# DOUBLE-TALK ROBUST ACOUSTIC ECHO CANCELLATION WITH CONTINUOUS NEAR-END ACTIVITY

*Toon van Waterschoot and Marc Moonen*

ESAT-SCD, Katholieke Universiteit Leuven
Kasteelpark Arenberg 10, B-3001 Leuven, Belgium
phone: +32 16 321927, fax: +32 16 321970, email: toon.vanwaterschoot@esat.kuleuven.ac.be
web: http://www.esat.kuleuven.ac.be/scd/

## ABSTRACT

In some acoustic echo cancellation scenarios, such as an automatic gain adjustment application, near-end noise may be continuously present. In this case a double-talk detector cannot be applied and the adaptive algorithm should behave in a robust way w.r.t. the disturbing near-end signal. From linear estimation theory it is known that the variance of the room impulse response estimate may be decreased by taking into account the near-end signal characteristics. From the expression for the best linear unbiased estimate, we derive a prediction error criterion from which the near-end signal model and the room impulse response can be estimated concurrently. We propose a new recursive identification algorithm for minimization of the proposed prediction error criterion. The proposed algorithm is in fact a variant of a prediction error identification algorithm that was developed recently for adaptive feedback cancellation. Simulation results indicate that indeed a fast converging echo cancellation algorithm may be obtained with the proposed method, as compared to ordinary RLS and NLMS adaptive algorithms.

## 1. INTRODUCTION

Acoustic echo cancellation (AEC) has been a popular research topic in acoustic signal processing, motivated mainly by the increasing demand for hands-free speech communication. A classical AEC scenario is shown in **Figure 1**. A speech signal $u(t)$ from the far-end side is broadcasted in an acoustic enclosure (the 'room') by means of a loudspeaker. A microphone is present in the room for recording a local signal $v(t)$ (the 'near-end signal') which is to be transmitted back to the far-end side. An acoustic echo path exists between the loudspeaker and the microphone such that the recorded microphone signal $y(t) = x(t) + v(t)$ contains an undesired echo component $x(t)$ in addition to the near-end signal component $v(t)$. If the echo path transfer function is modelled as a finite impulse response (FIR) filter $F(q,t) \triangleq f_0(t) + f_1(t)q^{-1} + \ldots + f_{n_F}(t)q^{-n_F}$, then the echo component can be considered as a filtered version of the loudspeaker signal: $x(t) = F(q,t)u(t)$. Here $q$ denotes the time shift operator, e.g. $q^{-k}u(t) = u(t-k)$. The main objective in AEC is to identify the unknown room impulse response (RIR) $F(q,t)$ and hence to subtract an estimate of the echo component from the microphone signal. In this way an echo-compensated signal $d(t) = y(t) - \hat{F}(q,t)u(t)$ is sent to the far-end side, with $\hat{F}(q,t)$ an estimate of $F(q,t)$.

Since $F(q,t)$ may be time-varying (e.g. due to people moving around the room), an adaptive algorithm is usually applied for the estimation of the RIR. It is a well-known problem in AEC that the
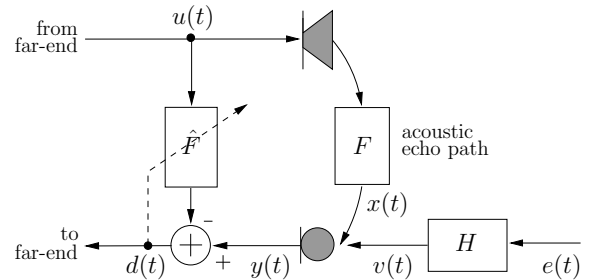
Figure 1: AEC scenario with AR modelling of the near-end signal.

convergence speed and hence the tracking capabilities of standard adaptive algorithms (like recursive least squares (RLS) or normalized least mean squares (NLMS)) may decrease severely when near-end noise is present ('double-talk' periods). A lot of effort has been spent on the design of efficient double-talk detectors (DTD), which are used to slow down or switch off the adaptation during double-talk periods [1]. Nevertheless in some scenarios near-end noise will be continuously present and the use of a DTD becomes futile. This may be the case for example in an automatic gain adjustment application.

