

## Sound Zones as an Optimal Filtering Problem

Nielsen, Jesper Kjær; Lee, Taewoong; Jensen, Jesper Rindom; Christensen, Mads Græsbøll

*Published in:*

2018 52nd Asilomar Conference on Signals, Systems, and Computers

*DOI (link to publication from Publisher):*

[10.1109/ACSSC.2018.8645268](https://doi.org/10.1109/ACSSC.2018.8645268)

*Creative Commons License*

Unspecified

*Publication date:*

2018

*Document Version*

Other version

[Link to publication from Aalborg University](#)

*Citation for published version (APA):*

Nielsen, J. K., Lee, T., Jensen, J. R., & Christensen, M. G. (2018). Sound Zones as an Optimal Filtering Problem. In *2018 52nd Asilomar Conference on Signals, Systems, and Computers* IEEE (Institute of Electrical and Electronics Engineers). <https://doi.org/10.1109/ACSSC.2018.8645268>

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

### Take down policy

If you believe that this document breaches copyright please contact us at [vbn@aub.aau.dk](mailto:vbn@aub.aau.dk) providing details, and we will remove access to the work immediately and investigate your claim.

# SOUND ZONES AS AN OPTIMAL FILTERING PROBLEM

*Jesper Kjær Nielsen, Taewoong Lee, Jesper Rindom Jensen, and Mads Græsbøll Christensen*

Audio Analysis Lab, CREATE  
Aalborg University, Denmark  
{jkn, tlee, jrj, mgc}@create.aau.dk

## ABSTRACT

The design of so-called sound zones is a fairly old idea which has recently gained some research attention. The idea is to control a loudspeaker array to reproduce a desired sound field in certain regions. In this paper, we show how the sound zone control problem can be solved using techniques from speech enhancement. Specifically, we first describe in detail the recently introduced variable span linear filtering (VSLF) framework which unifies the popular optimal filtering and subspace-based approaches to speech enhancement. We then show how the sound zone control problem can be solved using the VSLF framework and how a number of well-known sound zone control methods can be viewed as special cases of this solution. We also discuss in detail the differences between the speech enhancement problem and the sound zone control problem and argue that the VSLF framework is actually even better suited for controlling sound zones than for enhancing speech.

**Index Terms**— Speech enhancement, variable span linear filtering, sound zones, acoustic contrast control, pressure matching

## 1. INTRODUCTION

Sometimes multiple persons either wish or have to share the same acoustic environment while simultaneously enjoying different audio programs. In, e.g., a car cabin [1, 2], the parents in the front seats might wish to listen to music while the children in the rear seats might wish to watch a cartoon. Another example is an outdoor festival [3] where multiple concerts might happen in parallel, often within hearing range of a residential area. In these and other applications [4–6], simply using headphones to obtain the desired effect either hampers social interactions or is infeasible. An alternative, yet immature technology, is the sound zone concept where the desired audio programs are reproduced in spatially confined regions, i.e., sound zones, by controlling a loudspeaker array. This idea was originatively conceived two decades ago [7] and has since attracted some research attention, mainly in the field of acoustics. A number of different methods have been proposed based on various design principles (see [8] for a recent overview), and they all seek to reproduce the desired audio program in a given zone with the least amount of distortion while minimising the leakage to the other zones.

The dominant sound zone control strategy is to pre-filter the loudspeaker signals with fixed FIR filters. Designing sound zones, therefore, essentially boils down to designing these filters so that they simultaneously optimise metrics related to signal distortion and zone leakage. In [9], the acoustic contrast control (ACC) method was introduced, and it sought to minimise the zone leakage by maximising the ratio of acoustic potential energies between a so-called bright and dark zone. The introduced notions of bright and dark

zones were a convenient abstraction which allowed the control filters to be designed for each source and zone in isolation. That is, one zone was in turn considered to be the bright zone whereas the remaining zones were lumped together as a dark or quiet zone, and the multi-zone solution was obtained as a superposition of the individual solutions. The main disadvantage of the ACC approach is that the control filters are designed without penalising perturbations of the desired sound field in the bright zone. This was addressed in [10] with the introduction of the pressure matching (PM) method where the signal distortion in the bright zone was explicitly penalised in the control filter design. However, this improvement in the reconstruction quality of the desired sound field came at the expense of a much higher zone leakage, i.e., a smaller acoustic contrast between the bright and dark zone. Therefore, a number of methods [11–16] have sought to modify and/or combine ACC and PM so that a better trade-off between signal distortion and zone leakage is obtained.

