

Polyphonic music separation based on the Simplified Energy Splitter

KRISTÓF ACZÉL¹, ISTVÁN VAJK²

Department of Automation and Applied Informatics

Budapest University of Technology and Economics

3-9. Muegyetem rkp. , H-1111, Budapest

HUNGARY

aczekri@aut.bme.hu¹, vajk@aut.bme.hu²

Abstract: - In the past years many approaches have been developed that target the separation of polyphonic music material into independent source signals. Due to lack of information on the original signals it is currently practically impossible to extract the original waveforms from their mixture. Thus all of the approaches target the reconstruction of signals that are at least in some way close to the original. For that purpose common features of harmonic sounds are usually exploited. This paper proposes a system that uses frequency-domain instrument models as prior knowledge for reinserting information needed for the separation. The system provides over 18dB Signal to Distortion Ratio for two simultaneous notes, which slowly degrades as the level of polyphony increases. This makes the approach highly applicable both as a standalone separation tool and the ground of other signal manipulation methods.

Key-Words: - sound separation, instrument print, polyphonic music, energy split

1 Introduction

When we use the term ‘musical recording’, in technical terms we typically refer to a polyphonic material that is the composition of many individual notes originating from a number of instruments. The recording may or may not be made by using many microphones. However, even in the former case the output signals of the microphones are finally downmixed to fewer, typically one or two channels. Once the recording is downmixed, no correction can be made to the individual notes in it. Our long term interest in sound separation is motivated by the problem of correcting these recordings if they have incorrect notes (either in frequency, volume, intonation, length etc).

In this paper we propose a system architecture for the separation of polyphonic music to individual instrument note signals. Descriptions of the building blocks are also covered. We allow a reasonable amount of user input and processing time to achieve better separation quality. User input involves entering the musical score (note starting/ending times, frequencies, used instruments). Although this input will never be 100% accurate due to the nature of real-life music, it can be precise enough for getting a first estimate on the note parameters in the recording, thus we do not need to rely on musical transcription ([1], [2], [3], [4], [5]) and instrument recognition ([7], [8], [9], [10]) algorithms which are inferior to the performance of a human listener.

The complexity of separation can be attributed to the fact that the information to be retrieved is actually missing from the signal. This issue made researchers approach the problem in many different ways. In [11] a sound source separation algorithm is shown that requires no prior knowledge on the instrument notes in the recording, and performs the task of separation based purely on azimuth discrimination within the stereo field. The results are impressive. However, separating individual notes is not in the focus of the paper, only instrument groups are differentiated.

[13], [14], [15] describe a method which separates harmonic sounds by applying linear models for the overtone series of the sound. The method is based on a two-stage approach: after applying a multipitch estimator to find the initial sound parameters, more accurate sinusoidal parameters are estimated in an iterative procedure. Separating the spectra of concurrent musical sounds is based on the spectral smoothness principle [12].

Beamforming techniques [16], [17] along with the Independent Component Analysis framework offer a different way of separation. A large array of microphones is employed for recording an event. The travel time of the time and the difference of the recorded signals are used in calculations that increase the receiver sensitivity in the direction of signals of interest and decrease it in other directions. Beamforming relies on certain preliminary conditions and studio setup to achieve good results.

Several other approaches exist to sound source separation, (like Non-Negative Matrix Factorization [18], [19], sparse coding [5], [6], etc.). This article does not target to provide a complete list of these methods. As a good ‘state of the art’ a detailed overview can be found in [20].

The following sections guide the reader through the separation process that is proposed in this paper. Each section covers one important block of the algorithm. Section 2 shows an overview of these blocks and the signal flow in the system. The whole separation process is carried out in frequency domain, thus the first step is converting the signal. This step is covered in Section 3 along with the inverse transformation used at the end of the separation process. Section 4 describes the instrument model that is used in the system, while section 5 covers the details of the actual separation process, the Simplified Energy Split (SES) method. Missing information on note intonation (‘playmode’) and volume levels must be calculated prior to the actual separation. Section 6 deals with a simple solution for the problem.

2 Overview of the separation process

This section shows an overview of the separation process. Short descriptions of the building blocks are given which are discussed later in detail.

First, all time-domain signals are transformed to frequency domain. Windowed FFT is employed for this purpose along with Frequency Estimation [25] and Phase Memory methods.

The separation is aided by stored instrument samples. These samples are stored one by one in bandogram format, which is basically a spectrogram split to subbands, in which the energy is summed.

