



Research Paper

Speech signal enhancement in cocktail party scenarios by deep learning based virtual sensing of head-mounted microphones

Tim Fischer^{a,b}, Marco Caversaccio^{a,b}, Wilhelm Wimmer^{a,b,*}^a Hearing Research Laboratory, ARTORG Center for Biomedical Engineering Research, University of Bern, Bern 3008, Switzerland^b Department of ENT, Head and Neck Surgery, Inselspital, Bern University Hospital, University of Bern, Bern 3008, Switzerland

ARTICLE INFO

Article history:

Received 8 February 2021

Revised 31 May 2021

Accepted 7 June 2021

Available online 17 June 2021

Keywords:

Artificial intelligence

Selective hearing

Neural network

Beamformer

Hearing aid

Cochlear implant

ABSTRACT

The cocktail party effect refers to the human sense of hearing's ability to pay attention to a single conversation while filtering out all other background noise. To mimic this human hearing ability for people with hearing loss, scientists integrate beamforming algorithms into the signal processing path of hearing aids or implants' audio processors. Although these algorithms' performance strongly depends on the number and spatial arrangement of the microphones, most devices are equipped with a small number of microphones mounted close to each other on the audio processor housing. We measured and evaluated the impact of the number and spatial arrangement of hearing aid or head-mounted microphones on the performance of the established Minimum Variance Distortionless Response beamformer in cocktail party scenarios. The measurements revealed that the optimal microphone placement exploits monaural cues (pinna-effect), is close to the target signal, and creates a large distance spread due to its spatial arrangement. However, this microphone placement is impractical for hearing aid or implant users, as it includes microphone positions such as on the forehead. To overcome microphones' placement at impractical positions, we propose a deep virtual sensing estimation of the corresponding audio signals. The results of objective measures and a subjective listening test with 20 participants showed that the virtually sensed microphone signals significantly improved the speech quality, especially in cocktail party scenarios with low signal-to-noise ratios. Subjective speech quality was assessed using a 3-alternative forced choice procedure to determine which of the presented speech mixtures was most pleasant to understand. Hearing aid and cochlear implant (CI) users might benefit from the presented approach using virtually sensed microphone signals, especially in noisy environments.

© 2021 The Author(s). Published by Elsevier B.V.

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

1. Introduction

Following a conversation in a noisy setting is difficult. In literature, this phenomenon is referred to as the cocktail-party problem. It describes an acoustic scenario, where multiple speech and noise sources with different intensities and directions of incidence overlap (Cherry, 1953). For normal-hearing persons, the auditory system can handle conflicting sounds and focus on a specific conversation (Mesgarani and Chang, 2012; Peissig and Kollmeier, 1997). In hearing aids or CI audio processors, this separation of the conversational partner from a noise tangle is the goal of sophisticated

beamforming algorithms (Makino et al., 2007; Qian et al., 2018; Smaragdis, 1998; Veen and Buckley, 1988).

It is well known that the signal quality of beamforming algorithms increases with the number of available input microphones and their positioning with respect to the target source (Feng et al., 2011; Habets et al., 2009; Souden et al., 2009; Wimmer et al., 2015; 2017; Wouters and Berghe, 2001). Using numerical experiments, Feng et al. (2011) showed that the microphone positions play an essential role in the overall performance of beamforming algorithms. Jones et al. (2016) further showed for CI users that the microphone position at the ear canal versus behind the ear led to more detailed interaural level difference (ILD) information due to the frequency transformations of the pinna (Blauert, 1997; Fischer et al., 2021). In the specific case of unilateral CI users, it was demonstrated that an additional microphone positioned at the contralateral ear led to increased speech understand-

* Corresponding author at: Hearing Research Laboratory, ARTORG Center for Biomedical Engineering Research, University of Bern, Bern 3008, Switzerland.

E-mail address: wilhelm.wimmer@artorg.unibe.ch (W. Wimmer).

ing in noise (Arora et al., 2013; Dorman et al., 2018; Wimmer et al., 2017).

Since many conversations are held face to face (Wu and He, 2018), it is reasonable to assume that additional microphones in positions other than the contralateral ear canal, e.g., on the forehead, may further improve speech understanding. However, the additional placement of microphones on the head is impractical from the perspective of a hearing aid or CI user. One way of circumventing this limitation may be to place the microphones virtually rather than physically. The results of several virtual microphone sensing approaches suggest that estimating an additional microphone signal using information from the available microphones may improve the speech quality in a cocktail party scenario (Jinzai et al., 2019; Katahira et al., 2016; Yamaoka et al., 2019). The microphone array used to record the reference signals was similar in the studies and consisted of 2 microphones positioned in a straight line at a distance of 4 cm (Katahira et al., 2016; Yamaoka et al., 2019) or 3 cm (Jinzai et al., 2019) from each other. To generate virtual microphone signals, the phase was linearly interpolated (Katahira et al., 2016; Yamaoka et al., 2019) or extrapolated (Jinzai et al., 2019) using measurements of the real microphone signals. In Denk et al. (2018), functions transformed the sound pressure at a microphone positioned on a hearing aid to the pressure measured at the open eardrum. The basis for the determination of these functions were the relative transfer functions (RTFs) between the microphones, which in turn were determined by head related transfer functions (HRTFs) measurements using frequency sweeps in an anechoic chamber. Also using frequency sweeps, Corey et al. (2019) measured and evaluated impulse responses of 160 microphones spread across the body and affixed to wearable accessories. Their results suggest that microphone arrangements with large spatial distance spread across the body provided the best signal-to-noise ratio (SNR) values. Unlike microphones positioned on the head, the geometric arrangement of microphones placed on clothing may change according to posture. Likely, the quality of a beamforming algorithm defined for a specific microphone geometry suffers from the continually changing microphone geometries in everyday life (Himawan et al., 2008).

The tremendous progress in the field of machine learning leads to the expectation that in the future, the RTFs between microphones can be determined purely data-driven, i.e., without prior knowledge of the specific measurement setup. As a result, beamforming algorithms could be tuned to individual array geometries by simply providing sufficient reference data from the wearer without the need for anechoic chambers or knowledge of the sound sources' positions. In the Mic2Mic publication (Mathur et al., 2019) it was demonstrated that even with unlabeled and unpaired data, audio signals between different microphone domains could be translated. Based on the results, an additional virtual microphone at the head of a hearing aid or CI user generated or learned solely by data-driven rules seems like a realistic scenario. However, regardless of whether the microphones are placed virtually or physically on a subject's head, little is known about how their positioning affects beamforming.

To continue the discussion, the first objective of this work was to systematically investigate the speech signal quality in complex acoustic scenarios with varying head-mounted microphone arrangements and a minimum variance distortionless response (MVDR) beamformer as introduced by Souden et al. (2009). Based on these measurements, virtual microphone signals at specific positions were estimated using a deep neural network. Finally, subjective listening tests were conducted to investigate to what extent the virtually sensed microphone signals could improve the speech signal quality.

2. Methods

2.1. Linear observation model

In this work, recordings from $M = 16$ microphones attached to a human head were used. Each of the $i = 1 \dots M$ microphone signals $y_i(t)$ recorded varying acoustic cocktail party scenarios at time t . In the following, the cocktail party mixtures are described as the summation of the target speech source $s_i(t)$ and the noise $w_i(t)$ at microphone i :

$$y_i(t) = a_i s(t - \tau_i) + w_i(t)$$

where τ_i represents the time-delay of arrival and a_i is the amplitude modulation depending on the geometric arrangement of the microphones under the assumption of anechoic conditions. The noise $w_i(t)$ is assumed to be uncorrelated with the signal $s_i(t)$.

To enhance the perception of the target speech sources, the signals at each microphone can be combined using "beamforming" techniques. In this study, we used the widely studied MVDR beamformer (Capon, 1969; Habets et al., 2010), which is introduced in the following section.

2.2. MVDR beamforming

The MVDR beamformer minimizes the power of the beamformed signal while preserving the target signal, under the constraint of no distortion in the target signal (Souden et al., 2009). The MVDR is a filter-and-sum beamformer and as such it applies different phase weights $h_i(f)$ to the i input microphone channels in order to steer the main lobe of the directivity pattern to the direction of the target signal. The phase weights, or filters, are obtained in the frequency domain using (Erdogan et al., 2016):

$$\mathbf{h}_{ref}(f) = [h_{1,ref}(f), \dots, h_{M,ref}(f)]^T = \frac{1}{\lambda(f)} (\mathbf{G}(f) - \mathbf{I}_{M \times M}) \mathbf{e}_{ref} \quad (1)$$

Where \mathbf{I} is the identity matrix and $\mathbf{G}(f)$ can be obtained by $\mathbf{G}(f) = \Phi_{noise}^{-1}(f) \Phi_{obs}(f)$ with $\lambda(f) = \text{trace}(\mathbf{G}(f)) - M$ (Benesty et al., 2008; Souden et al., 2009). The spatial covariance matrices Φ can be computed by using time-frequency masks (Drude and Haeb-Umbach, 2017; Erdogan et al., 2016; Higuchi et al., 2017; Ito et al., 2016). However, in this work we focus on the impact of additional microphone channels on the MVDR beamformers performance and extract $\Phi_{noise}^{-1}(f)$, $\Phi_{obs}(f)$ and $\Phi_{target}(f)$ from the noise, observation and target recordings.

