

# Localization of Sound Sources in a Room with One Microphone

Helena Peić Tukuljac, Hervé Lissek and Pierre Vandergheynst

School of Computer and Communication Sciences  
École polytechnique fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland

## ABSTRACT

Estimation of the location of sound sources is usually done using microphone arrays. Such settings provide an environment where we know the difference between the received signals among different microphones in the terms of phase or attenuation, which enables localization of the sound source. In our solution we exploit the properties of the room transfer function in order to localize a sound source inside a room with only one microphone. Shape of the room and the position of the microphone are assumed to be known. The design guidelines and the limitations of the sensing matrix are given. Implementation is based on the sparsity in the terms of voxels in a room that are occupied by a source. What is especially interesting about our solution is that we provide localization of the sound sources not only in the horizontal plane, but in the terms of the 3D coordinates inside the room.

**Keywords:** resonant frequency, room impulse response, room mode, sparsity, sound source localization

## 1. INTRODUCTION

In the last decade the theory of compressed sensing [1, 2](#) has found its way in the domain of acoustic signal processing. There was always a need for finding a structure in the high dimensional acoustical data that was cumbersome to handle. In 2015 Boche et al. [3](#) have given a detailed state of the art for the application of compressed sensing in the domains of image and acoustic signal processing.

The origins of sparsity in acoustical data include, but are not limited to: voxels (directions of arrival) occupied by sound sources which is usually exploited for the localization of the sound sources in free field [4](#), or in rooms for the estimation of the sound pressure distribution [5](#). The sparsity in the image-source model has been mainly used for the estimation of the room shape [6](#), the sparsity of plane wave representation of the sound pressure is used for the characterization of the sound pressure inside the room [7](#). There exists another type of acoustic sparsity that is exploited in a few solutions - the sparsity of the room modes in the low-frequency part of the room transfer function (RTF) [8](#).

Our idea is to combine the sparsity that exists in the term of the voxels of a room occupied by the sound sources and the low-frequency room modes in the RTF. So, we will observe the transfer function below the so called Schröeder frequency, which is defined as:  $f_s = 2000\sqrt{\frac{RT_{60}}{V}}$ , where  $V$  is the volume of the room and  $RT_{60}$  is the reverberation time. This combination should result in a fast localization of sound sources by only one microphone as will be further explained.

To the best of our knowledge similar solutions have not been published so far.

The reminder of the paper is organized as follows: In Section [2](#) we discuss the sparsity that exists in the low frequency domain of the room transfer function. Section [3](#) gives a general introduction to compressed sensing and its application to the localization of sound sources. The design and the limitations of the sensing matrix for our case is given in Section [4](#) and final remarks and conclusions are given in Section [5](#).

---

Further author information:

Helena Peić Tukuljac: helena.peictukuljac@epfl.ch

Hervé Lissek: herve.lissek@epfl.ch

Pierre Vandergheynst: pierre.vandergheynst@epfl.ch

## 2. ROOM TRANSFER FUNCTION AND ITS LOW-FREQUENCY PROPERTIES

In the further development of our approach, we are going to rely on two facts: the room shape is known and the microphone position is known. These assumptions imply that we know the resonant frequencies of the room and the room modes related to the microphone's positions.

RTF between arbitrary two points  $\mathbf{r}'$  and  $\mathbf{r}''$  in a room in the Fourier domain is given by 9:

$$H_\omega(\mathbf{r}') = \rho_0 c^2 \omega Q \sum_n \frac{p_n(\mathbf{r}') p_n(\mathbf{r}'')}{K_n [2\delta_n \tilde{\omega}_n + i(\omega^2 - \tilde{\omega}_n^2)]} \quad (1)$$

where:  $\rho_0$  is the density of the propagating medium (air),  $c$  is the sound celerity,  $Q$  is a constant related to the loudspeaker (sound source),  $p_n(\cdot)$  are the room modes,  $K_n$  is the volume dependent scaling factor,  $\delta_n$  is the damping and  $\tilde{\omega}_n$  are the resonant frequencies. One of these points ( $\mathbf{r}'$  and  $\mathbf{r}''$ ) represents the position of the sound source and the other one represents the position of the microphone. We can notice an interesting underlying symmetry that exists in this equation: if we exchange the positions of the sound source and the microphone, the transfer function will remain the same.

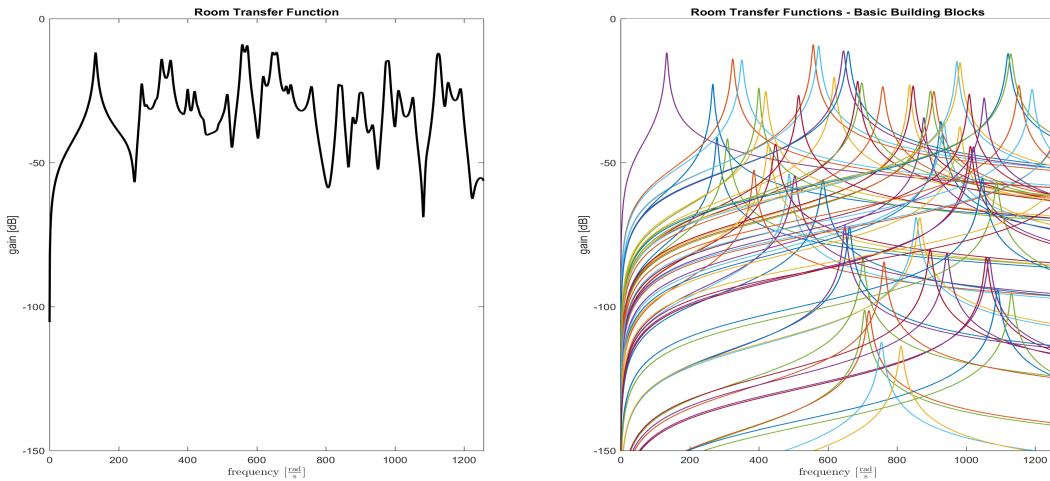


Figure 1. Basic components of the RTF are called room modes. As illustrated, room modes are simply second order bandpass filters.

