across multiple flows within the application can yield better CC decisions. Sharing state across even more flows beyond the application can yield even better CC decisions. The actual benefits and mechanisms of state sharing and coupled CC are described in [I-D.welzl-rmcat-coupled-cc].

Weighted Fairness: An important consideration in CC of multiple flows is their relative application-specified weights. Within an application, it is likely the different flows have different rate requirements, so equal bandwidth sharing may not be fair nor desirable, and weighted fairness may be required.

6. Acknowledgements

The RMCAT design team discussions contributed to this memo.

7. IANA Considerations

This memo includes no request to IANA.

8. Security Considerations

Amplification attacks often use UDP traffic to launch denial of service attacks. Attackers may attempt to subvert congestion control protocols in UDP applications to launch amplification attacks by

signaling more bandwidth than is actually available. For example, sending a victim a forged REMB or a few fast ACKs may result in the victim sending a high rate RTP stream. Attacks on conference servers could lead to further amplification if it distributes the streams to many others. One mitigation is to use SRTCP for congestion control messages where supported. Even if SRTCP is only authenticated not encrypted, SRTCP packets should always pass authentication checks before any message contents are interpreted. Non-secure RTCP should be avoided where possible.

9. References

9.1. Normative References

- [I-D.alvestrand-rmcat-remb]
Alvestrand, H., "RTCP message for Receiver Estimated Maximum Bitrate", draft-alvestrand-rmcat-remb-03 (work in progress), October 2013.
- [I-D.ietf-avtcore-rtp-circuit-breakers]
Perkins, C. and V. Singh, "Multimedia Congestion Control: Circuit Breakers for Unicast RTP Sessions", draft-ietf-avtcore-rtp-circuit-breakers-05 (work in progress), February 2014.
- [I-D.ietf-mmusic-sdp-bundle-negotiation]
Holmberg, C., Alvestrand, H., and C. Jennings, "Multiplexing Negotiation Using Session Description Protocol (SDP) Port Numbers", draft-ietf-mmusic-sdp-bundle-negotiation-05 (work in progress), October 2013.
- [I-D.ietf-rmcat-cc-requirements]
Jesup, R., "Congestion Control Requirements For RMCAT", draft-ietf-rmcat-cc-requirements-02 (work in progress), February 2014.
- [I-D.ietf-rmcat-eval-criteria]
Singh, V. and J. Ott, "Evaluating Congestion Control for Interactive Real-time Media", draft-ietf-rmcat-eval-criteria-00 (work in progress), January 2014.
- [I-D.welzl-rmcat-coupled-cc]
Welzl, M., Islam, S., and S. Gjessing, "Coupled congestion control for RTP media", draft-welzl-rmcat-coupled-cc-02 (work in progress), October 2013.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.

- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC3551] Schulzrinne, H. and S. Casner, "RTP Profile for Audio and Video Conferences with Minimal Control", STD 65, RFC 3551, July 2003.
- [RFC4566] Handley, M., Jacobson, V., and C. Perkins, "SDP: Session Description Protocol", RFC 4566, July 2006.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.
- [RFC5104] Wenger, S., Chandra, U., Westerlund, M., and B. Burman, "Codec Control Messages in the RTP Audio-Visual Profile with Feedback (AVPF)", RFC 5104, February 2008.
- [RFC5450] Singer, D. and H. Desineni, "Transmission Time Offsets in RTP Streams", RFC 5450, March 2009.
- [RFC5506] Johansson, I. and M. Westerlund, "Support for Reduced-Size Real-Time Transport Control Protocol (RTCP): Opportunities and Consequences", RFC 5506, April 2009.
- [RFC5761] Perkins, C. and M. Westerlund, "Multiplexing RTP Data and Control Packets on a Single Port", RFC 5761, April 2010.

9.2. Informative References

- [Nagy14] Nagy, M., Singh, V., Ott, J., and L. Eggert, "Congestion Control using FEC for Conversational Multimedia Communication", Proc. of 5th ACM International Conference on Multimedia Systems , 3 2014.
- [RFC0768] Postel, J., "User Datagram Protocol", STD 6, RFC 768, August 1980.
- [RFC4340] Kohler, E., Handley, M., and S. Floyd, "Datagram Congestion Control Protocol (DCCP)", RFC 4340, March 2006.
- [RFC5109] Li, A., "RTP Payload Format for Generic Forward Error Correction", RFC 5109, December 2007.

- [RFC5348] Floyd, S., Handley, M., Padhye, J., and J. Widmer, "TCP Friendly Rate Control (TFRC): Protocol Specification", RFC 5348, September 2008.
- [RFC5681] Allman, M., Paxson, V., and E. Blanton, "TCP Congestion Control", RFC 5681, September 2009.
- [RFC6330] Luby, M., Shokrollahi, A., Watson, M., Stockhammer, T., and L. Minder, "RaptorQ Forward Error Correction Scheme for Object Delivery", RFC 6330, August 2011.
- [RFC6865] Roca, V., Cunche, M., Lacan, J., Bouabdallah, A., and K. Matsuzono, "Simple Reed-Solomon Forward Error Correction (FEC) Scheme for FECFRAME", RFC 6865, February 2013.
- [Singh12] Singh, V., Ott, J., and C. Perkins, "Congestion Control for Interactive Media: Control Loops & APIs", Proc. of IAB /IRTF Workshop on Congestion Control for Interactive RTC , 7 2012.
- [Winstein13] Winstein,, K., Sivaraman,, A., and H. Balakrishnan, "Stochastic Forecasts Achieve High Throughput and Low Delay over Cellular Networks", Proc. of the 10th USENIX Symposium on Networked Systems Design and Implementation , 4 2013.

Authors' Addresses

Mo Zanaty
Cisco
Raleigh, NC
USA

Email: mzanaty@cisco.com

Varun Singh
Aalto University
Espoo, FIN
Finland

Email: varun@comnet.tkk.fi

Suhas Nandakumar
Cisco
San Jose, CA
USA

Email: snandaku@cisco.com

Zaheduzzaman Sarker
Ericsson AB
Luleae
Sweden

Email: zaheduzzaman.sarker@ericsson.com

RMCAT WG
Internet-Draft
Intended status: Informational
Expires: January 26, 2015

V. Singh
J. Ott
Aalto University
July 25, 2014

Evaluating Congestion Control for Interactive Real-time Media
draft-ietf-rmcat-eval-criteria-02

Abstract

The Real-time Transport Protocol (RTP) is used to transmit media in telephony and video conferencing applications. This document describes the guidelines to evaluate new congestion control algorithms for interactive point-to-point real-time media.

Status of This Memo

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

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

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

This Internet-Draft will expire on January 26, 2015.

Copyright Notice

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

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1.	Introduction	2
2.	Terminology	3
3.	Metrics	3
3.1.	RTP Log Format	4
4.	Guidelines	5
4.1.	Avoiding Congestion Collapse	5
4.2.	Stability	5
4.3.	Media Traffic	5
4.4.	Start-up Behaviour	6
4.5.	Diverse Environments	6
4.6.	Varying Path Characteristics	6
4.7.	Reacting to Transient Events or Interruptions	6
4.8.	Fairness With Similar Cross-Traffic	7
4.9.	Impact on Cross-Traffic	7
4.10.	Extensions to RTP/RTCP	7
5.	Minimum Requirements for Evaluation	7
6.	Status of Proposals	7
7.	Security Considerations	8
8.	IANA Considerations	8
9.	Contributors	8
10.	Acknowledgements	8
11.	References	8
11.1.	Normative References	8
11.2.	Informative References	9
Appendix A.	Application Trade-off	10
A.1.	Measuring Quality	10
Appendix B.	Change Log	10
B.1.	Changes in draft-ietf-rmcat-eval-criteria-02	10
B.2.	Changes in draft-ietf-rmcat-eval-criteria-01	10
B.3.	Changes in draft-ietf-rmcat-eval-criteria-00	10
B.4.	Changes in draft-singh-rmcat-cc-eval-04	10
B.5.	Changes in draft-singh-rmcat-cc-eval-03	11
B.6.	Changes in draft-singh-rmcat-cc-eval-02	11
B.7.	Changes in draft-singh-rmcat-cc-eval-01	11
Authors' Addresses		11

1. Introduction

This memo describes the guidelines to help with evaluating new congestion control algorithms for interactive point-to-point real time media. The requirements for the congestion control algorithm are outlined in [I-D.ietf-rmcat-cc-requirements]). This document builds upon previous work at the IETF: Specifying New Congestion Control Algorithms [RFC5033] and Metrics for the Evaluation of Congestion Control Algorithms [RFC5166].

The guidelines proposed in the document are intended to help prevent a congestion collapse, promote fair capacity usage and optimize the media flow's throughput. Furthermore, the proposed algorithms are expected to operate within the envelope of the circuit breakers defined in [I-D.ietf-avtcore-rtp-circuit-breakers].

This document only provides broad-level criteria for evaluating a new congestion control algorithm and the working group should expect a thorough scientific study to make its decision. The results of the evaluation are not expected to be included within the internet-draft but should be cited in the document.

2. Terminology

The terminology defined in RTP [RFC3550], RTP Profile for Audio and Video Conferences with Minimal Control [RFC3551], RTCP Extended Report (XR) [RFC3611], Extended RTP Profile for RTCP-based Feedback (RTP/AVPF) [RFC4585] and Support for Reduced-Size RTCP [RFC5506] apply.

3. Metrics

[RFC5166] describes the basic metrics for congestion control. Metrics that are of interest for interactive multimedia are:

- o Throughput.
- o Minimizing oscillations in the transmission rate (stability) when the end-to-end capacity varies slowly.
- o Delay.
- o Reactivity to transient events.
- o Packet losses and discards.
- o Section 2.1 of [RFC5166] discusses the tradeoff between throughput, delay and loss.

Each experiment is expected to log every incoming and outgoing packet (the RTP logging format is described in Section 3.1). The logging can be done inside the application or at the endpoints using pcap (packet capture, e.g., tcpdump, wireshark). The following are calculated based on the information in the packet logs:

1. Sending rate, Receiver rate, Goodput
2. Packet delay

3. Packet loss
4. If using, retransmission or FEC: residual loss
5. Packets discarded from the playout or de-jitter buffer
6. Fairness or Unfairness: Experiments testing the performance of an RMCAT proposal against any cross-traffic must define its expected criteria for fairness. The "unfairness" test guideline (measured at 1s intervals) is:
 1. Does not trigger the circuit breaker.
 2. No RMCAT stream achieves more than 3 times the average throughput of the RMCAT stream with the lowest average throughput, for a case when the competing streams have similar RTTs.
 3. RTT should not grow by a factor of 3 for the existing flows when a new flow is added.
For example, see the test scenarios described in [I-D.sarker-rmcat-eval-test].
7. Convergence time: The time taken to reach a stable rate at startup, after the available link capacity changes, or when new flows get added to the bottleneck link.
8. Bandwidth Utilization, defined as ratio of the instantaneous sending rate to the instantaneous bottleneck capacity. This metric is useful when an RMCAT flow is by itself or competing with similar cross-traffic.

From the logs the statistical measures (min, max, mean, standard deviation and variance) for the whole duration or any specific part of the session can be calculated. Also the metrics (sending rate, receiver rate, goodput, latency) can be visualized in graphs as variation over time, the measurements in the plot are at 1 second intervals. Additionally, from the logs it is possible to plot the histogram or CDF of packet delay.

[Open issue (1): Using Jain-fairness index (JFI) for measuring self-fairness between RTP flows? measured at what intervals? visualized as a CDF or a timeseries? Additionally: Use JFI for comparing fairness between RTP and long TCP flows?]

3.1. RTP Log Format

The log file is tab or comma separated containing the following details:

Send or receive timestamp (unix)
RTP payload type
SSRC
RTP sequence no
RTP timestamp
marker bit
payload size

If the congestion control implements, retransmissions or FEC, the evaluation should report both packet loss (before applying error-resilience) and residual packet loss (after applying error-resilience).

4. Guidelines

A congestion control algorithm should be tested in simulation or a testbed environment, and the experiments should be repeated multiple times to infer statistical significance. The following guidelines are considered for evaluation:

4.1. Avoiding Congestion Collapse

The congestion control algorithm is expected to take an action, such as reducing the sending rate, when it detects congestion. Typically, it should intervene before the circuit breaker [I-D.ietf-avtcore-rtp-circuit-breakers] is engaged.

Does the congestion control propose any changes to (or diverge from) the circuit breaker conditions defined in [I-D.ietf-avtcore-rtp-circuit-breakers].

4.2. Stability

The congestion control should be assessed for its stability when the path characteristics do not change over time. Changing the media encoding rate estimate too often or by too much may adversely affect the application layer performance.

4.3. Media Traffic

The congestion control algorithm should be assessed with different types of media behavior, i.e., the media should contain idle and data-limited periods. For example, periods of silence for audio, varying amount of motion for video, or bursty nature of I-frames.

The evaluation may be done in two stages. In the first stage, the endpoint generates traffic at the rate calculated by the congestion controller. In the second stage, real codecs or models of video

codecs are used to mimic application-limited data periods and varying video frame sizes.

4.4. Start-up Behaviour

The congestion control algorithm should be assessed with different start-rates. The main reason is to observe the behavior of the congestion control in different test scenarios, such as when competing with varying amount of cross-traffic or how quickly does the congestion control algorithm achieve a stable sending rate.

4.5. Diverse Environments

The congestion control algorithm should be assessed in heterogeneous environments, containing both wired and wireless paths. Examples of wireless access technologies are: 802.11, GPRS, HSPA, or LTE. One of the main challenges of the wireless environments for the congestion control algorithm is to distinguish between congestion induced loss and transmission (bit-error) loss. Congestion control algorithms may incorrectly identify transmission loss as congestion loss and reduce the media encoding rate by too much, which may cause oscillatory behavior and deteriorate the users' quality of experience. Furthermore, packet loss may induce additional delay in networks with wireless paths due to link-layer retransmissions.

4.6. Varying Path Characteristics

The congestion control algorithm should be evaluated for a range of path characteristics such as, different end-to-end capacity and latency, varying amount of cross traffic on a bottleneck link and a router's queue length. For the moment, only DropTail queues are used. However, if new Active Queue Management (AQM) schemes become available, the performance of the congestion control algorithm should be again evaluated.

In an experiment, if the media only flows in a single direction, the feedback path should also be tested with varying amounts of impairments.

The main motivation for the previous and current criteria is to identify situations in which the proposed congestion control is less performant.

4.7. Reacting to Transient Events or Interruptions

The congestion control algorithm should be able to handle changes in end-to-end capacity and latency. Latency may change due to route updates, link failures, handovers etc. In mobile environment the

end-to-end capacity may vary due to the interference, fading, handovers, etc. In wired networks the end-to-end capacity may vary due to changes in resource reservation.

4.8. Fairness With Similar Cross-Traffic

The congestion control algorithm should be evaluated when competing with other RTP flows using the same or another candidate congestion control algorithm. The proposal should highlight the bottleneck capacity share of each RTP flow.

4.9. Impact on Cross-Traffic

The congestion control algorithm should be evaluated when competing with standard TCP. Short TCP flows may be considered as transient events and the RTP flow may give way to the short TCP flow to complete quickly. However, long-lived TCP flows may starve out the RTP flow depending on router queue length.

The proposal should also measure the impact on varied number of cross-traffic sources, i.e., few and many competing flows, or mixing various amounts of TCP and similar cross-traffic.

4.10. Extensions to RTP/RTCP

The congestion control algorithm should indicate if any protocol extensions are required to implement it and should carefully describe the impact of the extension.

5. Minimum Requirements for Evaluation

The minimal requirements for RMCAT proposals is to produce or present results for the test scenarios described in Section 5 of [I-D.sarker-rmcat-eval-test] (Basic Test Cases).

6. Status of Proposals

Congestion control algorithms are expected to be published as "Experimental" documents until they are shown to be safe to deploy. An algorithm published as a draft should be experimented in simulation, or a controlled environment (testbed) to show its applicability. Every congestion control algorithm should include a note describing the environments in which the algorithm is tested and safe to deploy. It is possible that an algorithm is not recommended for certain environments or perform sub-optimally for the user.

[Editor's Note: Should there be a distinction between "Informational" and "Experimental" drafts for congestion control algorithms in RMCAT.

[RFC5033] describes Informational proposals as algorithms that are not safe for deployment but are proposals to experiment with in simulation/testbeds. While Experimental algorithms are ones that are deemed safe in some environments but require a more thorough evaluation (from the community).]

7. Security Considerations

Security issues have not been discussed in this memo.

8. IANA Considerations

There are no IANA impacts in this memo.

9. Contributors

The content and concepts within this document are a product of the discussion carried out in the Design Team.

Michael Ramalho provided the text for a specific scenario, which is now covered in [I-D.sarker-rmcat-eval-test].

10. Acknowledgements

Much of this document is derived from previous work on congestion control at the IETF.

The authors would like to thank Harald Alvestrand, Anna Brunstrom, Luca De Cicco, Wesley Eddy, Lars Eggert, Kevin Gross, Vinayak Hegde, Stefan Holmer, Randell Jesup, Karen Nielsen, Piers O'Hanlon, Colin Perkins, Michael Ramalho, Zaheduzzaman Sarker, Timothy B. Terriberry, Michael Welzl, and Mo Zanaty for providing valuable feedback on earlier versions of this draft. Additionally, also thank the participants of the design team for their comments and discussion related to the evaluation criteria.

11. References

11.1. Normative References

- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC3551] Schulzrinne, H. and S. Casner, "RTP Profile for Audio and Video Conferences with Minimal Control", STD 65, RFC 3551, July 2003.

- [RFC3611] Friedman, T., Caceres, R., and A. Clark, "RTP Control Protocol Extended Reports (RTCP XR)", RFC 3611, November 2003.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.
- [RFC5506] Johansson, I. and M. Westerlund, "Support for Reduced-Size Real-Time Transport Control Protocol (RTCP): Opportunities and Consequences", RFC 5506, April 2009.
- [I-D.ietf-rmcat-cc-requirements]
Jesup, R., "Congestion Control Requirements For RMCAT", draft-ietf-rmcat-cc-requirements-02 (work in progress), February 2014.
- [I-D.ietf-avtcore-rtp-circuit-breakers]
Perkins, C. and V. Singh, "Multimedia Congestion Control: Circuit Breakers for Unicast RTP Sessions", draft-ietf-avtcore-rtp-circuit-breakers-05 (work in progress), February 2014.

11.2. Informative References

- [RFC5033] Floyd, S. and M. Allman, "Specifying New Congestion Control Algorithms", BCP 133, RFC 5033, August 2007.
- [RFC5166] Floyd, S., "Metrics for the Evaluation of Congestion Control Mechanisms", RFC 5166, March 2008.
- [RFC5681] Allman, M., Paxson, V., and E. Blanton, "TCP Congestion Control", RFC 5681, September 2009.
- [I-D.sarker-rmcat-eval-test]
Sarker, Z., Singh, V., Zhu, X., and M. Ramalho, "Test Cases for Evaluating RMCAT Proposals", draft-sarker-rmcat-eval-test-00 (work in progress), February 2014.
- [SA4-EVAL]
R1-081955, 3GPP., "LTE Link Level Throughput Data for SA4 Evaluation Framework", 3GPP R1-081955, 5 2008.
- [SA4-LR]
S4-050560, 3GPP., "Error Patterns for MBMS Streaming over UTRAN and GERAN", 3GPP S4-050560, 5 2008.

[TCP-eval-suite]

Lachlan, A., Marcondes, C., Floyd, S., Dunn, L., Guillier, R., Gang, W., Eggert, L., Ha, S., and I. Rhee, "Towards a Common TCP Evaluation Suite", Proc. PFLDnet. 2008, August 2008.

Appendix A. Application Trade-off

Application trade-off is yet to be defined. see RMCAT requirements [I-D.ietf-rmcat-cc-requirements] document. Perhaps each experiment should define the application's expectation or trade-off.

A.1. Measuring Quality

No quality metric is defined for performance evaluation, it is currently an open issue. However, there is consensus that congestion control algorithm should be able to show that it is useful for interactive video by performing analysis using a real codec and video sequences.

Appendix B. Change Log

Note to the RFC-Editor: please remove this section prior to publication as an RFC.

B.1. Changes in draft-ietf-rmcat-eval-criteria-02

- o Incorporated fairness test as a working test.
- o Updated text on minimum evaluation requirements.

B.2. Changes in draft-ietf-rmcat-eval-criteria-01

- o Removed Appendix B.
- o Removed Section on Evaluation Parameters.

B.3. Changes in draft-ietf-rmcat-eval-criteria-00

- o Updated references.
- o Resubmitted as WG draft.

B.4. Changes in draft-singh-rmcat-cc-eval-04

- o Incorporate feedback from IETF 87, Berlin.
- o Clarified metrics: convergence time, bandwidth utilization.

- o Changed fairness criteria to fairness test.
- o Added measuring pre- and post-repair loss.
- o Added open issue of measuring video quality to appendix.
- o clarified use of DropTail and AQM.
- o Updated text in "Minimum Requirements for Evaluation"

B.5. Changes in draft-singh-rmcat-cc-eval-03

- o Incorporate the discussion within the design team.
- o Added a section on evaluation parameters, it describes the flow and network characteristics.
- o Added Appendix with self-fairness experiment.
- o Changed bottleneck parameters from a proposal to an example set.
- o

B.6. Changes in draft-singh-rmcat-cc-eval-02

- o Added scenario descriptions.

B.7. Changes in draft-singh-rmcat-cc-eval-01

- o Removed QoE metrics.
- o Changed stability to steady-state.
- o Added measuring impact against few and many flows.
- o Added guideline for idle and data-limited periods.
- o Added reference to TCP evaluation suite in example evaluation scenarios.

Authors' Addresses

Varun Singh
Aalto University
School of Electrical Engineering
Otakaari 5 A
Espoo, FIN 02150
Finland

Email: varun@comnet.tkk.fi
URI: <http://www.netlab.tkk.fi/~varun/>

Joerg Ott
Aalto University
School of Electrical Engineering
Otakaari 5 A
Espoo, FIN 02150
Finland

Email: jo@comnet.tkk.fi

RMCAT WG
Internet-Draft
Intended status: Informational
Expires: April 30, 2015

I. Johansson
Z. Sarker
Ericsson AB
October 27, 2014

Self-Clocked Rate Adaptation for Multimedia
draft-johansson-rmcat-scream-cc-03

Abstract

This memo describes a rate adaptation framework for conversational video services. The solution conforms to the packet conservation principle and uses a hybrid loss and delay based congestion control algorithm. The framework is evaluated over both simulated bottleneck scenarios as well as in a LTE (Long Term Evolution) system simulator and is shown to achieve both low latency and high video throughput in these scenarios.

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 30, 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
1.1. Wireless (LTE) access properties	3
2. Terminology	3
3. The adaptation framework	4
3.1. Congestion control	7
3.2. Transmission scheduling	8
3.3. Media rate control	8
4. Detailed description	8
4.1. Network congestion control	8
4.1.1. Congestion window update	9
4.1.1.1. Initial steps	9
4.1.1.2. Loss event is detected	11
4.1.1.3. If <code>in_exponential_start</code> = true and no loss event detected	11
4.1.1.4. If <code>in_exponential_start</code> = false and no loss event detected	11
4.1.1.5. Fairness enforcement	11
4.1.1.6. Final CWND adjustment step	12
4.1.1.7. Competing flows compensation, adjustment of <code>owd_target</code>	12
4.1.2. Transmission scheduling	13
4.1.2.1. Transmission decision	13
4.1.2.2. Next transmission attempt	14
4.2. Video rate control	14
4.2.1. Frame skipping	15
4.2.2. Rate change	16
4.2.2.1. Reduce rate	17
4.2.2.2. Increase rate	18
5. Conclusion	19
6. Open issues	20
7. Acknowledgements	20
8. IANA Considerations	20
9. Security Considerations	20
10. Change history	20
11. References	20
11.1. Normative References	21
11.2. Informative References	21
Authors' Addresses	22

1. Introduction

Rate adaptation is considered as an important part of a interactive realtime communication as the transmission channel bandwidth may vary over period of time. Wireless access such as LTE (Long Term Evolution), which is an integral part of the current Internet, increases the importance of rate adaptation as the channel bandwidth of a default LTE bearer [QoS-3GPP] can change considerably in a very short time frame. Thus a rate adaptation solution for interactive realtime media, such as WebRTC, in LTE system must be both quick and be able to operate over a large span in available channel bandwidth. This memo describes a solution that borrows the self-clocking principle of TCP and combines it with a new delay based rate adaptation algorithm, LEDBAT [RFC6817]. Because neither TCP nor LEDBAT was designed for interactive realtime media, a few extra features are needed to make the concept work well with in this context. This memo describes these extra features.

1.1. Wireless (LTE) access properties

[I-D.draft-sarker-rmcat-cellular-eval-test-cases] introduces the complications that can be observed in wireless environments. Wireless access such as LTE can typically not guarantee a given bandwidth, this is true especially for default bearers. The network throughput may vary considerably for instance in cases where the wireless terminal is moving around.

Unlike wireline bottlenecks with large statistical multiplexing it is not possible to try to maintain a given bitrate when congestion is detected with the hope that other flows will yield, this because there are generally few other flows competing for the same bottleneck. Each user gets its own variable throughput bottleneck, where the throughput depends on factors like channel quality, load and historical throughput. The bottom line is, if the throughput drops, the sender has no other option than to reduce the bitrate. In addition, the grace time, i.e. allowed reaction time from the time that the congestion is detected until a reaction in terms of a rate reduction is effected, is generally very short, in the order of one RTT (Round Trip Time).

2. Terminology

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

3. The adaptation framework

The adaptation framework has similarities to concepts like TFWC [TFWC]. One important property is self-clocking and compliance to the packet conservation principle. The packet conservation principle is described as an important key-factor behind the protection of networks from congestion [FACK].

The packet conservation principle is realized by including a vector of the sequence numbers of received packets in the feedback from the receiver back to the sender, the sender keeps a list of transmitted packets and their respective sizes. This information is then used to determine how many bytes can be transmitted. A congestion window puts an upper limit on how many bytes can be in flight, i.e. transmitted but not yet acknowledged. The congestion window is determined in a way similar to LEDBAT [RFC6817]. This ensures that the e2e latency is kept low. The basic functionality is quite simple, there are however a few steps to take to make the concept work with conversational media. These will be briefly described in sections Section 3.1 to Section 3.3.

The feedback is over RTCP [RFC3550] and is based on [RFC4585]. It is implemented as a transport layer feedback message, see proposed example in Figure 1. The feedback control information part (FCI) consists of the following elements.

- o **Timestamp:** A timestamp value indicating when the last packet was received which makes it possible to compute the one way (extra) delay (OWD).
- o **The ACK list (Highest received sequence number + ACK vector):** Makes it possible to detect lost packets and determine the number of bytes in flight.
- o **Source quench bit (Q):** Makes it possible to request the sender to reduce its congestion window. This is useful if WebRTC media is received from many hosts and it becomes necessary to balance the bitrates between the streams. The exact behavior and use for the source quench bit is T.B.D.
- o **ECE (Explicit Congestion Notification) echo:** Makes it possible to indicate if packets are ECN-CE (ECN Congestion Experienced) marked. The use for the ECN echo bits is T.B.D.

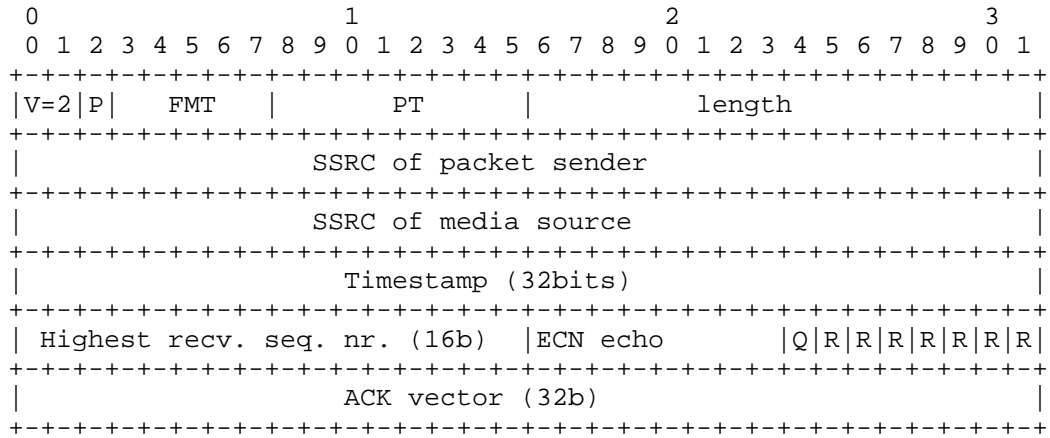


Figure 1: Transport layer feedback message

To make the feedback as frequent as possible, the feedback packets are transmitted as reduced size RTCP according to [RFC5506].

The timestamp clock time base is typically set to the same time base as the media source in question but as the protocol described here is not dependent on the media it can be set to a fixed value defined in this specification. The ACK vector is here a bit vector that indicates the reception of the last 1+32 = 33 RTP packets.

Section 4 describes the main algorithm details and how the feedback is used.

