

# SPATIAL MIXUP: DIRECTIONAL LOUDNESS MODIFICATION AS DATA AUGMENTATION FOR SOUND EVENT LOCALIZATION AND DETECTION

Ricardo Falc3n-P3rez<sup>\*\*</sup>, Kazuki Shimada<sup>†</sup>, Yuichiro Koyama<sup>†</sup>, Shusuke Takahashi<sup>†</sup>, Yuki Mitsufuji<sup>†</sup>

<sup>\*</sup>Aalto University, Espoo, Finland

<sup>†</sup>Sony Group Corporation, Tokyo, Japan

## ABSTRACT

Data augmentation methods have shown great importance in diverse supervised learning problems where labeled data is scarce or costly to obtain. For sound event localization and detection (SELD) tasks several augmentation methods have been proposed, with most borrowing ideas from other domains such as images, speech, or monophonic audio. However, only a few exploit the spatial properties of a full 3D audio scene. We propose Spatial Mixup, as an application of parametric spatial audio effects for data augmentation, which modifies the directional properties of a multi-channel spatial audio signal encoded in the ambisonics domain. Similarly to beamforming, these modifications enhance or suppress signals arriving from certain directions, although the effect is less pronounced. Therefore enabling deep learning models to achieve invariance to small spatial perturbations. The method is evaluated with experiments in the DCASE 2021 Task 3 dataset, where spatial mixup increases performance over a non-augmented baseline, and compares to other well known augmentation methods. Furthermore, combining spatial mixup with other methods greatly improves performance.

**Index Terms**— Sound event localization and detection, spatial audio, sound source localization, acoustic scene analysis, data augmentation

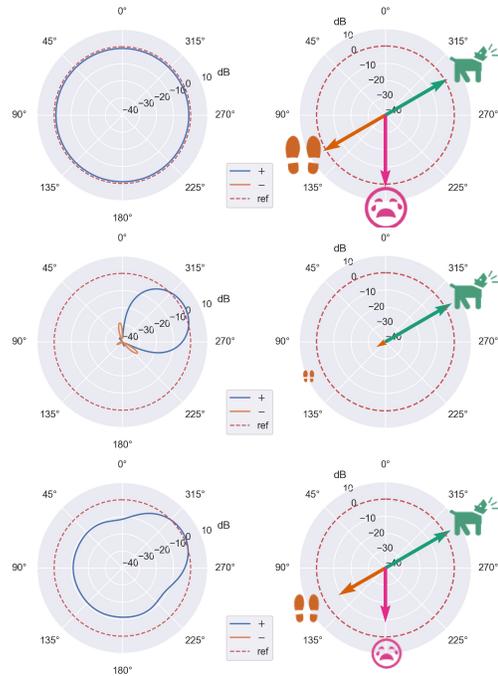
## 1. INTRODUCTION

A sound event localization and detection (SELD) task is a dual task where the goal is to classify the type of sounds present in an acoustic scene as well as estimate their direction of arrival (DOA) [1, 2, 3]. It is similar to image classification and image segmentation in the sense that the objective is to identify the content and location of elements in the scene. So SELD can be thought of as a machine listening problem that has attracted significant attention lately, as it is key for artificial intelligence methods to understand the world via sound [4].

In the current state of the art, most SELD tasks are solved using complex systems that are based on deep learning networks [5, 6, 7]. The input features are usually either log mel spectrograms (or other time-frequency transforms), or raw waveforms extracted from a soundfield recorded by a microphone array. The models vary in size and complexity, but include multiple architectures of convolutional networks with additional elements such as recurrent layers [8], transformer-based self-attention modules [9] or dense residual blocks [10]. Post-processing techniques such as weight averaging, ensemble methods and filtering are also common [5].

However, the labeled datasets available to train and evaluate such systems are generally very limited, as this requires labeling for

<sup>\*</sup>Work done during an internship at Sony Group Corporation.



