

A Multiple Listener Crosstalk Cancellation System Using Loudspeaker Dependent Regularization

Jacob Hollebon, Filippo Maria Fazi and Marcos F. Simón Gálvez

Institute of Sound and Vibration Research (ISVR)
University of Southampton, Southampton, SO17 1BJ, UK

Correspondence to J.Hollebon@soton.ac.uk

December 11, 2020

1 ABSTRACT

Binaural audio requires the use of crosstalk cancellation if reproduced using loudspeakers. A three-listener crosstalk cancellation system has been designed and built as part of this work. Simulations for different loudspeaker distributions are presented and a large system span is shown to increase low-frequency crosstalk cancellation performance whilst a denser loudspeaker distribution in front of a given listener increases mid to high-frequency crosstalk cancellation. The system's performance under perturbation of the listeners and loudspeakers is investigated and at high frequencies loudspeakers further away from any given listener are shown to be affected more by perturbations than those nearer the listener. To this issue, a novel implementation of weighting loudspeaker source strengths using loudspeaker dependent regularization is developed and optimised for use in this system. Hence, at high frequencies just loudspeakers close to the listener are used. This is shown to create a more robust solution than traditional crosstalk cancellation filter design when the system has undergone perturbations.

2 INTRODUCTION

The ability to reproduce a convincing 3D, or spatial, audio scene for a listener has received much attention following recent developments in the entertainment industry, such as with virtual reality or in cinema. One key technology for both recording and reproducing 3D audio is binaural audio. Binaural audio encompasses localization cues, such as interaural level differences and interaural time differences, that the brain uses to localize a sound source [1]. These localization cues can be modelled by recording using a binaural mannequin microphone or by synthesis using a set of Head Related Transfer Functions (HRTFs) [2]. The binaural stream is two channels, for the left and right ear of the listener respectively, which when reproduced correctly can convincingly reproduce any 3D audio scene.

A key requirement of binaural audio is that the signals must be reproduced at the listener's ears exactly to maintain the 3D effect, as any mixing of the signals would result in a loss of the interaural localization cues. This can be done simply using headphones. However, in many environments it may not be desirable to wear headphones. Crosstalk cancellation (CTC) systems have often been employed to reproduce binaural audio over loudspeakers. To ensure the binaural signals are correctly reproduced at the listener's ears inverse filtering is used to account for the undesired acoustical crosstalk. Ill-conditioning of the matrix undergoing inversion can cause solutions to be sensitive to small errors, as well as large loudspeaker gains. Hence, Tikhonov regularization is often used in the inversion process, which helps increase robustness of the system to errors, for example misalignment of the reproduction loudspeakers [3].

Early research focused on stereo CTC systems utilising a pair of loudspeakers with a 60 degree span [4, 5]. The mathematical theory to expand the problem to any number of loudspeakers or listeners was presented by Bauck and Cooper in [6, 7]. Further progress in CTC has considered more sophisticated loudspeaker geometries that improve the conditioning of the problem and allow for greater CTC, such as the stereo dipole [8, 9], optimal source distribution [10], and more recently uniform linear loudspeaker arrays (ULAs) which have been shown to be more robust to room reflections [11]. Traditionally CTC systems require the listener to remain static in a specific sweet-spot position, where even small misalignments of the listener or loudspeakers of just a few centimeters can drastically reduce the performance of the CTC

system [12, 13]. To allow listener movements, recent formulations have incorporated listener tracking to update the CTC filters in real-time adapting the sweet-spot position accordingly [14, 15, 16, 17, 18]. Similar problems occur when the listener is free to rotate their head, but this scenario may also be dynamically adjusted for [19]. Other issues such as mismatches between the HRTF of the listener and the HRTF used for the CTC filter design can introduce further perturbations, again degrading the ability of a system to provide CTC [20, 21].

Whilst much work has been done in CTC for a single listener, little consideration has been paid to the multiple listener problem. This is an important limiting factor in the field as a large range of practical applications for CTC, for example the entertainment industry, require multiple listener systems.

Initially, four-loudspeaker two-listener were trialed [22, 23] with particular focus on the positioning of the loudspeakers to minimize the ill-conditioning of the plant matrix [24, 25, 26]. However, more recent work has considered using a shared ULA between two listeners [27, 28] including systems adaptive to listener position [29, 30]. Notably, a static two-listener per loudspeaker array system has been built and installed in a commercial theme park ride [31]. This is one large install of many two listener CTC systems, hence each loudspeaker array attempts to control just the two listeners in front of the array, but does not control or steer nulls to any neighbouring listeners. Therefore leakage from the neighbouring CTC systems will introduce errors at any given listening position. This point has been expanded on in [32] where a control point configuration has been presented that allows multiple CTC systems to be placed side by side. This was achieved by focusing the output of each CTC system away from the neighbouring listener positions and towards just the listeners in front of each system. The setup considered was where each individual CTC system consisted of three ULAs for three listeners. Despite this, except for the demonstration of the mathematics and the simulation study in [32] no working CTC system for more than two listeners has yet been exhibited.

This paper is a continuation of the work described in [32] and presents the design and development of a static three-listener crosstalk cancellation system utilising personal uniform linear loudspeaker arrays for each listener, as well as the implementation of a novel loudspeaker dependent regularization technique in the filter design process. First, the theory of multiple listener multiple loudspeaker crosstalk cancellation is presented in Section 3, as well as the concept of loudspeaker dependent regularization. In Section 4 a range of different loudspeaker distributions for the system are numerically simulated using an analytical model and the best fitting system chosen as the final design, motivated by current state-of-the-art CTC loudspeaker systems. Next, the design is built as a real three-listener system. Section 5 considers the system's performance under realistic listener and loudspeaker position perturbations following similar examples in the literature. Motivated by the effect of these perturbations, a novel implementation of loudspeaker dependent regularization is derived and applied to the system such the loudspeakers furthest from each listener are regularized heavily and contribute less to controlling that listener position. Hence the system uses listener dependent variable loudspeaker regularization. This implementation is then optimised for the system and shown to increase the robustness of the system to perturbations.

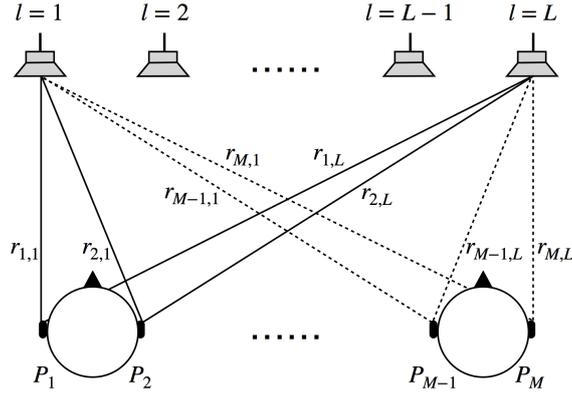


Figure 1: Geometry of the MIMO CTC problem, where L loudspeakers are used to control the pressure p_m at the M control points, which represent the M listeners' ears.

3 CROSSTALK CANCELLATION

3.1 MIMO Theory

The typical geometry of a Multiple-Input Multiple-Output (MIMO) CTC problem is shown in Fig. 1. L loudspeakers are used to control the pressure at M control points, which are placed at the listener's ears. Let \mathbf{p} , \mathbf{b} and \mathbf{v} be the vectors representing the reproduced pressure at the control points, the desired binaural signals to be replicated and the loudspeaker signals, respectively, such that

$$\begin{aligned}\mathbf{p} &= [p_1(\omega), p_2(\omega), \dots, p_M(\omega)]^T \\ \mathbf{b} &= [b_1(\omega), b_2(\omega), \dots, b_M(\omega)]^T \\ \mathbf{v} &= [v_1(\omega), v_2(\omega), \dots, v_L(\omega)]^T\end{aligned}\tag{1}$$

where $\omega = 2\pi f$ is the angular frequency of the reproduction signal. Then

$$\mathbf{p} = \mathbf{C}\mathbf{v}\tag{2}$$

where \mathbf{C} is the $M \times L$ plant matrix of acoustic transfer functions. The plant matrix can be measured for a specific loudspeaker and listener setup using a binaural mannequin microphone. However, it can also be modelled analytically under simple approximations. Assuming the loudspeakers radiate as ideal monopole sources in free field conditions, as well as using a shadowless head model, then

$$\mathbf{C} = \frac{1}{4\pi} \begin{bmatrix} \frac{e^{-jkr_{1,1}}}{r_{1,1}} & \dots & \frac{e^{-jkr_{1,L}}}{r_{1,L}} \\ \vdots & \ddots & \vdots \\ \frac{e^{-jkr_{M,1}}}{r_{M,1}} & \dots & \frac{e^{-jkr_{M,L}}}{r_{M,L}} \end{bmatrix}\tag{3}$$

where $k = \omega/c_0$ is the wavenumber, c_0 is the speed of sound and $r_{m,l}$ is the distance between loudspeaker l and control point m as shown in Fig. 1.

