

SUDO RM -RF: EFFICIENT NETWORKS FOR UNIVERSAL AUDIO SOURCE SEPARATION

Efthymios Tzinis[‡] Zhepei Wang[‡] Paris Smaragdis^{‡,§}

[‡] University of Illinois at Urbana-Champaign

[§] Adobe Research

ABSTRACT

In this paper, we present an efficient neural network for end-to-end general purpose audio source separation. Specifically, the backbone structure of this convolutional network is the SuDORM-RF as well as their aggregation which is performed through simple one-dimensional convolutions. In this way, we are able to obtain high quality audio source separation with limited number of floating point operations, memory requirements, number of parameters and latency. Our experiments on both speech and environmental sound separation datasets show that SuDORM-RF performs comparably and even surpasses various state-of-the-art approaches with significantly higher computational resource requirements.

Index Terms— Audio source separation, low-cost neural networks, deep learning

1. INTRODUCTION

The advent of the deep learning era has enabled the effective usage of neural networks towards single-channel source separation with mask-based architectures [1]. Recently, end-to-end source separation in time-domain has shown state-of-the-art results in a variety of separation tasks such as: speech separation [2, 3], universal sound separation [4, 5] and music source separation [6]. The separation module of ConvTasNet [2] and its variants [4, 5] consist of multiple stacked layers of depth-wise separable convolutions [7] which can aptly incorporate long-term temporal relationships. Building upon the effectiveness of a large temporal receptive field, a dual-path recurrent neural network (DPRNN) [3] has shown remarkable performance on speech separation. Demucs [6] has a refined U-Net structure [8] and has shown strong performance improvement on music source separation. Specifically, it consists of several convolutional layers in each a downsampling operation is performed in order to extract high dimensional features. A two-step approach has been introduced in [9] and showed that universal sound separation models could be further improved when working directly on the latent space and learning the ideal masks on a separate step.

Despite the dramatic advances in source separation performance, the computational complexity of the aforementioned methods might hinder their extensive usage across multiple devices. Specifically, many of these algorithms are not amenable to, e.g., embedded systems deployment, or other environments where computational resources are constrained. Additionally, training such systems is also an expensive computational undertaking which can amount to significant costs.

Several studies, mainly in the image domain, have introduced more efficient architectures in order to overcome the growing concern of large models with high computational requirements. Models with depth-wise separable convolutions [7] have shown strong potential for several image-domain tasks [10] while significantly reducing the computational requirements. Thus, several variants such as MobileNets [11] have been proposed for deep learning on edge-devices. However, convolutions with a large dilation factor might inject several artifacts and thus, lightweight architectures that combine several dilation factors in each block have been proposed for image tasks [12]. More recent studies propose meta-learning algorithms for optimizing architecture configurations given specific computational resource and accuracy requirements [13, 14].

Despite the recent success on low-resource architectures on the image domain, little progress has been made towards proposing efficient architectures for audio tasks and especially source separation. In [15] a WaveRNN is used for efficient audio synthesis in terms of floating point operations (FLOPs) and latency. Other studies have introduced audio source separation models with reduced number of trainable parameters [3, 16] and binarized models [17]. In this study, we propose a novel efficient neural network architecture for audio source separation while following a more holistic approach in terms of computational resources that we take into consideration (FLOPs, latency and total memory requirements). Our proposed model performs SuDORM-RF using depth-wise convolutions. By doing so, SuDORM-RF exploits the effectiveness of iterative temporal resampling strategies [18] and avoids the need of multiple stacked dilated convolutional layers [2]. We report a separation performance comparable or even better to several recent state-of-the-art models on speech and environmental sound separation tasks with significantly lower computational requirements. Our experiments suggest that SuDORM-RF models a) could be deployed on devices with limited resources, b) be trained significantly faster and achieve good separation performance and c) scale well when increasing the number of parameters. Our code is available online¹.

2. SUDO RM -RF NETWORK ARCHITECTURE

On par with many state-of-the-art approaches in the literature [2, 9, 3, 6], SuDORM-RF performs end-to-end audio source separation using a mask-based architecture with adaptive encoder and decoder basis. The input is the raw signal from a mixture $\mathbf{x} \in \mathbb{R}^T$ with T samples in the time-domain. First we feed the input mixture \mathbf{x} to an encoder \mathcal{E} in order to obtain a latent representation for the mixture $\mathbf{v}_{\mathbf{x}} = \mathcal{E}(\mathbf{x}) \in \mathbb{R}^{C_{\mathcal{E}} \times L}$. Consequently the latent mixture representation is fed through a separation module \mathcal{S} which estimates the

¹Supported by NSF grant #1453104. A Titan V used for this research was donated by the NVIDIA Corporation.