**Fig. 1:** Illustrative comparison of different spatial transformations for a 2D scene. (Top) shows the original omnidirectional response and (top right) the corresponding ACCDOA labels pointing to sources. (Middle left) shows a traditional beamformer and (middle right) the corresponding transformed ACCDOA labels, where high suppression effectively eliminates events outside the main lobe. (Bottom left) shows the proposed spatial mixup with a soft spherical cap and (bottom right) the optionally transformed labels.

both event class and DOA for each frame in the recording. Therefore, a significant component of modern SELD systems is the data augmentation strategy [11]. Although spectrograms can be interpreted as a 2D image representation of sound, and systems with those input features have successfully adapted some image data augmentation techniques (e.g. random cropping, scaling, etc.) [12], the best performing methods use augmentation techniques designed for audio content, such as mixup with equalization EMDA [13] or without equalization [14], spec augment [15], impulse response simulations [5], pitch and/or time shifting and stretching [16, 17], filtering, dynamic range compression [17], or spatial soundfield rotations [18].

Except for soundfield rotations and impulse response simulations, most of the aforementioned augmentation techniques were designed for monophonic or single channel audio, and only a few exploit the spatial characteristics of the input signals. To address this

issue, in this paper we propose Spatial Mixup, which uses a general parametric spatial audio effect as a data augmentation technique by applying a directional loudness modification to the audio data. This modification effectively transforms the spatial characteristics of the recorded soundfield, enhancing sound arriving from some directions while suppressing others. However, unlike a beamformer, the overall transformation is gentle, such that the overall content (class) and DOA of the all recorded events is preserved. Figure 1 shows a visualization of this concept compared with an omnidirectional 2D signal, and a traditional beamformer. The left column shows the omnidirectional responses, and the right column shows the activity-coupled Cartesian DOA (ACCDOA) vectors [19], which assign an event activity to the length of corresponding Cartesian DOA vectors.

## 2. RELATED WORK

Soundfield rotations in ambisonics domain such as swapping, arbitrary rotation, rotation over a single axis were first proposed by [18]. These generate new DOA labels, that might not exist in the dataset. However, the overall soundfield stays constant, where the acoustic environment (i.e the room) is rotated as well, so the relationship between direct sound, early reflections, and reverberation remains unchanged. The relative positions between events is also preserved. Nonetheless, this has proven successful, especially for the localization subtask, at least when there are few overlapping sources.

Expanding the concept of soundfield manipulations, [11] proposed audio channel swapping (ACS) and multi-channel simulation (MCS), where ACS applies a similar soundfield rotation as the channel swapping in [18], to both MIC and FOA signals. MCS aims to simulate new spatial information for specific events. To do this, a preprocessing step analyzes the recording to distinguish between noise and possible sources, then a beamformer extracts the direct sound while the spatial characteristics are extracted computing the spatial covariance matrix. This is comparable to a parametric decomposition into direct and diffuse components [20, 21]. Augmentation occurs by adding random perturbations to the spatial components, preserving the content, simulating new acoustical environments.

Finally, techniques that combine the content of multiple input signals have shown effectiveness too. First, mixup (referred in this paper as regular mixup) is a linear combination of an original signal with some other signal. In the case of time domain audio signals, this is equivalent to mixing two sound tracks together, which translates into two sound events occurring at the same time. In addition, the labels corresponding to the signals can be combined too if available. The regular mixup [14] can be expressed as

$$\hat{\mathbf{x}} = \lambda \mathbf{x} + (1 - \lambda) \mathbf{y}, \quad (1)$$

where  $\mathbf{x}$  is the single channel input audio signal,  $\mathbf{y}$  is the interfering signal,  $\hat{\mathbf{x}}$  is the augmented signal, and  $\lambda \sim \text{Beta}(\alpha, \beta)$  is a hyperparameter that controls the strength of the mixup. An alternative mixup is EMDA [13], applies random equalization to each signal, reducing the overlap in frequency domain, while preserving time mixing.

## 3. PROPOSED METHOD: SPATIAL MIXUP

For this paper, we assume that the input signals represent a soundfield encoded in ambisonics format [22]. This means that a full 3D sound scene has been captured by a microphone array and properly encoded into an orthonormal basis of spherical harmonics representing the full soundfield.

