

VOLUME I

INET'95 Conference Proceedings

Editor: Kilnam Chon

**International Networking Conference
Honolulu, Hawaii, USA
27-30 June 1995**

The Annual Conference of the Internet Society

To order copies, contact the ISOC Secretariat at:
Internet Society
12020 Sunrise Valley Dr., Suite 270
Reston, VA 22091
USA
isoc@isco.org

RPX Exhibit 1222

Reliable Audio for Use over the Internet

Vicky Hardman <v.hardman@cs.ucl.ac.uk>, Martina Angela Sasse <a.sasse@cs.ucl.ac.uk>, Mark Handley <m.handley@cs.ucl.ac.uk>, Anna Watson <a.watson@cs.ucl.ac.uk>

Abstract

This paper describes current problems found with audio applications over the MBONE (Multicast Backbone), and investigates possible solutions to the most common one - packet loss. The principles of packet speech systems are discussed, and how the structure allows the use of redundancy to design viable solutions to the problem. The paper proposes the use of synthetic speech coding algorithms (vocoders) to provide redundancy, since the algorithms produce a very low bit-rate stream, which only adds a small overhead to a packet. Preliminary experiments show that normal speech repaired with synthetic quality speech is intelligible, even at very high loss rates.

Introduction

The application of this work is multimedia conferencing over the MBONE (Multicast Backbone), an experimental overlay network of the Internet. The work has arisen from experiences in multi-way multimedia conferencing in Project MICE (Multimedia Integrated Conferencing for Europe) [1], is currently applied in Project ReLaTe (Remote Language Teaching over SuperJANET) [2], and includes formal experiments into the human perception of packet speech systems degraded by packet loss.

If multimedia conferencing is to become widely used in the Internet community, user must perceive the quality to be sufficiently good for most applications. Experience has shown that audio is almost always the most important component of multimedia conferencing. Whilst we have identified a number of problems which impair the quality of audio, the major one with audio over the MBONE is packet loss [3]. This paper attempts addresses the problem of packet loss over the MBONE.

Packet loss can occur for a number of reasons:

- congestion of routers and gateways, which lead to packet being discarded;
- delays in packet transmission, with packet arriving too late at the receiver to be played back;
- heavy loading of the workstations, leading to scheduling difficulties in multi-tasking operating systems.

Packet loss is a persistent problem, particularly given the increasing popularity, and therefore increasing lead, of the Internet. Possible ways of combatting congestion include bandwidth reservation and moves toward an integrated service management on the Internet. These would require wide-scale changed

to be agreed and implemented, so these solutions will be available in the short to medium term. Yet, the disruption of speech intelligibility even at low loss rates which we currently experience may convince a whole generation of users that multimedia conferencing over the Internet is not viable. We therefore propose a solution which renders the speech intelligible under current network conditions, and can be deployed in the short term. Such a solution will have to be at the application level, i.e. the multicast audio tools.

Current audio applications repair lost packets with silence, which leads to the speech clipping effects currently experienced by many users. Since comparatively large packets are used, even the loss of individual packet loss has a serious impact on the intelligibility of speech.

We propose a method of repairing damaged speech using cheap redundancy within the packets sent from the transmitter. The redundancy is synthetic speech, which, when split into packets, only adds a very small amount of overhead, and therefore does not add to the congestion at the network level. The redundancy for any given packet of speech is piggy-backed onto a later packet. This mechanism means that when the receiver suffers the loss of the primary speech information, it still has the possibility of substituting something sensible in the output stream of speech, provided that the redundancy can be received.

In order to establish the effectiveness of this solution, we have performed experiments into user perception of speech repaired with a synthetic substitute. The experiments subjectively measured speech intelligibility, and the results show that this technique is very successful at repairing speech with large packet sizes and for very high loss rates (results were taken up to 40%). The paper also describes how the proposed solution scales in the multicast environment.

