



UNIVERSITÄT PADERBORN
Die Universität der Informationsgesellschaft

FACULTY OF ELECTRICAL ENGINEERING
AND INFORMATION TECHNOLOGY
DEPARTMENT OF COMMUNICATIONS ENGINEERING

TECHNICAL REPORT

**Derivation of the Power Compensation Constant in
the Observation Model for Reverberant Speech in
the Logarithmic Mel Power Spectral Domain**

December 20, 2012

Authors:

Volker Leutnant Alexander Krueger Reinhold Haeb-Umbach
leutnant@nt.uni-paderborn.de krueger@nt.uni-paderborn.de haeb@nt.uni-paderborn.de

CONTENTS

I	Observation model for reverberant-only speech signals	2
II	Acoustic impulse response model	4
III	Derivation of power compensation constant	6
	References	11

The power compensation constant plays an important role in the observation model for reverberant speech signals [1], [2], whose derivation will be repeated here for convenience.

Note that to introduce the used notation and to set this technical report into proper context, Sec. I and Sec. II are directly taken from [1], which has recently been submitted to *IEEE Transactions on Audio, Speech And Language*. The experienced reader may thus skip these sections and directly proceed to Sec. III.

I. OBSERVATION MODEL FOR REVERBERANT-ONLY SPEECH SIGNALS

In the absence of background noise, the discrete-time microphone signal $\bar{s}(l)$ is given by the convolution of the clean speech signal $\bar{x}(l)$ with the acoustic impulse response (AIR) $h(l)$. According to [2], the corresponding STDFT $S(m, k)$ may be expressed as

$$S(m, k) = \sum_{k'=0}^{K-1} \sum_{m'=-L_{H,\ell}}^{L_H} X(m - m', k') h_{k,k'}(m') \quad (1)$$

with

$$h_{k,k'}(m') := \sum_{p'=0}^{L_h-1} h(p') \phi_{k,k'}(m'B - p') \quad (2)$$

and

$$\phi_{k,k'}(l) := e^{j\frac{2\pi}{K}k'l} \sum_{l'=0}^{L_w-1} w_A(l') w_S(l' + l) e^{-j\frac{2\pi}{K}(k-k')l'}. \quad (3)$$

The terms $h_{k,k'}(m')$ will in the following be referred to as cross-band filters for $k \neq k'$ and as band-to-band filters for $k = k'$, as in [3]. The lengths $L_{H,\ell}$ and L_H in (1) are defined by

$$L_{H,\ell} := \left\lfloor \frac{L_w - 1}{B} \right\rfloor, \quad L_H := \left\lfloor \frac{L_h + L_w - 2}{B} \right\rfloor. \quad (4)$$

Further, $w_S(l')$ denotes a synthesis window, which is bi-orthogonal to $w_A(l')$ [4] and has the same support as $w_A(l')$. For the power of $S(m, k)$ we now write

$$|S(m, k)|^2 = C_P \sum_{m'=0}^{L_H} |X(m - m', k)|^2 |h_{k,k}(m')|^2 + E^{(S)}(m, k). \quad (5)$$

The introduced constant C_P and the error term $E^{(S)}(m, k)$ thereby capture all terms incurring when the square of the sum given in (1) is approximated by the sum of the squares while also ensuring a causal relationship by dropping all negative frame indices. The constant C_P will be determined such that the error term $E^{(S)}(m, k)$ is zero-mean, or, equivalently,

$$\mathbb{E} \left[\left| \check{S}(m, k) \right|^2 \right] \stackrel{!}{=} \mathbb{E} \left[C_P \sum_{m'=0}^{L_H} \left| \check{X}(m - m', k) \right|^2 \left| \check{h}_{k,k}(m') \right|^2 \right]. \quad (6)$$

Note that in (6) and in the following, we use the breve mark ($\bar{\cdot}$) to distinguish a random variable from its realization. Due to the role of the constant C_P in (6), it is in the following referred to as the power compensation constant. Besides being additive in the targeted LMPSC domain, choosing a multiplicative constant C_P rather than an additive term to compensate for the bias introduced by the approximation is advantageous, since the desired compensation is made independent of the power of the clean speech signal and that of the AIR.

