# SEQUENTIAL MULTI-FRAME NEURAL BEAMFORMING FOR SPEECH SEPARATION AND ENHANCEMENT

*Zhong-Qiu Wang*[1,2,*], *Hakan Erdogan*[1], *Scott Wisdom*[1], *Kevin Wilson*[1],
*Desh Raj*[3], *Shinji Watanabe*[3], *Zhuo Chen*[4], *John R. Hershey*[1]

[1]Google Research, Cambridge, MA, [2]MERL, Cambridge, MA
[3]Johns Hopkins University, Baltimore, MD, [4]Microsoft, Seattle, WA

wang.zhongqiu41@gmail.com, {hakanerdogan, scottwisdom, kwwilson, johnhershey}@google.com,
draj@cs.jhu.edu, shinjiw@ieee.org, zhuc@microsoft.com

## Abstract

This work introduces sequential neural beamforming, which alternates between neural network based spectral separation and beamforming based spatial separation. Our neural networks for separation use an advanced convolutional architecture trained with a novel stabilized signal-to-noise ratio loss function. For beamforming, we explore multiple ways of computing time-varying covariance matrices, including factorizing the spatial covariance into a time-varying amplitude component and a time-invariant spatial component, as well as using block-based techniques. In addition, we introduce a multi-frame beamforming method which improves the results significantly by adding contextual frames to the beamforming formulations. We extensively evaluate and analyze the effects of window size, block size, and multi-frame context size for these methods. Our best method utilizes a sequence of three neural separation and multi-frame time-invariant spatial beamforming stages, and demonstrates an average improvement of 2.75 dB in scale-invariant signal-to-noise ratio and 14.2% absolute reduction in a comparative speech recognition metric across four challenging reverberant speech enhancement and separation tasks. We also use our three-speaker separation model to separate real recordings in the LibriCSS evaluation set into non-overlapping tracks, and achieve a better word error rate as compared to a baseline mask based beamformer.

## 1. Introduction

Audio source separation has many applications, for example as a front end for robust automatic speech recognition (ASR) and to improve voice quality for telephony. Leveraging multiple microphones has great potential to improve separation, since the spatial relationship among microphones provides complementary information to spectral cues exploited by monaural approaches. Multi-microphone processing can also improve the suppression of reverberation and diffuse background noise.

Recently, a new paradigm has emerged as a promising alternative to conventional beamforming approaches: neural beamforming, where the key advance is to utilize the non-linear modeling power of deep neural networks (DNN) to identify time-frequency (T-F) units dominated by each source for spatial covariance matrix computation [1, 2]. Unlike traditional approaches, neural beamforming methods have the potential to learn and adapt from massive training data, which improves their robustness to unknown positions and orientations of microphones and sources, types of acoustic sources, and room geometry. An initial success of neural beamforming approaches

---
*Work done during an internship at Google.

was improving time-invariant beamforming using T-F domain mask prediction, where predicted masks were used to obtain time-invariant spatial covariance matrices for all sources. This has proven useful in ASR tasks such as CHiME-3/4 [3]. Recent studies considered online or low-latency beamforming [4, 5] and time-varying beamforming [6] for improved performance in certain scenarios. In addition, spatial features such as inter-channel phase differences (IPD) [7], cosine and sine IPDs [8] and target direction compensated IPDs [9], which can encode spatial information, are utilized as additional network input to improve the mask estimation in masking-based beamforming. Other cues, such as visual information [10], location information [11] and speaker embeddings [12, 13], can also be used as additional inputs to improve target extraction and source separation in both single- and multi-microphone setups.

This paper explores alternating between spectral estimation using DNN-based masking and spatial separation using linear beamforming with a multichannel Wiener filter (MCWF), performing up to three applications of the neural separation network: separate, beamform, separate, beamform, and separate. By doing so, linear beamforming is effectively driven by DNN-based masking. This is inspired by the single-channel sequential network of [14], which we use as a baseline, and by the findings that better beamforming results can be used as extra features to improve spectral masking and vice versa. This sequential approach is related to iterative neural beamforming with postfiltering [15], which uses the same DNN repeatedly with only the beamformed signal as input for speech enhancement; in contrast, we train a different neural network for each stage, concatenating the mixture signal with the beamformed signals, applied to both speech enhancement and separation. For beamforming, we consider both time-invariant and time-varying ways of calculating covariance matrices to improve spatial separation. We also explore the effectiveness by incorporating multi-frame context during beamforming. Evaluation results on four challenging sound separation tasks demonstrate the effectiveness of the proposed algorithms. Also, our best three-speaker separation model achieves significantly improved word error rate on the LibriCSS dataset compared to a neural mask beamforming baseline separation system.

## 2. Contributions

The model we introduce has similarities with earlier mask-based beamforming models but has the following novel aspects which end up improving the performance significantly.

- We perform multiple repetitions of mask-based beamforming where each sequential application of a neural network

has its own parameters optimized separately.

- We use different window and hop sizes for the mask-prediction network and the beamformer, in contrast to previous works where the same STFT parameters were used for both. Our networks predict time-domain waveforms, and we take another independent STFT to perform beamforming. We achieve best results using a smaller window size in the mask network and a larger one for the beamformer.

- We use a state-of-the-art TDCN++ network [14] for mask-prediction as well as mixture consistency projection [16] to improve separation performance.

- We use a stabilized SNR loss function for training the mask-prediction neural network.

- Our best performing beamformer is a time-invariant multi-frame multichannel Wiener filter that improves the results significantly compared to previously-used single-frame beamformers.

- Our model is completely independent of microphone and room geometry, and we show that our three-speaker separation model works well on a mismatched real data set with different number of microphones and an unseen room geometry.

