

AVT
Internet-Draft
Intended status: Standards Track
Expires: January 12, 2011

J. Lennox
Vidyo
E. Iovov
SIP Communicator
E. Marocco
Telecom Italia
July 11, 2010

A Real-Time Transport Protocol (RTP) Header Extension for Client-to-Mixer Audio Level Indication
draft-lennox-avt-rtp-audio-level-exthdr-02

Abstract

This document defines a mechanism by which packets of Real-Time Transport Protocol (RTP) audio streams can indicate, in an RTP header extension, the audio level of the audio sample carried in the RTP packet. In large conferences, this can reduce the load on an audio mixer or other middlebox which wants to forward only a few of the loudest audio streams, without requiring it to decode and measure every stream that is received.

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 12, 2011.

Copyright Notice

Copyright (c) 2010 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.

This document may contain material from IETF Documents or IETF Contributions published or made publicly available before November 10, 2008. The person(s) controlling the copyright in some of this material may not have granted the IETF Trust the right to allow modifications of such material outside the IETF Standards Process. Without obtaining an adequate license from the person(s) controlling the copyright in such materials, this document may not be modified outside the IETF Standards Process, and derivative works of it may not be created outside the IETF Standards Process, except to format it for publication as an RFC or to translate it into languages other than English.

Table of Contents

1.	Introduction	4
2.	Terminology	4
3.	Audio Levels	4
4.	Signaling (Setup) Information	6
5.	Considerations on Use	6
6.	Limitations	6
7.	Security Considerations	7
8.	IANA Considerations	8
9.	References	8
9.1.	Normative References	8
9.2.	Informative References	8
Appendix A.	Open issues	9
Appendix B.	Changes From Earlier Versions	9
B.1.	Changes From Individual Submission Draft -01	9
B.2.	Changes From Individual Submission Draft -00	9
	Authors' Addresses	10

1. Introduction

In a centralized Real-Time Transport Protocol (RTP) [[RFC3550](#)] audio conference, an audio mixer or forwarder receives audio streams from many or all of the conference participants. It then selectively forwards some of them to other participants in the conference. In large conferences, it is possible that such a server might be receiving a large number of streams, of which only a few should be forwarded to the other conference participants.

In such a scenario, in order to pick the audio streams to forward, a centralized server needs to decode, measure audio levels, and possibly perform voice activity detection on audio data from a large number of streams. The need for such processing limits the size or number of conferences such a server can support.

As an alternative, this document defines an RTP header extension [[RFC5285](#)] through which senders of audio packets can indicate the audio level of the packets' payload, reducing the processing load for a server.

The header extension in this draft is different to, but complementary with, the one defined in [[I-D.ivov-avt-slic](#)], which defines a mechanism by which audio mixers can indicate to clients the levels of the contributing sources that made up the mixed audio.

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](#)] and indicate requirement levels for compliant implementations.

3. Audio Levels

The audio level header extension carries both the level of the audio carried in the RTP payload of the packet it is associated with, as well as an indication as to whether voice activity has been detected in the packet.

The form of the audio level extension block is as follows:

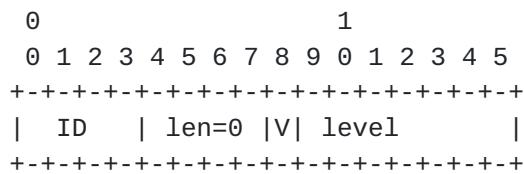


Figure 1

The length field takes the value 0 to indicate that 1 byte follows.

The audio level is defined in the same manner as is audio noise level in the RTP Comfort Noise [[RFC3389](#)] specification. In that specification, the overall magnitude of the noise level is encoded into the first byte of the payload, with spectral information about the noise in subsequent bytes. This specification's audio level parameter is defined so as to be identical to the comfort noise payload's noise-level byte.

The magnitude of the audio level is packed into the seven least significant bits of the single byte of the header extension, shown in Figure 1. The least significant bit of the audio level magnitude is packed into the least significant bit of the byte. The most significant bit of the byte is used as a separate flag bit "V", defined below.