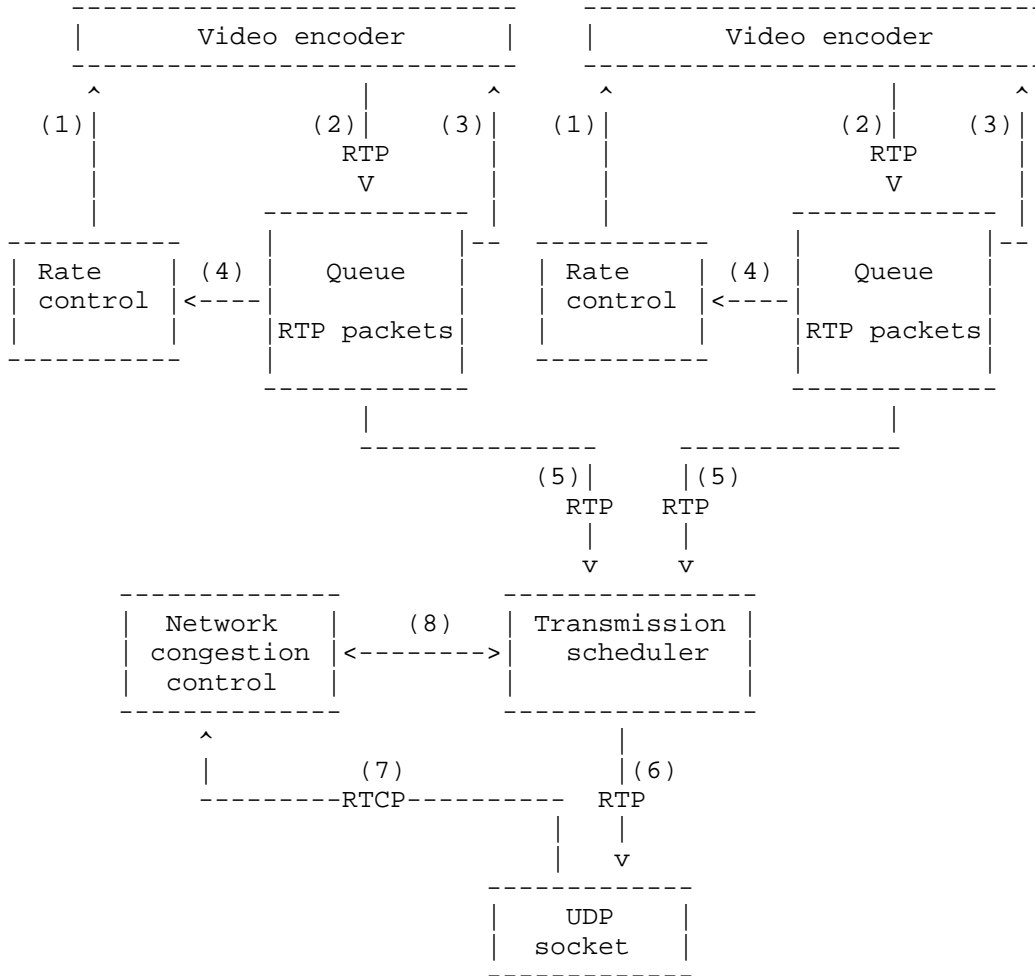


Figure 2: Rate adaptation framework

Figure 2 shows the functional overview of the adaptation framework. Each media type or source implements rate control and a queue, where RTP packets containing encoded media frames are temporarily stored for transmission, the figure shows the details for when two video sources are used. Video frames are encoded and forwarded to the queue (2). The media rate adaptation adapts to the age of the oldest RTP frame in the queue and controls the video bitrate (1). It is also possible to make the video encoder skip frames and thus temporarily reduce the frame rate if the queue age exceeds a given threshold (3). The RTP packets are picked from each queue based on some defined priority order or simply in a round robin fashion (5). A transmission scheduler takes care of the transmission of RTP packets, to be written to the UDP socket (6). In the general case all media must go through the packet scheduler and is allowed to be transmitted if the number of bytes in flight is less than the congestion window. However audio frames can be allowed to be transmitted as audio is typically low bitrate and thus contributes little to congestion, this is however something that is left as an implementation choice. RTCP packets are received (7) and the information about bytes in flight and congestion window is exchanged between the network congestion control and the transmission scheduler (8).

The rate adaptation solution constitutes three parts; congestion control, transmission scheduling and media rate adaptation.

3.1. Congestion control

The congestion control sets an upper limit on how much data can be in the network (bytes in flight); this limit is called CWND (congestion window) and is used in the transmission scheduling.

A congestion control method, similar to LEDBAT [RFC6817], measures the OWD (one way delay). The congestion window is allowed to increase if the OWD is below a predefined target, otherwise the congestion window decreases. The delay target is typically set to 50-100ms. This ensures that the OWD is kept low on the average. The reaction to loss events is similar to that of loss based TCP, i.e. an instant reduction of CWND.

LEDBAT is designed with file transfers as main use case which means that the algorithm must be modified somewhat to work with rate-limited sources such as video. The modifications are

- o Congestion window validation techniques. These are similar in action as the method described in [I-D.ietf-tcpm-newcwnd].

- o Fast start for bitrate increase. It makes the video bitrate ramp-up within 5 to 10 seconds. The behavior is similar to TCP slowstart. The fast start is exited when congestion is detected.
- o Adaptive delay target. This helps the congestion control to compete with FTP traffic to some degree.

3.2. Transmission scheduling

Transmission scheduling limits the output of data, given by the relation between the number of bytes in flight and the congestion window similar to TCP. Packet pacing is used to mitigate issues with coalescing that may cause increased jitter in the media traffic.

3.3. Media rate control

The media rate control serves to adjust the media bitrate to ramp up quickly enough to get a fair share of the system resources when link throughput increases.

The reaction to reduced throughput must be prompt in order to avoid getting too much data queued up in the sender frame queues. The queuing delay is determined and the media bitrate is decreased if it exceeds a threshold.

In cases where the sender frame queues increase rapidly such as the case of a RAT (Radio Access Type) handover it may be necessary to implement additional actions, such as discarding of encoded video frames or frame skipping in order to ensure that the sender frame queues are drained quickly. Frame skipping means that the frame rate is temporarily reduced. Discarding of old video frames is a more efficient way to reduce media latency than frame skipping but it comes with a requirement to repair codec state, frame skipping is thus to prefer as a first remedy.

4. Detailed description

This section describes the algorithm in more detail. It is split between the network congestion control and the video rate adaptation.

4.1. Network congestion control

This section explains the network congestion control, it contains two main functions

- o Computation of congestion window: Gives an upper limit to the number of bytes in flight i.e. how many bytes that have been transmitted but not yet acknowledged.

- o Transmission scheduling: RTP packets are transmitted if allowed by the relation between the number of bytes in flight and the congestion window

Unlike TCP, SCReAM is not a byte oriented protocol, rather it is an RTP packet oriented protocol. Thus it keeps a list of transmitted RTP packets and their respective sending times (wall-clock time). The feedback indicates the highest received RTP sequence number and a timestamp (wall-clock time) when it was received. In addition, an ACK list is included to make it possible to determine lost packets

4.1.1. Congestion window update

Below is described the actions when an acknowledgement is received.

4.1.1.1. Initial steps

Bytes in flight (`bytes_in_flight`) is computed as the sum of the sizes of the RTP packets ranging from the RTP packet most recently transmitted up to but not including the acknowledged packet with the highest sequence number. As an example: If RTP packet was sequence number SN with transmitted and the last ACK indicated SN-5 as the highest received sequence number then bytes in flight is computed as the sum of the size of RTP packets with sequence number SN-4, SN-3, SN-2, SN-1 and SN.

The congestion window is computed from the one way (extra) delay estimates (OWD) that are obtained from the send and received timestamp of the RTP packets. LEDBAT [RFC6817] explains the details of the computation of the OWD. An OWD sample is obtained for each received acknowledgement. No smoothing of the OWD samples occur, however some smoothing occurs naturally as the computation of the CWND is in itself a low pass filter function.

A variable `bytes_newly_acked` depicts the number bytes that was acknowledged with the last received acknowledgement.

`owd_mem` is an EWMA (Exponential Weighted Moving Average) filtered OWD

$$\text{owd_mem} = \max(\text{owd_mem} * 0.5 + \text{owd} * 0.5, \text{owd_mem} * 0.9 + \text{owd} * 0.1)$$

The OWD fraction is computed as

$$\text{owd_fraction} = \text{owd} / \text{owd_target}$$

where `owd_target` is the target (extra) delay, `owd_target` is typically set to `owd_target_lo=0.1s` but can in certain cases increase to `owd_target_hi=0.4s`. The OWD fraction is sampled every 50ms and the

last 20 samples are stored in a vector (`owd_fraction_hist`). This vector is used in the computation of an OWD trend that gives a value between 0.0 and 1.0 depending on how close to congestion it gets. The OWD trend is calculated as follows

Let $R(\text{owd_fraction_hist}, K)$ be the autocorrelation function of `owd_fraction_hist` at lag K . The 1st order prediction coefficient is formulated as

$$a = R(\text{owd_fraction_hist}, 1) / R(\text{owd_fraction_hist}, 0)$$

The prediction coefficient a has positive values if OWD shows an increasing trend, thus one get an indication of congestion before the OWD target is reached. The prediction coefficient is further multiplied with `owd_fraction` to reduce sensitivity to increasing OWD when OWD is small. The OWD trend is thus computed as

$$\text{owd_trend} = \max(0.0, \min(1.0, a * \text{owd_fraction}))$$

The `owd_trend` is utilized in the media rate control and to determine when to exit slow start.

An EWMA filtered version of `owd_trend` is computed

$$\text{owd_trend_ewma} = \max(\text{owd_trend}, \text{owd_trend_ewma} * (1.0 - \alpha) + \alpha * \text{owd_trend})$$

$$\alpha = (t_now - t_cwnd_update_prev) / 5000.0$$

`t_now` is the current wall clock time.

`owd_fraction_avg` is a lowpass filtered version of `owd_fraction`

$$\text{owd_fraction_avg} = 0.9 * \text{owd_fraction_avg} + 0.1 * \text{owd_fraction}$$

An off target value is computed as

$$\text{off_target} = (\text{owd_target} - \text{owd}) / \text{owd_target}$$

CWND is updated differently depending on whether the congestion control is in fast start or not and if a loss event is detected. A Boolean variable `in_exponential_start` (initialized to true) indicates if the congestion is in fast start.

A loss event indicates one or more lost RTP packets within an RTT. This is detected by means of inspection for holes in the sequence number space in the acknowledgements with some margin for possible packet reordering in the network.

4.1.1.2. Loss event is detected

If a loss event is detected then `in_exponential_start` is set to false and CWND is updated according to

$$\text{cwnd} = \max(\text{min_cwnd}, \text{cwnd} * 0.8) \text{ where } \text{min_cwnd} = 2 * \text{mss}$$

otherwise the CWND update continues

4.1.1.3. If `in_exponential_start` = true and no loss event detected

`in_exponential_start` is set to false if `owd_trend` \geq 0.2 and otherwise CWND is updated according to

$$\text{cwnd} = \text{cwnd} + \text{bytes_newly_acked}$$
4.1.1.4. If `in_exponential_start` = false and no loss event detected

Values of `off_target` $>$ 0.0 indicates that the congestion window can be increased. This is done according to the equations below (mss is the maximum RTP packet size).

$$\text{gain} = \text{gain_up} * (1.0 + \max(0.0, 1.0 - \text{owd_trend} / 0.2))$$
$$\text{cwnd} += \text{gain} * \text{off_target} * \text{bytes_newly_acked} * \text{mss} / \text{cwnd}$$

Values of `off_target` \leq 0.0 indicates congestion, CWND is then updated according to the equation

$$\text{cwnd} += \text{gain_down} * \text{off_target} * \text{bytes_newly_acked} * \text{mss} / \text{cwnd}$$

4.1.1.5. Fairness enforcement

Fairness enforcement is realized by reducing the congestion window by a fraction when a number of conditions are met. They are

- o `owd_target` $<$ `owd_target_lo` * 1.2 i.e no competing flows are compensated for
- o `owd_trend` $>$ 0.1 i.e. congestion is detected
- o more than `t_delta` since the congestion window was reduced the last time

`t_delta` is computed as

$$\text{t_delta} = 0.1 * \min(200.0, \max(20.0, 50.0e6 / \text{max_paced_bitrate}))$$

The bitrate is taken into account in the sense that the lower the bitrate, the more sparse the reductions in congestion window get.

If the above conditions are met then cwnd is adjusted according to

$$\text{cwnd} *= 0.8$$

4.1.1.6. Final CWND adjustment step

The congestion window is limited by the maximum number of bytes in flight over the last 1.0 seconds according to

$$\text{cwnd} = \min(\text{cwnd}, \text{max_bytes_in_flight} * \text{max_bytes_in_flight_head_room})$$

where `max_bytes_in_flight_head_room` = 1.1. This avoids possible over-estimation of the throughput after for example, idle periods-

Finally cwnd is set to ensure that it is at least `min_cwnd`

$$\text{cwnd} = \max(\text{cwnd}, \text{min_cwnd})$$

4.1.1.7. Competing flows compensation, adjustment of `owd_target`

In certain cases it becomes necessary to increase `owd_target`, one such case is where SCReAM competes with TCP based file transfer over a tail drop bottleneck link and the TCP congestion avoidance is loss based (for example Cubic or NewReno). The technique is to inhibit video long enough to make bytes in flight reach zero (no remaining RTP packets in flight) and then resume video. For the unfortunate case that the last RTP packet was lost, it is necessary to force video to resume after 1.0s as bytes in flight will never reach zero in this case.

This interruption is typically in the order of one RTT. Once video is resumed the average OWD (`owd_avg_c_flow`) is computed over the first 5 acknowledgements after video is resumed. If no competing flows exist then this average should be close to zero, otherwise `owd_avg_c_flow` has a value that corresponds roughly to the queuing delay caused by the competing flow. The `owd_target` is updated according to the value of `owd_avg_c_flow`.

The method above is executed if more than a given time since the last time video was inhibited (e.g. 20 seconds) and any of the two conditions below are fulfilled

- o `owd_mem > owd_target`
- o `owd_target > owd_target_lo`

The first condition indicates that another competing flows is possibly driving higher queuing delays in the network. The second condition indicates that the OWD target is increased and it should be determined if this can be lowered. Once `owd_avg_c_flow` is computed the `owd_target` is adjusted. The adjustment action depends on the value of `owd_avg_c_flow`

- o If `owd_avg_c_flow > owd_target_lo/2`:
Adjust the `owd_target` upwards according to
`owd_target = min(owd_target_hi, max(owd_target, owd_avg_c_flow *3.0))`
- o If `owd_avg_c_flow <= owd_target_lo/2`:
Adjust the `owd_target` downwards according to
`owd_target = 0.5*owd_target+ 0.5*Math.max(owd_target_lo, owd_avg_c_flow)`.
Furthermore `owd_target` is set to `owd_target_lo` if it is less than `owd_target_lo*1.2`.

4.1.2. Transmission scheduling

An RTP packet transmission attempt is scheduled at intervals given by `t_pace` that depends on the estimated throughput, the RTT and the size of the last transmitted RTP packet. This provides with packet pacing which is in some cases necessary in order to break up coalescing tendencies which can otherwise cause unwanted extra jitter or packet loss.

4.1.2.1. Transmission decision

The principle is to allow packet transmission of an RTP packet only if the number of bytes in flight is less than the congestion window. There are however two reasons why this strict rule will not work optimally

- o Bitrate variations. The video frame size is always varying to a larger or smaller extent, a strict rule as the one given above will have the effect that the video bitrate have difficulties to increase as the congestion window puts a too hard restriction on the video frame size variation, this further can lead to occasional queuing of RTP packets in the RTP packet queue that will prevent bitrate increase because of the increased queuing delay.
- o Reverse (feedback) path congestion. Especially in transport over buffer-bloated networks, the one way delay in the reverse direction may jump due to congestion. The effect of this is that the acknowledgements are delayed with the result that the self-

clocking is temporarily halted, even though the forward path is not congested.

Transmission of an RTP packet of size `rtp_size` is thus allowed when any of the following conditions is met.

- o If `owd > owd_target`:
Transmission is allowed if
`bytes_in_flight + rtp_size <= cwnd`.
This enforces a strict rule that helps to prevent further queue buildup.
- o If `owd <= owd_target`:
A helper variable
`x_cwnd=1.0+bytes_in_flight_slack*max(0.0, min(1.0,1.0-owd_trend/0.5))/100.0`
is computed. Transmission is allowed if
`bytes_in_flight+rtp_size <= max(cwnd*x_cwnd, cwnd+mss)` .
This gives a slack that reduces as congestion increases,
`bytes_in_flight_slack` is a maximum allowed slack in percent.
A large value such as 100% increases the robustness to bitrate variations in the source and congested feedback channel issues.
The possible drawback is increased delay or packet loss when forward path congestion occur. Recommended values are 20 to 50%.

4.1.2.2. Next transmission attempt

The interval until the next transmission attempt (`t_pace`) is set to 0.001s if no RTP packet was transmitted according to the decision in previous section. Otherwise it is calculated as

```
max_paced_bitrate = max (50000, cwnd* 8 / s_rtt)
```

```
t_pace = rtp_size * 8 / max_paced_bitrate
```

4.2. Video rate control

The video rate control is based on the queuing delay in the RTP packet queue and loss events. The video rate control function is executed for each video frame. The actual video rate adjustment may however be less frequent. The main reason is that there is typically a lag between the bitrate request and the actual bitrate from the video coder and this lag can be as much as 1 second. This makes it less efficient to try to react to congestion with prompt rate adjustments. The solution is to complement the rate reduction with frame skipping in order to keep the RTP queuing delay limited.

The queuing delay is sampled every frame period and the last N_a samples are stored in a vector `age_vec`.

An average queuing delay is computed as a weighted sum over the samples in `age_vec`. `age_avg` at the current time instant n is computed as

$$\text{age_avg}(n) = \text{SUM } \text{age_vec}(n-k) * w(k)$$

The sum is computed over $k=[0..N_a-1]$

$w(n)$ are weight factors arranged to give the most recent samples a higher weight.

N_a i.e. the number of samples that `avg_age` is computed over, depends on how slow the video encoder is to respond to video rate change requests. With a slow video encoder N_a is suggested to be set to

$$N_a = 1.0/\text{frame_period}$$

where `frame_peridod` is the video frame interval, 1.0 corresponds roughly to the time constant in the video coder rate control loop (1.0s).

If the video encoder is quicker to react to bitrate changes, N_a can be set to a lower value such as $N_a = 5$.

`avg_age` is used for rate adjustment instead of the current value, the reason is to avoid bitrate reduction because of temporal delay spikes. Instead the video rate control is a combination of slower rate adjustments and adjustments of the temporal frame rate by means of raw frame skipping on a shorter time scale. This is an adaptation to SCReAM as it works best when it has data to send because of its self-clocking properties. The concept also avoids very large rate reduction due to isolated delay spikes.

The change in `age_avg` is computed as

$$\text{age_d} = (\text{age_avg}(n) - \text{age_avg}(n-1))/\text{frame_period}$$

4.2.1. Frame skipping

Frame skipping is controlled by a flag `frame_skip` which, if set to 1, dictates that the video coder should skip the next video frame. The frame skipping intensity at the current time instant n is computed as

- o If `age_d > 0` and `age_avg > frame_period`:
The frame skip intensity is computed as

```
frame_skip_intensity = min(1.0, (age_vec(n)-frame_period)/(2*
frame_period))
```

- o Otherwise frame skip intensity is set to zero

Note that the frame skipping intensity is computed based on the current value of the queuing delay. Furthermore, frame skipping is enabled only if the average queue delay increases and is large enough.

The skip_frame flag is set depending on three variables

- o frame_skip_intensity
- o since_last_frame_skip, i.e the number of consecutive frames without frame skipping
- o consecutive_frame_skips, i.e the number of consecutive frame skips

The flag skip_frame is set to 1 if any of the conditions below is met.

- o $\text{age_vec}(n) > 0.2 \ \&\& \ \text{consecutive_frame_skips} < 5$
- o $\text{frame_skip_intensity} < 0.5 \ \&\& \ \text{since_last_frame_skip} \geq 1.0 / \text{frame_skip_intensity}$
- o $\text{frame_skip_intensity} \geq 0.5 \ \&\& \ \text{consecutive_frame_skips} < (\text{frame_skip_intensity} - 0.5) * 10$

The arrangement makes sure that no more than 4 frames are skipped in sequence, the rationale is to ensure that the input to the video encoder does not change too much, something that may give poor prediction gain.

4.2.2. Rate change

A variable target_bitrate is adjusted depending on the congestion state. The target bitrate can vary between a minimum value (target_bitrate_min) and a maximum value (target_bitrate_max).

First of all the target_bitrate is updated if a new loss event was indicated and the rate change procedure is exited.

```
target_bitrate = max(0.9* target_bitrate, target_bitrate_min)
```

If no loss event was indicated then the rate change procedure continues. Based on age_avg(n) and the time span since the last rate

reduction. A rate reduction condition is determined. This is evaluated differently depending on whether an ideal video coder is simulated for algorithm evaluation purposes or if the algorithm is executed in a real implement with a video coder that lags behind in the rate adjustment.

- o Ideal mode: `reduce_rate = age_avg(n) > frame_period/2 && t_now-t_last_rate_change >= rate_change_interval && t_now-t_last_rate_reduction > 0.5`
- o Non-ideal mode: `reduce_rate = age_avg(n) > frame_period*2 && t_now-t_last_rate_change >= rate_change_interval && t_now-t_last_rate_reduction > video_coder_time_constant`

`rate_change_interval` is set to 0.1s, `video_coder_time_constant` is set to a value that approximates the lag in the video coder rate change.

4.2.2.1. Reduce rate

If `reduce_rate` evaluates to true then the bitrate is reduced. First an inflection point is determined for later rate increase

```
target_bitrate_i = target_bitrate * 0.95
```

In addition, a restore point is determined for the case that false congestion was detected, for instance as an effect of congestion in the feedback path.

```
target_bitrate_restore_point = target_bitrate
```

A few variables are updated for future use

```
t_last_rate_change = t_now
```

```
max_owd_fraction = max(max_owd_fraction, owd_fraction_avg)
```

A rate reduction factor is determined

```
alpha = min(0.5, max(0.0, 0.9*age_d))
```

The target bitrate and `t_last_rate_reduction` are updated if `alpha > 0.0` according to

```
target_bitrate = max(target_bitrate_min, target_bitrate*(1.0-alpha))
```

4.2.2.2. Increase rate

A rate increase is allowed if two conditions are met

- o $t_{\text{now}} - t_{\text{last_rate_change}} \geq \text{rate_change_interval}$
- o $\text{age_avg}(n) \leq \text{frame_period}/2$

First the target bitrate is restored if false congestion was detected. This restoration is allowed if it is more than 2.0s since the last loss event and $\text{target_bitrate_restore_point} > 0.0$. Further, if an additional condition

$\text{do_restore} = \text{max_owd_fraction} < 0.4 \ \&\& \ \text{owd_trend_ewma} < 0.2$

evaluates to true then the target bitrate is restored as

$\text{target_bitrate} = \text{max}(\text{target_bitrate}, \text{target_bitrate_restore_point})$

Regardless of whether `do_restore` evaluates to true or false `target_bitrate_restore_point` is set to -1.0 and `max_owd_fraction` = 0.0. The target bitrate is increased, the increase rate depends on if the algorithm is in slow start or not, indicated by the variable `in_exponential_start`.

4.2.2.2.1. If `in_exponential_start` = true

The bitrate incremented is computed as

$\text{increment} =$
 $\text{target_bitrate_max} * \text{rate_change_interval} * \text{ramp_up_time_fast} *$
 $(1.0 - \text{min}(1.0, \text{owd_trend}/0.1))$

$\text{target_bitrate} = \text{min}(\text{target_bitrate_max}, \text{target_bitrate} + \text{increment})$

The target bitrate is allowed to reach the the highest bitrate within `ramp_up_time_fast` seconds if no congestion is detected. A recommended value for `ramp_up_time_fast` is 10.0s.

4.2.2.2.2. If `in_exponential_start` = false

The maximum allowed increment of the target bitrate is computed

$\text{increment_max} = \text{target_bitrate} * 0.2$

A variable gain factor is computed in a number of steps, first the gain factor is reduced if the target bitrate is close to the

inflection point i.e. the target bitrate when congestion was last detected.

```
gain = max(0.2,min(1.0, abs((target_bitrate - target_bitrate_i)/  
target_bitrate_i)*4.0))
```

Furthermore the gain is reduced if near (or past) congestion is detected

```
gain *= min(1.0, max(0.0,(1.0-owd_trend_ewma)))
```

The gain is increased if competing (potentially aggressive) flows are detected, this is indicated by that $owd_target/owd_target_lo > 1.0$

```
gain *= owd_target/owd_target_lo
```

A ramp-up speed is computed that is adjusted depending on the estimated congestion level

```
ramp_up_time = ramp_up_time_fast+(ramp_up_time_slow-  
ramp_up_time_fast)*  
max(0.0,Math.min(1.0, owd_trend_ewma /0.2))
```

A recommended value for `ramp_up_time_slow` is 20.0s. The increment is computed and the `target_bitrate` is updated

```
increment = min(target_bitrate_max*gain*rate_change_interval  
/(ramp_up_t),  
increment_max)
```

```
target_bitrate = min(target_bitrate_max, target_bitrate +increment)
```

5. Conclusion

This memo describes a congestion control framework for RMCAT that it is particularly good at handling the quickly changing condition in wireless network such as LTE. The solution conforms to the packet conservation principle and leverages on novel congestion control algorithms and recent TCP research, together with media bitrate determined by sender queuing delay and given delay thresholds. The solution has shown potential to meet the goals of high link utilization and prompt reaction to congestion. The solution is realized with a new RFC4585 transport layer feedback message.

6. Open issues

A list of open issues.

- o Describe use of Q bit
- o Describe how clock drift compensation is done
- o RTCP AVPF mode. Determine if AVPF early or immediate mode is to prefer
- o Determine format and use of ECN echo field

7. Acknowledgements

We would like to thank the following persons for their comments, questions and support during the work that led to this memo: Markus Andersson, Bo Burman, Tomas Frankkila, Laurits Hamm, Hans Hannu, Nikolas Hermanns, Stefan Haekansson, Erlendur Karlsson, Mats Nordberg, Jonathan Samuelsson, Rickard Sjoeborg, Magnus Westerlund.

8. IANA Considerations

A new RFC4585 transport layer feedback message needs to be standardized.

9. Security Considerations

The feedback can be vulnerable to attacks similar to those that can affect TCP. It is therefore recommended that the RTCP feedback is at least integrity protected.

10. Change history

A list of changes:

- o -02 to -03 : Added algorithm description with equations, removed pseudo code and simulation results
- o -01 to -02 : Updated GCC simulation results
- o -00 to -01 : Fixed a few bugs in example code

11. References

11.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.
- [RFC5506] Johansson, I. and M. Westerlund, "Support for Reduced-Size Real-Time Transport Control Protocol (RTCP): Opportunities and Consequences", RFC 5506, April 2009.
- [RFC6817] Shalunov, S., Hazel, G., Iyengar, J., and M. Kuehlewind, "Low Extra Delay Background Transport (LEDBAT)", RFC 6817, December 2012.

11.2. Informative References

- [FACK] "Forward Acknowledgement: Refining TCP Congestion Control", 2006.
- [I-D.alvestrand-rmcat-congestion]
Holmer, S., Cicco, L., Mascolo, S., and H. Alvestrand, "A Google Congestion Control Algorithm for Real-Time Communication", draft-alvestrand-rmcat-congestion-02 (work in progress), February 2014.
- [I-D.draft-sarker-rmcat-cellular-eval-test-cases]
Sarker, Z., "Evaluation Test Cases for Interactive Real-Time Media over Cellular Networks",
<<http://www.ietf.org/id/draft-sarker-rmcat-cellular-eval-test-cases-00.txt>>.
- [I-D.ietf-tcpm-newcwv]
Fairhurst, G., Sathiaseelan, A., and R. Secchi, "Updating TCP to support Rate-Limited Traffic", draft-ietf-tcpm-newcwv-07 (work in progress), September 2014.
- [QoS-3GPP]
TS 23.203, 3GPP., "Policy and charging control architecture", June 2011, <http://www.3gpp.org/ftp/specs/archive/23_series/23.203/23203-990.zip>.

[TFWC] University College London, "Fairer TCP-Friendly Congestion Control Protocol for Multimedia Streaming", December 2007, <<http://www-dept.cs.ucl.ac.uk/staff/M.Handley/papers/tfwc-conext.pdf>>.

Authors' Addresses

Ingemar Johansson
Ericsson AB
Laboratoriegraend 11
Luleae 977 53
Sweden

Phone: +46 730783289
Email: ingemar.s.johansson@ericsson.com

Zaheduzzaman Sarker
Ericsson AB
Laboratoriegraend 11
Luleae 977 53
Sweden

Phone: +46 761153743
Email: zaheduzzaman.sarker@ericsson.com

RMCAT WG
Internet-Draft
Intended status: Informational
Expires: September 3, 2015

I. Johansson
Z. Sarker
Ericsson AB
March 2, 2015

Self-Clocked Rate Adaptation for Multimedia
draft-johansson-rmcat-scream-cc-05

Abstract

This memo describes a rate adaptation algorithm for conversational video services. The solution conforms to the packet conservation principle and uses a hybrid loss and delay based congestion control algorithm. The algorithm is evaluated over both simulated Internet bottleneck scenarios as well as in a LTE (Long Term Evolution) system simulator and is shown to achieve both low latency and high video throughput in these scenarios.

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 September 3, 2015.

Copyright 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	2
1.1. Wireless (LTE) access properties	3
2. Terminology	3
3. Overview of SCReAM Algorithm	3
3.1. Congestion Control	4
3.2. Transmission Scheduling	5
3.3. Media Rate Control	5
4. Detailed Description of SCReAM	5
4.1. SCReAM Sender	5
4.1.1. Constants and Parameter values	7
4.1.2. Network congestion control	11
4.1.2.1. Congestion window update	12
4.1.2.2. Transmission scheduling	15
4.1.3. Video rate control	16
4.2. SCReAM Receiver	19
5. Feedback Message	20
6. Additional features	21
6.1. Packet pacing	21
6.2. Frame skipping	21
6.3. Q-bit semantics (source quench)	23
7. Discussion	23
8. Conclusion	24
9. Open issues	24
10. Source code	25
11. Acknowledgements	25
12. IANA Considerations	25
13. Security Considerations	25
14. Change history	25
15. References	26
15.1. Normative References	26
15.2. Informative References	26
Authors' Addresses	27

1. Introduction

Congestion in the internet is a reality and applications that are deployed in the internet must have congestion control schemes in place not only for the robustness of the service that it provides but also to ensure the function of the currently deployed internet. As the interactive realtime communication imposes a great deal of requirements on the transport, a robust, efficient rate adaptation for all access types is considered as an important part of interactive realtime communications as the transmission channel

bandwidth may vary over time. Wireless access such as LTE, which is an integral part of the current internet, increases the importance of rate adaptation as the channel bandwidth of a default LTE bearer [QoS-3GPP] can change considerably in a very short time frame. Thus a rate adaptation solution for interactive realtime media, such as WebRTC, must be both quick and be able to operate over a large span in available channel bandwidth. This memo describes a solution, named SCReAM, that is based on the self-clocking principle of TCP and uses techniques similar to what is used in a new delay based rate adaptation algorithm, LEDBAT [RFC6817]. Because neither TCP nor LEDBAT was designed for interactive realtime media, a few extra features are needed to make the concept work well within this context. This memo describes these extra features.

1.1. Wireless (LTE) access properties

[I-D.draft-sarker-rmcat-cellular-eval-test-cases] introduces the complications that can be observed in wireless environments. Wireless access such as LTE can typically not guarantee a given bandwidth, this is true especially for default bearers. The network throughput may vary considerably for instance in cases where the wireless terminal is moving around.