## 3. Methods

Assume an $M$-channel time-domain signal consisting of $S$ sources, $\mathbf{y}[n] = \sum_{s=1}^{S} \mathbf{x}^{(s)}[n]$, recorded in a reverberant environment. The short-time Fourier transform (STFT) of this multichannel signal can be written as $\mathbf{Y}_{t,f} = \sum_{s=1}^{S} \mathbf{X}_{t,f}^{(s)}$, where $\mathbf{Y}_{t,f}$ and $\mathbf{X}_{t,f}^{(s)} \in \mathbb{C}^M$ respectively represent the mixture and the reverberant image of source $s$ at time $t$ and frequency $f$. Our study proposes multiple algorithms to recover the constituent reverberant sources $X_{\text{ref}}^{(s)}$ from a reverberant mixture $Y_{\text{ref}}$ received by a reference microphone, with or without leveraging spatial information contained in $\mathbf{Y}$. We assume offline processing and non-moving sources throughout each utterance.

Figure 1 illustrates our proposed system. Each spectral masking stage uses an improved time-domain dilated convolutional neural network (TDCN++) [14]. The first stage performs single-channel processing to estimate each source via T-F masking. The estimated sources are then used to compute statistics for time-invariant or time-varying beamforming. The next masking stage combines spectral and spatial information by taking in the mixture and beamformed results for post-filtering. This sequence is then repeated several times.

As shown in Figure 1, we train through multiple iSTFT/STFT projection layers. These layers are helpful as they can effectively address the well-known phase inconsistency problem, a common issue of magnitude-domain masking [16, 17]. In addition, our masking networks operate at a typical 32 ms window size, but our system can use a larger window size for beamforming. This way, beamforming can be performed at a higher frequency resolution and produce finer separation. The iSTFT/STFT pairs are necessary here to change the window size back and forth during sequential processing. This strategy dramatically improves time-invariant MCWF in our experiments.

### 3.1. Spectral mask estimation for sound separation

For monaural speech enhancement and speaker separation, we use TDCN++ based T-F masking (see [18] for an overview) to produce source estimates $\hat{X}_{\text{MN}i}^{(s)} = \hat{A}_i^{(s)} \odot Y_{\text{ref}}$, where $\odot$
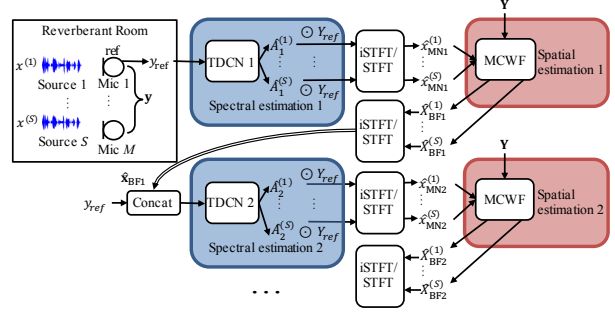


Figure 1: *System overview.*

denotes element-wise multiplication and $\hat{A}_i^{(s)}$ is the mask estimate produced by the $i$th TDCN++. Note that $i \in \{1, 2, ..., I\}$, meaning that there are $I$ (set to 3 in this study) stages in the sequence. The loss function at each stage maximizes a stabilized SNR in the time domain

$$\mathcal{L}_i = \min_{\pi \in \Pi} \sum_{s=1}^{S} -\text{SNR}_{\tau,\epsilon}\Big(\text{iSTFT}(\hat{X}_{\text{MN}i}^{(\pi(s))}), x_{\text{ref}}^{(s)}\Big), \quad (1)$$

where

$$\text{SNR}_{\tau,\epsilon}(\hat{\mathbf{x}}, \mathbf{x}) = 10 \log_{10}\left(\frac{\|\mathbf{x}\|^2}{\|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \tau\|\mathbf{x}\|^2 + \epsilon}\right), \quad (2)$$

$\Pi$ is the set of all permutations over $S$ sources, and $\pi(s)$ is the permuted source index $s$ under permutation $\pi$. The parameter $\tau$ acts as a soft threshold limiting the maximum SNR that can dampen the effect on the total loss from examples that are already well-separated. We use permutation invariance for speech separation, but not for speech enhancement. For $i > 1$, the network input is the concatenation of the mixture magnitude STFT features with those of all beamformed source estimates, $\hat{X}_{\text{BF}i-1}^{(s)}$.

### 3.2. Multi-frame multichannel Wiener filter

Inspired by the success of convolutional beamformers [19], we introduce multi-frame MCWF beamforming and investigate various context sizes. The rationale is that, by stacking multiple frames, the beamformer can have more contextual information and degrees of freedom for better noise suppression.

We first define a context-expanded observed signal $\bar{\mathbf{Y}}_{t,f} = \left[\mathbf{Y}_{t-a,f}^T, \ldots, \mathbf{Y}_{t,f}^T, \ldots, \mathbf{Y}_{t+b,f}^T\right]^T \in \mathbb{C}^{cM}$ which is a flattened complex vector including multiple frames around a T-F unit, where $a$ is the left and $b$ the right context size in frames, and $c = a + b + 1$. We treat all the contextual T-F units as if they are additional microphones in the subsequent beamforming formulations. In each stage $i$, estimated sources $\hat{X}_{\text{MN}i}^{(s)}$ from the TDCN++ are used to compute the spatial covariance of each source for a time-invariant MCWF (TI-MCWF)

$$\hat{\mathbf{w}}_{i,f}^{(s)} = (\hat{\boldsymbol{\Phi}}_f^{(y)})^{-1} \hat{\boldsymbol{\Phi}}_{i,f}^{(s)} \mathbf{u}_{\text{ref}}, \quad (3)$$

