

AN OPTIMIZED H.264-BASED VIDEO CONFERENCING SOFTWARE FOR MOBILE DEVICES

Hans L. Cycon*, Thomas C. Schmidt, Gabriel Hege†, Matthias Wählisch‡, Mark Palkow§

{h.cycon,hege}@fhtw-berlin.de, {t.schmidt,waehlich}@ieee.org, palkow@daviko.com

ABSTRACT

Mobile phones and related gadgets in networks are on the spot to deliver sufficient performance for rich multimedia applications and communication. In this report we introduce a video conferencing software, which seamlessly integrates mobile with stationary users into fully distributed multi-party conversations. Innovations related to this work are twofold. At first we report on a highly optimized realisation of a H.264 codec and our implementation experiences of the video conference *software* on a consumer mobile. Within the tight bounds of real-time requirements on mobiles, the coding software outperforms compatible H.264 realizations. At second we present an integrated peer-to-peer group communication solution, which scales well for medium-size conferences and accounts for the heterogeneous nature of mobile and stationary participants.

Index Terms— Mobile video coding, H.264/MPEG-4 AVC software codec, mobile conferencing, peer-to-peer group communication, distributed SIP conference management

1. INTRODUCTION

The idea of augmenting voice calls by video has been around for several decades, but only the flexibility of the Internet generated a noticeable deployment. As compared to audio, video processing places significantly higher demands on end system and network transmission capabilities. The rapid evolution of networks and processors have paved the way for realistic group conferences conducted at standard personal computers, combining about a dozen visual streams of Half-QVGA (240 x 160 pixel @ 15-30 fps) resolution.

Mobile phones and networked consumer portables are now on the spot to deliver sufficient performance for rich multimedia applications and communication, as well. Videoconferencing though, which requires simultaneous decoding and encoding in real-time, poses still a grand challenge to the mobile world. Limited and expensive wireless channels on the

one hand, high consumer demands on visual quality on the other, advise applications to take advantage of the latest standard for video coding H.264/AVC [1].

H.264/AVC provides gains in compression efficiency of up to 50 % over a wide range of bit rates and video resolutions compared to previous standards. While H.264/AVC decoding software has been successfully deployed on handhelds, high computational complexity still prevented pure software encoders in current mobile systems. There are however also fast hardware implementations available, which give rise to an increasing offer of device- and operator-bound video services.

In this work we first introduce a pure software solution for real-time video communication on standard smartphones in section 2. These mobile clients extend a lightweight, feature rich conferencing application developed for an infrastructure compliant use on standard PCs. In the second part we present the underlying peer-to-peer group communication scheme, which performs well for medium-size conferences and accounts for the heterogeneous nature of mobile and stationary participants, cf. section 3. This includes on the one hand SIP [2] standard compliant session signalling with respect to group communication, and on the other hand efficient, serverless media distribution, self-adjusting to the actual network infrastructure support. Conclusions and an outlook follow in the final section.

2. THE DAVIKO VIDEOCONFERENCING SOFTWARE

In this section we give an overview of our reference implementation, a digital audio-visual conferencing system, realised as a serverless multipoint video conferencing software without MCU developed by the authors [3]. It has been designed in a peer-to-peer model as a lightweight Internet conferencing tool aimed at email-like friendliness of use. The system is built upon a fast H.264/MPEG-4 AVC standard conformal video codec implementation [4] called DAVC. By controlling the coding parameters appropriately, the software permits scaling in bit rate from 48 to 1440 kbit/s on the fly.

Audio data is compressed using a 16 kHz speech-optimized variable bit rate codec [5] with extremely short latencies of 40 ms (plus network packet delay). All streams can be transmitted by unicast as well as by multicast protocols. Within

*The author is with FHTW Berlin, 10318 Berlin, Germany.

†Thomas and Gabriel are with HAW Hamburg, Dept. Informatik, Berliner Tor 7, 20099 Hamburg, Germany.

‡Matthias is with link-lab, Hönow Str. 35, 10318 Berlin, Germany and also with HAW Hamburg.

§Mark is with daViKo GmbH, Am Borsigturm 40, 13507 Berlin, Germany.

the application, audio streams are prioritized over video since user experience is usually more sensitive to losses in audio packets than those of video packets, which both may result from transmission errors or network congestions.