¹Code: https://github.com/etzinis/sudo_rm_rf

corresponding masks $\hat{\mathbf{m}}_i \in \mathbb{R}^{C_\mathcal{E} \times L}$ for each one of the N sources $\mathbf{s}_1, \dots, \mathbf{s}_N \in \mathbb{R}^T$ which constitute in the mixture. The estimated latent representation for each source in the latent space $\hat{\mathbf{v}}_i$ is retrieved by multiplying element-wise an estimated mask $\hat{\mathbf{m}}_i$ with the encoded mixture representation $\mathbf{v}_\mathbf{x}$. Finally, the reconstruction for each source $\hat{\mathbf{s}}_i$ is obtained by using a decoder \mathcal{D} to transform the latent-space $\hat{\mathbf{v}}_i$ source estimates back into the time-domain $\hat{\mathbf{s}}_i = \mathcal{D}(\hat{\mathbf{v}}_i)$. An overview of the SuDoRM-RF architecture is displayed in Figure 1. The encoder, separator and decoder modules are described in Sections 2.1, 2.2 and 2.3, respectively. For simplicity of our notation we will describe the whole architecture assuming that the processed batch size is one. Moreover, we are going to define some useful operators of the various convolutions which are used in SuDoRM-RF.

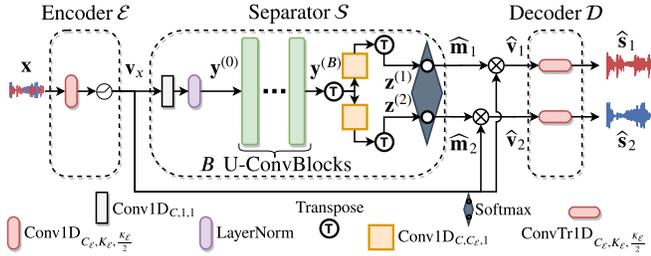


Fig. 1: SuDoRM-RF architecture for separating two sources.

Definition 2.1. $\text{Conv1D}_{C,K,S} : \mathbb{R}^{C_{in} \times L_{in}} \rightarrow \mathbb{R}^{C \times L}$ defines a kernel $\mathbf{W} \in \mathbb{R}^{C \times C_{in} \times K}$ and a bias vector $\mathbf{b} \in \mathbb{R}^C$. When applied on a given input $\mathbf{x} \in \mathbb{R}^{C_{in} \times L_{in}}$ it performs a one-dimensional convolution operation with stride equal to S as shown next:

$$\text{Conv1D}_{C,K,S}(\mathbf{x})_{i,l} = \mathbf{b}_i + \sum_{j=1}^{C_{in}} \sum_{k=1}^K \mathbf{W}_{i,j,k} \cdot \mathbf{x}_{j,S \cdot l - k}, \quad (1)$$

where the indices i, j, k, l denotes the output channel, the input channel, the kernel sample and the temporal index, respectively. Note that without loss of generality and performing appropriate padding, the last dimension of the output representation would be $L = \lfloor L_{in}/S \rfloor$.

Definition 2.2. $\text{ConvTr1D}_{C,K,S} : \mathbb{R}^{C_{in} \times L_{in}} \rightarrow \mathbb{R}^{C \times L}$ defines a one-dimensional transpose convolution. Since any convolution operation could be expressed as a matrix multiplication, transposed convolution can be directly understood as the gradient calculation for a regular convolution w.r.t. its input [19].

Definition 2.3. $\text{DWConv1D}_{C,K,S} : \mathbb{R}^{C_{in} \times L_{in}} \rightarrow \mathbb{R}^{C \times L}$ defines a one-dimensional depth-wise convolution operation [7]. In essence, this operator defines $G = C_{in}$ separate one-dimensional convolutions $\mathcal{F}_i = [\text{Conv1D}_{C_G,K,S}]_i$ with $i \in \{1, \dots, G\}$ where $C_G = \lfloor C/G \rfloor$. Given an input $\mathbf{x} \in \mathbb{R}^{C_{in} \times L_{in}}$ the i th one-dimensional convolution contributes to $C_G = \lfloor C/G \rfloor$ output channels by considering as input only the i th row of the input as described below:

$$\text{DWConv1D}_{C,K,S}(\mathbf{x}) = \text{Concat}(\{\mathcal{F}_i(\mathbf{x}_i), \forall i\}), \quad (2)$$

where $\text{Concat}(\cdot)$ performs the concatenation of all individual one-dimensional convolution outputs across the channel dimension.

2.1. Encoder

The encoder \mathcal{E} architecture consists of a one-dimensional convolution with kernel size $K_\mathcal{E}$ and stride equal to $K_\mathcal{E}/2$ similar to [2]. Each

convolved input audio-segment of $K_\mathcal{E}$ samples is transformed to a $C_\mathcal{E}$ -dimensional vector representation where $C_\mathcal{E}$ is the number of output channels of the 1D-convolution. We force the output of the encoder to be strictly non-negative by applying a rectified linear unit (ReLU) activation on top of the output of the 1D-convolution. Thus, the encoded input mixture representation could be expressed as:

$$\mathbf{v}_\mathbf{x} = \mathcal{E}(\mathbf{x}) = \text{ReLU}(\text{Conv1D}_{C_\mathcal{E},K_\mathcal{E},K_\mathcal{E}/2}(\mathbf{x})) \in \mathbb{R}^{C_\mathcal{E} \times L}, \quad (3)$$

where the activation $\text{ReLU}(\cdot)$ is applied element-wise.