In this paper we aim at developing a recursive identification algorithm that allows for continuous adaptation of the RIR estimate, yet behaves in a robust way in double-talk situations. From linear estimation theory [2], we know that the best (i.e. minimum variance) linear unbiased estimator (BLUE) for an unknown system depends on the characteristics of the noise acting upon the system. In the AEC context it is the near-end signal which acts as a noise signal to the RIR identification. Therefore we expect that by using knowledge of the near-end signal characteristics, the convergence properties of the RIR identification algorithm can be improved. However the near-end signal characteristics are typically unknown and may be highly time-varying. Therefore they need to be estimated concurrently with the unknown RIR.

The paper is organized as follows. We first review some results from linear estimation theory [2] in **Section 2** to indicate how the variance of the RIR estimate can be decreased. This leads to the expression for the best linear unbiased estimate (BLUE), from which we derive in **Section 3** a prediction error criterion. This criterion is a function of both the near-end signal model and the RIR. Then in **Section 4** a two-stage identification algorithm is described that makes use of the bilinearity of the prediction error. The algorithm comes in two flavours: a sliding window variant that follows naturally from the prediction error criterion and hopping window variant that exploits the quasistationary behaviour of audio signals and is computationally more efficient. The hopping window variant is equivalent to the PEM-AFROW algorithm proposed recently for adaptive feedback cancellation [3] (PEM-AFROW stands for prediction error method based adaptive filtering performing only row operations). Finally in **Section 5** both variants are compared by means of computer simulations, both for a Gauss-Newton and a stochastic gradient implementation.

## 2. BEST LINEAR UNBIASED ESTIMATE

Let us assume that a data record $\{u(k), y(k)\}_{k=1}^{t}$ of microphone and loudspeaker samples is available. Then the linear estimation problem at time $t$ can be written as

$$\begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(t) \end{bmatrix} = \begin{bmatrix} u(1) & \ldots & u(1-n_F) \\ u(2) & \ldots & u(2-n_F) \\ \vdots & \ddots & \vdots \\ u(t) & \ldots & u(t-n_F) \end{bmatrix} \cdot \begin{bmatrix} f_0 \\ \vdots \\ f_{n_F} \end{bmatrix} + \begin{bmatrix} v(1) \\ v(2) \\ \vdots \\ v(t) \end{bmatrix}$$
$$\Updownarrow$$
$$\mathbf{y} = \mathbf{U}\mathbf{f} + \mathbf{v}$$

where $\mathbf{f}$ is the $(n_F + 1) \times 1$ parameter vector containing the coefficients of $F(q,t)$ that are to be estimated.

Any linear estimate of parameter vector $\mathbf{f}$ can be written as a linear function of the data vector $\mathbf{y}$:

$$\hat{\mathbf{f}} = \mathbf{Z}^T \mathbf{y}. \tag{1}$$

For this estimate to be unbiased, the $t \times (n_F + 1)$ matrix $\mathbf{Z}$ should be subjected to two constraints:

$$\begin{cases} \mathbf{Z}^T \mathbf{U} & = & \mathbf{I}_{n_F+1} & (a) \\ E\mathbf{Z}^T \mathbf{v} & = & \mathbf{0}_{(n_F+1) \times 1} & (b) \end{cases} \tag{2}$$

Since $\mathbf{Z}$ is typically a function of loudspeaker Hankel matrix $\mathbf{U}$, constraint (2(b)) can be reduced to $E\mathbf{U}^T\mathbf{v} = \mathbf{0}$, which we assume to be fulfilled. In AEC this comes down to assuming that no closed signal loop is created due to an acoustic echo path in the far-end room.