The loudspeaker signals are given by

$$\mathbf{v} = \mathbf{H}\mathbf{b} \quad (4)$$

where \mathbf{H} is the $L \times M$ matrix of CTC filters. The aim is to find the set of CTC filters that minimize the cost function

$$\begin{aligned} J &= \|\mathbf{p} - \mathbf{b}\|_2^2 \\ &= \|\mathbf{C}\mathbf{v} - \mathbf{b}\|_2^2 \end{aligned} \quad (5)$$

where the operator $\|\cdot\|_2^2$ represents the l_2 norm. This corresponds to minimising the sum of the squared errors between the signals reproduced at the listener's ears and the desired binaural signals. For the underdetermined case when $L > M$, which is often the situation in CTC, there exists infinite solutions that minimise this cost function. However, the solution $\mathbf{H} = \mathbf{C}^+$ where $(\cdot)^+$ represents the pseudoinverse gives the minimum norm solution that minimises both J and the l_2 norm [3]. A modelling delay is also required to ensure causality of the filters, however is omitted here for conciseness.

Ill-conditioning of the plant matrix undergoing inversion can cause solutions that are extremely sensitive to errors, as well as very large loudspeaker gains; a large condition number of the plant matrix, κ , indicates ill-conditioning [33]. Hence, Tikhonov regularization is often used to improve the conditioning of the plant matrix at frequencies at which the system is ill-conditioned. Whilst regularization improves the robustness of the solution and can be used to limit the loudspeaker gains, it comes at the cost of intentionally introducing errors into the solution [34]. The cost function to be minimized now contains both an error term and a loudspeaker gain penalty term, such that

$$J = \|\mathbf{C}\mathbf{v} - \mathbf{b}\|_2^2 + \beta\|\mathbf{v}\|_2^2 \quad (6)$$

where β is the regularization parameter, a positive constant that determines the weight of the penalty. The least squares solution that minimizes this cost function is [3]

$$\mathbf{H} = \mathbf{C}^H [\mathbf{C}\mathbf{C}^H + \beta\mathbf{I}_M]^{-1} \quad (7)$$

where $(\cdot)^H$ denotes the hermitian transpose and \mathbf{I}_M is the $M \times M$ identity matrix.

3.2 Loudspeaker Dependent Regularization

The penalty term in the cost function can be modified to include a weight, such that a larger penalty term is applied to selected loudspeakers. Hence the contribution from any chosen loudspeaker to controlling the pressure at any of the control points can be reduced or completely suppressed if desired. Let $\tilde{\mathbf{\Gamma}}$ be a $L \times L$ diagonal weighting matrix, such that

$$\tilde{\mathbf{\Gamma}} = \begin{bmatrix} \tilde{\gamma}_1(\omega) & & & & \\ & \tilde{\gamma}_2(\omega) & & & \\ & & \ddots & & \\ & & & \tilde{\gamma}_{L-1}(\omega) & \\ & & & & \tilde{\gamma}_L(\omega) \end{bmatrix} \quad (8)$$

where $\tilde{\gamma}_l(\omega)$ is a real-valued, positive scalar. The cost function now becomes

$$J = \|\mathbf{C}\mathbf{v} - \mathbf{b}\|_2^2 + \beta \|\tilde{\mathbf{\Gamma}}\mathbf{v}\|_2^2 \quad (9)$$

and the corresponding least squares solution is [35]

$$\mathbf{H} = [\mathbf{C}^H\mathbf{C} + \beta\tilde{\mathbf{\Gamma}}^H\tilde{\mathbf{\Gamma}}]^{-1}\mathbf{C}^H. \quad (10)$$

A derivation of the above formula is shown in the Appendix. To simplify the procedure, the regularization terms $\beta\tilde{\mathbf{\Gamma}}^H\tilde{\mathbf{\Gamma}}$ can be replaced by an equivalent $L \times L$ diagonal weighting matrix, $\mathbf{\Gamma}$, such that

$$\beta\tilde{\mathbf{\Gamma}}^H\tilde{\mathbf{\Gamma}} = \mathbf{\Gamma} = \begin{bmatrix} \gamma_1(\omega) & & & & \\ & \gamma_2(\omega) & & & \\ & & \ddots & & \\ & & & \gamma_{L-1}(\omega) & \\ & & & & \gamma_L(\omega) \end{bmatrix} \quad (11)$$

where $\gamma_l(\omega)$ is a positive, real-valued regularization weight for loudspeaker l . Through this formulation, traditional regularization to combat ill-conditioning can be combined with selective loudspeaker weighting simply, through the relation $\gamma_l(\omega) = \beta\tilde{\gamma}_l(\omega)\tilde{\gamma}_l(\omega)$. Hence the CTC filters are

$$\mathbf{H} = [\mathbf{C}^H\mathbf{C} + \mathbf{\Gamma}]^{-1}\mathbf{C}^H. \quad (12)$$

It may be desirable to apply a different loudspeaker weight to each control point, in which case let $\mathbf{\Gamma}_m$ be the weighting matrix to be used for control point m and \mathbf{H}_m be the resulting set of CTC filters found using Eq. [12] for that given $\mathbf{\Gamma}_m$. To find the final loudspeaker signals just the row that relates to control point m must be extracted from \mathbf{H}_m , as all other rows relate to control the pressure at the other control points using the wrong weighting matrix. Mathematically we may write this extraction as

$$\begin{aligned} \mathbf{h}_m &= [\mathbf{C}^H\mathbf{C} + \mathbf{\Gamma}_m]^{-1}\mathbf{C}^H\mathbf{p}_T \\ &= \mathbf{H}_m\mathbf{p}_T \end{aligned} \quad (13)$$

where \mathbf{h}_m is the column vector of length L CTC filters for control point m and \mathbf{p}_T is a column vector of length M target pressures. \mathbf{p}_T has all entries equal to zero, except the m -th entry which is equal to one. In this case it may be thought of as asking a target pressure of 1 at the m -th control point, and 0 at all other control points. Thus \mathbf{p}_T selects the relevant column

from the matrix \mathbf{H}_m . Eq. 13 must be evaluated M times for each of the M different weighting matrices. However, it may be necessary to use the same weighting matrix for multiple control points - for example two control points represent the two ears of a single listener, so the same weighting might be applied to these two control points. Here, Eq. 13 must be performed only $M/2$ times and \mathbf{p}_T has entries of 1 at the indexes that correspond to the two control points in question and 0's elsewhere. In this case \mathbf{p}_T selects the linear combination of the two relevant columns of \mathbf{H}_m . Hence, the technique allows for the flexible application of regularization both per loudspeaker and per frequency, by simply changing one parameter per loudspeaker. The definition of the loudspeaker regularization terms can be chosen to suit any relevant design motivation. Furthermore, in different frequency regions, different loudspeaker regularization motivations can be applied and transitioned between easily.

Finally, the full matrix of CTC filters may be formed by stacking each of the \mathbf{h}_m columns in order such that

$$\mathbf{H} = [\mathbf{h}_1 \quad \mathbf{h}_2 \quad \dots \quad \mathbf{h}_M] \quad (14)$$

and the loudspeaker signals are given by Eq. 4 as previously.

An alternative solution has been proposed in 36 where the CTC filters are given by

$$\mathbf{H} = \mathbf{W}^{-1} \mathbf{C}^H [\mathbf{C} \mathbf{W}^{-1} \mathbf{C}^H + \beta \mathbf{I}_M]^{-1} \quad (15)$$

where here the weighting matrix, \mathbf{W} , is a diagonal matrix which only contains the loudspeaker weighting terms. Tikhonov regularization is applied in the traditional manner separately from the loudspeaker weighting matrix.

