

**Abstract.** More and more free multi-party video conferencing applications are readily available over the Internet and both Server-to-Client (S/C) or Peer-to-Peer (P2P) technologies are used. Investigating the mechanisms, analyzing their system performance, and measuring the quality are important objectives for researchers, developers and end users. In this paper, we take four representative video conferencing applications and reveal their characteristics and different aspects of Quality of Experience. Based on our observations and analysis, we recommend to incorporate the following aspects when designing video conferencing applications: 1) Traffic load control/balancing algorithms to better use the limited bandwidth resources and to have a stable conversation; 2) Use traffic shaping policy or adaptively re-encode streams in real time to limit the overall traffic.

This work is, to our knowledge, the first measurement work to study and compare mechanisms and performance of existing *free multi-party* video conferencing systems.

## 1 Introduction

The demand for video conferencing (VC) via the Internet is growing and services are provided in two different ways: (1) either utilizing a high-end VC room system with professional equipment and dedicated bandwidth, or (2) implementing a VC application on personal computers. The first category guarantee quality, but it is costly and limited to a fixed location, while the second category is often free of charge and easy to install and use, although the quality cannot be guaranteed.

In this paper, we focus on studying *free* applications that provide multi-party ( $\geq 3$  users) VC on the Internet, and focus on the following questions:

- *How do multi-party VC applications work?*
- *How much resources do they need?*
- *What is the Quality of Experience (QoE)?*
- *What is the bottleneck in providing multi-party VC over the Internet?*
- *Which technology and architecture offer the best QoE?*

The remainder of this paper is organized as follows. Section 2 presents the work. In Section 3, eighteen popular VC applications will be introduced and classified. Section 4 describes our QoE measurement scenario. Section 5 will show the measurement results obtained. Finally, we conclude in Section 6.

## 2 Related Work

Most research focuses on designing the network architectures, mechanisms and streaming technologies for VC. In this section we only discuss the work on comparing and comparing the mechanisms and performance of streaming applications.

Skype supports multi-party audio conferencing and 2-party video conferencing. Schulzrinne [1] analyzed key Skype functions such as login, call management, media transfer and audio conferencing and showed that Skype uses a centralized P2P network to support audio conferencing service. Cicco [2] measured Skype video responsiveness to bandwidth variations. Their results indicated that Skype video calls require a minimum of 40 kbps available bandwidth to start and are able to use as much as 450 kbps. A video flow is managed by a congestion control and an adaptive codec within that bandwidth interval.

Microsoft Office Live Meeting (Professional User License) uses a S/C architecture and has the ability to schedule and manage meetings with up to 100 participants. However, only few participants can be presenters who can share their videos and the others are non-active attendees.

Spier and Ventura [3] implemented IP multimedia subsystem (IMS) for VC systems with two different architectures, S/C and P2P, and measured signaling and data traffic overhead. The results show that S/C offers better network control together with a reduction in signaling and media overhead, whereas P2P allows flexibility, but at the expense of higher overhead.

Silver [4] discussed that applications built on top of web browsers are popular in the world of Internet applications today, but are fundamentally flawed. Problems listed include delays and discontinuities, confusion and error-prone interfacing and limited functionality.

Trueb and Lammers [5] analyzed the codec performance and security. They tested High Definition (HD) VC and Standard Definition (SD) VC characteristics and their corresponding video quality. In their results, HD provides a better video quality at good and acceptable network conditions. Under poor network conditions HD and SD have similar performance.

Few articles compare the different types of existing free multi-party VC applications or measure their QoE. In this paper, our aim is to provide such a comparison.

---

<sup>1</sup> <http://www.mebeam.com/>

<sup>2</sup> <http://www.qnext.com/>

<sup>3</sup> <http://www.vsee.com/>

<sup>4</sup> <http://www.nefsis.com/leads/free-trial.aspx>

frame rate they can support (the best video quality they can provide), the maximum number of simultaneous conference participants, and the category (S/C or P2P) they belong to in Table 1.