Unlike wireline bottlenecks with large statistical multiplexing it is not possible to try to maintain a given bitrate when congestion is detected with the hope that other flows will yield, this because there are generally few other flows competing for the same bottleneck. Each user gets its own variable throughput bottleneck, where the throughput depends on factors like channel quality, network load and historical throughput. The bottom line is, if the throughput drops, the sender has no other option than to reduce the bitrate. In addition, the grace time, i.e. allowed reaction time from the time that the congestion is detected until a reaction in terms of a rate reduction is effected, is generally very short, in the order of one RTT (Round Trip Time).

2. Terminology

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

3. Overview of SCReAM Algorithm

The core SCReAM algorithm has similarities to concepts like self-clocking used in TFWC [TFWC] and follows packet conservation principles. The packet conservation principle is described as an

important key-factor behind the protection of networks from congestion [FAACK].

The packet conservation principle is realized by including an indication of the highest received sequence number in the feedback, see Section 5, from the receiver back to the sender, the sender keeps a list of transmitted packets and their respective sizes. This information is then used to determine how many bytes can be transmitted. A congestion window puts an upper limit on how many bytes can be in flight, i.e. transmitted but not yet acknowledged. The congestion window is determined in a way similar to LEDBAT [RFC6817]. This ensures that the e2e latency is kept low. The basic functionality is quite simple, there are however a few steps to take to make the concept work with conversational media. These will be briefly described in sections Section 3.1 to Section 3.3.

The rate adaptation solution constitutes three parts- congestion control, transmission scheduling and media rate adaptation. All these three parts reside at the sender side. The receiver side algorithm is very simple in comparison as it only generates acknowledgements to received RTP packets.

3.1. Congestion Control

The congestion control sets an upper limit on how much data can be in the network (bytes in flight); this limit is called CWND (congestion window) and is used in the transmission scheduling.

The SCReAM congestion control method, uses LEDBAT [RFC6817] to measure the OWD (one way delay). The SCReAM sender calculates the congestion window based on the feedback from SCReAM receiver. The congestion window is allowed to increase if the OWD is below a predefined target, otherwise the congestion window decreases. The delay target is typically set to 50-100ms. This ensures that the OWD is kept low on the average. The reaction to loss events is similar to that of loss based TCP, i.e. an instant reduction of CWND.

LEDBAT is designed with file transfers as main use case which means that the algorithm must be modified somewhat to work with rate-limited sources such as video. The modifications are

- o Congestion window validation techniques. These are similar in action as the method described in [I-D.ietf-tcpm-newcwnd].
- o Fast start for bitrate increase. It makes the video bitrate ramp-up within 5 to 10 seconds. The behavior is similar to TCP slowstart. The fast start is exited when congestion is detected. The fast start state can be resumed if the congestion level is

low, this to enable a reasonably quick rate increase in case link throughput increases.

- o Adaptive delay target. This helps the congestion control to compete with FTP traffic to some degree.

3.2. Transmission Scheduling

Transmission scheduling limits the output of data, given by the relation between the number of bytes in flight and the congestion window similar to TCP. Packet pacing is used to mitigate issues with coalescing that may cause increased jitter and/or packet loss in the media traffic.

3.3. Media Rate Control

The media rate control serves to adjust the media bitrate to ramp up quickly enough to get a fair share of the system resources when link throughput increases.

The reaction to reduced throughput must be prompt in order to avoid getting too much data queued up in the RTP packet queues. The media bitrate is decreased if the RTP queue size exceeds a threshold.

In cases where the sender frame queues increase rapidly such as the case of a RAT (Radio Access Type) handover it may be necessary to implement additional actions, such as discarding of encoded video frames or frame skipping in order to ensure that the RTP queues are drained quickly. Frame skipping means that the frame rate is temporarily reduced. Discarding of old video frames is a more efficient way to reduce media latency than frame skipping but it comes with a requirement to repair codec state, frame skipping is thus to prefer as a first remedy. Frame skipping is described as an optional to implement feature in this specification.

4. Detailed Description of SCReAM

4.1. SCReAM Sender

This section describes the sender side algorithm in more detail. It is split between the network congestion control and the video rate adaptation.

Figure 1 shows the functional overview of a SCReAM sender. The RTP application interaction with congestion control is described in [I-D.ietf-rmcat-app-interaction]. Here we use a more decomposed version of the implementation model in the sense that the RTP packets may be queued up in the sender, the transmission of these RTP packets

is controlled by a transmission scheduler. A SCReAM sender implements rate control and a queue for each media type or source, where RTP packets containing encoded media frames are temporarily stored for transmission, the figure shows the details for when two video sources (a.k.a streams) are used.

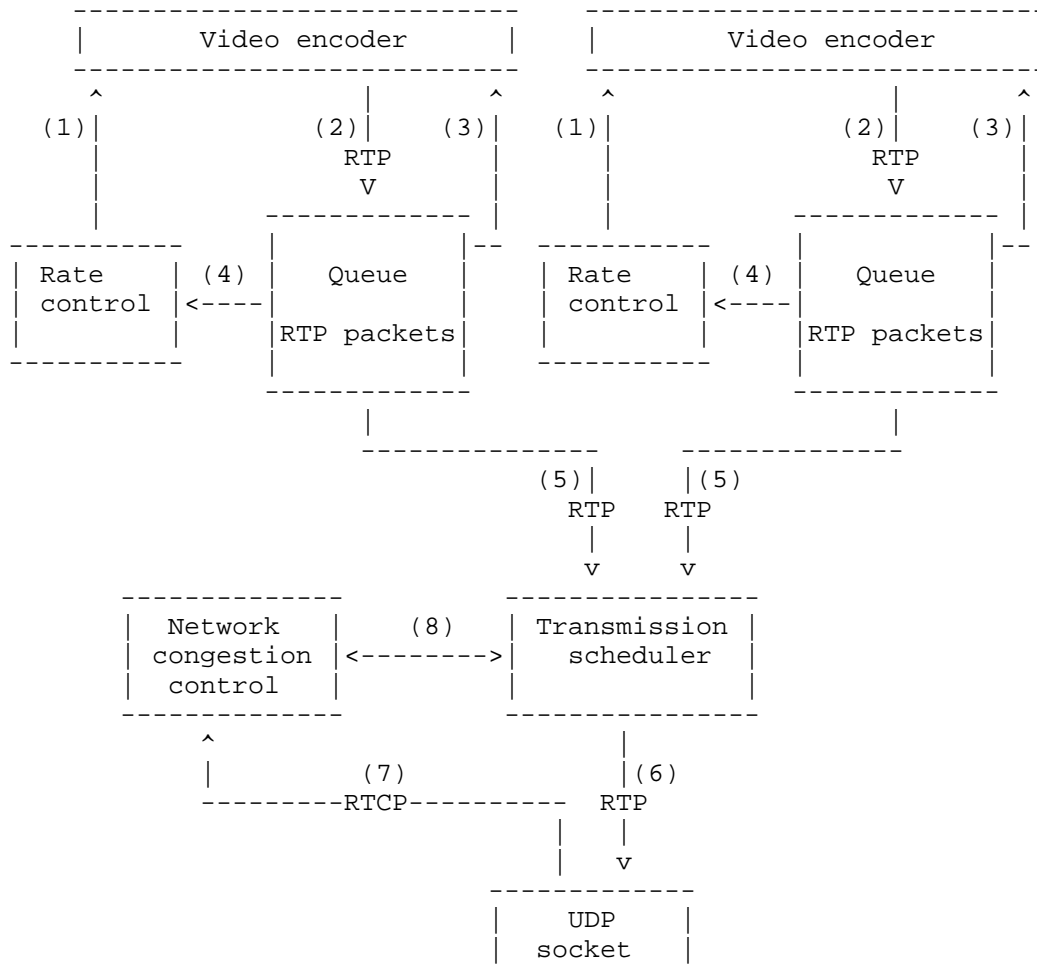


Figure 1: SCReAM sender functional view

Video frames are encoded and forwarded to the queue (2). The media rate adaptation adapts to the size of the RTP queue and controls the video bitrate (1). The RTP packets are picked from each queue based on some defined priority order or simply in a round robin fashion (5). A transmission scheduler takes care of the transmission of RTP

packets, to be written to the UDP socket (6). In the general case all media must go through the transmission scheduler and is allowed to be transmitted if the number of bytes in flight is less than the congestion window. Audio frames can however be allowed to be transmitted immediately as audio is typically low bitrate and thus contributes little to congestion, this is something that is left as an implementation choice. RTCP packets are received (7) and the information about bytes in flight and congestion window is exchanged between the network congestion control and the transmission scheduler (8).

4.1.1. Constants and Parameter values

A set of constants are defined in Table 1, state variables are defined in Table 2. And finally, local variables are described in Table 3.

An init value [] indicates an empty array.

Constant	Explanation	Value
OWD_TARGET_LO	Min OWD target	0.1s
OWD_TARGET_HI	Max OWD target	0.4s
MAX_BYTES_IN_FLIGHT_HEAD_ROOM	Headroom for limitation of CWND	1.1
GAIN	Gain factor for congestion window adjustment	1.0
BETA	CWND scale factor due to loss event	0.6
BETA_R	Target rate scale factor due to loss event	0.8
BYTES_IN_FLIGHT_SLACK	Additional slack [%] to the congestion window	10%
RATE_ADJUST_INTERVAL	Interval between video bitrate adjustments	0.1s
FRAME_PERIOD	Video coder frame period [s]	
TARGET_BITRATE_MIN	Min target_bitrate [bps]	
TARGET_BITRATE_MAX	Max target_bitrate [bps]	
RAMP_UP_TIME	Timespan [s] from lowest to highest bitrate	10s
PRE_CONGESTION_GUARD	Guard factor against early congestion onset. A higher value gives less jitter possibly at the expense of a lower video bitrate.	0.0..0.2
TX_QUEUE_SIZE_FACTOR	Guard factor against RTP queue buildup	0.0..2.0

Table 1: Constants

Variable	Explanation	Init value
owd_target	OWD target	OWD_TARGET_LO
owd_fraction_avg	EWMA filtered owd_fraction	0.0
owd_fraction_hist	Vector of the last 20 owd_fraction	[]
owd_trend	OWD trend, indicates incipient congestion	0.0
owd_norm_hist	Vector of the last 100 owd_norm	[]
mss	Maximum segment size = Max RTP packet size [byte]	1000
min_cwnd	Minimum congestion window [byte]	2*MSS
in_fast_start	True if in fast start state	true
cwnd	Congestion window [byte]	min_cwnd
cwnd_i	Congestion window inflection point	1
bytes_newly_acked	The number of bytes that was acknowledged with the last received acknowledgement i.e. bytes acknowledged since the last CWND update [byte]. Reset after a CWND update	0
send_wnd	Upper limit of how many bytes that can be transmitted [byte]. Updated when CWND is updated and when RTP packet is transmitted	0
t_pace	Approximate estimate of inter- packet transmission	0.001

age_vec	interval [s], updated when RTP packet transmitted A vector of the last 20 RTP packet queue delay samples	[]
frame_skip_intensity	Indicates the intensity of the frame skips	0.0
since_last_frame_skip	Number of video frames since the last skip	0
consecutive_frame_skips	Number of consecutive frame skips	0
target_bitrate	Video target bitrate [bps]	TARGET_BITRATE_MIN
target_bitrate_i	Video target bitrate inflection point i.e. the last known highest target_bitrate during fast start. Used to limit bitrate increase close to the last known congestion point	1
rate_transmit	Measured transmit bitrate [bps]	0.0
rate_acked	Measured throughput based on received acknowledgements [bps]	0.0
s_rtt	Smoothed RTT [s], computed similar to method depicted in [RFC6298]	0.0
rtp_queue_size	Size of RTP packets in queue [bits]	0
rtp_size	Size of the last transmitted RTP packets [byte]	0
frame_skip	Skip encoding of video frame if	false

```

|                                     | true                                     |
+-----+-----+-----+-----+-----+-----+

```

Table 2: State variables

Variable	Explanation
owd	OWD = One way delay with base delay subtracted [s]. This is an estimate of the network queueing delay.
owd_fraction	OWD as a fraction of the OWD target
owd_norm	OWD normalized to OWD_TARGET_LO
owd_norm_mean	Average OWD norm over the last 100 samples
owd_norm_mean_sh	Average OWD norm over the last 20 samples
owd_norm_var	OWD norm variance over the last 100 samples
off_target	Relation between OWD and OWD target
scl_i	A general scalefactor that is applied to the CWND or target_bitrate increase
x_cwnd	Additional increase of CWND, used when send_wnd is computed
pace_bitrate	The allowed RTP packet transmission rate, used in the computation of t_pace [bps]
age_avg	Average RTP queue delay [s]
increment	Allowed target_bitrate increase
current_rate	Max of rate_transmit and rate_acked

Table 3: Local temporary variables

4.1.2. Network congestion control

This section explains the network congestion control, it contains two main functions

- o Computation of congestion window at the sender: Gives an upper limit to the number of bytes in flight i.e. how many bytes that have been transmitted but not yet acknowledged.
- o Transmission scheduling at the sender: RTP packets are transmitted if allowed by the relation between the number of bytes in flight and the congestion window. This is controlled by the send window.

Unlike TCP, SCReAM is not a byte oriented protocol, rather it is an RTP packet oriented protocol. Thus it keeps a list of transmitted RTP packets and their respective sending times (wall-clock time). The feedback indicates the highest received RTP sequence number and a

timestamp (wall-clock time) when it was received. In addition, an ACK list is included to make it possible to determine lost packets.

4.1.2.1. Congestion window update

The congestion window is computed from the one way (extra) delay estimates (OWD) that are obtained from the send and received timestamp of the RTP packets. LEDBAT [RFC6817] explains the details of the computation of the OWD. An OWD sample is obtained for each received acknowledgement. No smoothing of the OWD samples occur, however some smoothing occurs anyway as the computation of the CWND is in itself a low pass filter function.

SCReAM uses the terminology "Bytes in flight (`bytes_in_flight`)" which is computed as the sum of the sizes of the RTP packets ranging from the RTP packet most recently transmitted down to but not including the acknowledged packet with the highest sequence number. As an example: If RTP packet was sequence number SN with transmitted and the last ACK indicated SN-5 as the highest received sequence number then bytes in flight is computed as the sum of the size of RTP packets with sequence number SN-4, SN-3, SN-2, SN-1 and SN.

CWND is updated differently depending on whether the congestion control is in fast start or not and if a loss event is detected. A Boolean variable `in_fast_start` indicates if the congestion is in fast start state.

A loss event indicates one or more lost RTP packets within an RTT. This is detected by means of inspection for holes in the sequence number space in the acknowledgements with some margin for possible packet reordering in the network. As an alternative, a timer for loss detection similar to TCP RACK may be used.

Below is described the actions when an acknowledgement from the receiver is received.

`bytes_newly_acked` is updated.

The OWD fraction and an average of it are computed as

```
owd_fraction = owd/owd_target
```

```
owd_fraction_avg = 0.9* owd_fraction_avg + 0.1* owd_fraction
```

The OWD fraction is sampled every 50ms and the last 20 samples are stored in a vector (`owd_fraction_hist`). This vector is used in the computation of an OWD trend that gives a value between 0.0 and 1.0

depending on how close to congestion it is. The OWD trend is calculated as follows

Let $R(\text{owd_fraction_hist}, K)$ be the autocorrelation function of owd_fraction_hist at lag K . The 1st order prediction coefficient is formulated as

$$a = R(\text{owd_fraction_hist}, 1) / R(\text{owd_fraction_hist}, 0)$$

The prediction coefficient a has positive values if OWD shows an increasing trend, thus an indication of congestion is obtained before the OWD target is reached. The prediction coefficient is further multiplied with owd_fraction_avg to reduce sensitivity to increasing OWD when OWD is very small. The OWD trend is thus computed as

$$\text{owd_trend} = \max(0.0, \min(1.0, a * \text{owd_fraction_avg}))$$

The owd_trend is utilized in the media rate control and to determine when to exit slow start.

An off target value is computed as

$$\text{off_target} = (\text{owd_target} - \text{owd}) / \text{owd_target}$$

A temporal variable scl_i is computed as

$$\text{scl}_i = \max(0.2, \min(1.0, (\text{abs}(\text{cwnd} - \text{cwnd}_i) / \text{cwnd}_i * 4)^2))$$

scl_i is used to limit the CWND increase when close to the last known max value, before congestion was last detected.

The congestion window update depends on whether a loss event has occurred, and if the congestion control is if fast start or not.

On loss event:

If a loss event is detected then in_fast_start is set to false and CWND is updated according to

$$\text{cwnd}_i = \text{cwnd}$$

$$\text{cwnd} = \max(\text{min_cwnd}, \text{cwnd} * \text{BETA})$$

otherwise the CWND update continues

`in_fast_start = true:`

`in_fast_start` is set to false and `cwnd_i=cwnd` if `owd_trend >= 0.2` and otherwise CWND is updated according to

`cwnd = cwnd + bytes_newly_acked*scl_i`

`in_fast_start = false:`

Values of `off_target > 0.0` indicates that the congestion window can be increased. This is done according to the equations below.

`gain = GAIN*(1.0 + max(0.0, 1.0 - owd_trend/ 0.2))`

The equation above limits the gain when near congestion is detected

`gain *= scl_i`

This equation limits the gain when CWND is close to its last known max value

`cwnd += gain * off_target * bytes_newly_acked * mss / cwnd`

Values of `off_target <= 0.0` indicates congestion, CWND is then updated according to the equation

`cwnd += GAIN*off_target*bytes_newly_acked*mss/cwnd`

The equations above are very similar to what is specified in [RFC6817]. There are however a few differences.

- o [RFC6817] specifies a constant GAIN, this specification however limits the gain when CWND is increased dependent on near congestion state and the relation to the last known max CWND value.
 - o [RFC6817] specifies that the CWND increased is limited by an additional function controlled by a constant ALLOWED_INCREASE. This additional limitation is removed in this specification.
-

A number of final steps in the congestion window update procedure are outlined below

Resume fast start:

Fast start can be resumed in order to speed up the bitrate increase in case congestion abates. The condition to resume fast start (`in_fast_start = true`) is that `owd_trend` is less than 0.2 for 1.0 seconds or more.

Competing flows compensation, adjustment of `owd_target`:

Competing flows compensation is needed to avoid that flows congestion controlled by SCReAM are starved out by flows that are more aggressive in their nature. The `owd_target` is adjusted according to the `owd_norm_mean_sh` whenever `owd_norm_var` is below a given value. The condition to update `owd_target` is fulfilled if `owd_norm_var < 0.16` (indicating that the standard deviation is less than 0.4). `owd_target` is then update as:

```
owd_target = min(OWD_TARGET_HI, max(OWD_TARGET_LO, owd_norm_mean_sh*
OWD_TARGET_LO*1.1))
```

Final CWND adjustment step:

The congestion window is limited by the maximum number of bytes in flight over the last 1.0 seconds according to

```
cwnd = min(cwnd, max_bytes_in_flight*MAX_BYTES_IN_FLIGHT_HEAD_ROOM)
```

This avoids possible over-estimation of the throughput after for example, idle periods.

Finally `cwnd` is set to ensure that it is at least `min_cwnd`

```
cwnd = max(cwnd, MIN_CWND)
```

4.1.2.2. Transmission scheduling

The principle is to allow packet transmission of an RTP packet only if the number of bytes in flight is less than the congestion window. There are however two reasons why this strict rule will not work optimally:

- o Bitrate variations: The video frame size is always varying to a larger or smaller extent, a strict rule as the one given above will have the effect that the video bitrate have difficulties to increase as the congestion window puts a too hard restriction on the video frame size variation, this further can lead to occasional queuing of RTP packets in the RTP packet queue that will prevent bitrate increase because of the increased RTP queue size.
- o Reverse (feedback) path congestion: Especially in transport over buffer-bloated networks, the one way delay in the reverse direction may jump due to congestion. The effect of this is that the acknowledgements are delayed with the result that the self-clocking is temporarily halted, even though the forward path is not congested.

Packets are transmitted at a pace given by the send window, computed below

The send window is computed differently depending on OWD and its relation to the OWD target.

- o If $owd > owd_target$:
The send window is computed as
 $send_wnd = cwnd - bytes_in_flight$
This enforces a strict rule that helps to prevent further queue buildup.
- o If $owd \leq owd_target$:
A helper variable
 $x_cwnd = 1.0 + BYTES_IN_FLIGHT_SLACK * \max(0.0, \min(1.0, 1.0 - owd_trend / 0.5)) / 100.0$
is computed. The send window is computed as
 $send_wnd = \max(cwnd * x_cwnd, cwnd + mss) - bytes_in_flight$
This gives a slack that reduces as congestion increases, BYTES_IN_FLIGHT_SLACK is a maximum allowed slack in percent. A large value increases the robustness to bitrate variations in the source and congested feedback channel issues. The possible drawback is increased delay or packet loss when forward path congestion occur.

4.1.3. Video rate control

The video rate control is operated based on the size of the RTP packet send queue and observed loss events. In addition, owd_trend is also considered in the rate control, this to reduce the amount of induced network jitter.

A variable `target_bitrate` is adjusted depending on the congestion state. The target bitrate can vary between a minimum value (`target_bitrate_min`) and a maximum value (`target_bitrate_max`).

For the overall bitrate adjustment, two network throughput estimates are computed :

- o `rate_transmit`: The measured transmit bitrate
- o `rate_acked`: The ACKed bitrate, i.e. the volume of ACKed bits per time unit.

Both estimates are updated every 200ms.

The current throughput `current_rate` is computed as the maximum value of `rate_transmit` and `rate_acked`. The rationale behind the use of `rate_acked` in addition to `rate_transmit` is that `rate_transmit` is affected also by the amount of data that is available to transmit, thus a lack of data to transmit can be seen as reduced throughput that may itself cause an unnecessary rate reduction. To overcome this shortcoming; `rate_acked` is used as well. This gives a more stable throughput estimate.

The bitrate is updated at regular intervals, given by `RATE_ADJUST_INTERVAL` and differently depending the fast start state

The rate change behavior depends on whether a loss event has occurred, and if the congestion control is if fast start or not.

On loss event:

First of all the `target_bitrate` is updated if a new loss event was indicated and the rate change procedure is exited.

```
target_bitrate_i = target_bitrate
```

```
target_bitrate = max(BETA_R* target_bitrate, TARGET_BITRATE_MIN)
```

If no loss event was indicated then the rate change procedure continues.

`in_fast_start = true:`

An allowed increment is computed based on the congestion level and the relation to `target_bitrate_i`

`scl_i = (target_bitrate - target_bitrate_i) / target_bitrate_i`

`increment = TARGET_BITRATE_MAX * RATE_ADJUST_INTERVAL / RAMP_UP_TIME * (1.0 - min(1.0, owd_trend / 0.1))`

`increment *= max(0.2, min(1.0, (scl_i * 4) ^ 2))`

`target_bitrate += increment`

`target_bitrate` is reduced further if congestion is detected.

`target_bitrate *= (1.0 - PRE_CONGESTION_GUARD * owd_trend)`

`target_bitrate = min(TARGET_BITRATE_MAX, max(TARGET_BITRATE_MIN, target_bitrate))`

`in_fast_start = false:`

`target_bitrate_i` is updated to the current value of `target_bitrate` if `in_fast_start` was true the last time the bitrate was updated.

A pre-congestion indicator is computed as

`pre_congestion = min(1.0, max(0.0, owd_fraction_avg - 0.3) / 0.7)`

`pre_congestion += owd_trend`

The target bitrate is computed as

`target_bitrate = current_rate * (1.0 - PRE_CONGESTION_GUARD * pre_congestion) - TX_QUEUE_SIZE_FACTOR * rtp_queue_size`

`target_bitrate = min(TARGET_BITRATE_MAX, max(TARGET_BITRATE_MIN, target_bitrate))`

4.2. SCReAM Receiver

The SCReAM receiver is very simple in its implementation. The task is to feedback acknowledgements of received packets. For that purpose a set of state variables are needed, these are explained in Table 4.

One set of state variables are maintained per stream.

Variable	Explanation	Init value
rx_timestamp	The wall clock timestamp when the latest RTP packet was received	0
highest_rtp_sequence_number	The highest received sequence number	0
ack_vector	A 16 bit vector that indicates received RTP packets with a sequence number lower than highest_rtp_sequence_number	0
n_loss	An 8 bit counter for the number of lost RTP packets, separate counters are maintained for each SSRC	0
n_ECN	An 8 bit counter for the number of ECN-CE marked RTP packets, separate counters are maintained for each SSRC	0
pending_feedback	Indicates that an RTP packet was received and that an RTCP packet can be generated when RTCP timing rules permit	false
last_transmit_t	Last time an RTCP packet was transmitted, this is used to ensure that RTCP feedback is generated fairly for all streams.	-1.0

Table 4: State variables

Upon reception of an RTP packet, the state variables in Table 4 should be updated and the RTCP processing function should be

notified. An RTCP packet is later generated based on the state variables, how often this is done depends on the RTCP bandwidth.

5. Feedback Message

The feedback is over RTCP [RFC3550] and is based on [RFC4585]. It is implemented as a transport layer feedback message (RTPFB), see proposed example in Figure 2. The feedback control information part (FCI) consists of the following elements.

- o Highest received RTP sequence number : The highest received RTP sequence number for the given SSRC
- o n_lost : Ackumulated number of lost RTP packets for the given SSRC
- o Timestamp: A timestamp value indicating when the last packet was received which makes it possible to compute the one way (extra) delay (OWD).
- o n_ECN : Ackumulated number of ECN-CE marked RTP packets for the given SSRC
- o Source quench bit (Q): Makes it possible to request the sender to reduce its congestion window. This is useful if WebRTC media is received from many hosts and it becomes necessary to balance the bitrates between the streams.

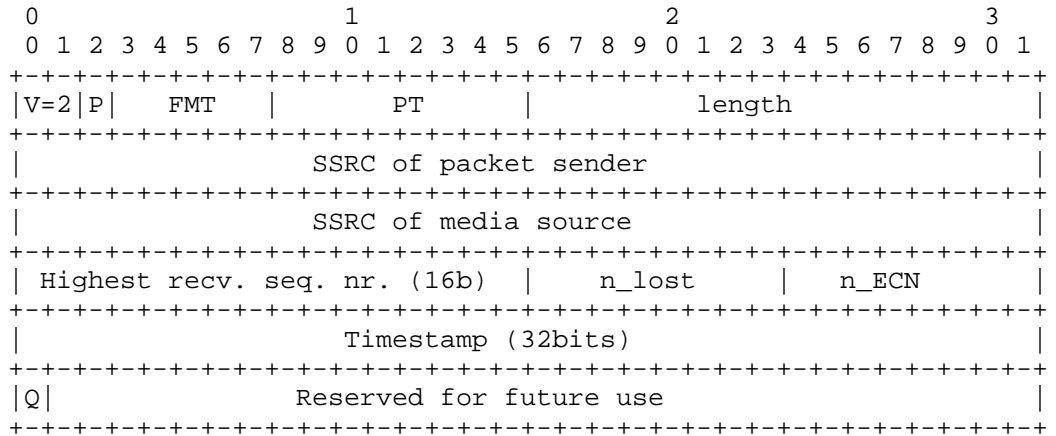


Figure 2: Transport layer feedback message

To make the feedback as frequent as possible, the feedback packets are transmitted as reduced size RTCP according to [RFC5506].

The timestamp clock time is recommended to be set to a fixed value such as 1000Hz, defined in this specification. The `n_lost` and `n_ECN` makes it possible to take necessary actions on the detection of lost and ECN marked packets.

Section 4 describes the main algorithm details and how the feedback is used.

6. Additional features

This section describes additional features. They are not required for the basic functionality of SCReAM but can improve performance in certain scenarios and topologies.

6.1. Packet pacing

Packet pacing is used in order to mitigate coalescing i.e. that packets are transmitted in bursts.

Packet pacing is enforced when `owd_fraction_avg` is greater than 0.1. The time interval between consecutive packet transmissions is then enforced to equal or higher than `t_pace` where `t_pace` is given by the equations below.

$$\text{pace_bitrate} = \max(50000, \text{cwnd} * 8 / \text{s_rtt})$$
$$\text{t_pace} = \text{rtp_size} * 8 / \text{pace_bitrate}$$

`rtp_size` is the size of the last transmitted RTP packet

6.2. Frame skipping

Frame skipping is a feature that makes it possible to reduce the size of the RTP queue in the cases that e.g. the channel throughput drops dramatically or even goes below the lowest possible video coder rate. Frame skipping is optional to implement as it can sometimes be difficult to realize e.g. due to lack of API function to support this.

Frame skipping is controlled by a flag `frame_skip` which, if set to 1 dictates that the video coder should skip the next video frame. The frame skipping intensity at the current time instant is computed according to the steps below

The queuing delay is sampled every frame period and the last 20 samples are stored in a vector `age_vec`

An average queuing delay is computed as a weighted sum over the samples in `age_vec`. `age_avg` at the current time instant is computed as

$$\text{age_avg}(n) = \text{SUM } \text{age_vec}(n-k) * w(k) \quad k = [0..20[$$

`w(n)` are weight factors arranged to give the most recent samples a higher weight.

The change in `age_avg` is computed as

$$\text{age_d} = \text{age_avg}(n) - \text{age_avg}(n-1)$$

The frame skipping intensity at the current time instant `n` is computed as

- o If `age_d > 0` and `age_avg > 2*FRAME_PERIOD`:
`frame_skip_intensity = min(1.0, (age_vec(n)-2*FRAME_PERIOD)/(4*FRAME_PERIOD)`
- o Otherwise frame skip intensity is set to zero

The `skip_frame` flag is set depending on three variables

- o `frame_skip_intensity`
- o `since_last_frame_skip`, i.e the number of consecutive frames without frame skipping
- o `consecutive_frame_skips`, i.e the number of consecutive frame skips

The flag `skip_frame` is set to 1 if any of the conditions below is met, otherwise it is set to 0.

- o `age_vec(n) > 0.2 && consecutive_frame_skips < 5`
- o `frame_skip_intensity < 0.5 && since_last_frame_skip >= 1.0 / frame_skip_intensity`
- o `frame_skip_intensity >= 0.5 && consecutive_frame_skips < (frame_skip_intensity - 0.5) * 10`

The arrangement makes sure that no more than 4 frames are skipped in sequence, the rationale is to ensure that the input to the video encoder does not change too much, something that may give poor prediction gain.

6.3. Q-bit semantics (source quench)

The Q bit in the feedback is set by a receiver to signal that the sender should reduce the bitrate. The sender will in response to this reduce the congestion window with the consequence that the video bitrate decreases. A typical use case for source quench is when a receiver receives streams from sources located at different hosts and they all share a common bottleneck, typically it is difficult to apply any rate distribution signaling between the sending hosts. The solution is then that the receiver sets the Q bit in the feedback to the sender that should reduce its rate, if the streams share a common bottleneck then the released bandwidth due to the reduction of the congestion window for the flow that had the Q bit set in the feedback will be grabbed by the other flows that did not have the Q bit set. This is ensured by the opportunistic behavior of SCReAM's congestion control. The source quench will have no or little effect if the flows do not share the same bottleneck.