where $\mathbf{u}_{\text{ref}}$ is a one-hot vector with the coefficient corresponding to the reference microphone at the center frame set to one, the multi-frame mixture covariance matrix is estimated as

$$\hat{\boldsymbol{\Phi}}_f^{(y)} = \frac{1}{T} \sum_{t=1}^{T} \bar{\mathbf{Y}}_{t,f} \bar{\mathbf{Y}}_{t,f}^H, \quad (4)$$

and $\hat{\boldsymbol{\Phi}}_{i,f}^{(s)}$ is the source covariance matrix computed as

$$\hat{\boldsymbol{\Phi}}_{i,f}^{(s)} = \frac{1}{T} \sum_{t=1}^{T} \hat{A}_{i,t,f}^{(s)} \bar{\mathbf{Y}}_{t,f} \bar{\mathbf{Y}}_{t,f}^{H}, \tag{5}$$

$$\hat{A}_{i,t,f}^{(s)} = \frac{|\hat{X}_{\mathrm{MN}i,t,f}^{(s)}|^2}{\sum_{s'=1}^{S} |\hat{X}_{\mathrm{MN}i,t,f}^{(s')}|^2}. \tag{6}$$

This approach follows recent developments in neural beamforming [3, 20, 2] and straightforwardly applies them to a multi-frame setup. The idea is to use T-F units dominated by source $s$ to compute its covariance matrix for beamforming. Here the Wiener-like mask $\hat{A}^{(s)}$, which can be derived based on a different window size, is recomputed in an alternate STFT domain from the reconstructed time-domain signal for source $s$ from the masking network. For convenience, the mask is considered the same across microphones, which is a reasonable approximation for compact arrays in far-field conditions. The beamforming result for source $s$ in stage $i$ is computed as

$$\hat{X}_{\mathrm{BF}i,t,f}^{(s)} = (\hat{\mathbf{w}}_{i,f}^{(s)})^{H} \bar{\mathbf{Y}}_{t,f}. \tag{7}$$

We also experimented with MVDR and MPDR beamformers [21, 22], but they did not perform as well as MCWF in terms of SI-SNR. This paper hence only reports results with MCWF.

### 3.3. Time-varying beamforming for spatial estimation

A TI-MCWF has limited power for separation, as it is only a linear time-invariant filter per-frequency. To obtain time-varying behavior, we experiment with a block-based approach, where we calculate TI-MCWF beamformers in half-overlapping blocks of frames with some windowing. We use windowed signals to calculate spatial covariance matrices and perform overlap-add for post-windowed beamformed signals. The Vorbis window [23] is used for this processing.

A frame-level way of computing a time-varying covariance matrix for each source is to factorize it as a product of a time-varying power spectral density (PSD) and a time-invariant coherence matrix [24, 25, 26]. The rationale is that for a non-moving source, its coherence matrix is time-invariant assuming that the beamforming STFT window is long enough to capture most of the reverberation. Unlike conventional methods, which typically use maximum likelihood estimation or non-negative matrix factorization to estimate the PSD and spatial coherence [24, 26], the proposed algorithm leverages estimated source signals produced by neural networks to compute these statistics. Mathematically,

$$\hat{\boldsymbol{\Phi}}_{i,t,f}^{(s)} = |\hat{X}_{\mathrm{MN}i,t,f}^{(s)}|^2 \hat{\boldsymbol{\Psi}}_{i,t,f}^{(s)} / \hat{D}_{i,t,f}, \tag{8}$$

where $|\hat{X}_{\mathrm{MN}i,t,f}^{(s)}|^2$ is the PSD estimate, $\hat{\boldsymbol{\Psi}}_{i,t,f}^{(s)}$ can be either $\hat{\boldsymbol{\Phi}}_{i,f}^{(s)}$ computed over all the frames in an utterance as in (5) for a time-invariant covariance matrix, or it could be a block-based one calculated over the frames in a block. $\hat{D}_{i,t,f} = \hat{d}_{i,t,f} \hat{d}_{i,t,f}^{T}$ with $\hat{d}_{i,t,f} = \mathrm{diag}(\hat{\boldsymbol{\Psi}}_{i,t,f}^{(s)})^{1/2}$ normalizes the spatial component to have a unit diagonal. In far-field conditions where level differences are negligible, $\hat{D}_{i,t,f} \approx (\hat{\boldsymbol{\Psi}}_{i,t,f}^{(s)})_{m,m} \mathbf{1}\mathbf{1}^{T}$ for a microphone index $m$. A time-varying factorized (TVF) MCWF is computed as

$$\hat{\mathbf{w}}_{i,t,f}^{(s)} = (\hat{\boldsymbol{\Phi}}_{i,t,f}^{(y)})^{-1} \hat{\boldsymbol{\Phi}}_{i,t,f}^{(s)} \mathbf{u}_{\mathrm{ref}}, \tag{9}$$

where $\hat{\boldsymbol{\Phi}}_{i,t,f}^{(y)} = \sum_{s'=1}^{S} \hat{\boldsymbol{\Phi}}_{i,t,f}^{(s')}$, and the beamformed result is

$$\hat{X}_{\mathrm{BF}i,t,f}^{(s)} = (\hat{\mathbf{w}}_{i,t,f}^{(s)})^{H} \bar{\mathbf{Y}}_{t,f}. \tag{10}$$

## 4. Data and models

### 4.1. Datasets

We use room impulse responses (RIRs) generated by an image-method room simulator with frequency-dependent wall filters. For each example, the RIR is created by sampling random locations for a cube-shaped microphone array and all sources within a room defined using a random size: width from 3 to 7 m, length from 4 to 8 m, and height between 2.13 and 3.05 m. The sides of the cube-shaped array was 20 cm long. During RIR generation, all source "image" locations are randomly perturbed by up to 8 cm in each direction to avoid the "sweeping echo" effect [27]. We generate 140,000 training, 20,000 validation and 20,000 test rooms which are used to generate train, validation and test data. Clean speech is from Libri-Light [28] and LibriTTS [29], and non-speech sounds are from freesound.org. We filtered out artificial sounds (such as synthesizer noises) based on user-annotated tags and used a sound classification network trained on AudioSet [30] to avoid clips likely containing speech. During training, sources are reverberated and mixed on the fly, and the validation and test sets consist of about 10 hours of mixture data each. Recipes for these datasets will be publicly released in the near future. We validate the proposed algorithms on 1, 2, 4, and 8-microphone setups.