In Figure 1, we can see a segment of RTF up to 200Hz and its decomposition into the room modes. The sharpness of the peaks of room modes is dependent on the damping properties of walls of the room. Peaks of the room modes are aligned with the resonant frequencies of the room.

The resonant frequencies of a rectangular room of size  $L_x \times L_y \times L_z$  are given by the expression:  $\tilde{\omega}_r = \pi c \sqrt{\left(\frac{n_x}{L_x}\right)^2 + \left(\frac{n_y}{L_y}\right)^2 + \left(\frac{n_z}{L_z}\right)^2}$  where  $(n_x, n_y, n_z) \in \mathbb{N}_0^3 \setminus (0, 0, 0)$ .

### 2.1 Room Transfer Function at Different Positions Across the Room

In 1985, Richardson et al. 10 have proposed a curve fitting algorithm for reconstruction of the RTF curve from discrete measurements by using room mode shaped functions as basic fitting elements. Every RTF is characterized by a set of parameters: resonant frequencies, which determine the position of the peaks of room modes, and damping, attenuation and phase of these room modes.

For different positions of the microphones/sound sources across the room, some parameters stay the same. These are the resonant frequencies which depend on the room shape, and the room mode damping which depends on the damping of the wall - *common parameters*. The attenuation and the phase of the room modes are position dependent parameters - *specific parameters*.

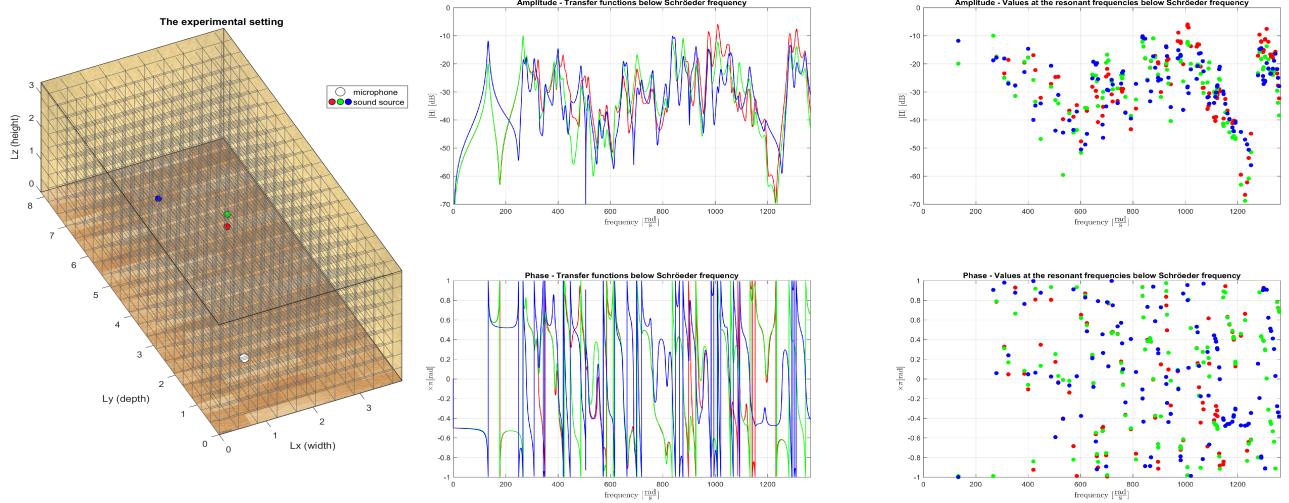


Figure 2. Values of the RTF across the room vary in the terms of attenuation and phase value at the resonant frequencies. We exploit only the difference in the magnitude because in our target experimental setting there exists only one microphone and the sources will emit white noise.

Figure 2. illustrates the difference between the attenuation and the phase of the RTFs across the room at the resonant frequencies. As can be observed, although that all the positions of the sound sources result in the peaks at the same set of frequencies (the resonant frequencies of the room), the set of the heights of these peaks seems unique (this will be further observed in the next section). This means that each pair of the positions of the sound source and the microphone could potentially result in a unique set of attenuation factors at the resonant frequencies.

Although that there exists uniqueness of phase for each room mode, since we plan to use only one microphone and white noise sources, this is irrelevant for our case but has a potential for other types of room characterization. We have decided to investigate the potential of unique representation of position of the sound source within the room with the set of attenuations of RTF at resonant frequencies.

Therefore we have established a valuable reasoning for the design of our sensing matrix.