Whilst Eq. 12 is the least squares solution to the cost function in Eq. 9 the solution in Eq. 15 is found by minimising the weighted norm of the the loudspeaker signals, such that the error between the reproduced and target binaural signals is zero, i.e.,

$$\min \|\mathbf{v}\|_W^2 \quad \text{s.t.} \quad \|\mathbf{C}\mathbf{v} - \mathbf{b}\|_2^2 = 0. \quad (16)$$

Practically, the two equations result in similar solutions when Tikhonov regularization is used, as in both the weighting equates to balancing the ratio of entries in the plant matrix to a regularization term, whether that be γ_l or β . In Eq. 12 this is achieved by diagonal loading, whilst in Eq. 15 this is done by minimising the diagonal terms. If a different weighting matrix is required for each control point, both solutions require the inversion to be calculated M times. The new proposed method allows for a simpler and more intuitive formulation of the regularization terms, such that the exact regularization parameter for each loudspeaker can be set, compared to Eq. 15. This makes it easier to directly control both the gain of each loudspeaker and the amount of regularization used when designing the CTC filters, through changing just one variable per loudspeaker. However, Eq. 15 may be better suited to adaptive CTC systems, as it requires inverting an $M \times M$ matrix whilst Eq. 12 requires inverting an $L \times L$ matrix which is more computationally intensive assuming, as with most CTC problems, $M < L$.

3.3 Performance Metrics

Two primary metrics are used to analyse the performance of the CTC system. The first is the normalised array effort, AE . The array effort is proportional to the energy required by

Distribution	Multiple or Single Arrays	Loudspeaker Separation, d	Total Length (m)	Regularization Parameter, β
ULA	Multiple	d	2.74	0.03350
Nested	Multiple	Increase by $d/2$ each iteration	3.10	0.03438
Extended Linear	Single	d , outer two pairs 0.25 m and 0.5 m respectively	3.42	0.03440
Extended Nested	Single	Increase by $d/2$ each iteration	7.80	0.03610

Table 1: Loudspeaker Distribution Parameters.

the system to perform crosstalk cancellation at a given frequency. Whilst controlling the array effort is not the same as controlling acoustic power, for a practical CTC system it is generally observed that a large array effort corresponds to a large acoustic power. Monitoring and limiting the energy required by the system is important to ensure the loudspeaker gains are not too large for a practical CTC system. In this sense, limiting the energy output of the CTC system can aid in robustness. Often where large energies are required in CTC this is for wave cancellation, which translates to comparatively little energy at the control points, thus this makes the CTC system extremely sensitive to perturbations. Here, in a practical scenario it is often best to sacrifice the system’s ability to provide CTC at these frequencies through using regularisation, which correspondingly increases the system’s robustness. It is also often desirable to minimise the energy output into the acoustical environment the CTC system is in, to reduce the effect room reflections and reverberation will have interfering with the reproduced signals at the listener’s ears which generally results in a loss of CTC performance [11]. If a maximum limit is set to the array effort, the regularization parameter, β , can be tuned such that this limit is not exceeded. This is important to ensure fair comparison between different CTC systems where different values of β may be required to match the same array effort limit. Consider $\mathbf{v}_{1,0}$, the loudspeaker signals to produce a binaural signal, $\mathbf{b}_{1,0}$, where

$$b_i(\omega) = \begin{cases} 1, & \text{if } i \text{ is odd} \\ 0, & \text{if } i \text{ is even} \end{cases} \quad (17)$$

corresponding to an impulse in the left ear and null in the right ear for all listeners. The array effort is defined as the norm of these loudspeaker signals divided by the norm of v_s ,

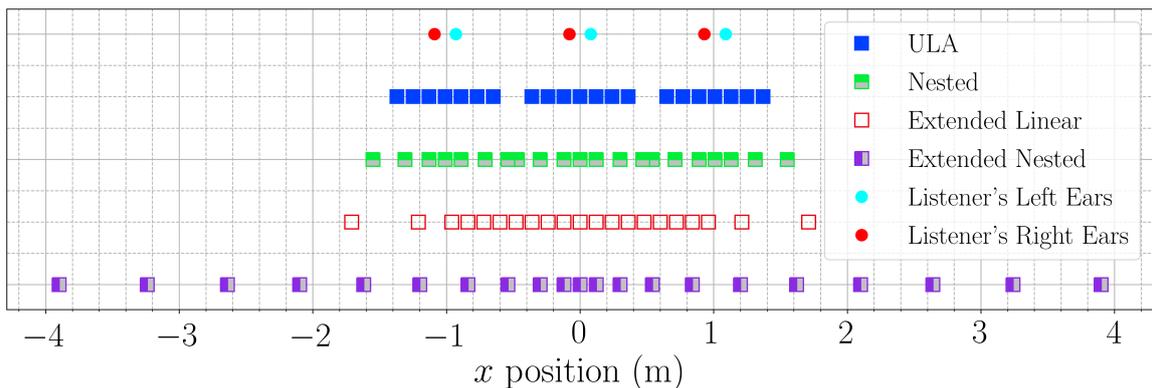


Figure 2: The four proposed loudspeaker distributions. Note the y axis does not denote y position.

the loudspeaker signal required to reproduce an impulse in a given ear by the most central loudspeaker. Hence [37]

$$AE = \frac{\mathbf{v}_{1,0}^H \mathbf{v}_{1,0}}{|v_s|^2}. \quad (18)$$

The array effort is a dimensionless quantity and typically presented on a decibel scale.

The second metric is the crosstalk cancellation spectrum. Whilst this has previously been used in [28], here an alternative form to account for the addition of multiple listeners is presented. Let $\mathbf{R} = \mathbf{C}\mathbf{H}$ be the $M \times M$ crosstalk matrix. The CTC spectrum, CTC , is then

$$CTC_i = \frac{|R_{i,i}|^2(M-1)}{\sum_{j=1, j \neq i}^M |R_{i,j}|^2} \quad (19)$$

such that for control point i the numerator is the reproduced signal and the denominator is the average of the squared crosstalk terms for the other control points. To simplify the metric for multiple listeners and hence many control points, the average CTC spectrum for a CTC system is

$$CTC_{avg} = \frac{\sum_{i=1}^M CTC_i}{M}. \quad (20)$$

The CTC spectrum is typically plotted in decibels.

4 SIMULATIONS FOR SYSTEM DESIGN

Numerical simulations were used to test the performance of four different loudspeaker distributions for the three-listener CTC system. For the simulations, an analytical plant matrix was built using Eq. [3] and the standard CTC filter approach in Eq. [7] was used. Hence all loudspeakers were set to control all the control points with equal weighting. A frequency independent regularization parameter, β , was set such that the array effort did not exceed 10 dB at any frequency - this meant the four loudspeaker distributions had the same relative amount of regularization applied during the simulations, even though slightly different values of β were required.

The geometry of the simulations represent a row of seated listeners. Three listeners, such that $M = 6$, are placed 1.1 m away from the CTC system. There is a spacing of 0.85 m between listeners and a head diameter of 0.16 m is used.

Due to the number and size of the specific loudspeakers available for a physical build of the CTC system, the simulations used a total number of $L = 21$ loudspeakers and the smallest spacing between loudspeakers, d , was 0.12 m. The four system designs are shown in Fig. [2] and detailed in Table [1]. Two different types of distributions, uniform linear arrays (ULAs) and nested arrays were tested. The ULA design consists of loudspeakers with a constant spacing, d . The nested design starts with an initial loudspeaker spacing, d , moving outwards from the central loudspeaker, and then increases the spacing by an additional $d/2$ for each consecutive loudspeaker iteration. ULAs were chosen as previous work in CTC has shown ULAs to be advantageous over traditional two-loudspeaker CTC systems [11]. However, the nested design

was also considered as it was theorised that by increasing the total system length this would improve the conditioning of the plant matrix at low frequencies, whilst using the same number of loudspeakers as a ULA. The ULA and nested designs were tested such that there were three smaller personal seven-loudspeaker arrays, each centred on one of the three listeners.