Note, that the approximation presented here is more general than that in [2], where we chose an empirical value of $(L_w/B)^2$.

By further introducing the mean of $|h_{k,k}(m')|^2$ over the q th mel band, i.e.,

$$\bar{\mathcal{H}}_{m',q} := \frac{1}{K_q^{(u)} - K_q^{(\ell)} + 1} \sum_{k=K_q^{(\ell)}}^{K_q^{(u)}} |h_{k,k}(m')|^2, \quad (7)$$

the MPSCs $\mathcal{S}_{m,q}$ of the reverberant speech signal can be written as

$$\mathcal{S}_{m,q} = C_P \sum_{m'=0}^{L_H} \bar{\mathcal{H}}_{m',q} \mathcal{X}_{m-m',q} + \mathcal{E}_{m,q}^{(S)} \quad (8)$$

$$=: \tilde{\mathcal{S}}_{m,q} + \mathcal{E}_{m,q}^{(S)}. \quad (9)$$

Hereby $\mathcal{X}_{m,q}$ denote the MPSCs of the clean speech signal and $\mathcal{E}_{m,q}^{(S)}$ the error resulting from the approximation of $\mathcal{S}_{m,q}$ by $\tilde{\mathcal{S}}_{m,q}$.

By introducing the logarithmic mel power spectral representation of the AIR

$$\bar{h}_{m',q} := \ln \{ \bar{\mathcal{H}}_{m',q} \} \quad (10)$$

and the LMPSC of the clean speech signal

$$x_{m,q} := \ln \{ \mathcal{X}_{m,q} \}, \quad (11)$$

we are now able to express the LMPSCs of the reverberant speech in terms of the underlying LMPSCs of the clean speech and the LMPSCs of the AIR, i.e.,

$$s_{m,q} := \ln \{ \mathcal{S}_{m,q} \} = \tilde{s}_{m,q} + v_{m,q}^{(s)}, \quad (12)$$

where

$$\tilde{s}_{m,q} := \ln \{ \tilde{\mathcal{S}}_{m,q} \} = \ln \left\{ C_P \sum_{m'=0}^{L_H} e^{x_{m-m',q} + \bar{h}_{m',q}} \right\}. \quad (13)$$

Thereby, employing the definition of $\tilde{\mathcal{S}}_{m,q}$ given in (9) instead of the equivalent formulation given in (13), the additive observation error in the LMPSC domain

$$v_{m,q}^{(s)} := s_{m,q} - \tilde{s}_{m,q} = \ln \{ \mathcal{S}_{m,q} \} - \ln \{ \tilde{\mathcal{S}}_{m,q} \} \quad (14)$$

$$= \ln \left\{ \frac{\mathcal{S}_{m,q}}{\sum_{m'=0}^{L_H} \bar{\mathcal{H}}_{m',q} \mathcal{X}_{m-m',q}} \right\} - \ln \{ C_P \} \quad (15)$$

captures the errors from the approximation of $\mathcal{S}_{m,q}$ by $\tilde{\mathcal{S}}_{m,q}$ in the MPSC domain. Note that the choice of C_P only affects the mean of the observation error $v_{m,q}^{(s)}$.

By introducing the observation mapping

$$f_s (\mathbf{x}_{m-L_H:m}, \bar{\mathbf{h}}_{0:L_H}) := \ln \left\{ C_P \sum_{m'=0}^{L_H} e^{\mathbf{x}_{m-m'} + \bar{\mathbf{h}}_{m'}} \right\}, \quad (16)$$

where the mathematical operations are understood to be performed on the vectors component-wise, the relationship (12) may compactly be formulated by

$$\mathbf{s}_m = f_s (\mathbf{x}_{m-L_H:m}, \bar{\mathbf{h}}_{0:L_H}) + \mathbf{v}_m^{(s)}, \quad (17)$$

where $\mathbf{x}_{m-L_H:m} := \{ \mathbf{x}_{m-L_H}, \dots, \mathbf{x}_m \}$ and $\bar{\mathbf{h}}_{0:L_H} := \{ \bar{\mathbf{h}}_0, \dots, \bar{\mathbf{h}}_{L_H} \}$ denote sequences of $L_H + 1$ Q -dimensional vectors.

