



MASTER 2
RÉSEAUX INFORMATIQUES ET SYSTÈMES EMBARQUÉS

Submitted by
Boris GROZEV
boris@jitsi.org
Strasbourg, June 14, 2014

FINAL REPORT

Media recording for multiparty video conferences based on WebRTC

Supervisor
Dr. Emil IVOV
emcho@jitsi.org

Hosting enterprise
BLUEJIMP



Contents

1	Introduction	1
1.1	<i>BlueJimp</i>	1
1.2	The <i>Jitsi</i> family	2
1.3	WebRTC	3
1.4	XMPP and Jingle	3
1.5	<i>Jitsi Meet</i>	4
1.6	Recording	5
2	Implementing support for the WebRTC transport layer	7
2.1	The RTP stack in <i>libjitsi</i>	8
2.2	RED	9
2.3	Uneven Level Protection Forward Error Correction	10
2.4	Retransmissions	12
2.5	RTCP compound packets	12
3	Recording video	12
3.1	The VP8 codec	12
3.2	Depacketization	13
3.3	Container format	15
3.4	Requesting keyframes	16
3.5	Jitter buffer	16
3.6	Overview	17
4	Recording audio	18
4.1	Opus	18
4.2	Recording a mix directly ("live" mixing)	18
4.3	Recording streams separately	19
4.4	Repairing gaps	19
4.5	Overview	20
5	Recording metadata	21
6	Synchronization	22
6.1	Implementation	23
6.2	Synchronization between different sources	25
7	Dominant speaker identification	25
7.1	Implementation	26
8	Conclusion	27
9	Future work	27

9.1	RTP Retransmissions	28
9.2	Recording audio in <i>ogg/opus</i> format	28
9.3	Support for payload type mappings	28
9.4	Improving post-processing	28
A	Jipopro	29
A.1	Video	29
A.2	Audio	30
A.3	Merging	30
A.4	Performance	30
B	The full contents of a metadata file	30

Abstract

This document describes the implementation of various media recording-related features in a modern video conferencing application based on WebRTC. We start by introducing the work environment and the various technological components involved in the project. We then proceed with a detailed description of the different stages of the process including the way audio and video data is transported over the network, the way it is persistently stored, the way multiple audio and video files are organized and combined in a single flat audio/video file, and how synchronization is ensured. The work on the project is still actively pursued and we therefore also present a number of planned future steps and improvements that will likely be implemented in the near future.

1 Introduction

The ever decreasing cost of bandwidth and processing resources have in the recent years made multi-party video conferencing over the internet viable for personal use. The advent of the WebRTC technology that adds audio/video communication capabilities to web browsers, has made the development of conferencing applications (or the addition of conferencing features to existing applications) simpler than ever before.

The work described in this document is about the development of video recording features within *Jitsi Meet*: an existing WebRTC video conferencing application. Throughout the rest of this section we introduce the working environment, the most important standards and software products used in *Jitsi Meet*, and we discuss the general concept of recording a conference. In sections 3 through 7 we examine specific parts of the recording process in detail. In section 8 we review accomplished work, and in section 9 we discuss new features and optimizations to existing features which we plan to develop in the near future.

1.1 *BlueJimp*

XXX expand?

*BlueJimp*¹ is a small company which offers support and development services mainly focused around the *Jitsi*² family of projects. The FLOSS (Free/Libre Open Source Software) nature of these projects makes for a slightly unusual business model. The company works with various kinds of customers who all have different use cases for *Jitsi* and need it adapted to their needs. While *BlueJimp* has no exclusivity on such adaptations, it is tightly involved in the development of the

¹<https://bluejimp.com>

²<https://jitsi.org>

project and some of the related technologies and standards. This has helped the company acquire significant credibility and offer advantageous price/quality ratios.

In addition to orders from customers, *BlueJimp* also often works on internal projects that aim to enrich *Jitsi* and make it more attractive to both users and enterprises.

BlueJimp is registered in Strasbourg, but the development team is international, with people working from different geographic locations. Most communication happens over the Internet using e-mail, instant messaging and audio/video calls.

My position in *BlueJimp* is that of a software developer. Apart from development, my tasks also involve a fair amount of research, experimentation and optimizations. I have worked on *Jitsi* previously and when my internship began, I was able to quickly get accustomed to the environment.

1.2 The *Jitsi* family

Jitsi is a feature-rich internet communications client. It is free software, licensed under the LGPL[2]. The project was started by Emil Ivov in the University of Strasbourg in 2003, and it was then known as SIP Communicator. In the beginning SIP Communicator was only a SIP client, but through the years it has evolved into a multi-protocol client (XMPP, AIM, Yahoo! and ICQ are now also supported) with a very wide variety of features: instant messaging (IM), video calls, desktop streaming, conferencing (multi-party, both audio and video), cross-protocol calls (SIP to XMPP), media encryption (SRTP), session establishment using ICE, file transfers and more. Most of the development is financed by *BlueJimp*.

Jitsi is written mostly in Java and runs on a variety of platforms (Windows, Linux, MacOSX and Android). The various projects comprise a massive codebase—over 700 000 lines of Java code alone.

A big part of the code which was originally in *Jitsi* is now split into a separate library—*libjitsi*. This allows it to be easily reused in other projects, such as *Jitsi Videobridge*. The code in *libjitsi* deals mainly with multimedia—capture and rendering, transcoding, encryption/decryption and transport over the network of audio and video data. It contains the RTP stack for *Jitsi* (partially implemented in *libjitsi* itself, partially in the external FMJ library).

Jitsi Videobridge is a server-side application which acts as a media relay and/or mixer. It allows a smart client to organize a conference using an existing technology (for example, SIP or XMPP/Jingle), outsourcing the bandwidth-intensive task of media relaying to a server. The organizing client controls *Jitsi Videobridge* over XMPP³, while the rest of the participating clients need not be aware that *Jitsi Videobridge* is in use (for example they can be simple SIP clients). Since the end of 2013 *Jitsi Videobridge* supports ICE and is WebRTC-compatible.

³Although a REST API is also available.

One of the latest additions to the *Jitsi* family is *Jitsi Meet*⁴. This is a WebRTC application, which runs completely within a browser and creates a video conference using a *Jitsi Videobridge* instance to relay the media. *Jitsi Meet* is discussed in more detail in 1.5.

1.3 WebRTC

WebRTC (Web Real-Time Communications) is a set of specifications currently in development, that allow browsers which implement them to open "peer connections". These are direct connections between two browsers (a webserver is used only to setup the connection), and can be used to send audio, video or application data. The specifications are open and are meant to be implemented in the browsers themselves (without the need for additional plug-ins).

WebRTC is divided in two main parts: the JavaScript standard APIs, being defined within the *WebRTC* working group[5] at W3C, and the on-the-wire protocols, being defined within the *RTCWEB* working group[4] at the IETF.

These standards provide web developers with a very powerful tool, which can be used to easily create rich real-time multimedia applications. There is also the possibility to pass arbitrary data in an application-defined format. These allow for some very interesting and more complicated use-cases.

The specifications are still being developed, but are already at an advanced stage. There is an open-source implementation of the network protocols, provided by Google with a BSD-like license. I will refer to this implementation as *webrtc.org* (which is the domain name of the project). Currently all browsers implementing WebRTC (Chrome/Chromium, Opera and Mozilla Firefox) use *webrtc.org* as their base. Because of this, *webrtc.org* is very important – it is used for all practical compatibility testing, making it the de-facto reference implementation.

1.4 XMPP and Jingle

Extensible Messaging and Presence Protocol (XMPP)[25] is a mature, XML-based protocol which allows endpoints to exchange messages in near real-time. The core XMPP protocol only covers instant messaging (that is, the exchange of text messages meant to be read by humans), but there are a variety of extensions that allow the protocol to cover a wide range of use cases. Many such extensions are published as XMPP Extension Protocols (XEPs), and there are XEPs for group chats, user avatars, file transfers, account registration, transporting XMPP over HTTP, discovery of node capabilities, management of wireless sensors (provisioning, control, data collection), and, most relevant here, internet telephony.

Jingle (defined in XEP-0166[14] and XEP-0167[15]) is a signalling protocol, serving a purpose similar to that of SIP: it uses an offer-answer model to setup an RTP

⁴<https://jitsi.org/Projects/JitsiMeet>

session. Many of the protocols often used with SIP, such as ICE[23], ZRTP[33], DTLS-SRTP[10], and RFC4575[24] can also be used with *Jingle*. Mappings have been defined between the two[26], which allow gateways or rich clients to organize cross-protocol calls.

COnferencing with LIghtweight BRIdging (COLIBRI[11]) is an extension developed mostly in *BlueJimp* for use with *Jitsi Videobridge*. It provides a way for a client to control a multi-media relay or mixer, such as *Jitsi Videobridge*. It works with the concept of a *conference*, which contains *channels*, separated in different *contents* (see fig.1). In the most common use case a client requests the creation of a *conference* with a specified number of channels. The mixer allocates

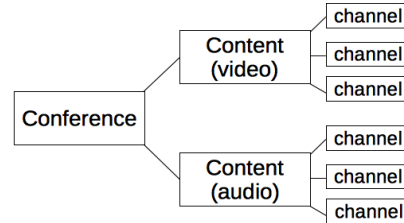


Figure 1: A COLIBRI conference.

local sockets for each *channel* and provides their addresses to the client. The client then uses these transport addresses as its own to establish, for example, a *Jingle* call. Instead of just allocating local sockets, the ICE protocol can be used, in which case the mixer provides a list of ICE candidates for each *channel*. The protocol works with a natural XML representation of a *conference*. After the *conference* is established, the client can add or remove channels from it, or change the parameters (such as the direction in which media is allowed to flow) of an existing *channel*.