The audio level is expressed in -dBov, with values from 0 to 127 representing 0 to -127 dBov. dBov is the level, in decibels, relative to the overload point of the system, i.e. the maximum-amplitude signal that can be handled by the system without clipping. (Note: Representation relative to the overload point of a system is particularly useful for digital implementations, since one does not need to know the relative calibration of the analog circuitry.) For example, in the case of u-law (audio/pcmu) audio [[ITU.G711.1988](#)], the 0 dBov reference would be a square wave with values +/- 8031. (This translates to 6.18 dBm0, relative to u-law's dBm0 definition in Table 6 of G.711.)

In addition, a flag bit (labeled V) indicates whether the encoder believes the audio packet contains voice activity (1) or does not (0). The voice activity detection algorithm is unspecified and left implementation-specific.

The audio level for digital silence (e.g. all-0 pcmu audio), for example for a muted audio source, MAY be represented as 127 (-127 dBov), regardless of the dynamic range of the encoded audio format.

When this header extension is used with RTP data sent using the RTP

Payload for Redundant Audio Data [[RFC2198](#)], the header's data describes the contents of the primary encoding.

4. Signaling (Setup) Information

The URI for declaring this header extension in an extmap attribute is "urn:ietf:params:rtp-hdrex:audio-level". There is no additional setup information needed for this extension (no extensionattributes).

5. Considerations on Use

Mixers and forwarders generally should not base audio forwarding decisions directly on packet-by-packet audio level information, but rather should apply some analysis of the audio levels and trends. This general rule applies whether audio levels are provided by endpoints (as defined in this document), or are calculated at a server, as would be done in the absence of this information. This section discusses several issues that mixers and forwarders may wish to take into account. (Note that this section provides design guidance only, and is not normative.)

First of all, audio levels should generally be measured over longer intervals than that of a single audio packet. In order to avoid false-positives for short bursts of sound (such as a cough or a dropped microphone), it is often useful to require that a participant's audio level be maintained for some period of time before considering it to be "real", i.e. some type of low-pass filter should be applied to the audio levels. Note, though, that such filtering must be balanced with the need to avoid clipping of the beginning of a speaker's speech.

Additionally, different participants may have their audio input set differently. It may be useful to apply some sort of automatic gain control to the audio levels. There are a number of possible approaches to achieving this, e.g. by measuring peak audio levels, by average audio levels during speech, or by measuring background audio levels (average audio level levels during non-speech).

6. Limitations

The audio levels carried by the extension header defined by this document are defined as dBov, decibels below system overload.

In principle, it could be more useful to have, instead, dB SPL, decibels of sound pressure level. In traditional telephony systems,

telephone handsets were calibrated such that a particular (e.g.) u-law audio level, or analog voltage, corresponded to a particular sound pressure level at the handset's mouthpiece.

However, in many environments, this information is not available. Notably, PC soundcard hardware can only determine the levels of mic- or line-in at the hardware input, and operating systems usually allow further adjustments of audio input levels without providing information about these transformations to applications. Furthermore, in many circumstances, such as speech synthesis or mixed audio, an "audio" signal may in fact never have actually existed as sound pressure at all.

Thus, while information about the correspondance between dB SPL and dBoV, or encoded audio, could be useful, this document does not attempt to define it. If there are circumstances in which this information would be useful, a separate header extension would be straightforward to define. (The information carried by such a header extension could indeed be useful independently from the information in the header extension defined by this document.)

7. Security Considerations

A malicious endpoint could choose to set the values in this header extension falsely, so as to falsely claim that audio or voice is or is not present. It is not clear what could be gained by falsely claiming that audio is not present, but an endpoint falsely claiming that audio is present could perform a denial-of-service attack on an audio conference, so as to send silence to suppress other conference members' audio. Thus, a device relying on audio level data from untrusted endpoints SHOULD periodically audit the level information transmitted, taking appropriate corrective action if endpoints appear to be sending incorrect data. (Note that endpoints MAY choose to measure audio levels prior to encoding, so some degree of discrepancy SHOULD be tolerated.)

In the Secure Real-Time Transport Protocol (SRTP) [[RFC3711](#)], RTP header extensions are authenticated but not encrypted. When this header extension is used, audio levels are therefore visible on a packet-by-packet basis to an attacker passively observing the audio stream. As discussed in [[I-D.perkins-avt-srtp-vbr-audio](#)], such an attacker might be able to infer information about the conversation, possibly with phoneme-level resolution. In scenarios where this is a concern, additional mechanisms SHOULD be used to protect the confidentiality of the header extension. One solution would be header extension encryption [[I-D.lennox-avt-srtp-encrypted-extension-headers](#)].