II. ACOUSTIC IMPULSE RESPONSE MODEL

The observation model (17) requires the logarithmic mel power spectral representation of the AIR $\bar{\mathbf{h}}_{0:L_H}$. In practice, however, this representation is usually unknown. To avoid a sensitive blind estimation of the AIR for a computation of $\bar{\mathbf{h}}_{0:L_H}$, we have proposed in [2] to employ a stochastic AIR model, which has previously been introduced in [5]. According to this model, the AIR is regarded to be a realization of a stochastic process $\check{h}(l)$ according to

$$\check{h}(l) = \sigma_h \check{v}_h(l) \chi_h(l) e^{-\frac{l}{\tau_h}}, \quad (18)$$

where $\check{v}_h(l)$ is a zero-mean white GAUSSIAN stochastic process of unit power. The indicator function

$$\chi_h(l) := \begin{cases} 1 & \text{for } 0 \leq l \leq L_h - 1 \\ 0 & \text{else} \end{cases} \quad (19)$$

assures the AIR to be causal having a finite length L_h . The term $e^{-\frac{l}{\tau_h}}$ causes an exponentially decaying envelope, where the decay constant τ_h depends on the reverberation time T_{60} through

$$\tau_h = \frac{T_{60}}{3 \ln(10) T_S}, \quad (20)$$

with T_S denoting the sampling duration. The constant σ_h may be used to control the AIR energy. The advantage of using this model is that it has only two parameters, i.e., τ_h and σ_h , which can be estimated more easily than the complete AIR.

Based on the AIR model (18) a reasonable length L_h may be determined in dependence on τ_h by

$$L_h = L_h(\tau_h) = \left\lceil -\frac{\tau_h}{2} \ln(\epsilon_h) \right\rceil, \quad (21)$$

which is obtained by minimizing the AIR length under the constraint that the relative energy of the neglected part of the AIR is smaller than ϵ_h [2].

It has been shown in [6] that the PDFs of the individual components $\check{h}_{m',q}$ of the logarithmic mel power spectral AIR representation $\check{\mathbf{h}}_{0:L_H}$ can be well modeled by GAUSSIANS.

Moreover, their respective means and variances have been found to only depend on the decay constant τ_h , the energy term σ_h and on parameters of the ETSI Standard Front-End.

As we have previously done in [2], we therefore propose to approximate the usually unknown logarithmic mel power spectral representation of the AIR $\bar{\mathbf{h}}_{0:L_H}$ by its mean $\boldsymbol{\mu}_{\check{\mathbf{h}}_{0:L_H}}$ under the AIR model (18). By doing so, the time-variance of the AIR due to, e.g., small movements of the speaker is absorbed through the observation error to a certain degree.

Having introduced the stochastic AIR model, we are able to compute the power compensation constant C_P from condition (6) by using the AIR model and two additional assumptions. We assume the AIR and the clean speech signal to be mutually independent and the latter to be a realization of a zero-mean white GAUSSIAN stochastic process. Under these assumptions, the derived C_P is given by

$$C_P = \frac{C_N}{C_D}, \quad (22)$$

where

$$C_N := K^2 \sum_{m', m'' = -L_H, \ell}^{L_H} \sum_{l=0}^{L_w-1} w_A(l) w_S(l) w_A(l + (m'' - m') B) w_S(l + (m'' - m') B) \cdot \sum_{l' = -L_w + 1}^{L_w - 1} \chi_h(m' B - l') e^{-\frac{2(m' B - l')}{\tau_h}} w_A^2(-l' + l), \quad (23)$$

$$C_D := \left[\sum_{m'=0}^{L_H} \sum_{l' = -L_w + 1}^{L_w - 1} w^2(-l') \chi_h(m' B - l') e^{-\frac{2(m' B - l')}{\tau_h}} \right] \left(\sum_{l=0}^{L_w - 1} w_A^2(l) \right) \quad (24)$$

and

$$w(l) := \sum_{l'=0}^{L_w-1} w_A(l') w_S(l' - l). \quad (25)$$

In practice, it is reasonable to precompute the power compensation constant for a set of predefined reverberation times prior to applying any of the observation models, e.g., to feature enhancement and successive ASR, and choose C_P based on an estimate of the reverberation time.