1.5 *Jitsi Meet*

Jitsi Meet uses the above-mentioned technologies to create a multi-party video conference. The endpoints of the conference are simply WebRTC-enabled browsers⁵ running the actual *Jitsi Meet* application. They all connect to an XMPP server and join a Multi-User Chat (MUC) chatroom. One of the participants (the first one to enter the chatroom) assumes the role of organizer (or focus).

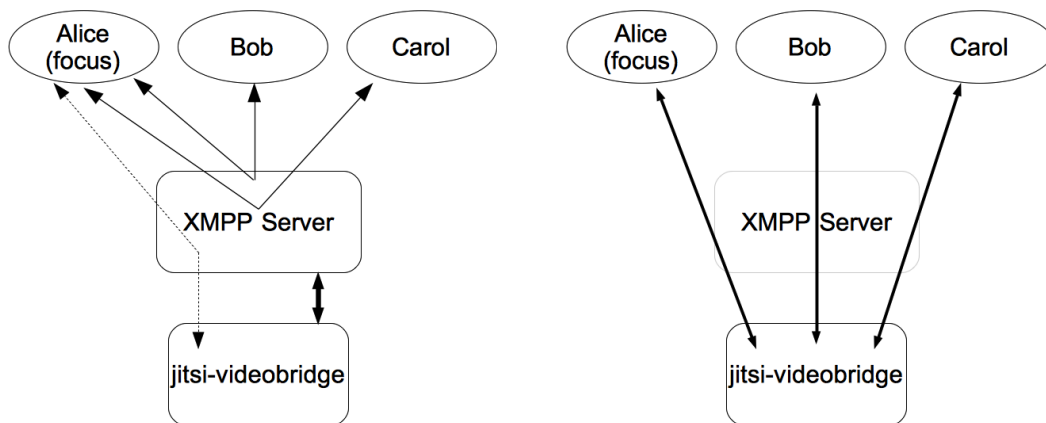
The focus creates a COLIBRI conference on a *Jitsi Videobridge* instance (*jvb*), and allocates two COLIBRI channels for each participant (one for audio, and one for video). Then, it initiates a separate *Jingle* session with each participant, using the transport information (i. e. the list of ICE candidates) obtained from *jvb* instead of its own. When the participants accept the Jingle sessions, they in effect perform ICE and establish direct RTP sessions with *jvb*.

The resulting connections for signalling and media are depicted in figures 2a and 2b respectively.

The *Jitsi Videobridge* instance runs as a relay (as opposed to a mixer) for both video and audio (meaning that it only passes RTP packets between the participants, without considering their payload).

On the user-interface end, *Jitsi Meet* aims to make it as easy as possible for a

⁵Although at the present time only Chrome/Chromium and Opera are supported.



(a) Signalling connections in a *Jitsi Meet* conference. The solid lines are XMPP/Jingle sessions, the dashed line is XMPP/COLIBRI. The thick line is an XMPP Component Connection.

(b) Media connections in a *Jitsi Meet* conference. The lines represent RTP/RTCP sessions.

person to enter or organize a conference. Entering a conference is accomplished by simply opening a URL such as <https://meet.jit.si/ConferenceID> (where *ConferenceID* can be chosen by the user). If *ConferenceID* doesn't exist, it is automatically created and the user assumes the role of focus, inviting anyone who enters later on. If *ConferenceID* exists, the user joins it (possibly after entering a password for the conference).

When in a conference, the interface has two main elements: one big video (taking all available space) and, overlayed on top of it, scaled down versions of the videos of all other participants. Figure 3 is a screenshot from the actual application. Current work is underway to use dominant speaker identification (see section 7) to change the video shown in full size to the person currently speaking.

In contrast to other products for video conferencing over the internet, the whole infrastructure needed to run *Jitsi Meet* can be installed in a custom environment. The only services which are needed are a web server (only serving static content), an XMPP server, and an instance of *Jitsi Videobridge*. This makes *Jitsi Meet* very suitable for businesses (or even individuals) who want full control over their conferencing solution.

1.6 Recording

By recording a multimedia conference in general we mean the following: all the audio and video from the conference is saved to disk in some format, in a way which allows the whole conference to be played back later on.

In our specific case, the recording of a conference has four main parts:

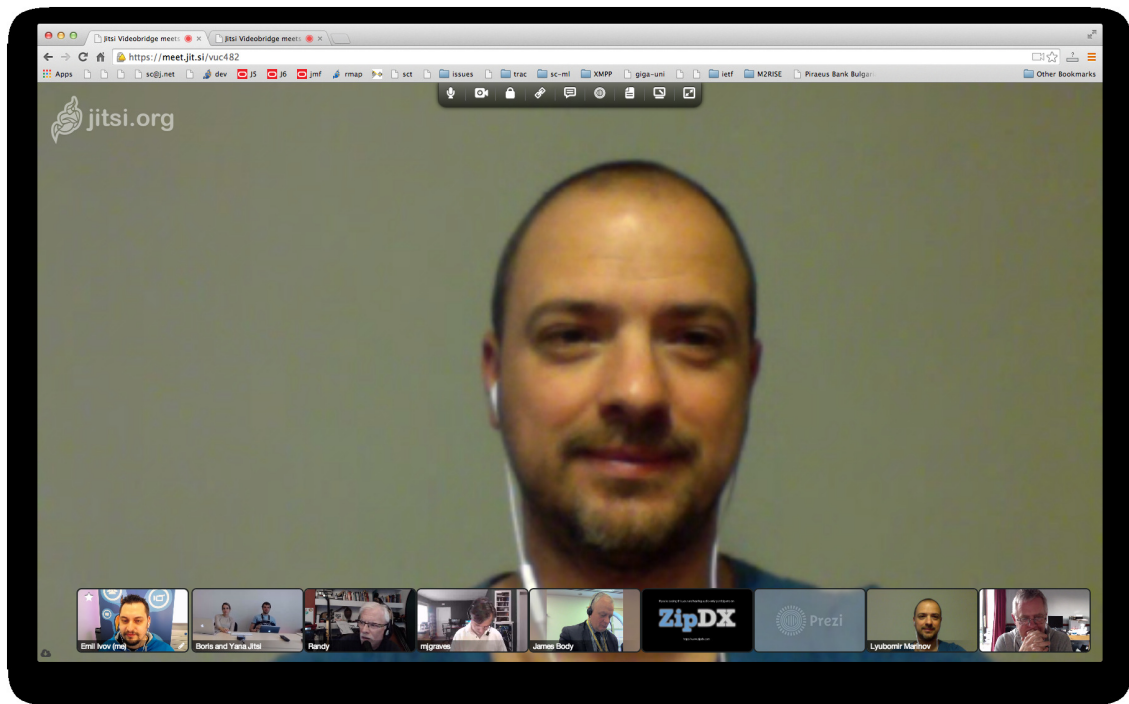


Figure 3: A screen capture from a *Jitsi Meet* conference.

- Recording video
- Recording audio
- Recording metadata
- Post-processing

Recording of the media is the process in which the audio and video RTP streams in the conference are converted to a convenient format and saved to disk. There are many different ways in which this can be done. Our final solution (and some of the ideas that didn't work) are discussed in detail in sections 3 and 4.

Recording of metadata means saving additional information (apart from the media itself) which is necessary to play back the conference later. This includes filenames, timing information, participant names and changes of the active speaker. There is a detailed discussion of the metadata that we use in section 5.

Post-processing in our case means taking all the recorded data and producing a single file with one audio track and one video track. The specifics depend on configuration, but generally all audio is mixed together, and the videos are combined in a way to resemble the *Jitsi Meet* interface. Appendix A discusses the post-processing application.

Where does recording happen? As can be seen in figure 2b, in a *Jitsi Meet* conference both *Jitsi Videobridge* and all the participants have access to the RTP streams, and so could potentially perform recording.

Since the participating clients are running an application within their browser, if we want one of them to do recording, we would need modifications to the browsers. This is inconvenient because users would need to use modified browsers, and because in most use-cases the recordings are going to be stored (and post-processed) on a server, so they would have to be transferred there somehow. To avoid this, a "fake" participant can be added, which does not actually participate in the conference (does not send audio or video), and runs on a server (and without the need for a browser at all). Still, it connects as a normal participant and establishes an RTP session with *Jitsi Videobridge* (see figure 4).

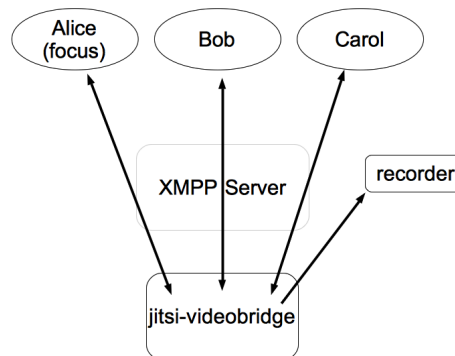


Figure 4: A recording application connected to a *Jitsi Meet* conference as a "fake" participant.

Recording directly on *Jitsi Videobridge* is more straightforward, so our initial implementation was focused on that. However, most of the code resides in *libjitsi*, allowing it to be easily reused.

Shunyang Li, a student from Peking University, is currently working, under my guidance and within context of the Google Summer of Code program, on *Jirecon*⁶—a standalone XMPP container for the recording application described here.

2 Implementing support for the WebRTC transport layer

The documents from the RTCWEB working group at the IETF specify how multimedia is to be transported between WebRTC endpoints. In the most part existing standards are reused.