The main idea of spatial mixup is to slightly modify the spatial characteristics of a recorded spatial audio signal to increase the robustness of neural networks models to these transformations. While regular mixup combines the content of two different sound signals, spatial mixup can be understood as applying the mixup operation to a spatially transformed version of the same signal. This operation is now defined as

$$\hat{\mathbf{X}} = \lambda \mathbf{X} + (1 - \lambda) \mathcal{T} \mathbf{X}, \quad (2)$$

where  $\hat{\mathbf{X}} \in \mathbb{R}^{n_{\text{out}}}$  is the augmented audio signal,  $n_{\text{out}} = (N_{\text{out}} + 1)^2$  for the output order  $N_{\text{out}}$ ,  $\mathbf{X} \in \mathbb{R}^{n_{\text{in}}}$  is the multi-channel spatial audio input signal,  $n_{\text{in}} = (N_{\text{in}} + 1)^2$  for the input order  $N_{\text{in}}$ , and  $\mathcal{T} \in \mathbb{R}^{n_{\text{out}} \times n_{\text{in}}}$  is a transformation matrix that performs linear combinations of the channels in  $\mathbf{X}$  via a matrix multiplication. The transformation matrix  $\mathcal{T}$  is further decomposed as

$$\mathcal{T} = \mathbf{Y}_{\text{grid}} \mathbf{G} \mathbf{W}, \quad (3)$$

where  $\mathbf{Y}_{\text{grid}} \in \mathbb{R}^{n_{\text{out}} \times n_{\text{grid}}}$  is a matrix of real spherical harmonics [23] of order  $N_{\text{out}}$  computed at the azimuth and elevation of a discrete set of points sampled from a unitary sphere;  $\mathbf{G} = \text{diag}[g_i] \in \mathbb{R}^{n_{\text{grid}} \times n_{\text{grid}}}$  is a diagonal matrix composed of gains for each point  $i$  in the grid defined for  $\mathbf{Y}_{\text{grid}}$ ; and  $\mathbf{W} \in \mathbb{R}^{n_{\text{grid}} \times n_{\text{in}}}$  is a beamforming matrix that couples the number of input channels to the grid directions. This effectively spatially decomposes the input soundfield into a discrete sampling. Although  $\mathbf{W}$  can be set to any beamforming matrix, in practice it is sufficient to set  $\mathbf{W} = \frac{1}{(N_{\text{in}} + 1)^2} \mathbf{Y}_{\text{grid}}^T$ , which corresponds to a hypercardioid beamforming to the grid points.

The setup presented for  $\mathcal{T}$  can be applied to many spatial transformations depending on the values of  $\mathbf{G}$ , including warping, compression, and acoustic zooms [24, 25, 21]. An additional rotation matrix  $\mathbf{R}$  can be added to Equation (3) if needed. In this paper we focus the analysis on a modification known as directional loudness.

### 3.1. Directional Loudness

A directional loudness modification corresponds to a spatial filter using some function. The spherical cap function divides the area of a unitary sphere into two sections, where the inner section is denoted as a spherical cap. This cap is parametrized by a central angle  $\Omega_c = (\theta, \phi)$  of azimuth and elevation, and a width  $\gamma_c$ . The spherical cap function [25, 24] is defined as the vector of gains

$$g_i = g_1 U(\Omega_c^T \Omega_i - \cos \frac{\gamma_c}{2}) + g_2 U(\cos \frac{\gamma_c}{2} - \Omega_c^T \Omega_i), \quad (4)$$

