



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH



MUSIC SOURCE SEPARATION USING DEEP NEURAL NETWORKS

A Degree Thesis
Submitted to the Faculty of the
Escola Tècnica d'Enginyeria de Telecomunicació de
Barcelona
Universitat Politècnica de Catalunya
by
Ferran López Soler

In partial fulfilment
of the requirements for the degree in
Telecommunications Technologies and Services
Engineering. Major in Audiovisual Systems

Advisor: Jose Adrián Rodríguez Fonollosa

Barcelona, June 2020

Abstract

Last years, Sound Source Separation (SSS) has been one of the most active fields within signal processing. The design of such algorithms seeks to recreate the human ability to identify individual sound sources.

In the music field, efforts are being made to isolate the main instruments from a single audio file with a mixture of stereo audio. The goal of these algorithms is to extract multiple audio files with specific instruments, such as bass, voice or drums.

This project focuses on analyzing the existing systems based on neural networks and their performance. In addition, it goes deeply into the Open-Unmix algorithm structure and tries to improve its results.

Resum

Aquests últims anys, la Separació de Fonts Sonores (SSS) ha estat un dels camps més actius dins del processat de senyal. El disseny d'aquest tipus d'algorismes intenta recrear l'habilitat humana d'identificar fonts sonores individuals.

En el camp de la música, es treballa per aïllar els principals instruments d'un únic fitxer amb una mescla d'àudio estèreo. Així doncs, l'objectiu d'aquests algorismes és obtenir diversos fitxers d'àudio amb instruments concrets, com ara el baix, la veu o la bateria.

Aquest treball se centra a analitzar les propostes existents de sistemes basats en les xarxes neuronals i el seu rendiment. A més, estudia a fons l'estructura proposada en l'algoritme *Open-Unmix* i tracta de millorar els seus resultats.

Resumen

En los últimos años, la Separación de Fuentes Sonoras (SSS) ha sido uno de los campos más activos dentro del procesado de señal. El diseño de estos algoritmos intenta recrear la habilidad humana de identificar fuentes sonoras individuales.

En el campo de la música, se trabaja para aislar los principales instrumentos de un único fichero con una mezcla de audio estéreo. Así pues, el objetivo de estos algoritmos es obtener varios archivos de audio con instrumentos concretos, como el bajo, la voz o la batería.

Este trabajo se centra en analizar las propuestas existentes de sistemas basados en las redes neuronales y su rendimiento. Además, estudia a fondo la estructura propuesta en el algoritmo *Open-Unmix* y trata de mejorar sus resultados.

Acknowledgements

I want to recognize the assistance and advice given by my supervisor Jose Adrián, my student colleagues and other people that helped me during the project process. Furthermore, I want to thank all the open source community for the importance of the task they are doing.

The covid-19 situation has been a difficult period by the presence of external episodes and motivation issues. I would not have finished this project without the support and motivation from my family and friends.

Revision history and approval record

Revision	Date	Purpose
0	03/06/2020	Document creation
1	15/06/2020	Document revision
2	26/06/2020	Document revision

DOCUMENT DISTRIBUTION LIST

Name	e-mail
Ferran López Soler	ferran.lopez.soler@gmail.com
José Adrián Rodríguez Fonollosa	jose.fonollosa@upc.edu

Written by:		Reviewed and approved by:	
Date	03/06/2020	Date	24/06/2020
Name	Ferran López	Name	José Adrián
Position	Project Author	Position	Project Supervisor

Table of contents

Introduction	9
Statement of purpose	9
Requirements and specifications	9
Methods and procedures	9
Work plan	10
Deviations from the initial plan	11
State of the art of the technology used or applied in this thesis	12
Music Source Separation (MSS)	12
Deep Learning and MSS	13
Machine Learning and Deep Learning Fundamentals	13
General description	13
Convolutional layers	15
LSTM modules	16
Batch Normalization	16
Iteration/Loss Training Analysis	17
MSS Algorithms	18
Waveform algorithms	18
Pre-processed waveform algorithms	19
Evaluation metrics	20
Programming Environment	21
PyTorch	21
Methodology / project development	22
Dataset	22
Baseline System	23
Overall structure	23
Specificities	26
Methodology	26
Parameter tuning	27
Improvement proposal	27
Results	29
Parameter tuning	30
Improvement results	33
Budget	37
Conclusions and future development	38

List of Figures

Figure 1.1. Breakdown work structure of the project	10
Figure 1.2. 20-week Time Plan of the Project	11
Figure 2.1. Representation of a music mixture in the time-frequency domain	12
Figure 2.2. Neural Network basic structure	14
Figure 2.3. Rosenblatt's Perceptron Unit	14
Figure 2.4. Kernel filter performance example	15
Figure 2.5. Unrolled RNN with time dependency data	16
Figure 2.6. LSTM implementation	16
Figure 2.7. Error/Iteration network training graphic example	17
Figure 2.8. Proposed Wave-U-Net structure with K sources and L layers	18
Figure 2.9. Deep U-Net convolutional structure	19
Figure 3.1. MUSDB18 dataset stem distribution	22
Figure 3.2. Open-Unmix model structure	23
Figure 3.3. Bidirectional LSTM scheme	24
Figure 4.1. First training results (default parameters)	30
Figure 4.2. Train/Test iteration/loss comparison under dimensionality reduction	32
Figure 4.3. Pre-masking model performance	34

List of Tables:

Table 1.1. Milestone chart distribution	10
Table 4.1. Open-Unmix most important training parameters	29
Table 4.2. Batch Size system performance comparison	31
Table 4.3. Regularization system performance comparison	31
Table 4.4. Dimensionality reduction system performance comparison	32
Table 4.5. Train/Test Iteration/Loss after bandwidth variance	33
Table 4.6. Learnable filter frequency response	34
Table 4.7. Learnable filter model performance	35
Table 5.1. Human resources budget	37
Table 5.2. Software resources budget	37

1. Introduction

1.1. Statement of purpose

This project consists on the usage of deep neural networks to perform music source separation (decomposing music into its constitutive components).

The main goals of the project are:

1. Understanding the configuration and performance of deep learning algorithms
2. Being able to make a diagnostic and purpose solutions from a concrete network training
3. Being able to extract voice from a mixed audio
4. Being able to set solid next steps to improve the system

1.2. Requirements and specifications

The project should perform well on the mixture separation from an audio source. It has to be able to extract voice from a final audio file. The software will be complemented with a detailed description of the network architecture used and the main reasons of the choice.

System performance will be evaluated using concrete metrics and will be compared with state-of-art of Music Source Separation, described in sections below.

1.3. Methods and procedures

The final implementation of the project is based on Open-Unmix Music Separation Software [1]. The training is performed using MUSDB18 Music Source Separation database [2].

The main study results are centered on the voice track. Focusing the problem to a single track extraction has allowed us to do a more consistent study and comparison of the results, especially when tuning the network hyperparameters. The main reasons of the choice and the overall base structure will be defined in next sections.

1.4. Work plan

The final work plan has not changed much from the initial planning. The work has been organized using different work blocks (packages) and following the structure below. As it can be seen, the work organization follows a natural path with the focus on getting knowledge. Software Development and Result Analysis were conceived in an iterative way, so that new results encourage changes to the software parameters or structure.

Breakdown structure

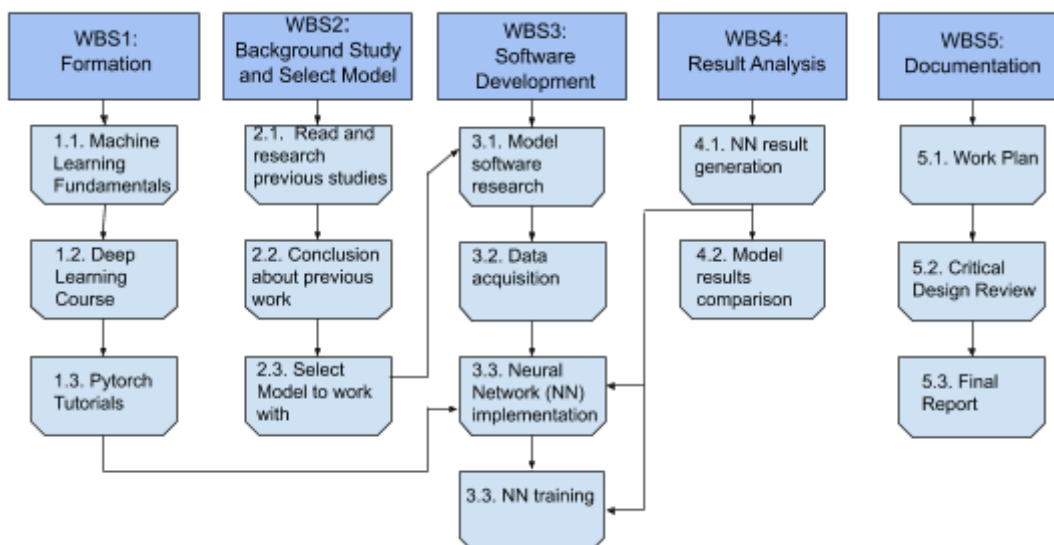


Figure 1.1. Breakdown work structure of the project

The monitoring of this structure has been done by setting milestones. Every milestone had a concrete objective and a fixed date. This is shown in the figure below.