**Table 1.** Popular video conferencing applications

	Max. frame rate (frames/second)	Max. # of simultaneous video participants	S/C or P2P
Eedo WebClass		6	web-based
IOMeeting	30	10	web-based
EarthLink	30	24	S/C
VideoLive	30	6	web-based
Himeeting	17	20	S/C
VidSoft	30	10	S/C
MegaMeeting	30	16	web-based
Smartmeeting	15	4	S/C
Webconference	15	10	web-based
Mebeam		16	web-based
Confest	30	15	S/C
CloudMeeting	30	6	S/C
Linktivity WebDemo	30	6	web-based
WebEx	30	6	web-based
Nefsis Free Trial	30	10	S/C
Lava-Lava	15	5	decentralized
Qnext		4	centralized
Vsee	30	8	decentralized

Even though there exist many free VC applications, many of them are unstable once installed. From Table 1, we observe that the maximum frame rate is 30 frames/s which corresponds to regular TV quality. All applications support only a very limited number of participants and the applications that support more than 10 simultaneous participants all use a centralized work structure.

Many other popular online chatting applications (like Skype, MSN messenger, Google talk, etc.) only support multi-party audio conference and 2-party video conference, and therefore are not considered here.

## 4 Experiments Set-Up

We have chosen four representative applications to study:

- *Mebeam*: web-browser based S/C with a single server center.
- *Qnext (version 4.0.0.46)*: centralized P2P. The node which hosts the application is the super node.

the four architectures under which all eighteen applications in Table 1 are classified.

We have performed two types of experiments: (1) local lab experiments posed of standard personal computers participating in a local video conference in order to investigate the login and call establishment process, as well as the protocol and packet distribution of the four VC applications; (2) global measurements, to learn more about the network topology, traffic load and QoS. A more realistic international video conference is carried out.

The global measurements were conducted during weekdays of May, under similar and stable conditions<sup>5</sup>:

- Client 1: 145.94.40.113; TUDelft, the Netherlands; 10/100 FastEthernet;
- Client 2: 131.180.41.29; Delft, the Netherlands; 10/100 FastEthernet;
- Client 3: 159.226.43.49; Beijing, China; 10/100 FastEthernet;
- Client 4: 124.228.71.177; Hengyang, China; ADSL 1Mbit/s.
- Client 1 always launches the video conference (as the host);
- Clients 1, 3 and 4 are behind a NAT.

To retrieve results, we used the following applications at each participant:

- *Jperf* to monitor the end-to-end available bandwidth during the process of each experiment. We observed that usually the network is quite stable and that the available end-to-end bandwidth is large enough for different applications and different participants.
- *e2eSoftVcam* to stream a stored video via a virtual camera at each participant. Each virtual camera is broadcasting in a loop a “New York” video (.avi file) with a bit rate of 910 Kbit/s, frame rate of 25 frames/s and resolution 480x270;
- *Camtasia Studio 6*. Because all applications use an embedded media player to display the Webcamera streaming content, we have to use a screen capture tool to capture the streaming content. The best software available to us was *Camtasia*, which could only capture at 10 frames/s. In order to have a fair comparison of the original video to the received video, we captured the streaming videos from all participants, but also the original stored video from the local virtual camera<sup>6</sup>.
- *Wireshark* to collect the total traffic at each participant.

---

<sup>5</sup> We have repeated the measurements in July, 2009 and obtained similar results to those obtained in May 2009.

<sup>6</sup> We assess the video quality using the full reference batch video quality metric (bVQM) which computes the quality difference of two videos. Capturing a video with frame rate of 25 frames/s may lead to a different bVQM score. However, because the video used has a stable content (there are only small changes in the person profile and background), we do not expect a significant variation in bVQM score with that of the 25 frames/s video. The results are similar for 10 frames/s videos.