Minimizing the variance $E(\hat{\mathbf{f}} - E\hat{\mathbf{f}})(\hat{\mathbf{f}} - E\hat{\mathbf{f}})^T$ of the estimate (1) under the unbiasedness constraint (2(a)) then yields the best linear unbiased estimate (BLUE):

$$\hat{\mathbf{f}}_{\mathbf{BLUE}} = (\mathbf{U}^T \mathbf{R}^{-1} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{R}^{-1} \mathbf{y}. \tag{3}$$

$\mathbf{R}$ represents the near-end signal correlation matrix, defined by

$$\mathbf{R} \triangleq E\mathbf{v}\mathbf{v}^T. \tag{4}$$

The BLUE covariance matrix is minimal among all linear unbiased estimates and given by

$$\text{cov}(\hat{\mathbf{f}}_{\mathbf{BLUE}}) = (\mathbf{U}^T \mathbf{R}^{-1} \mathbf{U})^{-1}.$$

Note that the BLUE in (3) cannot be calculated as such, because the near-end signal correlation matrix $\mathbf{R}$ is usually unknown. Nevertheless, from (3) we may derive a prediction error criterion from which both the RIR and the near-end signal characteristics may be estimated.

## 3. PREDICTION ERROR CRITERION

Let us first decompose the near-end signal correlation matrix $\mathbf{R}$ appearing in expression (3) for the BLUE. We therefore assume that the near-end signal $v(t)$ is generated as

$$v(t) = H(q,t)e(t) \quad \text{with} \quad Ee(t)e(t-k) = (k) \, _t^2.$$

The near-end excitation signal $e(t)$ is a white noise signal with a time-dependent variance $_t^2$, and $H(q,t)$ is a linear model with time-dependent coefficients. Expression (4) may then be rewritten as

$$\mathbf{R} = E\mathbf{H}\mathbf{e}\mathbf{e}^T \mathbf{H}^T \tag{5}$$

with $\mathbf{e} \triangleq [e(1) \quad \ldots \quad e(t)]^T$, and

$$\mathbf{H} = \mathbf{H}^T \triangleq \begin{bmatrix} H(q,1) & \ldots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \ldots & H(q,t) \end{bmatrix}.$$

If the near-end signal model $H(q,k)$, $k = 1 \ldots t$, is considered to be deterministic then the expectation operator in (5) can be shifted to the inner product $\mathbf{e}\mathbf{e}^T$:

$$\begin{aligned} \mathbf{R} & = & \mathbf{H}E\mathbf{e}\mathbf{e}^T\mathbf{H}^T \\ & = & \mathbf{H}\boldsymbol{\Lambda}\mathbf{H}^T \end{aligned}$$

with

$$\triangleq \begin{bmatrix} _1^2 & \ldots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \ldots & _t^2 \end{bmatrix}.$$

Hence the BLUE in (3) can be realized as

$$\hat{\mathbf{f}}_{\mathbf{BLUE}} = (\mathbf{U}^T \mathbf{H}^{-T} \boldsymbol{\Lambda}^{-1} \mathbf{H}^{-1} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{H}^{-T} \boldsymbol{\Lambda}^{-1} \mathbf{H}^{-1} \mathbf{y},$$

that is, by prefiltering and weighting the $k$-th row of $\mathbf{U}$ and $\mathbf{y}$ with the inverse near-end signal model $H^{-1}(q,k)$ and the inverse near-end excitation signal variance $_k^{-2}$ respectively.

If we impose an autoregressive (AR) model structure on the near-end signal, i.e.

$$H(q,t) = \frac{1}{A(q,t)} = \frac{1}{1 + a_1(t)q^{-1} + \ldots + a_{n_A}(t)q^{-n_A}},$$

then the prefilters $H^{-1}(q,k) = A(q,k)$, $k = 1 \ldots t$, turn out to be FIR filters of order $n_A$.

The BLUE can be seen to minimize at each time instant $t$ the prediction error criterion

$$V_{PE}(t,\mathbf{f}) = \frac{1}{2t} \sum_{k=1}^{t} \frac{1}{_k^2} \, ^2(k,\mathbf{f}), \tag{6}$$

with the prediction error defined as

$$(k,\mathbf{f}) = A(q,k)[y(k) - F(q,t)u(k)].$$

In **Section 4** we will derive a prediction error identification algorithm which minimizes the prediction error criterion in (6) recursively. However, in order to suit the application we have in mind, two modifications are made to the criterion in (6). First of all, we will allow the RIR $F(q,t)$ to vary with time, which is physically relevant as the acoustic environment may change. Therefore the parameter vector $\mathbf{f}(t)$ will be identified recursively and an exponential forgetting factor is included in the criterion. Secondly, up till now we have considered $A(q,k)$, $k = 1 \ldots t$, as a known, deterministic prefilter and $_k^{-2}$, $k = 1 \ldots t$, as a given weight. In practice, $A(q,t)$ and $_t^2$ have to be estimated concurrently with $F(q,t)$ at each time instant $t$. The modified prediction error criterion then looks like