Using this source data, the proposed models are evaluated on both speech separation and speech recognition. For speech separation evaluation, we construct three tasks: two-speaker separation, three-speaker separation and speech enhancement. For the speech enhancement task, a speech source is mixed with three directional noise sources, and the goal is to separate the speech from the noises. For each task, a random speech clip from clean source data is selected, and then each of the other sources is scaled to an SNR randomly drawn from $\mathcal{N}(0, 7)$ dB with respect to the speech clip. To better compare with previous arts, we used an additional two-speaker separation evaluation dataset introduced in [8], which is a multichannel reverberated version of WSJ0-2mix database simulated using a room simulator with random room configurations and microphone positions.

Besides these separation tasks, we evaluate our three-speaker separation model in terms of its ASR performance on the LibriCSS dataset which is a real meeting-like overlapping speech dataset [31]. This dataset has been collected by playing LibriSpeech utterances from loudspeakers and recording them in a room. Each loudspeaker takes a role of a single speaker and reads only utterances from that speaker [31].

### 4.2. Why use simulated data?

Our algorithms require multichannel data and we would like to make sure that we see a huge number of possible mixing configurations during training. Thus, we use simulation to generate training data with a multitude of possible source locations and microphone positions in random rooms. We have to use this type of simulated data since existing "real" multichannel recordings and room impulse response databases are nowhere close to the size that's needed to train a good separation model that generalizes to unseen conditions at test time. For example, the ACE database [32] and BUT ReverbDB [33] provide real room impulse responses, but they are limited in the number of rooms, possible source locations, and range of microphone geometries: ACE has 7 rooms, 1 source location per room, and one geometry each for 2, 3, 4, 5, 8, and 32 mics; BUT ReverbDB has 8 rooms, 2-11 source positions per room with a spherical 8-mic array in addition to 23 single microphones

with ad-hoc placements. Thus, these real RIR databases are not extensive enough for training generalizable multichannel separation models that can handle arbitrary numbers of sources in arbitrary rooms. For evaluation of separation systems, simulated test sets provide ground-truth source references that can be used to measure performance in terms of SI-SNR. Due to scarcity of common simulated multichannel evaluation sets, we generate our own development and evaluation sets. We plan to release our simulated databases to allow wide use of the academic community and serve as a benchmark for multichannel separation tasks. Real multichannel datasets with transcription but without ground-truth source reference signals are available. To provide evaluation on real data in this paper, we evaluate our best model with ASR word error rates on the LibriCSS dataset, which has real acoustic mixing and known transcriptions.

### 4.3. Networks

The network architecture of the TDCN++ networks is similar to the recently proposed Conv-TasNet [34] and includes a few improvements introduced in [14]. It consists of 4 repeats of 8 layers of convolutional blocks. Each block consists of a dilated separable convolution with feature-wise global layer normalization and a residual connection, where the dilation factor for the $k$th block is $2^k$. In contrast to Conv-TasNet, we utilize STFT basis with 32 ms windows rather than a learned basis with a very small window size, as initial results showed that the former leads to better performance. This is likely because an STFT with a larger window can better deal with room reverberation. The hop size is 8 ms. The sampling rate is 16 kHz. A 512-point FFT is used to extract 257-dimensional magnitude features for mask estimation. We use $\tau = 10^{-3}$ and $\epsilon = 10^{-8}$ with the soft-thresholded and stabilized negative SNR loss in (1). We intend to open-source our implementations of TDCN++ and sequential multichannel models. SI-SNR improvement (SI-SNRi) [35] over unprocessed speech is utilized as the evaluation metric. We also report *differential* word error rates (dWER) by assuming the speech recognition hypothesis of clean speech as the ground truth, and calculating the WER between the recognition output of enhanced/separated speech and that of clean speech. The ASR model we used for obtaining differential WERs on the simulated datasets is an attention-based encoder-decoder model with 16k word-piece output units trained on 960 hours of Librispeech data [36]. For LibriCSS experiments, we used a kaldi based hybrid-HMM ASR model and an end-to-end ASR model based on ESPNet.

As a single-channel baseline, we consider a sequential TDCN++ network [14], where no spatial information is used. This network consists of three masking networks. For the second and third networks, the separated time-domain outputs of the previous network are concatenated with the time-domain mixture signal as the input features to produce separated estimates. We report performance for the output of each stage. This model is trained with the negative stabilized SNR loss in (1) on the separated waveforms of all three stages.

## 5. Results

### 5.1. Results on simulated data

Figures 2 and 3 show the performance on the validation set of beamforming methods driven by a single-channel neural network under different conditions: using either TI or TV covariance estimation, for various block sizes, with 2, 4, or 8 microphones, for each of the four tasks, and with single (Figure 2) or

multiple context frames (Figure 3). We only display results for the best beamforming parameters over all tasks and numbers of microphones. We considered beamforming window sizes of 32, 64, 128 and 256 ms with half-sized hops. Frame context size were swept in powers of 2, and we chose the most promising frame context for each window. The best multi-frame TI result is achieved with window size 64 ms and 4 context frames (64 x 4), and the best multi-frame TV result uses window size 128 ms and 2 context frames (128 x 2). Our frame contexts are centered around the current frame, where for even context sizes, left context $a$ is one larger than the right context $b$.