*Mebeam*: We open the *Mebeam* official website to build a web video-conference room and all participants enter the room. The traces collected with *Wireshark* reveal that two computers located in the US with IP addresses 66.63.191.201 (Login Server) and 66.63.191.211 (Conference Server) are the servers of *Mebeam*. Each client first sends a request to the login server, and after getting a response, it sets up a connection with the single conferencing server center. When the client leaves the host, the host leaves from the conference room, the meeting can still continue. *Mebeam* uses TCP to transfer the signals, and RTMP<sup>7</sup> to transfer video and audio.

*Qnext*: The data captured by *Wireshark* reveals two login servers located in the US and Romania. Each client first sends packets to the login servers to join the network. After receiving a response, they use SSLv3 to set up a connection with the login servers. In the call establishment process, each client communicates encrypted handshake with 3 signaling servers located in the US and Romania and then uses SIP to set up a connection between the client and the signaling server. When client A sends a request to another client B to have a video conference and client B accepts A's request, they use UDP to transfer media data between each other. In a conference, there is one host and other clients can only communicate with the host. The host is the super node in the network. When the host leaves the meeting, the meeting ends. If another node leaves, the meeting will not be affected. *Qnext* uses SIP for signaling and UDP for video communication among participants.

*Vsee*: Each client uses UDP and TCP to communicate with the web servers and the login process. In the call establishment process, after receiving the invitation from the host, each client uses<sup>8</sup> T.38 to communicate with each other. *Vsee* has 3 web servers: during our experiment, one in the Netherlands, one in Canada, and one in the US. *Vsee* has a full-meshed P2P topology for video communication. However, only the host can invite other clients to participate in the conference. When the host leaves the meeting, the meeting cannot continue. Other clients can leave without disrupting the meeting. *Vsee* is a video-conferencing and collaboration service. The communication among users is usually of the P2P type, using UDP, with automatic tunneling through a relay if a direct connection is not available.

*Nefsis*: In the login process, the clients first use TCP and HTTP to connect to the Virtual Conference Servers (with IP addresses 128.121.149.212 in the Netherlands, 118.100.76.89 in Malaysia) and receive information about 5 other access points from the Virtual Conference Servers. These 5 access points are also Virtual Conference Servers owned by *Nefsis*, and they are located in the Netherlands (Amsterdam and Amsterdam), in the UK, India, Australia, and Singapore. After

---

<sup>7</sup> Real-Time Messaging Protocol (RTMP) is a protocol for streaming audio, video, and data over the Internet, between a Flash player and a server.

<sup>8</sup> T.38 is an ITU recommendation for fax transmission over IP networks in

set-up an end-to-end connection to communicate with each other directly. *Mebeam* uses TCP for signaling and delivering streaming data.

## 5.2 Packet Size Distribution and Traffic Load

To differentiate between non-data packets, video and audio packets formed three local experiments for each application. The first experiment used two computers with cameras and microphones to have a video conference. In the second experiment, two computers are only equipped with microphones without cameras (no video packets will be received). In the third experiment, two computers set-up a connection, both without microphones and cameras, so only non-data packets will be exchanged).

Based on *Wireshark* traces, we could distill for each VC application the packet size range as shown in Table 2:

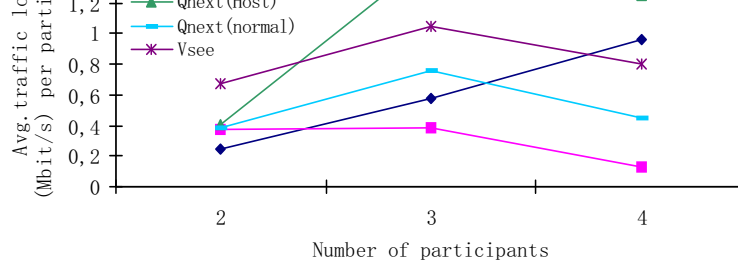
**Table 2.** The packet size distribution of *Mebeam*, *Qnext*, *Vsee* and *Nefsis*.

