

IDIAP RESEARCH REPORT



MODIFIED DISCRETE COSINE TRANSFORM FOR ENCODING RESIDUAL SIGNALS IN FREQUENCY DOMAIN LINEAR PREDICTION

Sriram Ganapathy Petr Motlicek
Hynek Hermansky

Idiap-RR-74-2008

DECEMBER 2008

MODIFIED DISCRETE COSINE TRANSFORM FOR ENCODING RESIDUAL SIGNALS IN FREQUENCY DOMAIN LINEAR PREDICTION

Sriram Ganapathy^{1,2}, Petr Motlicek¹, Hynek Hermansky^{1,2}

¹IDIAP Research Institute, Martigny, Switzerland

²École Polytechnique Fédérale de Lausanne (EPFL), Switzerland
{ganapathy,motlicek}@idiap.ch, hermansky@ieee.org

ABSTRACT

Frequency Domain Linear Prediction (FDLP) uses auto-regressive models to represent Hilbert envelopes of relatively long segments of speech/audio signals. Although the basic FDLP audio codec achieves good quality of the reconstructed signal at high bit-rates, there is a need for scaling to lower bit-rates without degrading the reconstruction quality. Here, we present a method for improving the compression efficiency of the FDLP codec by the application of the Modified Discrete Cosine Transform (MDCT) for encoding the FDLP residual signals. In the subjective and objective quality evaluations, the proposed FDLP codec provides competent quality of reconstructed signal compared to the state-of-the-art audio codecs for the 32 – 64 kbps range.

Index Terms— Audio coding, Frequency Domain Linear Prediction (FDLP), Modified Discrete Cosine Transform (MDCT), Perceptual Evaluation of Audio Quality (PEAQ)

1. INTRODUCTION

Conventional approaches to speech coding achieve signal compression with a linear source-filter model of speech production using the linear prediction (LP) [1]. The residual of this modeling process represents the source signal. While such approaches are commercially successful for toll quality conversational services, they do not perform well for mixed signals in many emerging multimedia services. On the other hand, perceptual codecs typically used for multi-media coding applications are not as efficient for speech content.

A new speech/audio coding technique based on modeling the temporal evolution of the spectral dynamics was proposed in [2, 3]. This technique is based on representing Amplitude Modulating (AM) signal using the Hilbert envelope estimate and Frequency Modulating (FM) signal using the Hilbert carrier. The technique exploits the predictability of slowly varying amplitude modulations for encoding speech/audio signals. Input signals are analyzed using a non-uniform Quadrature Mirror Filter (QMF) bank to decompose the signal into frequency sub-bands. For each sub-band signal, Hilbert envelopes are estimated using Frequency Domain Linear Prediction (FDLP), which is an efficient technique for auto-regressive (AR) modeling of the temporal envelopes of a signal [4]. The parameters of the AR model are transmitted to the decoder. The FDLP residual signals are transformed using Discrete Fourier Transform

(DFT) and the magnitude and phase components are quantized separately [3]. At the decoder, these steps are inverted to reconstruct the signal back.

The base-line FDLP codec provides good reconstruction signal quality at high bit-rates ~ 66 kbps. However, there is strong requirement for scaling to lower bit-rates while meeting the reconstruction quality constraints similar to those provided by the state-of-art codecs. The simple encoding set-up of using a DFT based processing for the FDLP residual signal ([3]) offers little freedom in reducing the bit-rates. This is mainly due to the fact that small quantization errors in DFT phase components of the sub-band FDLP residual signals (which consume 60 % of the bit-rate) give rise to significant coding artifacts in the reconstructed signal.

In this paper, we propose an encoding scheme for the FDLP residual signals using Modified Discrete Cosine Transform (MDCT). The MDCT, proposed in [5], outputs a set of critically sampled transform domain coefficients. Perfect reconstruction is provided by time domain alias cancelation and the overlapped nature of the transform. All these properties make the MDCT a potential candidate for application in many popular audio coding systems (for example Advanced Audio Coding (AAC) [6]).

In the proposed FDLP codec, MDCT is applied on short segments (50 ms) of the FDLP residual signals in each sub-band. These coefficients are Vector Quantized (VQ) and transmitted to the receiver along with the parameters of AR model. At the decoder, the MDCT coefficients of the residual are inverse transformed and are used to modulate the FDLP envelope for reconstructing the sub-band signal. Bit-rate scalability is provided by altering the number of VQ levels. The current version of the codec provides high-fidelity audio compression for speech/audio content operating in the bit-rate range of 32 – 64 kbps. In the objective and subjective quality evaluations, the proposed FDLP codec provides competitive results compared to the state-of-art codecs at similar bit-rates.