Milestones chart

WP#	Task #	Short title	Milestone / deliverable	Date (week)
1	1	Machine Learning Course	Coursera Official Title	01/02/2020
1	2	Deep Learning Course	Coursera Official Title	23/02/2020
1	3	Pytorch Tutorials	Pytorch Programming Skills	01/03/2020
2	1	Read previous studies	Various Model Study	15/02/2020
2	2	Conclusions about previous work	Model Viability	10/03/2020
2	3	Model Selection	Model Selection	30/03/2020
3	1	Model Software Research	First Model Implementation	01/04/2020
3	2	Data acquisition	Database	01/04/2020
3	3	NN Implementation/Training	Network Model Train	20/04/2020
3	4	Improving Implementations	System Improvements	10/06/2020

4	1	NN Diagnostic	Training graph (Loss/it)	30/04/2020
4	2	NN Evaluation (SDR)	SDR comparison	30/04/2020
5	1	Project Plan	Document	01/03/2020
5	2	Critical Review	Document	20/04/2020
5	3	Final Report	Document	20/06/2020

Table 1.1. Milestone chart distribution

Finally, the Time Plan shows the approximate time assignation of every task in a 20-week period.

Time Plan (Gantt Diagram)

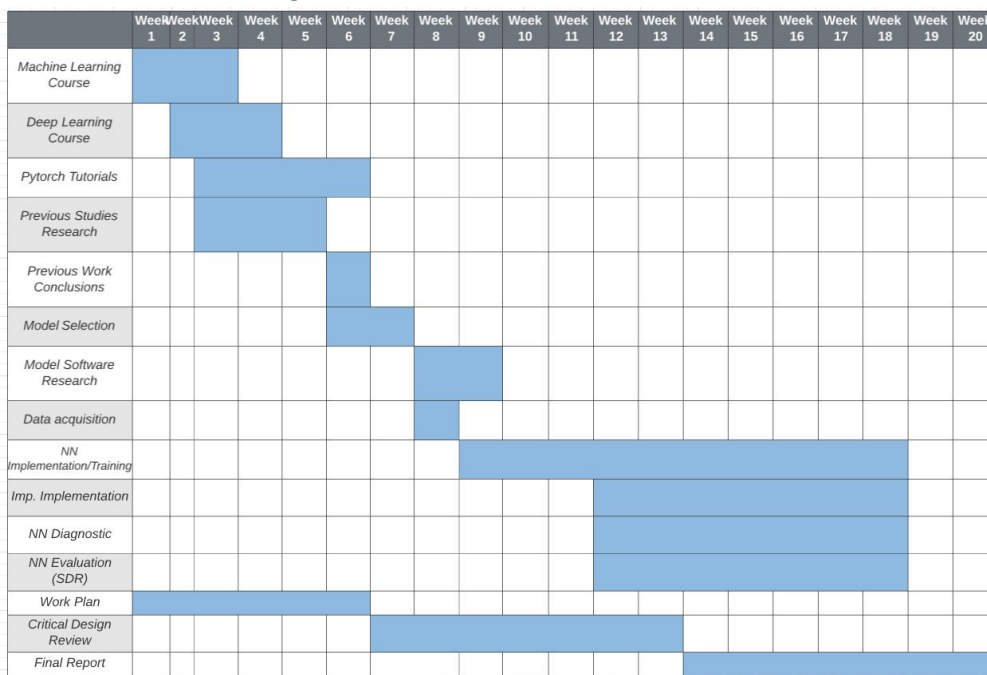


Figure 1.2. 20-week Time Plan of the Project

1.5. Deviations from the initial plan

The initial Work Plan has been affected by some factors related to the chosen source separation model and the installation of a virtual environment. First weeks, all the efforts were put on getting first training results and optimize system training: trying different parameters (like number of parallel workers) and using GPU. This problem affected directly the timeline and delayed obtaining the first results (from week 9 to week 11).

Other problems had appeared in the following weeks. When we needed to adapt the environment to achieve the specifications of the improvement proposal functions, we had several compatibility problems between libraries and Open-Unmix environment. Furthermore, some functions were not available due to problems with the server specifications.

2. State of the art of the technology used or applied in this thesis

2.1. Music Source Separation (MSS)

Last years, Sound Source Separation (SSS) algorithms have been one of the most active fields in signal processing. These algorithms try to recreate the human capacity to differentiate individual sound sources using spatial and frequency component (timbre) information.

Music Source Separation (MSS) is a concrete field into SSS. It tries to split a music audio clip into its constituent contributions (stems), such as the vocals, bass and drums. This separation can be used to re-mix, suppress or up-mix sources from a final recording mix.

Current algorithms are far from providing high quality individual audio stems. Usually, there may be a lot of musical instruments in a mono or stereo recording (only 1 or 2 channel recording), and the sources have been modified by signal filtering or reverberation addition in the mixing process, among others. All of these modifications make MSS a very complex problem. The nature of the problem can be seen in the time-frequency domain.

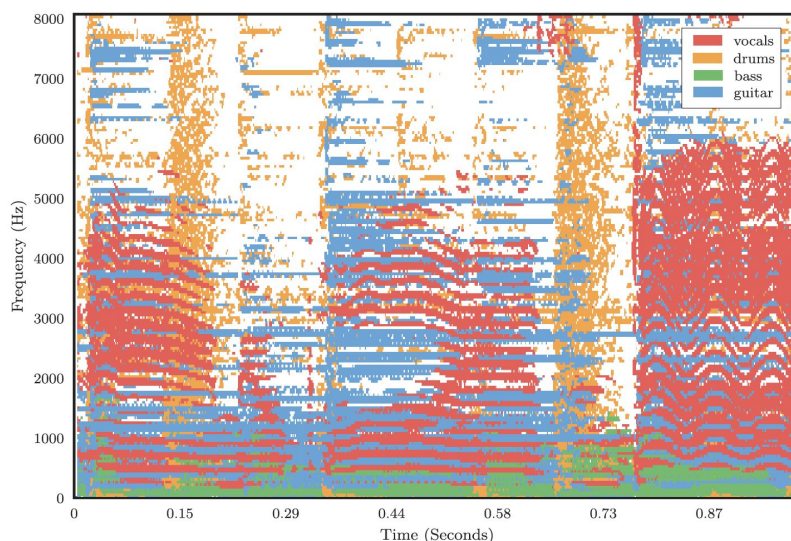


Figure 2.1. Representation of a music mixture in the time-frequency domain. Image from [1]

As it is shown above, the dominant musical source in each time-frequency bin is displayed with a different color. It can be seen that some instruments have a very wide frequency range. Besides, we can see some time-frequency masking problem, which will difficult source separation, especially in tracks like drums or guitar. In the following sections we will define the main strategies for MSS.

2.2. Deep Learning and MSS

Various algorithms have been implemented to achieve a high quality stem segmentation, such as Principal Component Analysis (PCA), defined in [2] or Non-Negative Matrix Factorization, explained in [3]. These algorithms have been lightly overshadowed by the success of Machine Learning (ML) and deep learning (DL) techniques.

Deep learning algorithms are able to model non-linearities and provide faster implementations than previous algorithms. These are usually formulated as a supervised learning problem using concrete targets (such as vocals or drums) and different cost functions. Most of the algorithms use Recurrent Neural Network structures such as LSTM modules to explore time dependencies. In the next section we will put the focus on those modules and other interesting concepts commonly used in MSS.

SiSEC MUS 18 [4] is a SSS evaluation challenge that introduce state-of-the-art algorithms in MSS and is used as starting point to analyze new algorithms.

2.2.1. Machine Learning and Deep Learning Fundamentals

In this section we will explain the key concepts to understand how MSS Deep Learning algorithms work. With this in mind, we will start with a short general description of how a neural network works and then we will explain three of the most common used modules in MSS. Finally, we will briefly explain the fundamentals of training diagnosis.