## 2.2 Relation Between Room Modes and Plane Waves

In a rectangular room, each room mode (eigenmode of the Laplacian operator) represents a sum of 8 plane waves that share a wave number:

$$p(\mathbf{k}_n, \mathbf{r}_m) = \sum_{i=1}^8 a_i e^{j(\mathbf{S}(:,j) \odot \mathbf{k}_n) \cdot \mathbf{r}_m} \quad (2)$$

where  $\odot$  is a Hadamard product,  $\mathbf{S}_{3 \times 8}$  is a sign matrix whose columns alternate from  $[1, 1, 1]^T$  to  $[-1, -1, -1]^T$ ,  $\mathbf{k}_n = (\frac{n_x \pi}{L_x}, \frac{n_y \pi}{L_y}, \frac{n_z \pi}{L_z})$ ,  $(n_x, n_y, n_z) \in \mathbb{N}_0^3 \setminus (0, 0, 0)$ , is a wave vector that uniquely defines the  $n^{\text{th}}$  room mode and  $\mathbf{r}_m$  is a position inside the room.

As can be seen in Figure 3., these vectors are just corners of a parallelepiped ( $\mathbf{k} = [\pm k_x, \pm k_y, \pm k_z]^T$ ). We can also notice the periodicity of the wave vector grid:  $\frac{\pi}{L_x}, \frac{\pi}{L_y}, \frac{\pi}{L_z}$  along each of the axes.

In a theoretical case in which all the walls are perfectly rigid all the plane waves have the same expansion coefficient ( $\forall i, a_i = a$ ), so our sum of the 8 plane waves can be represented as a product of cosine functions:

$$p(\mathbf{k}_n, \mathbf{r}_m) \sim \cos\left(\frac{n_x \pi}{L_x} \mathbf{r}_m(x)\right) \cos\left(\frac{n_y \pi}{L_y} \mathbf{r}_m(y)\right) \cos\left(\frac{n_z \pi}{L_z} \mathbf{r}_m(z)\right). \quad (3)$$

where  $a$  is a constant. An example of room mode for a room with rigid walls is given in Figure 4.

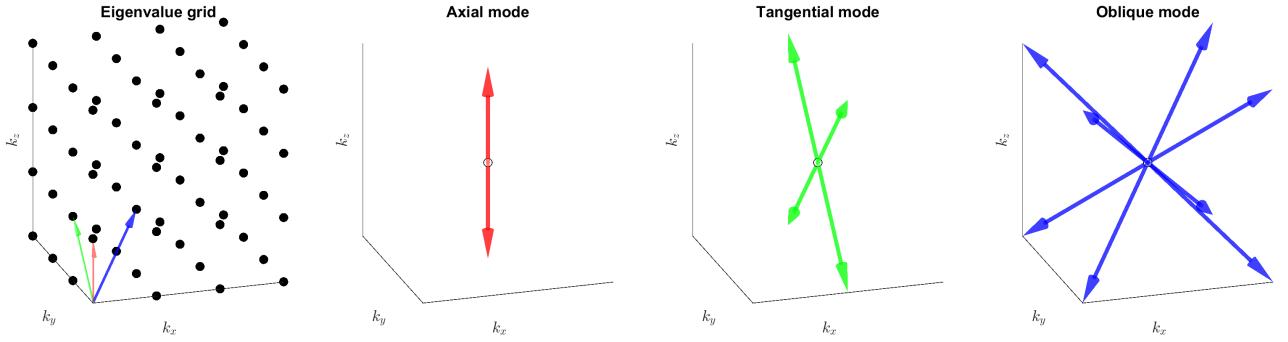


Figure 3. Eigenvalue space of a rectangular room with rigid walls. The left-hand side shows just one octant because of the symmetry that exists (there are 8 plane waves for each wave number). The length of the wave vector is proportional to the eigenvalue of the Laplacian.

