

# DEEP LEARNING ON COMPUTATIONAL ACCELERATORS 236605

## PAPER IMPLEMENTATION

### Monoaural Audio Source Separation Using Deep Convolutional Neural Networks

by P. Chandna, M. Miron, J. Janer, E. Gómez

Danny Priymak 307003434

Gloria Groysman 318372729

Technion IIT, spring 2019

## 1 Introduction

Source separation is the separation of a set of source signals from a set of mixed signals, without the aid of information (or with very little information) about the source signals or the mixing process. It is most commonly applied in digital signal processing and involves the analysis of mixtures of signals; the objective is to recover the original component signals from a mixture signal. Audio source separation has drawn the attention of researchers, with approaches varying from using timbre models [2], to those exploiting the repetitive nature of music [3]. While being an interesting problem in itself, the separation of sources from a mixture can serve as a intermediary step for other tasks such as automatic speech recognition and fundamental frequency estimation. Some applications, such as speech enhancement for cochlear implant users, require low-latency processing. The original paper we chose to focus on, by Chandna et al. [1] has tackled the audio signal source separation case, and has introduced a low-latency monaural source separation framework using a convolutional neural network (CNN). The authors used a CNN to estimate time-frequency soft masks, generated per-instrument, which in turn are applied to musical tracks in order to isolate the corresponding instrument's estimated, isolated audio. The training, as well as an evaluation of the neural network's performance was performed over a dataset containing fully mixed audio music tracks of four fixed instruments: vocals, drums, and bass, as well as a fourth instrument category representing all other remaining instruments, which vary from track to track. The authors' proposed architecture has been compared to a Multilayer Perceptron (MLP) architecture, and achieved on-par results in conjunction with a significant improvement in processing time. Our work intends to implement this architecture and provide possible enhancements.

## 2 Related Work

Several approaches have been proposed for the solution of this problem but development is currently still very much in progress. A more successful classical approach that does not involve machine learning concepts is non-negative matrix factorization (NMF). Additional popular approaches are principal components analysis and independent component analysis, which work well when there are no delays or echoes present; that is, the problem is simplified a great deal.

Approaches directly using deep neural networks for separation have been proposed as well. Nugraha et al. [4] adapt deep neural networks for multichannel source separation, using both phase and magnitude information. As for monaural source separation, Huang et al. [5] propose a method that utilizes a deep neural network, taking a single frame of the magnitude spectrogram of a mixture as an input feature to learn single-frame timbre features for each source. Temporal evolution is then modeled using a recurrent layer. Uhlich et al. [6] propose another method which takes multiple frames of the magnitude spectrogram of a mixture as input and consists of only fully connected layers. This method models timbre features across multiple time frames. While these approaches work well, they do not exploit completely local time-frequency features. Instead, they rely on global features across the entire frequency spectrum, over a longer period of time.

## 3 Method

Convolutional neural networks (CNNs), take advantage of small scale features present in data. CNNs require less memory and resources than regular fully connected neural networks, allowing for a faster, more efficient model. CNNs have proved to be successful in image processing for tasks such as image super-resolution and semantic segmentation of images. In the image processing field, CNNs take as input a two-dimensional vector of pixel intensities across the spatial dimension and exploit the local spatial correlation among input neurons to learn localized features. A similar two-dimensional representation is used in our model for audio mixtures, using the Short-Time Fourier Transform (STFT), which has frequency and time dimensions. Unlike 2D images, the STFT does not have symmetry across both axes, but a local symmetry can be found along each single axes. Therefore, the filters used in CNNs need to be adapted to the STFT representation of audio.

Figure 1 shows the block diagram for the proposed source separation framework. The STFT is computed on a segment of time context  $T$  of the mixture audio. The resulting magnitude spectrogram is then passed through the network, which outputs an estimate for each of the separated sources. The estimate is used to compute time-frequency soft masks, which are applied to the magnitude spectrogram of the mixture to compute final magnitude estimates for the separated sources. These estimates, along with the phase of the mixture, are used to obtain the audio signals corresponding to the separated sources.