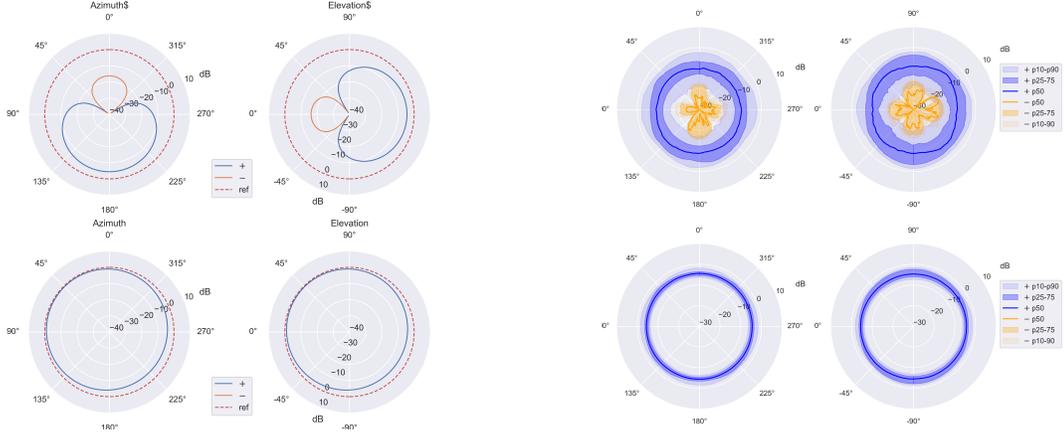
for all  $i \in n_{\text{grid}}$  grid directions in  $\mathbf{G}$ ,  $U$  is the unit step function,  $\gamma_c$  is the width of the cap in radians, and  $g_1, g_2$  are the gains for the region inside and outside the cap respectively.

Moreover, the aforementioned transformation matrix  $\mathcal{T}$  can also be applied to the SELD labels, especially if they are expressed as ACCDOA vectors. In this case, given that the ACCDOA labels are unitary activity vectors pointing to the DOA of a sound events, the application of spatial mixup applies  $\mathcal{T}$  to the labels is expressed as

$$\hat{\mathbf{Z}} = \lambda \mathbf{Z} + (1 - \lambda) \mathcal{T} \mathbf{Y}_0^0, \quad (5)$$

where  $\hat{\mathbf{Z}}$  are the augmented labels,  $\mathbf{Z}$  are the true labels expressed as ACCDOA vectors, and  $\mathbf{Y}_0^0$  are the spherical harmonics of first order and degree, computed for the direction of all active labels.

Generally, any function can be applied to  $\mathbf{G}$ , but we experimentally found that the success of the directional loudness as augmentation method relies heavily on a transformation that is gentle and not too extreme. For this reason, we select two hyperparameters sets



**Fig. 2:** Comparison of cross sections of the polar pattern responses for the omnidirectional channel of hard (top row) and soft (bottom row) spherical caps of first order. (Left column) A typical example directional loudness transform of each spherical cap. (Right) Distribution of the same responses for a sample of 500 patterns with randomized parameters for the spherical caps.

**Table 1:** Hyperparameters for the spherical caps.

$G_{type}$	Parameter	Distribution	Values
Sph-Cap (Soft)	Cap-center-azi	uniform	$[0, \pi]$
	Cap-center-ele	uniform	$[-\pi, \pi]$
	Cap-width	uniform	$[\pi/4, \pi]$
	$G_1$	exponential	$[0, -3]$
	$G_2$	uniform	$[-3, -6]$
Sph-Cap (Hard)	Cap-center-azi	uniform	$[0, \pi]$
	Cap-center-ele	uniform	$[-\pi, \pi]$
	Cap-width	uniform	$[\pi/4, \pi/2]$
	$G_1$	exponential	$[0, -6]$
	$G_2$	uniform	$[-6, -20]$

defined as *soft* and *hard* spherical caps, detailed in Table 1. Figure 2 compares these two types of spherical caps, for both a single example of a typical response as well as the coverage and distribution of a random sample of them. Overall, hard spherical caps cover a wider range and some patterns include negative phase regions, while soft spherical caps are more even, with smaller variations.