In the following discussion, we use MN$i$ to refer to the mask network output after step $i$ and BF$i$ to refer to the beamformer applied after step $i$, as consistent with the subscripts in equations in Section 3. For example, MN3 is the output of the separation model, and BF2 is the beamformed result achieved one step before that.

Figures 4 and 5 visualize the SI-SNRi performance on the validation and test sets of our best sequential neural beamforming models versus iteration, where the best multi-frame beamforming parameters are chosen from Figure 3. These plots also display the performance of single-channel baselines, including a single-channel iterative network and an oracle binary mask (OBM) both using the same STFT parameters. Notice that performance generally improves monotonically with iterations, with the outputs of the neural networks achieving better SI-SNRi compared to the beamforming outputs. Also, despite performing worse than TV on their own, TI beamforming performs best when used in a sequential setup. For all tasks, using more than one microphone improves performance.

The speech enhancement and speech separation datasets that we constructed from Libri-Light [28] and `freesound.org` have less overlap between sources compared to the WSJ0 speech separation dataset. To fairly compare results, figures 4 and 5 shows SI-SNRi computed only on fully-overlapping segments with darker colored sub-bars. Notice that the results are more comparable between our speech separation and WSJ0 speech separation in overlapped regions. Also, for speech enhancement and 2 speaker separation, our sequential neural beamformer exceeds the performance of the oracle binary mask (not shown) in fully overlapped segments.

Figures 6 and 7 display dWER for the validation and test sets. Our sequential neural beamforming models significantly decrease dWER, especially when more microphones are used. When using TI beamforming with 8 microphones, BF2 achieves the best dWER as opposed to MN3 since ASR models tend to work better with linear time-invariant processed signals. For two microphones though, MN3 output is the best, likely because two-microphone beamforming cannot achieve sufficient spatial separation. For the Libri-Light+Freesound speech enhancement and speech separation tasks, the best-performing outputs of our model achieve comparable or slightly better dWER than an oracle binary mask.

Table 1 presents the results of our best eight-microphone system as compared to a single-channel baseline, a multichannel baseline, and an oracle binary mask. We point out that our baselines are strong ones since we use a state-of-the-art neural network architecture and an improved SNR loss function. Also, for the beamformer, we use an optimal 128 ms window size which is typically not the case. Our methods obtain significantly better SI-SNR and dWER against these strong baselines.
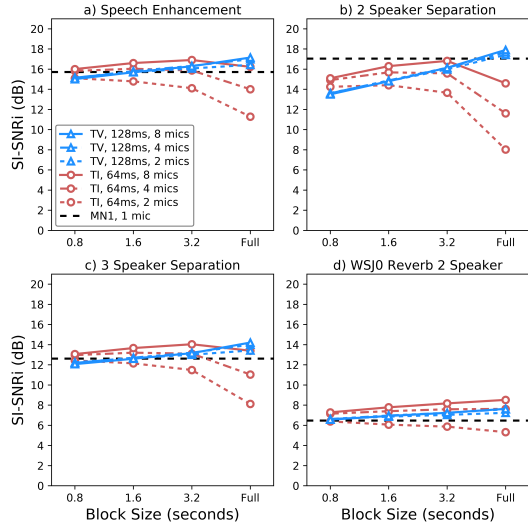
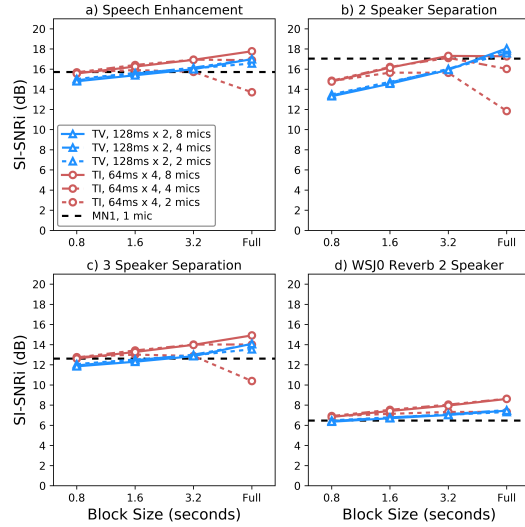Figure 2: *Single-frame beamforming versus block size.*



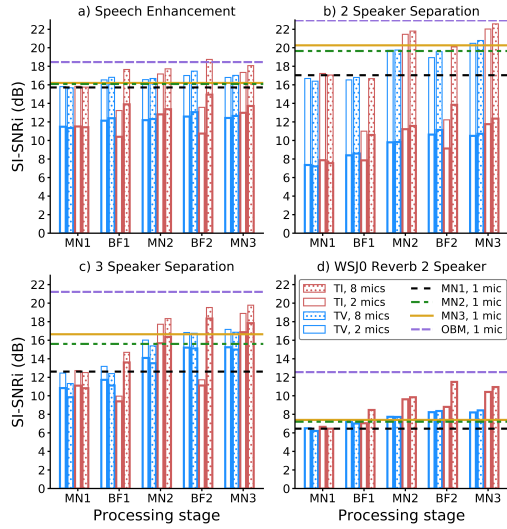Figure 3: *Multi-frame beamforming versus block size.*