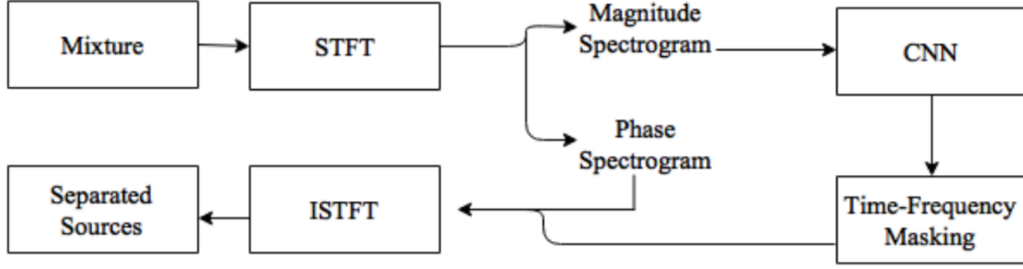


Figure 1: Data Flow

### 3.1 Architecture

The authors use a CNN which functions as a variation of an auto encoder architecture. The network is able to learn an end-to-end model for the separated sources by finding a compressed representation for the training data. The model proposed is shown in Figure 2.

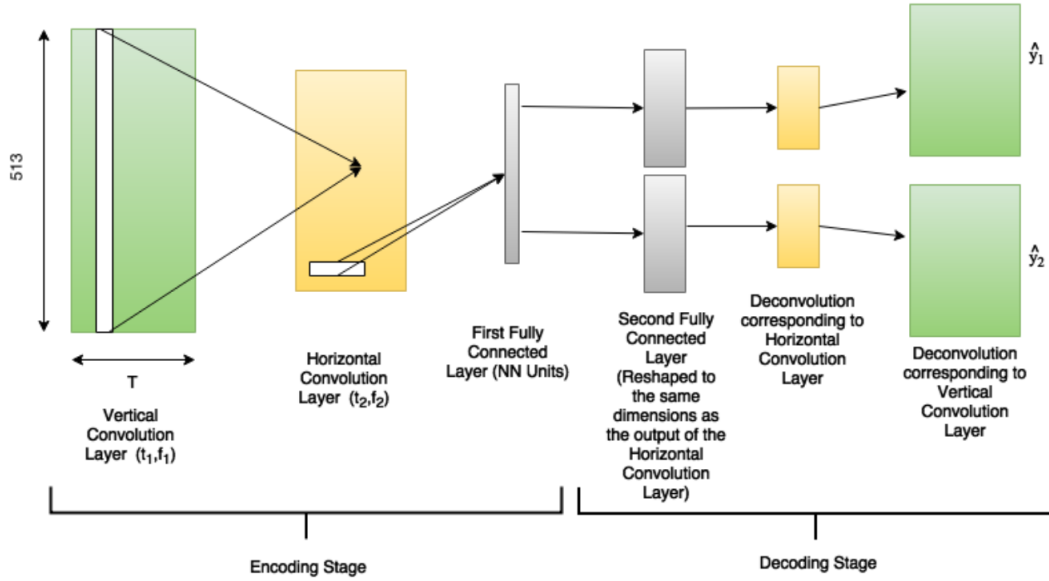


Figure 2: Network Architecture for Source Separation, Using Vertical And Horizontal Convolutions, Showing The Encoding And Decoding Stages

It uses a CNN with two stages: a convolution or encoding stage and the inverse operation, the deconvolution or decoding stage. The stages are explained in further details below.

**Encoding Stage.** This part of the network consists of two convolution layers and a fully connected dense layer, which acts as a bottleneck to compress information.

1. **Vertical convolution layer.** This convolution layer has a shape of  $(t_1, f_1)$ , spanning across a  $t_1$  time frame and taking into account  $f_1$  frequency bins. This layer tries to

capture local timbre information, allowing the model to learn timbre features, similar to the approach used in NMF algorithms for source separation. These features are shared among the sources to be separated, contrary to the NMF approach, where specific basis and activation gains are derived for each source. Therefore, the timbre features learned by this layer need to be robust enough to separate the required source across songs of different genres, where the type of instruments and singers might vary.  $N_1$  filters were used in this layer.