Bandograms originating from instrument notes are stored in the sample store. One instrument will have a number of samples differing by their frequency and playmode. This allows for the selection of the sample that best matches the note in the recording to help the separation process. The collection of bandogram samples from the same instruments is called an *instrument print*.

The playmode and volume detector receives the original recording, the musical score and the instrument prints. Its role is to select one sample from each instrument print that best fits the original instrument in the recording.

The Simplified Energy Splitter gets the spectrogram of the recording and the bandogram of the selected samples as input. It splits the energy in the spectrogram of the original recording to components that resemble the input bandograms.

Finally, the frequency-domain signals are transformed back to time domain.

Figure 1 shows the block-diagram of the sound separation process. The following notation is used for different representations of signals:

- **Simple waveform (W)**
- **Simple FFT (S):** spectrogram storing $c_{k,t}$ amplitude and $j_{k,t}$ phase for each bin.
- **Frequency estimated spectrogram (F):** $c_{k,t}$ amplitudes and $j_{k,t}$ phases remain the same as in simple FFT, but an $f_{k,t}^{true}$ true frequency value is stored in addition for each bin.
- **Bandogram (B):** A spectrogram split to subbands, in which the energy is summed. Only these sums are stored, no detailed information on bin amplitudes and no phase information either.

The following chapters will guide the reader step by step through the separation process. Each section will cover one block in detail. Due to size limitations, however, some of the blocks will not be exhaustively discussed in this paper.

3 Frequency domain transformation

This section proposes an easy, yet powerful algorithm that is able to generate a spectrogram of the recording that is much more precise for musical analysis than the conventional FFT spectrogram.

Earlier literature [21] and [22] covered different transformation methods in order to determine the best possible means for analysis of audio signals. Current research [23] has examined the analysis of polyphonic musical signals in particular.

In [25] a frequency estimation method is shown, that calculates true frequencies present in the original signal from subsequent phase values. For a frame starting at time t the FFT coefficients and phases are $c_{k,t}$ and $f_{k,t}$, respectively. In this paper the time index will be omitted in some of the equations for better readability. Two subsequent frames are needed by the algorithm for the calculation. Assuming that two subsequent frames start at t_1 and t_2 , a true frequency f_{k,t_2}^{true} can be computed for each bin. The nominal frequency of the k^{th} bin is

$$f_k = k \frac{\text{samplerate}}{\text{framesize}}. \quad (1)$$

The true frequency of each bin will deviate from this value, and can be expressed as:

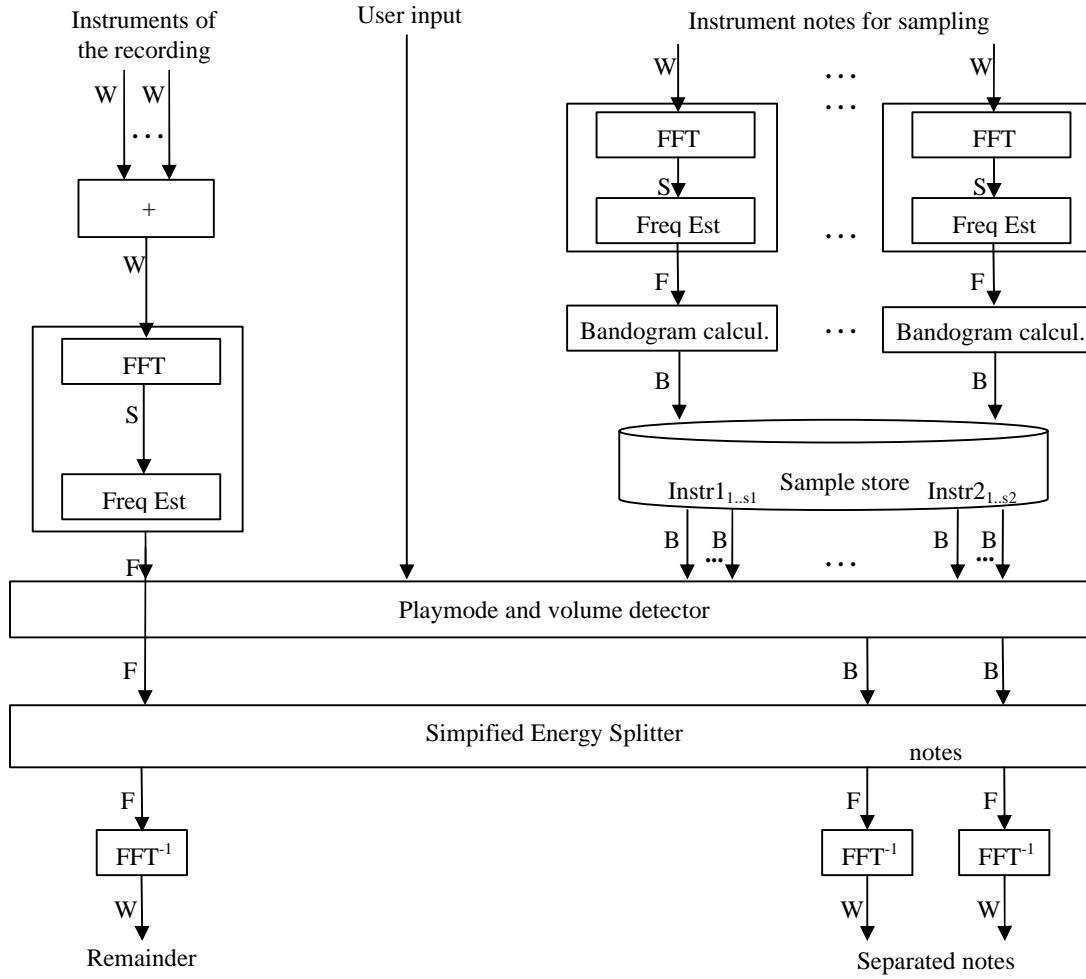


Figure 1: Signal flow and block diagram of the separation process

$$f_{k,t_2}^{true} = f_k + \frac{j_{k,t_2}^{dev}}{2p \cdot (t_2 - t_1)}, \quad (2)$$

with

$$\begin{aligned} j_{k,t_2}^{expt} &= j_{k,t_1} + (t_2 - t_1) \cdot 2p \cdot f_k \\ j_{k,t_2}^{dev} &= j_{k,t_2} - j_{k,t_2}^{expt} + l \cdot 2p \end{aligned}, \quad (3)$$

where j_{k,t_2}^{expt} is the phase of bin k in time t_2 ; j_{k,t_2}^{expt} is the expected phase; j_{k,t_2}^{dev} is the deviance between the expected and measured phase; f_{k,t_2}^{true} is the estimated true frequency of bin k in time t_2 and $l \in Z: -p < j_{k,t_2}^{dev} \leq +p$. The greater the time difference between the start of the frames the more precise the estimated value of f_{k,t_2}^{true} can be. On the other hand, large time differences limit the maximum detectable distance between f_{k,t_2}^{true} and f_k .

Sometimes, mainly for lower frequencies or complex signals with many components, f_{k,t_2}^{true} fluctuates around the real value that is present in the original signal. The true frequencies for periodic waves can be found more precisely by taking the weighted average of the last m , current and next m phase deviations. This extension will be referred to as Phase Memory (PM) and the new PM estimated true frequency can be calculated as

$$\hat{f}_{k,t_2}^{true} = f_k + \frac{\mathbf{J}_{k,t_2}^{dev}}{2p \cdot (t_2 - t_1)}, \quad (4)$$

where

$$\mathbf{J}_{k,t_x}^{dev} = \frac{\sum_{x=-\infty}^{\infty} j_{k,t_x}^{dev} \cdot c_{k,t_x} \cdot \mathbf{J}(x)}{\sum_{x=-\infty}^{\infty} c_{k,t_x} \cdot \mathbf{J}(x)}, \quad (5)$$

where $J(x)$ denotes the weighing function. Finding optimal $J(x)$ is out of the scope now.

It is important to mention that, although the proposed PM method is a very effective tool for signal analysis, it is not used in the transformation back to time domain. In a usual case, where the target is only the isolation of notes, even f_{k,t_2}^{true} can be omitted, a simple inverse FFT using c_{k,t_2} and j_{k,t_2} values will accomplish the task. However, in cases like pitch shifting, where the separated notes (and thus f_{k,t_2}^{true} frequencies) are altered, it may be necessary to recalculate the j_{k,t_2} phases from the new f_{k,t_2}^{true} values before the inverse FFT.

4 Instrument prints

The main complexity of sound separation lies in the paradox that we need to regain information from a signal that does not fully contain it. At some point we will definitely have to input additional information into the separation system to complete the missing data. Human listeners, who are known to be able to do the separation in their mind, use memories of instruments and memories of the notes in the musical piece being performed. This is their source of additional information. Copying nature has been proven to be the right approach many times. This section shows a way of implementing a memory of known instruments, trying to mimic the way the human brain works.

The representation used to store instrument features will be referred to as an instrument print. [26] presents experiments that examine the dynamic attributes of timbre evaluating the role of onsets in similarity judgments. It also gives an overview of researches tackling the challenge of the identification of the most important properties of instrument sounds that make a human listener able to distinguish between them. The instrument prints in this paper are partly based on these researches, in the sense that they contain the features that were found important in the experiments mentioned. However, separation purposes require more information on instruments than pure identification does.