$$V_{PE}(t,\mathbf{f}(t),\mathbf{a}(t),\,_t^2) = \frac{1}{2N} \sum_{k=1}^{t} \frac{t-k}{_k^2} \big( A(q,k)[y(k) - F(q,t)u(k)] \big)^2,$$

where $N = 1/(1 - )$ denotes the effective window length and $\mathbf{a}(t) \triangleq [a_1(t) \ldots a_{n_A}(t)]^T$ is the $n_A \times 1$ parameter vector containing the AR coefficients to be estimated at time $t$ (note that $a_0 = 1$ is not included in $\mathbf{a}(t)$).

## 4. PREDICTION ERROR IDENTIFICATION ALGORITHM

The prediction error $(t,\mathbf{f}(t),\mathbf{a}(t)) = A(q,t)[y(t) - F(q,t)u(t)]$ is nonlinear in the coefficients of $\mathbf{f}(t)$ and $\mathbf{a}(t)$. However, the prediction error has the property that if $\mathbf{a}(t)$ is assumed to be known, it is linear in $\mathbf{f}(t)$ and vice versa. The prediction error is said to be bilinear in $\mathbf{f}(t)$ and $\mathbf{a}(t)$ [4]. It is useful to exploit this property in the derivation of a prediction error identification algorithm, by performing the identification in two stages. We assume that at time

instant $t$ the estimates $\hat{\mathbf{a}}(t-1)$ and $\hat{\mathbf{f}}(t-1)$ are available from the previous recursion step.

In the first stage of the algorithm a linear prediction is performed on the echo-compensated signal $d(t,\hat{\mathbf{f}}(t-1))$, calculated using the previous estimate $\hat{\mathbf{f}}(t-1)$. The signal $d(t,\hat{\mathbf{f}}(t-1))$ is windowed with a rectangular sliding window of length $M$:

$$
\mathbf{d}(t)=\begin{bmatrix} y(t) \\ \vdots \\ y(t-M+1) \end{bmatrix} - \begin{bmatrix} u(t) & \dots & u(t-n_F) \\ \vdots & \ddots & \vdots \\ u(t-M+1) & \dots & u(t-M+1-n_F) \end{bmatrix}\hat{\mathbf{f}}(t-1).
$$

The autocorrelation functions $\hat{r}_{dd}(\tau)$, $\tau=0\dots n_A$, of $d(t,\hat{\mathbf{f}}(t-1))$ are then estimated using the autocorrelation method:

$$
\begin{bmatrix} \hat{r}_{dd}(0) \\ \hat{r}_{dd}(1) \\ \vdots \\ \hat{r}_{dd}(n_A) \end{bmatrix} = \begin{bmatrix} 0 & \dots & d(t) & \dots & d(t-M+1) \\ 0 & \dots & d(t-1) & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ d(t) & \dots & d(t-n_A) & \dots & 0 \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ d(t) \\ \vdots \\ d(t-M+1) \end{bmatrix}
$$

The near-end signal AR coefficients $\mathbf{a}(t)$ and the near-end excitation signal variance $\sigma_t^2$ are then estimated from $\hat{r}_{dd}(\tau)$, $\tau=0\dots n_A$, using the Levinson-Durbin recursion.

In the second stage of the identification algorithm, the microphone and loudspeaker data needed for the recursive update of the RIR estimate are prefiltered using the estimated coefficients $\hat{\mathbf{a}}(t)$ from the first stage:

$$
y_A(t) = [y(t) \quad \dots \quad y(t-n_A)]\begin{bmatrix} 1 \\ \hat{\mathbf{a}}(t) \end{bmatrix},
$$

$$
\mathbf{u_A}(t) = \begin{bmatrix} u(t) & \dots & u(t-n_A) \\ \vdots & \ddots & \vdots \\ u(t-n_F) & \dots & u(t-n_F-n_A) \end{bmatrix}\begin{bmatrix} 1 \\ \hat{\mathbf{a}}(t) \end{bmatrix}.
$$