The next two designs were variations of ULA and nested arrays. The extended nested design consists of one large format twenty-one-loudspeaker array which was centered on the middle listener. Finally, a large format extended linear design was also simulated. This distribution consisted of a twenty-one-loudspeaker ULA, however, the spacing of the outer two pairs of loudspeakers was increased by an extra 0.25 m and 0.50 m respectively, for the same low frequency motivation as with the nested designs.

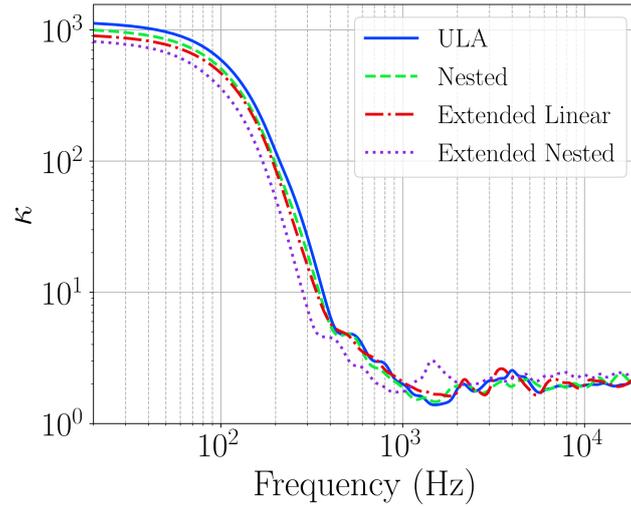
The condition number, κ , of the plant matrices corresponding to the four designs is shown in Fig. 3a. A large condition number at any given frequency indicates ill-conditioning and where regularization will have the strongest effect. This is clearly seen in the array effort of the CTC filters shown in Fig. 3b, where there is a low-frequency limit below which regularization dominates and also below which the condition number continues to increase. This corresponds to a lack in the system’s ability to provide CTC, as seen from Fig. 3c, which shows the CTC spectrum averaged across all three listener positions. Here, this trend is seen for all the designs with lack of CTC below 20 dB setting the low frequency limits between approximately 240–300 Hz the systems. However, for systems with a larger overall length (see Table 1) this low-frequency limit appears at a lower frequency, such that the extended nested designs performs the best at low-frequencies. This suggests that to achieve low-frequency CTC a system with a wider span is desirable and nested arrays could deliver this improvement whilst using the same number of loudspeakers as ULAs. However, this comes at the cost of a reduction in mid-frequency and high-frequency CTC, in which the ULAs are the strongest. The extended nested array also comes with an undesirably large increase in the length of the system, 5.06 m longer than the ULAs for only a 60 Hz drop in the low-frequency limit. The extended linear array appears a strong compromise where the larger system length lowers the low-frequency limit of the CTC whilst maintaining a strong performance in the mid-frequencies and high-frequencies, without making the system impractically large. Hence a large system span improves low-frequency CTC performance, whilst a more dense distribution of loudspeakers in front of each listener improves mid-frequency and high-frequency CTC control. However, when considering the design and size of the system span in light of low-frequency performance, it is also important to consider that at low-frequencies sound localisation can be poor [38] and an increase in CTC at these low-frequencies may not translate to a perceived increase in performance.

The final design chosen for the practical build was the personal ULAs. Whilst the simulations show there is a slight improvement in low-frequency CTC with the other designs, the ULAs perform strongly at all other frequencies. Furthermore, the ULA design consists of three identical arrays, hence the build of a practical system would be simpler and the size more suitable for practical use. Furthermore, the system is modular and for further work could be expanded for more listeners by adding further identical arrays to either end.

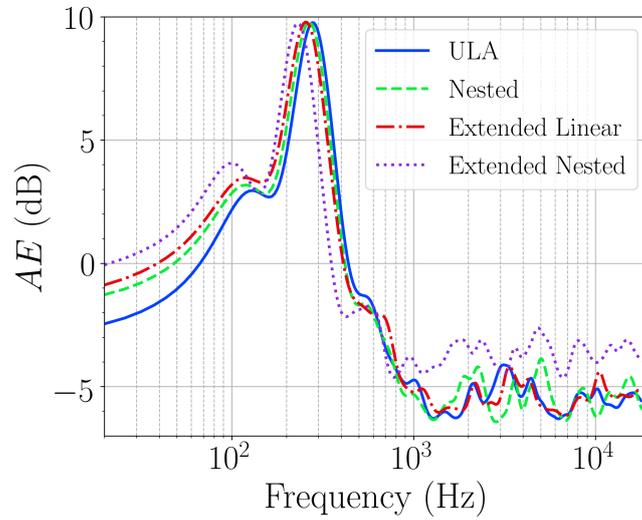
5 ROBUSTNESS ANALYSIS

5.1 Experimental Procedure

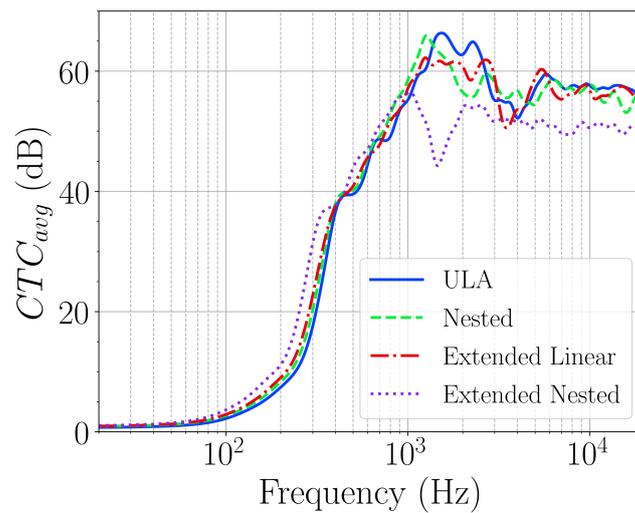
The three personal listener ULAs as described in Section 4 were built and implemented as a real three listener CTC system. As setting up the system assumes the loudspeaker arrays



(a) Condition number of the plant matrices.



(b) Array effort of the CTC filters.



(c) CTC spectrum averaged across all control points.

Figure 3: Simulated performance of the four CTC system designs.

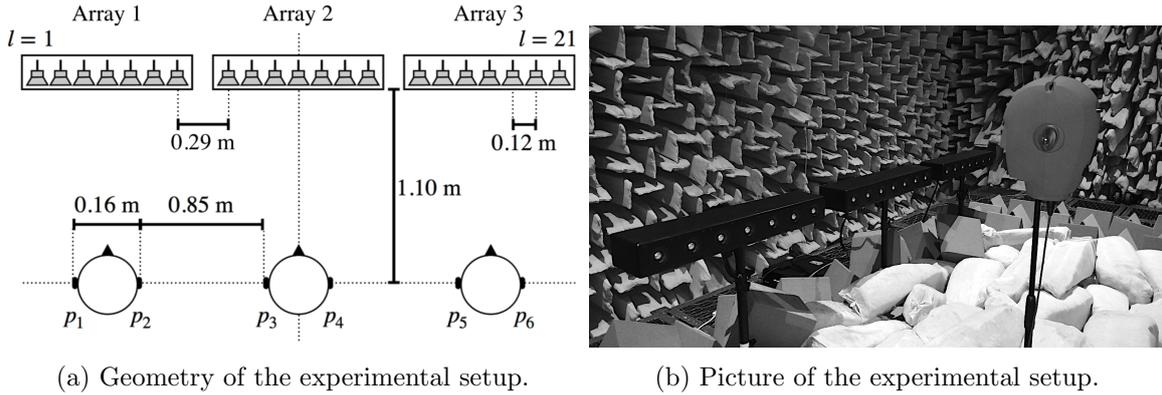


Figure 4: Experimental setup for the plant transfer function measurements.

and listeners are perfectly aligned, it is important to understand how well the system provides CTC under perturbations. Two main sources of perturbation are inexact positioning of the loudspeaker arrays and also the listener's head being positioned outside of the sweet-spot. Such errors are inevitable if the system is moved to a new location and re-setup. To experimentally verify the system's performance under such perturbations, the plant matrix of transfer functions for the real system was measured in an anechoic chamber and these measurements used to create a set of CTC filters. The system was then completely dismantled, and re-setup on a separate day to introduce some perturbations. The magnitude of the perturbations was approximately a distance of ± 0.05 m and an angular rotation of 5 degrees in each of the loudspeaker and listener positions. These values were chosen as it has been shown that the performance of CTC systems degrades substantially for perturbations larger than this magnitude [12]. Hence, to account for larger misalignments it is suggested that dynamic CTC utilising headtracking is required [15]. The transfer functions of the perturbed system were then measured. Offline simulations were used to find the CTC spectrum of the system such that

$$\mathbf{R} = \tilde{\mathbf{C}}\mathbf{H} \quad (21)$$

where $\tilde{\mathbf{C}}$ is the perturbed plant matrix and the CTC filters, \mathbf{H} , were created using the unperturbed plant matrix, \mathbf{C} . The CTC spectrum as detailed in Eq. 19 was then calculated.