Packet size	<i>Mebeam</i>	<i>Qnext</i>	<i>Vsee</i>	<i>Nefsis</i>
Audio packet	> 50 bytes	72 bytes	100 ~ 200 bytes	100 ~ 200 bytes
Video packet	> 200 bytes	50 ~ 1100 bytes	500 ~ 600 bytes	1000 ~ 1500 bytes
Signaling packet	50 ~ 200 bytes	50 ~ 400 bytes	50 ~ 100 bytes	50 ~ 100 bytes

Other interesting observations are: 1) If the person profile or background in the camera change/move acutely, a traffic peak is observed in our traces. The traffic does not necessarily increase as more users join the conference. Figure 1 shows the change of the average traffic load at each user when a new participant joins the conference<sup>9</sup>. The decreasing slope after 3 users indicates that *Qnext* and *Vsee* either re-encoded the videos or used traffic shaping in order to reduce/control the overall traffic load in the system. We can see from Figure 1 that only the traffic load at *Nefsis* clients does not decrease when the number of conferencing participants reaches to 4. Therefore, we introduced more participants into the video conference for *Nefsis*, and we found that the traffic at each user starts to decrease at 5 participants. Hence, we believe that in order to support more simultaneous conference participants, the overall traffic has to be controlled.

Fig. 1 illustrates that, compared with the traffic generated by *Nefsis* with the same coding technology and the same frame rate on the same video, *Qnext* and *Vsee* generate most traffic, especially the host client of *Qnext*. This is because *Qnext* and *Vsee* use P2P architectures where the signaling traffic overhead is more than the traffic generated by a S/C network with the same number of participants. The host client (super node) of *Qnext* generates 3 times more traffic than the other clients.

<sup>9</sup> We captured the packets after the meeting was set up and became stable.



**Fig. 1.** The average traffic load at an end-user when the number of conference participants increases from 2 to 4 (*Qnext* is limited to 4 participants)

other normal clients. Hence, for this architecture, a super-node selection is recommended to choose a suitable peer (with more resources, for example, a super node).

Fig. 1 also shows that *Mebeam* generates least traffic. Considering that all traffic load, which can be supported in a VC network, has an upper bound due to the limited users' bandwidth, and each *Mebeam* client generates much less traffic than the three other applications, it clarifies why *Mebeam* can support more simultaneous video users while *Nefsis* can only support 10 users, *Vsee* can support 10 users and *Qnext* can support 4 users.

### 5.3 Quality of Experience (QoE)

QoE can be measured through objective and subjective measurements. In this section, we assess via global measurements the QoE at the end user with respect to their video quality, audio-video synchronization, and level of interaction.

**Video Quality.** In the objective measurements, we use bVQM (Bandwidth Quality Metric) to analyze the VC's video quality off-line. bVQM takes the original video and the received video and produces quality scores that reflect the predicted fidelity of the impaired video with reference to its undistorted counterpart. The sampled video needs to be calibrated. The calibration consists of estimating and correcting the spatial and temporal shift of the processed video sequence with respect to the original video sequence. The final score is computed using a linear combination of parameters that describe perceptual changes in video quality by comparing features extracted from the processed video with those from the original video. The bVQM score scales from 0 to approximately 1, where a smaller score, the better the video quality.

<sup>10</sup> According to [7], bVQM scores may occasionally exceed 1 for video scenes that are extremely distorted.

**Table 3.** The video quality of *Mebeam*, *Qnext*, *Vsee* and *Nefsis* at 4 clients

VQM score	Client 1	Client 2	Client 3	Client 4	Average
<i>Mebeam</i> (Flash video, MPEG-4)	0.63	0.41	0.94	0.86	0.71
<i>Qnext</i> (MPEG-4, H.263, H.261)	1.05	0.94	0.63	0.83	0.86
<i>Vsee</i>	0.78	0.82	0.80	0.79	0.80
<i>Nefsis</i> (MPEG-4, H.263, H.263+)	0.34	0.61	0.61	0.87	0.61

Table 3 indicates that *Nefsis* features the best video quality among applications, although with an average bVQM score of 0.61 (its quality is on the lower end, which will be explained later with the subjective measurements). The highest bVQM score (the worst video quality) appears at Client 1 (the super node) of *Qnext*. Generally speaking, all four VC applications do not provide good quality<sup>12</sup>.