The RIR estimate $\hat{\mathbf{f}}(t-1)$ can then be updated recursively, either with the Gauss-Newton method:

$$
\hat{\mathbf{f}}(t) = \hat{\mathbf{f}}(t-1) + \frac{1}{\hat{\sigma}_t^2}\mathbf{R_f}^{-1}(t)\mathbf{u_A}(t)\varepsilon_p(t),
$$

$$
\mathbf{R_f}(t) = \mathbf{R_f}(t-1) + \frac{1}{\hat{\sigma}_t^2}\mathbf{u_A}(t)\mathbf{u_A}^T(t), \tag{7}
$$

or with the stochastic gradient method:

$$
\hat{\mathbf{f}}(t) = \hat{\mathbf{f}}(t-1) + \frac{\mathbf{u_A}(t)\varepsilon_p(t)}{\mathbf{u_A}^T(t)\mathbf{u_A}(t) + (n_F+1)\hat{\sigma}_t^2} \tag{8}
$$

where in both cases weighting is performed using the estimated variance $\hat{\sigma}_t^2$ from the first stage, and the a priori prediction error is calculated as

$$
\varepsilon_p(t) = \varepsilon(t,\hat{\mathbf{f}}(t-1),\hat{\mathbf{a}}(t)) = y_A(t) - \mathbf{u_A}^T(t)\hat{\mathbf{f}}(t-1).
$$

The complexity of the proposed algorithm as compared to an ordinary RLS or NLMS adaptive algorithm, may be reduced by taking into account that most audio signals exhibit a quasi-stationary behaviour. In this respect, if the near-end signal is assumed to behave stationary during time intervals with average length $P$, the first stage of the algorithm may be performed only every $P$-th time instant, instead of every time instant. In other words, the sliding window is replaced by a hopping window with hop size $P$.

In the hopping window variant of the prediction error identification algorithm, the first stage is only executed when $t/P \in \mathbb{Z}$. The linear prediction is then performed on a rectangular data window

that 'looks ahead' $P-1$ samples of the echo-compensated signal $d(t,\hat{\mathbf{f}}(t-1))$:

$$
\mathbf{d}(t)=\begin{bmatrix} y(t+P-1) \\ \vdots \\ y(t+P-M) \end{bmatrix} - \begin{bmatrix} u(t+P-1) & \dots & u(t+P-1-n_F) \\ \vdots & \ddots & \vdots \\ u(t+P-M) & \dots & u(t+P-M-n_F) \end{bmatrix}\hat{\mathbf{f}}(t-1).
$$

The estimated coefficients $\hat{\mathbf{a}}(t)$ and variance $\hat{\sigma}_t^2$ are then used in the second stage of the algorithm during $P$ recursive steps.

We conclude this section by noting that the hopping window prediction error identification algorithm is equivalent to an adaptive feedback cancellation algorithm proposed recently [3]. For convenience, we adopt the acronym PEM-AFROW, which stands for prediction error method based adaptive filtering applying only row operations (to the loudspeaker data matrix).

## 5. SIMULATION RESULTS

MATLAB simulations were performed to compare the convergence properties of both variants of the PEM-AFROW algorithm described in **Section 4**. A recursive least squares (RLS) and a normalized least mean squares (NLMS) algorithm were implemented as reference algorithms. At a sampling rate $f_s = 8kHz$, $F(q,t)$ was a fixed, realistic room impulse response of length $n_F + 1 = 1000$. In one series of experiments the AR model order was set to $n_A = 12$ which is a commonly used value in speech processing. In a second series the AR model order was raised to $n_A = 55$, a value high enough to predict also the pitch of the near-end excitation signal during voiced speech segments. The far-end signal $u(t)$ was a $1,5s$ male speech fragment and the near-end signal $v(t)$ a $1,5s$ female speech fragment. The near-end signal $v(t)$ was scaled such that the average echo-to-background ratio (EBR) was equal to $10dB$:

$$
EBR \triangleq \frac{\sum_{k=1}^{N}|x(k)|^2}{\sum_{k=1}^{N}|v(k)|^2} = 10dB.
$$