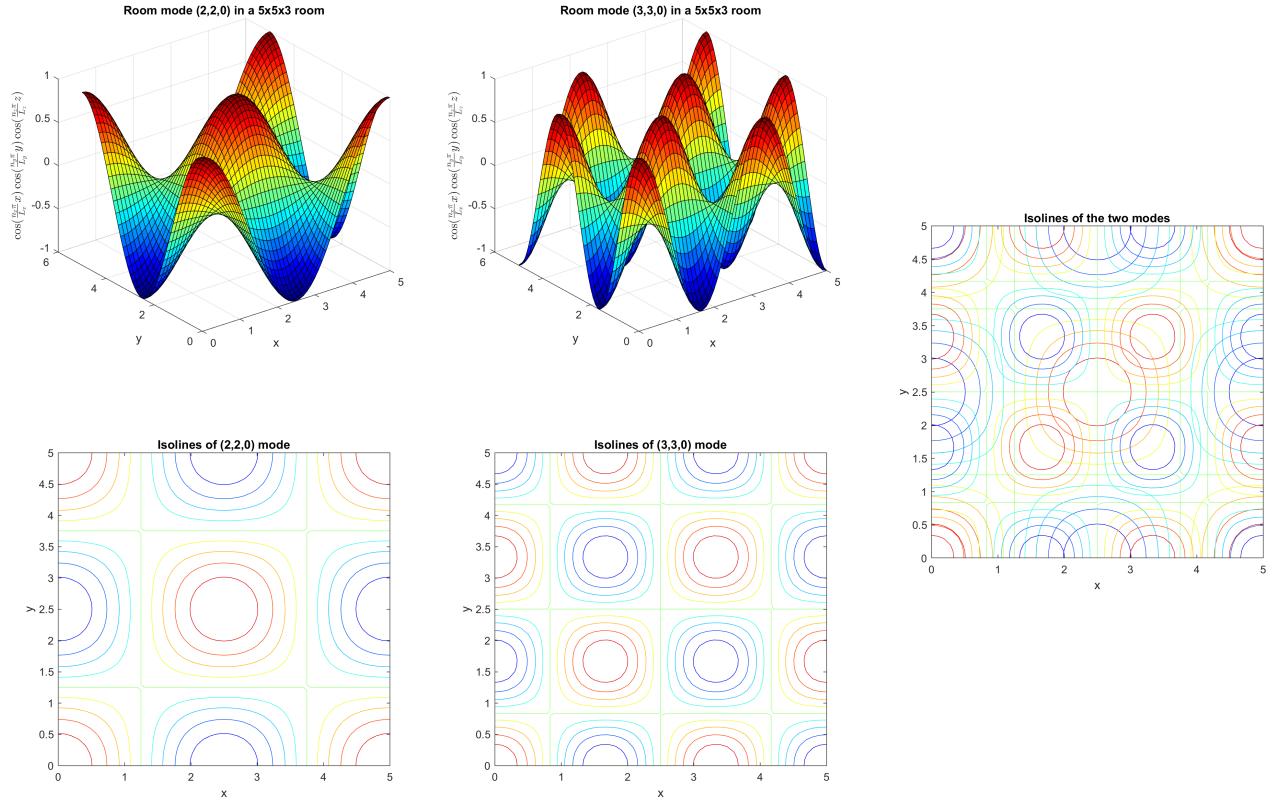


Figure 4. An example of  $(n_x, n_y, n_z) \in \{(2, 2, 0), (3, 3, 0)\}$  room modes in a  $5\text{m} \times 5\text{m} \times 3\text{m}$  room with rigid walls. We can notice that the isolines of different modes intersect in just a few locations, which supports our assumption of different height of sets of peaks in the RTF.

### 2.3 Ambiguities that Exist in the Terms of Uniqueness of the Attenuation Across the Room

We will observe the basic axial modes in Figure 5. Although the sound pressure value function is in 5D, we can visualize only values in 3D. First row shows the  $x$ - and  $y$ -axial modes (everything that will be said applies analogously to  $z$ -axial modes as well). We can see that these two modes form pairs of points that result in a unique location identifier. But, since we have decided to explore the special case with only one microphone, we need to neglect the phase of the RTF, therefore we can just observe the absolute value of the RTF. As seen in the second row of the same figure, this introduces ambiguity - there exists a unique representation, but only in  $\frac{1}{8}$  of the room.

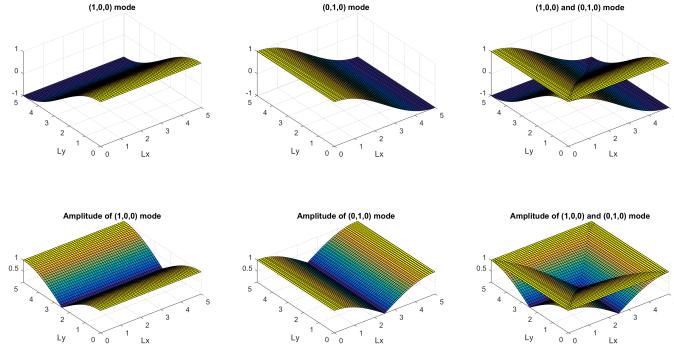


Figure 5. Basic modes and their attenuation values.

This ambiguity is illustrated in Figure 6. Here we see 3 modes of two different positions in a room. The small mode in the middle has the same amplitude and phase and the other two modes have an opposite phase, which can not be seen when we project it to neglect the phase. This means that we can not rely only on the basic axial modes for the sound source localization.

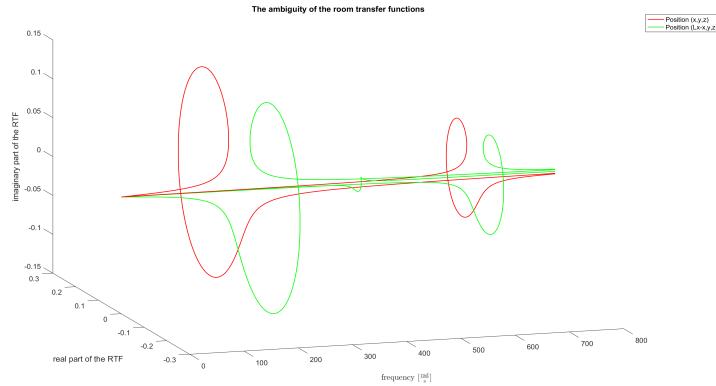


Figure 6. Ambiguities that exist in the term of the uniqueness of the RTF.

## 3. COMPRESSED SENSING AND SOUND SOURCE LOCALIZATION

### 3.1 Sparse Representation of the Position of Sources

In source localization problems the domain of interest is usually divided into an angular grid such that the sources occupy just a few of these angles. Since our sources are positioned inside a room, we will divide the room into

voxels and assume that the number of voxels occupied by a source is small. We recognize that this is a problem with underlying sparsity. These problems are usually solved by using the theory of compressed sensing.

### 3.2 Compressed Sensing

Our signal of interest  $\mathbf{y}$  is the measurement of sound pressure at a known location inside a known room:

$$\mathbf{y} = \Psi \Phi \mathbf{x} \quad (4)$$

where  $\mathbf{y} \in \mathbb{R}^N$  are the sound pressure measurements,  $\Psi_{N \times N}$  is the inverse Fourier Transform (represents the change of domain),  $\Phi_{N \times M}$  is a representational dictionary (columns of this matrix are called atoms) with the RTFs as columns and  $\mathbf{x} \in \mathbb{R}^M$  are the sparse expansion coefficients.  $\mathbf{x}$  is  $K$ -sparse, which means that it contains at most  $K$  non-zero elements and  $K \ll N$ . Since  $M > N$  we are facing an underdetermined system of equations. The problem in this form is non-convex. The product  $\mathbf{A} = \Psi \Phi$  is usually referred to as the sensing matrix .