Compared to sound zone control, speech enhancement is an older and much more mature field which is concerned with methods for retrieving a clean speech signal from a noisy mixture. Although this problem appears to be fundamentally different to that of designing sound zones, the main goal in speech enhancement is also to strike the balance between two conflicting requirements which are the speech distortion and the noise suppression [17]. More specifically, most speech enhancement methods involve the design of an FIR filter which filters out the noise from the noisy speech input signal while simultaneously preserving the clean speech signal in the output signal. This approach to doing speech enhancement is often referred to as optimal filtering [18] and, as the title of the present paper suggests, the optimal filtering principle can actually be adopted for the design of sound zones. We recently discovered this connection and used this insight to adapt the very flexible variable span linear filtering (VSLF) framework from speech enhancement [17, 19] to sound zone control [20]. Interestingly, the resulting sound zone control framework (VAST) has many of the existing control methods such as ACC, PM, and their variations as special cases. Moreover, the framework can be used to show theoretically that ACC gives the maximum acoustic contrast and distortion while PM gives the minimum acoustic contrast and distortion, something which had only been supported by empirical evidence prior to the introduction of VAST. In this paper, we further explore the connections between those speech enhancement and sound zone control methods which are formulated as optimal filtering problems. We start by a detailed description of the VSLF framework in the context of single channel speech enhancement in Sec. 2 since VAST becomes much easier to derive and understand when rooted in a good understanding of the VSLF framework. The VAST framework is described in Sec. 3 where we also describe its special cases. Finally, we discuss the differences between the VSLF and VAST frameworks in Sec. 4.

## 2. SPEECH ENHANCEMENT USING VSLF

Speech enhancement is an important, but also difficult problem which has attracted significant research attention in the speech processing community for many decades (see [21, 22] for fairly recent overviews). Although many variations of the problem exist, we here restrict our attention to its most basic form since this is sufficient to convey our main points. Specifically, we consider the single channel problem where we wish to extract the clean speech sample  $s(n) \in \mathbb{R}$  from a noisy mixture  $y(n) \in \mathbb{R}$ , i.e.,

$$y(n) = s(n) + e(n) \quad (1)$$

where  $e(n) \in \mathbb{R}$  is additive noise. There is a number of reasons for why estimating  $s(n)$  from  $y(n)$  is difficult. First, the problem is under-determined in the absence of any prior information since we have two latent variables for every observation. Second, specific prior information is often not available since the speech and noise characteristics can vary significantly from person to person and from noise environment to noise environments. Some noise types such as babble noise are even very speech-like since it is mostly composed of multiple talking persons. Third, reverberation results in that the noise component and the direct-path speech component become correlated which is difficult to take into account in general and is, therefore, often ignored. Fourth, and finally, speech and some noise types are highly non-stationary which means that segment-by-segment processing under a local stationarity assumption should be performed for short segments.

Although a number of different approaches to speech enhancement exist based on, e.g., spectral subtraction [23], statistical models [24], binary masking [25], subspace techniques [26], and non-negative matrix factorisations [27], the majority of enhancement methods can be classified as an optimal filtering approach [18] or has such an interpretation. The basic idea in the optimal filtering approach is to design an FIR filter  $\mathbf{h} \in \mathbb{R}^M$  which filters out the noise and preserves the clean speech from the noisy speech signal, i.e.,

$$\hat{s}(n) = \mathbf{h}^T \mathbf{y}(n) = \mathbf{h}^T \mathbf{s}(n) + \mathbf{h}^T \mathbf{e}(n) \quad (2)$$

where  $\hat{s}(n)$  is an estimate of the  $n$ th clean speech sample and  $\mathbf{y}(n) \in \mathbb{R}^M$  is defined as

$$\mathbf{y}(n) = [y(n) \quad y(n-1) \quad \cdots \quad y(n-M+1)]^T. \quad (3)$$

Note that  $\mathbf{s}(n)$  and  $\mathbf{e}(n)$  are defined analogously to  $\mathbf{y}(n)$ . From (2), we see that the filter  $\mathbf{h}$  should be defined so that the clean speech component is passed with as little distortion as possible, i.e.,  $\mathbf{h}^T \mathbf{s}(n) \approx s(n)$  and so that as much noise as possible is filtered out, i.e.,  $\mathbf{h}^T \mathbf{e}(n) \approx 0$ . Unfortunately, these two design criteria are in general conflicting since requiring no speech distortion leads to no noise reduction and, conversely, requiring complete noise reduction leads to that  $\hat{s}(n) = 0$ . Therefore, a compromise must be made and the various optimal filtering methods are different ways of trading-off speech distortion for noise suppression. Unfortunately, it is far from obvious how these errors should be measured in terms of perceptual meaningful and practical metrics [28] so they are usually simply measured using the mean squared error criterion for mathematical tractability. Specifically, if we denote the speech distortion by  $J_{SD}(\mathbf{h})$  and the noise suppression by  $J_{NS}(\mathbf{h})$ , their definitions are

$$\begin{aligned} J_{SD}(\mathbf{h}) &= \mathbb{E} \left[ \left( s(n) - \mathbf{h}^T \mathbf{s}(n) \right)^2 \right] \\ &= \mathbf{h}^T \mathbf{R}_s \mathbf{h} - 2\mathbf{h}^T \mathbf{r}_s + \boldsymbol{\iota}_1^T \mathbf{r}_s \end{aligned} \quad (4)$$