The standard unit vector of the reference microphone \mathbf{e}_{ref} , is selected by a maximum a posteriori expected SNR estimation. The reference microphone is chosen based on $ref = \underset{r}{\text{argmax}} \text{SNR}_{post,r}$ (Erdogan et al., 2016) and:

$$\text{SNR}_{post,r} = \frac{\sum_{f=0}^{F-1} \mathbf{h}_r^H(f) \Phi_{target}(f) \mathbf{h}_r(f)}{\sum_{f=0}^{F-1} \mathbf{h}_r^H(f) \Phi_{noise}(f) \mathbf{h}_r(f)}$$

Thus, the reference channel or microphone depends on $\mathbf{h}_r(f)$, which is the M -dimensional filter response (see Eq. (1)) at the discrete frequency index $f = 0, \dots, F - 1$, when \mathbf{e}_{ref} is set to \mathbf{e}_r . After the filters $\mathbf{h}_{ref}(f)$ are computed, the beamformed output $z_{t,f}$ is obtained by using the short-time Fourier transforms (STFTs) $y_{i,t,f}$ of the microphone signals $y_i(t)$:

$$z_{t,f} = \sum_{i=1}^M h_{i,ref}(f) y_{i,t,f}$$

For the MVDR beamformer, the input signals were down-sampled to 8kHz and a Blackman window was applied (Blackman and Tukey, 1959). Subsequently, an STFT (size =

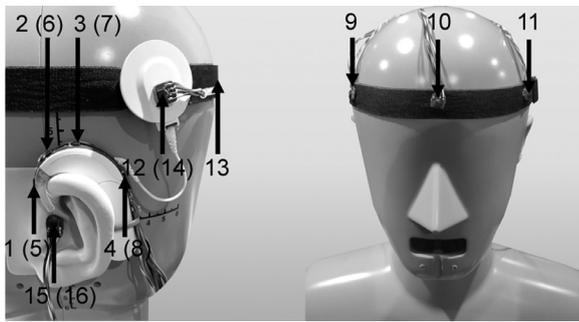


Fig. 1. Placement of the 16 microphones used for cocktail party scenario recordings. The IDs refer to the microphone signals assignment in the multi-channel recording audio files (Fischer et al., 2020b). Numbers in brackets refer to the contralateral (here: right side) assignment of the microphones. The sagittal plane is defined by a straight line between microphones 10 and 13 (front and back). A numeric description can be found in Table 1.

1	0	14	25	39	166	168	170	176	165	127	69	69	145	181	36	163
2	14	0	13	32	166	168	169	174	166	131	71	56	135	176	46	165
3	25	13	0	22	169	170	170	174	172	141	82	44	127	174	49	167
4	39	32	22	0	171	173	172	173	184	159	103	40	120	174	46	167
5	166	166	169	171	0	15	24	39	82	138	167	183	135	69	163	28
6	168	168	170	173	15	0	12	34	80	138	167	182	129	58	168	40
7	170	169	170	172	24	12	0	25	88	146	172	179	120	46	171	45
8	176	174	174	173	39	34	25	0	112	167	186	179	110	38	173	46
9	165	166	172	184	82	80	88	112	0	75	132	194	178	122	176	104
10	127	131	141	159	138	138	146	167	75	0	73	173	199	175	147	151
11	69	71	82	103	167	167	172	186	132	73	0	116	177	189	98	174
12	69	56	44	40	183	182	179	179	194	173	116	0	106	172	85	183
13	145	135	127	120	135	129	120	110	178	199	177	106	0	91	148	136
14	181	176	174	174	69	58	46	38	122	175	189	172	91	0	183	82
15	36	46	49	46	163	168	171	173	176	147	98	85	148	183	0	155
16	163	165	167	167	28	40	45	46	104	151	174	183	136	82	155	0

Fig. 2. Euclidean distances in millimeters between the microphones for the head and torso simulator measurements (Fischer et al., 2020b).

256 and shift = 128) was performed. To reconstruct the signal, an inverse short-time Fourier transform (ISTFT) with the overlap-add strategy was applied. The herein used MVDR beamformer to evaluate the benefits of virtual microphone signals is just one application scenario. Theoretically, any multi-channel speech-enhancement algorithm could have been used to assess the benefits of virtually sensed microphone signals.

2.3. Data

The Bern cocktail party (BCP) dataset is tailored to this work, as it contains multi-microphone recordings of hearing aid or CI users in cocktail party scenarios (Fischer et al., 2020b). For the recordings, 12 loudspeakers (Control 1 Pro, JBL, Northridge, USA) were aligned horizontally in a circle at the height of the ears (1.2 m) in an acoustic chamber (Fischer et al., 2020a; 2020d; Wimmer et al., 2017). For this work, we used the acoustic scenarios captured with 16 microphones (ICS-40619, TDK, Tokyo, Japan) attached to a head and torso simulator (Brüel & Kjaer, Type 4128, Nærum, Denmark) (see Figs. 1 and 2).

2.3.1. Test dataset

The results of this work were computed with an excerpt of 2400 samples from the BCP dataset (Fischer et al., 2020b). The duration of each sample was 1.5s, resulting in a total test dataset duration of 1h. The samples were randomly chosen under the constraint, that a majority of the recordings contain a target source

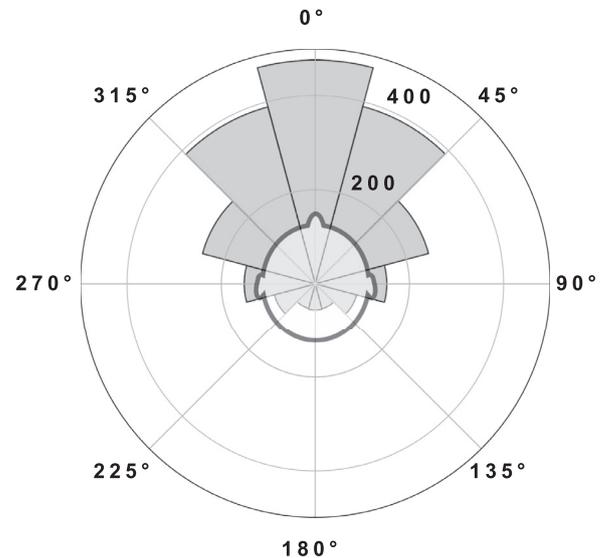


Fig. 3. Circular histogram of the frequency of occurrence of spatial source directions in relation to the head and torso simulator azimuth. The audio files were selected such that the directional distribution assumes a von-Mises distribution with $\mu = 0.0$ and $\kappa = 1.1$ (Fischer et al., 2020b).

azimuth inside the field of view (i.e., $\pm 45^\circ$), as this represents the most natural listening scenario (Wu et al., 2018) (see Fig. 3). All samples were randomly selected from an SNR distribution which covered conversational speech levels with 1 to 3 competing speakers and varying background noise types and intensities. The distribution of the audio mixture on the 12 output channels covered scenarios of spatially separated and non-separated speech and noise sources. The samples or audio mixtures had a mean SNR value of 1.2 dB with a standard deviation of 10.9 dB.

2.3.2. Training dataset

For the training and validation of the deep neural network 65h (78404 audio samples with 3s duration each) were randomly selected from the head and torso simulator recordings of the BCP dataset (Fischer et al., 2020b), excluding the test dataset (see Section 2.3.1). Ninety percent of the samples were used for training and 10% for validation. Because of the large size of the training and validation dataset, no cross-validation was performed.

2.4. Evaluation of microphone channel configurations

Various microphone channel configurations were evaluated by adding or omitting microphone channels with respect to a reference microphone channel configuration, as explained in detail later (Section 3, Tables 3–6). The results were computed by providing the MVDR beamformer (Souden et al., 2009) with the target and noise spatial covariance matrices Φ of the audio mixtures from the corresponding microphone configurations.

The reference microphone configurations were selected to cover reasonable microphone inputs of hearing aid devices or audio processors. Care was also taken to ensure that all microphones in the unilateral reference microphone configurations could technically be connected to the audio processor using an existing cable such as from the CI transmission coil to the audio processor.

To cover realistic use cases regarding the benefits of different microphone configurations, the results were divided into 4 categories rather than presenting all possible microphone channel combinations: subsets of unilateral CI microphone configurations (see Table 3), unilateral CI microphone configurations with additional ipsilateral microphones (Table 4), unilateral CI microphone

Table 1
Assignment of the 16 microphone positions to their respective IDs.

Microphone ID	Microphone position
{1}	Left audio processor. Facing forward.
{2}	Left audio processor. Facing to the top / forward.
{3}	Left audio processor. Facing to the top / backward.
{4}	Left audio processor. Facing back.
{5}	Right audio processor. Facing forward.
{6}	Right audio processor. Facing to the top / forward.
{7}	Right audio processor. Facing to the top / backward.
{8}	Right audio processor. Facing backward.
{9}	Right temple.
{10}	Front.
{11}	Left temple.
{12}	Left transmission coil.
{13}	Back.
{14}	Right transmission coil.
{15}	Left Ear. Entry of the ear canal.
{16}	Right Ear. Entry of the ear canal.

Table 2
Overview of all measured microphone configurations.

Unilateral microphone configurations	Bilateral microphone configurations
{1}	{1, 2, 3, 4, 9}
{2}	{1, 2, 3, 4, 14}
{3}	{1, 2, 3, 4, 16}
{4}	{1, 2, 3, 4, 5, 6, 7, 8}
{10}	{1, 2, 3, 4, 5, 6, 7, 8, 10}
{11}	{1, 2, 3, 4, 5, 6, 7, 8, 13}
{12}	{1, 2, 3, 4, 5, 6, 7, 8, 9, 11}
{13}	{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16}
{15}	{1, 2, 3, 4, 5, 6, 7, 8, 15, 16}
{1, 2}	{2, 3, 9}
{1, 2, 3, 4}	{2, 3, 14}
{1, 2, 3, 4, 10}	{2, 3, 16}
{1, 2, 3, 4, 11}	{2, 3, 6, 7}
{1, 2, 3, 4, 12}	{2, 3, 6, 7, 10}
{1, 2, 3, 4, 13}	{2, 3, 6, 7, 13}
{1, 2, 3, 4, 15}	{2, 3, 6, 7, 9, 11}
{1, 3}	{2, 3, 6, 7, 15, 16}
{1, 4}	{2, 3, 10, 13, 16}
{2, 3}	
{2, 3, 10}	
{2, 3, 11}	
{2, 3, 12}	
{2, 3, 13}	
{2, 3, 15}	
{2, 4}	
{3, 4}	

configurations with additional contralateral microphones (Table 5), symmetric bilateral CI configurations with additional microphones (Table 6). An overview of all measured microphone configurations can be found in Table 2.