These types of problems are usually solved using one of the 5 groups of approaches listed in 11, where the most common ones are the convex relaxation and the greedy pursuit. The convex relaxation 12,13, which is also known as the Basis Pursuit, relies on the relaxation of the minimization of the  $\ell_0$ -norm to  $\ell_1$ -norm which favors the sparse solutions, although at a cost of requiring higher number of measurements 14. There also exists a relaxation to  $\ell_2$ -norm, but this norm favors the minimization of the energy of the signal rather than the sparsity. In practice convex relaxation approach is usually used for smaller and medium size problems, because large scale data causes computational issues.

In our solution we will rely on the greedy approaches such as Orthogonal Matching Pursuit (OMP) 15 and Compressive Sampling Matching Pursuit (CoSaMP) 16. These methods select up to  $K$  atoms of a dictionary that give the least approximation error. CoSaMP is a faster contemporary method which works by selecting multiple atoms at every iteration.

Regardless of the approach, one of the main advantages of compressed sensing technique is the robustness to noise since we project our signal to the vectors that span the signal space, and therefore neglecting the residual related to the existing noise.

### 3.3 Conditions for Dictionary Design

#### 3.3.1 Spark and coherence of the dictionary

Spark of a matrix  $\Phi$  is the smallest number of linearly dependent columns of matrix  $\Phi$ . The requirement for the sensing matrix  $\Phi$  in compressed sensing is that the following holds:

$$\text{spark}(\Phi) > 2K \quad (5)$$

where  $K$  is the level of sparsity. In other words: To achieve an injective mapping we need to assure that there are no two  $k$ -sparse vectors that map to the same measurements. This implies that the rank of our sensing matrix has to be at least  $2k$  which is tightly related to the restriction on the coherence of the dictionary.

According to the theory of compressed sensing, we have to ensure the appropriate coherence parameter of a dictionary:  $\mu = \max_{1 \leq i < j \leq n} \frac{|\langle \varphi_i, \varphi_j \rangle|}{\|\varphi_i\|_2 \|\varphi_j\|_2} = \max_{1 \leq i < j \leq n} \cos \angle(\varphi_i, \varphi_j)$ , so coherence is the cosine of the acute angle between the closest pair of atoms in a given dictionary. We want our dictionary to be incoherent, so  $\mu$  should be the smallest possible. The ideal case is the case where we have orthogonal atoms with the coherence parameter equal to zero among different atoms of the dictionary.

#### 3.3.2 Restricted Isometry Property

The restricted isometry property guarantees that the distances (lengths) are preserved when moving from one space to another. Let  $\Phi$  be an  $M \times N$  matrix and let  $1 \leq K \leq N$  be an integer. Suppose that there exists a constant  $\delta_K \in (0, 1)$  such that, for every  $M \times K$  submatrix  $\Phi_K$  of  $\Phi$  and for every  $K$ -sparse vector  $\mathbf{y}$ ,

$$(1 - \delta_K) \|\mathbf{y}\|_2^2 \leq \|\Phi_K \mathbf{y}\|_2^2 \leq (1 + \delta_K) \|\mathbf{y}\|_2^2. \quad (6)$$

Then, the matrix  $\Phi$  is said to satisfy the  $K$ -restricted isometry property with restricted isometry constant  $\delta_K$ . In most cases it is hard to check whether this property holds or not.

### 3.4 Sound Source Localization in Underwater Acoustics

In 2014 Xenaki et al. [4] have used the theory of compressed sensing [1,2] for the estimation of the direction of arrival of sound sources in the free-field setting for underwater acoustics. The columns of their sensing matrix have the following form:

$$\mathbf{a}(\theta_i) = \frac{1}{\sqrt{M}} e^{j \frac{2\pi}{\lambda} \mathbf{r} \sin \theta_i} \quad (7)$$

where  $M$  is the number of microphones,  $\mathbf{r} = [r_1; \dots; r_M]^T$  are the locations of the microphones,  $\lambda$  is the wavelength and  $\theta_i \in [-90^\circ, 90^\circ]$  is the direction of arrival of associated plane wave with respect to the array axis. The elements of each column represent the propagation delay from the  $i^{\text{th}}$  potential source to each of the microphones. So the sensing matrix (dictionary) has the form:  $\mathbf{A} = [\mathbf{a}(\theta_1), \dots, \mathbf{a}(\theta_N)]$ . The geometry of the microphone array is assumed to be known.

Low coherence of the dictionary is shown by observing the maximum off-diagonal element of the Gram matrix:  $\mathbf{G}_{N \times N} = |\mathbf{A}^H \mathbf{A}|$ . This matrix is approximately equal to the identity matrix for appropriate microphone spacing and angular resolution.

### 3.5 Sound Source Localization in a Room

The following question rises: How to tailor a simple incoherent dictionary for fast localization of sources inside the room? In order to have a well-posed problem we introduce the following assumptions:

1. the shape of the room and the reverberation time are known,
2. the position of the microphone is known, and
3. all the sound sources have a flat spectrum in the observed frequency range.

The atoms of our localization dictionary are going to consist of the height of the peaks in the RTF at the resonant frequencies. In most cases the Restricted Isometry Property is hard to check. We know that random matrices, which were used as the dictionaries in the early stages of compressed sensing, satisfy this property. Therefore we will choose the potential position of the sources on uniformly at random on the regular grid.