III. DERIVATION OF POWER COMPENSATION CONSTANT

The constant C_P is introduced to obtain a tractable estimate of the power of the STDFT of the reverberant speech signal. We propose to choose it to satisfy the constraint (6), where the expectation is assumed to be taken over all possible clean speech signals and AIRs, which are possible for fixed constants T_{60} and σ_h^2 . The basis for the derivation is the stochastic AIR model (18). With it, the band-to-band filters $h_{k,k}(m)$ may be expressed by

$$h_{k,k}(m) = \sum_{p'=0}^{L_h-1} h(p') w(p' - mB) e^{j \frac{2\pi}{K} k(mB-p')} \quad (26)$$

using (2), (3) and (25). After applying the variable substitution $l = p' - mB$ to (26), we obtain

$$h_{k,k}(m) = \sum_{l=-L_w+1}^{L_w-1} h(l + mB) w(l) e^{-j \frac{2\pi}{K} kl}. \quad (27)$$

Assuming the AIR to be a realization of the stochastic process $\check{h}(l)$ defined in (18) whose auto-correlation function is given by

$$\mathbb{E} \left[\check{h}(l) \check{h}(l') \right] = \sigma_h^2 \delta(l - l') \chi_h(l) e^{-\frac{2l}{\tau_h}}, \quad (28)$$

with $\delta(l)$ denoting the DIRAC delta function, we can employ (27) and (28) to finally obtain

$$\mathbb{E} \left[\left| \check{h}_{k,k}(m) \right|^2 \right] = \mathbb{E} \left[\left| \sum_{l=-L_w+1}^{L_w-1} \check{h}(l + mB) w(l) e^{-j \frac{2\pi}{K} kl} \right|^2 \right] \quad (29)$$

$$= \sum_{l=-L_w+1}^{L_w-1} \sigma_h^2 \chi_h(l + mB) e^{-\frac{2(l+mB)}{\tau_h}} w^2(l) \quad (30)$$

For a feasible solution we further assume the clean speech signal $x(l)$ to be a realization of a white GAUSSIAN stochastic process $\check{x}(l)$ with power σ_x^2 , which is stochastically independent of the AIR. The auto-correlation function of its STDFT may thus be expressed as

$$\begin{aligned} & \mathbb{E} \left[\check{X}(m - m', k') \check{X}^*(m - m'', k'') \right] \\ &= \sigma_x^2 e^{j \frac{2\pi}{K} k''(m'' - m')} B \sum_{l=0}^{L_w-1} w_A(l) w_A(l + (m'' - m') B) e^{-j \frac{2\pi}{K} (k' - k'') l}. \end{aligned} \quad (31)$$

The auto-correlation function of the cross-band filters can be written with (2) and (28) as

$$\mathbb{E} \left[\check{h}_{k,k'}(m') \check{h}_{k,k''}^*(m'') \right] = \mathbb{E} \left[\sum_{l,l'=0}^{L_h-1} \check{h}(l) \check{h}(l') \phi_{k,k'}(m'B-l) \phi_{k,k''}^*(m''B-l') \right] \quad (32)$$

$$= \sigma_h^2 \sum_{l=0}^{L_h-1} \chi_h(l) \phi_{k,k'}(m'B-l) \phi_{k,k''}^*(m''B-l) e^{-\frac{2l}{\tau_h}}. \quad (33)$$

Exploiting the fact that the support of $\phi_{k,k'}(l)$ is given by $[-L_w + 1, L_w - 1]$, (33) may be reformulated using the variable substitution $l' = m'B - l$ by

$$\begin{aligned} & \mathbb{E} \left[\check{h}_{k,k'}(m') \check{h}_{k,k''}^*(m'') \right] \\ &= \sigma_h^2 \sum_{l'=-L_w+1}^{L_w-1} \phi_{k,k'}(l') \phi_{k,k''}^*(l' + (m'' - m')B) \chi_h(m'B - l') e^{-\frac{2(m'B-l')}{\tau_h}}. \end{aligned} \quad (34)$$