For the evaluation of the microphone configurations (i.e., real recordings and virtually sensed microphone channels), the following objective speech quality metrics were assessed: perceptual evaluation of speech quality (PESQ) (Rix et al., 2001), short-time objective intelligibility (STOI) (Taal et al., 2011) and scale-invariant speech to distortion ratio (SI-SDR) (Roux et al., 2019). The PESQ metric models the speech quality as perceived by human listeners. Analysis of speech-audio with the PESQ metric usually ranges from 1.0 (high distortion) to 4.5 (no distortion) (Rix et al., 2001). The values of STOI range from 0.0 (no word correctly understood) to 1.0 (all words correctly understood) and highly correlate with the intelligibility of degraded speech signals (Taal et al., 2011). The SI-SDR metric defines the energy ratio between the clean target signal and the acoustic distortions in decibel (dB). It is a slightly modified version of speech to distortion ratio (SDR), making it in-

Table 3
Values represent the mean difference in the performance of the unilateral cochlear implant (CI) microphone configurations compared to the mean performance of the reference channel configuration including channels positioned on the sagittal plane (see Fig. 1). The best result for each metric is marked in bold. All performance differences were statistically significant compared to the reference channel performance, except those marked with “†”.

Microphone IDs	Metric		
	PESQ	STOI	SI-SDR
Ref.: {1, 2, 3, 4}	1.77	0.48	-29.07
{1}	-0.28	-0.06	-2.95
{2}	-0.28	-0.06	-3.13
{3}	-0.29	-0.06	-3.13
{4}	-0.31	-0.07	-3.32
{10}	-0.24	-0.03	-2.77
{11}	-0.25	-0.04	-2.65
{12}	-0.30	-0.07	-3.24
{13}	-0.35	-0.08	-3.52
{15}	-0.29	-0.06	-3.19
{1, 2}	-0.17	-0.03	-1.25
{3, 4}	-0.13	-0.02	-0.86
{1, 3}	-0.15	-0.03	-0.97
{1, 4}	-0.08	-0.01†	-0.77
{2, 3}	-0.16	-0.03	-1.32
{2, 4}	-0.09	-0.01†	-0.89

Table 4
Values represent the mean difference in the performance of unilateral cochlear implant (CI) microphone configurations when additional ipsilateral, including sagittal plane, microphones were added (see Fig. 1). The performance difference is calculated in relation to the mean performance of the reference channel configuration. The best result for each metric is marked in bold. All performance differences were statistically significant compared to the reference channel performance, except those marked with “†”.

Microphone IDs	Metric		
	PESQ	STOI	SI-SDR
Ref.: {1, 2, 3, 4}	1.77	0.48	-29.07
Ref. + {10}	0.18	0.04	0.69
Ref. + {11}	0.20	0.04	0.59
Ref. + {12}	0.02	<0.01	0.14†
Ref. + {13}	0.11	0.03	0.64
Ref. + {15}	0.01	< 0.01†	-0.39†
Ref.: {2, 3}	1.61	0.45	-30.38
Ref. + {10}	0.22	0.06	1.38
Ref. + {11}	0.23	0.06	1.10
Ref. + {12}	0.12	0.03	0.81
Ref. + {13}	0.15	0.04	0.92
Ref. + {15}	0.03	<0.01	0.30

sensitive to power rescaling of the estimated signal (Roux et al., 2019).

For testing within a group of microphone configurations, the Friedman test was used (see Sections 3.1 and 3.2). To find the configurations that differed significantly after the Friedman test has rejected the null hypothesis, a post-hoc Nemenyi test was performed. In Section 3.3, two sets of paired samples were compared to each other with the two-sided Wilcoxon signed-rank test (no multiple testing). The significance level was chosen with $\alpha = 0.05$ for all statistical tests.

2.5. Virtual sensing of a microphone channel

The virtual sensing approach aimed to improve the speech quality in cocktail party scenarios by providing the beamformer

Table 5

Values represent the mean difference in the performance of unilateral cochlear implant (CI) microphone configurations when additional contralateral microphones were added (see Fig. 1). The performance difference is calculated in relation to the mean performance of the *reference channel configuration*. The best result for each metric is marked in bold. All performance differences were statistically significant compared to the reference channel performance.

Microphone IDs	Metric		
	PESQ	STOI	SI-SDR
Ref.: {1, 2, 3, 4}	1.77	0.48	-29.07
Ref. + {9}	0.12	0.03	0.41
Ref. + {14}	0.16	0.03	0.80
Ref. + {16}	0.19	0.04	0.42
Ref. + {5, 6, 7, 8}	0.30	0.05	0.61
Ref.: {2, 3}	1.61	0.45	-30.38
Ref. + {9}	0.18	0.04	1.30
Ref. + {14}	0.19	0.04	1.28
Ref. + {16}	0.21	0.05	1.13
Ref. + {6, 7}	0.26	0.06	1.44

Table 6

Values represent the mean difference in the performance of bilateral cochlear implant (CI) microphone configurations when additional microphones were added (see Fig. 1). The performance difference is calculated in relation to the mean performance of the *reference channel configuration*. The best result for each metric is marked in bold. All performance differences were statistically significant compared to the reference channel performance, except those marked with “†”.

Microphone IDs	Metric		
	PESQ	STOI	SI-SDR
Ref.: {1, 2, 3, 4, 5, 6, 7, 8}	2.07	0.54	-28.46
Ref. + {10}	0.11	0.01	0.02†
Ref. + {9, 11}	0.12	0.02	0.11
Ref.+ {15, 16}	-0.01	-0.01	-0.56
Ref.+ {13}	0.05	0.01	0.06†
Ref. + {9, 10, 11, 12, 13, 14, 15, 16}	0.19	0.01	-0.61†
Ref.: {2, 3, 6, 7}	1.87	0.51	-28.94
Ref. + {10}	0.16	0.03	0.49
Ref. + {13}	0.11	0.02	0.20
Ref. + {9, 11}	0.22	0.04	0.81
Ref. + {15, 16}	0.04	<0.01	-0.39†

with additional, virtually sensed, microphone signals. In this work, the estimation of the virtual microphone signals was realized by a purely data-driven deep learning approach on the raw-audio mixture without preprocessing (Stoller et al., 2018).

Most applications of deep neural networks in the domain of audio signal processing address the enhancement of speech signals by separating a target source (speech) from a mixture of interfering noise sources (Purwins et al., 2019). In the work presented here, however, no source separation was performed, but rather, in a transferred sense, a denoising of the reference signal, as explained in the following: Let the audio signal captured from a microphone inside the ear canal of the left ear be the reference signal and the audio signal inside the ear canal of the right ear the target signal. By trying to match the signal of the left ear to the right ear or denoise the left ear, we hypothesize that the network implicitly learns the RTF between the two microphone signals or, in other words, the “noise” to remove from the audio signal of the left ear. As a result, the network tries to virtually sense the right ear’s audio input by using the signal of the left ear. To evaluate the quality of the virtually sensed microphones, spatial covariance matrices Φ with and without virtually sensed microphone signals were provided as input for the MVDR beamformer (Souden et al., 2009).

The results were compared with the same metrics and statistics as with the real microphones measurements (see Section 2.4).

In this study, two microphone signals were used as reference signals, and three additional microphone signals were virtually sensed. The 2 reference signals consisted of the microphones {2, 3} and were chosen because their spatial arrangement corresponds to that of a conventional CI audio processor (see Figure 1 or Table 1). Motivated by the results of the head-mounted microphone measurements, the microphone on the forehead ({10}), the back ({13}) and inside the ear canal of the contralateral ear ({16}) were chosen as target signals for the virtual sensing approach. In the remainder of the manuscript, virtual channels are indicated by the subscript *v*. The resulting microphone configuration ({2, 3, 10_v, 13_v, 16_v}) provided the advantages as explained in the Discussion (Section 4.1): a high spatial spread of the microphone signals (Doclo et al., 2010), proximity to the target signal, and frequency transformations by the pinna and head shadow (Blauert, 1997).

2.5.1. Deep neural network architecture for the virtual sensing approach

The network architecture followed the U-Net adaption for end-to-end audio source separation in the time-domain (Stoller et al., 2018). The neural network operation on the raw-waveform in the time domain allowed to model the phase information of the audio signal, thus avoiding complex phase recovery algorithms (Griffin and Lim, 1984; Yatabe et al., 2018). The well known U-Net structure is composed as a convolutional autoencoder, and as such, consists of an encoder (contracting path), a bottleneck, and a decoder (expanding path) (Ronneberger et al., 2015). A diagram of our network’s architecture implementation is shown in Fig. 4.

In the encoder, an increasing number of higher-level features on coarser time scales were calculated, allowing the modeling of long-term dependencies in the audio signal. Our implementation of the encoder consisted of 5 levels, with each level working on half the time resolution and twice the number of feature maps as the previous one. In the bottleneck, the model was forced to learn a compression of the input data, containing only the relevant information (latent space) to construct the virtual microphone signal. The latent-space representation of the bottleneck layer was passed to the decoder, which tried to learn a mapping of the input data to match the desired virtual microphone signal. The decoder was the mirror image of the encoder and also consisted of 5 levels. Each level worked on double the time resolution and half the number of feature maps as the previous level. Based on the results of initial tests, transposed convolutions were used for the upsampling process. Each convolution was followed by group normalization, and a rectified linear unit (ReLU) activation function (Hahnloser et al., 2000; Wu and He, 2018). By introducing the skip connections in the encoder-decoder architecture, the encoder’s high-level features were concatenated with the local features computed during the upsampling block of the decoding. The result of this concatenation were multi-scale features that were fed in the output layer of the network (Ronneberger et al., 2015; Stoller et al., 2018). The output of the last convolutional layer was the estimation of the virtually sensed microphone signal.

The receptive field of the model was chosen to work with 2.1s (46077 samples), which provided an output vector with the desired test size of 1.5s (33797 samples).

Since no implicit zero padding was performed in the convolution operation, the neural network’s output sample size was smaller than the input sample size. Avoiding zero-padding allowed the convolutions to be performed in the correct audio context. As a result, audio artifacts in the results could be minimized, and the temporal continuity of the audio signal was better preserved (Stoller et al., 2018).