$$J_{NS}(\mathbf{h}) = \mathbb{E} \left[ \left( 0 - \mathbf{h}^T \mathbf{e}(n) \right)^2 \right] = \mathbf{h}^T \mathbf{R}_e \mathbf{h} \quad (5)$$

where  $\mathbb{E}[\cdot]$ ,  $\mathbf{R}_s = \mathbb{E}[\mathbf{s}(n)\mathbf{s}^T(n)]$ , and  $\mathbf{R}_e = \mathbb{E}[\mathbf{e}(n)\mathbf{e}^T(n)]$  are the expectation operator, the clean speech covariance matrix, and the noise covariance matrix, respectively. Moreover,  $\mathbf{r}_s = \mathbf{R}_s \boldsymbol{\iota}_1$  is a correlation vector with  $\boldsymbol{\iota}_1 = [1 \quad 0 \quad \cdots \quad 0]^T$ . Note that we have here assumed that the clean speech and noise components are zero mean and wide sense stationary (WSS) stochastic processes. Whereas the former assumption is reasonable for audio signals, the latter only approximately holds for a short speech segments of length  $N \geq M$ , say. Consequently, the optimal filtering methods are typically implemented on a segment-by-segment basis, and we here focus on the processing for just one such segment. We also limit ourselves to the case where the statistics is assumed known prior knowledge. For practical speech enhancement algorithms, this is, of course, unrealistic, but it turns out to be fulfilled for the sound zone control problem.

Based on the simple objective functions for speech distortion and noise suppression defined in (4) and (5), respectively, we can formulate a combined optimal filtering objective as

$$J_{OF}(\mathbf{h}) = J_{SD}(\mathbf{h}) + \mu J_{NS}(\mathbf{h}) \quad (6)$$

where  $\mu$  is a non-negative scalar. An interpretation of  $\mu$  and the combined objective  $J_{OF}(\mathbf{h})$  is that they are the Lagrange multiplier and Lagrangian, respectively, associated with the convex problem

$$\begin{aligned} &\underset{\mathbf{h} \in \mathbb{R}^M}{\text{minimise}} && J_{SD}(\mathbf{h}) \\ &\text{subject to} && J_{NS}(\mathbf{h}) \leq \beta \boldsymbol{\iota}_1^T \mathbf{r}_s \end{aligned} \quad (7)$$

where  $\beta \in [0, 1]$  is a user parameter which controls how important noise suppression is relative to speech distortion. Unfortunately, the relationship between  $\mu$  and  $\beta$  is not simple, except for in a few extreme special cases<sup>1</sup>. Instead of  $\beta$ , the Lagrange multiplier  $\mu$  is, therefore, often treated as a user parameter since this also turns the joint and iterative optimisation over both  $\mathbf{h}$  and  $\mu$  into a simple quadratic optimisation problem having a simple analytical solution.

Although  $\mu$  in (6) provides us with a handle for controlling the trade-off between speech distortion and noise suppression, we can obtain an even better control over this trade-off by essentially enforcing a low rank approximation to the clean speech covariance matrix  $\mathbf{R}_s$ . This is the main rationale behind the variable span linear filtering (VSLF) framework [17, 19] and also the principle which unifies the optimal filtering and subspace approaches to speech enhancement. To derive this, first consider the joint diagonalisation (also sometimes referred to as the generalised eigenvalue decomposition or matrix pencil) of the positive semi-definite clean speech covariance matrix  $\mathbf{R}_s$  and positive definite noise covariance matrix  $\mathbf{R}_e$  [29, p. 106]

$$\mathbf{B}^T \mathbf{R}_s \mathbf{B} = \boldsymbol{\Lambda} \quad (8)$$

$$\mathbf{B}^T \mathbf{R}_e \mathbf{B} = \mathbf{I}_M \quad (9)$$

where  $\mathbf{I}_M$  is the  $M \times M$  identity matrix,  $\boldsymbol{\Lambda}$  is a diagonal matrix containing the  $M$  non-negative and real-valued eigenvalues of  $\mathbf{R}_e^{-1} \mathbf{R}_s$  in descending order, and  $\mathbf{B} \in \mathbb{R}^{M \times M}$  is a non-singular matrix containing the eigenvectors of  $\mathbf{R}_e^{-1} \mathbf{R}_s$  sorted according to the eigenvalues in  $\boldsymbol{\Lambda}$ . Since  $\mathbf{B}$  is a non-singular matrix, any filter vector  $\mathbf{h}$  can be written as a linear combination of the eigenvectors in  $\mathbf{B}$  in (4), i.e.,

$$\mathbf{h} = \mathbf{B} \mathbf{a} \quad (10)$$

<sup>1</sup>For example, setting  $\mu = 0$  corresponds to setting  $\beta = 1$ , and letting  $\mu \rightarrow \infty$  corresponds to letting  $\beta \rightarrow 0^+$ . However, these special cases are uninteresting since they just produce the trivial solutions  $\mathbf{h} = \mathbf{0}$  and  $\mathbf{h} = \boldsymbol{\iota}_1$ , respectively.