Because no standard has been set for what level of bVQM score corresponds to what level of perceived quality of a VC service, we have also conducted subjective measurements. We use the average Mean Opinion Score (MOS) [8], a metric for user perceived quality, defined on a five-point scale<sup>13</sup>: 5 = *excellent*, 4 = *good*, 3 = *fair*, 2 = *poor*, 1 = *bad*.

We gave 7 different quality videos generated by VC applications to 20 participants, who gave a subjective MOS score independently. We also objectively measured their bVQM scores. Fig. 2 shows the correlation between the objective bVQM scores and the subjective MOS values.

We mapped between the bVQM scores and the average MOS score for 20 persons, and found that they have a linear correlation in the range 0.3 to 0.5. For a score  $\leq 1$ . Hence, the VC's video quality is predictable when using the objective metric bVQM.

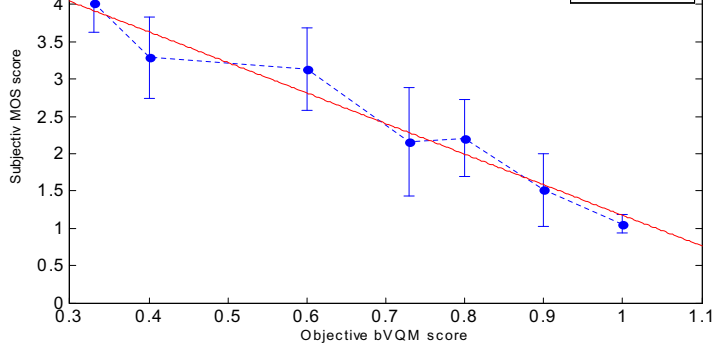
Compared with the video quality of a global P2PTV distribution service like P2P TV, which has an average MOS value of 4 [9], the video quality of a global VC service is poor (with an average bVQM score of 0.74 and MOS value of around 2.2), because the VC service requires end users to encode and upload their streams in real time.

<sup>11</sup> *VirtualDub* is a video capture and video processing utility for Microsoft Windows.

<sup>12</sup> We also objectively measured the audio quality using metric PESQ-LQ (Perceptual Evaluation of Speech Quality-Listening Quality) [6] [8] and found that the average score (scale from 1.0 to 4.5, where 4.5 represents an excellent audio quality) is 2.24, 2.68, 3.08 and 3.15 for *Mebeam*, *Qnext*, *Vsee*, and *Nefsis*, respectively.

<sup>13</sup> The threshold for acceptable TV quality corresponds to the MOS value 3.5.





**Fig. 2.** Relation between bVQM and MOS for video conferencing services.

Even the local uploaded video has a largely degraded quality although it is the best among all participants.

**Audio-Video Synchronization.** The relative timing of sound and images in the presentations of a streaming content may not be synchronized.

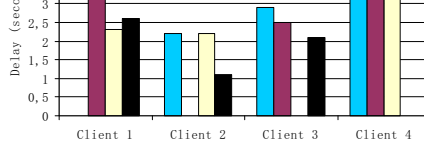
ITU [10] [11] has suggested that the viewer detection thresholds of audio-video synchronization lag are about +45 ms to −125 ms, and the acceptance thresholds are about +125 ms to −185 ms, for video broadcasting.

To analyze the A/V synchronization provided by each VC application, we generated an “artificially generated” video test sample, in which the video and audio tracks are temporally synchronized with markers. Similar to the experiments for testing the video quality, we captured at each end user the videos from the participants. When the audio and video tracks were extracted and compared side-by-side, there was an average difference in time between the two tracks of about 400 ms for *Mebeam*, 470 ms for *Qnext*, 400 ms for *Vsee* and 350 ms for *NetMeeting*. These large audio-video lags are mainly caused by a large amount of frame loss in the video stream, which lead to the low video quality mentioned already in Section 5.3.