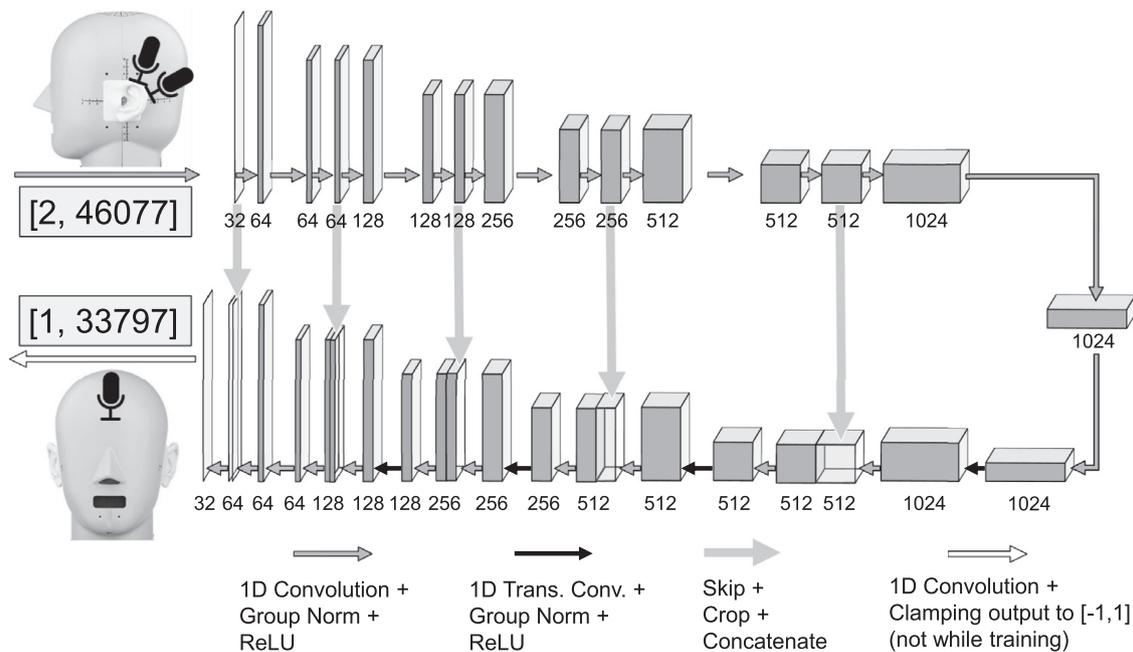


Fig. 4. The proposed deep neural network architecture for the virtual sensing of additional microphone channels based on the work of [Stoller et al. \(2018\)](#). The numbers below the blocks describe the input channel size of the following convolution. Shown is an example for the estimation of the microphone signal on the forehead ($\{10\}$) with the measurement data of 2 microphones as positioned in conventional cochlear implant (CI) audio processors (microphones $\{2, 3\}$). The network's input and output data blocks denoted with "[A, B]" describe the number of channels (A) and the number of samples (B). For an illustration of the microphone placement, please see [Figure 1](#).

2.5.2. Network training

To train the deep virtual sensing network, we extracted measurement data from the two reference channels ($\{2, 3\}$) and the microphone channel to be estimated. Due to the large size of the BCP training dataset (see [Section 2.3.2](#)), no data augmentation was necessary. In accordance with the original Wave-U-Net implementation ([Stoller et al., 2018](#)), the audio data of the BCP dataset ([Fischer et al., 2020b](#)) was downsampled to 22.05kHz. For evaluating the network's performance, the absolute differences between the actual value and the predicted value (L_1 loss) were used. To update the network weights iteratively based on training data, we applied the ADAM optimizer ([Kingma and Ba, 2014](#)) with the default decay rates of $\beta_1 = 0.9$, $\beta_2 = 0.999$ and a batch size of 16 ([Stoller et al., 2018](#)). Instead of monotonically decreasing the learning rate, cyclical learning rates ([Smith, 2017](#)) were used with upper and lower boundaries of 0.0002 and 0.00001, respectively. Early stopping was performed after 10 epochs with only minimal improvement on the validation loss. Afterward, the best model was fine-tuned with lower learning rate limits (0.000001 to 0.00001) and a batch size of 8, again until 10 epochs without improvement on the validation loss. The fine-tuned network was further used to predict the virtual channels. The test dataset to evaluate the virtually sensed microphone channels consisted of 2400 samples, which included the audio files described in [Section 2.3.1](#). Care was taken to ensure that none of the test samples were used to validate or train the network.

Since each virtual channel was estimated on a separate network, the networks were trained one after the other. The training time was reduced by successively using the previously trained network as a starting-point (transfer learning) ([Torrey and Shavlik, 2010](#)). All computations were performed with the open-source machine learning framework PyTorch version 1.6.0 ([Paszke et al., 2019](#)).

2.5.3. Subjective listening tests

Twenty normal hearing participants (6 female, 14 male, mean age in years = 29.8, SD = 3.6) performed a subjective listening test

to evaluate the benefit of the virtually sensed microphone signals on the speech quality. The test was performed in a quiet environment, and stimuli were presented via high definition insert earphones (Triple Driver, 1 More Inc. San Diego, CA) at the most comfortable loudness levels as selected by the subjects.

The questions of the subjective evaluation were twofold. First, we asked the subjects whether the signal processing applied by the MVDR beamformer lead to overall improved speech quality. Second, it was evaluated whether the beamformed signal based on the reference channels ($\{2, 3\}$) with additional virtual channels ($\{10_v\}$, $\{13_v\}$, $\{16_v\}$) outperforms the beamformed signal without virtual channels available, i.e. only the measured channels $\{2, 3\}$ were used (see [Fig. 1](#) or [Table 1](#) for a transcription of the channel IDs).

To answer these questions, the participants were asked to listen to 3 audio mixtures, all based on the same recording but either

- Beamformed based on the reference channels with additional virtual channels ($\{2, 3\} + \{10_v\}, \{13_v\}, \{16_v\}$)
- Beamformed based on the reference channels only ($\{2, 3\}$)
- The non-beamformed recording of the channels $\{2, 3\}$

The 3 audio mixtures were randomly assigned to 3 buttons on a graphical user interface (GUI). Since the beamformer's task was to enhance the speech quality for a predefined target signal, a fourth button on the GUI labeled "Target Signal" played back a recording of the corresponding target speech signal without interfering background noise. Finally, the participants' task was to select from the 3 audio mixtures the one in which the target signal was most comfortable to understand. Before the test started, trial runs were conducted until the participants confirmed that they understood the test procedure.

During the test and the trial runs, the participants were allowed to hear the 4 audio files (1 target signal and 3 audio mixtures) as many times as desired. The test stimuli consisted of 60 audio mixture quartets of 1.5 s length per file, ensuring that each file contained the utterance of at least one word. All audio mixtures were taken from the pool of the 2400 test files described in

Section 2.3 with distribution proportions as shown in Fig. 3. Evaluation of the presented audio files took about 20 min; no feedback was given during or after the test. After evaluating 30 of the 60 audio files, a pause of 3 min was taken during which the GUI was disabled. To minimize order bias, the 2 stimuli blocks that were evaluated before and after the pause were counter-balanced within the participants. The subjective listening evaluation was designed in accordance with the Declaration of Helsinki, written informed consent was obtained from all participants.

A Kruskal-Wallis test was used to determine if the frequency of choices within the 3 response options differed significantly from each other. After the Kruskal-Wallis test has rejected the null hypothesis, a post-hoc Nemenyi test was performed to investigate which of the response distributions differed significantly from each other. To determine whether the response distributions differed significantly from the chance level of the test (33 %), a chi-square test was applied. The significance level was chosen with $\alpha = 0.05$ for all statistical tests.

3. Results

3.1. MVDR beamforming with unilateral channel configurations

Table 3 shows the PESQ, STOI and SI-SDR performances of unilateral single microphone configurations compared to the performance with the reference configuration, i.e. a CI audio processor equipped with 4 microphones placed on top of the housing. For the PESQ and SI-SDR metric, the performances with single microphones were significantly worse than with the 4-channel reference configuration (all $p = 0.001$). The same was observed for STOI ($p = 0.001$) except for the microphones {1, 4} and {2, 4} (both $p = 0.9$). In all 3 metrics, the microphones that were facing the front (front {10}, left temple {11}, forward facing (audio processor) {1}, see Fig. 1 or Table 1) achieved the best results, whereas the performance differences between channels {10} and {11} were not statistically significant in terms of PESQ and SI-SDR ($p = 0.608$, $p = 0.9$) but for STOI ($p = 0.001$). Between the microphones {1} and {2} the metrics PESQ, STOI and SI-SDR did not differ significantly ($p = 0.408$, $p = 0.9$, $p = 0.115$) (a significance-matrix showing the results of the post-hoc Nemenyi tests for Table 3 can be found in the Appendix (Figures A.1-A.3)).

When the same 4-channel reference configuration (microphones {1, 2, 3, 4}) was extended by the aforementioned ipsilateral single microphone signals, again the front-facing microphones {10} and {11} (see Fig. 3) provided the greatest benefit (see Table 4). The performance differences for all metrics when channel {10} (front) was added did not differ significantly from the performance differences when channel {11} (left temple) was added to the reference configuration (PESQ: $p = 0.792$, STOI: $p = 0.736$, SI-SDR: $p = 0.9$) (a significance-matrix showing the results of the post-hoc Nemenyi tests for Table 4 can be found in the Appendix (Figures A.4-A.9)).

Since many CI audio processors record signals from 2 microphones positioned on top of the housing, the performance of different spatial arrangements of 2 microphones placed on the audio processor compared to the 4-channel reference configuration (microphones {1, 2, 3, 4}) was investigated and is shown in Table 3. The arrangement with the largest spatial distance between the 2 microphones, namely the microphones on top of the audio processors facing the front and back ({1, 4}), achieved the best performance (see Fig. 2 for a microphone distance matrix). The statistical analysis showed that the performance differences of the microphones {1, 4} did not differ significantly for PESQ and STOI from the results compared to the microphones on the audio processor facing the top and the back ({2, 4}) to the reference configuration ($p = 0.668$, $p = 0.9$). Both 2 channel microphone configurations did not differ significantly from the 4 channel reference configuration

in terms of STOI (both $p = 0.9$). For the SI-SDR metric, the differences when adding {1, 4} did not differ statistically significantly from any of the tested 2 channel configurations (all $p = 0.9$).