2.2. Separator

In essence, the separator \mathcal{S} module performs the following transformations to the encoded mixture representation $\mathbf{v}_\mathbf{x} \in \mathbb{R}^{C_\mathcal{E} \times L}$:

1. Projects the encoded mixture representation $\mathbf{v}_\mathbf{x} \in \mathbb{R}^{C_\mathcal{E} \times L}$ to a new channel space through a layer-normalization (LN) [20] followed by a point-wise convolution as shown next:

$$\mathbf{y}_0 = \text{Conv1D}_{C,1,1}(\text{LN}(\mathbf{v}_\mathbf{x})) \in \mathbb{R}^{C \times L}, \quad (4)$$

where $\text{LN}(\mathbf{v}_\mathbf{x})$ denotes a layer-normalization layer in which the moments used are extracted across the temporal dimension for each channel separately.

2. Performs repetitive non-linear transformations provided by B U-convolutional blocks (U-ConvBlocks) on the intermediate representation \mathbf{y}_0 . In other words, the output of the i th U-ConvBlock would be denoted as $\mathbf{y}_i \in \mathbb{R}^{C \times L}$ and would be used as input for the $(i+1)$ th block. Each U-ConvBlock extracts and aggregates information from multiple resolutions which is extensively described in Section 2.2.1.
3. Aggregates the information over multiple channels by applying a regular one-dimensional convolution for each source on the transposed feature representation $\mathbf{y}_B^T \in \mathbb{R}^{L \times C}$. Effectively, for the i th source we obtain an intermediate latent representation as shown next:

$$\mathbf{z}_i = \text{Conv1D}_{C,C_\mathcal{E},1}(\mathbf{y}_B^T)^T \in \mathbb{R}^{C_\mathcal{E} \times L} \quad (5)$$

This step has been introduced in [9] and empirically shown to make the training process more stable rather than using the activations from the final block \mathbf{y}_B to estimate the masks.

4. Combines the aforementioned latent codes for all sources $\mathbf{z}_i \forall i \in \{1, \dots, N\}$ by performing a softmax operation in order to get mask estimates $\hat{\mathbf{m}}_i \in [0, 1]^{C_\mathcal{E} \times L}$ which add up to one across the dimension of the sources. Namely, the corresponding mask estimate for the i th source would be:

$$\hat{\mathbf{m}}_i = \text{vec}^{-1} \left(\frac{\exp(\text{vec}(\mathbf{z}_i))}{\sum_{j=1}^N \exp(\text{vec}(\mathbf{z}_j))} \right) \in \mathbb{R}^{C_\mathcal{E} \times L}, \quad (6)$$

where $\text{vec}(\cdot) : \mathbb{R}^{K \times N} \rightarrow \mathbb{R}^{K \cdot N}$ and $\text{vec}^{-1}(\cdot) : \mathbb{R}^{K \cdot N} \rightarrow \mathbb{R}^{K \times N}$ denotes the vectorization of an input tensor and the inverse operation, respectively.

5. Estimates a latent representation $\hat{\mathbf{v}}_i \in \mathbb{R}^{C_\mathcal{E} \times L}$ for each source by multiplying element-wise the encoded mixture representation $\mathbf{v}_\mathbf{x}$ with the corresponding mask $\hat{\mathbf{m}}_i$:

$$\hat{\mathbf{v}}_i = \mathbf{v}_\mathbf{x} \odot \hat{\mathbf{m}}_i \in \mathbb{R}^{C_\mathcal{E} \times L}, \quad (7)$$

where $\mathbf{a} \odot \mathbf{b}$ is the element-wise multiplication of the two tensors \mathbf{a} and \mathbf{b} assuming that they have the same shape.

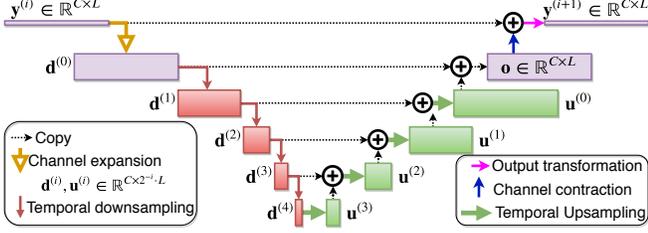


Fig. 2: U-ConvBlock architecture.