It is mandatory to use the Interactive Connectivity Establishment (ICE[23]) protocol to establish a session. This assures that an endpoint will not send any media before it has receive consent (in the form of a STUN message) from the remote side. This protects against possible traffic augmentation attacks, in which a malicious web-server causes browsers to send large amounts of data (e.g. a video stream) to a target.

After a connection is established using ICE, a DTLS-SRTP session is started.

⁶<https://github.com/jitsi/jirecon>

This means that the endpoints use Datagram TLS (DTLS[21]) to exchange key material, which is then used to generate session keys for a Secure Real-Time Protocol (SRTP[9]) session. The procedure is defined in RFC5763[10]. In a *Jitsi Meet* conference, each participant’s browser setups two SRTP sessions with *Jitsi Videobridge* in this way.

SRTP provides an unreliable transport. For this reason *webrtc.org* uses a couple of mechanisms on top of SRTP to improve the quality of the media. These mechanisms and their implementation in *libjitsi* are discussed in sections 2.2 to 2.4. Before that section 2.1 gives an overview of the RTP stack used in *libjitsi*.

2.1 The RTP stack in *libjitsi*

libjitsi makes heavy use of the Freedom for Media in Java (FMJ⁷) library. This is an open-source implementation of the Java Media Framework (JMF) API, and it is used in *libjitsi* for a variety of tasks: capture and playback of media, conversion (transcoding) of media, and for handling of basic RTP streams. FMJ is highly extensible, and many components (such as media codecs, capture devices and renderers) are written in *libjitsi*.

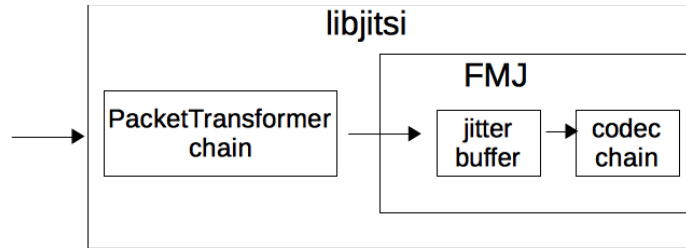


Figure 5: General scheme of the RTP stack in *libjitsi*.

The RTP stack used by FMJ lacks some features: notably support for SRTP and for asymmetric payload type mappings (i.e. sending and receiving a given format with two different RTP payload type numbers). In order for *libjitsi* to implement these features, it intercepts the RTP packets from the actual socket in use, and processes them before passing them on to FMJ. Specifically, packets go through a chain of *PacketTransformers*, which perform various tasks. Figure 5 illustrates this scheme.

PacketTransformers provide a convenient interface to intercept RTP and RTCP packets at different stages of their processing and perform additional operations on them. Despite their name, *PacketTransformers* don’t need to change the packets in any way, they can be used just for monitoring.

Figure 6 lists the currently used *PacketTransformers* in *libjitsi*. A packet which arrives from the network goes through the chain downwards (and when packets are sent, they go through the same chain but in the other direction). The transformer labelled ”RFC6464” is used to extract the audio level information from packets, which

⁷<http://sourceforge.net/projects/fmj/>

include this information in an RTP header extensions defined in RFC6464[12]. The audio levels are used for, among other things, performing dominant speaker identification (see section 7). This transformer also serves as a filter, dropping packets with audio marked to contain silence (in order to avoid unnecessary processing, which is why the transformer is first in the chain). The SRTP transformer decrypts SRTP packets. The "Override PT" transformer changes the payload type numbers of packets. It is used to implement asymmetric payload type mappings. The statistics transformer monitors RTCP packets and extracts statistics from them, making them available to other parts of the library.

The rest of the transformers (the ones in grey) were added with the implementation of the recording system, and will be discussed in the next sections.

2.2 RED

RFC2198[19] defines an RTP payload-type that allows the encapsulation of one or more "virtual" RTP packets in a single RTP packet. It is intended to be used with redundancy data. Its use is negotiated as a regular media format and it does not have a static payload-type number, so a dynamic number is assigned during negotiation.

In *webrtc.org*, RED is supported and used for video streams. In the case of a *Jitsi Meet* conference, it is negotiated between the clients, and *Jitsi Videobridge* has no way of affecting its use or its payload-type number, because it does not actively participate in the offer/answer procedure. This means that in order to record video, the recorder has to understand RED.

We decided that the best way to implement RED in *libjitsi* is as a *PacketTransformer*. There was one complication— *PacketTransformers* work with single packets (they take a single packet as input and produce a single packet as output), while a RED packet may contain multiple "virtual" RTP packets which would need to be output.

We modified *libjitsi*, so that all *PacketTransformers* work with multiple packets at a time— they take an array of packets as input and produce an array as output. This change was not easy, because we had to make sure that we didn't break existing code, but it proved useful later on when we added support for ULPFEC and RTCP compound packets.

We implemented a RED packet transformer following RFC2198 and inserted it in the transformer chain, after the SRTP transformer.

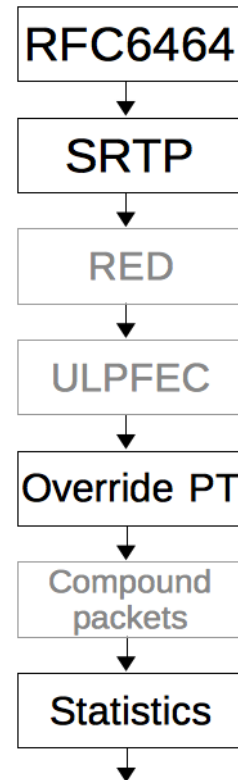


Figure 6: The chain of *PacketTransformers* in *libjitsi*. The shaded elements are new additions.

2.3 Uneven Level Protection Forward Error Correction

In general, Forward Error Correction (FEC) refers to a mechanism which allows lost data to be recovered without retransmission. It involves sending redundant data, in one way or another.

RFC5109[13] defines a specific RTP payload-type for redundant data called Uneven Level Protection FEC (ULPFEC). It is generic in the sense that it can be used with any media payload-type (audio or video, no matter what the codec is).

In *webrtc.org*, ULPFEC is used with video⁸, and while not strictly mandatory for our video-recording (as opposed to RED), it is important because by decreasing the number of irretrievably lost packets, it will improve the quality of the recordings.

ULPFEC (in the rest of the section we refer to it as simply FEC) is applied to an RTP stream (the "media stream", with "media packets") and adds additional packets to it ("FEC packets"). The basic idea is simple – take a set S of a few media packets and apply a parity operation (XOR) on it, resulting in a FEC packet f . If any one of the packets in S is lost, the receiver can use the rest of the packets in S together with f to reconstruct the lost packet⁹.

Along with the parity data, a FEC packet contains two fields which are used to describe the set S from which it was constructed: a "sequence number base" field, and a bitmap field that describes the sequence numbers of the packets in S using the base. The packet f is said to "protect" the packets in S . This scheme allows FEC to work without any additional signalling (apart from the payload-type number negotiated during session initialization).

The amount of FEC packets added to a stream can be controlled by changing the number of protected packets, and this can be done dynamically by the sender, to adapt to network conditions. This is the most common way to use FEC, and the one currently employed by *webrtc.org*: a given fraction of the configured bandwidth is allocated for FEC, and it is changed depending on the packet loss statistics received with RTCP. The aim is to mitigate the effects of packet loss without retransmissions.

Another way to use FEC is for probing the available bandwidth. When the sender detects stable network conditions, it wants to increase its sending bitrate, in order to improve quality. However, this risks causing a congestion, and therefore packet loss. The sender initially increases its sending rate by significantly increasing the amount of FEC. In this case, even if a congestion occurs, the receiver is more likely to be able to reconstruct the media packets (without the need for retransmissions). The sender then monitors the following RTCP reports. If they indicate a high percentage of packet loss, the sender goes back to the previous, lower rate. Otherwise, the sender can keep the total bitrate, but decrease the rate of FEC, using the available bitrate for the encoder instead, thus improving the video quality. This scheme is examined in [17].

⁸For audio, Opus' own FEC scheme which works differently than ULPFEC is used (and it is already supported in *libjitsi*).

⁹This is similar to how RAID5 works.

Implementation We decided to implement FEC as another *PacketTransformer*. This is how it works:

Two buffers of packets are kept: *bMedia* and *bFEC*. With every FEC packet f is associated the number $numMissing(f)$ of media packets protected by f which have not been received.

On reception of a media packet, the values $numMissing(f)$ are recalculated for all f in the *bFEC* buffer. Then, for all f in *bFEC*: if ($numMissing(f) == 0$), then f is removed from *bFEC*. If $numMissing(f) > 1$, then do nothing. If $numMissing(f) == 1$, use f and *bMedia* to reconstruct a media packet and remove f from *bFEC*.

On reception of a FEC packet f , $numMissing(f)$ is calculated, and the same as above is done according to its value.

bFEC is limited to a small size and if a new FEC packet arrives while it is full, the oldest FEC packet is dropped. This prevents "stale" FEC packets (for which $numMissing$ will always be > 1 , because more than one of their protected packets have been lost) to accumulate and cause needless computation.

RFC5109 does not place any restrictions on the placement of FEC packets within a stream, and in our architecture FEC packets are handled entirely in the FEC *PacketTransformer* and not passed on to the rest of the application. This presents a potential problem for the depacketizer (see section 3.2), because it cannot differentiate between a sequence number missing because a packet was lost and a sequence number missing because it was used by a FEC packet.

For this reason we initially implemented re-writing of the RTP sequence numbers of the media packets after they pass the *PacketTransformer*— we decreased their number by the number of FEC packets already received (see figure 7 for an illustration). This still leaves some problems, because a lost FEC packet might be incorrectly interpreted as a missing media packet, and because packets might be incorrectly re-numbered when a packet has arrived out of order.