The loudspeaker arrays and listener positions were set up in the same configuration as detailed in Section 4 for the numerical simulations, and are shown in Fig. 4a, whilst a picture of the experimental setup is shown in Fig. 4b. The transfer functions were measured using sine sweeps. The sine sweeps were played by a control computer through an RME MADIFace Pro audio interface, which was connected by MADI to an Innosonix MA32/LP amplifier, which drove each loudspeaker on a separate channel. A Neumann KU100 binaural microphone was used to record the sweep, which was sent to the MADIFace Pro and into MATLAB®. The anechoic chamber at the Institute of Sound and Vibration Research (ISVR), University of Southampton was used to ensure free field conditions.

5.2 Errors In The Plant Matrix

To quantify the effects of perturbation between the two sets of measured plant matrices, the error, e , is defined as

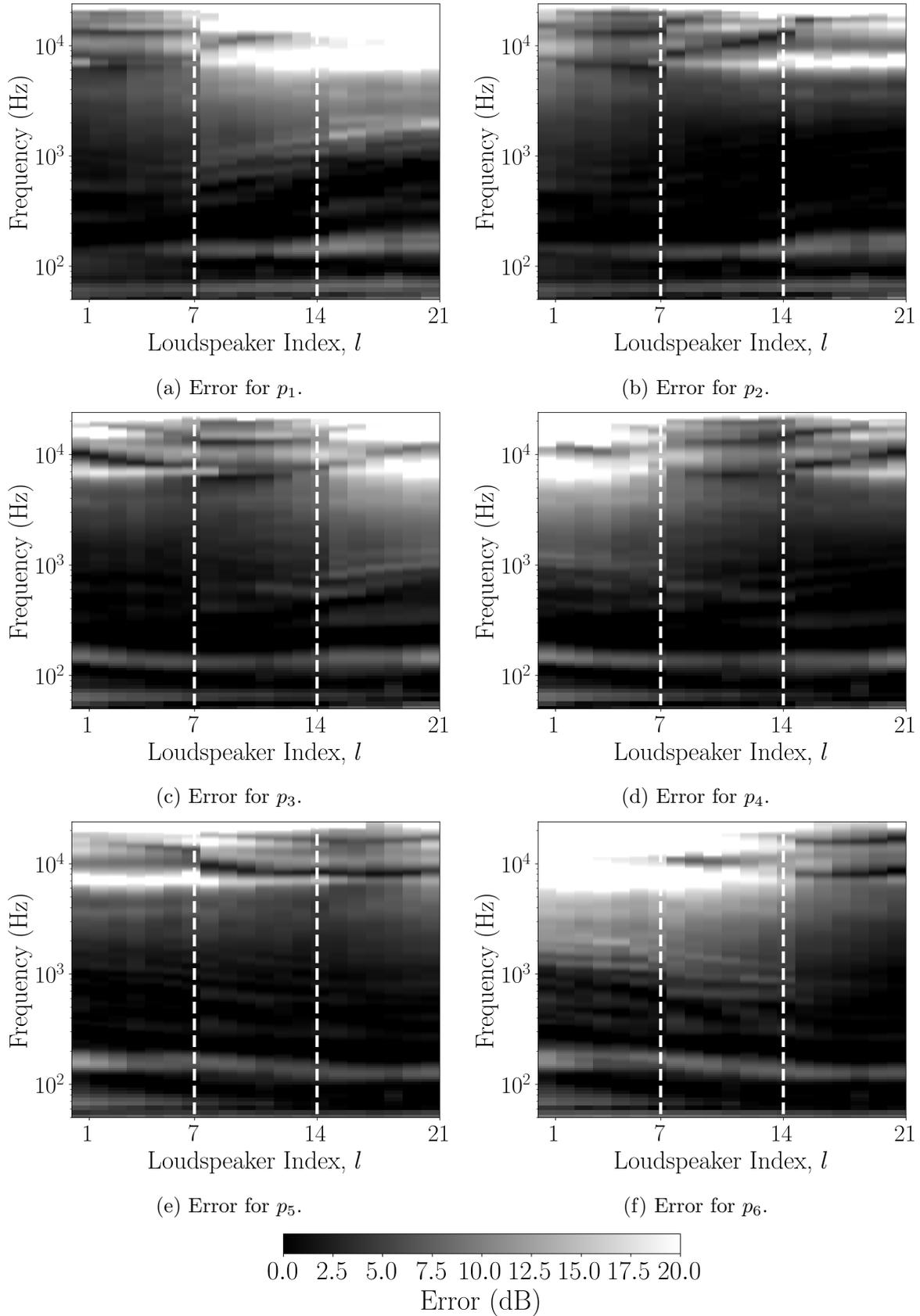


Figure 5: Errors in the plant acoustic transfer functions having undergone perturbations. The white dashed lines indicate the start and end of each of the loudspeaker arrays.

$$e_{m,l}(\omega) = \frac{|C_{m,l}(\omega) - \tilde{C}_{m,l}(\omega)|^2}{|C_{m,l}(\omega)|^2} \quad (22)$$

where $C_{m,l}$ denotes the transfer function between control point m and loudspeaker l . The error between the unperturbed and perturbed measured transfer functions is shown in Fig. 5 for all six control points. It is clear that the system is less robust at high frequencies where there are large errors between the entries of the two plant matrices. It is also clear the high frequency errors are smaller for loudspeakers closer to a given control point, and larger for loudspeakers further away. In particular, the errors are larger and across a wider frequency range when the loudspeakers are on the contralateral (opposite) side of the binaural mannequin head with regards to the control point in question. However at low frequencies, there is relatively little error regardless of the loudspeaker index. There is a notable error for all control points, at 170 Hz, which is at the resonance frequency of the loudspeakers used. As at the resonance frequency of the loudspeaker the phase changes quickly through frequency, here the error could be due to a change in temperature as the perturbed measurements were performed on a separate day, which causes a change in the speed of sound and hence a small variation in the measured resonance frequency of the loudspeaker. Overall, the perturbations introduce the largest error at high frequencies when the loudspeakers are further away from the control points. This suggests that the CTC system could be made more robust at high frequencies by using just the loudspeakers in front of each listener.

5.3 Loudspeaker Dependent Regularization

Motivated by the results of Section 5.2, the robustness of two different methods for calculating the CTC filters was tested using offline simulations. First, the standard inversion in Eq. 7 was used as a reference solution. When creating the filters using this method, a constant value of $\beta = 0.00399$ was used at all frequencies to ensure the array effort was never larger than 10 dB. Secondly, a solution utilising loudspeaker dependent regularization, as in Eq. 12, was implemented. The flexibility of Eq. 12 allows for system design choices to be implemented easily by changing the weighting matrix, $\mathbf{\Gamma}$.

For this specific system, it is not desirable to use any loudspeaker dependent weighting at low-frequencies as from Fig. 5 there is very little error in the transfer functions here. Furthermore, it has been shown in Section 4 that a large system length is required for strong low-frequency CTC, and applying a weighting would effectively reduce the system length. Therefore, the loudspeaker regularization should be set at low frequencies to use all the loudspeakers equally, whilst applying standard Tikhonov regularization to combat ill-conditioning. Hence at low frequencies below 900 Hz as dictated by Fig. 5, the diagonal elements were all set to the same value such that $\gamma_l = \beta = 0.00399 \forall l \in [1, L]$. This equates to using all loudspeakers with Tikhonov regularization as with the standard solution.