An instrument print contains samples from an instrument on different frequencies and with different intonations, *playmodes*. The term *playmode* refers to the way the instrument was played, e.g. the hardness of a piano key hit, the blowing strength of the flute or the intonation of a

saxophone note. One print can have more than one playmode dimensions, depending on the way an instrument can be played. These cannot always be defined by mathematical definitions, very often they can only be expressed by subjective terms (e.g. sharpness, warmth etc.). The instrument print is a collection of samples on different frequencies f and also with different values in the playmode space $\underline{\mathbf{M}} = [m_1, m_2, \dots, m_p]$. It can be regarded as a function

$$\underline{\mathbf{A}}(\underline{\mathbf{M}}, f_k, f_{base}, t) \quad (6)$$

showing how amplitudes through the frequency range change over time for a specific note at f_{base} frequency, played with a playmode $\underline{\mathbf{M}}$, with the conditions $t, m_x, f_{base} \in \mathbb{R}^+$, $0 < m_x < m_{x,max}$, $0 \leq t < \infty$ and $0 < f \leq 20000Hz$.

In reality, a sample will not store a continuous spectrogram, only a bandogram which represents the sum energy characteristics in certain frequency subbands that are aligned on a logarithmical frequency scale. One sample is calculated from a sound signal containing one note as

$$A_{\underline{\mathbf{M}}, f_{base}, b, t} = \sum_{f_{base} \cdot 2^{-\frac{b-0.5}{R}} < f_{k,t}^{true} < f_{base} \cdot 2^{\frac{b+0.5}{R}}} c_{k,t}, \quad (7)$$

where

$$b = \left\lceil \log_{\sqrt{2}} \frac{f_{base}}{\hat{f}_{k,t}^{true}} \right\rceil \quad (8)$$

identifies the specific subband, while R is an experimental value defining the resolution in frequency range, that is, the number of subbands per octave. Experiments showed that $R=12$ provides good enough resolution in log frequency.

The number of stored instrument samples is finite both in frequency and playmode spaces. Missing samples will be interpolated from the existing ones when needed.

5 The Simplified Energy Splitter

This section describes the heart of the separation process, the Simplified Energy Splitter (SES). First, the main issue of separation will be briefly presented. Since the original decomposition problem cannot be solved due to the lack of information, further on a certain simplification will be proposed that, although lowers the quality, makes it possible to carry out the separation even under these circumstances. By applying this change

the separation problem will be simplified to an energy split problem. Finally, the reader will be guided through the implementation of the split process itself.

The equation system of the original separation problem can be expressed as

$$\underline{c}_{rt} = \sum_{\forall i} \underline{s}_{i,rt}^{orig}, \quad (9)$$

where $\underline{c}_{rt} = [c_{rt,n} \cdot e^{g_{rt,n}}]$ is the mixed signal which is the input of the separation algorithm, and $\hat{\underline{s}}_{i,rt} = [s_{i,rt,n}^{orig} \cdot e^{s_{i,rt,n}}]$ are the original notes. Time is now represented as $t = rt$, where r stands for the current frame and t is the time difference between subsequent frames. The above undetermined system of equations cannot be solved unambiguously without any further constraints.

Unfortunately our knowledge on the original notes is rather limited. No precise information is available on the base frequency $f_{base_i}(rt)$ of the original notes, their starting/ending times and their playmode which also changes over time. Each original note will ‘resemble’ one instrument sample in our database to some extent, but there are no perfect matches ever. It is obvious that under these circumstances we will not be able to decompose the recording \underline{c} to an array of signals that are perfect replicas of the original \underline{s}^{orig} ones. The target is to decompose it to signals that resemble the original ones, or – lacking the original notes – at least the samples that are used in the separation.

The term *resemblance* is of course an expression taken from real life, not an exact mathematical measure. In the case of an automated algorithm, however, it must be defined in an exact manner in order to be able to interpret and validate the outcome of the separation algorithm. Due to space considerations the definition of resemblance is not discussed in this paper.

As the original separation problem cannot be solved, simplifications have to be made. The most obvious change is eliminating the unknown $\hat{s}_{i,0t,k}$ phases from the equation system:

$$\hat{s}_{i,rt,k} = g_{rt,k} \quad (10)$$

This rephrases the original problem to

$$\underline{c}_{rt} = \sum_{\forall i} \hat{\underline{s}}_{i,rt} + \hat{\underline{c}}_{rt} \quad (11)$$

where $\hat{\underline{s}}_i$ is separated note i and $\hat{\underline{c}}$ is the remainder of the recording. This perceptually motivated modification exploits the fact that the human ear does not differentiate by the phase of the heard sinusoids, we only hear magnitude differences.