The reduction in congestion window is proportional to the amount of SCReAM RTCP feedback with the Q bit set, the below steps outline how the sender should react to RTCP feedback with the Q bit set. The reduction is done once per RTT. Let :

- o n = Number of received RTCP feedback messages in one RTT
- o n_q = Number of received RTCP feedback messages in one RTT, with Q bit set.

The new congestion window is then expressed as:

$$cwnd = \max(\text{MIN_CWND}, cwnd * (1.0 - 0.5 * n_q / n))$$

Note that CWND is adjusted at most once per RTT. Furthermore The CWND increase should be inhibited for one RTT if CWND has been decreased as a result of Q bits set in the feedback.

The required intensity of the Q-bit set in the feedback in order to achieve a given rate distribution depends on many factors such as RTT, video source material etc. The receiver thus need to monitor the change in the received video bitrate on the different streams and adjust the intensity of the Q-bit accordingly.

7. Discussion

This section covers a few open discussion points

- o RTCP feedback overhead: SCReAM benefits from a relatively frequent feedback. Experiments have shown that a feedback rate roughly

equal to the frame rate gives a stable self-clocking and robustness against loss of feedback. With a maximum bitrate of 1500kbps the RTCP feedback overhead is in the range 10-15kbps with reduced size RTCP, including IP and UDP framing, in other words the RTCP overhead is quite modest and should not pose a problem in the general case. Other solutions may be required in highly asymmetrical link capacity cases. Worth notice is that SCReAM can work with as low feedback rates as once every 200ms, this however comes with a higher sensitivity to loss of feedback and also a potential reduction in throughput.

- o AVPF mode: The RTCP feedback is based on AVPF regular mode. The SCReAM feedback is transmitted as reduced size RTCP so save overhead, it is however required to transmit full compound RTCP at regular intervals, this interval can be controlled by `trr-int` depicted in [RFC4585].
- o BETA, CWND scale factor due to loss: The BETA value is recommended to be higher than 0.5. The reason behind this is that congestion control for multimedia has to deal with a source that is rate limited. A file transfer has "unlimited" source bitrate in comparison. The outcome is that SCReAM must be a little more aggressive than a file transfer in order to not be out competed.

8. Conclusion

This memo describes a congestion control algorithm for RMCAT that it is particularly good at handling the quickly changing condition in wireless network such as LTE. The solution conforms to the packet conservation principle and leverages on novel congestion control algorithms and recent TCP research, together with media bitrate determined by sender queuing delay and given delay thresholds. The solution has shown potential to meet the goals of high link utilization and prompt reaction to congestion. The solution is realized with a new RFC4585 transport layer feedback message.

9. Open issues

A list of open issues.

- o Describe how clock drift compensation is done
- o Describe how FEC overhead is accounted for in `target_bitrate` computation
- o Investigate the impact of more sparse RTCP feedback, for instance once per RTT

10. Source code

Source code for SCReAM is available in two formats :

- o C++ code, that is apt for experimentation. The code maintained as Visual Studio project. This code can possibly be included in simulators such as NS3. Available at <https://github.com/EricssonResearch/scream>
- o OpenWebRTC implementation : Work in progress, see <http://www.openwebrtc.io/> for information about the OpenWebRTC project

11. Acknowledgements

We would like to thank the following persons for their comments, questions and support during the work that led to this memo: Markus Andersson, Bo Burman, Tomas Frankkila, Frederic Gabin, Laurits Hamm, Hans Hannu, Nikolas Hermanns, Stefan Haekansson, Erlendur Karlsson, Daniel Lindstroem, Mats Nordberg, Jonathan Samuelsson, Rickard Sjoeborg, Robert Swain, Magnus Westerlund, Stefan Aelund.

12. IANA Considerations

A new RFC4585 transport layer feedback message needs to be standardized.

13. Security Considerations

The feedback can be vulnerable to attacks similar to those that can affect TCP. It is therefore recommended that the RTCP feedback is at least integrity protected.

14. Change history

A list of changes:

- o -04 to -05 : ACK vector is replaced by a loss counter, PT is removed from feedback, references to source code added
- o -03 to -04 : Extensive changes due to review comments, code somewhat modified, frame skipping made optional
- o -02 to -03 : Added algorithm description with equations, removed pseudo code and simulation results
- o -01 to -02 : Updated GCC simulation results

- o -00 to -01 : Fixed a few bugs in example code

15. References

15.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.
- [RFC5506] Johansson, I. and M. Westerlund, "Support for Reduced-Size Real-Time Transport Control Protocol (RTCP): Opportunities and Consequences", RFC 5506, April 2009.
- [RFC6298] Paxson, V., Allman, M., Chu, J., and M. Sargent, "Computing TCP's Retransmission Timer", RFC 6298, June 2011.
- [RFC6817] Shalunov, S., Hazel, G., Iyengar, J., and M. Kuehlewind, "Low Extra Delay Background Transport (LEDBAT)", RFC 6817, December 2012.

15.2. Informative References

- [FACK] "Forward Acknowledgement: Refining TCP Congestion Control", 2006.
- [I-D.draft-sarker-rmcat-cellular-eval-test-cases] Sarker, Z., "Evaluation Test Cases for Interactive Real-Time Media over Cellular Networks", <<http://www.ietf.org/id/draft-sarker-rmcat-cellular-eval-test-cases-00.txt>>.
- [I-D.ietf-rmcat-app-interaction] Zanaty, M., Singh, V., Nandakumar, S., and Z. Sarker, "RTP Application Interaction with Congestion Control", draft-ietf-rmcat-app-interaction-01 (work in progress), October 2014.

[I-D.ietf-tcpm-newcwg]

Fairhurst, G., Sathiaseelan, A., and R. Secchi, "Updating TCP to support Rate-Limited Traffic", draft-ietf-tcpm-newcwg-08 (work in progress), February 2015.

[QoS-3GPP]

TS 23.203, 3GPP., "Policy and charging control architecture", June 2011, <http://www.3gpp.org/ftp/specs/archive/23_series/23.203/23203-990.zip>.

[TFWC]

University College London, "Fairer TCP-Friendly Congestion Control Protocol for Multimedia Streaming", December 2007, <<http://www-dept.cs.ucl.ac.uk/staff/M.Handley/papers/tfwc-conext.pdf>>.

Authors' Addresses

Ingemar Johansson
Ericsson AB
Laboratoriegrend 11
Luleae 977 53
Sweden

Phone: +46 730783289
Email: ingemar.s.johansson@ericsson.com

Zaheduzzaman Sarker
Ericsson AB
Laboratoriegrend 11
Luleae 977 53
Sweden

Phone: +46 761153743
Email: zaheduzzaman.sarker@ericsson.com

RMCAT WG
Internet-Draft
Intended status: Experimental
Expires: April 27, 2015

V. Singh
M. Nagy
J. Ott
Aalto University
L. Eggert
NetApp
October 24, 2014

Congestion Control Using FEC for Conversational Media
draft-singh-rmcat-adaptive-fec-01

Abstract

This document describes a new mechanism for conversational multimedia flows. The proposed mechanism uses Forward Error Correction (FEC) encoded RTP packets (redundant packets) along side the media packets to probe for available network capacity. A straightforward interpretation is, the sending endpoint increases the transmission rate by keeping the media rate constant but increases the amount of FEC. If no losses and discards occur, the endpoint can then increase the media rate. If losses occur, the redundant FEC packets help in recovering the lost packets. Consequently, the endpoint can vary the FEC bit rate to conservatively (by a small amount) or aggressively (by a large amount) probe for available network capacity.

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	2
2. Terminology	3
3. Concept: FEC for Congestion Control	4
3.1. States	6
3.2. Framework	7
3.3. FEC Scheme	8
3.4. Applicability to other RMCAT Schemes	9
4. Security Considerations	9
5. IANA Considerations	10
6. Acknowledgements	10
7. References	10
7.1. Normative References	10
7.2. Informative References	11
Appendix A. Simulations	12
Authors' Addresses	12

1. Introduction

The Real-time Transport Protocol (RTP) [RFC3550] is widely used in voice telephony and video conferencing systems. Many of these systems run over best-effort UDP/IP networks, and are required to implement congestion to adapt the transmission rate of the RTP streams to match the available network capacity, while maintaining the user-experience [I-D.ietf-rmcat-cc-requirements]. The circuit breakers [I-D.ietf-avtcore-rtp-circuit-breakers] describe a minimal set of conditions when an RTP stream is causing severe congestion and should cease transmission. Consequently, the congestion control algorithm are expected to avoid triggering these conditions.

Conversational multimedia systems use Negative Acknowledgment (NACK), Forward Error Correction (FEC), and Reference Picture Selection (RPS) to protect against packet loss. These are used in addition to the codec-dependent resilience methods (for e.g., full intra-refresh and error-concealment). In this way, the multimedia system is anyway trading off part of the transmission rate for redundancy or retransmissions to reduce the effects of packet loss. An endpoint

often prefers using FEC in high latency networks where retransmissions may arrive later than the playout time of the packet (due to the size of the dejitter buffer) [Holmer13]. Therefore, the endpoint needs to adapt the transmission rate to best fit the changing network capacity and the amount of redundancy based on the observed/expected loss rate and network latency. Figure 1 shows the applicability of different error-resilience schemes based on the end-to-end latency and the observed packet loss [Devadoss08].

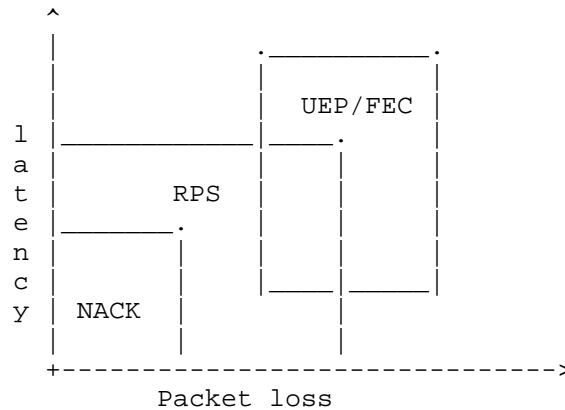


Figure 1: Applicability of Error Resilience Schemes based on the network delay and observed packet loss

In this document, we describe the use of FEC packets not only for error-resilience but also as a probing mechanism for congestion control (ramping up the transmission rate).

2. Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14, [RFC2119] and indicate requirement levels for compliant implementations.

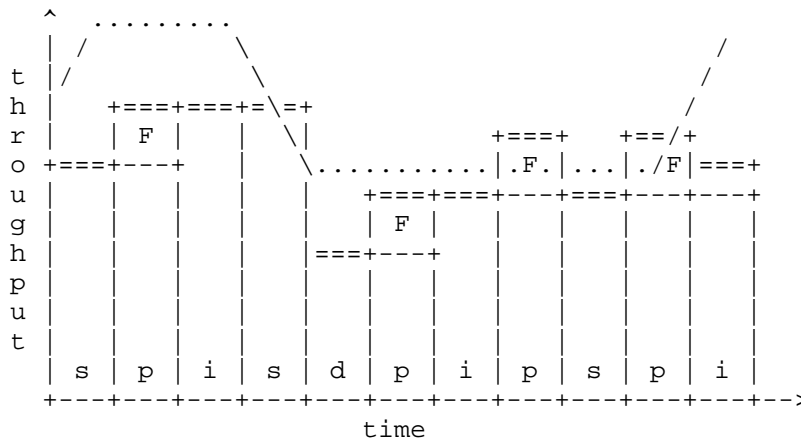
The terminology defined in RTP [RFC3550], RTP Profile for Audio and Video Conferences with Minimal Control [RFC3551], RTCP Extended Report (XR) [RFC3611], Extended RTP Profile for RTCP-based Feedback (RTP/AVPF) [RFC4585], RTP Retransmission Payload Format [RFC4588], Forward Error Correction (FEC) Framework [RFC6363], and Support for Reduced-Size RTCP [RFC5506] apply.

3. Concept: FEC for Congestion Control

FEC is one method for providing error-resilience, it improves reliability by adding redundant data to the primary media flow, which is used by received to recover packets that have been lost due to congestion or bit-errors. The congestion control algorithm on the other hand aims at maximizing the network path utilization, but risks over-estimating the available end-to-end network capacity leading to congestion (and therefore losses).

Figure 2 shows the timeline of enabling and disabling FEC. The main idea behind using FEC for congestion control is as follows: the sending endpoint chooses a high FEC rate to aggressively probe for available capacity and conversely chooses a low FEC rate to conservatively probe for available capacity. During the ramp up, if a packet is lost and the FEC packet arrives in time for decoding, the receiver is be able to recover the lost packet; if no packet is lost, the sender is able to increase the media encoding rate by swapping out a part of the FEC rate. This method can be especially useful when the transmission rate is close to the bottleneck link rate: by choosing an appropriate FEC rate, the endpoint is able to probe for available capacity without changing the target media rate and therefore not affecting the user-experience.

Hence, the congestion control algorithm is always able to probe for available capacity, as improved reliability compensates for possible errors resulting from probing for additional capacity (i.e., increase in observed latency and/or losses).



Key:

+====+ Media with minimal FEC for error protection

+====+
| F | Media with FEC for probing and error protection
+----+

.....
/ \ Available capacity

d,s,p,i: are the states: Decrease, Stay, Probe, Increase

Figure 2: Congestion Control enabling FEC.

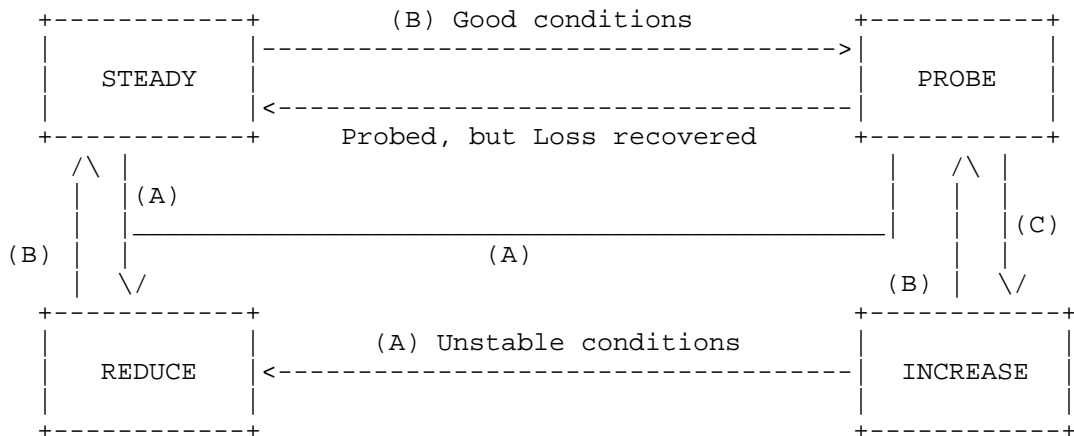


Figure 3: State machine of a Congestion Control enabling FEC.

3.1. States

The Figure 3 illustrates the the state machine of a congestion control algorithm incorporating FEC for probing. The state transitions occur based on the information reported in the feedback packet. In Figure 3 (A) indicates congestion, i.e., the congestion control observes increasing delay and/or packet loss, or any other congestion metric, and in response the congestion control reduces the transmission rate. In Figure 3 (B) occurs when the congestion cues report improvement in congestion metrics, and in response the congestion cue increases the transmission rate. For probing using FEC, the congestion control needs to map to the following 4 states: STEADY, PROBE, INCREASE, and REDUCE.

- o STEADY state: The congestion control keeps the same target media rate and no additional FEC packets are generated for probing. This is a transient state, after which the congestion control either attempts to increase the transmission rate, or observes congestion and reduces the transmission rate.
- o REDUCE state: The congestion control reduces the transmission rate based on the observed congestion cues, and generated no additional FEC packets than the minimum required for error-resilience. If in subsequent reports the conditions improve, the congestion control can directly transition to the PROBE state.
- o PROBE state: The congestion control observes no congestion for two reporting intervals (i.e., the transmission rate should be increased). The endpoint maintains the same target media bit rate, and instead increases the amount of FEC packets, thereby increasing the transmission rate.
- o INCREASE state: The endpoint is sending FEC packets and the congestion control observes no congestion (as reported in RTCP feedback), the media transmission rate is increased while maintaining minimal amount of FEC for error protection. In this case, the combined transmission rate (FEC+media) may remain the same as in the PROBE state. If the feedback reports packet loss, but some of these lost packets are recovered by the FEC packets, the congestion control can keep the same media bit rate and adjust the amount of FEC (compared to the previous PROBE state). If congestion is observed (the target rate calculated by the congestion control is much lower than the current media rate), the congestion control can transition to the REDUCE state and decrease the transmission rate.

3.2. Framework

The Figure 4 shows the interaction between the rate control module, the RTP and the FEC module.

At the sender, the rate control module calculates the new bit rate. If the new bit rate is higher than the previous than the previous bit rate indicates to the FEC module that the congestion control intends to probe. The FEC module depending on its internal state machine decides to add FEC for probing or not. Thereafter it indicates to the rate control module the bit rate remaining for the RTP media stream, which may be less than equal to the calculated bit rate.

At the reciver, the FEC module reconstructs lost packets in the primary stream from the packets received in the repair stream. If packets are repair it generates the post-repair loss report (discussed in Section 3.3) for the corresponding RTP packets.

At the sender, The FEC module also receives the RTCP Feedback related to the primary stream and any post-repair loss report. It uses the information from these RTCP reports to calculate the effectiveness of FEC for congestion control and is also the basis for changing its internal state.

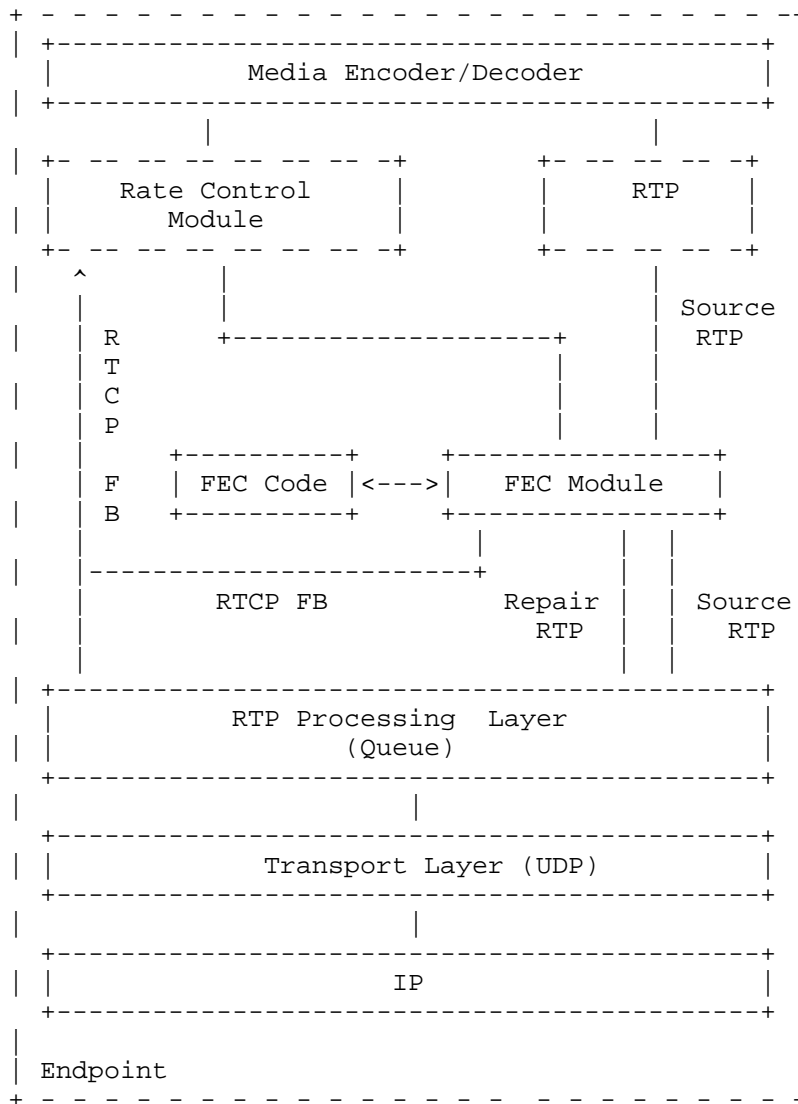


Figure 4: Interaction of Congestion Control and FEC Module.

3.3. FEC Scheme

[RFC6363] describes a framework for using Forward Error Correction (FEC) codes with RTP and allows any FEC code to be used with the framework. For this proposal, the FEC packets are created by XORing RTP media packets, the resulting redundant RTP packets are encoded

using the scheme defined in [I-D.singh-payload-rtp-ld2d-parity-scheme].

The endpoint MAY use a single-frame FEC (1-dimensional) or a multi-frame FEC (2-dimensional) for protecting the primary RTP stream. A single-frame FEC protects against a single packet loss and fails when burst loss occurs. Using multi-frame FEC helps mitigate these issues at the cost of higher overhead and latency in recovering lost packets. [Holmer13] shows examples of using a single- and multi-frame FEC.

The receiving endpoint may report the post-repair loss (or residual loss) using either the report block defined in [RFC5725] (Run-length encoding of packets repaired) or in [I-D.ietf-xrblock-rtcp-xr-post-repair-loss-count] (packet count of repaired packets).

Additionally, the receiving may report the occurrence of losses and discards via a run-length encoding (RLE) of lost [RFC3611] (Section 4.1), which enables the sender to detect the burst loss length and apply appropriate FEC scheme.

Packet that arrive too late to be played out by the receiver are discarded by the de-jitter buffer. Typically, the de-jitter buffer adjust the playout delay based on the observed frame inter-arrival delay, so that packets are played out smoothly. Reporting RLE of discarded packets [RFC7097] may further enable a sender to detect losses that occur after packet discards.

3.4. Applicability to other RMCAT Schemes

[Open issue: The current implementation is delay based and is documented in [Nagy14]. However, we would like to generalize the concept and apply it to different RMCAT algorithms for e.g., Google's Congestion Control algorithm [I-D.alvestrand-rmcat-congestion], SCReaM [I-D.johansson-rmcat-scream-cc], etc.]

4. Security Considerations

The security considerations of [RFC3550], RTP/AVPF profile for rapid RTCP feedback [RFC4585], circuit breaker [I-D.ietf-avtcore-rtp-circuit-breakers], and Generic Forward Error Correction [RFC5109] apply.

If non-authenticated RTCP reports are used, an on-path attacker can send forged RTCP feedback packets that can disrupt the operation of the underlying congestion control. Additionally, the forged packets can either indicate no packet loss causing the congestion control to

ramp-up quickly, or indicate high packet loss or RTT causing the circuit breaker to trigger.

5. IANA Considerations

There are no IANA impacts in this memo.

6. Acknowledgements

This document is based on the results published in [Nagy14].

The work of Varun Singh, and Joerg Ott has been partially supported by the European Institute of Innovation and Technology (EIT) ICT Labs activity RCLD 11882. The views expressed here are those of the author(s) only. Neither the European Commission nor the EITICT labs is liable for any use that may be made of the information in this document.

7. References

7.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.
- [RFC3551] Schulzrinne, H. and S. Casner, "RTP Profile for Audio and Video Conferences with Minimal Control", STD 65, RFC 3551, July 2003.
- [RFC3611] Friedman, T., Caceres, R., and A. Clark, "RTP Control Protocol Extended Reports (RTCP XR)", RFC 3611, November 2003.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.
- [RFC5506] Johansson, I. and M. Westerlund, "Support for Reduced-Size Real-Time Transport Control Protocol (RTCP): Opportunities and Consequences", RFC 5506, April 2009.

[I-D.ietf-avtcore-rtp-circuit-breakers]

Perkins, C. and V. Singh, "Multimedia Congestion Control: Circuit Breakers for Unicast RTP Sessions", draft-ietf-avtcore-rtp-circuit-breakers-05 (work in progress), February 2014.

[I-D.singh-payload-rtp-ld2d-parity-scheme]

Singh, V., Begen, A., and M. Zanaty, "RTP Payload Format for Non-Interleaved and Interleaved Parity Forward Error Correction (FEC)", draft-singh-payload-rtp-ld2d-parity-scheme-00 (work in progress), October 2014.

[I-D.ietf-xrblock-rtcp-xr-post-repair-loss-count]

Huang, R. and V. Singh, "RTP Control Protocol (RTCP) Extended Report (XR) for Post-Repair Loss Count Metrics", draft-ietf-xrblock-rtcp-xr-post-repair-loss-count-05 (work in progress), June 2014.

7.2. Informative References

[RFC4588] Rey, J., Leon, D., Miyazaki, A., Varsa, V., and R. Hakenberg, "RTP Retransmission Payload Format", RFC 4588, July 2006.

[RFC6363] Watson, M., Begen, A., and V. Roca, "Forward Error Correction (FEC) Framework", RFC 6363, October 2011.

[I-D.ietf-rmcat-cc-requirements]

Jesup, R., "Congestion Control Requirements For RMCAT", draft-ietf-rmcat-cc-requirements-02 (work in progress), February 2014.

[I-D.alvestrand-rmcat-congestion]

Holmer, S., Cicco, L., Mascolo, S., and H. Alvestrand, "A Google Congestion Control Algorithm for Real-Time Communication", draft-alvestrand-rmcat-congestion-02 (work in progress), February 2014.

[I-D.johansson-rmcat-scream-cc]

Johansson, I. and Z. Sarker, "Self-Clocked Rate Adaptation for Multimedia", draft-johansson-rmcat-scream-cc-02 (work in progress), June 2014.

[I-D.sarker-rmcat-eval-test]

Sarker, Z., Singh, V., Zhu, X., and M. Ramalho, "Test Cases for Evaluating RMCAT Proposals", draft-sarker-rmcat-eval-test-00 (work in progress), February 2014.

- [RFC5109] Li, A., "RTP Payload Format for Generic Forward Error Correction", RFC 5109, December 2007.
- [RFC5725] Begen, A., Hsu, D., and M. Lague, "Post-Repair Loss RLE Report Block Type for RTP Control Protocol (RTCP) Extended Reports (XRs)", RFC 5725, February 2010.
- [RFC7097] Ott, J., Singh, V., and I. Curcio, "RTP Control Protocol (RTCP) Extended Report (XR) for RLE of Discarded Packets", RFC 7097, January 2014.
- [Nagy14] Nagy, M., Singh, V., Ott, J., and L. Eggert, "Congestion Control using FEC for Conversational Multimedia Communication", Proc. of 5th ACM International Conference on Multimedia Systems (MMSys 2014) , 3 2014.
- [Devadoss08]
Devadoss, J., Singh, V., Ott, J., Liu, C., Wang, Y-K., and I. Curcio, "Evaluation of Error Resilience Mechanisms for 3G Conversational Video", Proc. of IEEE International Symposium on Multimedia (ISM 2008) , 3 2014.
- [Holmer13]
Holmer, S., Shemer, M., and M. Paniconi, "Handling Packet Loss in WebRTC", Proc. of IEEE International Conference on Image Processing (ICIP 2013) , 9 2013.

Appendix A. Simulations

This document is based on the results published in [Nagy14]. See the paper for ns-2 and testbed results; more results based on the scenarios listed in [I-D.sarker-rmcat-eval-test] will be published shortly.

Authors' Addresses

Varun Singh
Aalto University
School of Electrical Engineering
Otakaari 5 A
Espoo, FIN 02150
Finland

Email: varun@comnet.tkk.fi
URI: <http://www.netlab.tkk.fi/~varun/>

Marcin Nagy
Aalto University
School of Electrical Engineering
Otakaari 5 A
Espoo, FIN 02150
Finland

Email: marcin.nagy@aalto.fi

Joerg Ott
Aalto University
School of Electrical Engineering
Otakaari 5 A
Espoo, FIN 02150
Finland

Email: jo@comnet.tkk.fi

Lars Eggert
NetApp
Sonnenallee 1
Kirchheim 85551
Germany

Phone: +49 151 12055791
Email: lars@netapp.com
URI: <http://eggert.org/>

RMCAT WG
Internet-Draft
Intended status: Experimental
Expires: September 21, 2016

V. Singh
callstats.io
M. Nagy
J. Ott
Aalto University
L. Eggert
NetApp
March 20, 2016

Congestion Control Using FEC for Conversational Media
draft-singh-rmcat-adaptive-fec-03

Abstract

This document describes a new mechanism for conversational multimedia flows. The proposed mechanism uses Forward Error Correction (FEC) encoded RTP packets (redundant packets) along side the media packets to probe for available network capacity. A straightforward interpretation is, the sending endpoint increases the transmission rate by keeping the media rate constant but increases the amount of FEC. If no losses and discards occur, the endpoint can then increase the media rate. If losses occur, the redundant FEC packets help in recovering the lost packets. Consequently, the endpoint can vary the FEC bit rate to conservatively (by a small amount) or aggressively (by a large amount) probe for available network capacity.

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 September 21, 2016.

Copyright Notice

Copyright (c) 2016 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	2
2. Terminology	3
3. Concept: FEC for Congestion Control	4
3.1. States	6
3.2. Framework	7
3.3. FEC Scheme	8
3.4. Applicability to other RMCAT Schemes	9
4. Security Considerations	9
5. IANA Considerations	10
6. Acknowledgements	10
7. References	10
7.1. Normative References	10
7.2. Informative References	11
Appendix A. Simulations	13
Authors' Addresses	13

1. Introduction

The Real-time Transport Protocol (RTP) [RFC3550] is widely used in voice telephony and video conferencing systems. Many of these systems run over best-effort UDP/IP networks, and are required to implement congestion to adapt the transmission rate of the RTP streams to match the available network capacity, while maintaining the user-experience [I-D.ietf-rmcat-cc-requirements]. The circuit breakers [I-D.ietf-avtcore-rtp-circuit-breakers] describe a minimal set of conditions when an RTP stream is causing severe congestion and should cease transmission. Consequently, the congestion control algorithm are expected to avoid triggering these conditions.

Conversational multimedia systems use Negative Acknowledgment (NACK), Forward Error Correction (FEC), and Reference Picture Selection (RPS)

to protect against packet loss. These are used in addition to the codec-dependent resilience methods (for e.g., full intra-refresh and error-concealment). In this way, the multimedia system is anyway trading off part of the transmission rate for redundancy or retransmissions to reduce the effects of packet loss. An endpoint often prefers using FEC in high latency networks where retransmissions may arrive later than the playout time of the packet (due to the size of the dejitter buffer) [Holmer13]. Therefore, the endpoint needs to adapt the transmission rate to best fit the changing network capacity and the amount of redundancy based on the observed/expected loss rate and network latency. Figure 1 shows the applicability of different error-resilience schemes based on the end-to-end latency and the observed packet loss [Devadoss08].