Figure 4: *SI-SNRi of sequential neural BF (val). Dark bars indicate score of overlapping regions.*
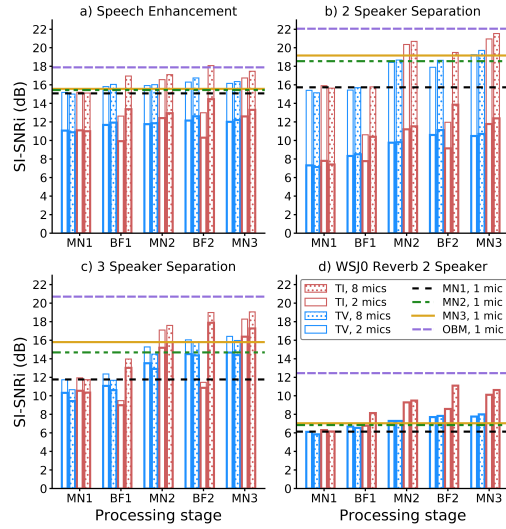


Figure 5: *SI-SNRi of sequential neural BF (test). Dark bars indicate score of overlapping regions.*

| Method | SI-SNRi (dB) | | | | | | | | dWER (%) | | | | | | | |
| | Speech Enhancement | | 2 Speaker Separation | | 3 Speaker Separation | | WSJ0 2 Spk. Separation | | Speech Enhancement | | 2 Speaker Separation | | 3 Speaker Separation | | WSJ0 2 Spk. Separation | |
| | val | tst | val | tst | val | tst | val | tst | val | tst | val | tst | val | tst | val | tst |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Noisy | - | - | - | - | - | - | - | - | 71.5 | 67.1 | 99.8 | 98.6 | 144.6 | 143.4 | 111.1 | 111.7 |
| MN3, 1 mic | 16.2 | 15.6 | 20.3 | 19.2 | 16.7 | 15.8 | 7.4 | 7.0 | 44.6 | 40.5 | 29.3 | 30.1 | 50.6 | 52.3 | 68.8 | 69.3 |
| BF1, 8 mic, TI 128ms x 1 | 15.7 | 15.2 | 16.1 | 15.5 | 14.4 | 13.7 | 8.6 | 8.3 | 31.1 | 26.6 | 27.8 | 27.5 | 52.0 | 52.6 | 57.1 | 58.2 |
| BF2, 8 mic, TI 64ms x 4 | **18.7\*** | **18.1\*** | 20.1 | 19.5 | 19.5 | **19.0** | **11.5** | **11.1** | **25.0\*** | **20.9\*** | **17.6\*** | **17.0\*** | 30.3 | 29.5 | 40.1 | **40.6** |
| MN3, 8 mic, TI 64ms x 4 | 18.1 | 17.5 | **22.6** | **21.5** | **19.8** | **19.0** | 10.9 | 10.6 | 26.1 | 21.9 | 18.5 | 18.4 | 32.5 | 32.7 | 43.3 | 42.8 |
| OBM, 1 mic, **oracle** mask | 18.5 | 17.9 | 23.0 | 22.1 | 21.2 | 20.7 | 12.6 | 12.4 | 25.8 | 22.8 | 18.5 | 18.6 | 28.3 | 28.7 | 35.2 | 33.6 |

Table 1: *SI-SNRi and differential WER (dWER) results of sequential neural beamforming on the validation and test data as compared to strong single-channel and 8-mic multichannel baselines, as well as an oracle binary mask (OBM) for four different tasks. \* indicates performance better than the oracle binary mask.*

## 5.2. Results on LibriCSS

We ran our three-speaker separation model on the LibriCSS evaluation dataset which contains 54 different 10 minute meetings with varying amounts of overlap recorded with a circular 7-microphone array [31]. We evaluate our separation model within a separation-diarization-recognition pipeline which is described in more detail in [37]. The model is applied in 8 second overlapping blocks with a 4 second shift. 7 microphone data is padded with another channel obtained by shifting the first microphone signal by one sample and adding white Gaussian noise with variance 1e-6. We used the BF2 output of the three-speaker separation model since it achieved better dWER in our experiments on simulated test data as can be seen in Ta-
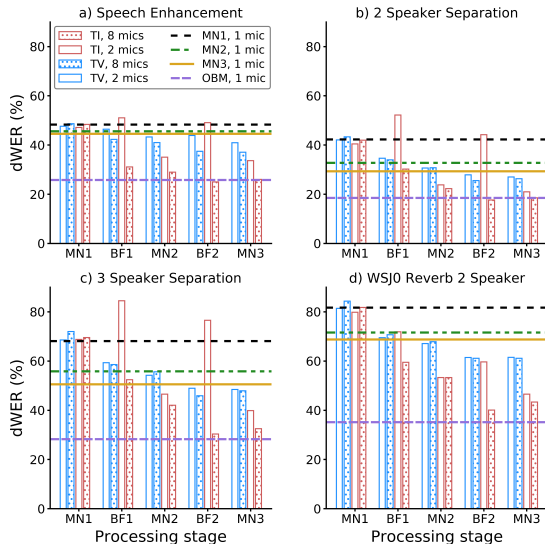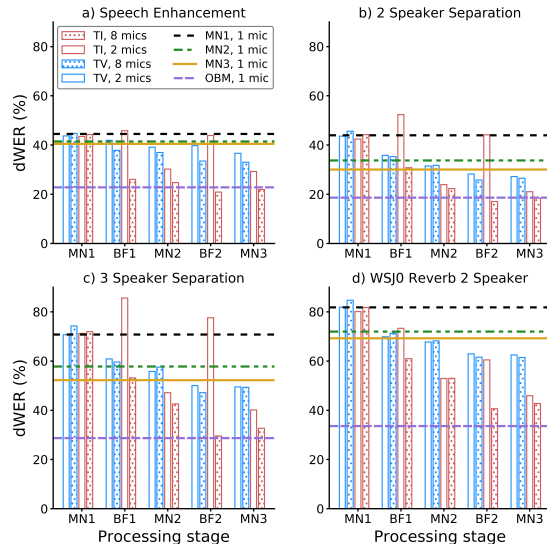
Figure 6: *Differential WER versus iteration (val).*