With (31) and (34), the power compensation constant C_P may now be computed by separately considering the left and right hand side of (6).

a) Left hand side of (6): Making use of the stochastic independence of the clean speech signal and the AIR, the left hand side of (6) may be written as

$$\begin{aligned} & \mathbb{E} \left[\left| \check{S}(m, k) \right|^2 \right] \\ &= \sum_{m', m''=-L_{H,\ell}}^{L_H} \sum_{k', k''=0}^{K-1} \mathbb{E} \left[\check{X}(m-m', k') \check{X}^*(m-m'', k'') \right] \mathbb{E} \left[\check{h}_{k,k'}(m') \check{h}_{k,k''}^*(m'') \right]. \end{aligned} \quad (35)$$

By plugging the found relationships (31) and (34) into (35), we arrive at

$$\begin{aligned} \mathbb{E} \left[\left| \check{S}(m, k) \right|^2 \right] &= \sigma_x^2 \sigma_h^2 \sum_{m', m''=-L_{H,\ell}}^{L_H} \sum_{l=0}^{L_w-1} w_A(l) w_A(l + (m'' - m')B) \\ &\quad \cdot \sum_{l'=-L_w}^{L_w-1} \chi_h(m'B - l') e^{-\frac{2(m'B-l')}{\tau_h}} \cdot \xi_{m''-m', l, l', k}, \end{aligned} \quad (36)$$

with

$$\xi_{m''-m', l, l', k} := \sum_{k', k''=0}^{K-1} \phi_{k,k'}(l') \phi_{k,k''}^*(l' + (m'' - m')B) \cdot e^{-j\frac{2\pi}{K}[k'l - k''(l + (m'' - m')B)]}. \quad (37)$$

By further substituting the definition of $\phi_{k,k'}(l)$ into (37) according to (3), we obtain

$$\begin{aligned} & \xi_{m''-m',l,l',k} \\ &= \sum_{k',k''=0}^{K-1} \left[\sum_{p'=0}^{L_w-1} w_A(p') w_S(p'+l') e^{j\frac{2\pi}{K}k'(p'+l')} e^{-j\frac{2\pi}{K}kp'} \right] e^{-j\frac{2\pi}{K}[k'l-k''(l+(m''-m')B)]} \\ & \quad \cdot \left[\sum_{p''=0}^{L_w-1} w_A(p'') w_S(p''+l'+(m''-m')B) e^{-j\frac{2\pi}{K}k''(p''+l'+(m''-m')B)} e^{j\frac{2\pi}{K}kp''} \right] \end{aligned} \quad (38)$$

$$= \sum_{p''=0}^{L_w-1} w_A(p'') w_S(p''+l'+(m''-m')B) e^{j\frac{2\pi}{K}kp''} \sum_{p'=0}^{L_w-1} w_A(p') w_S(p'+l') e^{-j\frac{2\pi}{K}kp'} \psi_{p',p'',l,l'} \quad (39)$$

with

$$\psi_{p',p'',l,l'} := \left[\sum_{k'=0}^{K-1} e^{j\frac{2\pi}{K}k'(p'+l'-l)} \right] \left[\sum_{k''=0}^{K-1} e^{-j\frac{2\pi}{K}k''(p''+l'-l)} \right]. \quad (40)$$

Considering the sum orthogonality

$$\frac{1}{K} \sum_{k=0}^{K-1} e^{j\frac{2\pi}{K}k\mu} = \sum_{\nu=-\infty}^{\infty} \delta(\mu - \nu K) \text{ for } \mu \in \mathbb{Z}, K \in \mathbb{N}, \quad (41)$$

the term $\psi_{p',p'',l,l'}$ can be expressed as

$$\psi_{p',p'',l,l'} = K^2 \left[\sum_{\nu'=-\infty}^{\infty} \delta(p'+l'-l-\nu'K) \right] \left[\sum_{\nu''=-\infty}^{\infty} \delta(p''+l'-l-\nu''K) \right]. \quad (42)$$