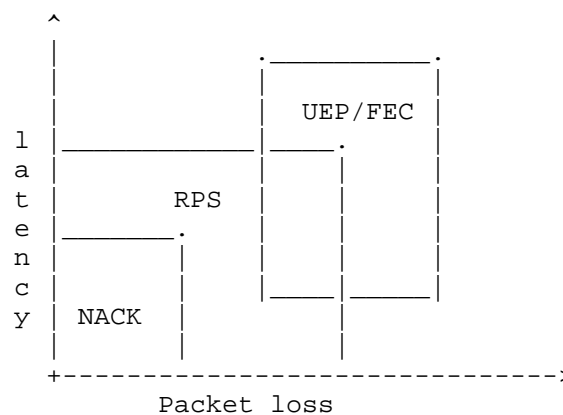


Figure 1: Applicability of Error Resilience Schemes based on the network delay and observed packet loss

In this document, we describe the use of FEC packets not only for error-resilience but also as a probing mechanism for congestion control (ramping up the transmission rate).

2. Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14, [RFC2119] and indicate requirement levels for compliant implementations.

The terminology defined in RTP [RFC3550], RTP Profile for Audio and Video Conferences with Minimal Control [RFC3551], RTCP Extended Report (XR) [RFC3611], Extended RTP Profile for RTCP-based Feedback (RTP/AVPF) [RFC4585], RTP Retransmission Payload Format [RFC4588],

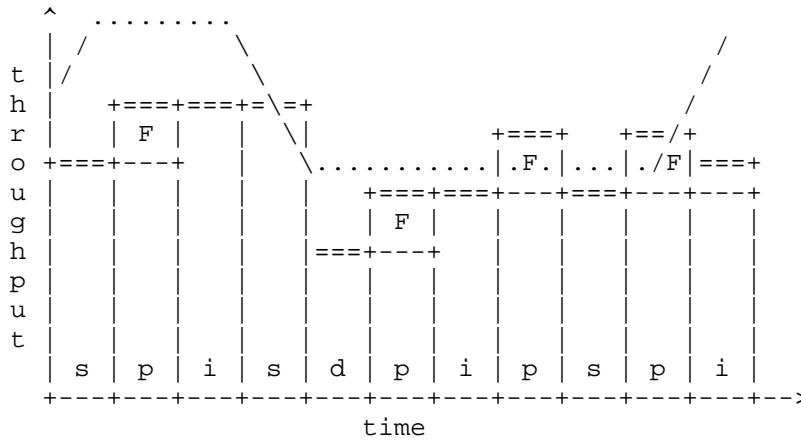
Forward Error Correction (FEC) Framework [RFC6363], and Support for Reduced-Size RTCP [RFC5506] apply.

3. Concept: FEC for Congestion Control

FEC is one method for providing error-resilience, it improves reliability by adding redundant data to the primary media flow, which is used by received to recover packets that have been lost due to congestion or bit-errors. The congestion control algorithm on the other hand aims at maximizing the network path utilization, but risks over-estimating the available end-to-end network capacity leading to congestion (and therefore losses).

Figure 2 shows the timeline of enabling and disabling FEC. The main idea behind using FEC for congestion control is as follows: the sending endpoint chooses a high FEC rate to aggressively probe for available capacity and conversely chooses a low FEC rate to conservatively probe for available capacity. During the ramp up, if a packet is lost and the FEC packet arrives in time for decoding, the receiver is be able to recover the lost packet; if no packet is lost, the sender is able to increase the media encoding rate by swapping out a part of the FEC rate. This method can be especially useful when the transmission rate is close to the bottleneck link rate: by choosing an appropriate FEC rate, the endpoint is able to probe for available capacity without changing the target media rate and therefore not affecting the user-experience.

Hence, the congestion control algorithm is always able to probe for available capacity, as improved reliability compensates for possible errors resulting from probing for additional capacity (i.e., increase in observed latency and/or losses).



Key:

====+ Media with minimal FEC for error protection

====+
| F | Media with FEC for probing and error protection
+----+

.....
/ \ Available capacity

d,s,p,i: are the states: Decrease, Stay, Probe, Increase

Figure 2: Congestion Control enabling FEC.

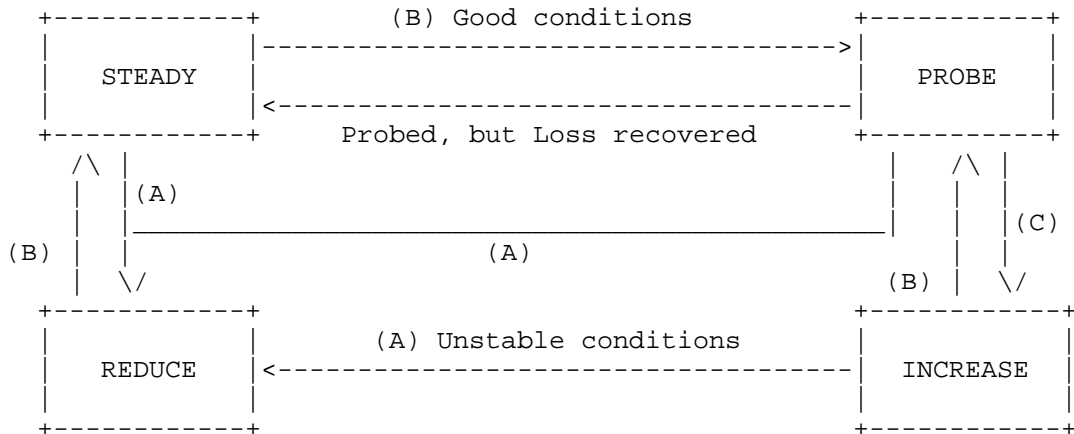


Figure 3: State machine of a Congestion Control enabling FEC.

3.1. States

The Figure 3 illustrates the the state machine of a congestion control algorithm incorporating FEC for probing. The state transitions occur based on the information reported in the feedback packet. In Figure 3 (A) indicates congestion, i.e., the congestion control observes increasing delay and/or packet loss, or any other congestion metric, and in response the congestion control reduces the transmission rate. In Figure 3 (B) occurs when the congestion cues report improvement in congestion metrics, and in response the congestion cue increases the transmission rate. For probing using FEC, the congestion control needs to map to the following 4 states: STEADY, PROBE, INCREASE, and REDUCE.

- o STEADY state: The congestion control keeps the same target media rate and no additional FEC packets are generated for probing. This is a transient state, after which the congestion control either attempts to increase the transmission rate, or observes congestion and reduces the transmission rate.
- o REDUCE state: The congestion control reduces the transmission rate based on the observed congestion cues, and generated no additional FEC packets than the minimum required for error-resilience. If in subsequent reports the conditions improve, the congestion control can directly transition to the PROBE state.
- o PROBE state: The congestion control observes no congestion for two reporting intervals (i.e., the transmission rate should be increased). The endpoint maintains the same target media bit rate, and instead increases the amount of FEC packets, thereby increasing the transmission rate.
- o INCREASE state: The endpoint is sending FEC packets and the congestion control observes no congestion (as reported in RTCP feedback), the media transmission rate is increased while maintaining minimal amount of FEC for error protection. In this case, the combined transmission rate (FEC+media) may remain the same as in the PROBE state. If the feedback reports packet loss, but some of these lost packets are recovered by the FEC packets, the congestion control can keep the same media bit rate and adjust the amount of FEC (compared to the previous PROBE state). If congestion is observed (the target rate calculated by the congestion control is much lower than the current media rate), the congestion control can transition to the REDUCE state and decrease the transmission rate.

3.2. Framework

The Figure 4 shows the interaction between the rate control module, the RTP and the FEC module.

At the sender, the rate control module calculates the new bit rate. If the new bit rate is higher than the previous than the previous bit rate indicates to the FEC module that the congestion control intends to probe. The FEC module depending on its internal state machine decides to add FEC for probing or not. Thereafter it indicates to the rate control module the bit rate remaining for the RTP media stream, which may be less than equal to the calculated bit rate.

At the reciver, the FEC module reconstructs lost packets in the primary stream from the packets received in the repair stream. If packets are repair it generates the post-repair loss report (discussed in Section 3.3) for the corresponding RTP packets.

At the sender, The FEC module also receives the RTCP Feedback related to the primary stream and any post-repair loss report. It uses the information from these RTCP reports to calculate the effectiveness of FEC for congestion control and is also the basis for changing its internal state.

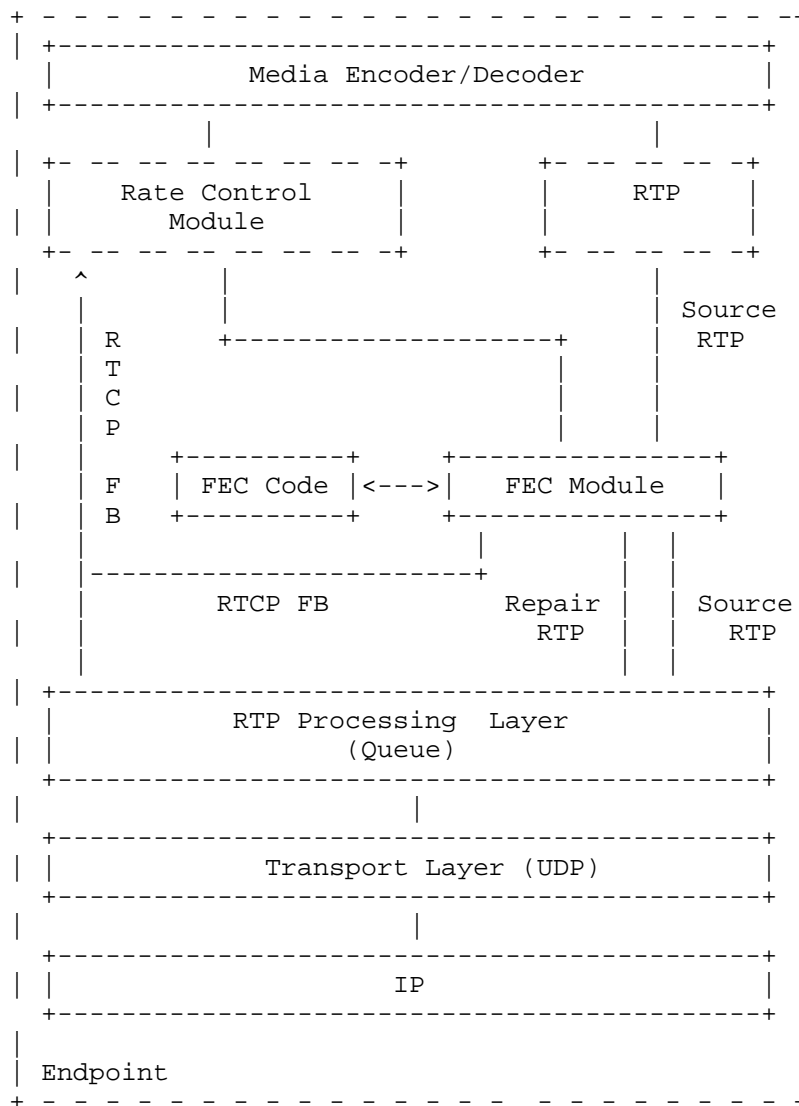


Figure 4: Interaction of Congestion Control and FEC Module.

3.3. FEC Scheme

[RFC6363] describes a framework for using Forward Error Correction (FEC) codes with RTP and allows any FEC code to be used with the framework. For this proposal, the FEC packets are created by XORing RTP media packets, the resulting redundant RTP packets are encoded using the scheme defined in [I-D.ietf-payload-flexible-fec-scheme].

The endpoint MAY use a single-frame FEC (1-dimensional) or a multi-frame FEC (2-dimensional) for protecting the primary RTP stream. A single-frame FEC protects against a single packet loss and fails when burst loss occurs. Using multi-frame FEC helps mitigate these issues at the cost of higher overhead and latency in recovering lost packets. [Holmer13] shows examples of using a single- and multi-frame FEC.

The receiving endpoint may report the post-repair loss (or residual loss) using either the report block defined in [RFC5725] (Run-length encoding of packets repaired) or in [RFC7509] (packet count of repaired packets).

Additionally, the receiving may report the occurrence of losses and discards via a run-length encoding (RLE) of lost [RFC3611] (Section 4.1), which enables the sender to detect the burst loss length and apply appropriate FEC scheme.

Packet that arrive too late to be played out by the receiver are discarded by the de-jitter buffer. Typically, the de-jitter buffer adjust the playout delay based on the observed frame inter-arrival delay, so that packets are played out smoothly. Reporting RLE of discarded packets [RFC7097] may further enable a sender to detect losses that occur after packet discards.

3.4. Applicability to other RMCAT Schemes

[Open issue: The current implementation is delay based and is documented in [Nagy14]. However, we would like to generalize the concept and apply it to different RMCAT algorithms for e.g., Google's Congestion Control algorithm [I-D.ietf-rmcat-gcc], SCReAM [I-D.ietf-rmcat-scream-cc], etc.]

4. Security Considerations

The security considerations of [RFC3550], RTP/AVPF profile for rapid RTCP feedback [RFC4585], circuit breaker [I-D.ietf-avtcore-rtcp-circuit-breakers], and Generic Forward Error Correction [RFC5109] apply.

If non-authenticated RTCP reports are used, an on-path attacker can send forged RTCP feedback packets that can disrupt the operation of the underlying congestion control. Additionally, the forged packets can either indicate no packet loss causing the congestion control to ramp-up quickly, or indicate high packet loss or RTT causing the circuit breaker to trigger.

5. IANA Considerations

There are no IANA impacts in this memo.

6. Acknowledgements

This document is based on the results published in [Nagy14].

The work of Varun Singh, and Joerg Ott has been partially supported by the European Institute of Innovation and Technology (EIT) ICT Labs activity RCLD 11882 (2012-2014). The views expressed here are those of the author(s) only. Neither the European Commission nor the EIT ICT labs is liable for any use that may be made of the information in this document.

Lars Eggert has received funding from the European Union's Horizon 2020 research and innovation programme 2014-2018 under grant agreement No. 644866. This document reflects only the authors' views and the European Commission is not responsible for any use that may be made of the information it contains.

7. References

7.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<http://www.rfc-editor.org/info/rfc2119>>.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, DOI 10.17487/RFC3550, July 2003, <<http://www.rfc-editor.org/info/rfc3550>>.
- [RFC3551] Schulzrinne, H. and S. Casner, "RTP Profile for Audio and Video Conferences with Minimal Control", STD 65, RFC 3551, DOI 10.17487/RFC3551, July 2003, <<http://www.rfc-editor.org/info/rfc3551>>.
- [RFC3611] Friedman, T., Ed., Caceres, R., Ed., and A. Clark, Ed., "RTP Control Protocol Extended Reports (RTCP XR)", RFC 3611, DOI 10.17487/RFC3611, November 2003, <<http://www.rfc-editor.org/info/rfc3611>>.

- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, DOI 10.17487/RFC4585, July 2006, <<http://www.rfc-editor.org/info/rfc4585>>.
- [RFC5506] Johansson, I. and M. Westerlund, "Support for Reduced-Size Real-Time Transport Control Protocol (RTCP): Opportunities and Consequences", RFC 5506, DOI 10.17487/RFC5506, April 2009, <<http://www.rfc-editor.org/info/rfc5506>>.
- [I-D.ietf-avtcore-rtp-circuit-breakers]
Perkins, C. and V. Varun, "Multimedia Congestion Control: Circuit Breakers for Unicast RTP Sessions", draft-ietf-avtcore-rtp-circuit-breakers-14 (work in progress), March 2016.
- [I-D.ietf-payload-flexible-fec-scheme]
Singh, V., Begen, A., Zanaty, M., and G. Mandyam, "RTP Payload Format for Flexible Forward Error Correction (FEC)", draft-ietf-payload-flexible-fec-scheme-01 (work in progress), October 2015.
- [RFC7509] Huang, R. and V. Singh, "RTP Control Protocol (RTCP) Extended Report (XR) for Post-Repair Loss Count Metrics", RFC 7509, DOI 10.17487/RFC7509, May 2015, <<http://www.rfc-editor.org/info/rfc7509>>.

7.2. Informative References

- [RFC4588] Rey, J., Leon, D., Miyazaki, A., Varsa, V., and R. Hakenberg, "RTP Retransmission Payload Format", RFC 4588, DOI 10.17487/RFC4588, July 2006, <<http://www.rfc-editor.org/info/rfc4588>>.
- [RFC6363] Watson, M., Begen, A., and V. Roca, "Forward Error Correction (FEC) Framework", RFC 6363, DOI 10.17487/RFC6363, October 2011, <<http://www.rfc-editor.org/info/rfc6363>>.
- [I-D.ietf-rmcat-cc-requirements]
Jesup, R. and Z. Sarker, "Congestion Control Requirements for Interactive Real-Time Media", draft-ietf-rmcat-cc-requirements-09 (work in progress), December 2014.

- [I-D.ietf-rmcat-gcc]
Holmer, S., Lundin, H., Carlucci, G., Cicco, L., and S. Mascolo, "A Google Congestion Control Algorithm for Real-Time Communication", draft-ietf-rmcat-gcc-01 (work in progress), October 2015.
- [I-D.ietf-rmcat-scream-cc]
Johansson, I. and Z. Sarker, "Self-Clocked Rate Adaptation for Multimedia", draft-ietf-rmcat-scream-cc-03 (work in progress), February 2016.
- [I-D.ietf-rmcat-eval-test]
Sarker, Z., Varun, V., Zhu, X., and M. Ramalho, "Test Cases for Evaluating RMCAT Proposals", draft-ietf-rmcat-eval-test-03 (work in progress), March 2016.
- [RFC5109] Li, A., Ed., "RTP Payload Format for Generic Forward Error Correction", RFC 5109, DOI 10.17487/RFC5109, December 2007, <<http://www.rfc-editor.org/info/rfc5109>>.
- [RFC5725] Begen, A., Hsu, D., and M. Lague, "Post-Repair Loss RLE Report Block Type for RTP Control Protocol (RTCP) Extended Reports (XRs)", RFC 5725, DOI 10.17487/RFC5725, February 2010, <<http://www.rfc-editor.org/info/rfc5725>>.
- [RFC7097] Ott, J., Singh, V., Ed., and I. Curcio, "RTP Control Protocol (RTCP) Extended Report (XR) for RLE of Discarded Packets", RFC 7097, DOI 10.17487/RFC7097, January 2014, <<http://www.rfc-editor.org/info/rfc7097>>.
- [Nagy14] Nagy, M., Singh, V., Ott, J., and L. Eggert, "Congestion Control using FEC for Conversational Multimedia Communication", Proc. of 5th ACM International Conference on Multimedia Systems (MMSys 2014) , 3 2014.
- [Devadoss08]
Devadoss, J., Singh, V., Ott, J., Liu, C., Wang, Y-K., and I. Curcio, "Evaluation of Error Resilience Mechanisms for 3G Conversational Video", Proc. of IEEE International Symposium on Multimedia (ISM 2008) , 3 2014.
- [Holmer13]
Holmer, S., Shemer, M., and M. Paniconi, "Handling Packet Loss in WebRTC", Proc. of IEEE International Conference on Image Processing (ICIP 2013) , 9 2013.

Appendix A. Simulations

This document is based on the results published in [Nagy14]. See the paper for ns-2 and testbed results; more results based on the scenarios listed in [I-D.ietf-rmcat-eval-test] will be published shortly.

Authors' Addresses

Varun Singh
Nemu Dialogue Systems Oy
Runeberginkatu 4c A 4
Helsinki 00100
Finland

Email: varun.singh@iki.fi
URI: <http://www.callstats.io/>

Marcin Nagy
Aalto University
School of Electrical Engineering
Otakaari 5 A
Espoo, FIN 02150
Finland

Email: marcin.nagy@aalto.fi

Joerg Ott
Aalto University
School of Electrical Engineering
Otakaari 5 A
Espoo, FIN 02150
Finland

Email: jo@comnet.tkk.fi

Lars Eggert
NetApp
Sonnenallee 1
Kirchheim 85551
Germany

Phone: +49 151 12055791
Email: lars@netapp.com
URI: <http://eggert.org/>

RTP Media Congestion Avoidance
Techniques (rmcat)
Internet-Draft
Intended status: Experimental
Expires: April 27, 2015

M. Welzl
S. Islam
S. Gjessing
University of Oslo
October 24, 2014

Coupled congestion control for RTP media
draft-welzl-rmcat-coupled-cc-04

Abstract

When multiple congestion controlled RTP sessions traverse the same network bottleneck, it can be beneficial to combine their controls such that the total on-the-wire behavior is improved. This document describes such a method for flows that have the same sender, in a way that is as flexible and simple as possible while minimizing the amount of changes needed to existing RTP applications.

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. Definitions 3
- 3. Limitations 4
- 4. Architectural overview 5
- 5. Roles 6
 - 5.1. SBD 6
 - 5.2. FSE 7
 - 5.3. Flows 7
 - 5.3.1. Example algorithm 1 - Active FSE 8
 - 5.3.2. Example algorithm 2 - Conservative Active FSE 9
- 6. Acknowledgements 10
- 7. IANA Considerations 10
- 8. Security Considerations 10
- 9. References 11
 - 9.1. Normative References 11
 - 9.2. Informative References 11
- Appendix A. Example algorithm - Passive FSE 12
 - A.1. Example operation (passive) 14
- Appendix B. Change log 19
 - B.1. Changes from -00 to -01 19
 - B.2. Changes from -01 to -02 19
 - B.3. Changes from -02 to -03 19
 - B.4. Changes from -03 to -04 19
- Authors' Addresses 19

1. Introduction

When there is enough data to send, a congestion controller must increase its sending rate until the path's capacity has been reached; depending on the controller, sometimes the rate is increased further, until packets are ECN-marked or dropped. This process inevitably creates undesirable queuing delay -- an effect that is amplified when multiple congestion controlled connections traverse the same network bottleneck. When such connections originate from the same host, it would therefore be ideal to use only one single sender-side congestion controller which determines the overall allowed sending rate, and then use a local scheduler to assign a proportion of this rate to each RTP session. This way, priorities could also be implemented quite easily, as a function of the scheduler; honoring user-specified priorities is, for example, required by rtcweb [rtcweb-usecases].

The Congestion Manager (CM) [RFC3124] provides a single congestion controller with a scheduling function just as described above. It is hard to implement because it requires an additional congestion controller and removes all per-connection congestion control functionality, which is quite a significant change to existing RTP based applications. This document presents a method that is easier to implement than the CM and also requires less significant changes to existing RTP based applications. It attempts to roughly approximate the CM behavior by sharing information between existing congestion controllers, akin to "Ensemble Sharing" in [RFC2140].

2. Definitions

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

Available Bandwidth:

The available bandwidth is the nominal link capacity minus the amount of traffic that traversed the link during a certain time interval, divided by that time interval.

Bottleneck:

The first link with the smallest available bandwidth along the path between a sender and receiver.

Flow:

A flow is the entity that congestion control is operating on. It could, for example, be a transport layer connection, an RTP session, or a subsession that is multiplexed onto a single RTP

session together with other subsessions.

Flow Group Identifier (FGI):

A unique identifier for each subset of flows that is limited by a common bottleneck.

Flow State Exchange (FSE):

The entity that maintains information that is exchanged between flows.

Flow Group (FG):

A group of flows having the same FGI.

Shared Bottleneck Detection (SBD):

The entity that determines which flows traverse the same bottleneck in the network, or the process of doing so.

3. Limitations

Sender-side only:

Coupled congestion control as described here only operates inside a single host on the sender side. This is because, irrespective of where the major decisions for congestion control are taken, the sender of a flow needs to eventually decide the transmission rate. Additionally, the necessary information about how much data an application can currently send on a flow is typically only available at the sender side, making the sender an obvious choice for placement of the elements and mechanisms described here.

When implementing a sender-side change to a congestion control mechanism such as TFRC [RFC5348], where receiver-side calculations make assumptions about the rate of the sender, the receiver also needs to be updated accordingly. Flows that have different senders but the same receiver, or different senders and different receivers can also share a bottleneck; such scenarios have been omitted for simplicity, and could be incorporated in future versions of this document. Note that limiting the number of flows on which coupled congestion control operates merely limits the benefits derived from the mechanism.

Shared bottlenecks do not change quickly:

As per the definition above, a bottleneck depends on cross traffic, and since such traffic can heavily fluctuate, bottlenecks can change at a high frequency (e.g., there can be oscillation between two or more links). This means that, when

flows are partially routed along different paths, they may quickly change between sharing and not sharing a bottleneck. For simplicity, here it is assumed that a shared bottleneck is valid for a time interval that is significantly longer than the interval at which congestion controllers operate. Note that, for the only SBD mechanism defined in this document (multiplexing on the same five-tuple), the notion of a shared bottleneck stays correct even in the presence of fast traffic fluctuations: since all flows that are assumed to share a bottleneck are routed in the same way, if the bottleneck changes, it will still be shared.

4. Architectural overview

Figure 1 shows the elements of the architecture for coupled congestion control: the Flow State Exchange (FSE), Shared Bottleneck Detection (SBD) and Flows. The FSE is a storage element that can be implemented in two ways: active and passive. In the active version, it initiates communication with flows and SBD. However, in the passive version, it does not actively initiate communication with flows and SBD; its only active role is internal state maintenance (e.g., an implementation could use soft state to remove a flow's data after long periods of inactivity). Every time a flow's congestion control mechanism would normally update its sending rate, the flow instead updates information in the FSE and performs a query on the FSE, leading to a sending rate that can be different from what the congestion controller originally determined. Using information about/from the currently active flows, SBD updates the FSE with the correct Flow State Identifiers (FSIs).

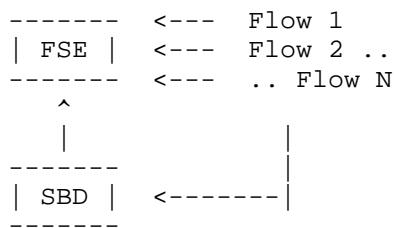


Figure 1: Coupled congestion control architecture

Since everything shown in Figure 1 is assumed to operate on a single host (the sender) only, this document only describes aspects that have an influence on the resulting on-the-wire behavior. It does,

for instance, not define how many bits must be used to represent FSIs, or in which way the entities communicate. Implementations can take various forms: for instance, all the elements in the figure could be implemented within a single application, thereby operating on flows generated by that application only. Another alternative could be to implement both the FSE and SBD together in a separate process which different applications communicate with via some form of Inter-Process Communication (IPC). Such an implementation would extend the scope to flows generated by multiple applications. The FSE and SBD could also be included in the Operating System kernel.

5. Roles

This section gives an overview of the roles of the elements of coupled congestion control, and provides an example of how coupled congestion control can operate.

5.1. SBD

SBD uses knowledge about the flows to determine which flows belong in the same Flow Group (FG), and assigns FGIs accordingly. This knowledge can be derived in three basic ways:

1. From multiplexing: it can be based on the simple assumption that packets sharing the same five-tuple (IP source and destination address, protocol, and transport layer port number pair) and having the same Differentiated Services Code Point (DSCP) in the IP header are typically treated in the same way along the path. The latter method is the only one specified in this document: SBD MAY consider all flows that use the same five-tuple and DSCP to belong to the same FG. This classification applies to certain tunnels, or RTP flows that are multiplexed over one transport (cf. [transport-multiplex]). In one way or another, such multiplexing will probably be recommended for use with rtcweb [rtcweb-rtp-usage].
2. Via configuration: e.g. by assuming that a common wireless uplink is also a shared bottleneck.
3. From measurements: e.g. by considering correlations among measured delay and loss as an indication of a shared bottleneck.

The methods above have some essential trade-offs: e.g., multiplexing is a completely reliable measure, however it is limited in scope to two end points (i.e., it cannot be applied to couple congestion controllers of one sender talking to multiple receivers). A measurement-based SBD mechanism is described in [sbd]. Measurements

can never be 100% reliable, in particular because they are based on the past but applying coupled congestion control means to make an assumption about the future; it is therefore recommended to implement cautionary measures, e.g. by disabling coupled congestion control if enabling it causes a significant increase in delay and/or packet loss. Measurements also take time, which entails a certain delay for turning on coupling (refer to [sbd] for details).

5.2. FSE

The FSE contains a list of all flows that have registered with it. For each flow, it stores the following:

- o a unique flow number to identify the flow
- o the FGI of the FG that it belongs to (based on the definitions in this document, a flow has only one bottleneck, and can therefore be in only one FG)
- o a priority P, which here is assumed to be represented as a floating point number in the range from 0.1 (unimportant) to 1 (very important). A negative value is used to indicate that a flow has terminated
- o The rate used by the flow in bits per second, FSE_R.

The FSE can operate on window-based as well as rate-based congestion controllers (TEMPORARY NOTE: and probably -- not yet tested -- combinations thereof, with calculations to convert from one to the other). In case of a window-based controller, FSE_R is a window, and all the text below should be considered to refer to window, not rates.

In the FSE, each FG contains one static variable S_CR which is meant to be the sum of the calculated rates of all flows in the same FG (including the flow itself). This value is used to calculate the sending rate.

The information listed here is enough to implement the sample flow algorithm given below. FSE implementations could easily be extended to store, e.g., a flow's current sending rate for statistics gathering or future potential optimizations.

5.3. Flows

Flows register themselves with SBD and FSE when they start, deregister from the FSE when they stop, and carry out an UPDATE function call every time their congestion controller calculates a new

sending rate. Via UPDATE, they provide the newly calculated rate and the desired rate (less than the calculated rate in case of application-limited flows, the same otherwise).

Below, two example algorithms are described. While other algorithms could be used instead, the same algorithm must be applied to all flows.

5.3.1. Example algorithm 1 - Active FSE

This algorithm was designed to be the simplest possible method to assign rates according to the priorities of flows. Simulations results in [fse] indicate that it does however not significantly reduce queuing delay and packet loss.

- (1) When a flow *f* starts, it registers itself with SBD and the FSE. FSE_R is initialized with the congestion controller's initial rate. SBD will assign the correct FGI. When a flow is assigned an FGI, it adds its FSE_R to S_CR.
- (2) When a flow *f* stops, its entry is removed from the list.
- (3) Every time the congestion controller of the flow *f* determines a new sending rate CC_R, the flow calls UPDATE, which carries out the tasks listed below to derive the new sending rates for all the flows in the FG. A flow's UPDATE function uses a local (i.e. per-flow) temporary variable S_P, which is the sum of all the priorities.