Figure 7: *Differential WER versus iteration (test).*

Table 2: *Performance of separation methods on LibriCSS eval set in terms of the resulting downstream diarization error rate (DER) (using spectral clustering), cpWER (using a hybrid HMM-DNN model) and cpWER (using an E2E model) results. Separation performance in terms of signal to distortion ratio (SDR) is reported on a different simulated LibriCSS-like eval set. For comparison, we also show results obtained on a "no separation" baseline.*

| Method | SDR (dB) | DER (%) | HMM-DNN cpWER (%) | E2E cpWER (%) |
|---|---|---|---|---|
| No separation | - | 18.28 | 31.04 | 27.11 |
| Mask-based MVDR [31] | 5.8 | **13.86** | 22.75 | 13.37 |
| Proposed: sequential multi-frame (BF2) | **14.1** | 14.07 | **19.28** | **12.70** |

ble 1. After separating each block into three speaker tracks, we stitch together the separated block-length tracks into meeting-length tracks. We use magnitude STFT domain mean-squared distance between common parts of neighboring blocks to find the best permutation between them. Diarization was done using x-vector based segment clustering from multiple separated tracks with some post-refinement [37]. Each diarized segment is recognized using either a hybrid HMM-DNN ASR model or an end-to-end (E2E) ASR model.

The results, taken from [37] are shown in Table 2. We compare with a baseline mask-based MVDR separation method that uses bidirectional LSTM layers [31, 37]. Separation performance was evaluated on a simulated test set since reference signals are required [37]. Separated tracks are first mapped to $N$-speaker tracks where $N$ is the number of participating speakers in a meeting, using a method we call "oracle track mapping" [37]. Average meeting-level SDR [38] is reported. Multi-frame MCWF beamformer (BF2) achieved a much better SDR value as compared to mask-based MVDR, however this may be expected since the MVDR beamformer does not attempt to reconstruct the target signal at the reference microphone directly.

On LibriCSS, diarization error rate (DER) was close between two separation methods. Our model achieved a concatenated minimum-permutation WER (cpWER) [39] of 19.28% on the LibriCSS eval set using a hybrid HMM-DNN model better than the baseline system. When using a superior E2E ASR model, our separation model achieved a cpWER of 12.70% as compared to a baseline of 13.37%. Note that our model is trained using a completely different microphone geometry and with mismatched data in terms of overlap amount. However,

our model has the advantage of training from on-the-fly mixtures and Libri-Light database which is quite large. We can see that the model generalizes to unseen conditions since it works well on data which contains a different microphone geometry, unseen source locations, different overlap amounts, and unseen reverberations with real recording conditions.

# 6. Conclusions

We have explored an alternating strategy between spectral estimation using a mask-based network and spatial estimation using beamformers. For spatial estimation, we introduced multi-frame beamforming and compared multiple ways of computing covariance matrices for time-invariant and time-varying beamforming. Evaluation results on four sound separation tasks suggest that, when combined with neural network based mask estimation, time-invariant multi-frame beamforming with a reasonably large window and context size produces the best separation performance for non-moving sources. Our best three-stage method demonstrates an average improvement of 2.75 dB in SI-SNR and an absolute reduction of 14.2% in dWER over several strong and representative baselines, across four challenging reverberant speech enhancement and separation tasks. Our three speaker separation model was used to separate tracks from LibriCSS evaluation dataset and ended up improving the constrained permutation word error rate as compared to a mask-based MVDR beamformer baseline.

# 7. References

[1] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural network based spectral mask estimation for acoustic beamforming," in *Proc. ICASSP*, 2016.

[2] H. Erdogan, J. R. Hershey, S. Watanabe, M. I. Mandel, and J. Le Roux, "Improved MVDR beamforming using single-channel mask prediction networks," in *Proc. Interspeech*, 2016.

[3] J. Barker, R. Marxer, E. Vincent, and S. Watanabe, "The third 'CHiME' speech separation and recognition challenge: Analysis and outcomes," *Computer Speech and Language*, vol. 46, 2017.

[4] C. Boeddeker, H. Erdogan, T. Yoshioka, and R. Haeb-Umbach, "Exploring practical aspects of neural mask-based beamforming for far-field speech recognition," in *Proc. ICASSP*, 2018.

[5] T. Yoshioka, Z. Chen, C. Liu, X. Xiao, H. Erdogan, and D. Dimitriadis, "Low-latency speaker-independent continuous speech separation," in *Proc. ICASSP*, 2019.

[6] Y. Kubo, T. Nakatani, M. Delcroix, K. Kinoshita, and S. Araki, "Mask-based MVDR beamformer for noisy multisource environments: Introduction of time-varying spatial covariance model," in *Proc. ICASSP*, 2019.

[7] T. Yoshioka, H. Erdogan, Z. Chen, and F. Alleva, "Multi-microphone neural speech separation for far-field multi-talker speech recognition," in *Proc. ICASSP*, 2018.

[8] Z.-Q. Wang, J. Le Roux, and J. R. Hershey, "Multi-channel deep clustering: Discriminative spectral and spatial embeddings for speaker-independent speech separation," in *Proc. ICASSP*, 2018.

[9] Z.-Q. Wang and D. Wang, "Combining spectral and spatial features for deep learning based blind speaker separation," *IEEE/ACM TASLP*, vol. 27, no. 2, 2018.

[10] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, "Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation," *ACM Transactions on Graphics*, 2018.

[11] Z. Chen, X. Xiao, T. Yoshioka, H. Erdogan, J. Li, and Y. Gong, "Multi-channel overlapped speech recognition with location guided speech extraction network," in *Proc. SLT*, 2018.