where  $\mathbf{a} \in \mathbb{R}^M$  contains the weights. By inserting this expression for  $\mathbf{h}$  in the objective functions for the speech distortion and noise suppression in (4) and (5), respectively, we readily obtain

$$J_{SD}(\mathbf{B}\mathbf{a}) = \mathbf{a}^T \mathbf{\Lambda} \mathbf{a} - 2\mathbf{a}^T \mathbf{\Lambda} \mathbf{B}^{-1} \mathbf{t}_1 + \mathbf{t}_1^T \mathbf{r}_s \quad (11)$$

$$J_{NS}(\mathbf{B}\mathbf{a}) = \mathbf{a}^T \mathbf{a}. \quad (12)$$

To understand how the joint diagonalisation allows us to trade off speech distortion for noise suppression, we now make a  $V (\leq M)$ -rank approximation to  $\mathbf{R}_s$  by forcing the  $M - V$  smallest eigenvalues to 0. Inserting this  $V$ -rank approximation in the speech distortion objective in (4) gives

$$J_{SD}(\mathbf{B}\mathbf{a}) \approx \mathbf{a}_V^T \mathbf{\Lambda}_V \mathbf{a}_V - 2\mathbf{a}_V^T \mathbf{\Lambda}_V \mathbf{U}_V^T \mathbf{t}_1 + \text{const.} \quad (13)$$

where  $\mathbf{\Lambda}_V$  and  $\mathbf{U}_V$  are the upper left  $V \times V$  submatrix of  $\mathbf{\Lambda}$  and the first  $V$  columns of  $\mathbf{U} = \mathbf{B}^{-T}$ , respectively. Note that the approximation is exact if  $\text{rank}(\mathbf{R}_s) = V$ . From (13), we see that the speech distortion objective only depends on the first  $V$  elements, denoted as  $\mathbf{a}_V$ , of the vector  $\mathbf{a}$ . As seen from (12), the remaining  $M - V$  elements only increase the noise suppression objective and should, therefore, be set to zero. Thus, if we define the filter vector  $\mathbf{h}$  as a linear combination of the first  $V$  eigenvectors, i.e.,

$$\mathbf{h} = \mathbf{B}_V \mathbf{a}_V, \quad (14)$$

we make the  $V$ -rank approximation to  $\mathbf{R}_s$  and obtain the solution for  $\mathbf{h}$  resulting in the largest noise suppression among all filter vectors satisfying the under-determined set of linear equations  $\mathbf{a}_V = \mathbf{U}_V^T \mathbf{h}$ .

We are now finally able to derive the VSLF optimal filter. To do that, we insert the expression for the filter vector in (14) into the combined objective function in (6) and obtain

$$\mathbf{a}_{VSLF}(V, \mu) = \underset{\mathbf{a}_V \in \mathbb{R}^M}{\text{argmin}} J_{OF}(\mathbf{B}_V \mathbf{a}_V) = (\mathbf{\Lambda}_V + \mu \mathbf{I}_V)^{-1} \mathbf{B}_V^T \mathbf{r}_s$$

from which it readily follows that

$$\mathbf{h}_{VSLF}(V, \mu) = \mathbf{B}_V \mathbf{a}_{VSLF}(V, \mu) = \sum_{v=1}^V \frac{\mathbf{b}_v^T \mathbf{r}_s}{\lambda_v + \mu} \mathbf{b}_v \quad (15)$$

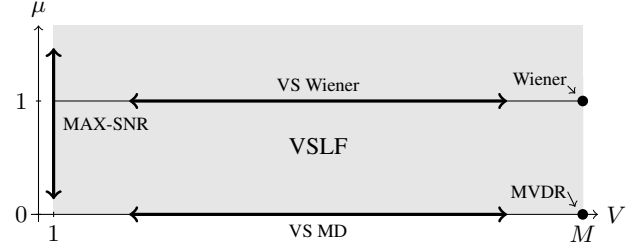
where  $\lambda_v$  and  $\mathbf{b}_v$  are the  $v$ th eigenvalue and eigenvector, respectively. Note that  $V$  and  $\mu$  are user-defined parameters which control the trade-off between the speech distortion and noise suppression in different ways. Moreover, note that  $\mathbf{b}_v^T \mathbf{r}_s$  is proportional to  $\lambda_v$ . Consequently, the elements in  $\mathbf{h}_{VSLF}(V, \mu)$  are automatically set to 0 if  $\lambda_v = 0$ , unless  $\mu = 0$ , which means that the VSLF solution is the same for all  $V \geq \text{rank}(\mathbf{R}_s)$ . Interestingly, we can also use the VSLF solution to derive simple expressions for the speech distortion and the noise suppression. By inserting the VSLF solution in (15) into the objectives in (4) and (5), we obtain

$$J_{SD}(\mathbf{h}_{VSLF}(V, \mu)) = \mathbf{t}_1^T \mathbf{r}_s - \sum_{v=1}^V \frac{\lambda_v + 2\mu}{(\lambda_v + \mu)^2} |\mathbf{b}_v^T \mathbf{r}_s|^2 \quad (16)$$

$$J_{NS}(\mathbf{h}_{VSLF}(V, \mu)) = \sum_{v=1}^V \frac{1}{(\lambda_v + \mu)^2} |\mathbf{b}_v^T \mathbf{r}_s|^2. \quad (17)$$