- (a) It updates S_CR.

$$S_CR = S_CR + CC_R - FSE_R(f)$$

- (b) It calculates the sum of all the priorities, S_P.

```
S_P = 0
for all flows i in FG do
    S_P = S_P + P(i)
end for
```

- (c) It calculates the sending rates for all the flows in an FG and distributes them.

```
for all flows i in FG do
    FSE_R(i) = (P(i)*S_CR)/S_P
```

```

        send FSE_R(i) to the flow i
    end for

```

5.3.2. Example algorithm 2 - Conservative Active FSE

This algorithm extends algorithm 1 to conservatively emulate the behavior of a single flow by proportionally reducing the aggregate rate on congestion. Simulations results in [fse] indicate that it can significantly reduce queuing delay and packet loss.

- (1) When a flow *f* starts, it registers itself with SBD and the FSE. FSE_R is initialized with the congestion controller's initial rate. SBD will assign the correct FGI. When a flow is assigned an FGI, it adds its FSE_R to S_CR.
- (2) When a flow *f* stops, its entry is removed from the list.
- (3) Every time the congestion controller of the flow *f* determines a new sending rate CC_R, the flow calls UPDATE, which carries out the tasks listed below to derive the new sending rates for all the flows in the FG. A flow's UPDATE function uses a local (i.e. per-flow) temporary variable S_P, which is the sum of all the priorities, and a local variable DELTA, which is used to calculate the difference between CC_R and the previously stored FSE_R. To prevent flows from either ignoring congestion or overreacting, a timer keeps them from changing their rates immediately after the common rate reduction that follows a congestion event. This timer is set to 2 RTTs of the flow that experienced congestion because it is assumed that a congestion event can persist for up to one RTT of that flow, with another RTT added to compensate for fluctuations in the measured RTT value.
 - (a) It updates S_CR based on DELTA.

```

if Timer has expired or not set then
    DELTA = CC_R - FSE_R(f)
    if DELTA < 0 then // Reduce S_CR proportionally
        S_CR = S_CR * CC_R / FSE_R(f)
        Set Timer for 2 RTTs
    else
        S_CR = S_CR + DELTA
    end if
end if

```

(b) It calculates the sum of all the priorities, S_P .

```
S_P = 0
for all flows i in FG do
  S_P = S_P + P(i)
end for
```

(c) It calculates the sending rates for all the flows in an FG and distributes them.

```
for all flows i in FG do
  FSE_R(i) = (P(i)*S_CR)/S_P
  send FSE_R(i) to the flow i
end for
```

6. Acknowledgements

This document has benefitted from discussions with and feedback from David Hayes, Andreas Petlund, and David Ros (who also gave the FSE its name).

This work was partially funded by the European Community under its Seventh Framework Programme through the Reducing Internet Transport Latency (RITE) project (ICT-317700).

7. IANA Considerations

This memo includes no request to IANA.

8. Security Considerations

In scenarios where the architecture described in this document is applied across applications, various cheating possibilities arise: e.g., supporting wrong values for the calculated rate, the desired rate, or the priority of a flow. In the worst case, such cheating could either prevent other flows from sending or make them send at a rate that is unreasonably large. The end result would be unfair behavior at the network bottleneck, akin to what could be achieved with any UDP based application. Hence, since this is no worse than UDP in general, there seems to be no significant harm in using this in the absence of UDP rate limiters.

In the case of a single-user system, it should also be in the

interest of any application programmer to give the user the best possible experience by using reasonable flow priorities or even letting the user choose them. In a multi-user system, this interest may not be given, and one could imagine the worst case of an "arms race" situation, where applications end up setting their priorities to the maximum value. If all applications do this, the end result is a fair allocation in which the priority mechanism is implicitly eliminated, and no major harm is done.

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.
- [RFC2140] Touch, J., "TCP Control Block Interdependence", RFC 2140, April 1997.
- [RFC3124] Balakrishnan, H. and S. Seshan, "The Congestion Manager", RFC 3124, June 2001.
- [RFC5348] Floyd, S., Handley, M., Padhye, J., and J. Widmer, "TCP Friendly Rate Control (TFRC): Protocol Specification", RFC 5348, September 2008.

9.2. Informative References

- [fse] Islam, S., Welzl, M., Gjessing, S., and N. Khademi, "Coupled Congestion Control for RTP Media", ACM SIGCOMM Capacity Sharing Workshop (CSWS 2014); extended version available as a technical report from <http://safiquili.at.ifi.uio.no/paper/fse-tech-report.pdf> , 2014.
- [rtcweb-rtp-usage] Perkins, C., Westerlund, M., and J. Ott, "Web Real-Time Communication (WebRTC): Media Transport and Use of RTP", draft-ietf-rtcweb-rtp-usage-18.txt (work in progress), October 2014.
- [rtcweb-usecases] Holmberg, C., Hakansson, S., and G. Eriksson, "Web Real-Time Communication Use-cases and Requirements", draft-ietf-rtcweb-use-cases-and-requirements-14.txt (work in progress), February 2014.

[sbd] Hayes, D., Ferlin, S., and M. Welzl, "Shared Bottleneck Detection for Coupled Congestion Control for RTP Media", draft-hayes-rmcat-sbd-00.txt (work in progress), October 2014.

[transport-multiplex] Westerlund, M. and C. Perkins, "Multiple RTP Sessions on a Single Lower-Layer Transport", draft-westerlund-avtcore-transport-multiplexing-07.txt (work in progress), October 2013.

Appendix A. Example algorithm - Passive FSE

Active algorithms calculate the rates for all the flows in the FG and actively distribute them. In a passive algorithm, UPDATE returns a rate that should be used instead of the rate that the congestion controller has determined. This can make a passive algorithm easier to implement; however, the resulting dynamics are not fully understood. The algorithm described below is to be considered as highly experimental and did not perform as well as the active variants in simulations.

This passive version of the FSE stores the following information in addition to the variables described in Section 5.2:

- o The desired rate DR. This can be smaller than the calculated rate if the application feeding into the flow has less data to send than the congestion controller would allow. In case of a bulk transfer, DR must be set to CC_R received from the flow's congestion module.

The passive version of the FSE contains one static variable per FG called TLO (Total Leftover Rate -- used to let a flow 'take' bandwidth from application-limited or terminated flows) which is initialized to 0. For the passive version, S_CR is limited to increase or decrease as conservatively as a flow's congestion controller decides in order to prohibit sudden rate jumps.

- (1) When a flow *f* starts, it registers itself with SBD and the FSE. FSE_R and DR are initialized with the congestion controller's initial rate. SBD will assign the correct FGI. When a flow is assigned an FGI, it adds its FSE_R to S_CR.
- (2) When a flow *f* stops, it sets its DR to 0 and sets P to -1.

(3) Every time the congestion controller of the flow f determines a new sending rate CC_R , assuming the flow's new desired rate new_DR to be "infinity" in case of a bulk data transfer with an unknown maximum rate, the flow calls `UPDATE`, which carries out the tasks listed below to derive the flow's new sending rate, $Rate$. A flow's `UPDATE` function uses a few local (i.e. per-flow) temporary variables, which are all initialized to 0: $DELTA$, new_S_CR and S_P .

(a) For all the flows in its FG (including itself), it calculates the sum of all the calculated rates, new_S_CR . Then it calculates the difference between $FSE_R(f)$ and CC_R , $DELTA$.

```

for all flows i in FG do
    new_S_CR = new_S_CR + FSE_R(i)
end for
DELTA = CC_R - FSE_R(f)

```

(b) It updates S_CR , $FSE_R(f)$ and $DR(f)$.

```

FSE_R(f) = CC_R
if DELTA > 0 then // the flow's rate has increased
    S_CR = S_CR + DELTA
else if DELTA < 0 then
    S_CR = new_S_CR + DELTA
end if
DR(f) = min(new_DR, FSE_R(f))

```

(c) It calculates the leftover rate TLO , removes the terminated flows from the FSE and calculates the sum of all the priorities, S_P .

```

for all flows i in FG do
    if P(i) < 0 then
        delete flow
    else
        S_P = S_P + P(i)
    end if
end for
if DR(f) < FSE_R(f) then
    TLO = TLO + (P(f)/S_P) * S_CR - DR(f)
end if

```

(d) It calculates the sending rate, Rate.

```
Rate = min(new_DR, (P(f)*S_CR)/S_P + TLO)

if Rate != new_DR and TLO > 0 then
    TLO = 0 // f has 'taken' TLO
end if
```

(e) It updates DR(f) and FSE_R(f) with Rate.

```
if Rate > DR(f) then
    DR(f) = Rate
end if
FSE_R(f) = Rate
```

The goals of the flow algorithm are to achieve prioritization, improve network utilization in the face of application-limited flows, and impose limits on the increase behavior such that the negative impact of multiple flows trying to increase their rate together is minimized. It does that by assigning a flow a sending rate that may not be what the flow's congestion controller expected. It therefore builds on the assumption that no significant inefficiencies arise from temporary application-limited behavior or from quickly jumping to a rate that is higher than the congestion controller intended. How problematic these issues really are depends on the controllers in use and requires careful per-controller experimentation. The coupled congestion control mechanism described here also does not require all controllers to be equal; effects of heterogeneous controllers, or homogeneous controllers being in different states, are also subject to experimentation.

This algorithm gives all the leftover rate of application-limited flows to the first flow that updates its sending rate, provided that this flow needs it all (otherwise, its own leftover rate can be taken by the next flow that updates its rate). Other policies could be applied, e.g. to divide the leftover rate of a flow equally among all other flows in the FGI.

A.1. Example operation (passive)

In order to illustrate the operation of the passive coupled congestion control algorithm, this section presents a toy example of two flows that use it. Let us assume that both flows traverse a common 10 Mbit/s bottleneck and use a simplistic congestion controller that starts out with 1 Mbit/s, increases its rate by 1 Mbit/s in the absence of congestion and decreases it by 2 Mbit/s in

the presence of congestion. For simplicity, flows are assumed to always operate in a round-robin fashion. Rate numbers below without units are assumed to be in Mbit/s. For illustration purposes, the actual sending rate is also shown for every flow in FSE diagrams even though it is not really stored in the FSE.

Flow #1 begins. It is a bulk data transfer and considers itself to have top priority. This is the FSE after the flow algorithm's step 1:

```
-----
| # | FGI | P | FSE_R | DR | Rate |
| 1 | 1 | 1 | 1 | 1 | 1 |
-----
```

S_CR = 1, TLO = 0

Its congestion controller gradually increases its rate. Eventually, at some point, the FSE should look like this:

```
-----
| # | FGI | P | FSE_R | DR | Rate |
| 1 | 1 | 1 | 10 | 10 | 10 |
-----
```

S_CR = 10, TLO = 0

Now another flow joins. It is also a bulk data transfer, and has a lower priority (0.5):

```
-----
| # | FGI | P | FSE_R | DR | Rate |
| 1 | 1 | 1 | 10 | 10 | 10 |
| 2 | 1 | 0.5 | 1 | 1 | 1 |
-----
```

S_CR = 11, TLO = 0

Now assume that the first flow updates its rate to 8, because the

total sending rate of 11 exceeds the total capacity. Let us take a closer look at what happens in step 3 of the flow algorithm.

CC_R = 8. new_DR = infinity.
 3 a) new_S_CR = 11; DELTA = 8 - 10 = -2.
 3 b) FSE_R(f) = 8. DELTA is negative, hence S_CR = 9;
 DR(f) = 8.
 3 c) S_P = 1.5.
 3 d) new sending rate = min(infinity, 1/1.5 * 9 + 0) = 6.
 3 e) FSE_R(f) = 6.

The resulting FSE looks as follows:

```
-----
```

#	FGI	P	FSE_R	DR	Rate
1	1	1	6	8	6
2	1	0.5	1	1	1

```
-----
```

S_CR = 9, TLO = 0

The effect is that flow #1 is sending with 6 Mbit/s instead of the 8 Mbit/s that the congestion controller derived. Let us now assume that flow #2 updates its rate. Its congestion controller detects that the network is not fully saturated (the actual total sending rate is 6+1=7) and increases its rate.

CC_R=2. new_DR = infinity.
 3 a) new_S_CR = 7; DELTA = 2 - 1 = 1.
 3 b) FSE_R(f) = 2. DELTA is positive, hence S_CR = 9 + 1 = 10;
 DR(f) = 2.
 3 c) S_P = 1.5.
 3 d) new sending rate = min(infinity, 0.5/1.5 * 10 + 0) = 3.33.
 3 e) DR(f) = FSE_R(f) = 3.33.

The resulting FSE looks as follows:

```
-----
```

#	FGI	P	FSE_R	DR	Rate
1	1	1	6	8	6
2	1	0.5	3.33	3.33	3.33

```
-----
```

S_CR = 10, TLO = 0

The effect is that flow #2 is now sending with 3.33 Mbit/s, which is close to half of the rate of flow #1 and leads to a total utilization of $6(\#1) + 3.33(\#2) = 9.33$ Mbit/s. Flow #2's congestion controller has increased its rate faster than the controller actually expected. Now, flow #1 updates its rate. Its congestion controller detects that the network is not fully saturated and increases its rate. Additionally, the application feeding into flow #1 limits the flow's sending rate to at most 2 Mbit/s.

CC_R=7. new_DR=2.

3 a) new_S_CR = 9.33; DELTA = 1.

3 b) FSE_R(f) = 7, DELTA is positive, hence S_CR = 10 + 1 = 11;
 DR = min(2, 7) = 2.

3 c) S_P = 1.5; DR(f) < FSE_R(f), hence TLO = $1/1.5 * 11 - 2 = 5.33$.

3 d) new sending rate = min(2, $1/1.5 * 11 + 5.33$) = 2.

3 e) FSE_R(f) = 2.

The resulting FSE looks as follows:

#	FGI	P	FSE_R	DR	Rate
1	1	1	2	2	2
2	1	0.5	3.33	3.33	3.33

S_CR = 11, TLO = 5.33

Now, the total rate of the two flows is $2 + 3.33 = 5.33$ Mbit/s, i.e. the network is significantly underutilized due to the limitation of flow #1. Flow #2 updates its rate. Its congestion controller detects that the network is not fully saturated and increases its rate.

CC_R=4.33. new_DR = infinity.
 3 a) new_S_CR = 5.33; DELTA = 1.
 3 b) FSE_R(f) = 4.33. DELTA is positive, hence S_CR = 12;
 DR(f) = 4.33.
 3 c) S_P = 1.5.
 3 d) new sending rate: $\min(\text{infinity}, 0.5/1.5 * 12 + 5.33) = 9.33$.
 3 e) FSE_R(f) = 9.33, DR(f) = 9.33.

The resulting FSE looks as follows:

#	FGI	P	FSE_R	DR	Rate
1	1	1	2	2	2
2	1	0.5	9.33	9.33	9.33

S_CR = 12, TLO = 0

Now, the total rate of the two flows is $2 + 9.33 = 11.33$ Mbit/s. Finally, flow #1 terminates. It sets P to -1 and DR to 0. Let us assume that it terminated late enough for flow #2 to still experience the network in a congested state, i.e. flow #2 decreases its rate in the next iteration.

CC_R = 7.33. new_DR = infinity.
 3 a) new_S_CR = 11.33; DELTA = -2.
 3 b) FSE_R(f) = 7.33. DELTA is negative, hence S_CR = 9.33;
 DR(f) = 7.33.
 3 c) Flow 1 has P = -1, hence it is deleted from the FSE.
 S_P = 0.5.
 3 d) new sending rate: $\min(\text{infinity}, 0.5/0.5*9.33 + 0) = 9.33$.
 3 e) FSE_R(f) = DR(f) = 9.33.

The resulting FSE looks as follows:

#	FGI	P	FSE_R	DR	Rate
2	1	0.5	9.33	9.33	9.33

S_CR = 9.33, TLO = 0

Appendix B. Change log

B.1. Changes from -00 to -01

- o Added change log.
- o Updated the example algorithm and its operation.

B.2. Changes from -01 to -02

- o Included an active version of the algorithm which is simpler.
- o Replaced "greedy flow" with "bulk data transfer" and "non-greedy" with "application-limited".
- o Updated new_CR to CC_R, and CR to FSE_R for better understanding.

B.3. Changes from -02 to -03

- o Included an active conservative version of the algorithm which reduces queue growth and packet loss; added a reference to a technical report that shows these benefits with simulations.
- o Moved the passive variant of the algorithm to appendix.

B.4. Changes from -03 to -04

- o Extended SBD section.
- o Added a note about window-based controllers.

Authors' Addresses

Michael Welzl
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Phone: +47 22 85 24 20
Email: michawe@ifi.uio.no

Safiqul Islam
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Phone: +47 22 84 08 37
Email: safiquli@ifi.uio.no

Stein Gjessing
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Phone: +47 22 85 24 44
Email: steing@ifi.uio.no

Network Working Group
Internet Draft
Intended Status: Informational
Expires: March 16, 2015

X. Zhu, R. Pan
M. A. Ramalho, S. M. de la Cruz
C. Ganzhorn, P. E. Jones
Cisco Systems
S. D'aronco
Ecole Polytechnique Federale de Lausanne
September 12, 2014

NADA: A Unified Congestion Control Scheme for Real-Time Media
draft-zhu-rmcat-nada-04

Abstract

This document describes a scheme named network-assisted dynamic adaptation (NADA), a novel congestion control approach for interactive real-time media applications, such as video conferencing. In the proposed scheme, the sender regulates its sending rate based on either implicit or explicit congestion signaling, in a unified approach. The scheme can benefit from explicit congestion notification (ECN) markings from network nodes. It also maintains consistent sender behavior in the absence of such markings, by reacting to queuing delays and packet losses instead.

We present here the overall system architecture, recommended behaviors at the sender and the receiver, as well as expected network node operations. Results from extensive simulation studies of the proposed scheme are available upon request.

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) 2012 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	4
2. Terminology	4
3. System Model	4
4. Network Node Operations	5
4.1 Default behavior of drop tail	5
4.2 ECN marking	5
4.3 PCN marking	6
4.4 Comments and Discussions	7
5. Receiver Behavior	7
5.1 Monitoring per-packet statistics	7
5.2 Aggregating congestion signals	8
5.3 Sending periodic feedback	8
5.4 Discussions on delay metrics	9
6. Sender Behavior	9
6.1 Video encoder rate control	10
6.2 Rate shaping buffer	11
6.3 Reference rate calculator	11
6.4 Video target rate and sending rate calculator	12
6.5 Start-up behavior	13
7. Incremental Deployment	13
8. Implementation Status	14
9. IANA Considerations	14
10. References	14
10.1 Normative References	14
10.2 Informative References	14

Authors' Addresses 15

1. Introduction

Interactive real-time media applications introduce a unique set of challenges for congestion control. Unlike TCP, the mechanism used for real-time media needs to adapt fast to instantaneous bandwidth changes, accommodate fluctuations in the output of video encoder rate control, and cause low queuing delay over the network. An ideal scheme should also make effective use of all types of congestion signals, including packet losses, queuing delay, and explicit congestion notification (ECN) markings.

Based on the above considerations, we present a scheme named network-assisted dynamic adaptation (NADA). The proposed design benefits from explicit congestion control signals (e.g., ECN markings) from the network, and remains compatible in the presence of implicit signals (delay or loss) only. In addition, it supports weighted bandwidth sharing among competing video flows.

This documentation describes the overall system architecture, recommended designs at the sender and receiver, as well as expected network nodes operations. The signaling mechanism consists of standard RTP timestamp [RFC3550] and standard RTCP feedback reports.

2. Terminology

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

3. System Model

The system consists of the following elements:

- * Incoming media stream, in the form of consecutive raw video frames and audio samples;
- * Media encoder with rate control capabilities. It takes the incoming media stream and encodes it to an RTP stream at a target bit rate R_v . Note that the actual output rate from the encoder R_o may fluctuate randomly around the target R_v . Also, the encoder can only change its rate at rather coarse time intervals, e.g., once every 0.5 seconds.
- * RTP sender, responsible for calculating the target bit rate R_n based on network congestion signals (delay or ECN marking reports from the receiver), and for regulating the actual sending rate R_s accordingly. A rate shaping buffer is employed

to absorb the instantaneous difference between video encoder output rate R_v and sending rate R_s . The buffer size L_s , together with R_n , influences the calculation of actual sending rate R_s and video encoder target rate R_v . The RTP sender also generates RTP timestamp in outgoing packets.

* RTP receiver, responsible for measuring and estimating end-to-end delay based on sender RTP timestamp. In the presence of packet losses and ECN markings, it also keeps track of packet loss and ECN marking ratios. It calculates the equivalent delay x_n that accounts for queuing delay, ECN marking, and packet losses, as well as the derivative (i.e., slope of change) of this congestion signal as x'_n . The receiver feeds both information (x_n and x'_n) back to the sender via periodic RTCP reports.

* Network node, with several modes of operation. The system can work with the default behavior of a simple drop tail queue. It can also benefit from advanced AQM features such as RED-based ECN marking, and PCN marking using a token bucket algorithm.

In the following, we will elaborate on the respective operations at the network node, the receiver, and the sender.

4. Network Node Operations

We consider three variations of queue management behavior at the network node, leading to either implicit or explicit congestion signals.

4.1 Default behavior of drop tail

In conventional network with drop tail or RED queues, congestion is inferred from the estimation of end-to-end delay and/or packet losses. Packet drops at the queue are detected at the receiver, and contributes to the calculation of the equivalent delay x_n . No special action is required at network node.

4.2 ECN marking

In this mode, the network node randomly marks the ECN field in the IP packet header following the Random Early Detection (RED) algorithm [RFC2309]. Calculation of the marking probability involves the following steps:

* upon packet arrival, update smoothed queue size q_{avg} as:

$$q_{avg} = \alpha * q + (1 - \alpha) * q_{avg}.$$

The smoothing parameter α is a value between 0 and 1. A value of $\alpha=1$ corresponds to performing no smoothing at all.

* calculate marking probability p as:

$$p = 0, \text{ if } q < q_{lo};$$

$$p = p_{max} * \frac{q_{avg} - q_{lo}}{q_{hi} - q_{lo}}, \text{ if } q_{lo} \leq q < q_{hi};$$

$$p = 1, \text{ if } q \geq q_{hi}.$$

Here, q_{lo} and q_{hi} corresponds to the low and high thresholds of queue occupancy. The maximum parking probability is p_{max} .

The ECN markings events will contribute to the calculation of an equivalent delay x_n at the receiver. No changes are required at the sender.

4.3 PCN marking

As a more advanced feature, we also envision network nodes which support PCN marking based on virtual queues. In such a case, the marking probability of the ECN bit in the IP packet header is calculated as follows:

* upon packet arrival, meter packet against token bucket (r, b) ;

* update token level b_{tk} ;

* calculate the marking probability as:

$$p = 0, \text{ if } b - b_{tk} < b_{lo};$$

$$p = p_{max} * \frac{b - b_{tk} - b_{lo}}{b_{hi} - b_{lo}}, \text{ if } b_{lo} \leq b - b_{tk} < b_{hi};$$

$$p = 1, \text{ if } b - b_{tk} \geq b_{hi}.$$

Here, the token bucket lower and upper limits are denoted by b_{lo} and b_{hi} , respectively. The parameter b indicates the size of the token bucket. The parameter r is chosen as $r = \gamma * C$, where $\gamma < 1$ is the

target utilization ratio and C designates link capacity. The maximum marking probability is p_{\max} .

The ECN markings events will contribute to the calculation of an equivalent delay x_n at the receiver. No changes are required at the sender. The virtual queuing mechanism from the PCN marking algorithm will lead to additional benefits such as zero standing queues.

4.4 Comments and Discussions

In all three flavors described above, the network queue operates with the simple first-in-first-out (FIFO) principle. There is no need to maintain per-flow state. Such a simple design ensures that the system can scale easily with large number of video flows and high link capacity.

The sender behavior stays the same in the presence of all types of congestion signals: delay, loss, ECN marking due to either RED/ECN or PCN algorithms. This unified approach allows a graceful transition of the scheme as the level of congestion in the network shifts dynamically between different regimes.

5. Receiver Behavior

The receiver periodically monitors end-to-end per-packet statistics in terms of delay, loss, and/or ECN marking ratios. It then aggregates all forms of congestion signals in terms of an equivalent delay, and periodically reports back to the sender.

5.1 Monitoring per-packet statistics

Upon receipt of each packet, the receiver calculates one-way delay as the difference between sender and receiver timestamps:

$$x_n = t_{r,n} - t_{s,n}.$$

It also maintains its estimate of baseline delay, d_f , as the minimum value of previously observed x_n 's over a relatively longer period. This assumes that that sending and receiving clocks are either well-synchronized, or have a relatively stable offset. In our implementation, the baseline delay estimation is updated once every 10 minutes.

Correspondingly, the queuing delay experienced by the packet n is estimated as:

$$d_n = x_n - d_f.$$

In addition, the receiver keeps track of both packet loss ratios as p_L via detection of gaps in the packet sequence numbers, and ECN marking ratios as p_M .

5.2 Aggregating congestion signals

The receiver aggregates all three forms of congestion signal in terms of an equivalent delay:

$$x_n = d_n + p_M d_M + p_L d_L, \quad (1)$$

where d_M is a prescribed fictitious delay value associated with ECN markings (e.g., $d_M = 200$ ms), and d_L is a prescribed fictitious delay value associated with packet losses (e.g., $d_L = 1$ second). By introducing a large fictitious delay penalty for ECN marking and packet losses, the proposed scheme leads to low end-to-end actual delays in the presence of such events.

While the value of d_M and d_L are fixed and predetermined in our current design, we also plan to investigate a scheme for automatically tuning these values based on desired bandwidth sharing behavior in the presence of other competing loss-based flows (e.g., loss-based TCP).

It should also be noted that in the absence of loss and marking information, the value of x_n falls back to the observed queuing delay d_n for packet n . Our algorithm operates in purely delay-based mode.

5.3 Sending periodic feedback

Periodically, the receiver sends back the most recent value of x_n in RTCP messages, to aid the sender in its calculation of target rate. It also calculates and sends the derivative of x_n as part of the RTCP message:

$$x'_n = \frac{x_n - x_{(n-k)}}{\text{delta}}. \quad (2)$$

Here, the interval between current and previous RTCP messages is denoted as delta , and the corresponding packet indices are n and $(n-k)$, respectively. Typically, the interval between adjacent RTCP receiver reports is on the order of sub-seconds (e.g., 100ms).

The size of acknowledgement packets are typically on the order of tens of bytes, and are significantly smaller than average video packet sizes. Therefore, the bandwidth overhead of the receiver acknowledgement stream is sufficiently low.

5.4 Discussions on delay metrics

Our current design works with relative OWD as the main indication of congestion. The value of the relative OWD is obtained by maintaining the minimum value of observed OWD over a longer time horizon and subtract that out from the observed absolute OWD value. Such an approach cancels out the fixed clock difference from the sender and receiver clocks, and has been widely adopted by other delay-based congestion control approaches such as LEDBAT [RFC6817]. As discussed in [RFC6817], the time horizon for tracking the minimum OWD needs to be chosen with care: long enough for an opportunity to observe the minimum OWD with zero queuing delay along the path, and sufficiently short so as to timely reflect "true" changes in minimum OWD introduced by route changes and other rare events.

The potential drawback in relying on relative OWD as the congestion signal is that when multiple flows share the same bottleneck, the flow arriving late at the network experiencing a non-empty queue may mistakenly account the standing queuing delay as part of the fixed path propagation delay. This will lead to slightly unfair bandwidth sharing amongst the flows.

Alternatively, one could move the per-packet statistical handling to the sender instead, and use RTT in lieu of OWD, assuming that per-packet ACKs are available. The main drawback of this latter approach, on the other hand, is that the scheme will be confused by congestion in the reverse direction.

Note that the adoption of either delay metric (relative OWD vs. RTT) involves no change in the proposed rate adaptation algorithm at the sender. Therefore, comparing the pros and cons regarding which delay metric to adopt can be kept as an orthogonal direction of investigation.

6. Sender Behavior

Figure 1 provides a more detailed view of the NADA sender. Upon receipt of an RTCP report from the receiver, the NADA sender updates its calculation of the reference rate R_n . It further adjusts both the target rate for the live video encoder R_v and the sending rate R_s over the network based on the updated value of R_n , as well as the size of the rate shaping buffer.

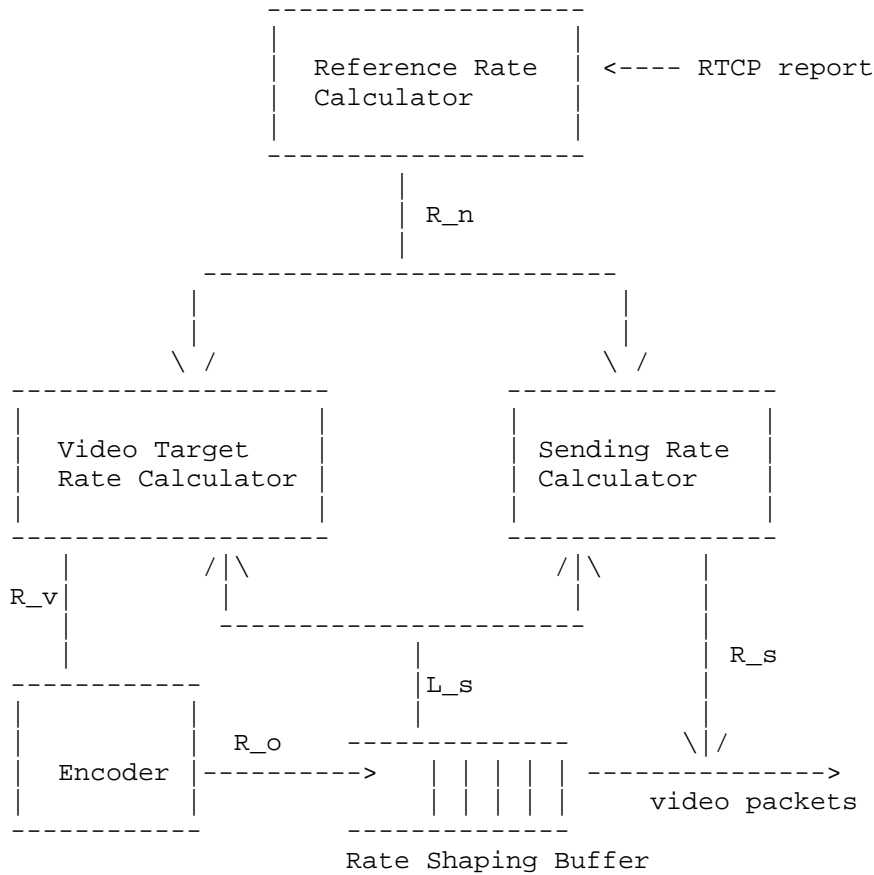


Figure 1 NADA Sender Structure

The following sections describe these modules in further details, and explain how they interact with each other.