However, at high frequencies for loudspeakers further away from the control points, the perturbations caused a maximal error. Therefore, the weighting matrix was set such that at high frequencies only the loudspeakers closest to a control point are used. The flexibility of the loudspeaker dependent regularization technique allows any definition for the weightings to be applied, dependant on any desired design choice. Thus to realise this idea, the weights were applied from 1100 Hz as a linear function of the path length between any given control point and loudspeaker. The weighting does not have to be applied in a linear manner, but this weighting definition was found to give satisfactory results. Hence, the entries for $\mathbf{\Gamma}_m$ are given

by

$$\gamma_l = \alpha r_{l,m} \quad (23)$$

where α is a constant that varies the magnitude of the weight and $r_{l,m}$ is the distance between control point m and loudspeaker l . A larger value of γ_l corresponds to a stronger weighting of the effort penalty to the loudspeaker, hence loudspeakers further away from the control point are penalised more.

A transition region was defined between 900 Hz and 1100 Hz, such that here the values for γ_l linearly transitioned from β to those set out in Eq. 23. This was implemented to minimise the effects of artefacts from transitioning between the two different frequency regions with different regularization motivations. Again for a fair comparison, it was enforced that the array effort did not exceed 10 dB at any frequency when deciding on the values for $\mathbf{\Gamma}$. Hence, a value of $\alpha = 0.00689$ was used. Thus the loudspeakers were weighted to contribute more or less heavily to specific control points in set frequency bands.

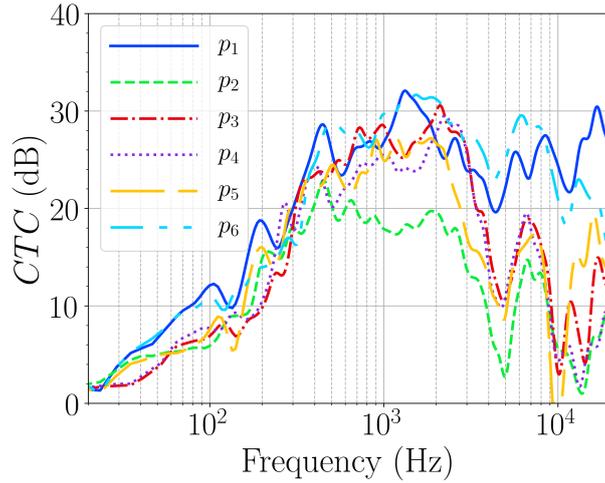
This particular implementation of weighting requires a different $\mathbf{\Gamma}_m$ for each listener, such that the correct loudspeakers are penalized for each control point. Eq. 12 must be computed for every instance of $\mathbf{\Gamma}_m$. However, the increase in computational load from both the multiple inversions and inverting an $L \times L$ matrix (as opposed to an $M \times M$ matrix in Eq. 7) is not an issue for this CTC system, as it is static hence the CTC filters can be calculated both offline and just once.

5.4 Results

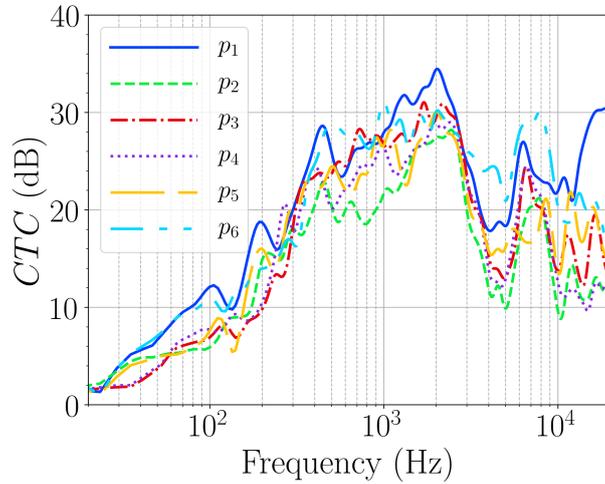
The CTC spectrum under perturbation for each of the six control points, using the CTC filters with and without loudspeaker dependent regularization, is shown in Fig. 6. A lower level of CTC is achieved than in the ideal simulations in Section 4. This is primarily due to the introduction of error through the inclusion of perturbations to the problem. From Fig. 6a for the unweighted CTC filters, it is clear that introducing these perturbations drastically reduces the amount of CTC and above 3000 Hz often below 20 dB of CTC is achievable. Furthermore, the CTC is extremely variable between each of the six control points. Notably, the two outer control points, $m = 1$ and 6, retain the greatest amount of CTC, whereas all other control points achieve approximately 15 dB less CTC in comparison. Hence, the system performs worse when perturbations are introduced, in a manner variable across all the different control points.

When using the loudspeaker dependent regularization, in Fig. 6b, there is a general increase in the CTC achieved above 1100 Hz. Furthermore, all positions have a much more similar response indicating more uniformity and robustness in the solution. Notably, $m = 1$ and 6 have a similar performance to all the other control points unlike with the standard CTC filter design. Hence, by weighting the loudspeakers to prioritise control points in front of them, the system has been made more robust to both loudspeaker and listener perturbations.

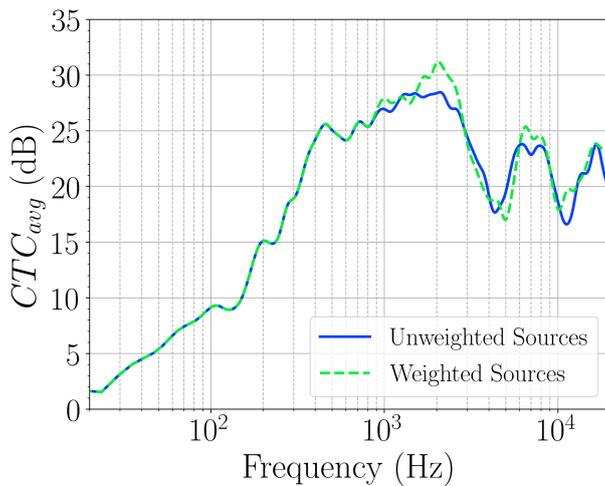
Fig. 6c shows the CTC averaged over all listener positions for the CTC filters with and without loudspeaker dependent regularization. It is clear that whilst the loudspeaker weighting method drastically limits the loudspeakers available to control each listener's ear, a similar level of CTC can still be achieved under perturbations to when using all the loudspeakers equally. Furthermore, using the loudspeaker dependent regularization the average CTC is often equal or slightly improved, compared to the standard solution. Hence, by applying loudspeaker



(a) CTC when using unweighted loudspeaker source strengths.



(b) CTC when using weighted loudspeaker source strengths.



(c) Average CTC across all control points for both weighted and unweighted loudspeaker source strengths.

Figure 6: CTC spectrum when the system is under perturbation, for the standard CTC filter design and that using loudspeaker dependent regularisation.

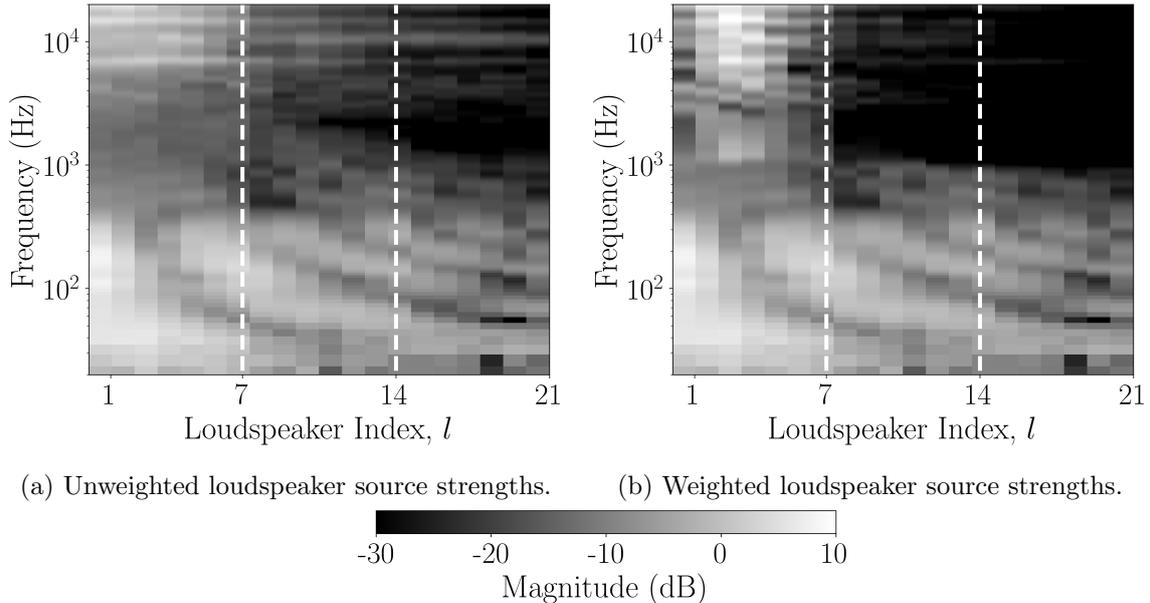


Figure 7: Loudspeaker source strengths when recreating a binaural signal of $\mathbf{b} = [1, 0, 0, 0, 0, 0]^T$. The white dashed lines indicate the start and end of each of the loudspeaker arrays.

dependent regularization the solution has been made more robust, whilst achieving a similar level of CTC on average even though less loudspeakers are required for each control point.