Algorithm 1: U-ConvBlock forward pass

Input: $\mathbf{y}^{(i)} \in \mathbb{R}^{C \times L}$
Output: $\mathbf{y}^{(i+1)} \in \mathbb{R}^{C \times L}$
// Expand channel dimensions
 $\mathbf{q} \leftarrow \text{PReLU}_{C_U} \left(\text{LN} \left(\text{Conv1D}_{C_U, 1, 1} \left(\mathbf{y}^{(i)} \right) \right) \right)$;
 $\mathbf{d}^{(0)} \leftarrow \text{PReLU}_{C_U} \left(\text{LN} \left(\text{DWConv1D}_{C_U, K_U, 1} \left(\mathbf{q} \right) \right) \right)$;
for $i = 1; i++$; **while** $i \leq Q$ **do**
// Successive depth-wise downsampling
 $\mathbf{d}^{(i)} \leftarrow \text{LN} \left(\text{DWConv1D}_{C_U, K_U, S_U} \left(\mathbf{d}^{(i-1)} \right) \right)$;
 $\mathbf{d}^{(i)} \leftarrow \text{PReLU}_{C_U} \left(\mathbf{d}^{(i)} \right)$;
end
 $\mathbf{u}^{(Q)} \leftarrow \mathbf{d}^{(Q)}$;
for $i = Q - 1; i--$; **while** $i \geq 0$ **do**
// Upsample and add resolutions
 $\mathbf{u}^{(i)} \leftarrow \mathbf{d}^{(i)} + \mathcal{I}_{S_U} \left(\mathbf{u}^{(i+1)} \right)$;
end
 $\mathbf{o} \leftarrow \text{LN} \left(\text{Conv1D}_{C, 1, 1} \left(\text{PReLU}_C \left(\text{LN} \left(\mathbf{u}^{(0)} \right) \right) \right) \right)$;
return $\text{PReLU}_C \left(\mathbf{y}^{(i)} + \mathbf{o} \right)$;

2.2.1. U-convolutional block (U-ConvBlock)

U-ConvBlock uses a block structure which resembles a depth-wise separable convolution [7] with a skip connection as in ConvTas-Net [2]. However, instead of performing a regular depth-wise convolution as shown in [10] or a dilated depth-wise which has been successfully utilized for source separation [2, 5, 9] our proposed U-ConvBlock extracts information from multiple resolutions using Q successive temporal downsampling and Q upsampling operations similar to a U-Net architecture [8]. More importantly the output of each block leaves the temporal resolution intact while increasing the effective receptive field of the network multiplicatively with each temporal sub-sampling operation [21]. An abstract view of the i th U-ConvBlock is displayed in Figure 2 while a detailed description of the operations is presented in Algorithm 1.

Definition 2.4. $\text{PReLU}_C : \mathbb{R}^{C \times L} \rightarrow \mathbb{R}^{C \times L}$ defines a parametric rectified linear unit (PReLU) [22] with C learnable parameters $\mathbf{a} \in \mathbb{R}^C$. When applied to an input matrix $\mathbf{y} \in \mathbb{R}^{C \times L}$ the non-linear transformation could be defined element-wise as:

$$\text{PReLU}_C(\mathbf{y})_{i,j} = \max(0, \mathbf{y}_{i,j}) + \mathbf{a}_i \cdot \min(0, \mathbf{y}_{i,j}) \quad (8)$$

Definition 2.5. $\mathcal{I}_M : \mathbb{R}^{C \times L} \rightarrow \mathbb{R}^{C \times M \cdot L}$ defines a nearest neighbor temporal interpolation by a factor of M . When applied on an input matrix $\mathbf{y} \in \mathbb{R}^{C \times L}$ this upsampling procedure could be formally expressed element-wise as: $\mathcal{I}_M(\mathbf{u})_{i,j} = \mathbf{u}_{i, \lfloor j/M \rfloor}$

2.3. Decoder

Our decoder module \mathcal{D} is the final step in order to transform the latent space representation $\hat{\mathbf{v}}_i$ for each source back to the time-domain. In our proposed model we follow a similar approach as in [9] where each latent source representation $\hat{\mathbf{v}}_i$ is fed through a different transposed convolution decoder $\text{ConvTr1D}_{C_E, K_E, K_E/2}$. The efficacy of dealing with different types of sources using multiple decoders has also been studied in [23]. Ignoring the permutation problem, for the i th source we have the following reconstruction in time:

$$\hat{\mathbf{s}}_i = \mathcal{D}_i(\hat{\mathbf{v}}_i) = \text{ConvTr1D}_{C_E, K_E, K_E/2}(\hat{\mathbf{v}}_i) \quad (9)$$

3. EXPERIMENTAL SETUP

3.1. Audio source separation tasks

Speech separation: We perform speech separation experiments in accordance with the publicly available WSJ0-2mix dataset [24] and other studies [3, 25, 26]. Speaker mixtures are generated by randomly mixing speech utterances with two active speakers at random signal to noise ratios (SNR)s between -5 and 5 dB from the Wall Street Journal (WSJ0) corpus [27].