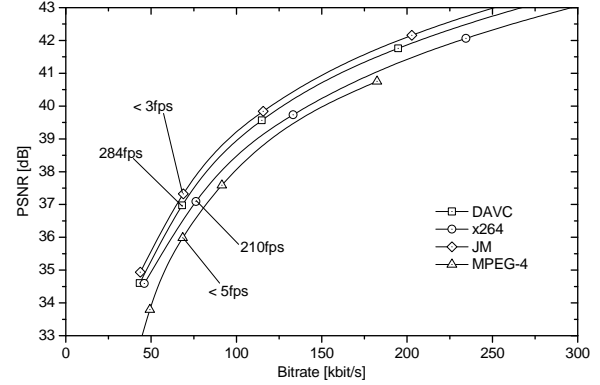
An application-sharing facility is included for collaboration and teleteaching. It enables participants to share or broadcast not only static documents, but also any selected dynamic PC actions like animations. All audio/video - streams including dynamic application sharing actions can be recorded on any site. This system is equally well suited to intranet and wireless video conferencing on a best effort basis, since the audio/video quality can be controlled to adapt the data stream to the available bandwidth.

The daViKo conferencing system is available for desktop computers running MS-Windows or Linux and on handhelds with Windows Mobile operating system.

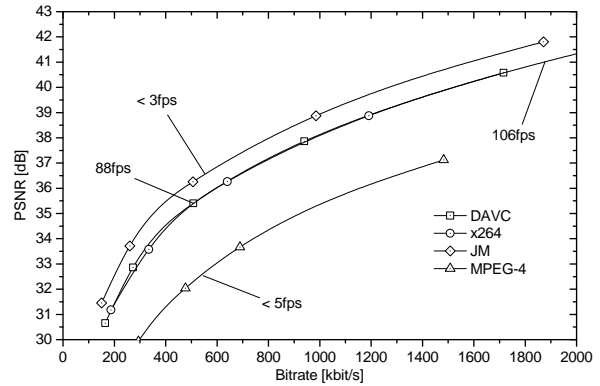
2.1. The DAVC Codec

DAVC, the core of the videoconferencing system, is a fast, highly optimized H.264/MPEG-4 AVC standard implementation. It realizes a Baseline profile, optimized for real-time encoding (as well as real-time decoding) by means of a fast motion estimation strategy including integer-pel diamond search as well as a fast subpel refinement strategy up to $\frac{1}{4}$ pel motion accuracy. Motion estimation includes the choice of several different macroblock (MB) partitions and multiple reference frames, as permitted by the H.264/MPEG-4 AVC standard. For choosing between different MB partitions for motion-compensated (i.e. temporal) prediction and MB-based intra (i.e. spatial) prediction modes, a fast rate-distortion (RD) based mode decision algorithm with early termination conditions has been employed.

In comparison to the well-known open source H.264/MPEG-4 AVC encoder implementation of x264 [6], our DAVC encoder implementation achieves up to 0.5 dB PSNR better RD performance and a considerable increase in encoding speed when using comparable encoder settings. For selected RD points we measured 284 encoded frames per second (fps) as compared to 210 fps for x264. In Figure 1, typical examples of such a comparison between x264 and DAVC are shown. In addition to the RD-performance of those two real-time encoder implementations, this plot also shows the RD behavior of two non real-time encoder implementations, as given by the H.264/MPEG-4 AVC Joint Model (JM) reference software (with Baseline profile settings) and a MPEG-4 (Part 2) Advanced Simple Profile implementation. The latter two encoders were operated using a high-complexity RD-based mode decision strategy for demonstrating the capabilities of both video coding standards when neglecting any real-time constraints. Figures 1(a) and 1(b) also contain the number of encoded frames per second (fps) for selected RD points as a measure for maximum encoding speed. Similar results were also achieved for other test sequences.



(a) Akiyo (cif, 300 frames at 30 Hz)



(b) Foreman (cif, 300 frames at 30 Hz)

Fig. 1. RD plot for test sequences in CIF resolution comparing three different H.264/MPEG-4 AVC encoder implementations as well as a RD-optimized MPEG-4 (Part 2) Advanced Simple Profile implementation.

Note that the DAVC codec along with the H.264/AVC design also includes some suitable mechanisms to quickly recover from video packet loss.