Finally, to demonstrate the loudspeaker weighting and therefore reduction in the magnitude of the loudspeaker signals, a binaural input of $\mathbf{b} = [1, 0, 0, 0, 0, 0]^T$ was used in Eq. 4. This binaural signal corresponds to an impulse in the left ear control point of the first listener who is situated in front of loudspeakers 3 and 4 (as in Fig. 4a), and nulls at all other control points. Using the standard CTC filters, it is expected that loudspeakers across all three arrays will be active. However, with the loudspeaker dependent regularization at high-frequencies the loudspeakers closest to the first listener should be the loudest.

Fig. 7 shows the magnitude of the loudspeaker signals for the unweighted and weighted solutions respectively. At low frequencies both methods have the same loudspeaker signals, as expected. When there is no weighting applied in the inversion, the solution still favours the loudspeakers closest to the control point. This is likely due to some natural weighting from the measured transfer functions, taking into account the longer path lengths between loudspeakers further away from the control point and also head shadowing from the binaural mannequin microphone. Despite this, there is a strong peak at 10 kHz that utilises all the loudspeakers. Applying the weighting forces the loudspeakers 1 – 7, which are in array 1 directly in front of the listener, to be used. Whilst this does increase the magnitude of the signal for the front loudspeakers, this does not exceed the magnitude that is required at low frequencies. Furthermore, the linear aspect of the weighting can be seen as the signal for loudspeakers in the middle of array 1 to have a larger magnitude than those at the edge. Loudspeakers 8 – 14 and 15 – 21, in array 2 and 3 respectively, are largely unused, with signals 30 – 40 dB smaller than those in array one at all frequencies above 1100 Hz.

6 CONCLUSION

A three-listener crosstalk cancellation system using uniform linear loudspeaker arrays has been designed, simulated, built, and experimentally verified. Furthermore, a novel loudspeaker dependent regularisation technique has been presented and optimised for use in the crosstalk cancellation system.

Four designs for the loudspeaker arrays were simulated using free-field numerical simulations. A larger span of the loudspeaker control sources is shown to increase the low-frequency crosstalk cancellation performance of the system. A more dense distribution of loudspeakers in front of each listener increases the mid-frequency and high-frequency crosstalk cancellation performance. Uniform linear loudspeaker arrays are demonstrated to be a good compromise between these two requirements for crosstalk cancellation systems whilst remaining a practical design for a build.

The crosstalk cancellation system was then built following the simulation results. Transfer functions between the three loudspeaker arrays to the three listener positions were measured twice, where the second time perturbations were introduced to the system. Perturbations such as loudspeaker and listener misalignments have been shown to drastically reduce the crosstalk cancellation performance of the system. In particular, it is shown that the error due to perturbations is larger for both loudspeakers further away from a given control point and as the frequency increases.

Motivated by these results, a novel loudspeaker dependent regularization technique was developed. This technique allows for individual choices of the regularization parameter for each loudspeaker, relevant to any given control point. Therefore flexible system design choices can be implemented by the weighting of each loudspeaker. This technique was implemented and optimised for this crosstalk cancellation system; the weighting is set such that at high frequencies only the loudspeakers closest to each control point are used for that given control point. However at low frequencies all loudspeakers are used, because following from the numerical simulation results, a large system span improves low-frequency CTC. The loudspeaker dependent regularization method is shown to increase the robustness of the solution under perturbation as well as achieving on average a slightly larger amount of crosstalk cancellation, even though less loudspeakers are used for each control point. Therefore, for multiple listener crosstalk cancellation it is advantageous at high-frequencies to only use loudspeakers close to a given listener.

For future work, the loudspeaker dependent regularization technique is to be validated further through the use of subjective listening tests. Furthermore, the performance of the system under greater perturbations than considered in this work may be investigated, as well as adding adaptive crosstalk cancellation through listener headtracking to aid in the performance under such conditions.

References

- [1] J. Blauert, *The Technology Of Binaural Listening*, pp. 1–32 (Springer, Berlin), 1st ed. (2012), DOI: <https://doi.org/10.1007/978-3-642-37762-4>.
- [2] H. Møller, “Fundamentals of binaural technology,” *Applied Acoustics*, vol. 36, no. 3, pp. 171 – 218 (1992), DOI: [https://doi.org/10.1016/0003-682X\(92\)90046-U](https://doi.org/10.1016/0003-682X(92)90046-U).
- [3] O. Kirkeby, P. A. Nelson, H. Hamada, F. Orduna-Bustamante, “Fast deconvolution of multichannel

- systems using regularization,” *IEEE Transactions on Speech and Audio Processing*, vol. 6, no. 2, pp. 189–194 (1998 March), DOI: <https://doi.org/10.1109/89.661479>.
- [4] M. R. Schroeder, B. S. Atal, “Computer Simulation of Sound Transmission in Rooms,” *Proceedings of the IEEE*, vol. 51, no. 3, pp. 536–537 (1963 March), DOI: <https://doi.org/10.1109/PROC.1963.2180>.
- [5] M. R. Schroeder, “Digital Simulation of Sound Transmission in Reverberant Spaces,” *The Journal of the Acoustical Society of America*, vol. 45, no. 1, pp. 303–303 (1969 Nov), DOI: <https://doi.org/10.1121/1.1971383>.
- [6] D. Cooper, J. Bauck, “Generalized Transaural Stereo and Applications,” *J. Audio Eng. Soc.*, vol. 44, no. 9, pp. 683–705 (1996 Sep).
- [7] D. Cooper, J. Bauck, “Prospects for Transaural Recording,” *J. Audio Eng. Soc.*, vol. 37, no. 1/2, pp. 3–19 (1989 Feb).
- [8] O. Kirkeby, P. A. Nelson, H. Hamada, “Local Sound Field Reproduction using Two Closely Spaced Loudspeakers,” *The Journal of the Acoustical Society of America*, vol. 104, no. 4, pp. 1973–1981 (1998), DOI: <https://doi.org/10.1121/1.423763>.
- [9] O. Kirkeby, P. A. Nelson, H. Hamada, “The ‘Stereo Dipole’: A Virtual Source Imaging System Using Two Closely Spaced Loudspeakers,” *J. Audio Eng. Soc.*, vol. 46, no. 5, pp. 387–395 (1998 Oct).
- [10] T. Takeuchi, P. A. Nelson, “Optimal Source Distribution for Binaural Synthesis over Loudspeakers,” *The Journal of the Acoustical Society of America*, vol. 112, no. 6, pp. 2786–2797 (2002 Dec), DOI: <https://doi.org/10.1121/1.1513363>.
- [11] M. F. Simón Gálvez, F. M. Fazi, “Loudspeaker Arrays For Transaural Reproduction,” presented at the *The 22nd International Congress of Sound and Vibration, Florence* (2015 July).
- [12] T. Takeuchi, P. A. Nelson, H. Hamada, “Robustness to head misalignment of virtual sound imaging systems,” *The Journal of the Acoustical Society of America*, vol. 109, no. 3, pp. 958–971 (2001 Mar).
- [13] Y. Parodi, P. Rubak, “Objective evaluation of the sweet spot size in spatial sound reproduction using elevated loudspeakers,” *The Journal of the Acoustical Society of America*, vol. 128, no. 3, pp. 1045–1055 (2010 Sep).
- [14] M. F. Simón Gálvez, T. Takeuchi, F. M. Fazi, “A Listener Adaptive Optimal Source Distribution System for Virtual Sound Imaging,” presented at the *Audio Engineering Society Convention 140* (2016 May).
- [15] M. F. Simón Gálvez, T. Takeuchi, F. M. Fazi, “Low-Complexity, Listener’s Position-Adaptive Binaural Reproduction Over a Loudspeaker Array,” *Acta Acustica united with Acustica*, vol. 103, no. 5, pp. 847–857 (2017 Sep), DOI: <https://doi.org/10.3813/AAA.919112>.
- [16] H. Kurabayashi, M. Otani, K. Itoh, M. Hashimoto, M. Kayama, “Development of dynamic transaural reproduction system using non-contact head tracking,” presented at the *2013 IEEE 2nd Global Conference on Consumer Electronics (GCCE)*, pp. 12–16 (2013 Oct), DOI: <https://doi.org/10.1109/GCCE.2013.6664768>.
- [17] X. Ma, C. Hohnerlein, J. Ahrens, “Concept and Perceptual Validation of Listener-Position Adaptive Superdirective Crosstalk Cancellation Using a Linear Loudspeaker Array,” *J. Audio Eng. Soc.*, vol. 67, no. 11, pp. 871–881 (2019 Nov).