Thus, the speech distortion and noise suppression are non-increasing or non-decreasing, respectively, for an increasing  $V$ , thus confirming the trade-off between speech distortion and noise suppression theoretically. Moreover, the combined objective reduces to

$$J_{OF}(\mathbf{h}_{VSLF}(V, \mu)) = \mathbf{t}_1^T \mathbf{r}_s - \sum_{v=1}^V \frac{1}{\lambda_v + \mu} |\mathbf{b}_v^T \mathbf{r}_s|^2 \quad (18)$$



**Fig. 1.** An illustration of how some of the traditional optimal filtering enhancement algorithms are special cases of the VSLF framework.

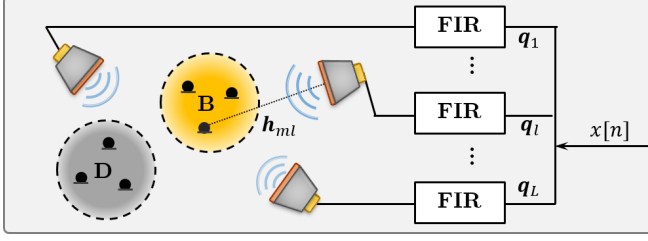
which is non-increasing with increasing  $V$ . A consequence of this is that we should always use all  $M$  eigenvectors if we want to minimise  $J_{OF}(\mathbf{h}_{VSLF}(V, \mu))$  over  $V$ , regardless of how  $\mu$  is selected. At first, this observation might seem to render the joint diagonalisation a mere academic exercise since the VSLF solution for  $V = M$  can be obtained directly as the minimiser of the objective in (6). However, it is important to remember that the above analysis is based on the unrealistic case of known statistics. Moreover, the objective function is not necessarily a good measure of the perceptual quality of the enhanced speech signal. Consequently, the joint diagonalisation is still a very useful tool for trading-off signal distortion for noise suppression when tuning a speech enhancement algorithm.

In summary, the key strength of the VSLF solution is the explicit control over the trade-off between the speech distortion and noise suppression through the user parameters  $V$  and  $\mu$ . Interestingly, many of the well-known optimal filtering methods such as the maximum SNR ( $V = 1$ ), Wiener ( $V = M$  and  $\mu = 1$ ), and MVDR ( $V = M$  and  $\mu = 0$ ) filters are also special cases of the VSLF solution, which is illustrated in Fig. 1. For a much more thorough discussion on the special cases of the VSLF solution, we refer the interested reader to [17, 19].

### 3. SOUND ZONE CONTROL USING VAST

On first sight, the discussion in the previous section on optimal filtering based speech enhancement might seem unrelated to the sound zone control problem. As alluded to in the introduction, however, one immediate parallel between the two problems is the trade-off between two conflicting design criteria which for the sound zone control problem are the signal distortion in the bright zone and the leakage to the dark zone. To define these criteria mathematically, we use the sketch in Fig. 2. Here, 'B' and 'D' denote the bright and dark zones, respectively, which are each formed by grouping a collection of microphones or control points. We index each of these control points using  $m$  and define the collection of bright and dark zone indices as  $\mathcal{M}_B$  and  $\mathcal{M}_D$ , respectively. In the bright zone, we wish to reproduce some desired signal  $d_m(n) \in \mathbb{R}$  where  $m \in \mathcal{M}_B$  whereas the control points in the dark zone, i.e.,  $m \in \mathcal{M}_D$ , should ideally measure 0. For a control point in either the bright or dark zone, the reproduced sound pressure is a summation of  $L$  signals which, for each loudspeaker and control filter index  $l \in \{1, \dots, L\}$ , are obtained as a convolution between the known input signal  $x(n) \in \mathbb{R}$ , the unknown control filter  $\mathbf{q}_l \in \mathbb{R}^J$ , and the known<sup>2</sup> room impulse

<sup>2</sup>We assume that these impulse response have been pre-measured or can be accurately simulated.