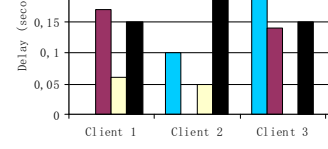
**Interactivity (communication delay).** During a video conference it is often required for participants to have large communication delay<sup>14</sup>. Large communication delay impacts the quality of real-time interactivity in our global multi-party VC experiments. We evaluated the video delays among participants by injecting in the network another video that mainly reproduced a timer with millisecond granularity.

In the video conference, this artificial “timer” video was uploaded via the local camera and transmitted among the participants via the different VC applications.

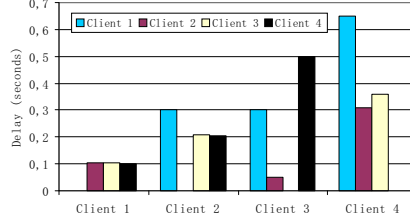
<sup>14</sup> In IP video conferencing scenarios, the maximum communication delay recommended by ITU is 400 ms [12].



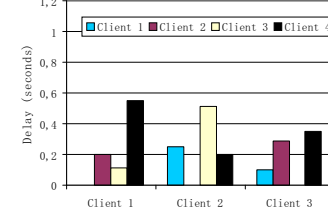
(a) Mebeam



(b) Qnext



(c) Vsee



(d) Nefsis

**Fig. 3.** The video delay between different participants

At each participant, we used a standard universal Internet time as reference. We displayed the “timer” videos of all participants in real time. After a 1-minute stable video conference, we cut the captured content at each participant and compared it with the content of the *RealDub* to compare the “timers” between any 2 participants. For each pair of participants, we took samples at 2 different times to calculate an average delay.

The video delays among participants are shown in Fig. 3. The  $x$  axis shows the 4 different clients. The  $y$  axis shows the video transmission delay from the participant on the  $x$  axis to the participant shown in the legend.

Fig. 3 shows that *Qnext* provides a video that is most synchronized among the clients. *Qnext*, *Vsee*, and *Nefsis* have a comparable level of average video delay, respectively 0.15 s, 0.27 s, and 0.41 s. However, *Mebeam* clients suffer a huge delay (2.77 s on average), because the processing time at the server is too long.

We also measured the audio delays among participants by injecting into the video an artificial DTMF (Dual-tone multi-frequency) tone. We sent and received audio at Client 1. Other participants kept their speaker and microphone muted and did not produce extra audio. Based on the recorded audio tracks, we compared the time the audio marker was sent from Client 1 and the time the same audio marker was heard again at Client 1 after the transmitted audio was played, recorded, and retransmitted by a client. The time difference is approximately twice the audio delay plus the processing delay at a client. Our results revealed

<sup>15</sup> <http://www.time.gov/timezone.cgi?Eastern/d/-5/java>

to-end delay including the transmission delay and the delay introduced by the application. In our experiment, the video delay represents the delay of a scene that was captured at the application interfaces of the sender and receiver, which does not include the time used for uploading the video to the server via applications. Hence, considering the audio delay, video delay, and the synchronization discussed in Section 5.3, we can conclude that the delay introduced by the application, when uploading, is large for *Qnext*.

## 6 Worst-Case Study

In another set of global experiments in June, 2009, our *Jperf* plots indicate that the end-to-end connections of clients 3 and 4 with the host were very unstable. We found that the two participants in China always passively disconnected from the conference or could not even log into *Mebeam*, *Nefsis* and *Qnext*. *Vsee* could still work, but the quality was awful, with bVQM scores close to 1.

In order to investigate the minimum bandwidth to support a video conference, we repeated many experiments adjusting the upload rate upper-bound (by *Netlimiter*) at each participant for a particular VC application to test the minimum upload bandwidth minimally required to launch a video conference.