2. **Horizontal convolution layer.** This layer models temporal evolution for different instruments from the features learned in the vertical convolution layer. This is particularly useful for modeling time-frequency characteristics of the different instruments present in the sources to be separated. The filter shape of this layer is  $(t_2, f_2)$  and  $N_2$  filters were used.
3. **Fully connected layer.** The output of the horizontal convolution layer is connected to a fully connected Rectified Linear Unit (ReLU) layer which acts as a bottleneck, achieving dimensional reduction. This layer consists of a non-linear combination of the features learned from the previous layers, with a ReLU non-linearity. The layer is chosen to have fewer elements to reduce the total number of network parameters and to ensure that the network is able to produce a robust representation of the input data. The number of nodes in this layer is represented as  $NN$ .

**Decoding Stage.** The output of the first fully connected layer is passed to another fully connected layer, with a ReLU non-linearity and the same size as the output of the second convolution layer. Thereafter, this layer is reshaped to the same dimensions as the horizontal convolution layer and passed through successive deconvolution layers, the inverse operations to the convolution stage.

### 3.2 Time-frequency masking

The authors wished integrate the computation of a softmask for each of the sources into the network. From the output of the network  $\hat{y}_n(f)$ , a soft mask  $m_n(f)$  can be computed as follows

$$m_n(f) = \frac{|\hat{y}_n(f)|}{\sum_{n=1}^N |\hat{y}_n(f)|}$$

where  $\hat{y}_n(f)$  represents the output of the network for the  $n$ -th source and  $N$  is the total number of sources to be estimated. The estimated mask is then applied to the input mixture signal to estimate the sources  $\tilde{y}_n$

$$\tilde{y}_n(f) = m_n(f) x(f)$$

where  $x(f)$  is the spectrogram of the input mixture signal.

### 3.3 Parameter learning

The neural network is trained to optimize parameters using a Stochastic Gradient Descent (SGD) in order to minimize the squared error between the estimate and the original source  $y_n$ . Formally

$$L_{sq} = \sum_{i=1}^N \|\tilde{y}_n - y_n\|^2$$

## 4 Training

### 4.1 Dataset

The authors used the **Demixing Secrets Dataset 100** (DSD100) for training and testing. This dataset consists of 100 professionally produced full track songs and is designed to evaluate signal source separation methods from music recordings. The dataset contains separate tracks for drums, bass, vocals and other instruments for each song in the set, present as stereo wav files with a sample rate of 44.1 KHz. The four source tracks are mixed using a professional Digital Audio Workstation to create the mixture track for each song. The dataset is divided into a dev set, used for training the network and a test set, which is used for testing the network. Both of these subsets consist of 50 songs each.

### 4.2 Training

During the training phase, the input mixture and the individual sources comprising the mixture were split into 20 seconds segments, and the short time Fourier transform (STFT) for each of these segments was computed. A Hanning window of length 1024 samples was used, which, at a sampling rate of 44.1 KHz corresponds to 23 milliseconds (ms), and a hopsize of 256 samples (5.8 ms), leading to an overlap of 75% across frames. The frames generated from this procedure were grouped into batches of  $T$  frames, representing the maximum time context that the network tries to model. The batches were fed to the model for training, with 30 batches being fed at each round. Thus, each batch consists of  $T$  frames of 513 frequency bins. A complete pass over the entire set is considered as one training epoch and the network is trained for 30 epochs.

## 5 Evaluation

For evaluation, the following measures were used: Source to Distortion Ratio (SDR), Source to Interference Ratio (SIR), Source to Artifacts Ratio(SAR), and Image to Spatial distortion Ratio (ISR). These measures are averaged for overlapping 30 second frames of each song in both the development and the test set.

## 6 Experiment Results

The following result tables represent the results achieved by our experiments in which we tweaked the original network’s parameters. The first 3 experiments investigated batch size

changes. The 4th and 5th experiments investigated different network layer sizes in conjunction with dropout to the fully connected fc0 layer. The 6h experiment investigated different learning rates.