The arrangement with the smallest inter-microphone distance (microphones {2, 3}, see Figs. 1 and 2), which is related to the conventional microphone positions of CI audio processors, achieved the lowest scores in 2 (STOI and SI-SDR) of the 3 evaluated objective metrics, even though for SI-SDR the differences of this configuration did not differ significantly from any of the tested 2 channel configurations (all $p = 0.9$). For the metrics PESQ and STOI no significant differences in the performances between the microphones {2, 3}, {1, 2} or {1, 3} were observed (PESQ: $p = 0.721$, $p = 0.601$, STOI: $p = 0.884$, $p = 0.134$). Table 4 shows the impact on the PESQ, STOI and SI-SDR metrics when additional ipsilateral, including those on the sagittal plane, microphones were added to the the conventional microphone arrangement ({2, 3}). The extension of the microphone arrangement ({2, 3}) with forward facing microphones (front {10} or left temple {11}) provided the greatest benefit. For none of the 3 tested metrics did the performance between adding the front ({10}) or left temple ({11}) microphone to the conventional microphone arrangement differ significantly (PESQ: $p = 0.067$, STOI: $p = 0.678$, SI-SDR: $p = 0.251$).

3.2. MVDR beamforming with bilateral channel configurations

Table 5 shows the PESQ, STOI and SI-SDR performances when additional bilateral microphones were added to the input signals of an unilateral CI audio processor equipped with 4 microphones placed on top of the housing (microphones {1, 2, 3, 4}, see Fig. 1 or Table 1). When a single contralateral microphone was added, it was not the microphone closest to the target source (microphone {9}, temple) that provided the greatest benefit in terms of the human perception-related objective metrics PESQ and STOI, but the contralateral ear canal microphone {16}. Compared to adding channels {9} or {14} (temple or contralateral CI transmission coil), the improvement of the PESQ and STOI metrics were significantly better when adding the contralateral ear-canal microphone (all $p = 0.001$) (a significance-matrix showing all results of the post-hoc Nemenyi test for Table 5 can be found in the Appendix (Figures A.10-A.15)). However, in terms of SI-SDR, the input from the microphone on the contralateral CI transmission coil (microphone {14}) achieved the best SI-SDR values and even outperformed the microphone configuration compared to an additional contralateral 4-channel CI audio processor. All differences in SI-SDR with the contralateral transmission coil microphone ({14}) compared to {9} (contralateral temple), {16} (contralateral ear canal) and Ref. ch. + {5, 6, 7, 8} were not statistically significant ($p = 0.362$, $p = 0.802$, $p = 0.409$). Since the cable connection between the CI transmission-coil and the audio processor could theoretically be exploited to transmit audio signals, a unilateral microphone configuration was also used as a reference, which included the coil signal ({12}) in addition to the 4 microphones on the audio processors. The results showed in Table 5 did differ only marginally and non significantly between the reference configuration with the CI transmission coil ({1, 2, 3, 4, 12}) and the reference configuration without the CI transmission coil microphone ({1, 2, 3, 4}). The small benefit of adding microphone {12} to the reference channel configuration is also indicated by the results of Table 4.

An analysis of the results with a reference microphone configuration based on the conventional spatial microphone arrangement in CI audio processors (microphones {2, 3}, see Fig. 1 or Table 1), lead to similar conclusions as with the 4-channel microphone configuration described above (see Table 5). Again, the overall result of a single additional microphone positioned at the contralateral ear-canal {16} was best, but only with respect to STOI and PESQ. For the PESQ metric, the performance with an addi-

tional microphone positioned in the contralateral ear canal differed non-significantly compared to the performance with an additional microphone on the temple ($\{9\}$) ($p = 0.763$). In terms of SI-SDR, the microphones on the contralateral side which were close to the sagittal plane (temple $\{9\}$ and transmission coil $\{14\}$) outperformed the contralateral ear-canal microphone $\{16\}$ when added to the microphone configuration $\{2, 3\}$ ($p = 0.006$, $p = 0.9$). An additional, identical, bilaterally connected processor with 2 microphones ($\{6, 7\}$) yielded significantly better values in all metrics than adding the single microphones shown in Table 5 (see Appendix Figure A.13-A.15 for p -values).

When a bilateral CI processor microphone configuration was taken as a reference (microphones $\{1, 2, 3, 4, 5, 6, 7, 8\}$, see Table 6), adding a microphone to the front ($\{10\}$) provided more benefit than adding a microphone facing the back ($\{13\}$) (PESQ and STOI: $p = 0.001$), but for SI-SDR not statistically significant ($p = 0.515$) (a significance-matrix showing all results of the post-hoc Nemenyi test for Table 6 can be found in the Appendix (Figures A.16-A.21)). The single front microphone achieved even similar and statistically not significantly differing STOI and SI-SDR values compared to the performance when adding 2 microphones at the left and right temple ($\{9,11\}$) (both metrics $p = 0.9$). For PESQ however, the performance with the additional 2 temple microphones ($\{9,11\}$) differed statistically significant compared to the additional microphone facing to the front ($\{10\}$) ($p = 0.001$). Adding the signals of the two in-ear microphones ($\{15, 16\}$) to the bilateral CI processor microphone configuration (microphones $\{1, 2, 3, 4, 5, 6, 7, 8\}$) did not provide any benefit, not even if only 2 bilateral ($\{2, 3, 6, 7\}$) instead of 4 ($\{1, 2, 3, 4, 5, 6, 7, 8\}$) bilateral processor microphones were used as a reference microphone configuration. The full 16-channel microphone configuration achieved the statistically significant best PESQ scores (all $p = 0.001$). However, in terms of STOI and SI-SDR the performance did barely, and for SI-SDR non significantly, differ compared to the 8-channel reference microphone configuration. Again, as with the unilateral 4-microphone CI audio processor configuration, adding the transmission-coil microphone signals ($\{12, 14\}$) to the bilateral microphone configurations ($\{1, 2, 3, 4, 5, 6, 7, 8\}$ or $\{2, 3, 6, 7\}$) did barely and statistically not significant influence the performance metrics shown in Table 6.

3.3. Virtual sensing of microphone channels

The bar graphs in Fig. 5 compare the performance in PESQ, STOI and SI-SDR (see Methods Section 2.4) between virtually sensed microphone signals and actually measured microphone signals placed at the same position on the head, i.e. the front ($\{10\}$), the back ($\{13\}$) and at the entry of the right external auditory canal ($\{16\}$) (see Fig. 1 or Table 1). For all 3 objective speech quality metrics tested, adding virtually sensed microphone signals to the input signals of the MVDR beamformer resulted in a significant improvement compared to the performance with microphone signals as used in conventional CI audio processors ($\{2, 3\}$) ($p < 0.001$).

The mean benefit in performance when additional virtual/measured microphone signals were used for beamforming was 0.24/0.34 units for PESQ, 0.06/0.07 units for STOI, and 1.17/1.25 dB for SI-SDR. For the PESQ and STOI metrics, the performance between the virtually sensed microphone signals and the measured microphones signals differed significantly ($p < 0.001$). In terms of SI-SDR, no significant difference between the two configurations were observed ($p = 0.998$).

An analysis of the performance of the neural networks with respect to each of the estimated channels $\{16\}$, $\{13\}$ and $\{10\}$ showed that the mean benefit when an additional virtual/measured microphone signal was used for beamforming was 0.154/0.211, 0.114/0.149, 0.178/0.219 for PESQ, 0.049/0.052, 0.028/0.032, 0.042/0.048 for STOI, and 1.000/1.057, 0.938/0.877,

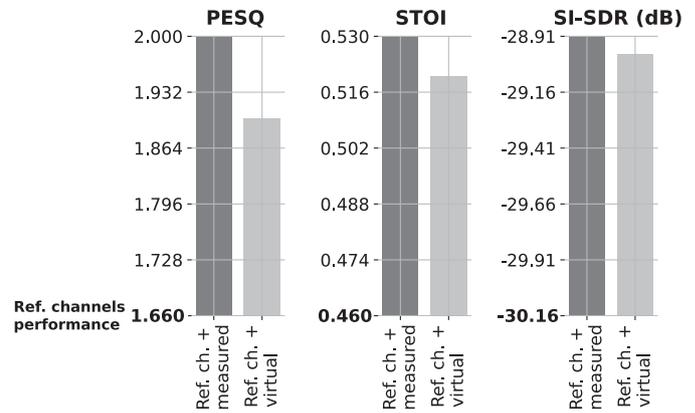


Fig. 5. Comparison of overall perceptual evaluation of speech quality (PESQ), short-time objective intelligibility (STOI) and scale-invariant speech to distortion ratio (SI-SDR) scores between 3 different microphone channel configurations used as input signals for the minimum variance distortionless response (MVDR) beamforming algorithm (Souden et al., 2009): 1) Reference channel configuration according to the conventional microphone placement on CI audio processors (microphone IDs $\{2, 3\}$) (bold letters); 2) Reference channel configuration with additional measured (real) microphones (microphone IDs $\{2, 3\} + \{10, 13, 16\}$) (dark grey bar); 3) Reference channel configuration with additional virtually sensed microphones (microphone IDs $\{2, 3\} + \{10_v, 13_v, 16_v\}$) (light grey bar). The dataset used to evaluate the microphone channel configurations consisted of 2400 cocktail party audio samples, as described in Section 2.3. Please see Fig. 1 or Table 1 for a description of the microphone IDs.

1.493/1.377 for SI-SDR. For the metrics PESQ and STOI the differences in performance between the additional virtually estimated microphone and the measured microphone were significant (all $p < 0.001$). For SI-SDR, the differences were significant only with respect to microphone channel $\{10\}$ ($p = 0.027$), but not for the channels $\{13\}$ and $\{16\}$ ($p = 0.244$, $p = 0.309$). The on average bad results for channel $\{16\}$, meaning the largest difference between the benefit of additional virtual/measured microphone signals, and the best results for channel $\{10\}$ were also reflected in the validation losses of the trained networks. For channel $\{16\}$, $\{13\}$ and $\{10\}$, the best L_1 -losses on the validation set were 2.1×10^{-4} , 1.5×10^{-4} and 1.4×10^{-4} , respectively.