For each of the potential positions of sound sources and a fixed position of the microphone we have one atom in the dictionary which consists out of the height of the peaks in the RTFs at the resonant frequencies. The height of the sensing matrix is proportional to the number of the resonant frequencies in the observed frequency range. The number of resonant frequencies below a given frequency  $f_s$  [9] can be computed by:  $N(f_s) = \frac{4}{3}\pi V\left(\frac{f_s}{c}\right)^3 + \frac{1}{4}\pi S\left(\frac{f_s}{c}\right)^2 + \frac{1}{2}L\frac{f_s}{c}$ , where:  $V = L_x L_y L_z$ ,  $S = 2(L_x L_y + L_y L_z + L_z L_x)$  and  $L = L_x + L_y + L_z$ . The width of the dictionary is proportional to the number of observation points on the predefined grid.

In order to localize the sources, we search for a subset of atoms that give the best fitting for the signal recorded by the microphone. Once we discover which atoms of our sensing matrix have the highest expansion coefficients in the sparse representation, we can easily recover the position of the sound sources in the room, because we know which atom corresponds to which position.

## 4. DESIGNING AN EFFICIENT SENSING MATRIX

### 4.1 Coherence

Coherence of a dictionary can be seen from the maximum off-diagonal element of the coherence Gram matrix  $\mu = \max_{i \neq j} |\mathbf{G}_{ij}|$ . In our case where  $\Psi_{N \times N}$  is the inverse Fourier Transform and  $\Phi_{N \times M}$  is the matrix with the RTF coefficients, the Gram matrix has the following form:

$$\mathbf{G} = |\mathbf{A}^H \mathbf{A}| = |(\Psi \Phi)^H \Psi \Phi| = |\Phi^H \Psi^H \Psi \Phi|. \quad (8)$$

Since the Fourier matrix has orthonormal atoms up to a scaling constant  $\Psi^H \Psi = \frac{1}{N} \mathbf{I}$ , we have:

$$\mathbf{G} = \frac{1}{N} \Phi^H \Phi. \quad (9)$$

Therefore we observe the coherence of the sensing matrix by focusing on the discretization of the room transfer function.

## 4.2 Battle of the Grids

Our problem has two degrees of freedom and both of them represent a selection process of the nodes on a grid. We have a grid of wave vectors - *features* and a grid of potential positions of sound sources - *samples*. In Figure 7 the grid on the left-hand side repeats in all 6 directions and the one on the right-hand side repeats in 3 directions.

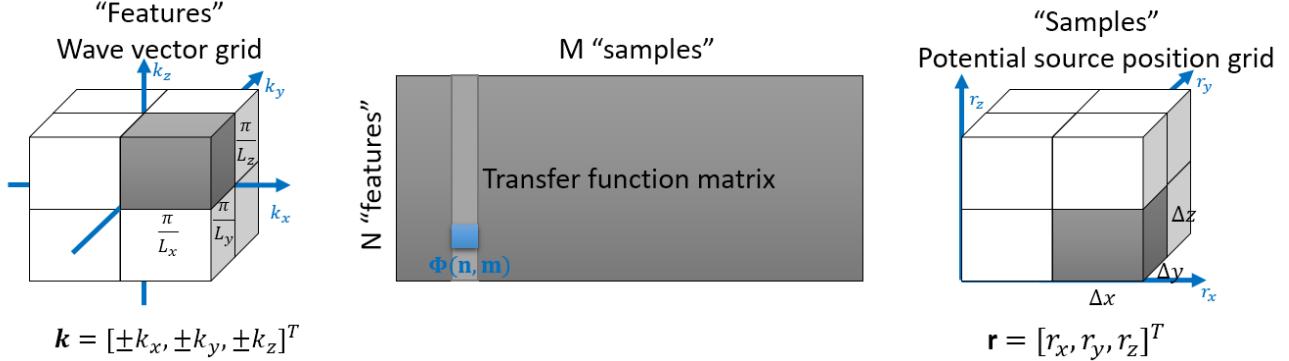


Figure 7. Two grids that represent two degrees of freedom that we have for designing an incoherent sensing matrix.

We will observe the room transfer function in a matrix form at the resonant frequencies. If we go back to equation (1) and introduce  $\omega = \tilde{\omega}_n$ , we get that each of the elements of our sensing matrix  $\Phi_{N \times M}$  is of the following form:

$$\Phi(n, m) = \frac{\rho_0 c^2 Q_m}{2 K_n \delta_n} p(\mathbf{k}_n, \mathbf{r}_{\text{mic}}) p(\mathbf{k}_n, \mathbf{r}_m) \quad (10)$$

which corresponds to  $n^{\text{th}}$  wave vector and  $m^{\text{th}}$  potential sound source position. The only coefficients that differ among the atoms of the dictionary are represented in blue. The difference due to the quality factor of the sound source  $Q$ , will not affect our approach. It has an effect only on the expansion coefficients of the sparse representation. Therefore we focus on the sound sources' position that produces different attenuation of room modes.