Figure 7: Re-writing sequence numbers after removal of FEC packets. The marked packets are FEC. The line above shows the sequence numbers before FEC is removed, the line below— after.

Upon further research we found that *webrtc.org* restricts the placement of FEC packets in the stream by only adding them at the end of a VP8 frame, after the RTP packet with the M -bit set (see section 3). This restriction allows the depacketizer to

distinguish between the case of a lost part of a frame and a sequence number being used for FEC, and makes re-numbering of the media packets unnecessary.

2.4 Retransmissions

RFC4585[18] defines a type of RTCP Feedback Message (called NACK), with which a receiver can indicate to a sender that a specific RTP packet (or a set of packets) has not been received. Upon receipt of a NACK, a sender might attempt to retransmit the lost packets.

In *webrtc.org* NACKs and retransmissions are used for video streams. Current versions do retransmissions by sending the exact same RTP packets (without even re-encrypting, which causes some SRTP implementations to falsely detect a replay attack), but there's planned switch to using the payload format defined in RFC4588[22] to encapsulate retransmitted packets.

Currently *libjitsi* does not support RFC4588, but we plan to implement it (as another *PacketTransformer*). Retransmitted packets are not handled in a special way— we just ensure that we have buffers of sufficient size so that retransmitted packets are not dropped as arriving too late (see section 3.5).

When the recording application runs on the *Jitsi Videobridge*, requesting retransmissions with NACKs is not very important, because all RTP packets go through the bridge, and if the bridge is missing a packet, then so are the rest of the participants. The recorder can rely on the participants for sending NACKs, and just make use of the retransmissions themselves. However, when the recorder runs in a separate application (as a "fake" participant), this approach doesn't work, because packets might be lost between the bridge and the recorder. Our implementation does not yet support sending NACKs, but we plan to introduce it.

2.5 RTCP compound packets

RFC3550 specifies that two or more RTCP packets can be combined into a compound RTCP packet. The format is very simple – the packets are just concatenated together and the length fields in their headers allow their later reconstruction.

Because of the lack of support for such packets in FMJ (and because *webrtc.org* makes use of them), we implemented it in *libjitsi* (as a *PacketTransformer*).

3 Recording video

3.1 The VP8 codec

The WebRTC standards do not define a video codec which has to be supported by all clients. There has been a very long discussion at the IETF about whether to make a codec mandatory to implement (MTI), and if so, which one. The four

main options suggested were (i) make the VP8 codec MTI; (ii) make the H264 codec MTI; (iii) have no MTI codecs; (iv) make both VP8 and H264 MTI. No consensus has been reached. Nevertheless, currently VP8 is the de-facto standard codec for WebRTC, because it is the only codec supported by *webrtc.org* (and therefore *Jitsi Meet*).

VP8 is a video compression format, defined in RFC6386[8]. It was originally developed by *On2 Technologies*, which was acquired by *Google* in 2010. *Google* published the specification and released a reference implementation (*libvpx*) under a BSD-like opensource license. They also provided a statement granting permission for royalty-free use of any of their patents used in *libvpx*¹⁰.

Both VP8 in general and *libvpx* work exclusively with the I420 raw (uncompressed) image format[7]. A component called a VP8 encoder, which takes as input an I420 image and produces a "VP8 Compressed Frame". Similarly, a decoder reads VP8 compressed frames and produces I420 images.

A separate specification[32] defines how to transport a VP8 compressed frame over RTP. In short, the process involves optionally splitting a VP8 compressed in parts, prefixing each part with a structure called a "VP8 Payload Descriptor", and then encapsulating each part in RTP. This process is referred to as packetization, and the reverse process (of collecting RTP packets and constructing VP8 compressed frames)—depacketization. Figure 8 provides a high-level overview of the use of VP8 with RTP.

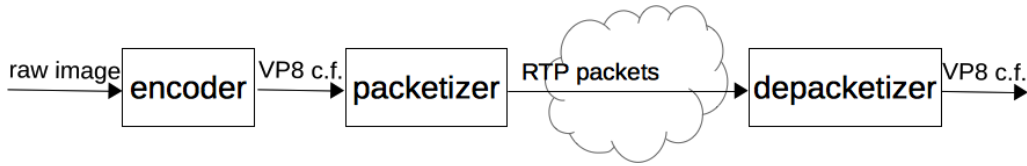


Figure 8: Using VP8 over RTP (c.f. stands for "compressed frame")

libjitsi already has a VP8 implementation, which is used in the *Jitsi* client and consists of four parts: an encoder and decoder (wrappers around *libvpx*), a packetizer and a depacketizer. For the purposes of recording, we need only a depacketizer. We found that the existing depacketizer is not compliant with the specification, and also not compatible with *webrtc.org*. We decided to re-write it from scratch.

3.2 Depacketization

The VP8 Payload Descriptor can be thought of as an extension to the RTP header. It has variable length (between 1 and 6 bytes) and contains, among others, the following fields:

- *S*-bit: start of VP8 partition, set only if the first byte of the payload of the packet is the first byte of a VP8 partition.

¹⁰<http://www.webmproject.org/license/additional/>

- *PID*: Partition ID, specifies the ID of the VP8 partition to which the first byte of the payload of the packet belongs.
- *PictureID*: A running index of frames, incremented by 1 for each subsequent VP8 frame.

Apart from this, the depacketizer also uses the following fields from the RTP header:

- *Timestamp*: a 32-bit field specifying a generation timestamp of the payload. For VP8, all RTP packets from a given frame have the same timestamp.
- *Sequence number*: An index of RTP packets.
- *M-bit*: set for the last RTP packet of a frame (and only for it).

We implemented the algorithm suggested in the specification: we buffer RTP packets locally, until we receive a packet from a new frame. At this point, we check whether the buffer contains a full VP8 frame, and if it does we output it. Otherwise, we drop the buffer and start to collect packets for the next frame. In *libjitsi*, the depacketizer is part of the FMJ codec chain (see fig. 5).

In order to decide whether a received packet is from a new frame or not, the RTP timestamp and the *PictureID* fields are used (if either don't match, then it's a new frame).

The *PID* and *S* fields from the VP8 Payload Descriptor allow us to detect the first packet from a frame – the first, and only the first packet will have both the *S*-bit set and *PID* set to 0.

This allows us to easily check whether we have a full packet in the buffer or not: we have a full packet if: (i) *we have the beginning of a frame*; (ii) *we have the end of a frame (a packet with the M-bit set)*; and (iii) *we have all RTP sequence numbers inbetween*.

The following pseudo-code outlines the procedure.

```

receive(Packet p){
    if (!belongsToBuffer(p))
        flush();
    push(p);

    if (haveFullFrame())
        outputFrame();
}

belongsToBuffer(p){
    if (bufferEmpty())
        return true;

```

```

        else if (bufferRtpTimestamp == p.RtpTimestamp
                && bufferPictureID == p.PictureID)
            return true;
        return false;
    }

haveFullFrame(){
    if (! (buffer.first.S && buffer.first.PID == 0))
        return false;
    if (!buffer.last.M)
        return false;
    for (int i=buffer.first.seq; i<=buffer.last.seq; i++)
        if (!buffer.contains(i))
            return false;
    return true;
}

```

3.3 Container format

After depacketization, we are left with a stream of VP8 Compressed Frames. We needed to decide how to store them on disk. We considered three options:

- Use the *ivf* container format
- Define and use our own container format
- Use the *webm* container format

The *ivf* format is a very simple video-only, VP8-only storage format. It was developed with *libvpx* for the purposes of testing the implementation. It precedes each frame with a fixed-size header containing just the length of the frames and a presentation timestamp. The only advantage of using this format is the relative simplicity of its implementation. The disadvantages include that not many players support it (for example browsers don't support it), and that its lack of extensibility.

Defining our own container format has one advantage, and that's the possibility to design it in a way that allows partial VP8 frames. The *libvpx* decoder has a mode which allows it to decode a frame even if parts of it are missing. In order to use this API, however, the decoder needs to be provided with information about which parts (which VP8 partitions) are missing, and this information would be lost if we use *ivf* or *webm*. The disadvantages of this approach are the complexity that it brings, and the inflexibility with regards to the players – no of the other tools which we use (like *ffmpeg*) will be able to handle it, and we would need to implement our own decoder.

The *webm*[6] format is a subset of *matroska*[3]. It has been designed specifically to contain VP8 (possibly interleaved with audio) and to be played in browsers. It

allows for much more flexibility than *ivf*. The main advantage is that it can be played by many players, and that it supports many features, so we can later extend our implementation if needed.

We found a small library (written in C++, but with Java mappings already available) with a simple API that would allow us to write *webm* files easily, so we decided to ignore *ivf*, postpone the potential definition of a new format, and use *webm*.

We adapted the library to our needs, and implemented a *WebmDataSink* class which takes as input a stream of VP8 compressed frames and saves them in a *webm* file.

A VP8 stream transported over RTP does not have a strictly defined frame rate. Each VP8 frame has its own, independent timestamp, generated by the source (usually at the time of capture from a camera, before VP8 encoding has taken place), and this timestamp gets translated into an RTP timestamp and is carried in RTP.

When we save VP8 in *webm* format, we use the RTP timestamp in order to calculate a presentation timestamp. This is simple— for the first recorded frame we use a presentation timestamp of 0 (in milliseconds) and for subsequent frames we calculate the difference from the first frame.

3.4 Requesting keyframes

VP8 has two types of frames: I-frames (or keyframes) which can be decoded without any other context, and P-frames, whose decoding depends on previously decoded frames. In order to start working, a VP8 decoder needs to first decode a keyframe. Therefore, we want the recorded *webm* files to start with a keyframe.