2.2.1.1. General description

Deep Learning is a subfield of machine learning concerned with algorithms inspired by the structure and function of the brain called artificial neural networks. Furthermore, a Neural Network is combination of several layers of neurons or perceptrons (the basic elementary unit). NNs are made of input and output layers/dimensions, and in most cases, they also have a hidden layer consisting of units that transform the input into something that the output layer can use.

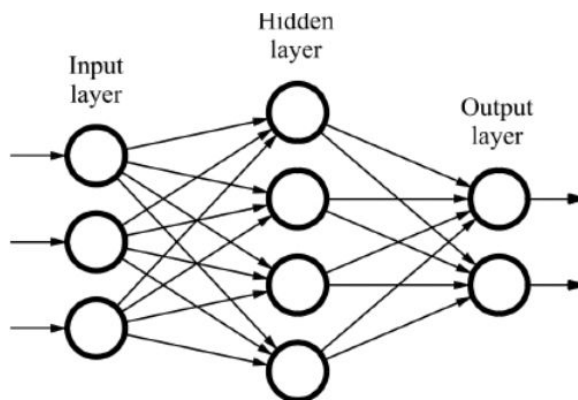


Figure 2.2. Neural Network basic structure. Image from [5]

In a single neuron unit the input data is pondered and transformed in a linear or non-linear way to generate a single output.

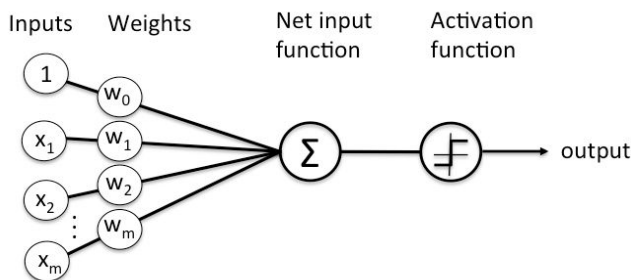


Figure 2.3. Rosenblatt's Perceptron Unit. Image from [6]

The weights are transformed in every iteration until the network outputs are close to the desired outputs. This weight actualization is called back-propagation and it uses the difference between the actual output and the ideal output of every layer (minimization of loss function), from the output layer to the input layer. Some parameters are used in order to achieve a good network performance, such as learning rate, that controls how quickly is the weights adaption, or regularization, which control how much we want to penalize the flexibility of our model by shrinking the coefficients towards zero.

2.2.1.2. Convolutional layers

Convolutional Neural Networks (CNN) were originally designed for image processing. A convolutional layer was designed to reduce the images into a form that is easier to process (dimensionality reduction) without losing features which are essential for a good prediction.

CNN are able to successfully capture the spatial and temporal dependencies through the application of relevant filters (Kernel).

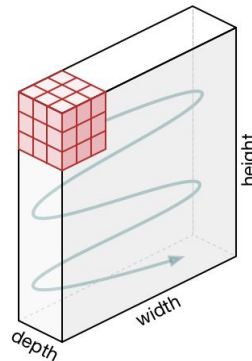


Figure 2.4. Kernel filter performance example. Image from [7]

In the image, we can see that the Kernel (red block) moves through the entire image with a certain Stride Value (positions to be moved after the first position in every iteration). The output in every position is the sum of the element-wise products between the filter and the multidimensional input. During the network training the filter coefficients change to optimize results.

The capacity for dimensionality reduction and capturing spatial and temporal dependence makes CNN a good choice in audio processing. Depending on the input dimensionality, 1D or 2D layers can be used (e.g. raw audio signal or spectrogram).

2.2.1.3. LSTM modules

Long-Short Time Memory (LSTM) modules are special kind of Recurrent Neural Networks (RNN), networks with loops in them, allowing information to persist. LSTMs are capable of learning long-term dependencies between input samples. Furthermore, they are great remembering information for long time periods.

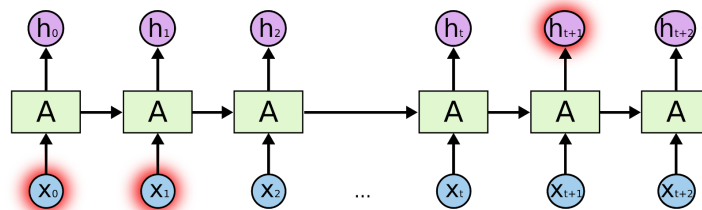


Figure 2.5. Unrolled RNN with time dependency data. Image from [8]

LSTMs have a chain like structure with a repeating module. In one of the most common implementations, each module is composed by a cell, an input gate, an output gate and a forget gate. The cell remembers values over arbitrary time intervals and the three gates regulate the flow of information into and out of the cell.

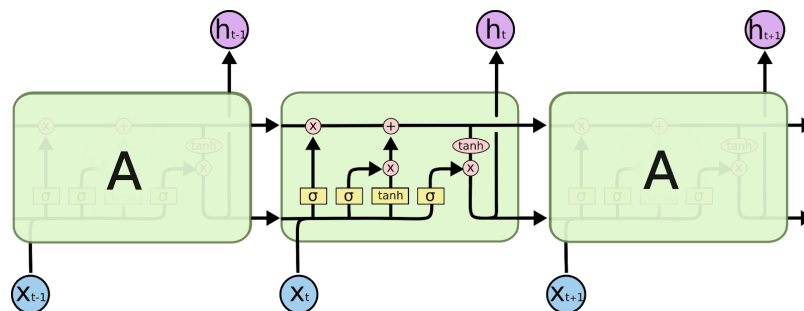


Figure 2.6. LSTM implementation. Image from [8]

2.2.1.4. Batch Normalization

Batch normalization is used to increase the stability of a neural network. The process consists on normalizing the output of a previous activation layer by subtracting the batch mean and dividing by the batch standard deviation. It makes sure that there's no activation that's gone really high or really low and reduces overfitting by adjusting the parameterization of a model in order to make the loss surface smoother.

2.2.1.5. Iteration/Loss Training Analysis

To analyze how good a model performance is, we use the relation between loss function and number of iterations (epochs). A lower loss implies a better data modeling. To understand how good the model generalizing is for new data, the loss is calculated on training and validation (test) and its interpretation is how well the model is doing for these two sets.

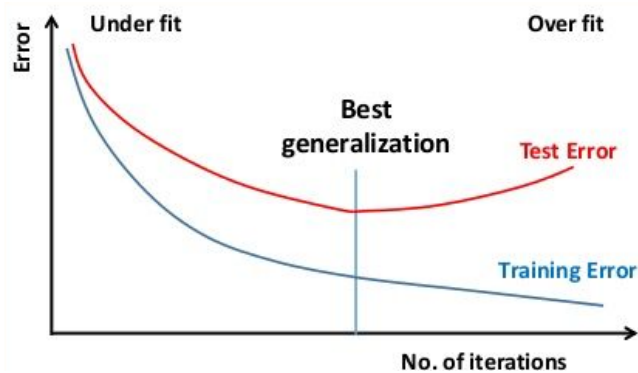


Figure 2.7. Error/iteration network training graphic example. Image from [9]

For network diagnosis we need to define two terms:

Bias

Bias measures whether the average predicted values are far from the actual values. It measures whether the model does or does not capture the complexity of data. It is referred to the training sample.

Variance

Variance measures the difference between the model performance in train and test datasets. High-variance problem tells us that the model is not generalizing well for new data.

When training a model we can get stuck as the model "memorizes" the training examples and becomes kind of ineffective for the test set (over-fitting). Over-fitting also occurs in cases where you have a very complex model or the dataset is not large enough.

On the other hand, under-fitting happens when a model is not able to accurately capture relationships between dataset features and a target variable.

Some measures can be applied to reduce both problems. To reduce over-fitting we can try to get more training examples, reduce the number of features or increase regularization. To solve under-fitting we can try to get additional features or increase regularization.

2.2.2. MSS Algorithms

MSS algorithms are divided in two families. Depending on the input, we can distinguish two types of algorithms: waveform and pre-processed waveform (e.g.: spectrogram) algorithms. We will explain their main differences and most common structures below.

2.2.2.1. Waveform algorithms

Waveform domain algorithms are designed as end-to-end systems. Working in this domain allows modeling phase information, avoiding fixed spectral transformations and low latency calculations.

Most of these algorithms use multiple 1D convolution layers and encoding-decoding structure using downsampling and upsampling. This kind of structure is called *U-Net* structure. It was introduced in biomedical imaging to improve precision and localization of microscopic images. One example could be the *Wave-U-Net* neural network, defined in [10].