The rest of the paper is organized as follows. Sec. 2 describes the FDLP technique for AR modelling of Hilbert Envelopes. The basic structure of the proposed FDLP codec is described in Sec. 3. The objective and subjective evaluations are reported in Sec. 4.

2. FREQUENCY DOMAIN LINEAR PREDICTION

Typically, auto-regressive (AR) models have been used in speech applications for representing the envelope of the power spectrum of the signal by performing the operation of Time Domain Linear Prediction (TDLP) [7]. This paper utilizes AR models for obtaining smoothed, minimum phase, parametric models of temporal rather than spectral envelopes. The duality between the time and frequency domains means that AR modeling can be applied equally well to dis-

This work was partially supported by grants from ICSI Berkeley, USA; the Swiss National Center of Competence in Research (NCCR) on “Interactive Multi-modal Information Management (IM)2”; managed by the IDIAP Research Institute on behalf of the Swiss Federal Authorities. The authors also thank Marios Athineos for MDCT code fragments.

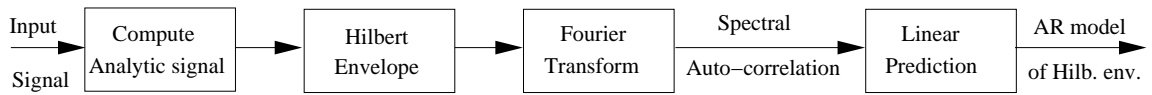


Fig. 1. Steps involved in the FDLP technique for AR modelling of Hilbert Envelopes.

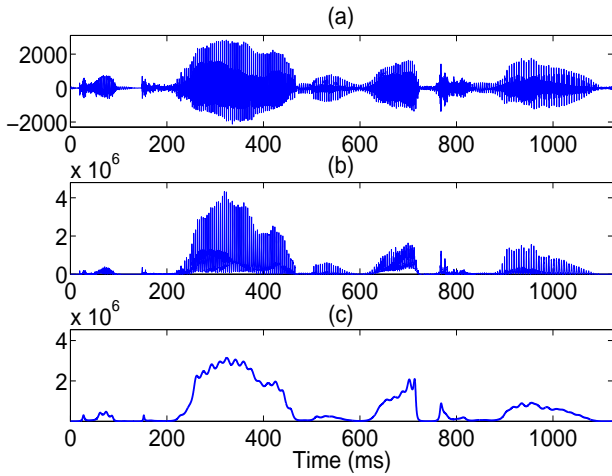


Fig. 2. Illustration of the AR modelling property of FDLP. (a) a portion of speech signal, (b) its Hilbert envelope and (c) all-pole model obtained using FDLP.

crete spectral representations of the signal instead of time-domain signal samples.

The block schematic showing the steps involved in deriving the AR model of Hilbert envelope is shown in Fig. 1. The first step is to compute the analytic signal for the input signal. For a discrete time signal, the analytic signal can be obtained using the Fourier Transform [8]. Hilbert envelope (squared magnitude of the analytic signal) and spectral auto-correlation function form Fourier transform pairs [9]. This relation is used to derive the auto-correlation of the spectral components of the signal which are then used for deriving the FDLP models (in manner similar to the computation of the TDLP models from temporal autocorrelations [7]).

For the FDLP technique, the squared magnitude response of the all-pole filter approximates the Hilbert envelope of the signal. An illustration of the all-pole modelling property of the FDLP technique

is shown in Fig. 2, where we plot a portion of speech signal, its Hilbert envelope computed using Fourier transform [8] and the AR model fit to the Hilbert envelope provided by FDLP.

3. SPEECH/AUDIO CODEC BASED ON FDLP

The block schematic of the FDLP based encoder is given in Fig. 3.

3.1. FDLP analysis

Long temporal segments (1000 ms) of the input speech/audio signals are decomposed into 32 non-uniform QMF sub-bands which approximate the critical band decomposition in the auditory system. In each sub-band, the FDLP analysis is applied to obtain a set of AR model parameters and the FDLP residual signal. The FDLP envelope coefficients are converted to Line Spectral Frequencies (LSF) which approximate the sub-band temporal envelopes. These LSF parameters are quantized using Vector Quantization (VQ).