Non-speech sound separation: For our non-speech sound separation experiments we follow the exact same setup as in [9] and utilize audio clips from the environmental sound classification (ESC50) data collection [28] which consists of a wide variety of sounds (non-speech human sounds, animal sounds, natural soundscapes, interior sounds and urban noises). For each data sample, two audio sources are mixed with a random SNR between -2.5 and 2.5 dB where each source belongs to a distinct sound category from a total of 50.

3.2. Data preprocessing and generation

We follow the same data augmentation process which was firstly introduced in [9] and it has been show beneficial in other recent studies [25]. The process for generating a mixture is the following: A) random choosing two sound classes (for universal sound separation) or speakers (for speech separation) B) random cropping of 4sec segments from two sources audio files C) mixing the source segments with a random SNR (as specified in Section 3.1). For each epoch, 20,000 new training mixtures are generated. Validation and test sets are generated once with each one containing 3,000 mixtures. Moreover, we downsample each audio clip to 8kHz, subtract its mean and divide with the standard deviation of the mixture.

3.3. Training and evaluation details

All models are trained for 120 epochs using a batch size equal to 4. As a loss function we use the negative permutation-invariant [29] scale-invariant signal to distortion ratio (SI-SDR) [30] which is defined between the clean sources \mathbf{s} and the estimates $\hat{\mathbf{s}}$ as:

$$\mathcal{L} = -\text{SI-SDR}(\mathbf{s}^*, \hat{\mathbf{s}}) = -10 \log_{10} \left(\frac{\|\alpha \mathbf{s}^*\|^2}{\|\alpha \mathbf{s}^* - \hat{\mathbf{s}}\|^2} \right), \quad (10)$$

where \mathbf{s}^* denotes the permutation of the sources that maximizes SI-SDR and $\alpha = \hat{\mathbf{s}}^\top \mathbf{s}^* / \|\mathbf{s}\|^2$ is just a scalar. In order to evaluate the performance of our models we use the SI-SDR improvement (SI-SDRi) which is the gain that we get on SI-SDR measure using estimated signal instead of the mixture signal.

3.4. SuDoRM-RF configurations

For the encoder \mathcal{E} and decoder modules \mathcal{D} we use a kernel size $K_{\mathcal{E}} = 21$ corresponding to 2.625ms and a number of basis equal to $C_{\mathcal{E}} = 512$. For the configuration of each U-ConvBlock we set the input number of channels equal to $C = 128$, the number of successive resampling operations equal to $Q = 4$ and the expanded number of channels equal to $C_U = 512$. In each subsampling operation we reduce the temporal dimension by a factor of 2 and all depth-wise separable convolutions have a kernel length of $K_U = 5$ and a stride of $S_U = 2$. We propose 3 different models which are configured through the number B of U-ConvBlocks inside the separator module \mathcal{S} . Namely, SuDoRM-RF 1.0x, SuDoRM-RF 0.5x, SuDoRM-RF 0.25x consist of 16, 8 and 4 blocks, respectively. During training, we use the Adam optimizer [31] with an initial learning rate set to 0.001 and we decrease it by a factor of 5 every 50 epochs.

3.5. Literature models configurations

We compare against the best configurations of some of the latest state-of-the-art approaches for speech [2, 3], universal [9] and music [6] source separation. For a fair comparison with the aforementioned models we use the authors original code, the best performing configurations for the proposed models as well as the suggested training process. For Demucs [6], 80 channels are used instead of 100 in order to be able to train it on a single graphical processing unit (GPU).

3.6. Measuring computational resources

One of the main goals of this study is to propose a model for audio source separation which could be trained using limited computational resources and deployed easily on a mobile or edge-computing device [32]. Specifically, we consider the following aspects which might cause a computational bottleneck during inference or training:

1. Number of executed floating point operations (FLOPs).
2. Number of trainable parameters.
3. Memory allocation required on the device for a single pass.
4. Time for completing each process.

We are using various sampling profilers in Python for tracing all the requirements on an Intel Xeon CPU E5-2695 v3 @ 2.30GHz CPU and a Nvidia Tesla K80 GPU.

4. RESULTS & DISCUSSION

In Table 1, we show the separation performance alongside computational requirements for some of the most recent state-of-the-art models in the literature and the proposed SuDoRM-RF configurations. It is easy to see that the proposed models can match and even outperform the separation performance of other several state-of-the-art models using orders of magnitude less computational requirements. A better visualization for understanding the Pareto efficiency of the proposed architectures is displayed in Figure 3 where we show for each model its performance on non-speech sound separation vs a specific computational requirement. We do not show the same plots for speech separation on the WSJ dataset as the patterns were similar.