8. IANA Considerations

This document defines a new extension URI to the RTP Compact Header Extensions subregistry of the Real-Time Transport Protocol (RTP) Parameters registry, according to the following data:

Extension URI: urn:ietf:params:rtp-hdext:audio-level
Description: Audio Level
Contact: jonathan@vidyo.com
Reference: RFC XXXX

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.
- [RFC2198] Perkins, C., Kouvelas, I., Hodson, O., Hardman, V., Handley, M., Bolot, J., Vega-Garcia, A., and S. Fosse-Parisis, "RTP Payload for Redundant Audio Data", [RFC 2198](#), September 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.
- [RFC5285] Singer, D. and H. Desineni, "A General Mechanism for RTP Header Extensions", [RFC 5285](#), July 2008.

9.2. Informative References

- [I-D.ivov-avt-slic]
Ivov, E. and E. Marocco, "A Real-Time Transport Protocol (RTP) Extension Header for Mixer-to-client Audio Level Indication", [draft-ivov-avt-slic-02](#) (work in progress), October 2009.
- [I-D.lennox-avt-srtp-encrypted-extension-headers]
Lennox, J., "Encryption of Header Extensions in the Secure Real-Time Transport Protocol (SRTP)", [draft-lennox-avt-srtp-encrypted-extension-headers-01](#) (work in progress), March 2010.
- [I-D.perkins-avt-srtp-vbr-audio]
Perkins, C. and J. Valin, "Guidelines for the use of Variable Bit Rate Audio with Secure RTP",

[draft-perkins-avt-srtp-vbr-audio-04](#) (work in progress),
July 2010.

[ITU.G711.1988]

International Telecommunications Union, "Pulse code modulation (PCM) of voice frequencies", ITU-T Recommendation G.711, November 1988.

[ITU.P56.1993]

International Telecommunications Union, "Objective Measurement of Active Speech Level", ITU-T Recommendation P.56, March 1988.

[RFC3389] Zopf, R., "Real-time Transport Protocol (RTP) Payload for Comfort Noise (CN)", [RFC 3389](#), September 2002.

[RFC3711] Baugher, M., McGrew, D., Naslund, M., Carrara, E., and K. Norrman, "The Secure Real-time Transport Protocol (SRTP)", [RFC 3711](#), March 2004.

[Appendix A.](#) Open issues

- o In order to more accurately determine signal-to-noise ratio, it would be useful for a sender to also send its estimate of its current audio noise floor. If so, it's unclear whether this would be better as a separate header extension element, or added to this header extension element.
- o It has been suggested to reference ITU P.56 [[ITU.P56.1993](#)] for level measurement. This needs to be investigated.

[Appendix B.](#) Changes From Earlier Versions

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

[B.1.](#) Changes From Individual Submission Draft -01

- o This version is primarily a document refresh.
- o Emil Ivov and Enrico Marocco have been added as co-authors.
- o Additional open issues listed.

[B.2.](#) Changes From Individual Submission Draft -00

- o The draft name has been changed to clarify that this document defines Client-To-Mixer Audio Levels, to more clearly distinguish it from [[I-D.ivov-avt-slic](#)].

- o The header extension format has been changed from a two-byte to a one-byte payload, eliminating the 7 reserved bits and the one must-be-zero bit.
- o The sections Considerations on Use ([Section 5](#)) and Limitations ([Section 6](#)) have been added.
- o It has been noted that senders MAY indicate -127 dBov for digital silence, and that level measurement MAY be done prior to encoding audio.
- o A reference to [[I-D.lennox-avt-srtp-encrypted-extension-headers](#)] has been added to the security considerations.
- o The term "header extension" is now used consistently throughout the document (as opposed to "extension header").

Authors' Addresses

Jonathan Lennox
Vidyo, Inc.
433 Hackensack Avenue
Seventh Floor
Hackensack, NJ 07601
US

Email: jonathan@vidyo.com

Emil Ivov
SIP Communicator
Strasbourg 67000
France

Email: emcho@sip-communicator.org

Enrico Marocco
Telecom Italia
Via G. Reiss Romoli, 274
Turin 10148
Italy

Email: enrico.marocco@telecomitalia.it