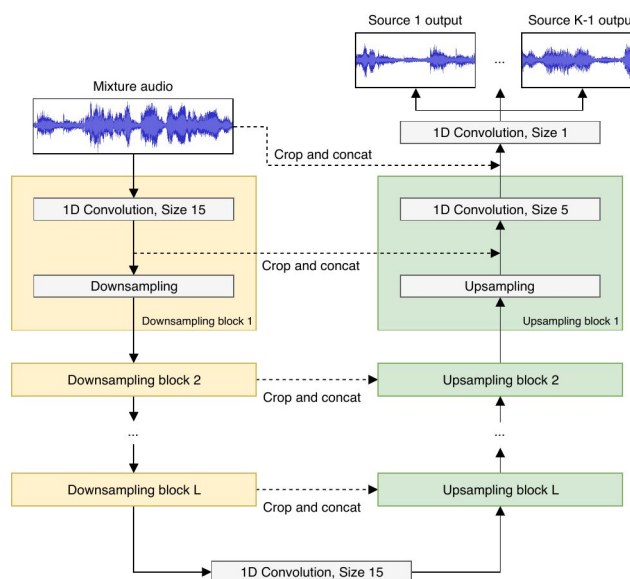


Figure 2.8. Proposed Wave-U-Net structure with K sources and L layers. Image from [10]

We can see that the initial audio mixture is used in the final step through a neural network skipped connection. This provides the magnitude and phase information about the original audio, which is used to obtain final source outputs by applying the mask weights obtained from the network. Most famous algorithms are *Demucs* [11], *Wave-U-Net* and *Conv-Tasnet* [12]. Despite not being the most extended algorithms, last results show the potential of end-to-end systems in MSS and place these models in the forefront of the MSS list on *SISEC18*.

2.2.2.2. Pre-processed waveform algorithms

Most common pre-processed waveform algorithms operate on the spectrograms generated by the Short-Time Fourier Transform (STFT). They produce a mask on the magnitude spectrums for each frame and each source. The output audio is generated by running an inverse STFT on the masked spectrograms. However, the STFT output depends on many parameters, such as the size and overlap of audio frames, and it can affect the time and frequency resolution.

As in waveform-domain, these algorithms usually explore *U-Net* structure networks. But in this case using 2D convolutional layers (e.g. using three dimensional tensors with number of channels, time-steps and frequency bins). One example could be *Spleeter* [13] by using a U-Net Deep Convolutional Layer from [14] as a baseline.

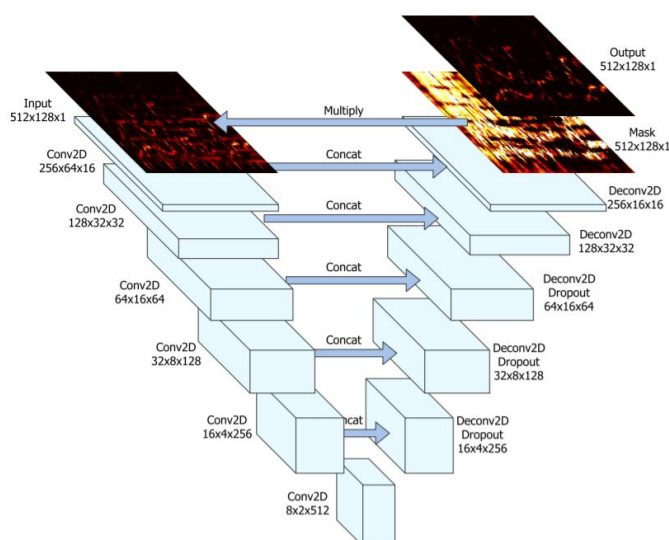


Figure 2.9. Deep U-Net convolutional structure from [14]

Most famous algorithms are *Spleeter* and *Open-Unmix* [15]. *Open-Unmix* will be our proposed approach for the project and it will be explained in detail in Section 3.

2.2.3. Evaluation metrics

Evaluating the results of a music source separation model is a difficult task. However, several objective metrics were developed to facilitate evaluating Blind Audio Source Separation (BASS) algorithms [16]. To this end, the estimated sources \hat{s}_j (with $j = 1 \dots J$) are decomposed as:

$$\hat{s}_j = s_{target} + e_{interf} + e_{noise} + e_{artif}$$

Where s_{target} is a version of the original source, e_{interf} is the interference coming from not desired sources, e_{noise} is the sensor noise and e_{artif} refers to the musical noise self-generated for the separation algorithm.

This decomposition has established the objective metrics and the values we want to maximize:

- SDR (Source to Distortion Ratio), defined as:

$$SDR := 10 * \log_{10} \frac{\|s_{target}\|^2}{\|e_{interf} + e_{noise} + e_{artif}\|^2}$$

- SIR (Source to Interference Ratio), defined as:

$$SIR := 10 * \log_{10} \frac{\|s_{target}\|^2}{\|e_{interf}\|^2}$$

- SAR (Source to Artefact Ratio), defined as:

$$SAR := 10 * \log_{10} \frac{\|s_{target} + e_{interf} + e_{noise}\|^2}{\|e_{artif}\|^2}$$

2.3. Programming Environment

The programming and development of the project has been made using Python and PyTorch.

2.3.1. PyTorch

PyTorch is an open source machine learning framework that accelerates the path from prototyping to production deployment. The PyTorch Torch package contains data structures for multi-dimensional tensors (data container) and mathematical operations over these are defined. Its principal modules are torch.nn or torch.optim. The nn modules in PyTorch provide us a higher level API to build and train deep network.

3. Methodology / project development

The project development has been centered into the optimization of the Open-unmix software performance. The system has been trained using the MUSDB18 dataset. The main reasons of the choice will be explained in detail in the next sections.

3.1. Dataset

The SigSep MUSDB18 data set [17] consists of a total of 150 full-track songs of different styles and it includes both the stereo mixtures and the original sources, divided into a training subset and a test subset. All the tracks are stereophonic and encoded at 44,1kHz.

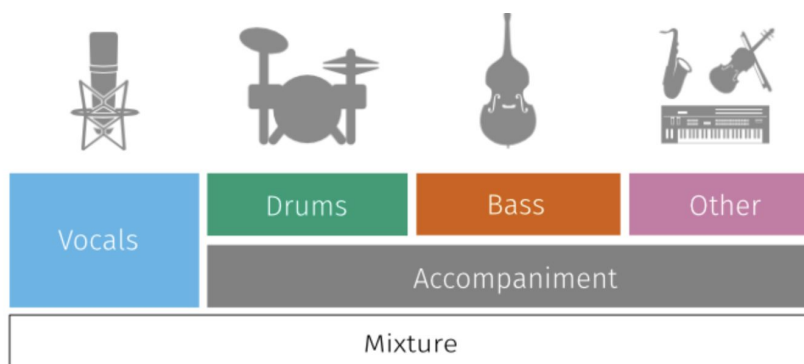


Figure 3.1. MUSDB18 dataset stem distribution. Image from [17]

Its purpose is to serve as a reference database for the design and the evaluation of source separation algorithms. The objective of such signal processing methods is to estimate one or more sources from a set of mixtures, e.g. for karaoke applications. It has been used as the official dataset in the professionally-produced music recordings task for SiSEC 2018 [4].

The database is presented in compressed and uncompressed (high quality) format. The compressed files have a bandwidth limited to 16 kHz. But nevertheless, it seems not affecting to the evaluation performance.

MUSDB18 is actually the most complete database for MSS. Furthermore, all the systems from state-of-art have been using this dataset in order to compare and share their results.

3.2. Baseline System

The proposal approach is based on Open-Unmix Architecture [15] [18]. We will describe its main advantages and the main reasons of this choice below.

3.2.1. Overall structure

Open-Unmix is an open source MSS software that works with audio spectrograms. Its main advantages are simplicity and fast training.

Its architecture is very simple. It is based on LSTM modules. We will describe its main modules using the Open-Unmix paper description and some considerations.