**Fig. 2.** A sketch of how a loudspeaker array is controlled via a set of FIR control filters to produce a bright and a dark zone.

response  $\mathbf{h}_{ml} \in \mathbb{R}^K$ , i.e.,

$$p_m(n) = \sum_{l=1}^L \sum_{k=0}^{K-1} \sum_{j=0}^{J-1} x(n-k-j) h_{ml}(k) q_l(j) \quad (19)$$

$$= \sum_{l=1}^L \mathbf{q}_l^T \mathbf{y}_{ml}(n) = \mathbf{q}^T \mathbf{y}_m(n) \quad (20)$$

where we have defined

$$\mathbf{y}_{ml}(n) = \mathbf{X}(n) \mathbf{h}_{ml} \quad (21)$$

$$\mathbf{X}(n) = \begin{bmatrix} x(n) & \cdots & x(n-K+1) \\ \vdots & \ddots & \vdots \\ x(n-J+1) & \cdots & x(n-K-J+2) \end{bmatrix} \quad (22)$$

$$\mathbf{y}_m(n) = [\mathbf{y}_{m1}^T(n) \cdots \mathbf{y}_{mL}^T(n)]^T \quad (23)$$

$$\mathbf{q} = [\mathbf{q}_1^T \cdots \mathbf{q}_L^T]^T. \quad (24)$$

Note that  $\mathbf{y}_{ml}(n) \in \mathbb{R}^J$  is known and what we refer to as the uncontrolled pressures at control point  $m$  originating from loudspeaker  $l$ . The control filters in  $\mathbf{q}$ , however, are unknown and what we are trying to design. To do that, we define the signal distortion objective  $J_B(\mathbf{q})$  for the bright zone and the residual error objective  $J_D(\mathbf{q})$  as

$$J_B(\mathbf{q}) = \frac{1}{N|\mathcal{M}_B|} \sum_{n=0}^{N-1} \sum_{m \in \mathcal{M}_B} \|d_m(n) - \mathbf{q}^T \mathbf{y}_m(n)\|_2^2 \quad (25)$$

$$J_D(\mathbf{q}) = \frac{1}{N|\mathcal{M}_D|} \sum_{n=0}^{N-1} \sum_{m \in \mathcal{M}_D} \|0 - \mathbf{q}^T \mathbf{y}_m(n)\|_2^2. \quad (26)$$

If we now define the positive semi-definite spatial covariance matrices  $\mathbf{R}_B$  and  $\mathbf{R}_D$  as well as the spatial correlation vector  $\mathbf{r}_B$  as

$$\mathbf{R}_C = \frac{1}{N|\mathcal{M}_C|} \sum_{n=0}^{N-1} \sum_{m \in \mathcal{M}_C} \mathbf{y}_m(n) \mathbf{y}_m^T(n) \quad \text{for } C \in \{B, D\} \quad (27)$$

$$\mathbf{r}_B = \frac{1}{N|\mathcal{M}_B|} \sum_{n=0}^{N-1} \sum_{m \in \mathcal{M}_B} \mathbf{y}_m(n) d_m(n), \quad (28)$$

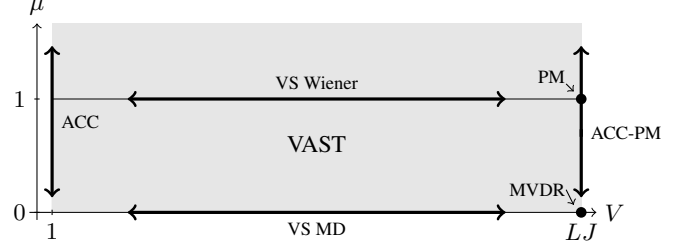
we readily obtain that

$$J_B(\mathbf{q}) = \mathbf{q}^T \mathbf{R}_B \mathbf{q} - 2\mathbf{q}^T \mathbf{r}_B + \sigma_d^2 \quad (29)$$

$$J_D(\mathbf{q}) = \mathbf{q}^T \mathbf{R}_D \mathbf{q} \quad (30)$$

where we have also defined the constant

$$\sigma_d^2 = \frac{1}{N|\mathcal{M}_B|} \sum_{n=0}^{N-1} \sum_{m \in \mathcal{M}_B} d_m^2(n). \quad (31)$$



**Fig. 3.** An illustration of how some of the traditional sound zone methods are special cases of the VAST framework.

By comparing the objectives in (29) and (30) to the objectives in (4) and (5), we see that they have exactly the same quadratic form. Thus, the sound zone control filter design can be viewed as an optimal filtering problem and be solved using the VSLF framework, provided that  $\mathbf{R}_D$  has full rank. From (27), we see that  $N|\mathcal{M}_D|$  must be at least greater than or equal to  $LJ$  to satisfy this.