3.2. Encoding FDLP residual signals using MDCT

The sub-band FDLP residual signals are split into relatively short frames (50 ms) and transformed using the MDCT. We use the sine window with 50 % overlap for the MDCT analysis as this was experimentally found to provide the best reconstruction quality (based on objective quality scores). Since a full-search VQ in the MDCT domain with good resolution would be computationally infeasible, the split VQ approach is employed. Although the split VQ approach is suboptimal, it reduces the computational complexity and memory requirements to manageable limits without severely degrading the VQ performance. The VQ codebooks are trained on a large audio database using the LBG algorithm. The quantized levels are Huffman encoded for further reduction of bit-rates (bit-rate reduction of about 10 %). Quantization of the MDCT coefficients using the split VQ consumes around 80% of the bit-rate. The MDCT coefficients for the lower frequency sub-bands are quantized using higher number of VQ levels compared to those from the higher bands. For the

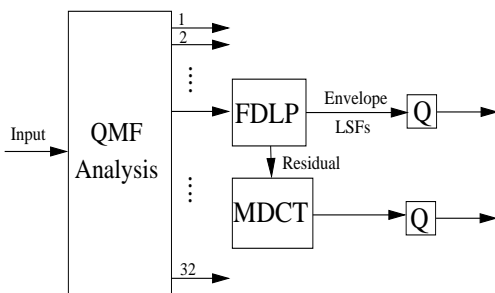


Fig. 3. Scheme of the FDLP encoder.

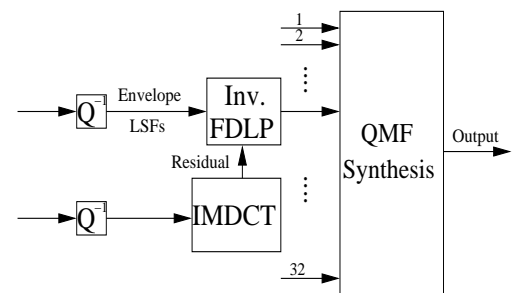


Fig. 4. Scheme of the FDLP decoder.

ODG Scores	Quality
0	imperceptible
-1	perceptible but not annoying
-2	slightly annoying
-3	annoying
-4	very annoying

Table 1. PEAQ scores and their meanings.

bit-rate [kbps]	64	64	66	64
Codec	LAME	AAC	FDLP-DFT	FDLP
PEAQ	-1.6	-0.8	-1.2	-0.7
bit-rate [kbps]	48	48	48	48
Codec	LAME	AAC	FDLP	FDLP
PEAQ	-2.5	-1.1	-1.2	-1.2
bit-rate [kbps]	32	32	32	32
Codec	LAME	AAC	AMR	FDLP
PEAQ	-3.0	-2.4	-2.2	-2.4

Table 2. Average PEAQ scores for 28 speech/audio files at 64, 48 and 32 kbps.

purpose of scaling the bit-rates, all the sub-bands are treated uniformly. The current version of the codec follows a simple signal independent bit assignment for the MDCT coefficients.

3.3. FDLP based decoder

In the decoder, shown in Fig. 4, quantized MDCT coefficients of the FDLP residual signals are reconstructed and transformed back to the time-domain using Inverse MDCT (IMDCT). The reconstructed FDLP envelopes (from LSF parameters) are used to modulate the corresponding sub-band residual signals. Finally, sub-band synthesis is applied to reconstruct the full-band signal.

4. QUALITY EVALUATIONS

The subjective and objective evaluations of the proposed audio codec are performed using audio signals (sampled at 48 kHz) present in the framework for exploration of speech and audio coding [10]. This database is comprised of speech, music and speech over music recordings. The music samples contain a wide variety of challenging audio samples ranging from tonal signals to highly transient signals.

The objective and subjective quality evaluations of the following codecs are considered:

1. The proposed FDLP codec with MDCT based residual signal processing, at 32, 48 and 64 kbps, denoted as FDLP.
2. The previous version of the FDLP codec [3], at 66 kbps, denoted as FDLP-DFT.
3. LAME MP3 (MPEG 1, layer 3) [13], at 32, 48 and 64, kbps denoted as LAME.
4. MPEG-4 HE-AAC, v1, at 32, 48 and 64 kbps [6], denoted as AAC. The HE-AAC coder is the combination of Spectral Band Replication (SBR) [14] and Advanced Audio Coding (AAC) [15].
5. AMR-WB plus standard [16], at 32 kbps, denoted as AMR.

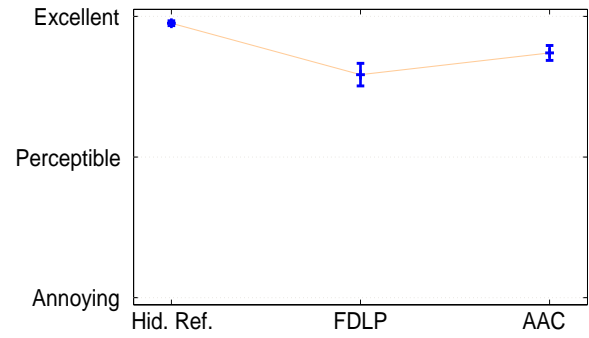


Fig. 5. BS.1116 results for 5 speech/audio samples using two coded versions at 64 kbps (FDLP-MDCT (FDLP), MPEG-4 HE-AAC (AAC)), hidden reference (Hid. Ref.).