[12] M. Delcroix, K. Zmolikova, K. Kinoshita, A. Ogawa, and T. Nakatani, "Single channel target speaker extraction and recognition with speaker beam," in *Proc. ICASSP*, 2018.

[13] Q. Wang, H. Muckenhirn, K. Wilson, P. Sridhar, Z. Wu, J. R. Hershey, R. A. Saurous, R. J. Weiss, Y. Jia, and I. L. Moreno, "VoiceFilter: Targeted voice separation by speaker-conditioned spectrogram masking," in *Proc. Interspeech*, 2019.

[14] I. Kavalerov, S. Wisdom, H. Erdogan, B. Patton, K. Wilson, J. Le Roux, and J. R. Hershey, "Universal sound separation," *Proc. WASPAA*, 2019.

[15] X. Zhang, Z.-Q. Wang, and D. Wang, "A speech enhancement algorithm by iterating single-and multi-microphone processing and its application to robust asr," in *Proc. ICASSP*. IEEE, 2017.

[16] S. Wisdom, J. R. Hershey, K. Wilson, J. Thorpe, M. Chinen, B. Patton, and R. A. Saurous, "Differentiable consistency constraints for improved deep speech enhancement," in *Proc. ICASSP*, 2019.

[17] Z.-Q. Wang, K. Tan, and D. Wang, "Deep learning based phase reconstruction for speaker separation: A trigonometric perspective," in *Proc. ICASSP*, vol. 2019, 2019.

[18] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM TASLP*, vol. 26, Aug. 2018.

[19] T. Nakatani and K. Kinoshita, "A unified convolutional beamformer for simultaneous denoising and dereverberation," *IEEE Signal Processing Letters*, vol. 26, no. 6, 2019.

[20] J. Heymann, L. Drude, A. Chinaev, and R. Haeb-Umbach, "BLSTM supported GEV beamformer front-end for the 3rd CHiME challenge," in *ASRU*, 2015.

[21] M. Souden, J. Benesty, and S. Affes, "On optimal frequency-domain multichannel linear filtering for noise reduction," *IEEE Transactions on audio, speech, and language processing*, vol. 18, no. 2, 2009.

[22] L. Ehrenberg, S. Gannot, A. Leshem, and E. Zehavi, "Sensitivity analysis of mvdr and mpdr beamformers," in *IEEE Convention of Electrical and Electronics Engineers in Israel*, 2010.

[23] J.-M. Valin, G. Maxwell, T. B. Terriberry, and K. Vos, "High-quality, low-delay music coding in the Opus codec," *Proc. of AES Convention*, 2016.

[24] N. Duong, E. Vincent, and R. Gribonval, "Under-determined reverberant audio source separation using a full-rank spatial covariance model," *IEEE TASLP*, vol. 18, no. 7, 2010.

[25] T. Higuchi, N. Ito, T. Yoshioka, and T. Nakatani, "Robust MVDR beamforming using time-frequency masks for online/offline ASR in noise," in *Proc. ICASSP*, 2016.

[26] K. Shimada, Y. Bando, M. Mimura, K. Itoyama, K. Yoshii, and T. Kawahara, "Unsupervised speech enhancement based on multi-channel NMF-informed beamforming for noise-robust automatic speech recognition," *IEEE/ACM TASLP*, vol. 27, no. 5, 2019.

[27] E. De Sena, N. Antonello, M. Moonen, and T. Van Waterschoot, "On the modeling of rectangular geometries in room acoustic simulations," *IEEE/ACM TASLP*, vol. 23, no. 4, 2015.

[28] J. Kahn, M. Rivière, W. Zheng, E. Kharitonov, Q. Xu, P.-E. Mazaré, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen *et al.*, "Libri-Light: A benchmark for ASR with limited or no supervision," in *Proc. ICASSP*, 2020.

[29] H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, and Y. Wu, "LibriTTS: A corpus derived from LibriSpeech for text-to-speech," in *Proc. Interspeech*, 2019.

[30] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in *Proc. ICASSP*, 2017.

[31] Z. Chen, T. Yoshioka, L. Lu, T. Zhou, Z. Meng, Y. Luo, J. Wu, X. Xiao, and J. Li, "Continuous speech separation: dataset and analysis," 2020.

[32] J. Eaton, N. D. Gaubitch, A. H. Moore, and P. A. Naylor, "The ace challenge—corpus description and performance evaluation," in *2015 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2015.

[33] I. Szoke, M. Skacel, L. Mosner, J. Paliesek, and J. Cernocky, "Building and evaluation of a real room impulse response dataset," *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 4, Aug 2019.

[34] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing ideal time-frequency magnitude masking for speech separation," *IEEE/ACM TASLP*, vol. 27, no. 8, 2019.

[35] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR–half-baked or well done?" in *Proc. ICASSP*, 2019.

[36] K. Irie, R. Prabhavalkar, A. Kannan, A. Bruguier, D. Rybach, and P. Nguyen, "On the choice of modeling unit for sequence-to-sequence speech recognition," in *Proc. Interspeech*, 2019.

[37] D. Raj, P. Denisov, Z. Chen, H. Erdogan, Z. Huang, M. He, S. Watanabe, J. Du, T. Yoshioka, Y. Luo, N. Kanda, J. Li, S. Wisdom, and J. Hershey, "Integration of speech separation, diarization, and recognition for multi-speaker meetings: System description, comparison, and analysis," in *submitted to 2020 IEEE Workshop on Spoken Language Technology (SLT)*. IEEE, 2020.

[38] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE transactions on audio, speech, and language processing*, vol. 14, no. 4, 2006.

[39] S. Watanabe, M. Mandel, J. Barker, and E. Vincent, "CHiME-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings," *ArXiv*, vol. abs/2004.09249, 2020.