### Experiment 1

Our configuration			Authors' Results
(Train Set, Validation Set, Test Set, Batch Size) = (80%, 10%, 10%, 30)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513, 12, 1, 50, 30, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = (200, $1 \times 10^{-3}$ , [60, 120])			
(Training Loss, Validation Loss) = (20538.27, 111510.24)			
SDR [dB]	Bass	$-21.270 \pm 3.08$	$0.9 \pm 2.7$
	Drums	$-25.027 \pm 1.82$	$2.4 \pm 2.0$
	Others	$-22.087 \pm 1.78$	$0.8 \pm 1.5$
	Vocals	$-21.70 \pm 2.21$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-1.45 \pm 3.23$	$4.6 \pm 4.4$
	Drums	$-4.44 \pm 2.24$	$9.1 \pm 4.3$
	Others	$-2.57 \pm 1.64$	$3.8 \pm 4.0$
	Vocals	$-1.41 \pm 2.05$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.14 \pm 1.49$	$6.9 \pm 2.3$
	Drums	$-19.09 \pm 0.68$	$7.0 \pm 2.8$
	Others	$-17.47 \pm 1.31$	$2.8 \pm 2.4$
	Vocals	$-17.77 \pm 1.35$	$5.3 \pm 2.9$

### Experiment 2

Our configuration			Authors' Results
(Train Set, Validation Set, Test Set, Batch Size) = (80%, 10%, 10%, 5)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513, 12, 1, 50, 30, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = (200, $1 \times 10^{-3}$ , [60, 120])			
(Training Loss, Validation Loss) = (21876.123, 37201.086)			
SDR [dB]	Bass	$-21.60 \pm 2.97.$	$0.9 \pm 2.7$
	Drums	$-24.74 \pm 1.39$	$2.4 \pm 2.0$
	Others	$-22.76 \pm 1.68$	$0.8 \pm 1.5$
	Vocals	$-23.58 \pm 2.086$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-1.41 \pm 2.72$	$4.6 \pm 4.4$
	Drums	$-4.81 \pm 2.087$	$9.1 \pm 4.3$
	Others	$-3.476 \pm 1.81$	$3.8 \pm 4.0$
	Vocals	$-4.32 \pm 2.02$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.584 \pm 1.65$	$6.9 \pm 2.3$
	Drums	$-18.54 \pm 1.00$	$7.0 \pm 2.8$
	Others	$-17.54 \pm 1.08$	$2.8 \pm 2.4$
	Vocals	$-17.74 \pm 1.33$	$5.3 \pm 2.9$

### Experiment 3

Our configuration			Authors' Results
(Train Set, Validation Set, Test Set, Batch Size) = (80%, 10%, 10%, 128)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513, 12, 1, 50, 30, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = $(200, 1 \times 10^{-3}, [60, 120])$			
(Training Loss, Validation Loss) = (14560.074, 144205.705)			
SDR [dB]	Bass	$-21.61 \pm 2.53.$	$0.9 \pm 2.7$
	Drums	$-24.37 \pm 1.23$	$2.4 \pm 2.0$
	Others	$-23.22 \pm 1.53$	$0.8 \pm 1.5$
	Vocals	$-23.43 \pm 2.08$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-1.59 \pm 2.77$	$4.6 \pm 4.4$
	Drums	$-4.44 \pm 1.83$	$9.1 \pm 4.3$
	Others	$-2.97 \pm 1.22$	$3.8 \pm 4.0$
	Vocals	$-4.32 \pm 2.02$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.480 \pm 1.29$	$6.9 \pm 2.3$
	Drums	$-18.46 \pm 0.87$	$7.0 \pm 2.8$
	Others	$-18.39 \pm 1.10$	$2.8 \pm 2.4$
	Vocals	$-18.07 \pm 1.37$	$5.3 \pm 2.9$