3.3.1. Subjective listening tests

Figure 6 shows that the participants preferred the audio mixture that was beamformed using the additional virtual channels (Mean=65%, SD=8%) compared to a beamformed signal generated using only the microphones as placed in CI audio processors (Mean=23%, SD=4%). This difference in selection frequency was statistically significant with $p < 0.001$.

The non-beamformed signal was rarely selected as the signal that was easiest to understand (Mean=13%, SD=7%). The beamformed signal based on the reference channel only and the beamformed signal based on additional virtual channels differed significantly to the non-beamformed audio mixture selection frequency ($p = 0.002$, $p < 0.001$).

For all of the presented signal configurations, the distribution of the frequency of choices differed significantly from the chance level of the test (all $p < 0.001$).

To investigate if the subjects' choice of the ãsignal most comfortable to understandg was dependent on the SNR of the original or raw audio mixture, the SNRs of the corresponding raw audio mixtures were compared. It was observed that the subjects preferred the beamformed signal with additional virtual channels if the SNRs of the raw audio mixture were low (Mean=2.4, SD=9.3) compared to the raw audio mixtures' SNRs when the beamformed signal based on the reference channels only was selected (Mean=5.2, SD=8.0, $p = 0.001$). The SNRs of the raw audio mix-

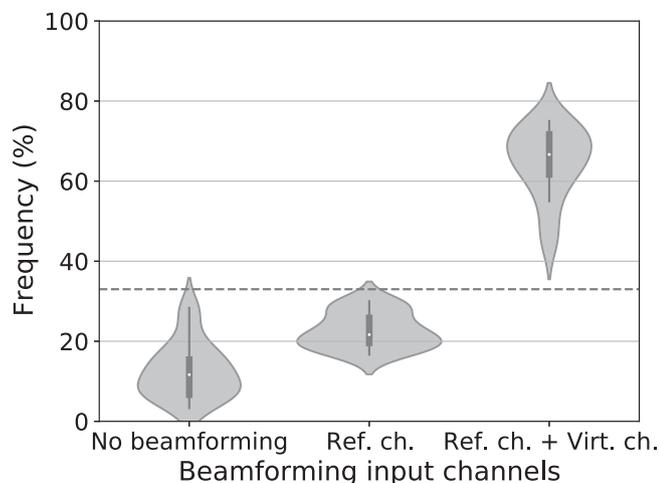


Fig. 6. Violin plots (Hintze and Nelson, 1998) of the frequency of choices in the subjective listening test. The data represents the choices for the non-beamformed signal, the beamformed signal with the measured reference channel configuration as input channels (microphone IDs {2, 3}) and the beamformed signal with additional virtually sensed microphone signals as input channels (microphone IDs {2, 3} + {10_v, 13_v, 16_v}) (see Fig. 1 or Table 1). The dashed horizontal line indicates the chance level of the test. The probability of observations taking a given value (Frequency (%)) is indicated by the violin's width, while each violin is normalized to have the same area. The thick black bar in the center of the violin represents the interquartile range. The thin black line extended from it represents the 95% confidence intervals, and the white dot represents the median.

tures when the non-beamformed signal was selected (Mean=2.1, SD=9.2) was not significantly different from the SNRs of the raw audio mixtures when the beamformed signal with additional virtual channels was selected ($p = 0.987$). However, it was significantly different from the SNRs of the raw audio mixtures when the beamformed signal based on the reference channels was chosen ($p = 0.029$).

4. Discussion

Herein, we presented a comprehensive comparison of different head-mounted microphone configurations and their effect on the output of an MVDR beamforming algorithm. The results showed that microphone positions, such as placing a microphone on the forehead, would be desirable for better speech understanding. Since these microphone positions are not practicable in reality, we proposed and evaluated a purely data-driven virtual sensing technique.

4.1. Association of the speech quality and the microphone positioning

Our measurements of varying head-mounted microphone arrangements in cocktail party scenarios confirmed that the performance of beamforming algorithms and thus the speech quality improves with additional microphone signals (Doclo et al., 2010). Single-microphone speech-enhancement algorithms can only exploit temporal and spectral information cues, whereas multi-microphone beamformers can additionally exploit the spatial information of the sound sources (Doclo et al., 2010; Souden et al., 2009).

However, a high number of microphones alone does not necessarily lead to a better speech quality (Souden et al., 2009). In the case of bilaterally placed microphones (Table 6), we observed saturation in terms of speech signal enhancement with additional microphones that were placed close to the reference microphones. In particular, the SI-SDR metric showed that noise from additional microphone signals can dominate compared to the redundant information in the audio signal used for speech enhancement. As

also shown by Corey et al. (2019), the microphone arrangement's spatial diversity played a significant role in the quality of the acoustic beamforming. The herein performed measurements confirmed this finding since no improvements were observed when additional microphones were placed at a distance of about 5cm to the reference microphones. It was assumed that even for low frequencies, these microphones were too closely spaced to provide inter-microphone information for the beamforming algorithm (Corey et al., 2019). Besides, the microphones' distance was too small for an effect of the acoustic head shadow (Blauert, 1997). With the same reasoning, the slightly worse result of the unilateral, conventional microphone configuration ({2, 3}) and the good result of the arrangement with the largest inter-microphone distance (front and back facing {1, 4}) compared to other 2-channel microphone arrangements on the audio processor can be argued.

Although adding a microphone with a high Euclidean distance to the reference microphone configuration is a good rule of thumb to improve acoustic beamforming, other microphone positioning factors, such as exploiting the acoustic head shadow (Blauert, 1997), may be just as important. In the unilateral configuration (see Table 4), we observed that the proximity to the most likely target source with an additional microphone on the temple ({2, 3}+{11}) was more important than the spatial diversity of the microphones with an additional microphone placed on the back of the head ({2, 3}+{13}). In addition to the proximity to the target signal and the microphone distance, our measurements confirmed that the pinna's directional frequency transformation provided relevant information for improving the quality of the beamforming algorithm (Blauert, 1997; Fischer et al., 2021; Wimmer et al., 2016). We observed that the most useful additional contralateral microphone was neither the one closest to the target signal ({11}, temple) nor the one with the highest Euclidean distance to the reference microphone configuration ({14}, CI transmission coil). It was the contralateral microphone placed in the ear canal facing away from the target signal ({16}).

4.2. Virtual sensing of head-mounted microphone signals

In this work, we presented and evaluated a method for virtual sensing of microphone signals to improve the speech quality of hearing aid and CI users in noisy environments. The proposed methodology enabled to capture microphone signals at positions on the head, including but not limited to the forehead, where a physical placement of microphones is impractical. Our objective measurements showed, that adding strategically positioned virtual microphones on the head significantly improved the speech quality compared to the speech quality as obtained with a microphone arrangement found in conventional CI audio processors. This result was also confirmed in human listening tests using a 3-alternative forced-choice procedure with the task of selecting the speech mixture that was most comfortable to understand.

In addition to the general assumption that adding microphone signals to hearing aid applications can increase the performance of beamforming algorithms (Doclo et al., 2010), we hypothesized and confirmed that replacing real microphone signals with virtual microphone signals can also increase beamformer performance. In contrast to the work presented in (Jinzai et al., 2019; Katahira et al., 2016; Yamaoka et al., 2019), our entirely data-driven approach showed that explicit knowledge of the real microphone's positioning might not be necessary to enhance the speech quality with virtual microphone channels. The mathematical reasoning for the success of our deep learning-based approach is the subject of ongoing research (Fan et al., 2020; Yu and Principe, 2019).

In the measurements with the reference microphone configuration according to conventional CI audio processors ({2, 3}), we observed that an additional microphone on the forehead produced

similar improvements in speech quality as an additional microphone placed at the entry of the contralateral ear canal. However, due to the poor estimation of the contralateral ear signal by the neural network, a higher benefit was obtained with the virtual microphone channel estimating the signal at the forehead. Therefore the estimation of optimal microphone positions for neural network-based beamforming approaches requires further investigation.

The subjective feedback of the 20 participants significantly showed that the additional virtual microphone signals were preferred, especially in cocktail party scenarios with low SNRs. On the other hand, the participants' choices also showed that in low SNRs scenarios, the MVDR beamforming, either with real or real and additional virtual channels, might degrade the subjective speech signal quality instead of enhancing it. This finding confirmed that although MVDR beamformers aim to keep the target signal undistorted (Veen and Buckley, 1988), there was a trade-off between noise reduction and speech signal distortion (Souden et al., 2009).

4.3. Limitations and outlook

Although the virtually sensed microphones significantly improved the speech quality within this study, further research is needed before the methodology can be used in hearing aids or CI audio processors.

Due to the input data size of 2 s, the delay of the proposed network architecture is too long to be applicable in a real hearing aid application. However, this paper's main objective was to demonstrate a proof of concept for purely data-driven virtual channel estimations in hearing aids or CIs. Tackling the problem of latency and neural network complexity in online speech enhancement is ongoing research (Koyama et al., 2020; Luo and Mesgarani, 2019; Pandey and Wang, 2019; Takeuchi et al., 2020) with promising results and input frame lengths as little as 2s (Luo and Mesgarani, 2019). Future research should investigate whether the significant reduction in network time delay required for an application in hearing devices affects the performance of the presented approach. In addition to progress in reducing the computational costs, substantial progress is continuously being made in other areas of speech signal enhancement with artificial neural networks relevant for the methodology of this work, such as in blind source separation (BSS) (Drude and Haeb-Umbach, 2019; Drude et al., 2019; Lu et al., 2019), acoustic scene classification (ASC) (Kong et al., 2020; Koutini et al., 2019; Phaye et al., 2019), domain shift (Borsos et al., 2020; Mathur et al., 2019) and the usage of loss functions to optimize the parameters of the network based on the human perception of speech (Fu et al., 2019; Koyama et al., 2020). The results of Drude et al. (2019) indicated, that the benefit of the presented approach when using estimated coherence matrices may be different from the benefit achieved with the oracle matrices. For computational time reasons, no sophisticated optimization of the presented network's architecture was performed. Further research may investigate the optimal number and size of hidden layers for the presented approach.