4.1. Objective Evaluations

The objective measure employed is the Perceptual Evaluation of Audio Quality (PEAQ) distortion measure [17]. In general, the perceptual degradation of the test signal with respect to the reference signal is measured, based on the ITU-R BS.1387 (PEAQ) standard. The output combines a number of Model Output Variables (MOV's) into a single measure, the Objective Difference Grade (ODG) score, which is an impairment scale with meanings shown in Table 1. The mean PEAQ score for the 28 speech/audio files in [10] is used as the objective quality measure.

The first set of results given in Table 2 compare the objective quality scores of the proposed FDLP codec at 64 kbps with the FDLP-DFT codec at 66 kbps. The objective quality scores for AAC and LAME codecs at 64 kbps are also shown. This table shows the advantage of using the MDCT for encoding the FDLP residuals instead of using the DFT.

The next set of results in Table 2 show the average PEAQ scores for the proposed FDLP codec with AAC and LAME codecs at 48 kbps and the scores for these codecs along with the AMR codec at 32 kbps. The objective scores for the proposed FDLP codec at these bit-rates follow a similar trend to that of the state-of-the-codecs.

4.2. Subjective Evaluations

Since the encoded audio signals at 64 kbps have small impairments compared to the original, we perform the BS.1116 methodology of subjective evaluation [18]. As this subjective evaluation is time consuming, only two coded versions (FDLP and AAC) are compared at 64 kbps along with the hidden reference. The subjective results with 7 listeners using 5 speech/audio samples from the database is shown in Fig. 5. Here, the mean scores are plotted with 95% confidence interval. The proposed FDLP codec at 64 kbps is judged to be similar to the AAC codec at the same bit-rate.

For the audio signals encoded at 48 kbps and 32 kbps, the MUSHRA (Multiple Stimuli with Hidden Reference and Anchor) methodology for subjective evaluation is employed. It is defined by ITU-R recommendation BS.1534 [19]. We perform the MUSHRA tests on 6 speech/audio samples from the database with 6 listeners. The mean MUSHRA scores (with 95% confidence interval) for the subjective listening tests at 48 kbps and 32 kbps (given in Fig. 6 and Fig. 7 respectively) show that the subjective quality of the proposed codec is slightly poorer than the AAC codec but better than

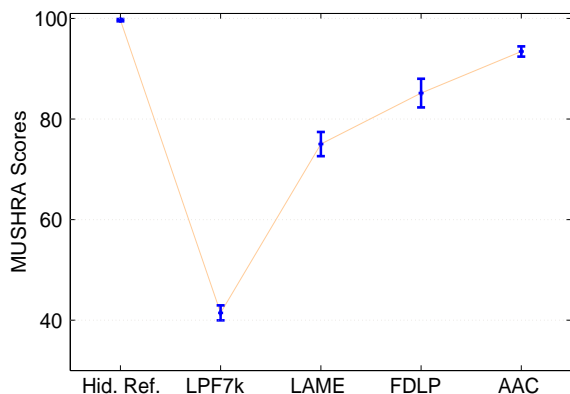


Fig. 6. MUSHRA results for 6 speech/audio samples using three coded versions at 48 kbps (FDLP-MDCT (FDLP), MPEG-4 HE-AAC (AAC) and LAME-MP3 (LAME)), hidden reference (Hid. Ref.) and 7 kHz low-pass filtered anchor (LPF7k).

the LAME codec.

5. CONCLUSIONS

In order to improve the compression efficiency of audio codecs based on spectral dynamics, we propose a new method of encoding the FDLP residual signals by the application of MDCT. This new technique offers the advantage of bit-rate scalability similar to the state-of-the-art codecs. Objective evaluations justify the improvement provided by the use of MDCT as compared to the use of DFT in the previous versions of the FDLP codec. The current version of the codec provides subjective results which are competitive to the state of the art codecs in the bit-rate range of 32-64 kbps. Furthermore, this performance is achieved without utilizing standard modules like psycho-acoustic modelling and signal adaptive windowing. The inclusion of these techniques form part of the future work.