Table 1: First 3 experiments investigating batch size changes

### Experiment 4

Our configuration			Authors' Results
(Train Set, Validation Set, Test Set, Batch Size) = (80%, 10%, 10%, 30)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513/3, 12, 1, 50, 30/5, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = (200, $1 \times 10^{-3}$ , [60, 120])			
Convolutional layer kernel size changed to (513/3, 1)			
Added maxpool with kernel size $(2 \cdot (513/3), 1)$			
Added unpool			
(Training Loss, Validation Loss) = (17156.322, 114269.075)			
SDR [dB]	Bass	$-22.32 \pm 2.26.$	$0.9 \pm 2.7$
	Drums	$-24.92 \pm 0.92$	$2.4 \pm 2.0$
	Others	$-23.01 \pm 1.25$	$0.8 \pm 1.5$
	Vocals	$-23.16 \pm 1.88$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-2.25 \pm 1.68$	$4.6 \pm 4.4$
	Drums	$-5.68 \pm 1.58$	$9.1 \pm 4.3$
	Others	$-3.44 \pm 1.59$	$3.8 \pm 4.0$
	Vocals	$-3.66 \pm 1.94$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.91 \pm 1.44$	$6.9 \pm 2.3$
	Drums	$-18.09 \pm 0.92$	$7.0 \pm 2.8$
	Others	$-17.83 \pm 0.94$	$2.8 \pm 2.4$
	Vocals	$-17.80 \pm 1.13$	$5.3 \pm 2.9$

### Experiment 5

Our configuration			Authors' Results
(Train Set, Validation Set, Test Set, Batch Size) = (80%, 10%, 10%, 30)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513, 12, 1, 50, 30, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = (200, $1 \times 10^{-3}$ , [60, 120])			
Same as experiment 4, but added dropout to fully connected layer fc0			
(Training Loss, Validation Loss) = (22073.190, 117374.315)			
SDR [dB]	Bass	$-21.50 \pm 2.80$	$0.9 \pm 2.7$
	Drums	$-25.07 \pm 1.63$	$2.4 \pm 2.0$
	Others	$-22.71 \pm 1.84$	$0.8 \pm 1.5$
	Vocals	$-23.75 \pm 2.10$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-1.32 \pm 2.72$	$4.6 \pm 4.4$
	Drums	$-5.10 \pm 2.51$	$9.1 \pm 4.3$
	Others	$-3.17 \pm 1.78$	$3.8 \pm 4.0$
	Vocals	$-4.52 \pm 2.05$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.53 \pm 1.54$	$6.9 \pm 2.3$
	Drums	$-18.61 \pm 1.14$	$7.0 \pm 2.8$
	Others	$-17.69 \pm 1.17$	$2.8 \pm 2.4$
	Vocals	$-17.76 \pm 1.20$	$5.3 \pm 2.9$

Table 2: Experiments 4 and 5 investigating different network layer sizes and adding dropout to the fully connected fc0 layer

### Experiment 6

Our configuration			Authors' Results
(Train Set, Validation Set, Test Set, Batch Size) = (80%, 10%, 10%, 30)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513, 12, 1, 50, 30, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = (200, $1 \times 10^{-5}$ , [60, 120])			
(Training Loss, Validation Loss) = (21492.434, 115899.219)			
SDR [dB]	Bass	$-21.82 \pm 2.37$	$0.9 \pm 2.7$
	Drums	$-24.99 \pm 1.67$	$2.4 \pm 2.0$
	Others	$-22.21 \pm 1.19$	$0.8 \pm 1.5$
	Vocals	$-23.08 \pm 1.90$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-1.60 \pm 2.57$	$4.6 \pm 4.4$
	Drums	$-5.62 \pm 2.56$	$9.1 \pm 4.3$
	Others	$-2.97 \pm 1.62$	$3.8 \pm 4.0$
	Vocals	$-4.11 \pm 1.94$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.72 \pm 1.10$	$6.9 \pm 2.3$
	Drums	$-18.12 \pm 1.06$	$7.0 \pm 2.8$
	Others	$-17.34 \pm 1.26$	$2.8 \pm 2.4$
	Vocals	$-17.39 \pm 1.35$	$5.3 \pm 2.9$