$N = 12000$ denotes the number of data points used for simulations with the Gauss-Newton method. For the stochastic gradient simulations, $N = 480000$ and the far-end and near-end speech fragments were repeated 40 times. The exponential forgetting factor in (7) was set to $\lambda = 0.9997$ and the step size in (8) to $\mu = 0.5$. The length of the rectangular window for linear prediction was set to $M = 215$ for all experiments. For the hopping window variant, the hop size was set to $P = M - n_A$. The performance measure used for comparison was the logarithmic normalized bias, defined as

$$
\nu(t) = 20\log_{10}\frac{\|\hat{\mathbf{f}}(t)-\mathbf{f}\|}{\|\mathbf{f}\|}.
$$

The convergence curves for the sliding window (SW) and hopping window (HW) PEM-AFROW algorithm using the Gauss-Newton method are shown in **Figures 2** and **4** respectively. It is clear that for both AR model orders the HW variant outperforms the SW variant (compare the dashed curves), which may come as a surprise. It turns out that keeping the AR coefficients fixed during several consecutive recursion steps (as is the case in the HW variant) prevents the algorithm from converging to a local minimum of the prediction error criterion. Moreover, even when the AR model is identified on the true near-end signal (see the dotted curves), the HW variant shows a faster convergence than the SW variant. However in this case the prediction error criterion has no local minima. So it appears that the variance of the RIR estimate is lower in case the AR coefficients are not estimated at each time instant.

It can be seen from the dotted curves that both PEM-AFROW variants show a potential convergence improvement between $10dB$ and $20dB$ compared to an ordinary RLS algorithm. When knowledge of the true near-end signal is not used, the improvement of the HW variant compared to the RLS algorithm is still $5dB$ to $10dB$ if the AR model order is set high enough.

In **Figures 3** and **5** the convergence curves for the stochastic gradient implementation of both PEM-AFROW variants are shown.
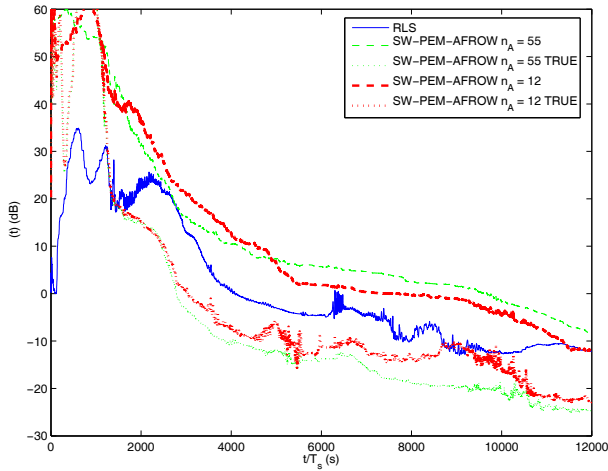
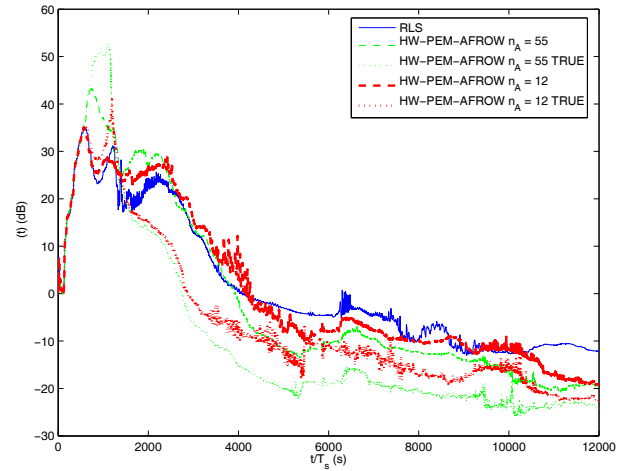Figure 2: Convergence curves of sliding window PEM-AFROW using the Gauss-Newton method.
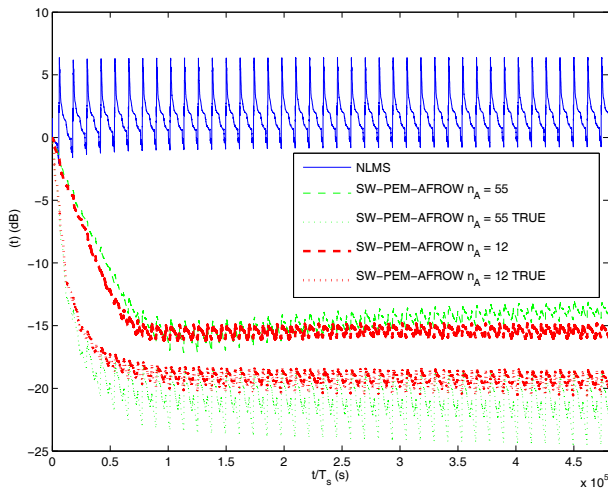