4.1. Floating point operations (FLOPs)

Different devices (CPU, GPU, mobiles, etc.) have certain limitations on their FLOPs throughput capacity. In the case of an edge device,

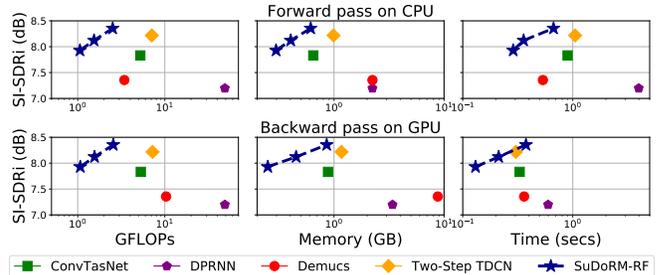


Fig. 3: SI-SDRi non-speech sound separation performance on ESC50 vs computational resources with an input audio of 8000 samples for all models. (Top row) computational requirements for a single forward pass on CPU (Bottom) for a backward pass on GPU. All x-axis are shown in log-scale while the 3 connected blue stars correspond to the three SuDoRM-RF configurations from Section 3.4.

the computational resource one might be interested in is the number of FLOPs required during inference. On the other hand, training on cloud machines might be costly if a huge number of FLOPs is needed in order to achieve high separation performance. As a result, it is extremely important to be able to train and deploy models which require a low number of computations [11]. We see from the first column of Figure 3 that SuDoRM-RF models scale well as we increase the number of U-ConvBlocks B from 4 \rightarrow 8 \rightarrow 16. Furthermore, we see that for both forward and backward passes the family of the proposed SuDoRM-RF models appear more Pareto efficient in terms of SI-SDRi performance vs Giga-FLOPs (GFLOPs) and time required compared to the other state-of-the-art models which we take into account. Specifically, the DPRNN model [3] which performs sequential matrix multiplications (even with a low number of parameters) requires at least 45 times more FLOPs for a single pass compared to SuDoRM-RF 0.25x while performing worse when trained for the same number of epochs.

4.1.1. Cost-efficient training

Usually one of the most detrimental factors for training deep learning models is the requirement of allocating multiple GPU devices for several days or weeks until an adequate performance is obtained on the validation set. This huge power consumption could lead to huge cloud services rental costs and carbon dioxide emissions [14]. In Figure 4, we show the validation SI-SDRi performance for the speech separation task which is obtained by each model versus the total amount of FLOPs performed. For each training epoch all models perform updates while iterating over 20,000 audio mixtures. Notably, SuDoRM-RF models outperform all other models in terms of cost-efficient training as they obtain better separation performance while requiring significantly less amount of training FLOPs. For instance, SuDoRM-RF 1.0x obtains ≈ 16 dB in terms of SI-SDRi compared to ≈ 10 dB of DPRNN [3] which manages to complete only 3 epochs given the same number of training FLOPs.

4.2. Trainable parameters

From Table 1 it is easy to see that SuDoRM-RF architectures are using orders of magnitude fewer parameters compared to the U-net architectures like Demucs [6] where each temporal downsampling is followed by a proportional increase to the number of channels. Moreover, the upsampling procedure inside each U-ConvBlock

| Model | SI-SDRi (dB) | | Parameters (millions) | GFLOPs | | Memory (GB) | | Time (sec) | |
|-------------------|-------------------|-----------------------|-----------------------|--------|-------|-------------|------|------------|------|
| | Speech separation | Non-speech separation | | I | B | I | B | I | B |
| ConvTasNet [2] | 15.30* | 7.74 | 5.05 | 5.23 | 5.30 | 0.65 | 0.88 | 0.90 | 0.33 |
| Demucs [6] | 12.12 | 7.23 | 415.09 | 3.43 | 10.34 | 2.24 | 8.77 | 0.53 | 0.36 |
| DPRNN [3] | 18.80* | 7.20 | 2.63 | 48.89 | 48.90 | 2.23 | 3.40 | 3.98 | 0.60 |
| Two-Step TDCN [9] | 16.10* | 8.22 | 8.63 | 7.09 | 7.23 | 0.99 | 1.17 | 1.05 | 0.30 |
| SuDoRM-RF 1.0x | 17.02 | 8.35 | 2.66 | 2.52 | 2.56 | 0.61 | 0.86 | 0.67 | 0.38 |
| SuDoRM-RF 0.5x | 15.37 | 8.12 | 1.42 | 1.54 | 1.56 | 0.40 | 0.45 | 0.36 | 0.21 |
| SuDoRM-RF 0.25x | 13.39 | 7.93 | 0.79 | 1.06 | 1.07 | 0.30 | 0.25 | 0.29 | 0.13 |

Table 1: SI-SDRi separation performance for all models on both separation tasks (speech and non-speech) alongside their computational requirements for performing inference on CPU (I) and a backward update step on GPU (B) for one second of input audio or equivalently 8000 samples. * We assign the maximum SI-SDRi performance obtained by our runs and the reported number on the corresponding paper.