## 4. EXPERIMENTS

### 4.1. Experimental Setup

To evaluate the performance of our method, we used the TAU-NIGENS Spatial Sound Events 2021, introduced in the DCASE 2021 Task3 challenge [26]. This dataset includes 600 one minute long clips of spatial audio recordings, presented in two formats: the raw microphone array signals (MIC) and first order ambisonics (FOA), sampled at 24 kHz. The dataset is split into training, validation, and test subsets, consisting of 400, 100 and 100 minutes respectively. Each clip contains a dynamic acoustic scene, where specific sound events, are mixed together in a simulated acoustical environment. In total, there are 12 sound event classes, with examples such as footsteps or a dog barking. Each scene is generated with a collection of up to 3 concurrent sound events, that can be either spatially stationary, or follow a trajectory inside the sound scene. In addition, interference noise along with background noise are also present, where the former are localized sound events that do not belong to any of the classes, and the latter are continuous multi-channel recordings of ambient noise naturally present in the acoustical environments where the impulses were collected.

The evaluation tasks hence consists of the classification of sound events as well as the estimation of DOAs for full clips. For this purpose we used the same classification ( $ER_{LD} \downarrow$ ,  $F_{LD} \uparrow$ ) and localization ( $LE_{CD} \downarrow$ ,  $LR_{CD} \uparrow$ ) as explained in the official DCASE challenge [27]. We also adopted an aggregated SELD error ( $\mathcal{E}_{SELD} \downarrow$ ), as

$$\mathcal{E}_{SELD} = \frac{ER_{LD} + (1 - F_{LD}) + \frac{LE_{CD}}{\pi} + (1 - LR_{CD})}{4}. \quad (6)$$

The goal of the experiments was to compare the impact of the augmentation methods fairly, and not necessarily to get the best possible performance in the task. Consequently, the training setup was the same for all. All experiments were trained for 100,000 iterations of batch size 32, with validation every 10,000 steps, using the official subset split. We minimize the MSE loss, using Adam optimizer and a learning rate scheduler with a warmup stage starting at learning rate  $1e-4$ , reaching  $1e-3$  after 5 validation steps, followed by a reduce on plateau scheduler (monitoring the validation SELD error) with patience of 3 validation steps and a decay rate of 0.9. For each experiment, we report the test subset results, from the model of the best validation step of multiple runs.

### 4.2. Features and models

The experiments are conducted using two different systems, a low complexity model with standard architecture, and another with a more sophisticated model with a large amount of parameters. The systems are:

1. **Basic system** - We use the CRNN10 model proposed by [28], which consists of 2d convolutional layers with batch normalization and increasing number of channels. The inputs are linear amplitude STFT and interchannel differences using only the FOA input signals, for a total of 7 input channels. The spectrograms are computed using frame size of 512, hop size of 240, and total input length for the network is 1.27 seconds.
2. **Sophisticated system** - We use the RD3Net [19], which consists of a series of densely connected blocks of 2d, dilated convolutions. Each block is followed by a down sampling module and finally, a gated recurrent unit (GRU), fully-connect (FC) layer, and up sampling operation as outputs. The input features are the same as in the basic system.

**Table 2:** Performance of Spatial Mixup with different directional loudness  $\mathbf{G}$  matrix types in the *basic* system.

System	ER <sub>LD</sub> ↓	F <sub>LD</sub> ↑	LE <sub>CD</sub> ↓	LR <sub>CD</sub> ↑	$\mathcal{E}_{\text{SELD}}$ ↓
Baseline	0.689	40.5	20.7	44.4	0.489
Random	0.776	24.2	26.9	32.5	0.590
Identity	0.668	<b>42.2</b>	19.5	42.9	0.481
Sph-Cap (Hard)	0.693	39.1	22.1	<b>45.6</b>	0.492
Sph-Cap (Soft)	<b>0.664</b>	42.1	<b>19.4</b>	43.2	<b>0.480</b>

**Table 3:** Performance of Spatial Mixup with different directional loudness  $\mathbf{G}$  matrix types in the larger *sophisticated* system.