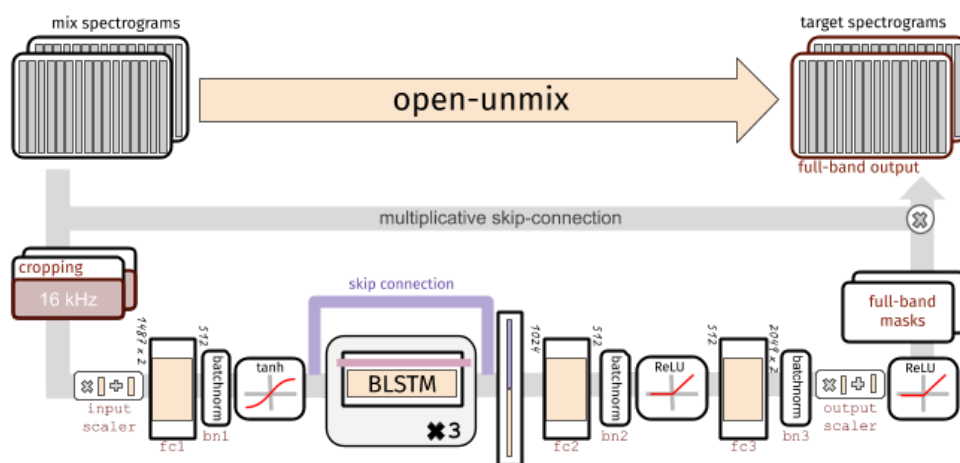


Figure 3.2. Open-Unmix model structure. Image from [18]

Input Stage

Open-Unmix operates in the time-frequency domain to perform its prediction. The input of the model is either:

- Time domain signal tensor (later transformed with STFT)
- Magnitude spectrogram tensor (e.g. when pre-computed and loaded from disk)

First of all, the audio input is chunked into 6 seconds excerpts. It guarantees the audio correlation between samples without compromising the computational cost. In one iteration (epoch), 64 samples in random position are selected from each track to ensure a balanced track sampling.

The input spectrogram is cropped at 16kHz (in relation with the compressed MUSDB 18 bandwidth) and standardized using the global mean and standard deviation for every frequency bin across all frames. Furthermore, batch normalization is applied in multiple stages of the model to make the training more robust against gain variation.

Dimensionality reduction

The LSTM is not operating on the original input spectrogram resolution. In the first step, the network learns to compress the frequency and channel axis of the model to reduce redundancy and make the model converge faster. Fully connected time-distributed layers are used for dimensionality reduction and augmentation, thus encoding/decoding the input and output.

Bidirectional LSTM

The core of open-unmix is a three layer bidirectional LSTM network. The figure below shows the structure of a single bidirectional LSTM layer.

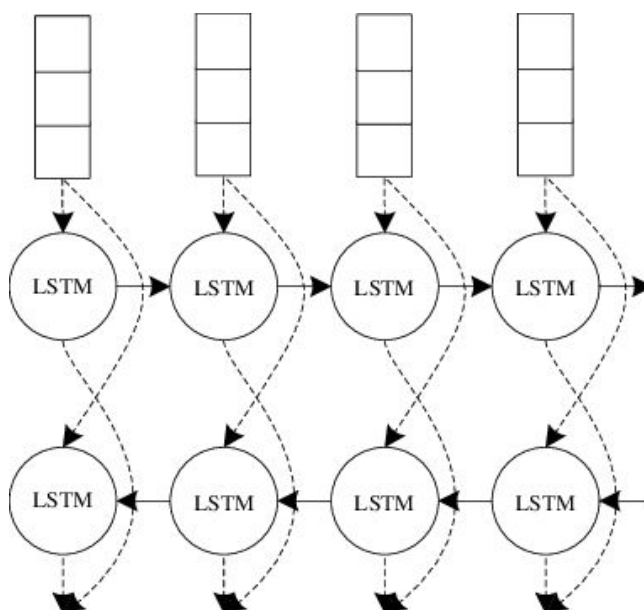


Figure 3.3. Bidirectional LSTM scheme. Image from [19]

As it is shown above, the information of every input sample is used in LSTM cells to work with the previous or next LSTM cell. Since the model takes information from past and future simultaneously, the model cannot be used in an online/real-time manner.

Skipped Connection

Skipped Connections are used in two ways:

- The output to recurrent layers are augmented with their input, and this proved to help convergence.
- The output spectrogram is computed as an element-wise multiplication of the input. This means that the system has to learn how much each TF bin does belong to the target source.

Output Stage

After applying the LSTM, the signal is decoded back to its original input dimensionality. In the last steps the output is multiplied with the input magnitude spectrogram, so that the models are asked to learn a mask. In this final step the output signal is synthesized by inverse STFT.

For inference, the signal is post-processed with an implementation of a multichannel Wiener filter that is a very popular way of filtering multichannel audio for several applications. This filtering method assumes you have some way of estimating power or magnitude spectrograms for all the audio sources composing a mixture. Norbert allows us to create a residual model when working with individual sources.

As described in the technical details on Open-Unmix SigSep documentation, different activation functions are used.

- Rectified linear units (ReLU) allow intermediate layers to comprise nonnegative activations.
- Tanh are necessary for good training of LSTM model, notably because they avoid exploding input and output.
- Sigmoid activation is chosen to mimic the way legacy systems take the output as a filtering of the input.

3.2.2. Specificities

In this section we will describe the main advantages of the Open-Unmix software and the main reasons why we have chosen it.

The design has been oriented to reach two main objectives:

- To have state-of-art performance
- To be easily understandable (for research purposes)

Furthermore, other specificities have been accomplished. The code is:

- Simple to extend: The pre/post-processing, data-loading, training and models part of the code is isolated and easy to replace/update.
- Not a package: The software is composed of largely independent and self-containing parts, keeping it easy to use and easy to change.
- Hackable (MNIST like): Open-unmix mimics the famous MNIST example, including the ability to instantly start training on a dataset that is automatically downloaded.
- Reproducible: Releasing Open-Unmix attempts to provide a reliable implementation sticking to established programming practice.

The determinant factors for choosing this architecture were its training speed, which allowed us to train multiple times and have results within two or three days, and its simplicity, which gave us a concrete vision of all the network hyperparameters and its functions.

3.3. Methodology

The methodology followed during the project has been based on an alternation between result analysis and improvement proposal and development. After the generation of the first training results the work has been focused in two main blocks: network parameter and hyperparameter tuning and network improvements proposal. All the changes were tested in the same target to really evaluate its impact.

The network training has been done using the UPC Calcula service. This allowed us to streamline the training by using GPU.

We have decided to center the study of the systems in terms of Training/Test loss instead of the metrics explained in 2.2.3. This decision is based on the computational cost and time cost of the metrics calculation and it is also based on the directly proportional relation between SDR and best epoch test loss. This decision has allowed us to increase the number of experiments and has helped especially in parameter tuning.

The training methodology and the Open-Unmix data processing show some result differences between two identical parameter trainings. For this reason, we have trained each model three times and we have selected the best performance.

3.3.1. Parameter tuning

As it has been seen before, *Open-Unmix* architecture has some parameters that are interesting to study.

After analyzing the first results, a network diagnostic has shown possible parameter changes. We have studied the algorithm performance under variations of the batch size, the dimensionality reduction in the fully-connected layer (input of the LSTM) and the weight decay for regularization.

We have centered the study on the most common parameters or the parameters that we thought to be most important for the training performance.

3.3.2. Improvement proposal

The improvement proposal and implementation have been conditioned by the Open-Unmix main specificities. We have tried to look for modifications that would not affect directly on the computational cost or the understandability of the architecture. For this reason, the work has been centered on searching pre-processing tools to improve the overall system performance.

In pre-processing we have worked to emphasize the selected target in the input of the network. We have tried fixed filters and learnable filters to explore how equalization can affect to a concrete target extraction. We also tested to reduce the frequency threshold of the maximum bandwidth processed by the LSTM (input stage frequency cropping).

On post processing, we have explored but not implemented the use of phase information in the frequency-time conversion, using some algorithms for phase reconstruction such as Griffin&Limm algorithm. We have decided not to put into post-processing because the evaluation metrics previously defined in 2.2.3 are not efficient to compare obtained signals sample by sample, and other metrics would be necessary. Furthermore, post-processing did not have repercussion to the whole training performance.

4. Results

This section will include the data analysis and findings. We will start with an analysis of the first obtained train results. Then, these results will be used for parameter tuning orientation. Finally, we will analyze the post-processing improvement proposal results.

The default most important training parameters and hyper-parameters are shown in the table below. Other network parameters can be found in [18].

Argument	Description	Default
<code>--batch-size <int></code>	Batch size has influence on memory usage and performance of the LSTM layer	16
<code>--seq-dur <int></code>	Sequence duration in seconds of chunks taken from the dataset. A value of ≤ 0.0 results in full/variable length	6.0
<code>--hidden-size <int></code>	Hidden size parameter of dense bottleneck layers	512
<code>--lr <float></code>	learning rate	0.001
<code>--weight-decay <float></code>	weight decay for regularization	0.00001
<code>--bandwidth <int></code>	maximum bandwidth in Hertz processed by the LSTM. Input and Output is always full bandwidth!	16000
<code>--nfft <int></code>	size of fft	4096

Table 4.1. Open-Unmix most important training parameters

After the first training attempts, the first results were generated. We defined next steps after analyzing its performance. Concretely, the iteration/loss graphic.