Having in mind the identity

$$\delta(l-\mu)\delta(l-\mu') = \delta(l-\mu)\delta(\mu-\mu') \text{ for } l, \mu, \mu' \in \mathbb{Z}, \quad (43)$$

we may simplify $\psi_{p',p'',l,l'}$ further by

$$\psi_{p',p'',l,l'} = K^2 \left[\sum_{\nu'=-\infty}^{\infty} \delta(p'+l'-l-\nu'K) \right] \left[\sum_{\nu''=-\infty}^{\infty} \delta(p''-p'-(\nu''-\nu')K) \right]. \quad (44)$$

Since the difference $p''-p'$ lies within the interval $[-L_w+1, L_w-1]$ and $K > L_w$ holds, the argument of the second DIRAC delta function in (44) may become zero only for $\nu'' = \nu'$. For that reason

$$\psi_{p',p'',l,l'} = K^2 \sum_{\nu'=-\infty}^{\infty} \delta(p'+l'-l-\nu'K) \delta(p''-p') \quad (45)$$

holds. By substituting (45) into (39), it follows that

$$\begin{aligned} & \xi_{m''-m',l,l',k} \\ &= K^2 \sum_{\nu'=-\infty}^{\infty} \sum_{p'=0}^{L_w-1} w_A(p') w_S(p'+l') \delta(p'+l'-l-\nu'K) \\ & \quad \cdot \sum_{p''=0}^{L_w-1} w_A(p'') w_S(p''+l'+(m''-m')B) \delta(p''-p') e^{-j\frac{2\pi}{K}k(p'-p'')} \end{aligned} \quad (46)$$

$$= K^2 \sum_{\nu'=-\infty}^{\infty} \sum_{p'=0}^{L_w-1} w_A^2(p') w_S(p'+l') \delta(p'+l'-l-\nu'K) w_S(p'+l'+(m''-m')B) \quad (47)$$

$$= K^2 \sum_{\nu'=-\infty}^{\infty} w_A^2(-l'+l+\nu'K) w_S(l+\nu'K) w_S(l+\nu'K+(m''-m')B), \quad (48)$$

From this expression it can be seen that $\xi_{m''-m',l,l',k}$ does not depend on k . Since for $l \in [-L_w+1, L_w-1]$ the equality $w_S(l+\nu'K) = 0 \forall \nu' \neq 0$ holds, we have

$$\xi_{m''-m',l,l',k} = K^2 w_A^2(-l'+l) w_S(l) w_S(l+(m''-m')B). \quad (49)$$

By substituting (49) into (36), we obtain

$$\mathbb{E} \left[\left| \check{S}(m, k) \right|^2 \right] = \sigma_x^2 \sigma_h^2 \cdot C_N \quad (50)$$

with

$$\begin{aligned} C_N := K^2 & \sum_{m', m'' = -L_H, \ell}^{L_H} \sum_{l=0}^{L_w-1} w_A(l) w_S(l) w_A(l+(m''-m')B) w_S(l+(m''-m')B) \\ & \cdot \sum_{l'=-L_w+1}^{L_w-1} \chi_h(m'B-l') e^{-\frac{2(m'B-l')}{\tau_h}} w_A^2(-l'+l). \end{aligned} \quad (51)$$

b) Right hand side of (6): The expression on the right hand side of (6) may be simplified using (31) and (34) to be

$$\begin{aligned} & \mathbb{E} \left[C_P \cdot \sum_{m'=0}^{L_H} \left| \check{X}(m-m', k) \right|^2 \left| \check{h}_{k,k}(m') \right|^2 \right] \\ &= C_P \sum_{m'=0}^{L_H} \mathbb{E} \left[\left| \check{X}(m-m', k) \right|^2 \right] \mathbb{E} \left[\left| \check{h}_{k,k}(m') \right|^2 \right] \end{aligned} \quad (52)$$

$$= C_P \sum_{m'=0}^{L_H} \left(\sigma_x^2 \sum_{l=0}^{L_w-1} w_A^2(l) \right) \left(\sigma_h^2 \sum_{l'=-L_w+1}^{L_w-1} |\phi_{k,k}(l')|^2 \chi_h(m'B-l') e^{-\frac{2(m'B-l')}{\tau_h}} \right) \quad (53)$$

$$= C_P \cdot \sigma_x^2 \sigma_h^2 \cdot C_D. \quad (54)$$

Employing (3) for $k=k'$ and (25), C_D may be expressed as

$$C_D := \left[\sum_{m'=0}^{L_H} \sum_{l'=-L_w+1}^{L_w-1} w^2(-l') \chi_h(m'B - l') e^{-\frac{2(m'B-l')}{\tau_h}} \right] \left(\sum_{l=0}^{L_w-1} w_A^2(l) \right). \quad (55)$$