Our approach follows a two-step procedure to estimate a virtual microphone channel that is used as an additional input to the beamformer. We chose this procedure to improve the compatibility with existing beamforming technology in current devices. However, the entire approach could be replaced by an end-to-end single-network artificial intelligence solution for hearing devices.

One of the biggest challenges of the presented methodology to be applicable in a real-world application will be to ensure the robustness of the network's predictions in acoustic environments with high reverberation (Chen et al., 2017; Feng and Kowalski, 2019; Inoue et al., 2019; Xie et al., 2019). In the context of this work, the first step in this direction would be the use of more

challenging acoustic training data, for example, by simulating conditions with higher reverberation (Cuevas-Rodríguez et al., 2019) or the use of dynamically moving sound sources (Fischer et al., 2020c; 2020d). Another possibility would be to record acoustic scenarios using a portable microphone array (Kiselev et al., 2017). In a real-world application, this data could be collected as part of an audiological fitting routine. In both cases, whether the data was simulated or recorded in real environments for each subject, the additional recordings and the personalization of the network through transfer learning would most likely increase the robustness of applied neural network solutions (Zhuang et al., 2021). To account for the different head geometries and thus varying inter-microphone features, the information of 3D head scans as provided in Fischer et al. (2020b) could be fed into a neural network architecture that allows metadata injection.

Although the speech quality may improve by applying the proposed measures, binaural cues would still be discarded, resulting in a low spatial quality of the perceived sounds (Blauert, 1997). It remains unclear whether the findings of this study will also hold for current state-of-the-art beamformers with binaural output. To preserve the binaural cues and thus improve the spatial quality of the MVDR beamforming algorithm (Souden et al., 2009), adaptations such as those proposed by Marquardt et al. (2015) or Marquardt and Doclo (2018) could yield improvements in this regard while still enhancing the speech quality (Gößling et al., 2020).

5. Conclusions

In this work, real and virtual microphone signals were combined as input for an MVDR beamformer to investigate the effects on speech quality for hearing aid or CI users in cocktail party scenarios. The measurements with respect to the number and spatial arrangement of real microphones indicated that, optimally, microphones should be placed as close as possible to the target source, encode monaural cues, and produce a large distance spread by their spatial arrangement. In reality, however, it is inconvenient to place the microphones according to these criteria. To overcome this problem, virtual microphone signals were estimated using a deep neural network without explicit knowledge of the spatial microphone arrangement. The results of 3-alternative forced choice subjective listening tests and objective speech quality metrics suggest that hearing aid or CI users might benefit from virtually sensed microphone signals, especially in challenging cocktail party scenarios.

Funding

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Appendix A. Additional Figures

Please see appendix_A.pdf for significance-matrices of the post-hoc Nemenyi tests concerning the data in Tables 3–6.

Supplementary material

Supplementary material associated with this article can be found, in the online version, at [10.1016/j.heares.2021.108294](https://doi.org/10.1016/j.heares.2021.108294).

CRediT authorship contribution statement

Tim Fischer: Conceptualization, Methodology, Software, Formal analysis, Writing - original draft, Investigation, Data curation, Visualization. **Marco Caversaccio:** Resources, Supervision. **Wilhelm Wimmer:** Conceptualization, Validation, Writing - review & editing, Supervision, Methodology, Project administration.

References

- Arora, R., Amoodi, H., Stewart, S., Friesen, L., Lin, V., Nedzelski, J., Chen, J., 2013. The addition of a contralateral routing of signals microphone to a unilateral cochlear implant system a prospective study in speech outcomes. *Laryngoscope* 123, 746–751.
- Benesty, J., Chen, J., Huang, Y., 2008. *Microphone array signal processing*. Springer Science & Business Media.
- Blackman, R.B., Tukey, J.W., 1959. Particular pairs of windows. *The Measurement of Power Spectra, From the Point of View of Communications Engineering* 98–99.
- Blauert, J., 1997. *Spatial Hearing: The Psychophysics of Human Sound Localization*. MIT Press.
- Borsos, Z., Li, Y., Gfeller, B., Tagliasacchi, M., 2020. MicAugment: one-shot microphone style transfer. *arXiv preprint arXiv:2010.09658*.
- Capon, J., 1969. High-resolution frequency-wavenumber spectrum analysis. *Proc. IEEE* 57, 1408–1418.
- Chen, Z., Li, J., Xiao, X., Yoshioka, T., Wang, H., Wang, Z., Gong, Y., 2017. Cracking the cocktail party problem by multi-beam deep attractor network. In: 2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), pp. 437–444. doi:10.1109/ASRU.2017.8268969.
- Cherry, E.C., 1953. Some experiments on the recognition of speech, with one and with two ears. *J. Acoust. Soc. Am.* 25, 975–979. doi:10.1121/1.1907229. <http://asa.scitation.org/doi/10.1121/1.1907229>
- Corey, R.M., Tsuda, N., Singer, A.C., 2019. Acoustic impulse responses for wearable audio devices. In: ICASSP 2019–2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, pp. 216–220.
- Cuevas-Rodríguez, M., Picinali, L., González-Toledo, D., Garre, C., de la Rubia-Cuevas, E., Molina-Tanco, L., Reyes-Lecuona, A., 2019. 3D tune-in toolkit: an open-source library for real-time binaural spatialisation. *PLOS ONE* 14, e0211899. doi:10.1371/journal.pone.0211899. <http://dx.plos.org/10.1371/journal.pone.0211899>
- Denk, F., Ernst, S.M., Ewert, S.D., Kollmeier, B., 2018. Adapting hearing devices to the individual ear acoustics: database and target response correction functions for various device styles. *Trends Hearing* 22. doi:10.1177/2331216518779313.
- Doclo, S., Gannot, S., Moonen, M., Spriet, A., Haykin, S., Liu, K.R., 2010. Acoustic beamforming for hearing aid applications. In: *Handbook on Array Processing and Sensor Networks*, pp. 269–302.
- Dorman, M.F., Natale, S.C., Agrawal, S., 2018. The value of unilateral CIs, CI-CROS and bilateral CIs, with and without beamformer microphones, for speech understanding in a simulation of a restaurant environment. *Audiol. Neurotol.* 23, 270–276.
- Drude, L., Haeb-Umbach, R., 2017. Tight integration of spatial and spectral features for BSS with deep clustering embeddings. In: *Interspeech*, pp. 2650–2654.
- Drude, L., Haeb-Umbach, R., 2019. Integration of neural networks and probabilistic spatial models for acoustic blind source separation. *IEEE J. Sel. Top. Signal Process.* 13, 815–826.
- Drude, L., Hasenklever, D., Haeb-Umbach, R., 2019. Unsupervised training of a deep clustering model for multichannel blind source separation. In: ICASSP 2019–2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, pp. 695–699.
- Erdogan, H., Hershey, J.R., Watanabe, S., Mandel, M.I., Roux, J.L., 2016. Improved MVDR beamforming using single-channel mask prediction networks. In: *Interspeech*, pp. 1981–1985.
- Fan, F., Xiong, J., Wang, G., 2020. On interpretability of artificial neural networks. Preprint at <https://arxiv.org/abs/2001.02522>.
- Feng, F., Kowalski, M., 2019. Underdetermined reverberant blind source separation: Sparse approaches for multiplicative and convolutive narrowband approximation. *IEEE/ACM Trans. Audio Speech Lang. Process.* 27, 442–456. doi:10.1109/TASLP.2018.2881925.
- Feng, Z.G., Yiu, K.F.C., Nordholm, S.E., 2011. Placement design of microphone arrays in near-field broadband beamformers. *IEEE Trans. Signal Process.* 60, 1195–1204.
- Fischer, T., Caversaccio, M., Wimmer, W., 2020. A front-back confusion metric in horizontal sound localization: the FBC score. In: *ACM Symposium on Applied Perception* 2020, pp. 1–5.
- Fischer, T., Caversaccio, M., Wimmer, W., 2020. Multichannel acoustic source and image dataset for the cocktail party effect in hearing aid and implant users. *Sci. Data* 7. doi:10.1038/s41597-020-00777-8.
- Fischer, T., Caversaccio, M., Wimmer, W., 2020c. System for combined hearing and balance tests of a person with moving sound source devices. <https://patents.google.com/patent/WO2020254462A1WO Patent WO2020254462A1>.
- Fischer, T., Kompis, M., Mantokoudis, G., Caversaccio, M., Wimmer, W., 2020. Dynamic sound field audiometry: Static and dynamic spatial hearing tests in the full horizontal plane. *Appl. Acoust.* 166, 107363. doi:10.1016/j.apacoust.2020.107363.
- Fischer, T., Schmid, C., Kompis, M., Mantokoudis, G., Caversaccio, M., Wimmer, W., 2021. Pinna-imitating microphone directionality improves sound localization and discrimination in bilateral cochlear implant users. *Ear Hearing* 42, 214.
- Fu, S.-W., Liao, C.-F., Tsao, Y., Lin, S. D., 2019. MetricGAN: generative adversarial networks based black-box metric scores optimization for speech enhancement. *arXiv preprint arXiv:1905.04874*.
- Gößling, N., Marquardt, D., Doclo, S., 2020. Perceptual evaluation of binaural MVDR-based algorithms to preserve the interaural coherence of diffuse noise fields. *Trends Hearing* 24, 10.1177/2331216520919573.
- Griffin, D., Lim, J., 1984. Signal estimation from modified short-time fourier transform. *IEEE Trans. Acoust. Speech Signal Process.* 32, 236–243.
- Habets, E.A., Benesty, J., Gannot, S., Cohen, I., 2010. The MVDR beamformer for speech enhancement. In: *Speech Processing in Modern Communication*. Springer, pp. 225–254.
- Habets, E.A.P., Benesty, J., Cohen, I., Gannot, S., Dmochowski, J., 2009. New insights into the MVDR beamformer in room acoustics. *IEEE Trans. Audio Speech Lang. Process.* 18, 158–170.
- Hahnloser, R.H., Sarpeshkar, R., Mahowald, M.A., Douglas, R.J., Seung, H.S., 2000. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature* 405, 947–951.
- Higuchi, T., Kinoshita, K., Delcroix, M., Nakatani, T., 2017. Adversarial training for data-driven speech enhancement without parallel corpus. In: 2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU). IEEE, pp. 40–47.
- Himawan, I., Sridharan, S., McCowan, I., 2008. Dealing with uncertainty in microphone placement in a microphone array speech recognition system. In: 2008 IEEE International Conference on Acoustics, Speech and Signal Processing. IEEE, pp. 1565–1568.
- Hintze, J.L., Nelson, R.D., 1998. Violin plots: a box plot-density trace synergism. *Am. Stat.* 52, 181–184.
- Inoue, S., Kameoka, H., Li, L., Seki, S., Makino, S., 2019. Joint separation and dereverberation of reverberant mixtures with multichannel variational autoencoder. In: ICASSP 2019 – 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 96–100. doi:10.1109/ICASSP.2019.8683497.
- Ito, N., Araki, S., Nakatani, T., 2016. Complex angular central gaussian mixture model for directional statistics in mask-based microphone array signal processing. In: 2016 24th European Signal Processing Conference (EUSIPCO). IEEE, pp. 1153–1157.
- Jinzai, R., Yamaoka, K., Matsumoto, M., Makino, S., Yamada, T., 2019. Wave-length proportional arrangement of virtual microphones based on interpolation/extrapolation for underdetermined speech enhancement. In: 2019 27th European Signal Processing Conference (EUSIPCO). IEEE, pp. 1–5.
- Jones, H.G., Kan, A., Litovsky, R.Y., 2016. The effect of microphone placement on interaural level differences and sound localization across the horizontal plane in bilateral cochlear implant users. *Ear Hearing* 37, e341.
- Katahira, H., Ono, N., Miyabe, S., Yamada, T., Makino, S., 2016. Nonlinear speech enhancement by virtual increase of channels and maximum SNR beamformer. *EURASIP J. Adv. Signal Process.* 2016, 1–8.
- Kingma, D. P., Ba, J., 2014. Adam: a method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kiselev, I., Ceolini, E., Wong, D., De Cheveigne, A., Liu, S., 2017. WHISPER: wirelessly synchronized distributed audio sensor platform. In: 2017 IEEE 42nd Conference on Local Computer Networks Workshops (LCN Workshops), pp. 35–43. doi:10.1109/LCN.Workshops.2017.62.
- Kong, Q., Cao, Y., Iqbal, T., Wang, Y., Wang, W., Plumbley, M.D., 2020. PANNs: large-scale pretrained audio neural networks for audio pattern recognition. *IEEE/ACM Trans. Audio Speech Lang. Process.* 28, 2880–2894.
- Koutini, K., Eghbal-zadeh, H., Dorfer, M., Widmer, G., 2019. The receptive field as a regularizer in deep convolutional neural networks for acoustic scene classification. In: 2019 27th European Signal Processing Conference (EUSIPCO), pp. 1–5. doi:10.23919/EUSIPCO.2019.8902732.
- Koyama, Y., Vuong, T., Uhlich, S., Raj, B., 2020. Exploring the best loss function for DNN-based low-latency speech enhancement with temporal convolutional networks. *arXiv preprint arXiv:2005.11611*.
- Lu, J., Cheng, W., He, D., Zi, Y., 2019. A novel underdetermined blind source separation method with noise and unknown source number. *J. Sound Vib.* 457, 67–91.
- Luo, Y., Mesgarani, N., 2019. Conv-TasNet: surpassing ideal time-frequency magnitude masking for speech separation. *IEEE/ACM Trans. Audio Speech Lang. Process.* 27, 1256–1266.
- Makino, S., Lee, T.-W., Sawada, H., 2007. *Blind Speech Separation*. Springer.
- Marquardt, D., Doclo, S., 2018. Interaural coherence preservation for binaural noise reduction using partial noise estimation and spectral postfiltering. *IEEE/ACM Trans. Audio Speech Lang. Process.* 26, 1261–1274. doi:10.1109/TASLP.2018.2823081.
- Marquardt, D., Hohmann, V., Doclo, S., 2015. Interaural coherence preservation in multi-channel wiener filtering-based noise reduction for binaural hearing aids. *IEEE/ACM Trans. Audio Speech Lang. Process.* 23, 2162–2176. doi:10.1109/TASLP.2015.2471096.
- Mathur, A., Isopoussu, A., Kawsar, F., Berthouze, N., Lane, N.D., 2019. Mic2Mic: Using cycle-consistent generative adversarial networks to overcome microphone variability in speech systems. In: *Proceedings of the 18th International Conference on Information Processing in Sensor Networks*, pp. 169–180.
- Mesgarani, N., Chang, E.F., 2012. Selective cortical representation of attended speaker in multi-talker speech perception. *Nature* 485, 233–236.
- Pandey, A., Wang, D., 2019. TCNN: temporal convolutional neural network for real-time speech enhancement in the time domain. In: ICASSP 2019–2019 IEEE