So the RTF matrix has the following decomposition:

$$\Phi = \frac{\rho_0 c^2}{2} \begin{bmatrix} \frac{p(\mathbf{k}_1, \mathbf{r}_{\text{mic}})}{K_1 \delta_1} & \dots & \frac{p(\mathbf{k}_1, \mathbf{r}_{\text{mic}})}{K_1 \delta_1} \\ \vdots & \ddots & \vdots \\ \frac{p(\mathbf{k}_N, \mathbf{r}_{\text{mic}})}{K_N \delta_N} & \dots & \frac{p(\mathbf{k}_N, \mathbf{r}_{\text{mic}})}{K_N \delta_N} \end{bmatrix} \odot \begin{bmatrix} Q_1 p(\mathbf{k}_1, \mathbf{r}_1) & \dots & Q_M p(\mathbf{k}_1, \mathbf{r}_M) \\ \vdots & \ddots & \vdots \\ Q_1 p(\mathbf{k}_N, \mathbf{r}_1) & \dots & Q_M p(\mathbf{k}_N, \mathbf{r}_M) \end{bmatrix}. \quad (11)$$

Just to recall, our rigid wall room modes are of the form:  $p(\mathbf{k}_n, \mathbf{r}_m) = \sum_{i=1}^8 e^{j(\mathbf{S}(:, j) \odot \mathbf{k}_n) \cdot \mathbf{r}_m}$ , where  $\mathbf{k}_n$  belongs to the positive octant of the left-hand side grid.

For a uniform case, the off-diagonal elements of our Gram matrix are of the following form:

$$\mathbf{G}_{ij} = \cos(k_x r_x) \cos(k_x(r_x \pm m \Delta x)) + \cos(k_y r_y) \cos(k_x(r_y \pm n \Delta y)) + \cos(k_z r_z) \cos(k_x(r_z \pm o \Delta z)). \quad (12)$$

Due to the smoothness of cosine function, the points on the potential sound source position grid that lay close result in similar heights of the peaks in RIR.

## 4.3 The Uniqueness of the Atoms of the Dictionary and Results

We can restate our problem in a following way:  $\mathbf{S}_{rf} \Psi^* y = \mathbf{S}_{rf} \mathbf{S}_{sp} \Phi x$  where  $y$  is the measured signal,  $\mathbf{S}_{rf}$  is a resonant frequency selector (defines which points on the wave vector grid we observe),  $\mathbf{S}_{sp}$  is a sound source position selector (defines which points on the potential source position grid we observe),  $\Psi^*$  is the Fourier transform and  $x$  is sparse. Support of  $x$  shows which of the positions on the grid are the most probable positions of the sources.

Since our exponentials in the plane wave representation are not equidistant, we can not apply the Dirichlet kernel sum to our case. Therefore our selectors will select the nodes from the grid in a uniformly at random way which should reduce the coherence of our dictionary due to the Steinhaus theorem for sequences 17.

Candès et al. 18 discuss the potential of recovery of data that has a sparse representation in a coherent dictionary. Coherent dictionaries can give guarantees only on the recovery of the sparse signal, but not on the recovery of the index of the atoms in sparse representation. That is because if we have pairs of atoms that are extremely coherent (almost collinear), we can not tell which one of them will be used for our sparse representation when projecting to a lower-dimensional space.

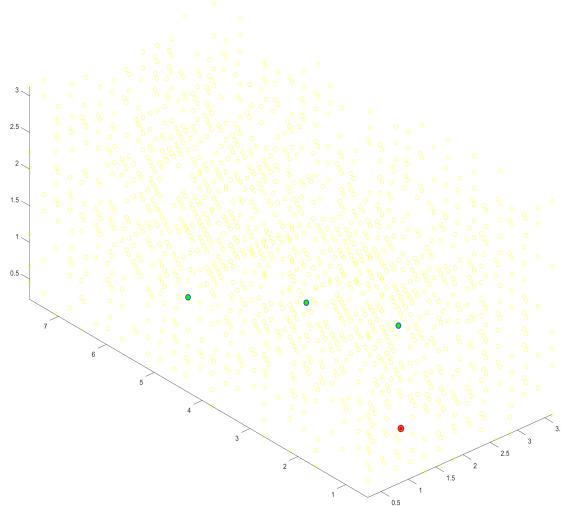


Figure 8. These are the results for localization of 3 sound sources inside the room.

Figure 8. shows a reconstruction example for a case with 3 sound sources. Yellow circles are the potential positions taken into account, blue rings are the true positions and green points are the reconstructed positions. The red point represents the known position of the sound source. Here is the description of the algorithm:

---

#### **Algorithm 1** Localization of sound sources in a room with one microphone

---

**Input :** Size of the room ( $L_x \times L_y \times L_z$ ), position of the sound source  $r_{\text{mic}}$ , ground truth positions of the sound sources, uniform grid of potential points of the sound sources  $\mathbf{P}$ , uniform grid of the wave vectors  $\mathbf{W}$ , dictionary dimension -  $N$  and  $M$ .

**Output:** Reconstructed positions of the sound sources.

**do**

Select  $N$  random wave vectors and  $M$  random potential positions of the sound sources from the uniform grids.

Generate a room mode dictionary  $\Phi_{N \times M}$ .

Try to estimate the positons of the sound sources.

**while** CoSaMP 16 sparse representation does not converge (has a residual close to zero)

---

#### 4.4 Precision and Basis Mismatch

Due to the smoothness of the room mode functions, there is a small variation in the value between the close points. This supports the idea of similarity of the atoms of the dictionary of the spatially close positions.

Compressed sensing usually assumes the existence of a grid with finite density and our signals of interests can fail to coincide with the nodes of the predefined grid, especially in the case of moving sources. As shown in 19 this can cause that sparse signals appear as incompressible. The work we have observed before 4 has an extension to a continuous case 20 by applying the semi-definite programming. In our observations we have assumed that our grid of the potential positions of the sources is dense enough to avoid the spectral leakage and continuous approaches will be left for future work.