To differentiate the VSLF solution for sound zones from the VSLF solution for speech enhancement, we introduced the VAST acronym in [20]. Analogously to (15), the VAST filter is given by

$$\mathbf{q}_{\text{VAST}}(V, \mu) = \sum_{v=1}^V \frac{\mathbf{b}_v^T \mathbf{r}_B}{\lambda_v + \mu} \mathbf{b}_v \quad (32)$$

which can also be used to show that  $J_B(\mathbf{q})$ ,  $J_D(\mathbf{q})$ , and  $J_B(\mathbf{q}) + \mu J_D(\mathbf{q})$  are non-increasing, non-decreasing, and non-increasing, respectively, for an increasing  $V$ . Moreover, we can also use the VAST framework to show that the acoustic contrast defined as

$$\gamma(\mathbf{q}) = \frac{|\mathcal{M}_D| \mathbf{q}^T \mathbf{R}_B \mathbf{q}}{|\mathcal{M}_B| \mathbf{q}^T \mathbf{R}_D \mathbf{q}}, \quad (33)$$

which is an important metric for sound zone control, is non-increasing for an increasing  $V$  when the VAST solution is used. This observation is extremely interesting since the ACC and PM solutions are two extreme special cases of the VAST solution for  $V = 1$  and for  $V = LJ$  and  $\mu = 1$ , respectively (see also Fig. 3 for some special cases of VAST). Thus, the VAST framework can be used to show theoretically that you not only trade-off signal distortion in the bright zone for the residual error in the dark zone, but also trade-off signal distortion in the bright zone for the acoustic contrast between the bright and dark zones.

#### 4. DISCUSSION

So far, we have focused on the similarities between the speech enhancement and sound zone control problem. However, there are also some very important differences which have consequences for VAST. First, the main difficulties in speech enhancement, we initially listed in Sec. 2, are actually not (or much smaller) problems for the sound zone control. That is, we have direct access to the individual signals and, therefore know the exact statistics defined in (27) and (28). Consequently, non-stationary signals do not lead to major problems since we do not have to estimate their statistics from short segments. Second, most (if not all) existing sound zone control methods assume Oracle knowledge of the room impulse responses (RIRs) from the loudspeakers to the control points. In practice, these are typically pre-measured in a separate calibration step so that the zones do not have to contain physical microphones during play-back. However, having Oracle knowledge of the RIRs is typically an optimistic assumption in practice since the RIRs cannot be

measured without errors [30] and are normally time-varying due to, e.g., temperature changes [31]. Third, the correlation vectors  $\mathbf{r}_s$  and  $\mathbf{r}_B$  have slightly different interpretations, unless the desired signal  $d_m(n)$  in the bright zone is defined to be the uncontrolled pressures  $y_m(n)$ . Usually, however, the desired signal is defined as the pressure originating from a virtual source, and this means a distortion of 0 cannot always be obtained, even for the MVDR solution which results in zero distortion for the speech enhancement problem. Fourth, whereas the clean speech covariance matrix is typically much more sparse than the noise covariance matrix in speech enhancement, a similar relationship between the spatial covariance matrices in the bright and dark zones seems not to exist. How this influences the choice of  $V$  is still an open question, but we anticipate that a larger  $V$  should generally be used for VAST than for VSLF. Fifth, the vector and matrix dimensions in the sound zone control problem are usually much larger than in the speech enhancement problem. Consequently, the computational cost of computing the control filters is much larger. The typical approach to circumvent this is to compute the control filters offline for generic signals such as white Gaussian noise sources or Dirac's delta functions. This means that all frequencies are given equal weight in the optimisation which is clearly suboptimal when highly coloured signals such as speech and audio are played back.

Despite that the VSLF framework was originally developed with the speech enhancement problem as its main application, the above discussion actually reveals that it is even better suited for the sound zone control problem, primarily because we there have direct access to the individual signals and, therefore, also the statistics. The main disadvantage of applying the VSLF framework to the sound zone control problem is the computational cost and memory requirements. However, the problem is highly structured, and we believe that this can be exploited to develop fast algorithms which either solve the problem exactly or approximately.