6. REFERENCES

- [1] Schroeder M. R. and Atal B. S., "Code-excited linear prediction (CELP): high-quality speech at very low bit rates", *Proc. of ICASSP*, Vol. 10, pp. 937-940, Apr. 1985.
- [2] Motlicek P., Ganapathy S., Hermansky H., and Garudabri H., "Frequency Domain Linear Prediction for QMF Sub-bands and Applications to Audio coding", *Proc. of MLMI*, LNCS Series, Springer-Verlag, 2007.
- [3] Ganapathy S., Motlicek P., Hermansky H., and Garudabri H., "Autoregressive Modelling of Hilbert Envelopes for Wideband Audio Coding", *Audio Engg. Soc.*, 124th Convention, May 2008.
- [4] Athineos M., and Ellis D., "Autoregressive Modeling of Temporal Envelopes", *IEEE Trans. on Signal Proc.*, Vol. 55, pp. 5237 - 5245, Nov. 2007.
- [5] Princen J., Johnson A., Bradley A., "Subband/Transform Coding Using Filter Bank Designs Based on Time Domain Aliasing Cancellation", *Proc. of ICASSP*, Vol. 87, pp 2161-2164, May 1987.

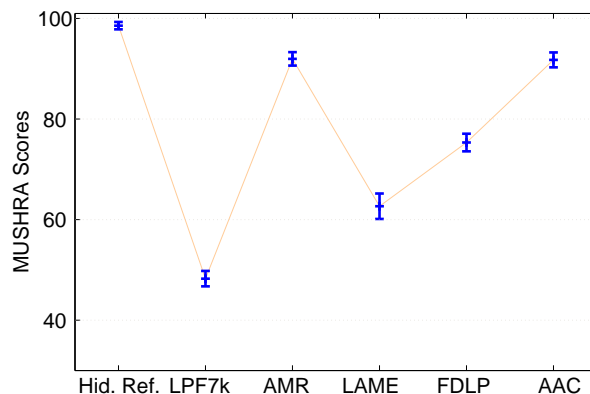


Fig. 7. MUSHRA results for 6 speech/audio samples using four coded versions at 32 kbps (AMR-WB+ (AMR), FDLP-MDCT (FDLP), MPEG-4 HE-AAC (AAC) and LAME-MP3 (LAME)), hidden reference (Hid. Ref.) and 7 kHz low-pass filtered anchor (LPF7k).

- [6] 3GPP TS 26.401: Enhanced aacPlus general audio codec; General Description.
- [7] Makhoul J., "Linear Prediction: A Tutorial Review", *Proc. of the IEEE*, Vol. 63, pp. 561-580, 1975.
- [8] Marple L.S., "Computing the Discrete-Time Analytic Signal via FFT", *IEEE Trans. on Acoust., Speech and Sig. Proc.*, Vol. 47, pp. 2600-2603, 1999.
- [9] Herre J. and Johnston J. D., "Enhancing the Performance of Perceptual Audio Coders by using Temporal Noise Shaping (TNS)", *Audio Engg. Soc.*, 101st Convention, pp. 1-24, 1996.
- [10] ISO/IEC JTC1/SC29/WG11, "Framework for Exploration of Speech and Audio Coding", *MPEG2007/N9254*, July 2007.
- [11] LAME MP3 codec: <http://lame.sourceforge.net>
- [12] Martin Dietz, Lars Liljeryd, Kristofer Kjorling and Oliver Kunz, "Spectral Band Replication, a novel approach in audio coding", *Audio Engg. Soc.*, 112th Convention, May 2002.
- [13] LAME MP3 codec: <http://lame.sourceforge.net>
- [14] Martin Dietz, Lars Liljeryd, Kristofer Kjorling and Oliver Kunz, "Spectral Band Replication, a novel approach in audio coding", *Audio Engg. Soc.*, 112th Convention, May 2002.
- [15] M. Bosi, K. Brandenburg, S. Quackenbush, L. Fielder, K. Akagiri, H. Fuchs, M. Dietz, J. Herre, G. Davidson, Y. Oikawa, "ISO/IEC MPEG-2 Advanced Audio Coding", *J. Audio Eng. Soc.*, Vol. 45, pp. 789-814, Oct. 1997.
- [16] "Extended AMR Wideband codec", <http://www.3gpp.org/ftp/Specs/html-info/26290.htm>
- [17] ITU-R Recommendation BS.1387, "Method for objective psychoacoustic model based on PEAQ to perceptual audio measurements of perceived audio quality", Dec. 1998.
- [18] ITU-R Recommendation BS.1116, "Methods for the subjective assessment of small impairments in audio systems including multichannel sound systems", Oct. 1997.
- [19] ITU-R Recommendation BS.1534, "Method for the subjective assessment of intermediate audio quality", June 2001.