The desired constant C_P is finally obtained from comparing expressions (54) and (50) to be

$$C_P = \frac{C_N}{C_D}. \quad (56)$$

From (51) and (55), it can be seen that the constant C_P only depends on the parameters employed for feature extraction and on the reverberation time.

The power compensation constant obtained for the parameter values according to the ETSI Standard Front-End [7] are given in Fig. 1 for different reverberation times. The required AIR parameters for different reverberation times resulting from the choice $\epsilon_h = 10^{-3}$ are listed in Tab. I It can be seen that

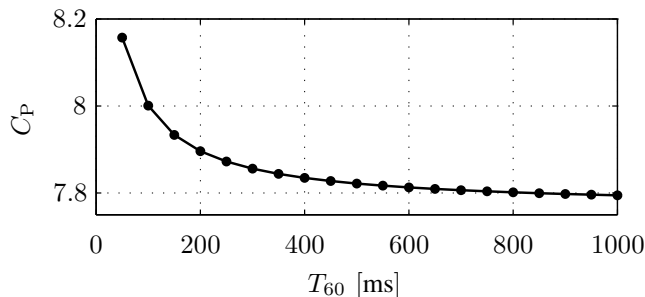


Fig. 1. Power compensation constant for different reverberation times.

the constant depends only on the parameters for the extraction, i.e. the analysis and synthesis windows, the number of frequency bins etc., and on the reverberation time and that its value is about 8 for a large range of practically relevant reverberation times.

In practice, it is reasonable to precompute the power compensation constant for a set of predefined reverberation times prior to applying any of the observation models, e.g., to feature enhancement and successive ASR, and choose C_P based on an estimate of the reverberation time.

TABLE I

AIR PARAMETERS FOR DIFFERENT REVERBERATION TIMES RESULTING FROM THE CHOICE $\epsilon_h = 10^{-3}$.

\hat{T}_{60}	250 ms	350 ms	450 ms	550 ms	650 ms
AIR length \hat{L}_h	1000	1400	1800	2200	2600
LMPSC length \hat{L}_H	14	19	24	29	34

REFERENCES

- [1] V. Leutnant, A. Krueger, and R. Haeb-Umbach, "A new observation model in the logarithmic mel power spectral domain for the automatic recognition of noisy reverberant speech," Sep. 2012, submitted for publication in *IEEE Trans. Audio, Speech, Lang. Process.*
- [2] A. Krueger and R. Haeb-Umbach, "Model-based feature enhancement for reverberant speech recognition," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 7, pp. 1692–1707, 2010.
- [3] Y. Avargel and I. Cohen, "System identification in the short-time Fourier transform domain with crossband filtering," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 15, no. 4, pp. 1305–1319, May 2007.
- [4] S. Qian and D. Chen, "Discrete Gabor transform," *IEEE Trans. Signal Process.*, vol. 41, no. 7, pp. 2429–2438, Jul. 1993.
- [5] J. Polack, "La transmission de l'énergie sonore dans les salles," Dissertation, Université du Maine, 1988.
- [6] A. Krüger, "Modellbasierte Merkmalsverbesserung zur robusten automatischen Spracherkennung in Gegenwart von Nachhall und Hintergrundstörungen," Ph.D. dissertation, University of Paderborn, Germany, Dec. 2011.
- [7] ETSI, *ETSI standard document, Speech Processing, Transmission and Quality Aspects (STQ); Distributed speech recognition; Front-end feature extraction algorithm; Compression algorithms, ETSI ES 201 108 V1.1.3 (2003-09)*.