Because keyframes are rarely sent, and because we want to be able to start recording a conference at any time, we needed a way to trigger the generation of a keyframe. One way to do this, which is supported by *webrtc.org*, is by the use of RTCP Full-Intra Request (FIR) feedback messages (defined in RFC5104[31], section 3.5.1).

We added support for FIR messages in *libjitsi* and made use of them to request keyframes in the beginning of a recording. We faced a difficulty because FIR messages contain an "SSRC of packer sender" field, and *webrtc.org* clients only accept messages from SSRCs that they know about (that is, that have been specifically added via signalling). We had to make *Jitsi Videobridge* generate and use its own SSRC, which it also announces to the focus via COLIBRI.

3.5 Jitter buffer

A jitter buffer is normally used in a real-time multimedia application to introduce a certain amount of delay in the playback of media, mitigating the effects of the

varying delay of packets on the network (the jitter). A buffer which is too small gets emptied quickly when packets are delayed and playback has to be paused. A buffer which is too big adds unnecessary delay to the playback. For this reason adaptive jitter buffers are used, which change their size according to network conditions.

In the case of video with WebRTC, apart from just jitter on the network, there might be packet retransmissions triggered by the receiver (which necessarily arrive at least one RTT later than the originally transmitted packets), and it is beneficial if the buffer is large enough to accept them (as opposed to dropping them because they have arrived too late).

For the purposes of recording, a relatively long delay (on the order of a few seconds) is acceptable, provided that the buffer can be emptied at the end of the recording, without packets being discarded. For this reason we decided to use a fixed-size jitter buffer.

FMJ already includes a jitter buffer in its chain (see figure 5), but for technical reasons it was hard to implement a way to empty it without discarding its contents. We decided to implement our own buffer, in the form of a *PacketTransformer*. For simplicity, we limited its size to a given number of packets, and not to a given length of time. We used a default size of 300 packets, since we observed that this usually corresponds to between 3 and 10 seconds.

3.6 Overview

The components discussed above are pieced together in a *libjitsi* class named *RecorderRtpImpl*¹¹, which we implemented specifically for recording video¹².

RecorderRtpImpl takes its input in the form of RTP packets. It first demultiplexes the packets by their SSRC, and puts them in a jitter buffer, which is placed as the last element of the *PacketTransformer* chain in figure 5.

After they exit the jitter buffer, all packets are passed to an instance of FMJ which is configured to transcode from the "VP8/RTP" format, to the "VP8" (i.e. to do depacketization). FMJ does its own demultiplexing and creates a VP8 depacketizer for each stream. The depacketizer is part of the FMJ codec chain (again, see figure 5), and the internal FMJ jitter buffer is practically not used. The output of FMJ is in the form of a stream of VP8 compressed frames, which is passed to a *WebmDataSink* instance, which produces the final *webm* file.

As soon as a new VP8 stream is detected, before its packets enter the jitter buffer, a keyframe is requested, by sending an RTCP FIR message.

When the recording of a stream is stopped (either because of a user request via the API, or because an RTCP BYE message is received, or because of a timeout), the jitter buffer for the specific SSRC is emptied, its contents being processed.

¹¹Because it implements the *Recorder* interface, using RTP as opposed to a capture device as input

¹²Although we later adapted it to record audio as well.

Using this process, each SSRC stream (with the VP8 format) is saved in a separate *webm* file.

4 Recording audio

The WebRTC standards define two mandatory to implement audio codecs: *G711*[1] and *Opus*[29]. All WebRTC endpoints are required to implement them.

G711 is an audio codec originally developed in the 1970s by the ITU for use in the telephone network. It works on an input PCM signal with a sampling rate of 8000Hz (which already limits the sound quality) and produces an encoded signal with a constant bitrate of 64kbps. It is included in WebRTC for interoperation with legacy devices.

4.1 Opus

Opus is a modern audio codec, whose definition was published as an RFC in 2012. It has a variable bitrate and works with fullband sound (sampled at 48000Hz). The average bitrate is configurable, and at 32kbps (the default in many applications which encode voice) it produces sound with remarkable quality.

In a *Jitsi Meet* conference, since all participants use full-fledged browsers, *Opus* is always the codec being used.

Opus supports frames of different size, but almost all applications (including *webrtc.org*) use frames of 20ms.

Opus has a build-in FEC mechanism which works differently than the RFC5109 mechanism discussed in section 2.3. If enabled for a packet, it encodes the packet at a lower quality, and attaches this version to the subsequent packet. This allows a single lost packet to be recovered (in lower quality), but if two consecutive packets are missing, only one can be recovered. It also has a Packet Loss Concealment (PLC) function, which can be used in the absence of a packet to produce a low level of noise in such a way as to reduce audible crackling.

libjitsi has support for *Opus* and for its FEC and PLC features.

4.2 Recording a mix directly ("live" mixing)

When we started work on recording audio, the most straightforward solution, given the existing code, was to create a mix on-the-spot (in the recording application) and record a single file. This was because *libjitsi* has an advanced audio mixer component, which is used in the *Jitsi* client and *Jitsi Videobridge* (if it is so configured) to mix the audio for a conference. The API allows the addition of a basic RTP stream, and handles all the processing (decoding, resampling, if needed) automatically (see figure 9). Also, the audio mixer serves as a capture device, and can thus be recorded with the existing *libjitsi* recorder.

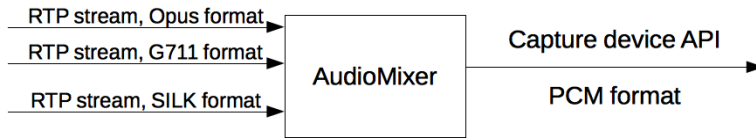


Figure 9: The API of the *libjitsi* audio mixer

Indeed, the implementation proved easy to do, but we faced a problem when we started to think about synchronizing the recorded audio and video. In order to do this, we would at the very least need to know the exact moment where a stream has been added to the mix. This was very hard to do with the audio mixer API.

We would also need a way to guarantee that the spacing between two audio packets in the mix exactly matches the spacing between the packets at the source. We had no way of doing this with the mixer.

We decided that we needed to give up using the *libjitsi* audio mixer and implement code which would allow us finer control over the whole process.

4.3 Recording streams separately

Since we were going to implement the code to handle the decoding and recording of audio streams anew, doing mixing in the application would require additional work, and would introduce substantial complexity. Additionally, there are already many tools which allow the mixing of already recorded files. For these reasons, we decided to record streams individually.

We already had a *RecorderRtpImpl* (see section 3.6) which we used for video, and which performed many of the things needed for audio as well. We adapted *RecorderRtpImpl* so that it could handle streams with different formats, including audio formats.

The existing code in *libjitsi* supports recording only in *mp3* format. We decided to use this, for the time being, because it is relatively easy to transcode it to whatever format is needed in post-processing.

For an audio stream, we use FMJ to decode (to a PCM format), and then encode it and save it in *mp3* format.

4.4 Repairing gaps

In contrast to video in a *webm* file, in which frames have their individual presentation timestamps, an audio stream consists of just sound samples to be played one after the other. This means that if some samples are missing for some reason (e.g. because of a packet lost in the network), the resulting stream will be shorter than the original.

When recording a conference, which might last a relatively long time (possibly hours), the missing samples accumulate, and, when audio and video are merged in the end, this leads to audio and video drifting apart. In order to solve this problem,

any gaps in the audio streams need to be repaired.

The *libjitsi* implementation of *Opus* already repairs gaps on two levels while decoding. First, it tries to use FEC, which can be used when only a single packet is missing. Then, in case of 3 or less missing packets, it uses *Opus*' PLC. But gaps of more than three packets are not repaired.

Apart from lost packets, there is another cause of gaps in the audio streams. When participants in a conference use the "mute" functionality, their browsers continue to send audio packets. However, these packets are specifically marked as containing silence using RFC6464, and *Jitsi Videobridge* drops them. When a participant unmutes, this results in a gap, which is possibly many minutes long.

Repairing such a gap by adding silence samples would require non-negligible amount of computation, because all these samples will then need to be encoded. In such cases it would be better to re-start the recording for the stream in a new file, giving it its own timestamp in the metadata, so that it can be properly mixed in post-processing.

To solve this problem, we decided to introduce an element in the FMJ codec chain, which would be placed between the *Opus* decoder and the *mp3* encoder, and would detect and handle gaps. We called this the *SilenceEffect*.

It works by monitoring the timestamps of the buffers which pass through it. It handles small gaps (the length of gaps to consider short is configurable, with a default of 3 seconds) by inserting the necessary amount of samples (accurate up to a single sample). For gaps longer than this, it signals to the *RecorderRtpImpl*, which restarts the recording into a new file.

This scheme allows us to know, for each recorded audio file, the RTP timestamp of the first RTP packet, the audio from which is contained in the file. This piece of information allows us to ensure the synchronization of the resulting audio and video files (see section 6).

4.5 Overview

Initially we started with recording a mix of all audio streams, because that was the easiest solution. We later changed to recording streams individually. Each incoming audio stream (which is usually *Opus*, but any of the codecs supported in *libjitsi* can be used as well), is first decoded (and re-sampled to 48000Hz, if necessary), then gaps are repaired by inserting silence. Then the stream is encoded and saved in an *mp3* file. See figure 10 for an illustration. If a gaps of more than 3 seconds is detected, we start recording in a new file.