6.1 Video encoder rate control

The video encoder rate control procedure has the following characteristics:

- * Rate changes can happen only at large intervals, on the order of seconds.
- * Given a target rate R_o , the encoder output rate may randomly fluctuate around it.

* The encoder output rate is further constrained by video content complexity. The range of the final rate output is $[R_{\min}, R_{\max}]$. Note that it's content-dependent, and may change over time.

Note that operation of the live video encoder is out of the scope of our design for a congestion control scheme in NADA. Instead, its behavior is treated as a black box.

6.2 Rate shaping buffer

A rate shaping buffer is employed to absorb any instantaneous mismatch between encoder rate output R_o and regulated sending rate R_s . The size of the buffer evolves from time $t-\tau$ to time t as:

$$L_s(t) = \max [0, L_s(t-\tau) + (R_o - R_s) \cdot \tau].$$

A large rate shaping buffer contributes to higher end-to-end delay, which may harm the performance of real-time media communications. Therefore, the sender has a strong incentive to constrain the size of the shaping buffer. It can either deplete it faster by increasing the sending rate R_s , or limit its growth by reducing the target rate for the video encoder rate control R_v .

6.3 Reference rate calculator

The sender initializes the reference rate R_n as R_{\min} . Upon receipt of a new receiver RTCP reports containing values of x_n and x'_n , it updates the rate as:

$$R_n \leftarrow R_n + \frac{\kappa \cdot \delta_s}{\tau_o^2} * (\theta - (R_n - R_{\min}) * \hat{x}) \quad (3)$$

where

$$\theta = w * (R_{\max} - R_{\min}) * x_{\text{ref}}. \quad (4)$$

$$\hat{x} = x_n + \eta * \tau_o * x'_n \quad (5)$$

In (3), Δ_s refers to the time interval between current and previous rate updates. Note that Δ_s is the same as the RTCP report interval at the receiver (see Δ from (2)) when the backward path is uncongested.

In (4), R_{\min} and R_{\max} denote the content-dependent rate range the encoder can produce. The weight of priority level is w . The reference congestion signal x_{ref} is chosen so that the maximum rate of R_{\max} can be achieved when $\hat{x} = w \cdot x_{\text{ref}}$.

Proper choice of the scaling parameters η and κ in (3) and (5) can ensure system stability so long as the RTT falls below the upper bound of τ_o . In our design, τ_o is chosen as 500ms.

The final target rate R_n is clipped within the range of $[R_{\min}, R_{\max}]$.

Note that the sender does not need any explicit knowledge of the management scheme inside the network. Rather, it reacts to the aggregation of all forms of congestion indications (delay, loss, and marking) via the composite congestion signals x_n and x'_n from the receiver in a coherent manner.

6.4 Video target rate and sending rate calculator

The target rate for the live video encoder is updated based on both the reference rate R_n and the rate shaping buffer size L_s , as follows:

$$R_v = R_n - \beta_v * \frac{L_s}{\tau_v}. \quad (6)$$

Similarly, the outgoing rate is regulated based on both the reference rate R_n and the rate shaping buffer size L_s , such that:

$$R_s = R_n + \beta_s * \frac{L_s}{\tau_v}. \quad (7)$$

In (6) and (7), the first term indicates the rate calculated from network congestion feedback alone. The second term indicates the influence of the rate shaping buffer. A large rate shaping buffer nudges the encoder target rate slightly below -- and the sending rate slightly above -- the reference rate R_n .

Intuitively, the amount of extra rate offset needed to completely drain

the rate shaping buffer within the same time frame of encoder rate adaptation τ_v is given by L_s/τ_v . The scaling parameters β_v and β_s can be tuned to balance between the competing goals of maintaining a small rate shaping buffer and deviating the system from the reference rate point.

6.5 Start-up behavior

The rate adaptation algorithm specified by (3)--(5) naturally leads to a linear rate increase at start-up, when queuing delay stays at zero in the beginning:

$$R_n \leftarrow R_n + \frac{\kappa \cdot \delta_s}{\tau_o^2} * \theta \quad (8)$$

Given that $\theta = w \cdot (R_{\max} - R_{\min}) \cdot x_{\text{ref}}$, the speed of increase scales with the value of κ , weight of priority w , and dynamic range of the flow ($R_{\max} - R_{\min}$).

In practice, one may desire a more aggressive ramp-up behavior during the start-up period, e.g., by doubling the rate upon the receipt of each new RTCP message which reports on near-zero values of x_n and x'_n .

We note here that design of the start-up behavior can be kept orthogonal to the design of the steady-state rate adaptation behavior. This topic is worthy of further investigation separately.

7. Incremental Deployment

One nice property of proposed design is the consistent video end point behavior irrespective of network node variations. This facilitates gradual, incremental adoption of the scheme.

To start off with, the proposed encoder congestion control mechanism can be implemented without any explicit support from the network, and rely solely on observed one-way delay measurements and packet loss ratios as implicit congestion signals.

When ECN is enabled at the network nodes with RED-based marking, the receiver can fold its observations of ECN markings into the calculation of the equivalent delay. The sender can react to these explicit congestion signals without any modification.

Ultimately, networks equipped with proactive marking based on token bucket level metering can reap the additional benefits of zero standing queues and lower end-to-end delay and work seamlessly with existing senders and receivers.

8. Implementation Status

The proposed NADA scheme has been implemented in the ns-2 simulation platform [ns2]. Extensive simulation evaluations of an earlier version of the draft are documented in [Zhu-PV13]. Evaluation results of current draft over several test cases in [I-D.draft-sarker-rmcat-eval-test] have been presented at the recent IETF meeting [IETF-90].

The scheme has also been implemented in Linux and has been evaluated in a lab setting also described in [IETF-90]. Evaluation results of NADA in single-flow and multi-flow scenarios from this testbed will be disclosed soon.

9. IANA Considerations

There are no actions for IANA.

10. References

10.1 Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.

10.2 Informative References

- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.
- [RFC2309] Braden, B., Clark, D., Crowcroft, J., Davie, B., Deering, S., Estrin, D., Floyd, S., Jacobson, V., Minshall, G., Partridge, C., Peterson, L., Ramakrishnan, K., Shenker, S., Wroclawski, J., and L. Zhang, "Recommendations on Queue Management and Congestion Avoidance in the Internet", RFC 2309, April 1998.
- [RFC6187] S. Shalunov, G. Hazel, J. Iyengar, and M. Kuehlewind, "Low Extra Delay Background Transport (LEDBAT)", RFC 6817, December 2012
- [ns2] "The Network Simulator - ns-2", <http://www.isi.edu/nsnam/ns/>

- [Zhu-PV13] Zhu, X. and Pan, R., "NADA: A Unified Congestion Control Scheme for Low-Latency Interactive Video", in Proc. IEEE International Packet Video Workshop (PV'13). San Jose, CA, USA. December 2013.
- [I-D.draft-sarker-rmcat-eval-test] Sarker, Z., Singh, V., Zhu, X., and Ramalho, M., "Test Cases for Evaluating RMCAT Proposals", draft-sarker-rmcat-eval-test-01 (work in progress), June 2014.
- [IETF-90] Zhu, X. et al., "NADA Update: Algorithm, Implementation, and Test Case Evaluation Results", presented at IETF 90, <https://tools.ietf.org/agenda/90/slides/slides-90-rmcat-6.pdf>

Authors' Addresses

Xiaoqing Zhu
Cisco Systems,
12515 Research Blvd.,
Austin, TX 78759, USA
Email: xiaoqzhu@cisco.com

Rong Pan
Cisco Systems
510 McCarthy Blvd,
Milpitas, CA 95134, USA
Email: ropan@cisco.com

Michael A. Ramalho
6310 Watercrest Way Unit 203
Lakewood Ranch, FL, 34202, USA
Email: mramalho@cisco.com

Sergio Mena de la Cruz
EPFL, Quartier de l'Innovation
Batiment E
Ecublens, Vaud 1015, Switzerland
Email: semen@cisco.com

Charles Ganzhorn
7900 International Drive
International Plaza, Suite 400
Bloomington, MN 55425, USA
Email: cganzhor@cisco.com

Paul E. Jones
7025 Kit Creek Rd.
Research Triangle Park, NC 27709, USA
Email: paulej@packetizer.com

Stefano D'Aronco
EPFL STI IEL LTS4
ELD 220 (Batiment ELD), Station 11
CH-1015 Lausanne, Switzerland
Email: stefano.daronco@epfl.ch

Network Working Group
Internet Draft
Intended Status: Informational
Expires: September 27, 2015

X. Zhu, R. Pan
M. A. Ramalho, S. Mena
C. Ganzhorn, P. E. Jones
Cisco Systems
S. De Aronco
Ecole Polytechnique Federale de Lausanne
March 26, 2015

NADA: A Unified Congestion Control Scheme for Real-Time Media
draft-zhu-rmcat-nada-06

Abstract

Network-Assisted Dynamic Adaptation (NADA) is a novel congestion control scheme for interactive real-time media applications, such as video conferencing. In NADA, the sender regulates its sending rate based on either implicit or explicit congestion signaling in a consistent manner. As one example of explicit signaling, NADA can benefit from explicit congestion notification (ECN) markings from network nodes. It also maintains consistent sender behavior in the absence of explicit signaling by reacting to queuing delay and packet loss.

This document describes the overall system architecture for NADA, as well as recommended behavior at the sender and the receiver.

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) 2012 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. Terminology	3
3. System Model	3
4. NADA Receiver Behavior	4
4.1 Estimation of one-way delay and queuing delay	4
4.2 Estimation of packet loss/marketing ratio	5
4.3 Non-linear warping of delay	6
4.4 Aggregating congestion signals	7
4.5 Estimating receiving rate	7
4.6 Sending periodic feedback	7
4.7 Discussions on delay metrics	8
5. NADA Sender Behavior	9
5.1 Reference rate calculation	10
5.1.1 Accelerated ramp up	10
5.1.2 Gradual rate update	11
5.2 Video encoder rate control	12
5.3 Rate shaping buffer	12
5.4 Adjusting video target rate and sending rate	12
6. Incremental Deployment	13
7. Implementation Status	13
8. IANA Considerations	14
9. References	14
9.1 Normative References	14
9.2 Informative References	14
Appendix A. Network Node Operations	15
A.1 Default behavior of drop tail	16
A.2 ECN marking	16
A.3 PCN marking	16
Authors' Addresses	17

1. Introduction

Interactive real-time media applications introduce a unique set of challenges for congestion control. Unlike TCP, the mechanism used for real-time media needs to adapt quickly to instantaneous bandwidth changes, accommodate fluctuations in the output of video encoder rate control, and cause low queuing delay over the network. An ideal scheme should also make effective use of all types of congestion signals, including packet loss, queuing delay, and explicit congestion notification (ECN) [RFC3168] markings.

Based on the above considerations, this document describes a scheme called network-assisted dynamic adaptation (NADA). The NADA design benefits from explicit congestion control signals (e.g., ECN markings) from the network, yet also operates when only implicit congestion indicators (delay and/or loss) are available. In addition, it supports weighted bandwidth sharing among competing video flows.

This documentation describes the overall system architecture, recommended designs at the sender and receiver, as well as expected network node operations. The signaling mechanism consists of standard RTP timestamp [RFC3550] and standard RTCP feedback reports.

2. Terminology

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

3. System Model

The overall system consists of the following elements:

- * Source media stream, in the form of consecutive raw video frames and audio samples;
- * Media encoder with rate control capabilities. It takes the source media stream and encodes it to an RTP stream at a target bit rate R_v . Note that the actual output rate from the encoder R_o may fluctuate around the target R_v . Also, the encoder can only change its rate at rather coarse time intervals, e.g., once every 0.5 seconds.
- * RTP sender, responsible for calculating the target bit rate R_n based on network congestion indicators (delay, loss, or ECN marking reports from the receiver), for updating the video encoder with a new target rate R_v , and for regulating the

actual sending rate R_s accordingly. A rate shaping buffer is employed to absorb the instantaneous difference between video encoder output rate R_v and sending rate R_s . The buffer size L_s , together with R_n , influences the calculation of actual sending rate R_s and video encoder target rate R_v . The RTP sender also generates RTP timestamp in outgoing packets.

* RTP receiver, responsible for measuring and estimating end-to-end delay based on sender RTP timestamp. In the presence of packet loss and ECN markings, it keeps track of packet loss and ECN marking ratios. It calculates the equivalent delay x_n that accounts for queuing delay, ECN marking, and packet loss, as well as the derivative (i.e., rate of change) of this congestion signal as x'_n . The receiver feeds both pieces of information (x_n and x'_n) back to the sender via periodic RTCP reports.

* Network node, with several modes of operation. The system can work with the default behavior of a simple drop tail queue. It can also benefit from advanced AQM features such as RED-based ECN marking, and PCN marking using a token bucket algorithm. Note that network node operation is out of scope for the design of NADA.

In the following, we will elaborate on the respective operations at the NADA receiver and sender.

4. NADA Receiver Behavior

The receiver continuously monitors end-to-end per-packet statistics in terms of delay, loss, and/or ECN marking ratios. It then aggregates all forms of congestion indicators into the form of an equivalent delay and periodically reports this back to the sender. In addition, the receiver tracks the receiving rate of the flow and includes that in the feedback message.

4.1 Estimation of one-way delay and queuing delay

The delay estimation process in NADA follows a similar approach as in earlier delay-based congestion control schemes, such as LEDBAT [RFC6817]. NADA estimates the forward delay as having a constant base delay component plus a time varying queuing delay component. The base delay is estimated as the minimum value of one-way delay observed over a relatively long period (e.g., tens of minutes), whereas the individual queuing delay value is taken to be the difference between one-way delay and base delay.

In mathematical terms, for packet n arriving at the receiver, one-way delay is calculated as:

$$x_n = t_{r,n} - t_{s,n},$$

where $t_{s,n}$ and $t_{r,n}$ are sender and receiver timestamps, respectively. A real-world implementation should also properly handle practical issues such as wrap-around in the value of x_n , which are omitted from the above simple expression for brevity.

The base delay, d_f , is estimated as the minimum value of previously observed x_n 's over a relatively long period. This assumes that the drift between sending and receiving clocks remains bounded by a small value.

Correspondingly, the queuing delay experienced by the packet n is estimated as:

$$d_n = x_n - d_f.$$

The individual sample values of queuing delay should be further filtered against various non-congestion-induced noise, such as spikes due to processing "hiccup" at the network nodes. We denote the resulting queuing delay value as $d_{\hat{n}}$.

Our current implementation employs a simple 5-point median filter over per-packet queuing delay estimates, followed by an exponential smoothing filter. We have found such relatively simple treatment to suffice in guarding against processing delay outliers observed in wired connections. For wireless connections with a higher packet delay variation (PDV), more sophisticated techniques on de-noising, outlier rejection, and trend analysis may be needed.

Like other delay-based congestion control schemes, performance of NADA depends on the accuracy of its delay measurement and estimation module. Appendix A in [RFC6817] provides an extensive discussion on this aspect.

4.2 Estimation of packet loss/marketing ratio

The receiver detects packet losses via gaps in the RTP sequence numbers of received packets. It then calculates instantaneous packet loss ratio as the ratio between the number of missing packets over the number of total transmitted packets in the given time window (e.g., during the most recent 500ms). This instantaneous value is passed over an exponential smoothing filter, and the filtered result is reported back to the sender as the observed packet loss ratio p_L .

We note that more sophisticated methods in packet loss ratio calculation, such as that adopted by TFRC [Floyd-CCR00], will likely be beneficial. These alternatives are currently under investigation.

Estimation of packet marking ratio p_M , when ECN is enabled at bottleneck network nodes along the path, will follow the same procedure as above. Here it is assumed that ECN marking information at the IP header are somehow passed along to the transport layer by the receiving endpoint.

4.3 Non-linear warping of delay

In order for a delay-based flow to hold its ground and sustain a reasonable share of bandwidth in the presence of a loss-based flow (e.g., loss-based TCP), it is important to distinguish between different levels of observed queuing delay. For instance, a moderate queuing delay value below 100ms is likely self-inflicted or induced by other delay-based flows, whereas a high queuing delay value of several hundreds of milliseconds may indicate the presence of a loss-based flow that does not refrain from increased delay.

Inspired by the delay-adaptive congestion window backoff policy in [Budzisz-TON11] -- the work by itself is a window-based congestion control scheme with fair coexistence with TCP -- we devise the following non-linear warping of estimated queuing delay value:

$$\begin{aligned}
 d_{\text{tilde}_n} &= (d_{\text{hat}_n}), & \text{if } d_{\text{hat}_n} < d_{\text{th}}; \\
 d_{\text{tilde}_n} &= d_{\text{th}} \frac{(d_{\text{max}} - d_{\text{hat}_n})^4}{(d_{\text{max}} - d_{\text{th}})^4}, & \text{if } d_{\text{th}} < d_{\text{hat}_n} < d_{\text{max}}; \\
 d_{\text{tilde}_n} &= 0, & \text{if } d_{\text{hat}_n} > d_{\text{max}}.
 \end{aligned}$$

Here, the queuing delay value is unchanged when it is below the first threshold d_{th} ; it is discounted following a non-linear curve when its value falls between d_{th} and d_{max} ; above d_{max} , the high queuing delay value no longer counts toward congestion control.

When queuing delay is in the range $(0, d_{\text{th}})$, NADA operates in pure delay-based mode if no losses/markings are present. When queuing delay exceeds d_{max} , NADA reacts to loss/markings only. In between d_{th} and d_{max} , the sending rate will converge and stabilize at an operating point with a fairly high queuing delay and non-zero packet loss ratio.

In our current implementation d_{th} is chosen as 50ms and d_{max} is chosen as 400ms. The impact of the choice of d_{th} and d_{max} will be investigated in future work.

4.4 Aggregating congestion signals

The receiver aggregates all three forms of congestion signal in terms of an equivalent delay:

$$x_n = d_{\tilde{n}} + p_M d_M + p_L d_L, \quad (1)$$

where d_M is a prescribed fictitious delay value associated with ECN markings (e.g., $d_M = 200$ ms), and d_L is a prescribed fictitious delay value associated with packet losses (e.g., $d_L = 1$ second). By introducing a large fictitious delay penalty for ECN marking and packet loss, the proposed scheme leads to low end-to-end actual delay in the presence of such events.

While the value of d_M and d_L are fixed and predetermined in the current design, a scheme for automatically tuning these values based on desired bandwidth sharing behavior in the presence of other competing loss-based flows (e.g., loss-based TCP) is being studied.

In the absence of ECN marking from the network, the value of x_n falls back to the observed queuing delay d_n for packet n when queuing delay is low and no packets are lost over a lightly congested path. In that case the algorithm operates in purely delay-based mode.

4.5 Estimating receiving rate

Estimation of receiving rate of the flow is fairly straightforward. NADA maintains a recent observation window of 500ms, and simply divides the total size of packets arriving during that window over the time span. The receiving rate is denoted as R_r .

4.6 Sending periodic feedback

Periodically, the receiver feeds back a tuple of the most recent values of $\langle d_{\hat{n}}, x_n, x'_n, R_r \rangle$ in RTCP feedback messages to aid the sender in its calculation of target rate. The queuing delay value $d_{\hat{n}}$ is included along with the composite congestion signal x_n so that the sender can decide whether the network is truly underutilized (see Sec. 6.1.1 Accelerated ramp-up).

The value of x'_n corresponds to the derivative (i.e., rate of change) of the composite congestion signal:

$$x'_n = \frac{x_n - x_{(n-k)}}{\text{delta}}, \quad (2)$$

where the interval between consecutive RTCP feedback messages is denoted as δ . The packet indices corresponding to the current and previous feedback are n and $(n-k)$, respectively.

The choice of target feedback interval needs to strike the right balance between timely feedback and low RTCP feedback message counts. Through simulation studies and frequency-domain analysis, it was determined that a feedback interval below 250ms will not break up the feedback control loop of the NADA congestion control algorithm. Thus, it is recommended to use a target feedback interval of 100ms. This will result in a feedback bandwidth of 16Kbps with 200 bytes per feedback message, less than 0.1% overhead for a 1Mbps flow.

4.7 Discussions on delay metrics

The current design works with relative one-way-delay (OWD) as the main indication of congestion. The value of the relative OWD is obtained by maintaining the minimum value of observed OWD over a relatively long time horizon and subtract that out from the observed absolute OWD value. Such an approach cancels out the fixed difference between the sender and receiver clocks. It has been widely adopted by other delay-based congestion control approaches such as LEDBAT [RFC6817]. As discussed in [RFC6817], the time horizon for tracking the minimum OWD needs to be chosen with care: it must be long enough for an opportunity to observe the minimum OWD with zero queuing delay along the path, and sufficiently short so as to timely reflect "true" changes in minimum OWD introduced by route changes and other rare events.

The potential drawback in relying on relative OWD as the congestion signal is that when multiple flows share the same bottleneck, the flow arriving late at the network experiencing a non-empty queue may mistakenly consider the standing queuing delay as part of the fixed path propagation delay. This will lead to slightly unfair bandwidth sharing among the flows.

Alternatively, one could move the per-packet statistical handling to the sender instead and use RTT in lieu of OWD, assuming that per-packet ACKs are available. The main drawback of this latter approach is that the scheme will be confused by congestion in the reverse direction.

Note that the choice of either delay metric (relative OWD vs. RTT) involves no change in the proposed rate adaptation algorithm at the sender. Therefore, comparing the pros and cons regarding which delay metric to adopt can be kept as an orthogonal direction of investigation.

5. NADA Sender Behavior

Figure 1 provides a detailed view of the NADA sender. Upon receipt of an RTCP report from the receiver, the NADA sender updates its calculation of the reference rate R_n . It further adjusts both the target rate for the live video encoder R_v and the sending rate R_s over the network based on the updated value of R_n , as well as the size of the rate shaping buffer.

In the following, we describe these modules in further detail, and explain how they interact with each other.

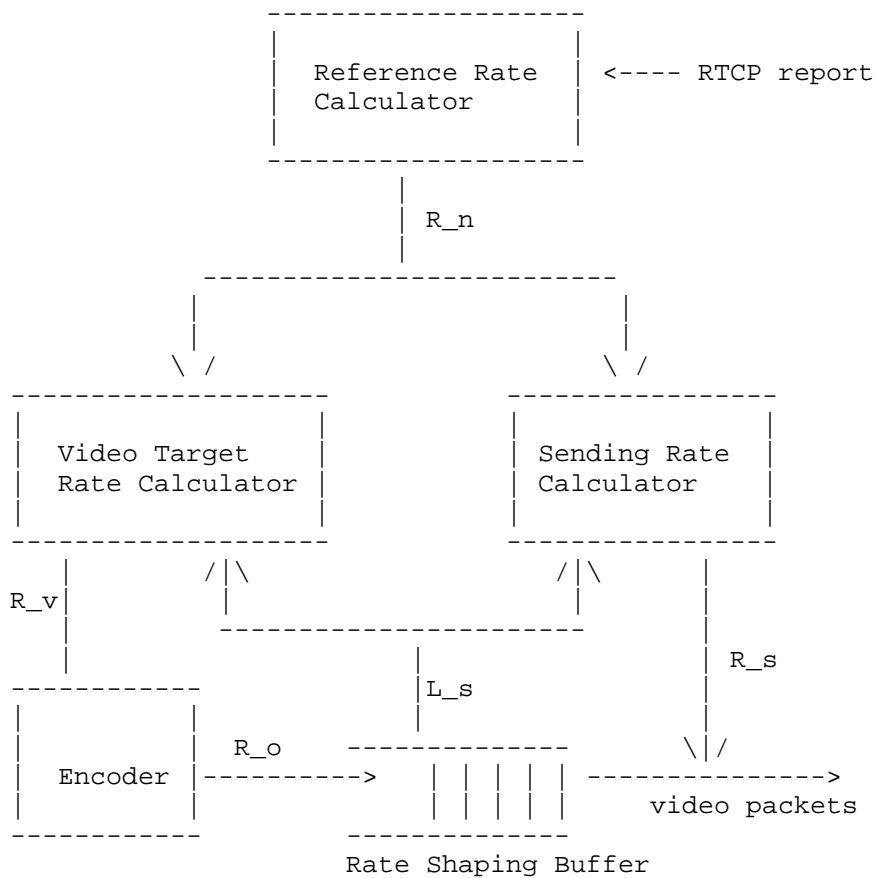


Figure 1 NADA Sender Structure

5.1 Reference rate calculation

The sender initializes the reference rate R_n as R_{\min} by default, or to a value specified by the upper-layer application. [Editor's note: should proper choice of starting rate value be within the scope of the CC solution?]

The reference rate R_n is calculated based on receiver feedback information regarding queuing delay $d_{\tilde{n}}$, composite congestion signal x_n , its derivative x'_n , as well as the receiving rate R_r . The sender switches between two modes of operation:

- * Accelerated ramp up: if the reported queuing delay is close to zero and both values of x_n and x'_n are close to zero, indicating empty queues along the path of the flow and, consequently, underutilized network bandwidth; or

- * Gradual rate update: in all other conditions, whereby the receiver reports on a standing or increasing/decreasing queue and/or composite congestion signal.

5.1.1 Accelerated ramp up

In the absence of a non-zero congestion signal to guide the sending rate calculation, the sender needs to ramp up its estimated bandwidth as quickly as possible without introducing excessive queuing delay. Ideally the flow should inflict no more than T_{th} milliseconds of queuing delay at the bottleneck during the ramp-up process. A typical value of T_{th} is 50ms.

Note that the sender will be aware of any queuing delay introduced by its rate increase after at least one round-trip time. In addition, the bottleneck bandwidth C is greater than or equal to the receive rate R_r reported from the most recent "no congestion" feedback message. The rate R_n is updated as follows:

$$\gamma = \min \left[\gamma_0, \frac{T_{th}}{RTT_0 + \delta_0} \right] \quad (3)$$

$$R_n = (1 + \gamma) R_r \quad (4)$$

In (3) and (4), the multiplier γ for rate increase is upper-bounded by a fixed ratio γ_0 (e.g., 20%), as well as a ratio which depends

on T_{th} , base RTT as measured during the non-congested phase, and target ACK interval δ_0 . The rationale behind this is that the rate increase multiplier should decrease with the delay in the feedback control loop, and that $RTT_0 + \delta_0$ provides a worst-case estimate of feedback control delay when the network is not congested.

5.1.2. Gradual rate update

When the receiver reports indicate a standing congestion level, NADA operates in gradual update mode, and calculates its reference rate as:

$$R_n \leftarrow R_n + \frac{\kappa * \delta_s}{\tau_o^2} * (\theta - (R_n - R_{min}) * x_{hat}) \quad (5)$$

where

$$\theta = w * (R_{max} - R_{min}) * x_{ref}. \quad (6)$$

$$x_{hat} = x_n + \eta * \tau_o * x'_n \quad (7)$$

In (5), δ_s refers to the time interval between current and previous rate updates. Note that δ_s is the same as the RTCP report interval at the receiver (see δ from (2)) when the backward path is uncongested.

In (6), R_{min} and R_{max} denote the content-dependent rate range the encoder can produce. The weighting factor reflecting a flow's priority is w . The reference congestion signal x_{ref} is chosen so that the maximum rate of R_{max} can be achieved when $x_{hat} = w * x_{ref}$.

Proper choice of the scaling parameters η and κ in (5) and (7) can ensure system stability so long as the RTT falls below the upper bound of τ_o . The recommended default value of τ_o is chosen as 500ms.

For both modes of operations, the final reference rate R_n is clipped within the range of $[R_{min}, R_{max}]$. Note also that the sender does not need any explicit knowledge of the management scheme inside the network. Rather, it reacts to the aggregation of all forms of congestion indications (delay, loss, and explicit markings) via the composite congestion signals x_n and x'_n from the receiver in a coherent manner.

5.2 Video encoder rate control

The video encoder rate control procedure has the following characteristics:

- * Rate changes can happen only at large intervals, on the order of seconds.
- * The encoder output rate may fluctuate around the target rate R_v .
- * The encoder output rate is further constrained by video content complexity. The range of the final rate output is $[R_{min}, R_{max}]$. Note that it is content-dependent and may vary over time.

The operation of the live video encoder is out of the scope of the design for the congestion control scheme in NADA. Instead, its behavior is treated as a black box.

5.3 Rate shaping buffer

A rate shaping buffer is employed to absorb any instantaneous mismatch between encoder rate output R_o and regulated sending rate R_s . The size of the buffer evolves from time $t-\tau$ to time t as:

$$L_s(t) = \max [0, L_s(t-\tau) + (R_o - R_s) \cdot \tau].$$

A large rate shaping buffer contributes to higher end-to-end delay, which may harm the performance of real-time media communications. Therefore, the sender has a strong incentive to constrain the size of the shaping buffer. It can either deplete it faster by increasing the sending rate R_s , or limit its growth by reducing the target rate for the video encoder rate control R_v .

5.4 Adjusting video target rate and sending rate

The target rate for the live video encoder is updated based on both the reference rate R_n and the rate shaping buffer size L_s , as follows:

$$R_v = R_n - \beta_v * \frac{L_s}{\tau_v}. \quad (8)$$

Similarly, the outgoing rate is regulated based on both the reference rate R_n and the rate shaping buffer size L_s , such that:

$$R_s = R_n + \beta_s * \frac{L_s}{\tau_v}. \quad (9)$$

In (8) and (9), the first term indicates the rate calculated from network congestion feedback alone. The second term indicates the influence of the rate shaping buffer. A large rate shaping buffer nudges the encoder target rate slightly below -- and the sending rate slightly above -- the reference rate R_n .

Intuitively, the amount of extra rate offset needed to completely drain the rate shaping buffer within the same time frame of encoder rate adaptation τ_v is given by L_s/τ_v . The scaling parameters β_v and β_s can be tuned to balance between the competing goals of maintaining a small rate shaping buffer and deviating the system from the reference rate point.

6. Incremental Deployment

One nice property of NADA is the consistent video endpoint behavior irrespective of network node variations. This facilitates gradual, incremental adoption of the scheme.

To start off with, the encoder congestion control mechanism can be implemented without any explicit support from the network, and relies solely on observed one-way delay measurements and packet loss ratios as implicit congestion signals.

When ECN is enabled at the network nodes with RED-based marking, the receiver can fold its observations of ECN markings into the calculation of the equivalent delay. The sender can react to these explicit congestion signals without any modification.

Ultimately, networks equipped with proactive marking based on token bucket level metering can reap the additional benefits of zero standing queues and lower end-to-end delay and work seamlessly with existing senders and receivers.