2.2. Mobile Video Codec Performance

In ongoing work, the DAVC codec has been adapted to sustain real-time performance on mobile devices. The mobile codec version operates at reduced complexity for motion compensation with a highly optimized code base for the target platform. Motion compensation has been limited to work on 16 x 16 pixel blocks, only. The code tuning includes the efficient use of the wireless MMX instruction set available at the target system. Portability is sustained by an ANSI compliant C version, to be augmented incrementally by platform specific injections.

The application was tested on a 520 MHz Xscale processor built in an Asus P735 system. Thereon it can reliably encode and decode a QCIF video stream in parallel at 5/10 fps, without CPU exhaustion or frame dropping. Real-time encoding rate increases up to 10 fps for moderate video com-

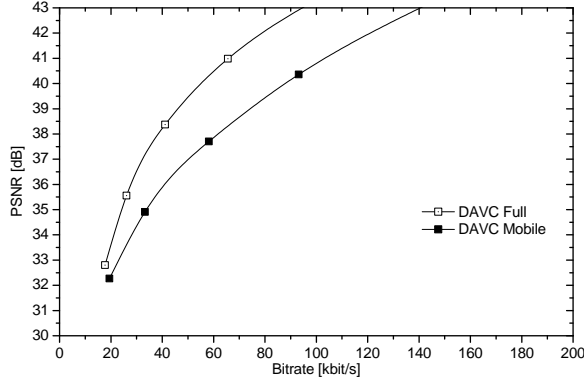


Fig. 2. RD plot for test sequence “Akiyo” in QCIF resolution at 10 fps, comparing the DAVC mobile encoder to the DAVC fully optimized implementation.



Fig. 3. The mobile video application.

plexity. QCIF @ 15 fps is the maximal image feed that can be obtained from the front camera in our test equipment. Performance values for the mobile encoder are displayed in figure 2 and compared to the results for the full DAVC.

Reduced coding complexity results in an enhanced data rate send by the mobile, but the gross total rate for a bidirectional video exchange at 10 fps complies to 3GPP/UMTS bandwidths constraints. Note that experimental conditions are not fully compatible: The image sequence obtained from the front camera of the mobile is significantly more noisy than our standard USB cameras connected to the desktop, which increases the image complexity and thereby the data rate.

3. DISTRIBUTED POINT-TO-MULTIPOINT CONFERENCING

Our application aims for simple, flexible, and cost-efficient ad-hoc conferencing functions, which scale appropriately well, but avoid any infrastructure assistance. Such a solution re-

quires group session management and media distribution at peers, which for the sake of standard compliance we realize with group conferencing functions in SIP, cf. [7, 8, 9]. Implemented as pure software on standard personal devices, user agent peers are exposed to severe restrictions in real-world deployments: Often they are located behind NATs and firewalls with network capacities confined to asymmetric DSL or wireless links. Capacity constraints and resilience to node failures require peer-managed ad-hoc conferences to organize in a distributed multi-party model. As a key component, the heterogeneity of clients must be accounted for, whereas the range of scalability is limited to about a dozen parties in videoconferences.

3.1. P2P Adaptive Architecture

A peer-to-peer conferencing system faces the grand challenge to be robust w.r.t. the infrastructure. The role a user agent is able to attain in a distributed scenario needs to be adaptively determined according to constraints of its device and current network attachment. In a simplified scenario, clients may be divided into two groups, distinguished by their ability to act as a SIP conference focus or not. A focus must be globally addressable and have access to necessary processing and network resources.

This elementary adaptation scheme can be based on individual decisions of user agents and gives rise to a hybrid architecture of super peers, chosen from potential focus nodes, and remaining leaf nodes. To decide on its potential role of building a focus, a client at first needs to determine NATs and firewalls. Aside from address evaluation, this is done by a simple probe packet exchange. As the implementation is CPU-type aware, processing restrictions are easily evaluated, as well. However, an a priori judgement on available network bandwidth is not easily obtained. An evaluation of the local link capacity is frequently misleading, as wireless devices may be located behind wired transmitters of lower, asymmetric capacity such as in ADSL. Current experiments to quickly retrieve reasonable estimates of up- and downstream capacity are ongoing on the basis of variable packet size, nonintrusive estimators, cf. [10]. Note that network capacity detection is of vital use for temporal adaptation of the video codecs, as well.