This solution involves re-encoding and that makes it suboptimal – it is more computationally expensive than necessary and reduces the sound quality of the recording. A better solution would be, whenever possible, to record in a container format that does can accommodate the original format, and thus does not require transcoding or even decoding (as we do for video). For *Opus*, a viable solution would

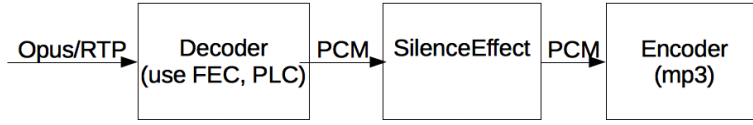


Figure 10: The chain used to transcode audio for the purposes of recording.

be to use the ogg/opus format and we plan to implement it , see section 9.2.

5 Recording metadata

In order for the post-processing application to correctly process the recorded audio and video files, it needs some additional information— a list of file names, at the very least. In the context of our recording system, we call this information metadata.

We decided to keep all metadata in a single file, and to use a JSON format to store it. This would allow to easily extend it with new fields if the need arises.

We defined a JSON object, which we call a *RecorderEvent*, and the metadata file consist of two arrays of such objects: one array for audio and one for video.

All *RecorderEvents* have two mandatory fields: an event type, and an instant. The instant is an integer that represents the time at which the event took place, in milliseconds. The instant has no global interpretation, it is only used to provide timing information relative between the events in a single metadata file. There are three types of events:

RECORDING_STARTED events signify that a recording began in a specific file (given by the "filename" field), at a specific time. They also contain an "ssrc" field, giving the SSRC associated with the recording. This field is used by the post-processing application to associate different files with a single participant, and is necessary in case recording of a stream is split among more than one file. These events also contain a "mediaType" field used to distinguish between audio and video, and two optional fields for additional information about the participant: "participantName" and "participantDescription". This is used in post-processing to overlay some text identifying the participant on top of their video.

RECORDING_ENDED events indicate that a particular recording ends at a specific time. These events have "filename", "ssrc" and "mediaType" fields with the same meaning as for RECORDING_STARTED events. The post-processing application uses these events to decide when to remove a certain audio or video stream from the final mix, but they are optional, because it can determine this from the actual media file.

SPEAKER_CHANGED events are used to signal that the dominant speaker in the conference has changed (at a specific time). They include an "audioSsrc" field which specifies the SSRC of the audio stream of the new dominant speaker, and an "ssrc" field which specifies the SSRC of the video stream associated with the new

dominant speaker. The post-processing application uses `SPEAKER_CHANGED` events to change the video shown in full size.

A few examples of *RecorderEvents* follow, and appendix B has the full contents of a metadata file from an actual recording of a conference.

```
{
  "instant" : 1395767658389,
  "type" : "RECORDING_STARTED",
  "filename" : "3360910907.webm",
  "ssrc" : 3360910907,
  "mediaType" : "video",
  "aspectRatio" : "16_9",
  "participantName" : "Jane Doe",
  "participantDescription" : "
}

{
  "instant" : 1395767704521,
  "type" : "RECORDING_ENDED",
  "filename" : "1277672956-2.mp3",
  "ssrc" : 1277672956,
  "mediaType" : "audio"
}

{
  "instant" : 1395767658527,
  "type" : "SPEAKER_CHANGED",
  "ssrc" : 500778727,
  "audioSsrc" : 1277672956
}
```

6 Synchronization

As explained in the previous section, for each recorded audio and video file, we save a *RecorderEvent*, and each such event has an "instant" field containing an integer specifying milliseconds. This field specifies the relative time at which the audio or video should be included in the final mix. This section describes how the values of the "instant" fields are calculated, in order to provide proper synchronization between the audio and video of a participant and between different participants in the final result produced by the post-processing application.

6.1 Implementation

If we use a simple and naive method, for example, saving the local time when a recording starts, this would leave the audio and the video stream of a participant not synchronized, because of network jitter and differences in the time used for local processing. Therefore we would like to somehow use timing information, which comes from a place as close as possible to the source of the media.

According to its definition, the timestamp field in an RTP packet "reflects the sampling instant of the first octet in the RTP data packet"[27]. When we record video, we know the RTP timestamp of the first frame written, and similarly for audio we know the RTP timestamp of the first packet which is included in the recording¹³.

We implemented a component, named *Synchronizer*, which is responsible for translating RTP timestamps (coming from different RTP streams) into a common format.

As the common format, we used the time provided by *Java's* *System.currentTimeMillis()*, which gives the "difference, measured in milliseconds, between the current time and midnight, January 1, 1970 UTC)". We can then use this format directly when writing the *RecorderEvents*' instants.

The RTP timestamp is a 32bit unsigned integer field, which uses a frequency specific to the media format. For example, *Opus* uses a frequency of 48000Hz, and *VP8* uses 90000Hz. This means that two *Opus* packets, generated exactly one second apart will have RTP timestamps differing by 48000. The initial RTP timestamp should be chosen at random.

When two RTP streams (for example one audio and one video) come from the same source (the same participant in a conference), their RTP timestamps are not directly related, but they share the same wall clock from which the RTP timestamps are generated.

Saving mappings RTCP Sender Report packets contain two timestamp fields. One is an RTP timestamp, and the other is a representation of the wall clock in NTP timestamp format[16]. Both fields represent the same instant, so we can think of the pair as a mapping. Having this mapping allows us to calculate the relative difference between RTP timestamps from different streams (provided that they come from the same source, of course).

We save one such mapping of an RTP timestamp to and NTP timestamp for each SSRC.

Then, for each source, we save a pair mapping an NTP timestamp, to a time on the local clock. It is important to only use one such mapping, because any consecutive mapping made are likely to not be accurate because of network jitter. To create such a mapping, we use the time of reception of an RTCP Sender Report

¹³The initial implementation of the recording didn't expose this information, we had to add it for the purpose of synchronization.

and the NTP timestamp in said report.

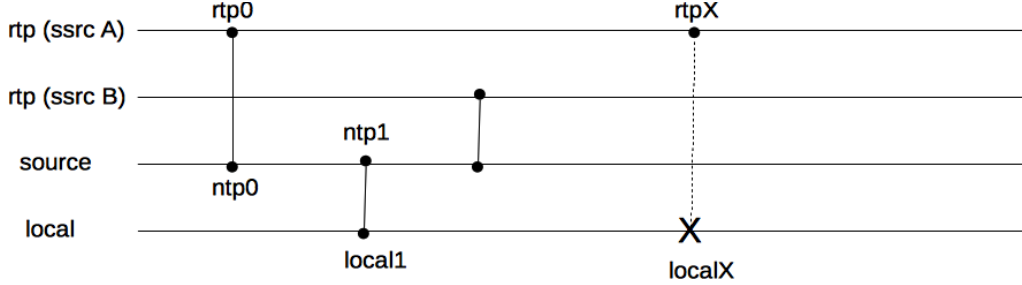


Figure 11: The mappings saved for a single source

The calculation After we have these two necessary mappings— from the RTP clock of a specific SSRC to the source’s clock, and from the source’s clock to the local clock – we can perform translation from the RTP clock to the local clock. The procedure is illustrated on figure 11. The calculation is not complex. Using the notation from the figure, to translate the RTP timestamp $rtpX$ coming for SSRC A to its corresponding local time $localX$ we calculate:

$$localX = local1 + (ntpX - ntp1)$$

where $ntpX$ is the source’s wallclock time corresponding to $rtpX$, for which we have

$$ntpX = ntp0 + (rtpX - rtp0)$$

Simplifying:

$$localX = local1 + (ntp0 - ntp1) + (rtpX - rtp0)$$

Of course we have to take into account the different formats of the timestamps: RTP timestamps are integers with an SSRC-specific rate which are 32-bit and wrap, the local timestamp is an integer in milliseconds, and the NTP timestamps are 64-bit fixed-point numbers (with 32 bits for the whole part and 32 bits for the fraction) in seconds, which we locally store in a *double*.

So the final calculation would be something like

$$localX = local1 + 1000 * (ntp0 - ntp1) + 1000 * diff(rtpX, rtp0) / rateA$$

where $rateA$ is the RTP clock rate for the SSRC A, and $diff$ calculates the difference taking into account the wrap at 2^{32} .

Identifying sources By specification, two SSRCs can be recognized to come from the same source and thus share the same wallclock by looking at the CNAME field in RTCP SDES packets. This would allow for a very simple API for the *Synchronizer*

in which it is fed RTCP packets, and recognizes which SSRCs are grouped together automatically.

However, there is a bug in Chrome¹⁴, which makes it use different CNAME for each SSRC, making it impossible to deduce which SSRCs come from the same source by only looking at RTP and RTCP packets. We implemented a workaround, in which we explicitly specify to the *Synchronizer* identifiers for the source of each SSRC (and we get this information from the signalling path).

We support the CNAME mode of operation and plan to switch to it if the bug in Chrome is fixed.

Overview The scheme just presented allows for audio and video to be very accurately synchronized, because it takes into account timestamps sampled at the source of the streams. It has the drawback that it needs to be initialized before it can perform translation – it usually needs an RTCP Sender Report packet. Before that, the calculation simply cannot be performed. This makes the usage of the *Synchronizer* slightly inconvenient.

6.2 Synchronization between different sources

When streams come from different sources we have no mechanism which allows us to synchronize them. As explained in the previous section, we keep a single mapping of a time of the source’s wall clock to a local time, for each source. We save these mappings independently. We perform the measurement of the local time as soon as possible (as soon as the packet is received by the recorder), in order to minimize the effects of varying time used for local processing.

These same limitations apply for clients participating in a conference, trying to render multiple incoming streams together.

7 Dominant speaker identification

Dominant speaker identification (DSI) refers to the process of continuously analysing a set of audio streams (which are assumed to contain human voice) and keeping track of the one which comes from the person currently speaking (or the person currently speaking "the most"). Obviously the problem does not have a single solution, and so it is hard to measure objectively how well a system performs. A good description of the general problem and a proposed solution is available in [30].