System	ER <sub>LD</sub> ↓	F <sub>LD</sub> ↑	LE <sub>CD</sub> ↓	LR <sub>CD</sub> ↑	$\mathcal{E}_{\text{SELD}}$ ↓
Baseline	0.678	43.5	21.8	53.5	0.458
Random	0.744	27.0	29.3	41.8	0.555
Identity	0.643	46.9	22.1	<b>56.0</b>	0.434
Sph-Cap (Hard)	0.660	44.7	22.1	55.7	0.445
Sph-Cap (Soft)	<b>0.615</b>	<b>48.9</b>	<b>18.7</b>	54.2	<b>0.422</b>

### 4.3. Results

#### 4.3.1. Effects of Directional Gain $\mathbf{G}$

Table 2 shows the performance of the proposed augmentation method with different types of  $\mathbf{G}$  compared with a baseline without any data augmentation for the basic model. In these experiments, the spatial mixup was applied only to the input signals  $\mathbf{X}$ , the output order was 1, and the  $\mathbf{Y}_{\text{grid}}$  was computed for a t-design [29] of degree 3 (giving 6 total directions). For most metrics, the performance is better for all  $\mathbf{G}$  types except random (uniformly random diagonal matrix), which is significantly worse. Surprisingly, the identity matrix and the soft spherical caps show similar results, while the hard spherical caps are not as good. The former can be explained in part because the full transformation matrix is not truly orthogonal, due to the limited spatial resolution in the grid that generates some small non-zero values. The latter is most likely because the hard spherical caps can sometimes generate extreme patterns that resemble a beamformer, rather than a gentle manipulation of the soundfield. These patterns are possibly too extreme for the model to learn invariance.

The same comparison of different  $\mathbf{G}$  types was explored for the sophisticated model in Table 3. Here, the spatial mixup setup was the same as the previous experiments, except that the  $\mathbf{Y}_{\text{grid}}$  was computed with a larger t-design of degree 7, for 24 directions, giving better spatial resolution. The results show a similar trend as Table 3, where the main difference is that the  $\mathbf{G}$  type identity was not quite as good as the soft spherical caps. This suggests that the small non zero values are not as strong here, due to the larger number of directions in the grid. In addition, the higher capacity of the model is able to accommodate for the soft spherical caps, learning better invariance. In summary, a smooth  $\mathbf{G}$  function, with a sufficiently large grid works best.

#### 4.3.2. Other effects

The general method allows to utilize a different order for the outputs than the input signals. However, even if a higher output order is selected, it is not possible to generate spatial information that is not already there. That said, for data augmentation purposes, having a higher order might enable a different architecture for neural network models, (e.g. with more input channels), as well as more nuanced directional loudness modifications. In addition, the labels can also be augmented as described in Equation (5). Cursory experimental

**Table 4:** Performance of common data augmentation method compared with the Spatial Mixup using the *basic* system.

System	ER <sub>LD</sub> ↓	F <sub>LD</sub> ↑	LE <sub>CD</sub> ↓	LR <sub>CD</sub> ↑	$\mathcal{E}_{\text{SELD}}$ ↓
Baseline (B1)	0.689	40.5	20.7	44.4	0.489
B1+Mixing	0.649	45.7	20.4	51.9	0.447
B1+Rotation	0.633	<b>46.5</b>	20.4	51.1	0.442
B1+SpecAugment	0.702	37.6	23.4	45.2	0.501
B1+EQ	0.675	42.5	20.9	44.6	0.480
B1+All	0.652	46.2	22.4	<b>57.3</b>	<b>0.435</b>
B1+Sph-cap (Soft)	0.662	42.8	<b>20.1</b>	45.7	0.472
B1+All+Sph-cap (Soft)	<b>0.628</b>	46.3	<b>20.1</b>	50.8	0.442

results of both effects showed very small differences to results already presented, so these are not explicitly included in this paper. Nevertheless, it is possible that other tasks, in particular tasks with high order ambisonics data might see more significant performance gains when using spatial mixup.