Background

Speech Coding

Speech coding schemes have been standardised for use over telephone networks; a variety of speech coding algorithms exist for a single target quality of service (QoS), and at a very few discrete bit-rates: Pulse Code Modulation (PCM) operates at 64 kbps, Adaptive Differential Pulse Code Modulation (ADPCM) operates at 32 kbps, and Code Excited Linear Prediction (LD-CELP) operates at 16 kbps. The target QoS is 'toll' (or telephone) quality, and each algorithm available at this QoS produces a different bit-rate: the improvement in bit-rate being obtained for

increasing complexity in the coding algorithms. Another standard speech coding algorithm is Groupe Speciale Mobile (GSM), which was designed for use over cellular telephone networks. The target QoS is consequently slightly less than toll, but the algorithm is popular, since it operates at the same bit-rate as CELP, but is much less complex. A fuller discussion of toll quality speech coding algorithms can be found in [4].

A second class of coding algorithms exist, which operate at the 'communications' or synthetic QoS. These algorithms operate at very low bit-rates (approx. 4.8kbps and below), and produce very mechanical sounding speech. Perhaps the most important method of this class is Linear Predictive Coding (LPC), since the principle is also an integral part of both the CELP and GSM coders. A fuller description of which can be found in [5].

Packet Speech Systems

Packet speech systems usually employ the standard speech coding algorithms, and group the emerging stream of codewords into packets for transmission over the network. At the receiver, the packets may be delivered: out of order, not at all, or at non-uniform intervals. Consequently, a reconstruction delay must be used at the receiver to repair the network effects; this enables sample play-out to be smoothed.

In a packet speech system, the end-to-end delay is always a critical factor in the usability of a real-time voice system, and should be kept below 600ms in the absence of echoes (The figure may be in fact be less than this - 400ms) [6], if conversation patterns are not to break down. The size of the packets (in ms) chosen for a packet speech system directly impacts the end-to-end delay. A delay equal to the size of one packet is incurred at the transmitter, since the samples in the packet have to be collected before a packet can be sent. At the receiver, a rough estimate of the reconstruction delay required to smooth out packet arrival times is two packets worth in ms [7] [8], although the true value may be substantially in excess of this rule of thumb. Consequently, a minimum of three packets worth of delay is incurred on an end-to-end basis, before the network propagation delay has been taken into account.

The delay introduced will be enough to receive most of the packets, but some will always arrive too late to be played back, and can be considered 'lost'. Furthermore, the network itself may lose packets. In such situations, the speech 'stream' must be repaired, and a dummy packet inserted in place of the lost one, so that the correct timing relationship is maintained between the transmitter and receiver. The presence of the dummy packet is usually discernible to the listener, and unfortunately, the perceptibility of the loss increases with increasing packet size, as well as with increasing loss rate.

The impact of the two factors identified above, (delay and loss), is such that small packet sizes are required for real-time voice links. However, the use of

small packets increases the overhead of packet headers, and any processing incurred at network nodes, and therefore increases the likelihood of congestion and loss. Consequently, a trade-off exists between the requirements of the network, and the requirements of real-time voice connections.

Voice Reconstruction Techniques

Repair methods for packet loss are known as voice reconstruction mechanisms. The aim is to construct a suitable dummy packet at the receiver, so that the loss is as imperceptible as possible. With compressed speech, voice reconstruction mechanisms not only have to produce a suitable fill-in packet, but also have to maintain the decoder tracking, since the algorithms transmit difference information. Voice reconstruction techniques can be split into two categories; receiver only, and combined source and channel techniques.

Table 1: Voice Reconstruction Techniques

Receiver-Only	Combined Source and Channel
Silence	Embedded Coding
White Noise	Redundancy
Waveform Substitution	
Sample Interpolation	

Receiver-only techniques are those that try to reconstruct the missing segment of speech solely at the receiver, possibly from correctly received packets preceding that which was lost. Combined source and channel techniques are those that try to make the system robust to loss by either arranging for the transmitter to code the speech in such a way as to be robust to packet loss, or by transmitting extra information to help with reconstruction.

Receiver-Only Techniques

The original voice reconstruction techniques were receiver-only, and used either silence, white noise, repetition of part of the last correctly received speech waveform, or sample interpolation as the substitute.

Silence substitution is favoured because it is simple to implement, and it gives adequate performance for small packet sizes (<16ms), and up to 1% loss [9] [10].

It is well known that other methods give substantially better results than those obtained from silence substitution. Warren [11] investigated the human perception of speech interrupted by silence compared to noises, such as coughs. The results show that phonemic restoration (the ability of the human brain to subconsciously repair the missing segment of speech with the correct sound) occurs for the noise situation,

and does not occur for silence substitution.