In a *Jitsi Meet* conference there are two use-cases for DSI: for "live" use in a conference and for recording. For "live" use, the purpose is to have the client interface change during a conference, according to who is the dominant speaker (i.e. the video shown in full changes, following the dominant speaker). In this case DSI is

¹⁴<https://code.google.com/p/webrtc/issues/detail?id=3431>

performed on *Jitsi Videobridge*, which uses WebRTC data channels¹⁵ to notify the clients of the change.

For use in recording, the purpose is for the post-processing application to take into account the dominant speaker when combining the videos (in general, it follows the *Jitsi Meet* interface and renders the dominant speaker in full size). In this case, DSI is performed by the recording application¹⁶, and changes of the dominant speaker are saved as metadata.

7.1 Implementation

Our first task was to design and implement an API and a simple, non-optimized algorithm, which we could later improve upon.

We came up with a very simplified scheme, implemented it and tweaked its parameters until it seemed to work well in a conversation with two people. It only takes into account the audio levels (that is, the loudness of the sound) of the different streams. It works like this:

At all times one stream is considered active, while the rest are considered 'competing'. All streams have an associated score. When new audio levels are measured for a stream its score is recomputed, and the scores of all streams are examined in order to determine if one of the competing streams should replace the then active stream.

In order to be 'eligible' to replace the active stream, a competing stream has to:

1. Have a score at least C times as much as the score of the currently active stream.
2. Have score at least MIN_SCORE.
3. Have been active in the last MIN_ACTIVE milliseconds.

The first rule helps to avoid often changing the active when there are two streams with similar levels. The second prevents the active speaker from being changed during a pause of his speech, while no one else is speaking either (but someone else generates higher levels during "silence" than the speaker). The third is to make sure that when participants leave a conference, they aren't mistakenly chosen as active (this is due to implementation details)

The values for the parameters are very dependent on the exact rules for scoring. We chose the values based on a few uncontrolled experiments we performed. We

¹⁵These are channels which carry application data (and not audio or video), but are negotiated and established in the same way as the media channels.

¹⁶Although, if the recording application is running as a participant, it can potentially also open a WebRTC data channel to the *Jitsi Videobridge* and use the events received there, instead of doing DSI itself.

used $\text{MIN_ACTIVE} = 1000$ and $C = 1.15$ (MIN_SCORE is too arbitrary to deserve mention).

Scores are computed as follows: $C_{recent} * \text{avg}(0, N_1) + C_{older} * \text{avg}(N_1, N_2)$ where $\text{avg}(X_1, X_2)$ is the average audio level for the interval $[\text{now} - X_2, \text{now} - X_1]$ (in milliseconds).

Currently the values of the parameters are: $C_{recent} = 2, C_{older} = 1, N_1 = 250, N_2 = 1250$.

We found that this solution works sufficiently well, at least for the purposes of testing recording and post-processing. One problem that we face is when one of the participants is generating constant noise (caused for example by a spinning fan near the microphone). In this case the algorithm tends to select this participant, even if they are not speaking. One suggestion for a fix is to measure the variance of the sound levels for a participant, and penalize the score if the variance is low. This depends on the assumption that a person speaking would produce sound with varying loudness (and that this variation can be detected with the 20ms resolution that we use). However, we have not yet tried to solve this problem.

Lyubomir Marinov has now taken over this task, and has implemented an improved algorithm that more closely resembles the one proposed in [30]. However, his implementation still relies only on "audio levels" and not on other analysis of the audio samples. This offers a significant advantage in our use case, because this information (the "audio level" for a particular RTP packet) is already calculated by the sending client, and is included in the RTP packet in the format specified in RFC6464[12]. This allows *Jitsi Videobridge* to do DSI very efficiently, without the need to decode the received audio streams.

8 Conclusion

During the last six months we have successfully managed to implement a viable recording solution. We have also learned the shortcoming of some of our initial ideas, such as live audio mixing and post-processing via. optimal quality and scalability characteristics.

We have gathered valuable experience.

- de ce que ça a apporté à l'entreprise (realisation) - de ce que ça vous a apporté (acquisition de compétences)

9 Future work

The recording system examined in this document is now fully functional, but there's room for many improvements and optimizations. This section discusses features planned for the near future.

9.1 RTP Retransmissions

There are two things related to retransmissions which need to be implemented. Handling of packets retransmitted using the format defined in RFC4588 needs to be added, because the next release of Chrome/Chromium is expected to use it. We also need to be able to request retransmission by use of RTCP NACK packets, in order to improve video quality when the recorder is run as a separate application (as opposed to on *Jitsi Videobridge*).

9.2 Recording audio in *ogg/opus* format

The current procedure for recording and post-processing audio transcodes the stream twice. First, the stream is transcoded from *Opus* to *mp3* while recording. Then in post-processing it is transcoded to *vorbis*. We could avoid one reencoding by recording audio directly in a container format which supports *Opus*.

We have planned to implement recording in the *ogg/opus* format. One difficulty that we face, which we don't yet know how to solve is handling of FEC in this case. In the current implementation we control the *Opus* decoder and we repair gaps using FEC, whenever possibly. The *ogg/opus* format supports filling gaps, but the specifications[20, 28] do not define a way to handle FEC, and the currently available tools don't support it.

9.3 Support for payload type mappings

The current implementation uses predefined mappings for the RTP payload type numbers. We need to implement support for set them dynamically. One complication arises, because we currently do not support re-writing of the payload type numbers for encapsulated formats such as RED and RTX.

9.4 Improving post-processing

The post-processing application (see appendix A) has limited performance, and our examinations so far show that it would require significant changes to allow the

A Jipopro

Jipopro (short for Jitsi Post-Processing) is an application which was created specifically for post-processing recordings for *Jitsi Meet*. It was developed in parallel with the recording implementation discussed in this document, and is closely tied with the formats for media and metadata produced by that system. *Jipopro* is written in Java, but uses external applications for many of its tasks. It's main author is Vladimir Marinov.

On a very high level, *jipopro* does four things:

1. Reads a metadata file.
2. Processes video, producing a single file.
3. Processes audio, producing a single file.
4. Merges the audio and video files together.

The audio and video processing are independent.

A.1 Video

For all video handling, *jipopro* uses *ffmpeg*¹⁷.

First every input *webm* file is transcoded to an MJPEG file with a static frame rate (25fps by default). Then, a timeline of events is constructed according to the metadata. The events divide the timeline into sections. Figure 12 shows an example.

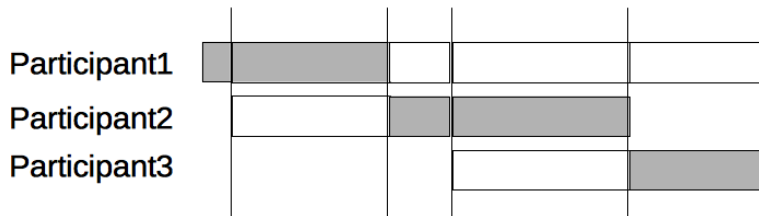


Figure 12: A representation of the process of combining videos. There are 5 sections, separated by vertical lines. The shared regions represent the dominant speaker.

For each section, a list of participants, one of whom is active, is maintained. The sections are then processed: a "slice" of video is taken (with a specific offset) from the MJPEG file for every participant, and the obtained videos (all with the same length) are decoded, combined together (overlaid) and encoded in an MJPEG file again.

After all sections have been processed, they are concatenated together, resulting in a single MJPEG file.

Then the file is transcoded to a *webm* file.

¹⁷<https://ffmpeg.org>

A.2 Audio

For audio handling we use *sox*¹⁸. First all *mp3* files are converted to *wav* format and padded, if necessary. All the resulting files are mixed together in a single *wav* file.

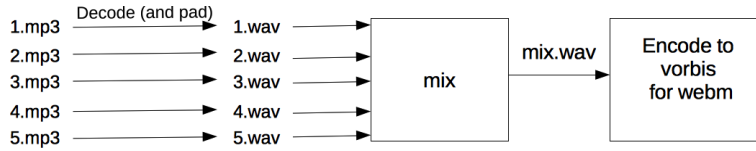


Figure 13: Post-processing for audio.

A.3 Merging

The resulting audio (a single *wav* file) and video (a single *webm* file) are combined using *ffmpeg*. The audio is automatically transcoded to *vorbis*, because it is the codec supported by *webm*¹⁹.

A.4 Performance

The running time of the application depends on a variety of parameters: the number of participants, the chosen frame rate and resolution, the number of `SPEAKER_CHANGED` events.

The process was optimized by running certain tasks– the initial ”decoding”, and the section processing and concatenation– in parallel (in different operating system processes).

We observed that for a particular recording with 4 participants the running time, after all optimizations, was about 1.3 times longer than the length of the conference.

The process is not very efficient and performs at least 2 redundant transcodings. We found no other way to optimizing it without a major re-design, which we are planning to do in the future.