#### 4.3.3. Comparison to other data augmentation

Table 4 compares spatial mixup to other common data augmentation methods including mixing [19], spec augment [15], FOA soundfield rotations [18], random equalization [13] and a combination of all four. These results show that from the common augmentations, FOA rotations and mixing are the best performers, achieving a faster convergence and better metrics in general. While random EQ increases the performance slightly, it is surprising that spec augment is markedly worse than the baseline. More importantly, spatial mixup with soft spherical caps (using a grid with t-design of degree 7) shows considerable performance gains over the baseline, better than EQ but slightly less than mixing. Lastly, when comparing the combinations with and without spatial mixup, it seems that adding spherical caps reduces errors, but also reduces recall for events. However, larger systems might exploit this better.

A possible explanation for these results is that the DCASE 2021 Task3 benefits the most from methods that improve equivariance rather than invariance, given that both FOA rotations and mixing increase the coverage of the labels. In contrast, spatial mixup modifies the relative levels of certain directions, increasing sound level diversity for all events. Nonetheless, spatial mixup shows good results.

## 5. CONCLUSIONS

In this paper, we proposed a data augmentation method for sound event localization and detection (SELD) tasks, based on the application of spatial audio parametric effects, in a process we call Spatial Mixup. This enables modifications to the spatial characteristics of audio signals encoded in ambisonics format, by applying a transformation matrix to the time domain input signals. This matrix is obtained by spatially sampling the original soundfield, and applying a directional loudness gain modification. The method was evaluated in the DCASE2021 Task3 dataset, which includes complex sound scenes with overlapping and non-stationary sources, as well as interference and background noise. The method proved effective when using gentle modifications known as soft spherical caps. It improves all metrics when compared to a non-augmented baseline, and shows similar advantages compared to well known augmentation methods. Future research could explore the application of the method for other machine learning tasks with spatial audio, as well as analyze further transforms such as audio warping or acoustic zoom.