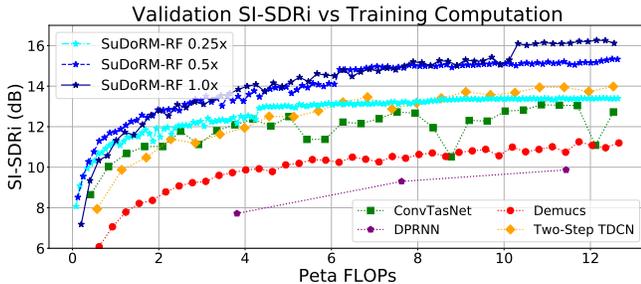


Fig. 4: Validation SI-SDRi separation performance for speech separation vs the number of FLOPs executed during training. All models are trained using batches of 4 mixtures with 32,000 time-samples each. Each point corresponds to a completed training epoch.

does not require any additional parameters. The SuDoRM-RF models seem to increase their effective receptive field with significantly fewer parameters compared to dilated convolutional architectures [2, 9]. Notably, our largest model SuDoRM-RF 1.0x matches the relatively low number of parameters of the DPRNN [3] model which is based on stacked RNN layers.

4.3. Memory requirements

In most of the studies where efficient architectures are introduced [11, 10, 12, 13] authors are mainly concerned with the total number of trainable parameters of the network. The same applies to efficient architectures for source separation [2, 3, 16]. However, the trainable parameters is only a small portion of total amount of memory required for a single forward or backward pass. The space complexity could easily blow up by the storage of intermediate representations. The latter could become even worse when multiple skip connections are present, gradients from multiple layers have to be stored or implementations require augmented matrices (dilated, transposed convolutions, etc.). In Figure 3, we see that SuDoRM-RF models are more pareto-efficient in terms of the memory required compared to the dilated convolutional architectures of ConvTasNet [2] and Two-Step TDCN [9] where they require an increased network depth in order to increase their receptive field. Although SuDoRM-RF models do not perform downsampling in every feature extraction step as Demucs [6] does, we see that the proposed models require orders of magnitude less memory especially during a backward update step as the number of parameters in Demucs is significantly higher. Finally, SuDoRM-RF models have a smaller memory footprint because the encoder \mathcal{E} performs a temporal downsampling by a factor

| $K_{\mathcal{E}}$ | C | B | Q | Norm | Mask Act. | Dec. | SI-SDRi |
|-------------------|-----|-----|-----|------|-----------|------|---------|
| 21 | 128 | 16 | 4 | LN | Softmax | 2 | 16.0 |
| 17 | 128 | 16 | 4 | LN | ReLU | 1 | 15.9 |
| 17 | 128 | 16 | 4 | GLN | ReLU | 1 | 16.8 |
| 21 | 256 | 20 | 4 | GLN | ReLU | 1 | 17.7 |
| 41 | 256 | 32 | 4 | GLN | ReLU | 1 | 17.1 |
| 41 | 256 | 20 | 4 | GLN | ReLU | 1 | 16.8 |
| 21 | 512 | 18 | 7 | GLN | ReLU | 1 | 18.0 |
| 21 | 512 | 20 | 2 | GLN | ReLU | 1 | 17.4 |
| 21 | 512 | 34 | 4 | GLN | ReLU | 1 | 18.9 |

Table 2: SI-SDRi separation performance on WSJ0-2mix for various parameter configurations of SuDoRM-RF models. Mask Act. corresponds to the activation function before the mask estimation and Dec. specifies the number of decoders we are using before reconstructing the time-domain signals. GLN corresponds to the global layer normalization as described in [2]. All the other parameters have the same values as described in Section 3.4

of $\text{div}(K_{\mathcal{E}}, 2) = 10$ compared to DPRNN [3] which does not reduce the temporal resolution at all.

4.4. Ablation study on WSJ0-2mix

We perform a small ablation study in order to show how different parameter choices in SuDoRM-RF models affect the separation performance. In order to be directly comparable with the numbers reported by several other studies [2, 3, 25, 26], we train our models for 200 epochs and test them using the given data splits from WSJ0-2mix dataset [24]. The results are shown in Table 2.

5. CONCLUSIONS

In this study, we have introduced the SuDoRM-RF network, a novel architecture for efficient universal sound source separation. The proposed model is capable of extracting multi-resolution temporal features through successive depth-wise convolutional downsampling of intermediate representations and aggregates them using a non-parametric interpolation scheme. In this way, SuDoRM-RF models are able to significantly reduce the required number of layers in order to effectively capture long-term temporal dependencies. We show that these models can perform similarly or even better than recent state-of-the-art models while requiring significantly less computational resources in FLOPs, memory and time. In the future, we aim to use SuDoRM-RF models for real-time low-cost source separation.