B The full contents of a metadata file

```
{
  "audio" : [
    { "ssrc" : 2473760775,
      "filename" : "2473760775.mp3",
      "type" : "RECORDING_STARTED",
      "instant" : 1401807483964,
```

¹⁸<http://sox.sourceforge.net/sox.html>

¹⁹Although the specification is being extended to allow for Opus

```

    "mediaType" : "audio" },
  { "ssrc" : 2452647913,
    "filename" : "2452647913.mp3",
    "type" : "RECORDING_STARTED",
    "instant" : 1401807484421,
    "mediaType" : "audio" },
  { "ssrc" : 1107392327,
    "filename" : "1107392327.mp3",
    "type" : "RECORDING_STARTED",
    "instant" : 1401807501709,
    "mediaType" : "audio" },
  { "ssrc" : 1565260729,
    "filename" : "1565260729.mp3",
    "type" : "RECORDING_STARTED",
    "instant" : 1401807524062,
    "mediaType" : "audio"},
  { "ssrc" : 1107392327,
    "filename" : "1107392327-1.mp3",
    "type" : "RECORDING_STARTED",
    "instant" : 1401807527409,
    "mediaType" : "audio"},
  { "ssrc" : 2452647913,
    "filename" : "2452647913-1.mp3",
    "type" : "RECORDING_STARTED",
    "instant" : 1401807532821,
    "mediaType" : "audio"}
],

"video" : [
  { "ssrc" : 2612617218,
    "filename" : "2612617218.webm",
    "aspectRatio" : "16_9",
    "type" : "RECORDING_STARTED",
    "instant" : 1401807484011,
    "mediaType" : "video"},

  { "ssrc" : 9413050,
    "filename" : "9413050.webm",
    "aspectRatio" : "16_9",
    "type" : "RECORDING_STARTED",
    "instant" : 1401807484013,
    "mediaType" : "video"},

```

```

{ "ssrc" : 2612617218,
  "audioSsrc" : 2473760775,
  "type" : "SPEAKER_CHANGED",
  "instant" : 1401807484078,
  "mediaType" : "video"},

{ "ssrc" : 9413050,
  "audioSsrc" : 2452647913,
  "type" : "SPEAKER_CHANGED",
  "instant" : 1401807485275,
  "mediaType" : "video"},

{ "ssrc" : 1190123626,
  "filename" : "1190123626.webm",
  "aspectRatio" : "16_9",
  "type" : "RECORDING_STARTED",
  "instant" : 1401807500796,
  "mediaType" : "video"},

{ "ssrc" : 2612617218,
  "audioSsrc" : 2473760775,
  "type" : "SPEAKER_CHANGED",
  "instant" : 1401807503515,
  "mediaType" : "video"},

{ "ssrc" : 2612617218,
  "filename" : "2612617218.webm",
  "type" : "RECORDING_ENDED",
  "instant" : 1401807509187,
  "mediaType" : "video"},

{ "ssrc" : 9413050,
  "audioSsrc" : 2452647913,
  "type" : "SPEAKER_CHANGED",
  "instant" : 1401807512310,
  "mediaType" : "video"},

{ "ssrc" : 3879530045,
  "filename" : "3879530045.webm",
  "aspectRatio" : "16_9",
  "type" : "RECORDING_STARTED",

```

```
    "instant" : 1401807522745,  
    "mediaType" : "video"},  
  
    { "ssrc" : 3879530045,  
      "audioSsrc" : 1565260729,  
      "type" : "SPEAKER_CHANGED",  
      "instant" : 1401807524818,  
      "mediaType" : "video"},  
  
    { "ssrc" : 1190123626,  
      "filename" : "1190123626.webm",  
      "type" : "RECORDING_ENDED",  
      "instant" : 1401807534861,  
      "mediaType" : "video"},  
  ]  
}
```

References

- [1] G.711 : Pulse code modulation (PCM) of voice frequencies. <http://www.itu.int/rec/T-REC-G.711/e>.
- [2] GNU Lesser General Public Licence. <https://www.gnu.org/licenses/lgpl.html>.
- [3] Matroska. <http://matroska.org/technical/specs/index.html>.
- [4] Real-time communication in web-browsers working group. <http://tools.ietf.org/wg/rwcweb/>.
- [5] Web real-time communications working group. <http://www.w3.org/2011/04/webrtc/>.
- [6] The webm project. <http://www.webmproject.org/about/>.
- [7] Yuv pixel formats. <http://www.fourcc.org/yuv.php#IYUV>.
- [8] J. Bankoski, J. Koleszar, L. Quillio, J. Salonen, P. Wilkins, and Y. Xu. VP8 Data Format and Decoding Guide, Internet Engineering Task Force Request for Comments (RFC) 5245, November 2011.
- [9] M. Baugher, D. McGrew, M. Naslund, E. Carrara, and K. Norrman. The Secure Real-time Transport Protocol (SRTP), Internet Engineering Task Force Request for Comments (RFC) 3711, March 2004.
- [10] J. Fischl, H. Tschofenig, and E. Rescorla. Framework for Establishing a Secure Real-time Transport Protocol (SRTP) Security Context Using Datagram Transport Layer Security (DTLS), Internet Engineering Task Force Request for Comments (RFC) 5763, May 2010.
- [11] Emil Ivov, Lyubomir Marinov, and Philipp Hancke. CONferences with Lightweight BRIdging (COLIBRI), XMPP Standards Foundation XEP-0340, January 2014. <http://xmpp.org/extensions/xep-0340.html>.
- [12] J. Lennox, E. Ivov, and E. Marocco. A Real-time Transport Protocol (RTP) Header Extension for Client-to-Mixer Audio Level Indication, Internet Engineering Task Force Request for Comments (RFC) 6464, December 2011.
- [13] A. Li. RTP Payload Format for Generic Forward Error Correction, Internet Engineering Task Force Request for Comments (RFC) 5109, December 2007.
- [14] Scott Ludwig, Joe Beda, Peter Saint-Andre, Robert McQueen, Sean Egan, and Joe Hildebrand. Jingle, XMPP Standards Foundation XEP-0166, December 2009. <http://xmpp.org/extensions/xep-0166.html>.
- [15] Scott Ludwig, Peter Saint-Andre, Sean Egan, Robert McQueen, and Diana Cionoiu. Jingle RTP Sessions, XMPP Standards Foundation XEP-0167, December 2009. <http://xmpp.org/extensions/xep-0167.html>.

- [16] D. Mills, U. Delaware, J. Martin, J. Burbank, and W. Kasch. Network Time Protocol Version 4: Protocol and Algorithms Specification, Internet Engineering Task Force Request for Comments (RFC) 5905, June 2010.
- [17] Marcin Nagy, Varun Singh, Jörg Ott, and Lars Eggert. Congestion control using fec for conversational multimedia communication. *CoRR*, abs/1310.1582, 2013.
- [18] J. Ott, S. Wenger, N. Sato, C. Burmeister, and J. Rey. Extended RTP Profile for Real-time Transport Control Protocol (RTCP)-Based Feedback (RTP/AVPF), Internet Engineering Task Force Request for Comments (RFC) 4585, July 2006.
- [19] C. Perkins, I. Kouvelas, O. Hodson, V. Hardman, M. Handley, J.C. Bolot, A. Vega-Garcia, and S. Fosse-Parisis. RTP Payload for Redundant Audio Data, Internet Engineering Task Force Request for Comments (RFC) 2198, September 1997.
- [20] S. Pfeiffer. The Ogg Encapsulation Format Version 0, Internet Engineering Task Force Request for Comments (RFC) 3533, May 2003.
- [21] E. Rescorla and N. Modadugu. Datagram Transport Layer Security, Internet Engineering Task Force Request for Comments (RFC) 4347, April 2006.
- [22] J. Rey, D. Leon, A. Miyazaki, V. Varsa, and R. Hakenberg. RTP Retransmission Payload Format, Internet Engineering Task Force Request for Comments (RFC) 4588, July 2006.
- [23] J. Rosenberg. Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols, Internet Engineering Task Force Request for Comments (RFC) 5245, April 2010.
- [24] J. Rosenberg, H. Schulzrinne, and O. Levin. A Session Initiation Protocol (SIP) Event Package for Conference State, Internet Engineering Task Force Request for Comments (RFC) 4575, August 2006.
- [25] P. Saint-Andre. Extensible Messaging and Presence Protocol (XMPP): Core, Internet Engineering Task Force Request for Comments (RFC) 6120, March 2011.
- [26] P. Saint-Andre, S. Ibarra, and E. Iovov. Interworking between the Session Initiation Protocol (SIP) and the Extensible Messaging and Presence Protocol (XMPP): Media Sessions, Internet Engineering Task Force Internet Draft, March 2014. <https://tools.ietf.org/html/draft-ietf-stox-media-04>.
- [27] H. Schulzrinne, S. Casner, R. Frederick, and V. Jacobson. RTP: A Transport Protocol for Real-Time Applications, Internet Engineering Task Force Request for Comments (RFC) 3550, July 2003.
- [28] T. Terriberry, R. Lee, and R. Giles. RTP Payload Format for Opus Speech and Audio Codec, Internet Engineering Task Force Internet Draft, February 2014. <http://tools.ietf.org/html/draft-ietf-codec-oggopus-03>.

- [29] JM. Valin, K. Vos, and T. Terriberry. Definition of the Opus Audio Codec, Internet Engineering Task Force Request for Comments (RFC) 6716, September 2012.
- [30] Ilana Volfin and Israel Cohen. Dominant speaker identification for multipoint video-conferencing. *Computer Speech Language*, 27(4):895 – 910, 2013.
- [31] S. Wenger, U. Chandra, M. Westerlund, and B. Burman. Codec Control Messages in the RTP Audio-Visual Profile with Feedback (AVPF), Internet Engineering Task Force Request for Comments (RFC) 5104, February 2008.
- [32] P. Westin, H. Lundin, M. Glover, J. Uberti, and F. Galligan. RTP Payload Format for VP8 Video, Internet Engineering Task Force Internet Draft, February 2014. <http://tools.ietf.org/html/draft-ietf-payload-vp8-11>.
- [33] P. Zimmermann, A. Johnston, and J. Callas. ZRTP: Media Path Key Agreement for Unicast Secure RTP, Internet Engineering Task Force Request for Comments (RFC) 6189, April 2011.