## 6. REFERENCES

- [1] Hendrik Purwins, Bo Li, Tuomas Virtanen, Jan Schlüter, Shuo-Yiin Chang, and Tara Sainath, “Deep learning for audio signal processing,” *IEEE Journal of Selected Topics in Signal Processing*, 2019.
- [2] Sharath Adavanne, Archontis Politis, Joonas Nikunen, and Tuomas Virtanen, “Sound event localization and detection of overlapping sources using convolutional recurrent neural networks,” *IEEE Journal of Selected Topics in Signal Processing*, 2018.
- [3] Sharath Adavanne, Archontis Politis, and Tuomas Virtanen, “Direction of arrival estimation for multiple sound sources using convolutional recurrent neural network,” in *Proc. of EUSIPCO*, 2018.
- [4] Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicens Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman, “Soundspaces: Audio-visual navigation in 3D environments,” in *Proc. of ECCV*, 2020.
- [5] Kazuki Shimada, Naoya Takahashi, Yuichiro Koyama, Shusuke Takahashi, Emiru Tsunoo, Masafumi Takahashi, and Yuki Mitsufuji, “Ensemble of ACCDOA- and EINV2-based systems with D3Nets and impulse response simulation for sound event localization and detection,” in *Tech. report of DCASE Challenge*, 2021.
- [6] Thi Ngoc Tho Nguyen, Karn Watcharasupat, Ngoc Khanh Nguyen, Douglas L. Jones, and Woon Seng Gan, “DCASE 2021 Task 3: Spectrotemporally-aligned features for polyphonic sound event localization and detection,” in *Tech. report of DCASE Challenge*, 2021.
- [7] Sang-Hoon Lee, Jung-Wook Hwang, Sang-Buem Seo, and Hyung-Min Park, “Sound event localization and detection using cross-modal attention and parameter sharing for DCASE2021 challenge,” in *Tech. report of DCASE Challenge*, 2021.
- [8] Emre Cakir, Sharath Adavanne, Giambattista Parascandolo, Konstantinos Drossos, and Tuomas Virtanen, “Convolutional recurrent neural networks for bird audio detection,” in *Proc. of EUSIPCO*, 2017.
- [9] Koichi Miyazaki, Tatsuya Komatsu, Tomoki Hayashi, Shinji Watanabe, Tomoki Toda, and Kazuya Takeda, “Weakly-supervised sound event detection with self-attention,” in *Proc. of IEEE ICASSP*, 2020.
- [10] Naoya Takahashi and Yuki Mitsufuji, “Densely connected multi-dilated convolutional networks for dense prediction tasks,” in *Proc. of IEEE/CVF CVPR*, 2021.
- [11] Qing Wang, Jun Du, Hua-Xin Wu, Jia Pan, Feng Ma, and Chin-Hui Lee, “A four-stage data augmentation approach to ResNet-Conformer based acoustic modeling for sound event localization and detection,” *arXiv:2101.02919*, 2021.
- [12] Dan Hendrycks and Thomas Dietterich, “Benchmarking neural network robustness to common corruptions and perturbations,” in *Proc. of ICLR*, 2018.
- [13] Naoya Takahashi, Michael Gygli, and Luc Van Gool, “AENet: learning deep audio features for video analysis,” *IEEE Trans. on Multimedia*, 2017.
- [14] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz, “mixup: Beyond empirical risk minimization,” in *Proc. of ICLR*, 2018.
- [15] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” in *Proc. of INTERSPEECH*, 2019.
- [16] Jan Schlüter and Thomas Grill, “Exploring data augmentation for improved singing voice detection with neural networks,” in *Proc. of ISMIR*, 2015.
- [17] Justin Salamon and Juan Pablo Bello, “Deep convolutional neural networks and data augmentation for environmental sound classification,” *IEEE Signal Processing Letters*, 2017.
- [18] Luca Mazzon, Yuma Koizumi, Masahiro Yasuda, and Noboru Harada, “First order Ambisonics domain spatial augmentation for DNN-based direction of arrival estimation,” in *Proc. of DCASE Workshop*, 2019.
- [19] Kazuki Shimada, Yuichiro Koyama, Naoya Takahashi, Shusuke Takahashi, and Yuki Mitsufuji, “ACCDOA: Activity-coupled Cartesian direction of arrival representation for sound event localization and detection,” in *Proc. of IEEE ICASSP*, 2021, pp. 915–919.
- [20] Leo McCormack, Archontis Politis, and Ville Pulkki, “Parametric spatial audio effects based on the multi-directional decomposition of Ambisonic sound scenes,” in *Proc. of DAFX*, 2021.
- [21] Archontis Politis, Tapani Pihlajamäki, and Ville Pulkki, “Parametric spatial audio effects,” in *Proc. of DAFX*, 2012.
- [22] Hannes Pomberger, Franz Zotter, and A Sontacchi, “An Ambisonics format for flexible playback layouts,” in *Proc. of Ambisonics Symposium*, 2009.
- [23] Boaz Rafaely, *Fundamentals of Spherical Array Processing*, Springer, 2015.
- [24] Franz Zotter and Matthias Frank, *Ambisonics: A Practical 3D Audio Theory for Recording, Studio Production, Sound Reinforcement*, Springer, 2019.
- [25] Matthias Kronlachner and Franz Zotter, “Spatial transformations for the enhancement of Ambisonic recordings,” in *Proc. of International Conference on Spatial Audio*, 2014.
- [26] Archontis Politis, Sharath Adavanne, Daniel Krause, Antoine Deleforge, Prerak Srivastava, and Tuomas Virtanen, “A dataset of dynamic reverberant sound scenes with directional interferers for sound event localization and detection,” *arXiv:2106.06999*, 2021.
- [27] Archontis Politis, Annamaria Mesaros, Sharath Adavanne, Toni Heittola, and Tuomas Virtanen, “Overview and evaluation of sound event localization and detection in DCASE 2019,” *IEEE/ACM Trans. on ASLP*, vol. 29, pp. 684–698, 2020.
- [28] Yin Cao, Qiuqiang Kong, Turab Iqbal, Fengyan An, Wenwu Wang, and Mark D Plumbley, “Polyphonic sound event detection and localization using a two-stage strategy,” in *Proc. of DCASE Workshop*, 2019.
- [29] Manuel Gräf and Daniel Potts, “On the computation of spherical designs by a new optimization approach based on fast spherical Fourier transforms,” *Numerische Mathematik*, vol. 119, 2011.