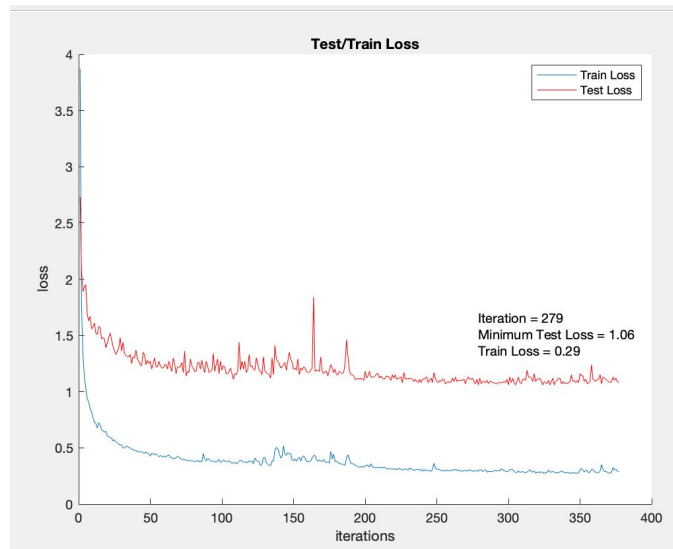


Figure 4.1. First training results (default parameters)

As it can be observed in the figure above, we can assume that the network is suffering a high-variance (overfitting) problem. We can also appreciate a high bias problem, but we decided reducing Test Loss as the main objective, and usually reduce bias involve an increment of variance. As it has been explained before, some actions can be taken to reduce this problem. We have centered on understanding the network parameters and its functions.

In the next sections we will see the performance of parameter tuning and improvement proposal using the network training performance.

4.1. Parameter Tuning

In this section we have tested the performance of Open-Unmix model by changing some interesting parameters. The changes were directly oriented to reduce overfitting and to try to reduce the test loss or make test/train loss closer.

Batch Size

There is a tension between batch size and the speed and stability of the learning process. It is also related with the estimation accuracy, so we decided to compare Open-Unmix performances on different batch sizes.

Batch Size	Total epochs trained*	Best Epoch Test Loss
16	481	1.06
32	517	0.98
64	530	0.99
128	327	1.1

Table 4.2. Batch Size system performance comparison

In this case we cannot see relevant difference between trainings in terms of converge speed or stability. Despite of that, we can see an important difference in time per epoch, so increasing batch size accelerates the whole training process. We also have certain advances in terms of Test Loss using 32 or 64 batch size values. Using too low values may affect to the gradient estimation. In contrast, using too high values can produce a degradation in the model generalization ability.

Regularization parameter tuning

Regularization is commonly used in NN trainings to reduce overfitting. We tried to increase the weight decay regularization in order to reduce the high-variance problem.

Regularization	Train Loss	Test Loss	SDR
0.00002	0.24	1.02	5.514
0.00004	0.27	1.03	5.466
0.0001	0.35	1.08	5.124
0.001	0.8	1.50	0.548

Table 4.3. Regularization system performance comparison

We can see that the distance between Train Loss and Test (Validation) Loss reduces as regularization increase. The training bias also increase, so regularization tuning it is not useful as an independent measure but it could be interesting when training the system with an augmented dataset, for example. We decided to show SDR values to emphasise that too high regularization values can slash SDR metric.

Dimensionality reduction

If we analyze the dimensionality reduction in the fully-connected layer (input of the LSTM), we can see that information compression to reduce redundancy and make the model converge faster is not optimal in terms of training results.

LSTM input size	Total epochs trained*	Best Epoch Train Loss	Best Epoch Test Loss
64	433	0.43	1.32
128	638	0.32	1.11
256	649	0.25	1.03
1024	483	0.21	0.97

Table 4.4. Dimensionality reduction system performance comparison

This results show that dimensionality reduction is not affecting or does not seem relevant in terms of convergence velocity. When comparing the Test Loss between the different models, we can affirm that higher input size is a good choice to reduce variance. We can also see that bias decreases when LSTM input size increases.

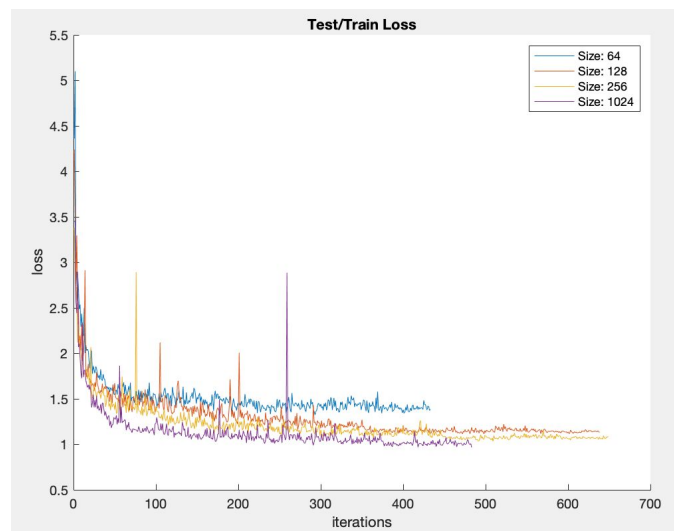


Figure 4.2. Train/Test iteration/loss comparison under dimensionality reduction

4.2. Improvement results

As it has been explained before, the improvements have been focused on filtering for source enhancement. We have tried learnable filters of different number of coefficients in the pre-processing stage of the architecture. We implemented a pre-masking stage by assigning learnable weights to the every spectrogram frequency bin. We also analyzed an internal parameter of the network, the frequency bandwidth of the input. We decided to put its results here because it is directly related with spectrum pre-processing.

We have studied the importance of the frequency cropping pre-processing stage. For the voice track we thought that high frequencies were essential for the source extraction. However, results show that we can reduce a lot the frequency bandwidth and inclusive improving Test performance. In the figures below we can see the network performance reducing the bandwidth to 15 kHz, 13 kHz, 11 kHz and 9 kHz (from upper-left corner to lower-right corner).