Experience from the MICE project has shown that listeners can become very frustrated with MBONE speech. Their frustration stems from a variety of audio problems:

- interrupted speech due to packet loss - the speech sounds at first bubbly, and then individual interruptions are apparent;
- talking into a 'dead' channel, and the lack of automatic gain control for listeners;
- high levels of background noise cutting in and out;
- distortion due to overloading the microphone when speakers shout, or mis-matched levels.

Packet loss is the most frustrating problem, and one users cannot cure of their own accord. The frustration with interrupted speech can be explained by considering the linguistic construct of a sentence, which includes a pause (of duration > a phoneme) at the end of the sentence. Since the size of packets used over the MBONE are often comparable to the length of a phoneme, the interruptions in the speech flow sometimes occur at inappropriate points, which sends ambiguous signals to brain, as to whether speech is continuing or not [11].

White noise was shown to give a subjective performance improvement over silence by Miller [12] when contextual information in speech was removed, and an intelligibility improvement by Warren [11] when the contextual information was present. Consequently, silence substitution is not a suitable means of voice reconstruction, since white noise is known to give improvements, and is as easy to generate as silence.

Other receiver-only voice reconstruction techniques rely on the assumption that the speech characteristics have not changed from a preceding segment of speech, and use this preceding segment information to reconstruct the missing part; a simple example of this sort of voice reconstruction would be to repeat the last correctly received packet. The mechanisms fail when the packet sizes are large, and the loss rate is high (packets are more likely to be lost in twos or threes, than singularly). A fuller explanation of existing receiver-only techniques can be found in [13].

Combined Source and Channel Techniques

Combined source and channel techniques generally show significant improvement over receiver only techniques. The techniques either transmit extra information within the speech packets (to help with reconstruction at the receiver), or alter the speech coding algorithm and network operation (to make system as a whole more robust to packet loss).

Embedded speech coding techniques used with adaptive differential pulse code modulation (ADPCM [14]), such as those by [15], [16], and code excited linear prediction (CELP) [17], have shown significant

performance improvements during packet loss. Embedded speech coding techniques allow the bit-rate can be adjusted from 40 to 32 or 23 kbps, without the introduction of large amounts of noise; essentially the feed-back loops in the encoder and decoder operate at a lower resolution than usual. The standard was designed to ease the problem of packet loss in packet networks; the codewords are segmented into high and low priority bits, and then placed in different packets. The mechanism relies on arranging for the network to drop packets containing LSBs only, which means that the mechanism is not applicable to networks which do not provide this support, such as today's Internet.

Lara-Barron [15] investigated embedded ADPCM coding techniques at 16-32 ms, and reported success for up to 40% loss (no reduction in speech quality for up to 6% loss).

The significant improvement resulting from the use of this mechanism is mostly due to the preservation of the decoder adaption logic [16].

Speech Quality

Speech quality may be assessed by either subjective or objective means, although it is well known that subjective assessment methods provide more accurate results.

Subjective assessment is usually made by performing listening tests using a large number of subjects. The material used, and the measurements made, depend upon the likely degree of distortion expected.

Toll quality speech coding algorithms are usually assessed by mean opinion scores (MOS) [18], where encoding distortion and noise are the likely type of degradation suffered. The technique involves the listener making a category rating after listening to a passage of speech.

Synthetic quality speech coding algorithms result in speech that has far greater degradation than found in toll quality systems; intelligibility is usually only adequate at best. Consequently, the MOS method is not suitable, and communications quality systems are assessed using comprehension or intelligibility tests. There is a wide range of speech material available, ranging from a sequence of syllables (the listeners transcribe what they hear) to passages (the comprehension of which is ascertained by asking a series of questions) [19].

The speech material is chosen based on the required sensitivity of the results, desired experimental control, and range of human faculties included in the test.

Reliable Audio for Use over the Internet

We have developed a new voice reconstruction scheme that uses redundancy to improve voice reconstruction at the receiver.

The redundant information is the output of a synthetic quality speech coding algorithm (LPC), which

is very low bit-rate (4.8kbps). LPC is generally considered to contain about 60% of the information content of the speech signal, as the overall shape of the frequency spectrum is preserved at the expense of short-term amplitude and pitch variations. This technique is exactly what is required for successful voice reconstruction; the gap will be filled with a sound that is expected, and phonemic restoration should improve the situation further.