Figure 3: Convergence curves of sliding window PEM-AFROW using the stochastic gradient method.



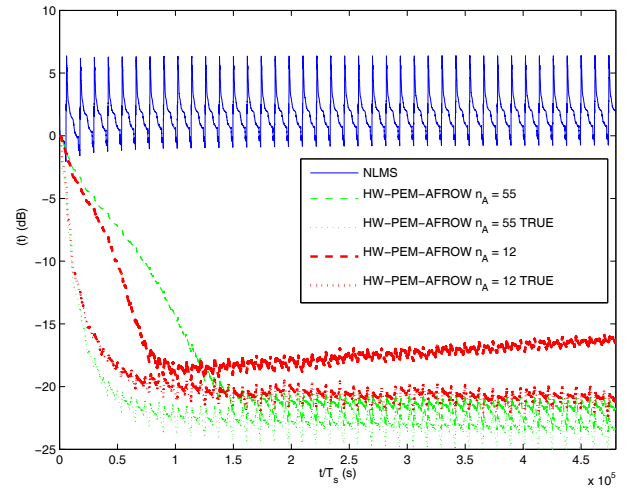Figure 4: Convergence curves of hopping window PEM-AFROW using the Gauss-Newton method.



Figure 5: Convergence curves of hopping window PEM-AFROW using the stochastic gradient method.

It is clear that, whereas an RLS algorithm still performs relatively robust with respect to double-talk, the NLMS algorithm does not converge at all in a continuous double-talk situation. The proposed PEM-AFROW algorithm may outperform the NLMS algorithm with as much as $25dB$. Again the HW variant performs on average somewhat better than the SW variant, but the performance gap is not so large as with the Gauss-Newton method. We also note that some of the PEM-AFROW convergence curves tend to diverge after initial convergence. This is again due to convergence to a local minimum of the prediction error criterion.

## 6. CONCLUSIONS AND FURTHER WORK

We have proposed a new way of coping with a continuous double-talk situation in acoustic echo cancellation. Inspired by linear estimation theory, we have suggested to lower the variance of the RIR estimate by taking into account the near-end signal characteristics. These may be estimated concurrently with the RIR using a two-stage prediction error identification algorithm, by using either a sliding window or a hopping window for linear prediction of the near-end signal. The hopping window variant outperforms the sliding window variant and is computationally cheaper. The proposed method has the potential of improving the echo canceller's convergence during double-talk with $10dB$ resp. $20dB$ as compared to an

ordinary RLS resp. NLMS algorithm.

## REFERENCES

[1] J. Benesty, T. Gänsler, D.R. Morgan, M.M. Sondhi, and S.L. Gay, *Advances in Network and Acoustic Echo Cancellation*, Springer-Verlag, Berlin, Germany, 2001.

[2] S. M. Kay, *Fundamentals of statistical signal processing: estimation theory*, Prentice-Hall Inc., Upper Saddle River, New Jersey, USA, 1993.

[3] G. Rombouts, T. van Waterschoot, K. Struyve, and M. Moonen, "Acoustic feedback cancellation for long acoustic paths using a nonstationary source model," in *Proceedings of the 13th European Signal Processing Conference (EUSIPCO-2005)*, Antalya, Turkey, September 4-8, 2005.

[4] L. Ljung, *System Identification: Theory for the User*, Prentice-Hall Inc., Englewood Cliffs, New Jersey, USA, 1987.