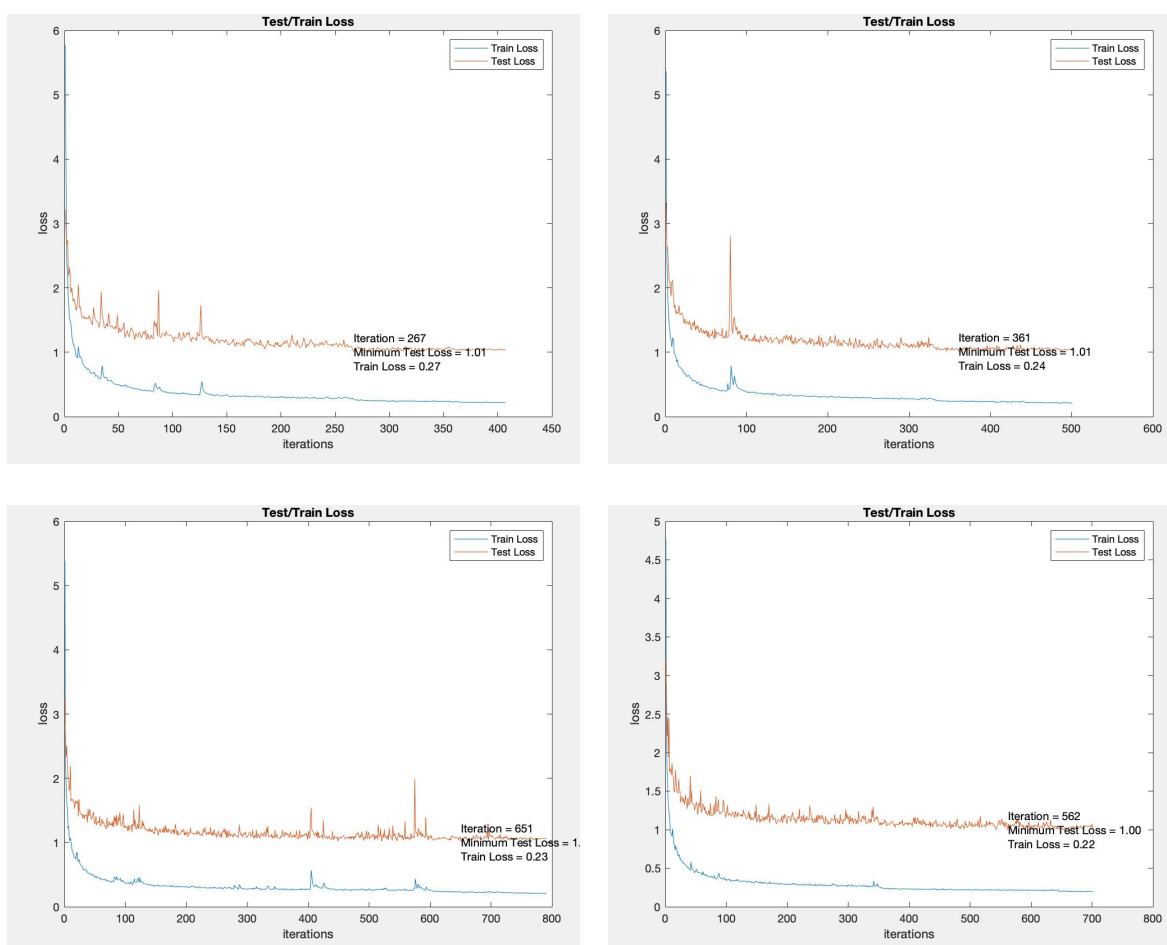


Table 4.5. Train/Test Iteration/Loss after bandwidth variance. Bandwidth: 15 kHz, 13 kHz, 11 kHz and 9 kHz (from upper-left corner to lower-right corner)

Before implementing the filters we designed a learnable pre-masking stage. We defined learnable weights applied to every frequency bin of the fft input. As it can be seen below, the training results are similar to other results showed in this section, however the training convergence velocity increases, so we decided not to get stuck into this implementation.

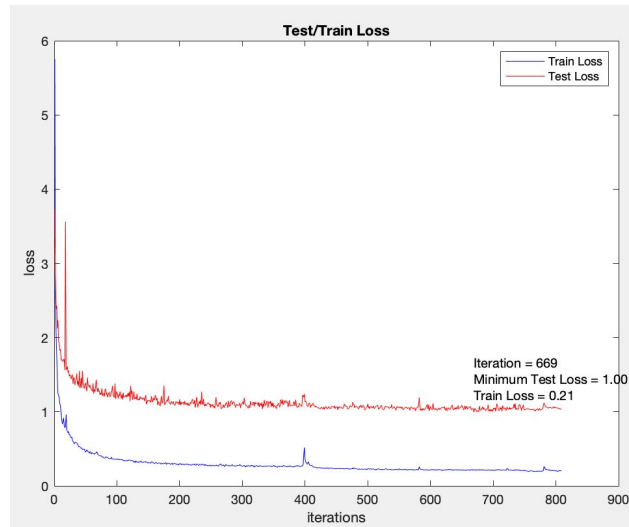
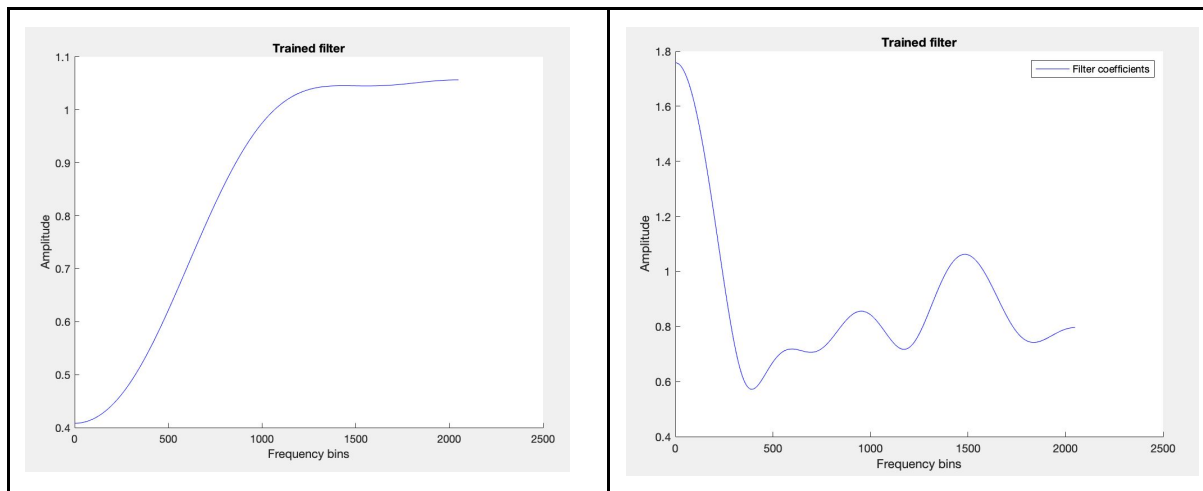


Figure 4.3. Pre-masking model performance

As it is explained in other sections, we have designed a pre-processing learnable filter to try to enhance a concrete target. We have studied its performance with different filter coefficient numbers.



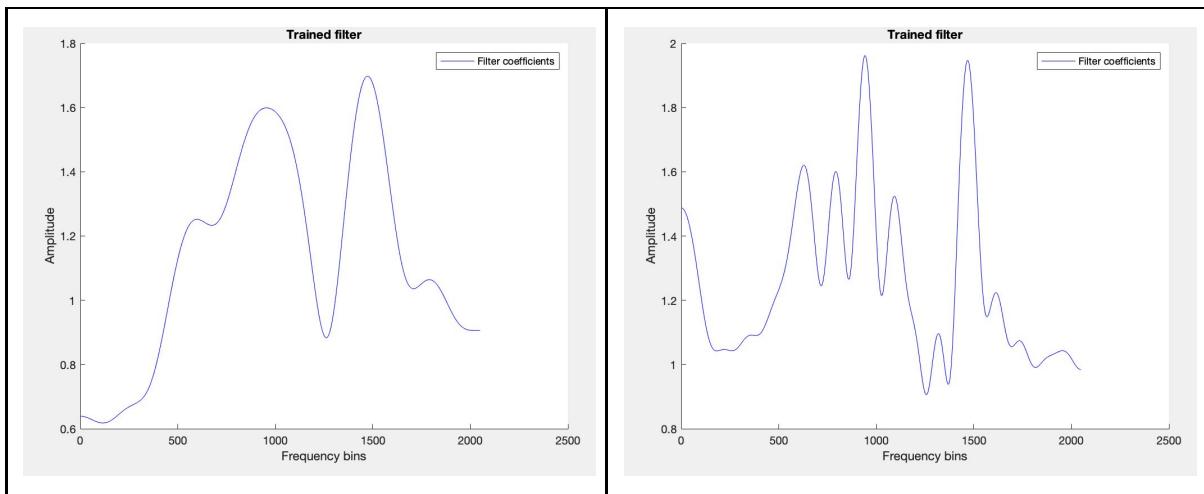


Table 4.6 Learnable filter frequency response of 5, 10, 15 and 35 coefficients (from upper-left corner to lower-right corner)

Number of filter coefficients	Best Epoch Train Loss	Best Epoch Test Loss
5	0.29	1.10
10	0.27	1.07
15	0.21	1.00
35	0.21	1.00

Table 4.7. Learnable filter model performance

From the results showed above, we can make some interesting considerations. As the number of coefficients increase, the system seems to have a slightly better performance. Furthermore, increasing a lot the number of coefficients does not provide relevant improvement.

If we analyze the frequency response and of the filters, we can assert that high-pass filtering could be a good option as pre-processing to extract the voice track. This results are comparable with filtering proposal from professional mixing techniques as [20], concretely with the 15 coefficient filter response.

As we can see, the filter response variates a lot between trainings. However, we can see a gap around 12 kHz in all the figures. This frequency range could be directly related with drum cymbals or other instruments, concluding that its presence is not relevant to extract the voice target.

As this filter was applied before the frequency bandwidth cropping, this gap can also be considered as a high frequency attenuator. Furthermore, we can see that the filters also have a downward trend starting at bin 1500, approximately 16 kHz frequency, which it is probably learned automatically from the network specifications during the training.

5. Budget

To analyze the budget of this project, we consider human and software resources. Both resources can be seen in detail in the tables below.

HUMAN RESOURCES

Description	Cost per hour	Number of hours	Total cost (20 weeks)
Junior Engineer*	10 €/h	18 hours / week	3600 €
Advisoring	20 €/h	1 hour / week	400 €
TOTAL			4000 €

Table 5.1. Human resources budget

*For this cost estimation we have considered the salary of a non-Master graduated student

SOFTWARE RESOURCES

In this project we mostly used open source software options, reducing the license amortization costs.