For *Qnext*, the threshold is 50 Kbit/s. If an end user's available upload bandwidth is < 50 Kbit/s, (s)he cannot launch *Qnext*. For *Vsee*, the threshold is 10 Kbit/s; for *Nefsis* it is 28 Kbit/s; and for *Mebeam* it is 5 Kbit/s, which values are the minimally supported streaming bit rates set by the applications.

## 7 Summary and Conclusions

Through a series of local and global experiments with four representative video conferencing systems, we examined their behavior and mechanisms, analyzed their login process, the call establishment process, the packet scheduling, transfer protocols, traffic load, delivery topology, and different aspects of Quality of Experience.

Our main conclusions from the measurement results on the traffic characteristics of four different video conferencing systems are: (1) The QoE of multi-party video conferencing is very sensitive to bandwidth fluctuations, especially on the bottleneck link. Hence, an adaptive bit rate/frame rate policy should be deployed; (2) as the number of participants increases, the traffic load at each participant always increase correspondingly (see Fig. 1), suggesting that re-encoding of video or a traffic shaping policy take place to control the overall traffic of the system.

---

<sup>16</sup> *NetLimiter* is an Internet traffic control and monitoring tool designed for controlling download/upload transfer rate limits for applications.

video and audio quality, large audio-video lag, and long communication times in some cases); (2) Only a limited number of multimedia participants are supported and rare high definition webcam streaming is supported due to the limited available bandwidth or the limited processing capability; (3) The existing systems are not reliable in the worst cases. When the network is unstable or the available upload bandwidth is very limited (thresholds have been found), none of the systems work properly.

It seems that the Server-to-Client architecture with many servers located all over the world is currently the best architecture for providing video conferencing via the Internet, because it introduces the least congestion at both servers and clients. Load balancing and load control algorithms help the overall performance of the system and the codec used is important for the quality that end users receive. The bottleneck to support video conferencing with more participants is high definition streams is the overhead traffic generated by them. To support simultaneous participants in a single conferencing session, the traffic load must be controlled/limited by using traffic shaping policy or re-encoding the video streams.

We have chosen four representative video conferencing systems for comparison, but the measurement methodologies mentioned in this paper can also be applied to other video conferencing applications, which could be compared with the results in the future.

## Acknowledgements

We would like to thank Rob Kooij for a fruitful discussion on measurement of delay. This work has been supported by TRANS (<http://www.trans-res.nl>).

## References

1. Baset, S.A., Schulzrinne, H.: An Analysis of the Skype Peer-to-Peer Interference Protocol. In: INFOCOM '06, Barcelona, Spain (April 2006)
2. De Cicco, L., Mascolo, S., Palmisano, V.: Skype Video Responsiveness to Bandwidth Variations. In: NOSSDAV '08, Braunschweig, Germany (May 2008)
3. Spiers, R., Ventura, N.: An Evaluation of Architectures for IMS Based Video Conferencing, Technical Report of University of Cape Town (2008)
4. Silver, M.S.: Browser-based applications: popular but flawed? *Information Systems and E-Business Management* 4(4) (October 2006)
5. Trueb, G., Lammers, S., Calyam, P.: High Definition Videoconferencing: Confidentiality, Performance, Security, and Collaboration Tools, REU Report, Ohio Supercomputer Center, USA (2007)
6. Rix, A.W.: A new PESQ-LQ scale to assist comparison between P.862 PESQ and subjective MOS, ITU-T SG12 COM12-D86 (May 2003)
7. Pinson, M.H., Wolf, S.: A New Standardized Method for Objectively Measuring Video Quality. *IEEE Transactions on Broadcasting* 50(3), 312–322 (2004)

- Experience of SopCast. *International Journal of Internet Protocol Technol* 11–23 (2009)
10. ITU BT.1359-1, Relative timing of sound and vision for broadcasting (1990)
  11. Lias, J.L.: HDMI's Lip Sync and audio-video synchronization for broadcast home video, Simplay Labs, LLC (August 2008)
  12. Bartoli, I., Iacovoni, G., Ubaldi, F.: A synchronization control scheme for streaming services. *Journal of multimedia* 2(4) (August 2007)