#### 4.5 Requirements and Limitations

In a setting where we have multiple sound sources and a microphone, the sound received is equal to the linear combination of the convolution of sounds emitted by the sound sources and the transfer functions that correspond to their positions. Therefore we need to introduce another assumption: we can efficiently localize sources which have a flat spectrum in the observed frequency band, since they result in a nearly constant Fourier coefficients of emitted signals' spectrum. Otherwise we have to know upfront the signals that will be emitted by sources.

### 5. CONCLUSION

Unlike most of the localization algorithms, this approach guarantees the localization in 3D, without neglecting the elevation angle, which is rarely estimated. Although that we have mainly focused on the shoe-box shaped rooms, this method is applicable to any room shape as long as we know the form of the room modes. The simplicity of our solution lays in the low required prior knowledge about the room - only the height of the peaks in the RTF at the resonant frequencies for modes should be known.

Our solution has the potential of being applied to the optimization of the quality of the hearing aids as well as in robotics for monoaural localization.

Future work will include estimation *off the grid* in order to avoid the basis mismatch and the challenging computational costs.

### 6. SUPPLEMENTARY MATERIALS

*matlab* code used for generating each of the figures in this paper as well as the acoustical room mode framework is available for download on the following link: [https://github.com/epfl-lts2/room\\_transfer\\_function\\_toolkit](https://github.com/epfl-lts2/room_transfer_function_toolkit). *python* version of the toolkit is also available.

### ACKNOWLEDGMENTS

The work of H. Peić Tukuljac was supported by the Swiss National Science Foundation under Grant No. 200021\_169360 for the project “*Compressive Sensing applied to the Characterization and the Control of Room Acoustics*”.

### REFERENCES

- [1] Donoho, D. L., “Compressed sensing,” *IEEE Trans. Information Theory* **52**(4), 1289–1306 (2006).
- [2] Candes, E. J., Romberg, J., and Tao, T., “Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information,” *IEEE Trans. Inf. Theor.* **52**, 489–509 (Feb. 2006).
- [3] Boche, H., Calderbank, R., Kutyniok, G., and Vybral, J., [*Compressed Sensing and Its Applications: MATH-EON Workshop 2013*], Birkhäuser Basel, 1st ed. (2015).
- [4] Xenaki, A., Gerstoft, P., and Mosegaard, K., “Compressive beamforming,” *The Journal of the Acoustical Society of America* **136**(1), 260–271 (2014).
- [5] Kitić, S., Bertin, N., and Gribonval, R., “Hearing behind walls: Localizing sources in the room next door with cosparsity,” in [*2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*], 3087–3091 (May 2014).
- [6] Dokmanić, I., Parhizkar, R., Walther, A., Lu, Y. M., and Vetterli, M., “Acoustic echoes reveal room shape,” *Proceedings of the National Academy of Sciences* **110**(30), 12186–12191 (2013).

- [7] Koyano, Y., Yatabe, K., Ikeda, Y., and Oikawa, Y., “Physical-model based efficient data representation for many-channel microphone array,” in [*2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*], 370–374 (March 2016).
- [8] Mignot, R., Chardon, G., and Daudet, L., “Low frequency interpolation of room impulse responses using compressed sensing,” *IEEE/ACM Trans. Audio, Speech and Lang. Proc.* **22**, 205–216 (Jan. 2014).
- [9] Kuttruff, H. and Mommertz, E., [*Room Acoustics*], 239–267, Springer Berlin Heidelberg, Berlin, Heidelberg (2013).
- [10] Richardson, M. H. and Formenti, D. L., “Global curve fitting of frequency response measurements using the rational fraction polynomial method,” (1985).
- [11] Tropp, J. A. and Wright, S. J., “Computational methods for sparse solution of linear inverse problems,” *Proceedings of the IEEE* **98**, 948–958 (June 2010).
- [12] Tropp, J. A., “Just relax: convex programming methods for identifying sparse signals in noise,” *IEEE Transactions on Information Theory* **52**, 1030–1051 (March 2006).
- [13] Boyd, S. and Vandenberghe, L., [*Convex Optimization*], Cambridge University Press, New York, NY, USA (2004).
- [14] Candès, E. J., Wakin, M. B., and Boyd, S. P., “Enhancing sparsity by reweighted l1 minimization,” *Journal of Fourier Analysis and Applications* **14**(5), 877–905 (2008).
- [15] Tropp, J. A., Gilbert, A. C., and Strauss, M. J., “Algorithms for simultaneous sparse approximation: Part i: Greedy pursuit,” *Signal Process.* **86**, 572–588 (Mar. 2006).
- [16] Needell, D. and Tropp, J., “Cosamp: Iterative signal recovery from incomplete and inaccurate samples,” *Applied and Computational Harmonic Analysis* **26**(3), 301 – 321 (2009).
- [17] Rauhut, H., “Compressive sensing and structured random matrices,” *Theoretical Foundations and Numerical Methods for Sparse Recovery* **9**(4), 1–92 (2010).
- [18] Candès, E. J., Eldar, Y. C., and Needell, D., “Compressed sensing with coherent and redundant dictionaries,” *CoRR* **abs/1005.2613** (2010).
- [19] Chi, Y., Pezeshki, A., Scharf, L., and Calderbank, R., “Sensitivity to basis mismatch in compressed sensing,” in [*2010 IEEE International Conference on Acoustics, Speech and Signal Processing*], 3930–3933 (March 2010).
- [20] Xenaki, A. and Gerstoft, P., “Grid-free compressive beamforming,” *CoRR* **abs/1504.01662** (2015).