The Characteristics of the MBONE / Internet

The Internet, and its multicast overlay (MBONE) is a unique 'shared' packet network, that offers scalable multi-way communication. Such a network has traditionally not been considered suitable for speech applications, because of the large end-to-end delay that is commonly experienced over the network, and the potentially high probability of packet loss, (with relatively large segments of speech being lost). Current audio tools such as vat [20], and nevot [21] have, however, demonstrated that these problems are not prohibitive to successful voice communications, since use of these tools is very widespread.

The Internet provides variable length packets, a feature which has the potential for fine-grained control over the trade-off between network and speech performance requirements. The 'per-packet', rather than 'size-of-packet' network penalty for small packets coupled with the ability to have variable length packets also means that the state information from the speech coding algorithms can be transmitted in the packet, which substantially improves the perception of packet loss. Current audio tools transmit coding algorithm state information in each packet, but replace lost packets with silence.

The Loss Characteristics of the MBONE

Current research by a MICE partner, Bolot [22], is investigating the number of consecutive losses found over the MBONE. The results show that for light and intermediate loads, losses are essentially non-consecutive for an audio stream, and for heavy loads, the behaviour is similar, but consecutive losses are more prevalent.

These results suggest that a model where the redundancy is positioned in the packet after speech from the primary coding algorithm is suitable for light and medium network loads, and a model with the redundancy positioned a number of packets later is suitable for heavy loads.

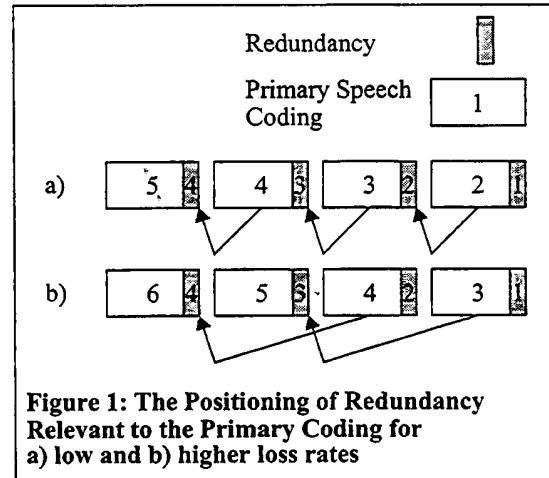
Voice Reconstruction for the MBONE

LPC as the redundant information adds only a small amount of overhead to an RTP [23] packet (12 bytes per 160 bytes of PCM (or per 80 bytes of ADPCM)). The information is piggy-backed to the packet following that containing the primary speech code-words; that the loss of an individual packet can be repaired using the redundant information in the following packet. This mechanism is unique to packet networks, and is only feasible because of the recon-

struction delay introduced at the receiver.

The use of this redundancy technique means an increase in the reconstruction delay by the time equivalent of the distance of the redundancy component after the primary component; this implies an extra delay of one packet for light and medium loading conditions.

Multiple multicast receivers in a single conference may experience a variety of the characteristics reported in [22]. Consequently, the reconstruction mechanism may occasionally have more than one instance of the redundancy after the primary coding scheme packet. In this way, the heavy loading characteristics seen by one site do not affect the performance of the majority.



The provision of LPC redundant information for use in voice reconstruction is intended to be used with per-packet state information; this prevents decoder mistracking in the case of loss. When a packet has been lost, the receiver decodes the redundant information, and feeds the samples to the audio hardware. Consequently, the output speech waveform consists of periods of toll quality speech, interspersed with periods of synthetic quality speech.

While LPC is a fairly complex speech coding algorithm, it should be noted that linear predictive analysis and synthesis are an essential part of all new higher compression schemes: GSM and CELP both use these techniques as a first step in their algorithms. LPC also has the potential to be used elsewhere in the system; as an improvement to the silence detection function, which is an integral part of most packet speech systems.

Experimental Design

Voice reconstruction experiments to date have usually been performed with packet sizes of 16-32ms, or less. The packet sizes used over the MBONE are usually greater than these values (40ms is commonly used). Consequently, little information exists about the degradation commonly experienced in voice con-

Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

Real-Time Litigation Alerts



Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

Advanced Docket Research



With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

Analytics At Your Fingertips



Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.