7. Implementation Status

The NADA scheme has been implemented in the ns-2 simulation platform [ns2]. Extensive simulation evaluations of an earlier version of the draft are documented in [Zhu-PV13]. Evaluation results of the current draft over several test cases in [I-D.draft-sarker-rmcat-eval-test] have been presented at recent IETF meetings [IETF-90][IETF-91].

The scheme has also been implemented and evaluated in a lab setting as described in [IETF-90]. Preliminary evaluation results of NADA in single-flow and multi-flow scenarios have been presented in [IETF-91].

8. IANA Considerations

There are no actions for IANA.

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.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.

9.2 Informative References

- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.
- [RFC2309] Braden, B., Clark, D., Crowcroft, J., Davie, B., Deering, S., Estrin, D., Floyd, S., Jacobson, V., Minshall, G., Partridge, C., Peterson, L., Ramakrishnan, K., Shenker, S., Wroclawski, J., and L. Zhang, "Recommendations on Queue Management and Congestion Avoidance in the Internet", RFC 2309, April 1998.
- [RFC6817] Shalunov, S., Hazel, G., Iyengar, J., and Kuehlewind, M., "Low Extra Delay Background Transport (LEDBAT)", RFC 6817, December 2012
- [Floyd-CCR00] Floyd, S., Handley, M., Padhye, J., and Widmer, J., "Equation-based Congestion Control for Unicast Applications", ACM SIGCOMM Computer Communications Review, vol. 30. no. 4., pp. 43-56, October 2000.
- [Budzisz-TON11] Budzisz, L. et al., "On the Fair Coexistence of Loss- and Delay-Based TCP", IEEE/ACM Transactions on Networking, vol. 19, no. 6, pp. 1811-1824, December 2011.

- [ns2] "The Network Simulator - ns-2", <http://www.isi.edu/nsnam/ns/>
- [Zhu-PV13] Zhu, X. and Pan, R., "NADA: A Unified Congestion Control Scheme for Low-Latency Interactive Video", in Proc. IEEE International Packet Video Workshop (PV'13). San Jose, CA, USA. December 2013.
- [I-D.draft-sarker-rmcat-eval-test] Sarker, Z., Singh, V., Zhu, X., and Ramalho, M., "Test Cases for Evaluating RMCAT Proposals", draft-sarker-rmcat-eval-test-01 (work in progress), June 2014.
- [IETF-90] Zhu, X. et al., "NADA Update: Algorithm, Implementation, and Test Case Evaluation Results", presented at IETF 90, <https://tools.ietf.org/agenda/90/slides/slides-90-rmcat-6.pdf>
- [IETF-91] Zhu, X. et al., "NADA Algorithm Update and Test Case Evaluations", presented at IETF 91 Interim, <https://datatracker.ietf.org/meeting/91/agenda/rmcat/>

Appendix A. Network Node Operations

NADA can work with different network queue management schemes and does not assume any specific network node operation. As an example, this appendix describes three variations of queue management behavior at the network node, leading to either implicit or explicit congestion signals.

In all three flavors described below, the network queue operates with the simple first-in-first-out (FIFO) principle. There is no need to maintain per-flow state. Such a simple design ensures that the system can scale easily with a large number of video flows and high link capacity.

NADA sender behavior stays the same in the presence of all types of congestion indicators: delay, loss, ECN marking due to either RED/ECN or PCN algorithms. This unified approach allows a graceful transition of the scheme as the network shifts dynamically between light and heavy congestion levels.

A.1 Default behavior of drop tail

In a conventional network with drop tail or RED queues, congestion is inferred from the estimation of end-to-end delay and/or packet loss. Packet drops at the queue are detected at the receiver, and contributes to the calculation of the equivalent delay x_n . No special action is required at network node.

A.2 ECN marking

In this mode, the network node randomly marks the ECN field in the IP packet header following the Random Early Detection (RED) algorithm [RFC2309]. Calculation of the marking probability involves the following steps:

* upon packet arrival, update smoothed queue size q_{avg} as:

$$q_{avg} = \alpha * q + (1 - \alpha) * q_{avg}.$$

The smoothing parameter α is a value between 0 and 1. A value of $\alpha=1$ corresponds to performing no smoothing at all.

* calculate marking probability p as:

$$p = 0, \text{ if } q < q_{lo};$$

$$p = p_{max} * \frac{q_{avg} - q_{lo}}{q_{hi} - q_{lo}}, \text{ if } q_{lo} \leq q < q_{hi};$$

$$p = 1, \text{ if } q \geq q_{hi}.$$

Here, q_{lo} and q_{hi} corresponds to the low and high thresholds of queue occupancy. The maximum marking probability is p_{max} .

The ECN markings events will contribute to the calculation of an equivalent delay x_n at the receiver. No changes are required at the sender.

A.3 PCN marking

As a more advanced feature, we also envisage network nodes which support PCN marking based on virtual queues. In such a case, the marking probability of the ECN bit in the IP packet header is calculated as follows:

- * upon packet arrival, meter packet against token bucket (r,b);
- * update token level b_{tk};
- * calculate the marking probability as:

$p = 0$, if $b - b_{tk} < b_{lo}$;

$p = p_{max} * \frac{b - b_{tk} - b_{lo}}{b_{hi} - b_{lo}}$, if $b_{lo} \leq b - b_{tk} < b_{hi}$;

$p = 1$, if $b - b_{tk} \geq b_{hi}$.

Here, the token bucket lower and upper limits are denoted by b_{lo} and b_{hi} , respectively. The parameter b indicates the size of the token bucket. The parameter r is chosen as $r = \gamma * C$, where $\gamma < 1$ is the target utilization ratio and C designates link capacity. The maximum marking probability is p_{max} .

The ECN markings events will contribute to the calculation of an equivalent delay x_n at the receiver. No changes are required at the sender. The virtual queuing mechanism from the PCN marking algorithm will lead to additional benefits such as zero standing queues.

Authors' Addresses

Xiaoqing Zhu
Cisco Systems,
12515 Research Blvd.,
Austin, TX 78759, USA
Email: xiaoqzhu@cisco.com

Rong Pan
Cisco Systems
510 McCarthy Blvd,
Milpitas, CA 95134, USA
Email: ropan@cisco.com

Michael A. Ramalho
6310 Watercrest Way Unit 203
Lakewood Ranch, FL, 34202, USA
Email: mramalho@cisco.com

Sergio Mena de la Cruz
Cisco Systems
EPFL, Quartier de l'Innovation, Batiment E
Ecublens, Vaud 1015, Switzerland
Email: semena@cisco.com

Charles Ganzhorn
7900 International Drive
International Plaza, Suite 400
Bloomington, MN 55425, USA
Email: charles.ganzhorn@gmail.com

Paul E. Jones
7025 Kit Creek Rd.
Research Triangle Park, NC 27709, USA
Email: paulej@packetizer.com

Stefano D'Aronco
EPFL STI IEL LTS4
ELD 220 (Batiment ELD), Station 11
CH-1015 Lausanne, Switzerland
Email: stefano.daronco@epfl.ch

Network Working Group
Internet-Draft
Intended status: Informational
Expires: April 29, 2015

X. Zhu
S. Mena
Cisco Systems
Z. Sarker
Ericsson AB
October 26, 2014

Modeling Video Traffic Sources for RMCAT Evaluations
draft-zhu-rmcat-video-traffic-source-00

Abstract

This document describes two reference video traffic source models for evaluating RMCAT candidate algorithms. The first model statistically characterizes the behavior of a live video encoder in response to changing requests on target video rate. The second model is trace-driven, and emulates the encoder output by scaling the pre-encoded video frame sizes from a widely used video test sequence. Both models are designed to strike a balance between simplicity, repeatability, and authenticity in modeling the interactions between a video traffic source and the congestion control module.

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 29, 2015.

Copyright Notice

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

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of

publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

- 1. Introduction 2
- 2. Terminology 3
- 3. Desired Behavior of Synthetic Video Traffic Model 3
- 4. Interactions Between Synthetic Video Traffic Source and Congestion Control 4
- 5. A Statistical Reference Model 6
 - 5.1. Time-damped response to target rate update 6
 - 5.2. Temporary burst/oscillation during transient 6
 - 5.3. Output rate fluctuation at steady state 7
 - 5.4. Rate range limit imposed by video content 7
- 6. A Trace-Based Model 7
 - 6.1. Choosing the video sequence and generating the traces 8
 - 6.2. Using the traces in the syntetic codec 9
 - 6.2.1. Main algorithm 9
 - 6.2.2. Notes to the main algorithm 11
 - 6.3. Varying frame rate and resolution 11
- 7. Combining The Two Models 12
- 8. IANA Considerations 12
- 9. References 12
 - 9.1. Normative References 12
 - 9.2. Informative References 13
- Authors' Addresses 14

1. Introduction

When evaluating candidate congestion control algorithms designed for real-time interactive media -- as chartered by the RMCAT Working Group -- it is important to account for the characteristics of traffic patterns generated from a live video encoder. Unlike synthetic traffic sources that can conform perfectly to the rate changing requests from the congestion control module, a live video encoder can be sluggish in reacting to such changes. Output rate of a live video encoder also typically deviates from the target rate due to uncertainties in the encoder rate control process. Consequently, end-to-end delay and loss performance of a RMCAT flow can be further impacted by such rate variations introduced by the live encoder.

On the other hand, evaluation results of a candidate RMCAT algorithm should mostly reflect performance of the congestion control module,

and be somewhat decoupled from the choice of a specific video codec. It is also desirable that the evaluation tests are repeatable, and be easily duplicated across different candidate algorithms.

One way to strike a balance between the above considerations is to evaluate RMCAT algorithms using a video traffic source model that is slightly more sophisticated than perfectly conforming CBR traffic, but rather captures the key characteristics of the behavior of a live video encoder. To this purpose, this draft presents two reference models based on two different approaches: statistical modelling or trace driven, respectively.

2. Terminology

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

The terminology defined in RTP [RFC3550], RTP Profile for Audio and Video Conferences with Minimal Control [RFC3551], RTCP Extended Report (XR) [RFC3611], Extended RTP Profile for RTCP-based Feedback (RTP/AVPF) [RFC4585], Support for Reduced-Size RTCP [RFC5506], and RTP Circuit Breaker Algorithm [I-D.ietf-avtcore-rtp-circuit-breakers] apply.

3. Desired Behavior of Synthetic Video Traffic Model

A live video encoder employs encoder rate control to meet a target rate by varying its encoding parameters, such as quantization step size, frame rate, and picture resolution, based on its estimate of the video content (e.g., motion and scene complexity). In practice, however, several factors prevent the output video rate from perfectly conforming to the input target rate.

Due to uncertainties in the captured video scene, the output rate typically deviates from the specified target. In the presence of a significant change in target rate, it sometimes takes several frames before the encoder output rate converges to the new target. Finally, while most of the frames in a live session are encoded in predictive mode, the encoder can occasionally generate a large intra-coded frame (or a frame partially containing intra-coded blocks) in an attempt to recover from losses or re-sync with the receiver, or during the transient period of responding to target rate changes.

Hence, a synthetic video source should have -

- o ability to change bitrate. This includes ability to change the framerate and/or the resolution, skip frames when required.

- o ability to fluctuate around the target bitrate set by the congestion control module.
- o ability to add delay in convergence to the target bitrate.
- o ability to produce Intra-coded frames on demand.

While there exists many different approaches in developing a synthetic video traffic model, it is desirable that the outcome follows a few common characteristics, as outlined below.

- * **Low Computational Complexity:** The model should be computationally lightweight, otherwise it defeats the whole purpose of serving as a substitute for a live video encoder.
- * **Temporal Pattern Similarity:** The individual traffic trace instances generated by the model should mimic the temporal pattern of those from a real video encoder.
- * **Statistical Accuracy:** The synthetic traffic should match the outcome of the real video encoder in terms of statistical characteristics, such as the mean, variance, peak, and autocorrelation coefficients of the bitrate. It is also important that the statistical resemblance should hold across different time scales, ranging from tens of milliseconds to sub-seconds.
- * **Wide Range of Coverage:** The model should be easily configurable to cover a wide range of codec behaviors (e.g., with either fast or slow reaction time in live encoder rate control) and video content variations (e.g, ranging from high-motion to low-motion).

These distinct behavior features can be characterized via simple statistical models, or a trace-driven approach. In the next three sections, we present an example of each.

4. Interactions Between Synthetic Video Traffic Source and Congestion Control

Figure 1 illustrates how the synthetic video traffic source module interacts with the congestion control module at the sender. Both reference models described later in Sections 5. and 6. assume the same set of interactions between encoder rate control and congestion control, as well as the underlying packet transport module.

We model the synthetic video encoder to take in raw video frames captured by the camera along with a set of requests from the

congestion control module. It then dynamically generates a sequence of encoded video frames with varying size and interval. These encoded frames are segmented and packetized into RTP packets by the RTP stack, and encapsulated over UDP/IP before they are transmitted to the network interface. Upon the receipt of an updated RTCP report from the receiver, the congestion control module may further revise its request to the synthetic video encoder, which in turn updates the size and interval of encoded video frames at its output.

In our model, the key notion of "congestion control requests" --- marked as (a) in the figure --- comprises several options:

- o Target rate $R_v(t)$: requested at time t from the congestion control module to the encoder. Depending on the congestion control algorithm in use, the update requests can either be periodic (e.g., once per 1 second), or on-demand (e.g., only when drastic bandwidth change over the network is observed).

- o Target frame rate $FPS(t)$: the instantaneous frame rate measured in frames-per-second at time t . This depends on the native camera capture frame rate as well as the target/preferred frame rate configured by the application or user.

- o Instant frame skipping: the request from the congestion control module to skip the encoding of one or several captured video frames, typically when a drastic decrease in available network bandwidth is detected.

- o On-demand generation of intra (I) frame: the request to encode another I frame to avoid further error propagation at the receiver, if severe packet losses are observed. Strictly speaking, this request should come from the error control module, not the congestion control module in the sender.

Optionally, the synthetic video encoder can inform the congestion control module of the dynamic range of its output rate for the current video contents: $[R_{min}, R_{max}]$. Here, R_{min} and R_{max} are meant to capture the dynamic rate range the encoder is capable of outputting. This typically depends on the video content complexity and/or display type (e.g., higher R_{max} for video contents with higher motion complexity, or for displays of higher resolution). Therefore, these values will not change with R_v , but may change over time if the content is changing.

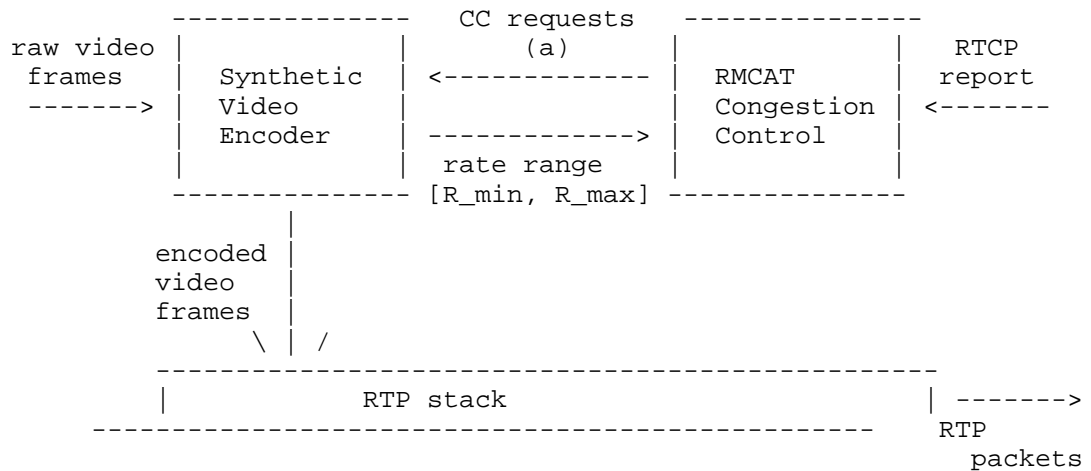


Figure 1: Interaction between synthetic video encoder, congestion control, and packet transport module.

5. A Statistical Reference Model

In this section, we describe one simple statistical model of the live video encoder traffic source. A more complete survey of popular methods can be found in [Tanwir2013].

5.1. Time-damped response to target rate update

While the congestion control module can update its target rate request $R_v(t)$ at any time, our model dictates that the encoder will only react to such changes after τ_v seconds from a previous rate transition. In other words, when the encoder has reacted to a rate change request at time t , it will simply ignore all subsequent rate change requests until time $t + \tau_v$.

5.2. Temporary burst/oscillation during transient

The output rate R_o during the period $[t, t + \tau_v]$ is considered to be in transient. Based on observations from video encoder output data, we model the transient behavior of an encoder upon reacting to a new target rate request in the form of largely varying output sizes. It is assumed that the overall average output rate R_o during this period matches the target rate R_v . Consequently, the occasional burst of large frames are followed by smaller-than average encoded frames.

This temporary burst is characterized by two parameters:

- o burst duration K_d : number frames in the burst event; and
- o burst size K_r : ratio of a burst frame and average frame size at steady state.

It can be noted that these burst parameters can also be used to mimic the inserion of a large on-demand I frame in the presence of severe packet losses. The values of K_d and K_r are fitted to reflect the typical ratio between I and P frames for a given video content.

5.3. Output rate fluctuation at steady state

We model output rate R_o as randomly fluctuating around the target rate R_v after convergence. There are two variants in modeling the random fluctuation $R_e = R_o - R_v$:

- o As normal distribution: with a mean of zero and a standard deviation $SIGMA$ specified in terms of percentage of the target rate. A typical value of $SIGMA$ is 10 percent of target rate.
- o As uniform distribution bounded between $-DELTA$ and $DELTA$. A typical value of $DELTA$ is 10 percent of target rate.

The distribution type (normal or uniform) and model parameters ($SIGMA$ or $DELTA$) can be learned from data samples gathered from a live encoder output.

5.4. Rate range limit imposed by video content

The output rate R_o is further clipped within the dynamic range $[R_{min}, R_{max}]$, which in reality are dictated by scene and motion complexity of the captured video content. In our model, these parameters are specified by the user.

6. A Trace-Based Model

We now present the second approach to model a video traffic source. This approach is based on running an actual live video encoder offline on a set of chosen raw video sequences and using the encoder's output traces for constructing a synthetic live encoder. With this approach, the recorded video traces naturally exhibit temporal fluctuations around a given target rate request $R_v(t)$ from the congestion control module.

The following list summarizes this approach's main steps:

- 1) Choose one or more representative raw video sequences.

- 2) Using an actual live video encoder, encode the sequences at various bitrates. Keep just the sequences of frame sizes for each bitrate.
- 3) Construct a data structure that contains the output of the previous step. The data structure should allow for easy bitrate lookup.
- 4) Upon a target bitrate request $R_v(t)$ from the controller, look up the closest bitrates among those previously stored. Use the frame size sequences stored for those bitrates to approximate the frame sizes to output.
- 5) The output of the synthetic encoder contains "encoded" frames with random contents but with realistic sizes.

Section 6.1 explains steps 1), 2), and 3), Section 6.2 elaborates on steps 4) and 5). Finally, Section 6.3 briefly discusses the possibility to extend the model for supporting variable frame rate and/or variable frame resolution.

6.1. Choosing the video sequence and generating the traces

The first step we need to perform is a careful choice of a set of video sequences that are representative of the use cases we want to model. Our use case here is video conferencing, so we must choose a low-motion sequence that resembles a "talking head", for instance a news broadcast or a video capture of an actual conference call.

The length of the chosen video sequence is a tradeoff. If it is too long, it will be difficult to manage the data structures containing the traces we will produce in the next steps. If it is too short, there will be an obvious periodic pattern in the output frame sizes, leading to biased results when evaluating congestion controller performance. In our experience, a one-minute-long sequence is a fair tradeoff.

Once we have chosen the raw video sequence, denoted S , we use a live encoder, e.g. [H264] or [HEVC] to produce a set of encoded sequences. As discussed in Section 3, a live encoder's output bitrate can be tuned by varying three input parameters, namely, quantization step size, frame rate, and picture resolution. In order to simplify the choice of these parameters for a given target rate, we assume a fixed frame rate (e.g. 25 fps) and a fixed resolution (e.g., 480p). See section 6.3 for a discussion on how to relax these assumptions.

Following these simplifications, we run the chosen encoder by setting a constant target bitrate at the beginning, then letting the encoder vary the quantization step size internally while encoding the input video sequence. Besides, we assume that the first frame is encoded as an I-frame and the rest are P-frames. We further assume that the encoder algorithm does not use knowledge of frames in the future so as to encode a given frame.

We define R_{min} and R_{max} as the minimum and maximum bitrate at which the synthetic codec is to operate. We divide the bitrate range between R_{min} and R_{max} in $n_s + 1$ bitrate steps of length $l = (R_{max} - R_{min}) / n_s$. We then use the following simple algorithm to encode the raw video sequence.

```

r = R_min
while r <= R_max do
    Traces[r] = encode_sequence(S, r, e)
    r = r + l

```

where function `encode_sequence` takes as parameters, respectively, a raw video sequence, a constant target rate, and an encoder algorithm; it returns a vector with the sizes of frames in the order they were encoded. The output vector is stored in a map structure called `Traces`, whose keys are bitrates and values are frame size vectors.

The choice of a value for n_s is important, as it determines the number of frame size vectors stored in map `Traces`. The minimum value one can choose for n_s is 1, and its maximum value depends on the amount of memory available for holding map `Traces`. A reasonable value for n_s is one that makes the steps' length $l = 200$ kbps. We will further discuss step length l in the next section.

6.2. Using the traces in the syntethic codec

The main idea behind the trace-based synthetic codec is that it mimics a real live codec's rate adaptation when the congestion controller updates the target rate $R_v(t)$. It does so by switching to a different frame size vector stored in the map `Traces` when needed.

6.2.1. Main algorithm

We maintain two variables `r_current` and `t_current`:

- * `r_current` points to one of the keys of the map `Traces`. Upon a change in the value of $R_v(t)$, typically because the congestion controller detects that the network conditions have changed, `r_current` is updated to the greatest key in `Traces` that is less than

or equal to the new value of $R_v(t)$. For the moment, we assume the value of $R_v(t)$ to be clipped in the range $[R_{\min}, R_{\max}]$.

```

r_current = r
such that
  ( r in keys(Traces) and
    r <= R_v(t) and
    (not(exists) r' in keys(Traces) such that r < r' &lt;= R_v(t)) )

```

* t_{current} is an index to the frame size vector stored in $\text{Traces}[r_{\text{current}}]$. It is updated every time a new frame is due. We assume all vectors stored in Traces to have the same size, denoted size_traces . The following equation governs the update of t_{current} :

```

if t_current < SkipFrames then
  t_current = t_current + 1
else
  t_current = ((t_current+1-SkipFrames) % (size_traces- SkipFrames) )
              + SkipFrames

```

where operator $\%$ denotes modulo, and SkipFrames is a predefined constant that denotes the number of frames to be skipped at the beginning of frame size vectors after t_{current} has wrapped around. The point of constant SkipFrames is avoiding the effect of periodically sending a (big) I-frame followed by several smaller-than-normal P-frames. We typically set SkipFrames to 20, although it could be set to 0 if we are interested in studying the effect of sending I-frames periodically.

We initialize r_{current} to R_{\min} , and t_{current} to 0.

When a new frame is due, we need to calculate its size. There are three cases:

- a) $R_{\min} \leq R_v(t) < R_{\max}$: In this case we use linear interpolation of the frame sizes appearing in $\text{Traces}[r_{\text{current}}]$ and $\text{Traces}[r_{\text{current}} + 1]$. The interpolation is done as follows:

```

size_lo = Traces[r_current][t_current]
size_hi = Traces[r_current + 1][t_current]
distance_lo = ( R_v(t) - r_current ) / 1
framesize = size_hi * distance_lo + size_lo * (1 - distance_lo)

```

- b) $R_v(t) < R_{\min}$: In this case, we scale the trace sequence with the lowest bitrate, in the following way:

```
factor = R_v(t) / R_min
framesize = max(1, factor * Traces[R_min][t_current])
```

c) $R_v(t) \geq R_{max}$: We also use scaling for this case. We use the trace sequence with the greatest bitrate:

```
factor = R_v(t) / R_max
framesize = factor * Traces[R_max][t_current]
```

In case b), we set the minimum to 1 byte, since the value of factor can be arbitrarily close to 0.

6.2.2. Notes to the main algorithm

* Reacting to changes in target bitrate. Similarly to the statistical model presented in Section 5, the trace-based synthetic codec has a time bound, τ_v , to reacting to target bitrate changes. If the codec has reacted to an update in $R_v(t)$ at time t , it will delay any further update to $R_v(t)$ to time $t + \tau_v$. Note that, in any case, the value of τ_v cannot be chosen shorter than the time between frames, i.e. the inverse of the frame rate.

* I-frames on demand. The synthetic codec could be extended to simulate the sending of I-frames on demand, e.g., as a reaction to losses. To implement this extension, the codec's API is augmented with a new function to request a new I-frame. Upon calling such function, $t_{current}$ is reset to 0.

* Variable length l of steps defined between R_{min} and R_{max} . In the main algorithm's description, the step length l is fixed. However, if the range $[R_{min}, R_{max}]$ is very wide, it is also possible to define a set of steps with a non-constant length. The idea behind this modification is that the difference between 400 kbps and 600 kbps as bitrate is much more important than the difference between 4400 kbps and 4600 kbps. For example, one could define steps of length 200 Kbps under 1 Mbps, then length 300 kbps between 1 Mbps and 2 Mbps, 400 kbps between 2 Mbps and 3 Mbps, and so on.

6.3. Varying frame rate and resolution

The trace-based synthetic codec model explained in this section is relatively simple because we have fixed the frame rate and the frame resolution. The model could be extended to have variable frame rate, variable frame resolution, or both.

When the video quality for a given bitrate is low, one can decrease the frame rate (if the video sequence is currently low motion) or the frame resolution in order to attempt an improvement in the quality-

of-experience (QoE). On the other hand, if the bitrate increases to a point where there is no longer a perceptible improvement in the QoE, then we might afford to increase the frame resolution or the frame rate (useful if the video is currently high motion).

Many techniques have been proposed to choose over time the best values for these two parameters, together with the quantization step size, in order to maximize the quality of live video codecs [Ozer2011], [Hu2010]. In the future, we will consider extending the trace-based codec to be able to use variable frame rate and/or resolution.

From the perspective of congestion control performance, varying the frame resolution will not impact the outcome of a synthetic video codec: the resulting encoded video frames bear the same data size regardless of resolution choice. On the other hand, different choices of frame rates lead to different levels of burstiness in the encoded video traffic trace: e.g., many small packets at a high frame rate vs. sparsely spaced large packets at a low frame rate. Such difference in traffic profiles may affect the performance of congestion control differently, especially when outgoing packets are not paced at the transport module. We leave the investigation of varying frame rate to future work.

7. Combining The Two Models

This section discusses the pros and cons of the two reference models, as well as how one may combine them for evaluation of RMCAT candidates.

[EDITOR'S NOTE: Will add more details after sollicitating initial feedback on the draft from mailing list and in-meeting presentation.]

8. IANA Considerations

There are no IANA impacts in this memo.

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.
- [RFC3550] Schulzrinne, H., Casner, S., Frederick, R., and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", STD 64, RFC 3550, July 2003.

- [RFC3551] Schulzrinne, H. and S. Casner, "RTP Profile for Audio and Video Conferences with Minimal Control", STD 65, RFC 3551, July 2003.
- [RFC3611] Friedman, T., Caceres, R., and A. Clark, "RTP Control Protocol Extended Reports (RTCP XR)", RFC 3611, November 2003.
- [RFC4585] Ott, J., Wenger, S., Sato, N., Burmeister, C., and J. Rey, "Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF)", RFC 4585, July 2006.
- [RFC5506] Johansson, I. and M. Westerlund, "Support for Reduced-Size Real-Time Transport Control Protocol (RTCP): Opportunities and Consequences", RFC 5506, April 2009.
- [I-D.ietf-avtcore-rtp-circuit-breakers]
Perkins, C. and V. Singh, "Multimedia Congestion Control: Circuit Breakers for Unicast RTP Sessions", draft-ietf-avtcore-rtp-circuit-breakers-05 (work in progress), February 2014.
- [I-D.ietf-rmcat-eval-criteria]
Singh, V. and J. Ott, "Evaluating Congestion Control for Interactive Real-time Media", draft-ietf-rmcat-eval-criteria-01 (work in progress), March 2014.
- [I-D.ietf-rmcat-cc-requirements]
Jesup, R., "Congestion Control Requirements For RMCAT", draft-ietf-rmcat-cc-requirements-04 (work in progress), April 2014.
- [H264] ITU-T Recommendation H.264, "Advanced video coding for generic audiovisual services",
<<http://www.itu.int/rec/T-REC-H.264-201304-I>>.
- [HEVC] ITU-T Recommendation H.265, "High efficiency video coding", .

9.2. Informative References

- [I-D.ietf-rtcweb-use-cases-and-requirements]
Holmberg, C., Hakansson, S., and G. Eriksson, "Web Real-Time Communication Use-cases and Requirements", draft-ietf-rtcweb-use-cases-and-requirements-14 (work in progress), February 2014.

- [RFC5033] Floyd, S. and M. Allman, "Specifying New Congestion Control Algorithms", BCP 133, RFC 5033, August 2007.
- [Hu2010] Hu, H., Ma, Z., and Y. Wang, "Optimization of Spatial, Temporal and Amplitude Resolution for Rate-Constrained Video Coding and Scalable Video Adaptation", inproceedings in Proc. 19th IEEE International Conference on Image Processing (ICIP'12), September 2012.
- [Ozer2011] Ozer, J., "Video Compression for Flash, Apple Devices and HTML5", ISBN ISBN-13:978-0976259503, 2011.
- [Tanwir2013] Tanwir, S. and H. Perros, "A Survey of VBR Video Traffic Models", journal IEEE Communications Surveys and Tutorials, vol. 15, no. 5, pp. 1778-1802., October 2013.

Authors' Addresses

Xiaoqing Zhu
Cisco Systems
12515 Research Blvd., Building 4
Austin, TX 78759
USA

Email: xiaoqzhu@cisco.com

Sergio Mena de la Cruz
Cisco Systems
EPFL, Quartier de l'Innovation, Batiment E
Ecublens, Vaud 1015
Switzerland

Email: semena@cisco.com

Zaheduzzaman Sarker
Ericsson AB
Luleae, SE 977 53
Sweden

Phone: +46 10 717 37 43
Email: zaheduzzaman.sarker@ericsson.com