## 5. REFERENCES

- [1] J. Cheer, S. J. Elliott, and M. F. S. Gálvez, "Design and implementation of a car cabin personal audio system," *J. Audio Eng. Soc.*, vol. 61, no. 6, pp. 412–424, 2013.
- [2] X. Liao, J. Cheer, S. J. Elliott, and S. Zheng, "Design array of loudspeakers for personal audio system in a car cabin," in *Proc. Int. Congr. Sound Vib.*, 2016.
- [3] F. Heuchel, D. Caviedes Nozal, and F. T. Agerkvist, "Sound field control for reduction of noise from outdoor concerts," in *Conv. Audio Eng. Soc.* Audio Engineering Society, 2018.
- [4] J.-H. Chang, C.-H. Lee, J.-Y. Park, and Y.-H. Kim, "A realization of sound focused personal audio system using acoustic contrast control," *J. Acoust. Soc. Am.*, vol. 125, no. 4, pp. 2091–2097, 2009.
- [5] S. J. Elliott, J. Cheer, H. Murfet, and K. R. Holland, "Minimally radiating sources for personal audio," *J. Acoust. Soc. Am.*, vol. 128, no. 4, pp. 1721–1728, 2010.
- [6] J.-M. Lee, T. Lee, J.-Y. Park, and Y.-H. Kim, "Generation of a private listening zone; Acoustic parasol," in *Proc. Int. Congr. Acoust.*, 2010.
- [7] W. F. Druyvesteyn and J. Garas, "Personal sound," *J. Audio Eng. Soc.*, vol. 45, no. 9, pp. 685–701, 1997.
- [8] T. Betlehem, W. Zhang, M. A. Poletti, and T. D. Abhayapala, "Personal sound zones: Delivering interface-free audio to multiple listeners," *IEEE Signal Process. Mag.*, vol. 32, no. 2, pp. 81–91, 2015.
- [9] J. Choi and Y. Kim, "Generation of an acoustically bright zone with an illuminated region using multiple sources," *J. Acoust. Soc. Am.*, vol. 111, no. 4, pp. 1695–1700, Apr. 2002.
- [10] M. Poletti, "An investigation of 2-D multizone surround sound systems," in *Conv. Audio Eng. Soc.*, 2008, number 125.
- [11] J.-H. Chang and F. Jacobsen, "Sound field control with a circular double-layer array of loudspeakers," *J. Acoust. Soc. Am.*, vol. 131, no. 6, pp. 4518–4525, 2012.
- [12] M. B. Møller, M. Olsen, and F. Jacobsen, "A hybrid method combining synthesis of a sound field and control of acoustic contrast," in *Conv. Audio Eng. Soc.*, 2012.
- [13] M. F. Simón-Gálvez, S. J. Elliott, and J. Cheer, "Time domain optimization of filters used in a loudspeaker array for personal audio," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 11, pp. 1869–1878, 2015.
- [14] Y. Cai, M. Wu, L. Liu, and J. Yang, "Time-domain acoustic contrast control design with response differential constraint in personal audio systems," *J. Acoust. Soc. Am.*, vol. 135, no. 6, pp. 252–257, 2014.
- [15] M. B. Møller and M. Olsen, "Sound zones: On performance prediction of contrast control methods," in *AES Int. Conf. Sound Field Control*, 2016.
- [16] D. H. M. Schellekens, M. B. Møller, and M. Olsen, "Time domain acoustic contrast control implementation of sound zones for low-frequency input signals," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2016, pp. 365–369.
- [17] J. Benesty, M. G. Christensen, and J. R. Jensen, *Signal enhancement with variable span linear filters*, vol. 7, Singapore, Singapore: Springer Science+Business Media Singapore Pte Ltd, 2016.
- [18] J. Benesty, M. G. Christensen, J. R. Jensen, and J. Chen, "A brief overview of speech enhancement with linear filtering," *EURASIP J. on Advances in Signal Process.*, vol. 2014, no. 1, pp. 162, 2014.
- [19] J. R. Jensen, J. Benesty, and M. G. Christensen, "Noise reduction with optimal variable span linear filters," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 4, pp. 631–644, 2016.
- [20] T. Lee, J. K. Nielsen, J. R. Jensen, and M. G. Christensen, "A unified approach to generating sound zones using variable span linear filters," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2018.
- [21] P. C. Loizou, *Speech Enhancement: Theory and Practice*, Boca Raton, FL, USA: CRC Press, 2. edition, 2013.
- [22] T. Gerkmann, M. Krawczyk-Becker, and J. Le Roux, "Phase processing for single-channel speech enhancement: History and recent advances," *IEEE Signal Process. Mag.*, vol. 32, no. 2, pp. 55–66, 2015.
- [23] S. Boll, "Suppression of acoustic noise in speech using spectral subtraction," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 27, no. 2, pp. 113–120, 1979.
- [24] Y. Ephraim and D. Malah, "Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 32, no. 6, pp. 1109–1121, 1984.
- [25] G. Kim, Y. Lu, Y. Hu, and P. C. Loizou, "An algorithm that improves speech intelligibility in noise for normal-hearing listeners," *J. Acoust. Soc. Am.*, vol. 126, no. 3, pp. 1486–1494, 2009.
- [26] J. Benesty, J. R. Jensen, M. G. Christensen, and J. Chen, *Speech enhancement: A signal subspace perspective*, Elsevier, 2014.
- [27] N. Mohammadiha, P. Smaragdis, and A. Leijon, "Supervised and unsupervised speech enhancement using nonnegative matrix factorization," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 21, no. 10, pp. 2140–2151, 2013.
- [28] P. C. Loizou and G. Kim, "Reasons why current speech-enhancement algorithms do not improve speech intelligibility and suggested solutions," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 19, no. 1, pp. 47–56, 2011.
- [29] J. N. Franklin, *Matrix theory*, Mineola, New York, USA: Dover Publications Inc., 2000.
- [30] M. B. Møller, J. K. Nielsen, E. Fernandez-Grande, and S. K. Olesen, "On the influence of transfer function noise on low frequency pressure matching for sound zones," in *Proc. IEEE Sensor Array and Multich. Signal Process. Workshop*, 2018, pp. 331–335.
- [31] M. Olsen and M. B. Møller, "Sound zones: on the effect of ambient temperature variations in feed-forward systems," in *Conv. Audio Eng. Soc.*, 2017.