### Experiment 7

Our configuration			Authors' Results
(Train Set, Validation Set, Test Set, Batch Size) = (80%, 10%, 10%, 30)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513, 12, 1, 50, 30, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = (200, $1 \times 10^{-5}$ , [60, 120])			
Changed learning rate to $1 \times 10^{-6}$ after 10 epochs followed by $1 \times 10^{-7}$ after 50 epochs			
(Training Loss, Validation Loss) = (29001.09, 127572.20)			
SDR [dB]	Bass	$-22.45 \pm 1.89$	$0.9 \pm 2.7$
	Drums	$-25.11 \pm 1.41$	$2.4 \pm 2.0$
	Others	$-22.63 \pm 1.32$	$0.8 \pm 1.5$
	Vocals	$-23.49 \pm 1.94$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-2.57 \pm 2.01$	$4.6 \pm 4.4$
	Drums	$-6.31 \pm 2.28$	$9.1 \pm 4.3$
	Others	$-3.48 \pm 1.69$	$3.8 \pm 4.0$
	Vocals	$-4.55 \pm 1.99$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.82 \pm 1.00$	$6.9 \pm 2.3$
	Drums	$-17.73 \pm 1.03$	$7.0 \pm 2.8$
	Others	$-17.41 \pm 1.16$	$2.8 \pm 2.4$
	Vocals	$-17.49 \pm 1.21$	$5.3 \pm 2.9$

Table 3: Experiments 6 and 7 investigating varying learning rates

### Experiment 8

Our configuration (Same as authors)			Authors' Results
(Train Set, Test + Validation Set, Batch Size) = (50%, 50%, 30)			
$(t_1, f_1, t_2, f_2, N_1, N_2, NN) = (1, 513, 12, 1, 50, 30, 128)$			
$(\alpha, \beta, \beta_{\text{vocals}}) = (0.001, 0.01, 0.03)$			
(Epochs, Learning Rate, Milestones) = (200, $1 \times 10^{-5}$ , [60, 120])			
(Training Loss, Validation Loss) = (18473.046, 97306.92)			
SDR [dB]	Bass	$-22.17 \pm 2.40$	$0.9 \pm 2.7$
	Drums	$-25.37 \pm 1.62$	$2.4 \pm 2.0$
	Others	$-22.25 \pm 2.01$	$0.8 \pm 1.5$
	Vocals	$-22.73 \pm 2.51$	$1.3 \pm 2.4$
SIR [dB]	Bass	$-2.10 \pm 2.72$	$4.6 \pm 4.4$
	Drums	$-5.10 \pm 1.67$	$9.1 \pm 4.3$
	Others	$-2.39 \pm 2.05$	$3.8 \pm 4.0$
	Vocals	$-3.30 \pm 2.17$	$7.2 \pm 3.6$
SAR [dB]	Bass	$-17.74 \pm 1.54$	$6.9 \pm 2.3$
	Drums	$-18.99 \pm 1.23$	$7.0 \pm 2.8$
	Others	$-17.72 \pm 1.50$	$2.8 \pm 2.4$
	Vocals	$-17.59 \pm 1.69$	$5.3 \pm 2.9$

Table 4: Experiment 8 investigating a 50%/50% training and test sets distribution of the dataset, as the original authors did

## 7 Limitations and Proposed Ways to Alleviate Them

### References

- [1] P. Chandna, M. Miron, J. Janer, E. G´omez. Monoaural audio source separation using deep convolutional neural networks. In 13th International Conference on Latent Variable Analysis and Signal Separation (LVAICA2017), 02/2017 2017
- [2] Durrieu, J., Ozerov, A., and F´evotte, C. (2009). Main instrument separation from stereophonic audio signals using a source/filter model. 17th European Signal Processing Conference.
- [3] Rafii, Z., Liutkus, A., and Pardo, B. (2014). REPET for Background/Foreground Separation in Audio, pages 395–411. Springer Berlin Heidelberg, Berlin, Heidelberg.
- [4] Nugraha, A. A., Liutkus, A., and Vincent, E. (2016). Multichannel audio source separation with deep neural networks. Technical report.
- [5] Huang, P.-S., Kim, M., Hasegawa-Johnson, M., and Smaragdis, P. (2014). Deep Learning for Monaural Speech Separation. Acoustics, Speech and Signal Processing (ICASSP), pages 1562–1566.
- [6] Uhlich, S., Giron, F., and Mitsufuji, Y. (2015). Deep neural network based instrument extraction from music. 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 2135–2139. IEEE.