- [18] M. F. Simón Gálvez, D. Menzies, F. M. Fazi, “Dynamic Audio Reproduction with Linear Loudspeaker Arrays,” *J. Audio Eng. Soc.*, vol. 67, no. 4, pp. 190–200 (2019 April).
- [19] M. F. Simón Gálvez, E. Hamdan, D. Menzies, F. M. Fazi, “A Study of the Effect of Head Rotation on Transaural Reproduction,” presented at the *Audio Engineering Society Convention 145* (2018 Oct).
- [20] M. A. Akeroyd, J. Chambers, D. Bullock, A. R. Palmer, A. Q. Summerfield, P. A. Nelson, S. Gatehouse, “The binaural performance of a cross-talk cancellation system with matched or mismatched setup and playback acoustics,” *The Journal of the Acoustical Society of America*, vol. 121, no. 2, pp. 1056–1069 (2007 Feb).
- [21] J. Hollebon, E. Hamdan, F. M. Fazi, “A Comparison Of The Performance Of HRTF Models In Inverse Filter Design For Crosstalk Cancellation,” presented at the *Proceedings of the Institute of Acoustics: Reproduced Sound*, vol. 41 (2019 Nov).
- [22] Y. Kahana, P. A. Nelson, O. Kirkeby, H. Hamada, “Objective and Subjective Assessment of Systems for the Production of Virtual Acoustic Images for Multiple Listeners,” presented at the *Audio Engineering Society Convention 103* (1997 Sep).
- [23] J.-S. Lim, C. Kyriakakis, “Virtual Loudspeaker Rendering for Multiple Listeners,” presented at the *Audio Engineering Society Convention 109* (2000 Sep).
- [24] Y. Kim, O. Deille, P. A. Nelson, “Crosstalk Cancellation in Virtual Acoustic Imaging Systems for Multiple Listeners,” *Journal of Sound and Vibration*, vol. 297, no. 1, pp. 251 – 266 (2006), DOI: <https://doi.org/10.1016/j.jsv.2006.03.042>
- [25] B. Masiero, X. Qiu, “Two Listeners Crosstalk Cancellation System Modelled by Four Point Sources and Two Rigid Spheres,” *Acta Acustica united with Acustica*, vol. 95, no. 2, pp. 379–385 (2009 Mar), DOI: <https://doi.org/10.3813/AAA.918160>.
- [26] B. Masiero, “Source Positioning in a Two Listener Crosstalk Cancellation System,” presented at the *NAG/DAGA International Conference on Acoustics, Rotterdam* (2009 March).
- [27] M. F. Simón Gálvez, F. M. Fazi, “A Loudspeaker Array For 2 People Transaural Reproduction,” presented at the *24th International Congress on Sound and Vibration* (2017 July).
- [28] C. House, et al., “Development of a Loudspeaker Array for Multi-Listener Transaural Reproduction in a Vehicle,” *Proceedings of the Institute of Acoustics*, vol. 39, no. 2 (2017).
- [29] M. F. Simón Gálvez, F. M. Fazi, “Listener Adaptive Filtering Strategies for Personal Audio Reproduction over Loudspeaker Arrays,” presented at the *Audio Engineering Society Conference: 2016 AES International Conference on Sound Field Control* (2016 July).
- [30] H. Kurabayashi, M. Otani, M. Hashimoto, M. Kayama, “Development of dynamic crosstalk cancellation system for multiple-listener binaural reproduction,” *Acoustical Science and Technology*, vol. 36, no. 6, pp. 537–539 (2015 Nov), DOI: <https://doi.org/10.1250/ast.36.537>.
- [31] P. Otto, E. Hamdan, “Bridging Near and Far Acoustical Fields: a Hybrid Systems Approach to Improved Dimensionality in Multi-Listener Spaces,” presented at the *Audio Engineering Society Conference: 2016 AES International Conference on Sound Field Control* (2016 July).
- [32] J. Hollebon, M. F. Simón Gálvez, F. M. Fazi, “Multiple Listener Crosstalk Cancellation Using Linear Loudspeaker Arrays For Binaural Cinematic Audio,” presented at the *Audio Engineering Society Conference: 2019 AES International Conference on Immersive and Interactive Audio* (2019 March).

- [33] S. J. Elliott, *Signal Processing for Active Control* (Academic Press, London), 1st ed. (2001).
- [34] F. M. Fazi, E. Hamdan, “Stage Compression in Transaural Audio,” presented at the *Audio Engineering Society Convention 144* (2018 May).
- [35] M. Poletti, “Robust Two-Dimensional Surround Sound Reproduction for Nonuniform Loudspeaker Layouts,” *J. Audio Eng. Soc.*, vol. 55, no. 7/8, pp. 598–610 (2007 July).
- [36] B. Masiero, M. Vorländer, “A Framework for the Calculation of Dynamic Crosstalk Cancellation Filters,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 9, pp. 1345–1354 (2014 Sep.), DOI: <https://doi.org/10.1109/TASLP.2014.2329184>.
- [37] M. F. Simón Gálvez, S. J. Elliott, J. Cheer, “A Superdirective Array of Phase Shift Sources,” *The Journal of the Acoustical Society of America*, vol. 132, no. 2, pp. 746–756 (2012 Aug), DOI: <https://doi.org/10.1121/1.4733556>.
- [38] J. Borenus, “Perceptibility of Direction and Time Delay Errors in Subwoofer Reproduction,” presented at the *Audio Engineering Society Convention 79* (1985 Oct).

APPENDIX

The aim is to minimize the cost function given in Eq. 9. This can be rearranged to give

$$J = \mathbf{b}^H \mathbf{b} + \mathbf{v}^H \mathbf{C}^H \mathbf{C} \mathbf{v} - 2\mathbf{v}^H \mathbf{C}^H \mathbf{b} + \beta \mathbf{v}^H \tilde{\Gamma}^H \tilde{\Gamma} \mathbf{v}. \quad (24)$$

The cost function is minimized by finding the derivative with respect to \mathbf{v} and setting the result equal to zero. Hence

$$\frac{\partial J}{\partial \mathbf{v}} = 2(\mathbf{C}^H \mathbf{C} \mathbf{v} - \mathbf{C}^H \mathbf{b} + \beta \tilde{\Gamma}^H \tilde{\Gamma} \mathbf{v}) = 0. \quad (25)$$

Rearranging gives

$$\begin{aligned} \mathbf{v} &= [\mathbf{C}^H \mathbf{C} + \beta \tilde{\Gamma}^H \tilde{\Gamma}]^{-1} \mathbf{C}^H \mathbf{b} \\ &= \mathbf{H} \mathbf{b} \end{aligned} \quad (26)$$

and finally

$$\mathbf{H} = [\mathbf{C}^H \mathbf{C} + \beta \tilde{\Gamma}^H \tilde{\Gamma}]^{-1} \mathbf{C}^H. \quad (27)$$