Leaf nodes attach to super peers in subordinate position, whereas potential focus nodes may be assigned to be super peers or leaves. Super peers provide global connectivity among each other and NAT traversal assistance to leaves, while leaf nodes experience super peers in different roles: A leaf nodes sees its next hop super peer as the conference focus, while the remote super peers act as proxies on the path to the leaves behind.¹ This set-up corresponds to the well known architecture

¹This architecture relies on the presence of at least one globally addressable, sufficiently powerful peer. As there are many scenarios, where this is likely to fail, we advise for and offer a permanently deployed ‘silent’ relay-

of Gnutella 0.6 and successive hybrid unstructured peer-to-peer systems, cf. [11]. Despite its architectural analogy, our routing layer for real-time group applications follows a different, next-hop design.

4. CONCLUSIONS & OUTLOOK

We have presented a peer-to-peer software for high-quality videoconferencing on mobiles, admitting utmost flexibility with respect to end systems, operators and network provisioning. To the best of our knowledge, this is the first software implementation of an H.264 video encoder that operates in real-time on mobile phones. An adaptive, fully distributed conference management scheme with SIP has been developed as part of the multi-party scenario. This hybrid peer-to-peer model accounts for client capabilities as well as network attachment, and does scale well beyond standard use.

In future work we will concentrate on further optimization and generalization of the video coding software to make it available for a wider variety of platforms. Network adaptation and capacity evaluation will require further work to arrive at estimates that reliably serve the needs in real world environments, as well.

Additional research will target at benefits possibly inherited from key-based routing. As common application layer multicast schemes, which rely on dedicated shared or source specific trees, are significantly sensitive to client departure and of insufficient performance in medium size groups, and as conference routing actually can be seen as an application layer broadcasting problem, new and highly optimized structured broadcast algorithms are desirable. The bidirectional shared tree approach introduced in [12] may be a promising point to start at.

Acknowledgement

This work is supported by the German Bundesministerium für Bildung und Forschung within the project *Moviecast* (<http://moviecast.realmv6.org>).

5. REFERENCES

- [1] ITU-T Recommendation H.264 & ISO/IEC 14496-10 AVC, "Advanced Video Coding for Generic Audiovisual Services," ITU, Tech. Rep., 2005, draft Version 3.
- [2] J. Rosenberg, H. Schulzrinne, G. Camarillo, A. Johnston, J. Peterson, R. Sparks, M. Handley, and E. Schooler, "SIP: Session Initiation Protocol," IETF, RFC 3261, June 2002.
- [3] M. Palkow, "The daViKo homepage," 2008, <http://www.daviko.com>.
- [4] J. Ostermann, J. Bormans, P. List, D. Marpe, N. Narroschke, F. Pereira, T. Stockhammer, and T. Wedi, "Video Coding with H.264/AVC: Tools, Performance and Complexity," *IEEE Circuits and Systems Magazine*, vol. 4, no. 1, pp. 7–28, April 2004.
- [5] "The Speex projectpage," <http://www.speex.org>, 2007.
- [6] "VideoLan: x264 - a free h264/avc encoder," <http://www.videolan.org/developers/x264.html>, 2007.
- [7] A. Johnston and O. Levin, "Session Initiation Protocol (SIP) Call Control - Conferencing for User Agents," IETF, RFC 4579, August 2006.
- [8] R. Mahy, R. Sparks, J. Rosenberg, D. Petrie, and A. Johnston, "A Call Control and Multi-party usage framework for the Session Initiation Protocol (SIP)," IETF, Internet Draft - work in progress 9, November 2007.
- [9] T. C. Schmidt and M. Wählisch, "Group Conference Management with SIP," in *SIP Handbook: Services, Technologies, and Security*, S. Ahson and M. Ilyas, Eds. Boca Raton, FL, USA: CRC Press, November 2008, to appear, on invitation.
- [10] R. Prasad, C. Dovrolis, M. Murray, and K. Claffy, "Bandwidth Estimation: Metrics, Measurement Techniques, and Tools," *IEEE Network*, vol. 17, no. 6, pp. 27–35, November–December 2003.
- [11] R. Steinmetz and K. Wehrle, Eds., *Peer-to-Peer Systems and Applications*, ser. LNCS. Berlin Heidelberg: Springer-Verlag, 2005, vol. 3485.
- [12] M. Wählisch and T. C. Schmidt, "Between Underlay and Overlay: On Deployable, Efficient, Mobility-agnostic Group Communication Services," *Internet Research*, vol. 17, no. 5, pp. 519–534, November 2007.