Description	Cost per licence	Number of licences	Total cost
Matlab Student Subscription*	70 €	1	70 €
TOTAL			70 €

Table 5.2. Software resources budget

*As Matlab Student Subscription is personal and intransferible, it is excluded from amortization. It is an annual subscription simulated with a student doing its last year.

6. Conclusions and future development

Improving Open-Unmix performance has been a difficult task to manage. All the modifications were carefully studied and chosen in order to improve the system performance without compromising the training speed and the overall system simplicity and understandability. Therefore, finding and implementing modifications has been limited by the system main advantages.

The result comparison with the state-of-art is affected by the performance of the system post-processing Norbert Wiener Filter. As explained in 3.2.1, Norbert filter works with all the individual stems obtained from the mix (vocals, bass, drums and others). When the filter does not have other stem information, it uses an internal approximation of the mix without the voice. This adaptation makes the metrics differ from the metrics found in the state-of-art of the project. So centering the trainings on the voice track has make the performance of the Norbert filter poorer. This is the one of the reasons because of we have used the training loss/iteration metric to compare the different system performances.

Some implementations and purposes have been excluded from the results because it directly affected to the characteristics explained above. A clear example would be a pre-processing temporal filtering design that dragged the training time from two days to more than a month and a half. The sequential filtering was specially slow in GPU. We tried to move it into the pre-processing stage but it did not perform as expected in terms of training loss and speed.

Moreover, parameter tuning is a very complex task. Finding the optimal values of the parameters chosen has been a challenging task. Parallel training has helped to reduce the waiting time. Nevertheless, it has been difficult to control the performance when changing the different parameters (two or more parameters) in the same training example. Helpfully, Open-Unmix is a relatively small network and the internal network parameters to tune were not too hard to study separately. We have centered the study on the most common parameters or the parameters that we thought to be most important for the training performance.

The results show a slightly upgrade in source separation performance by using filters. In spite of that, more accurate filtering or parametric equalization could work better for this concrete problem. As it is mentioned in 4.2, reducing the frequency bandwidth could help with source extraction in tracks like bass or voice, but with tracks with a very wide frequency range such as drums it will not. Other alternatives such as panning coefficients [21] or other stereo audio based source enhancement methods could be implemented for a better separation. Music repetitive structure has also been used to help source separation, as in [22].

Filtering would be an interesting option to try in the post-processing stage to reduce final audio artificial interference. Furthermore, different spectrogram signal reconstruction, such as Griffin&Lim, or more complex structures, such as wavenets [23], could perform better obtaining the final unique source audio.

Up to now, end-to-end systems seem to perform better in MSS than pre-processing waveform algorithms. This could change introducing some improvements or exploring different features. Concretely, in spectrogram based algorithms efforts are made to use the FFT phase information as an input or extra feature of the system to help with the source extraction.

It is well known that this kind of problems depend a lot on how much data are available to train the system. So database augmentation is a must when we think on the next steps. Various techniques have been followed in the state-of-art of this project, such as pitch modification or time stretching. Another option would be a multiple dataset adaptation or directly obtaining new tracks in a further research task.

The main objectives of the project have been accomplished. Focusing on the voice track has helped in the result comparison but has also compromised the global vision of source separation algorithms on other tracks like bass or drums. Taking a chance on open source algorithm implementations has been a good choice. Open source community has helped a lot in terms of research and information accessibility.

Bibliography:

- [1] Cano, E., Fitzgerald, D., Liutkus, A., Plumbley, M., Robert-Stöter, F. (2019). Musical Source Separation: An Introduction. *IEEE Signal Processing Magazine*, 36 (1), 31-40.
- [2] Abdi. H., Williams, L.J. (2010). Principal component analysis. *Wiley Interdisciplinary Reviews: Computational Statistics*, 2 (4), 433-459.
- [3] Inderjit S. Dhillon and Suvrit Sra. (2005). Generalized nonnegative matrix approximations with Bregman divergences. In *Proceedings of the 18th International Conference on Neural Information Processing Systems (NIPS'05)* (283–290). Cambridge, MA, USA: MIT Press.
- [4] Stöter, E., Liutkus, A., Ito, N. (2018). *14th International Conference on Latent Variable Analysis and Signal Separation (LVA/ICA 2018)*.
- [5] Quiza, R., Davim, J. (2009). *Computational modeling of machining systems*.
- [6] Nicholson, C. (n.d.). *A Beginner's Guide to Neural Networks and Deep Learning*. Available: <https://pathmind.com/wiki/neural-network> [Accessed: 1 June 2020]
- [7] Saha, S. (2016). *A Comprehensive Guide to Convolutional Neural Networks — the ELI5 way*. Available: <https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53> [Accessed: May 2020]
- [8] Olah, C. (2015). *Understanding LSTM Networks*. Available: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/> [Accessed: May 2020]
- [9] Microsoft Ventures. (2014). *Intro to Machine Learning*. Available: <https://www.slideshare.net/microsoftventures/microsoft-ventures-workshop> [Accessed: 10 June 2020]
- [10] Stoller, D., Ewert, S., Dixon, S. (2018). Wave-u-net: A multi-scale neural network for end-to-end audio source separation. *Proc. Int. Soc. Music Inf. Retrieval*, (334–340).
- [11] Défossez, A., Usunier, N., Bottou, L., Bach, F. (2019). *Music Source Separation in the Waveform Domain*. Available: <https://hal.archives-ouvertes.fr/hal-02379796/document> [Accessed: April 2020]
- [12] Luo, Y., Mesgarani, N. (2019). Conv-TasNet: Surpassing Ideal Time–Frequency Magnitude Masking for Speech Separation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27.8 (1256–1266).
- [13] Hennequin, R., Khlif, A., Voituret, F., Moussallam, M. (2019). *SPLEETER: A FAST AND STATE-OF-THE-ART MUSIC SOURCE SEPARATION TOOL WITH PRE-TRAINED MODELS*. Available: <https://archives.ismir.net/ismir2019/latebreaking/000036.pdf> [Accessed: April 2020].
- [14] Jansson A., Humphrey E. J., Montecchio, N., Bittner R. M., Kumar, A., Weyde, T. (2017). Singing Voice Separation with Deep U-Net Convolutional Networks. *ISMIR*, (745-751).
- [15] Stöter F., Uhlich, S., Liutkus, A., Mitsufuji, Y. (2019). *Open-Unmix - A Reference Implementation for Music Source Separation*. Available : <https://sigsep.github.io/open-unmix/> [Accessed: March-June 2020]
- [16] Vincent, E., Gibbonval R., Fevotte, C. (2006). Performance measurement in blind audio source separation," in *IEEE Transactions on Audio, Speech, and Language Processing*, 14-4 (1462-1469).

- [17] Rafii, Z., Liutkus, A., Stöter, F., Mimitakis, S. I., Bittner, R. (2017) *MUSDB18 - a corpus for music separation*. <https://sigsep.github.io/datasets/musdb.html> [Accessed: March-June 2020]
- [18] Stöter, F., Liutkus, A., Inria and LIRMM. (2019). *Open-Unmix - A Reference Implementation for Music Source Separation*. Available: <https://sigsep.github.io/open-unmix/#using-the-pytorch-version> [Accessed: March-June 2020]
- [19] Hu, Z., Ma, X., Liu, Zhengzhong H., Xing, E. (2016). *Harnessing Deep Neural Networks with Logic Rules*. <https://arxiv.org/pdf/1603.06318.pdf> [Accessed: April 2020]
- [20] Moxey, J. (n.d.). *Songstuff Music Resources - EQ Frequencies*. Available: https://www.songstuff.com/recording/article/eq_frequencies/ [Accessed: June 2020]
- [21] Avendano, C. (2003). Frequency-domain source identification and manipulation in stereo mixes for enhancement, suppression and re-panning applications. *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*. (55-58).
- [22] Rafii, Z., Pardo, B. (2011). A simple music/voice separation method based on the extraction of the repeating musical structure. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. (221-224).
- [23] Oord, A. V. D., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., Kavukcuoglu, K. (2016). *WaveNet: a generative model for raw audio*. <https://arxiv.org/pdf/1609.03499.pdf> [Accessed: May 2020]

Glossary

SSS	Sound Source Separation
MSS	Music Source Separation
NN	Neural Network
ML	Maschine Learning
DL	Deep Learning
PCA	Principal Component Analysis
LSTM	Long-Short Time Memory
CNN	Convolutional Neural Networks
RNN	Recurrent Neural Networks
STFT	Short-Time Fourier Transform
SDR	Source to Distortion Ratio
SIR	Source to Interference Ratio
SAR	Source to Artefact Ratio
ReLU	Rectified Linear Unit