6. REFERENCES

- [1] Po-Sen Huang, Minje Kim, Mark Hasegawa-Johnson, and Paris Smaragdis, “Deep learning for monaural speech separation,” in *Proc. ICASSP*, 2014, pp. 1562–1566.
- [2] Yi Luo and Nima Mesgarani, “Conv-tasnet: Surpassing ideal time–frequency magnitude masking for speech separation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [3] Yi Luo, Zhuo Chen, and Takuya Yoshioka, “Dual-path rnn: efficient long sequence modeling for time-domain single-channel speech separation,” in *Proc. ICASSP*, 2020.
- [4] Ilya Kavalero, Scott Wisdom, Hakan Erdogan, Brian Patton, Kevin Wilson, Jonathan Le Roux, and John R Hershey, “Universal sound separation,” in *Proc. WASPAA*, 2019, pp. 175–179.
- [5] Efthymios Tzinis, Scott Wisdom, John R Hershey, Aren Jansen, and Daniel PW Ellis, “Improving universal sound separation using sound classification,” in *Proc. ICASSP*, 2020.
- [6] Alexandre Défossez, Nicolas Usunier, Léon Bottou, and Francis Bach, “Music source separation in the waveform domain,” *arXiv preprint arXiv:1911.13254*, 2019.
- [7] Laurent Sifre and Stéphane Mallat, “Rigid-motion scattering for image classification,” *Ph. D. thesis*, 2014.
- [8] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [9] Efthymios Tzinis, Shrikant Venkataramani, Zhepei Wang, Cem Subakan, and Paris Smaragdis, “Two-step sound source separation: Training on learned latent targets,” in *Proc. ICASSP*, 2020.
- [10] François Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *Proc. CVPR*, 2017, pp. 1251–1258.
- [11] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam, “Mobilenets: Efficient convolutional neural networks for mobile vision applications,” *arXiv preprint arXiv:1704.04861*, 2017.
- [12] Sachin Mehta, Mohammad Rastegari, Linda Shapiro, and Hananeh Hajishirzi, “Espnetv2: A light-weight, power efficient, and general purpose convolutional neural network,” in *Proc. CVPR*, 2019, pp. 9190–9200.
- [13] Jiahui Yu, Linjie Yang, Ning Xu, Jianchao Yang, and Thomas Huang, “Slimmable neural networks,” in *Proc. ICLR*, 2019.
- [14] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han, “Once for all: Train one network and specialize it for efficient deployment,” in *Proc. ICLR*, 2020.
- [15] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aaron Oord, Sander Dieleman, and Koray Kavukcuoglu, “Efficient neural audio synthesis,” in *International Conference on Machine Learning*, 2018, pp. 2410–2419.
- [16] Alejandro Maldonado, Caleb Rascon, and Ivette Velez, “Lightweight online separation of the sound source of interest through blstm-based binary masking,” *arXiv preprint arXiv:2002.11241*, 2020.
- [17] Minje Kim and Paris Smaragdis, “Efficient source separation using bitwise neural networks,” in *Audio Source Separation*, pp. 187–206. Springer, 2018.
- [18] Muhammad Haris, Gregory Shakhnarovich, and Norimichi Ukita, “Deep back-projection networks for super-resolution,” in *Proc. CVPR*, 2018, pp. 1664–1673.
- [19] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, “Deep inside convolutional networks: Visualising image classification models and saliency maps,” in *Workshop Proc. ICLR*, 2014.
- [20] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.
- [21] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel, “Understanding the effective receptive field in deep convolutional neural networks,” in *Advances in neural information processing systems*, 2016, pp. 4898–4906.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proc. CVPR*, 2015, pp. 1026–1034.
- [23] G. Brunner, N. Naas, S. Palsson, O. Richter, and R. Wattenhofer, “Monaural music source separation using a resnet latent separator network,” in *Proc. ICTAI*, 2019, pp. 1124–1131.
- [24] John R Hershey, Zhuo Chen, Jonathan Le Roux, and Shinji Watanabe, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *Proc. ICASSP*, 2016, pp. 31–35.
- [25] Neil Zeghidour and David Grangier, “Wavesplit: End-to-end speech separation by speaker clustering,” *arXiv preprint arXiv:2002.08933*, 2020.
- [26] Yuzhou Liu and DeLiang Wang, “Divide and conquer: A deep casa approach to talker-independent monaural speaker separation,” *arXiv preprint arXiv:1904.11148*, 2019.
- [27] Douglas B. Paul and Janet M. Baker, “The design for the wall street journal-based CSR corpus,” in *Speech and Natural Language: Proceedings of a Workshop Held at Harriman, New York, February 23-26, 1992*, 1992.
- [28] Karol J Piczak, “Esc: Dataset for environmental sound classification,” in *Proc. ACM International Conference on Multimedia*, 2015, pp. 1015–1018.
- [29] Dong Yu, Morten Kolbæk, Zheng-Hua Tan, and Jesper Jensen, “Permutation invariant training of deep models for speaker-independent multi-talker speech separation,” in *Proc. ICASSP*, 2017, pp. 241–245.
- [30] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R Hershey, “Sdr-half-baked or well done?,” in *Proc. ICASSP*, 2019, pp. 626–630.
- [31] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [32] Nicholas D Lane, Sourav Bhattacharya, Petko Georgiev, Claudio Forlivesi, Lei Jiao, Lorena Qendro, and Fahim Kawsar, “DeepX: A software accelerator for low-power deep learning inference on mobile devices,” in *Proc. IPSN*, 2016, pp. 1–12.