- International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, pp. 6875–6879.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S., 2019. PyTorch: an imperative style, high-performance deep learning library. In: *Advances in Neural Information Processing Systems* 32. Curran Associates, Inc., pp. 8024–8035. <http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- Peissig, J., Kollmeier, B., 1997. Directivity of binaural noise reduction in spatial multiple noise-source arrangements for normal and impaired listeners. *J. Acoust. Soc. Am.* 101, 1660–1670.
- Phaye, S.S.R., Benetos, E., Wang, Y., 2019. SubSpectralNet using sub-spectrogram based convolutional neural networks for acoustic scene classification. In: *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 825–829. doi:10.1109/ICASSP.2019.8683288.
- Purwins, H., Li, B., Virtanen, T., Schlüter, J., Chang, S.-Y., Sainath, T., 2019. Deep learning for audio signal processing. *IEEE J. Sel. Top. Signal Process.* 13, 206–219.
- Qian, Y.-m., Weng, C., Chang, X.-k., Wang, S., Yu, D., 2018. Past review, current progress, and challenges ahead on the cocktail party problem. *Front. Inf. Technol. Electron. Eng.* 19, 40–63.
- Rix, A.W., Beerends, J.G., Hollier, M.P., Hekstra, A.P., 2001. Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs. In: *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No. 01CH37221)*, vol. 2. IEEE, pp. 749–752.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: convolutional networks for biomedical image segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, pp. 234–241.
- Roux, J.L., Wisdom, S., Erdogan, H., Hershey, J.R., 2019. SDR-half-baked or well done? In: *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 626–630.
- Smaragdīs, P., 1998. Blind separation of convolved mixtures in the frequency domain. *Neurocomputing* 22, 21–34.
- Smith, L.N., 2017. Cyclical learning rates for training neural networks. In: *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, pp. 464–472.
- Souden, M., Benesty, J., Affes, S., 2009. On optimal frequency-domain multichannel linear filtering for noise reduction. *IEEE Trans. Audio Speech Lang. Process.* 18, 260–276.
- Stoller, D., Ewert, S., Dixon, S., 2018. Wave-u-net: A multi-scale neural network for end-to-end audio source separation. *arXiv preprint arXiv:1806.03185*.
- Taal, C.H., Hendriks, R.C., Heusdens, R., Jensen, J., 2011. An algorithm for intelligibility prediction of time-frequency weighted noisy speech. *IEEE Trans. Audio Speech Lang. Process.* 19, 2125–2136.
- Takeuchi, D., Yatabe, K., Koizumi, Y., Oikawa, Y., Harada, N., 2020. Real-time speech enhancement using equilibrated RNN. In: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 851–855.
- Torrey, L., Shavlik, J., 2010. Transfer learning. In: *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques*. IGI global, pp. 242–264.
- Veen, B.D.V., Buckley, K.M., 1988. Beamforming: a versatile approach to spatial filtering. *IEEE ASSP Mag.* 5, 4–24.
- Wimmer, W., Caversaccio, M., Kompis, M., 2015. Speech intelligibility in noise with a single-unit cochlear implant audio processor. *Otol. Neurotol.* 36, 1197–1202.
- Wimmer, W., Kompis, M., Stieger, C., Caversaccio, M., Weder, S., 2017. Directional microphone contralateral routing of signals in cochlear implant users: a within-subjects comparison. *Ear Hearing* 38, 368–373.
- Wimmer, W., Weder, S., Caversaccio, M., Kompis, M., 2016. Speech intelligibility in noise with a pinna effect imitating cochlear implant processor. *Otol. Neurotol.* 37, 19–23.
- Wouters, J., Berghe, J.V., 2001. Speech recognition in noise for cochlear implantees with a two-microphone monaural adaptive noise reduction system. *Ear Hearing* 22, 420–430.
- Wu, Y., He, K., 2018. Group normalization. In: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 3–19.
- Wu, Y.-H., Stangl, E., Chipara, O., Hasan, S.S., Welhaven, A., Oleson, J., 2018. Characteristics of real-world signal-to-noise ratios and speech listening situations of older adults with mild-to-moderate hearing loss. *Ear Hearing* 39, 293.
- Xie, Y., Xie, K., Yang, J., Wu, Z., Xie, S., 2019. Underdetermined reverberant audio-source separation through improved expectation-maximization algorithm. *Circuits Syst. Signal Process.* 38, 2877–2889.
- Yamaoka, K., Li, L., Ono, N., Makino, S., Yamada, T., 2019. Cnn-based virtual microphone signal estimation for MPDR beamforming in underdetermined situations. In: *2019 27th European Signal Processing Conference (EUSIPCO)*, pp. 1–5. doi:10.23919/EUSIPCO.2019.8903040.
- Yatabe, K., Masuyama, Y., Oikawa, Y., 2018. Rectified linear unit can assist Griffin-Lim phase recovery. In: *2018 16th International Workshop on Acoustic Signal Enhancement (IWAENC)*. IEEE, pp. 555–559.
- Yu, S., Principe, J.C., 2019. Understanding autoencoders with information theoretic concepts. *Neural Netw.* 117, 104–123.
- Zhuang, F., Qi, Z., Duan, K., Xi, D., Zhu, Y., Zhu, H., Xiong, H., He, Q., 2021. A comprehensive survey on transfer learning. *Proceedings IEEE* 109, 43–76. doi:10.1109/JPROC.2020.3004555.