Of course any modifications will have a smaller or greater impact on the quality of the separation, causing artifacts in the output. The modified equation does not handle periodic signals with closely located frequencies well. If two or more signals cancel each other, this effect will also appear in the separated notes. Experiments showed that this tradeoff is acceptable in most cases.

In the energy split step bandograms of the right samples will be used to recreate spectrograms of the notes that are to be separated from the remaining part of the recording. Semi-linear decomposition will be used. Assuming that we know the exact frequency, volume and playmode of a certain note that we want to separate, the following iterative algorithm can be proposed to divide the energy between the target notes. We start out with the original Frequency Estimated FFT image of the recording. In each step a fraction of the energy of the selected samples is transferred from the FFT of the recording to the FFT of the separated notes. This ensures a fair division of the energy of the recording.

Let $\hat{\underline{c}}_{[0],0,rt} = [c_{[0],0,k,rt}] = \underline{c}$ denote the initial energy residing in the recording. $\hat{\underline{s}}_{[d],i,rt}$ values denote the current energy in the separated note i , being $\hat{\underline{s}}_{[0],i,rt} = \underline{0}$ initially. Each step $[d]$ contains i substeps, in which an a fraction of the energy in the reference sample A_i is transferred from the current remaining energy $\hat{\underline{c}}_{[d],i,rt}$ to the separated note $\hat{\underline{s}}_{[d+1],i,rt}$ if still possible, as in

$$c_{[d],i,k,rt} = \begin{cases} f_{base} \cdot 2^{\frac{b-0.5}{R}} < \hat{f}_{k,rt}^{true} < f_{base} \cdot 2^{\frac{b+0.5}{R}} : \\ \quad \mathbf{d} \left(\hat{c}_{[d],i-1,k,rt} \left(1 - \frac{\mathbf{a}_{[d],b,i,rt}}{D} \right) \right) \\ \text{otherwise :} \\ \quad \hat{c}_{[d],i,k,rt} \end{cases} \quad (12)$$

with

$$\mathbf{a}_{[d],b,i,rt} = \frac{A_i (M_i, f_{base_i}, b, rt - T_{i,start})}{\sum_{f_{base} \cdot 2^{\frac{b-0.5}{R}} < \hat{f}_{k,rt}^{true} < f_{base} \cdot 2^{\frac{b+0.5}{R}}} \hat{c}_{[0],0,k,rt}} \quad (13)$$

The current energy in note i (which is being isolated) can be calculated in step $[d]$ as

$$\hat{\underline{s}}_{[d+1],i,rt} = \hat{\underline{s}}_{d,i,rt} + (\hat{\underline{c}}_{d,i-1,rt} - \hat{\underline{c}}_{d,i,rt}) \quad (14)$$

which is the starting value of the next $[d+1]$ step for separated notes, while

$$\hat{\underline{c}}_{[d+1],0,rt} = \hat{\underline{c}}_{d,I,rt} \quad (15)$$

is the starting value of the next step for the remaining energy in the recording, where I is the number of instruments in the time frame. With D denoting the number of steps, $\underline{c}_{[D],0,rt}$ is the remaining part of the recording after the separation and $\hat{\underline{s}}_{[D],i,rt}$ will represent the coefficients of instrument i after the separation.

6 Playmode detection

In the previous section all note parameters were assumed to be known. However, this is not the usual scenario. While the user can input the location of the instrument notes in frequency and time, they may not be capable of entering either the playmode matrix \underline{M} or the volume. Further on the automatic detection of the playmode and volume will be covered.

To carry out the energy split step an optimal playmode matrix \underline{M} must be found. For the sake of convenience the volume will also be incorporated in \underline{M} from now on. \underline{M} is by definition perfect if the separated notes are the perfect replicas of the parent instrument samples that were used in the energy split, and the remaining part is zero. In general, an \underline{M} matrix is considered good provided the energy split step that uses \underline{M} generates notes that 'resemble' their parent sample while letting $\underline{c}_{[D]}$ get as close to zero as possible. All combinations of \underline{M} matrices will be tested for separation error, and the one causing the least error will be selected. Finding the most suitable error definition is not covered here. This article will later use the Signal to Distortion Ratio (SDR), which will be introduced in Section 7.

Depending on the size and possible values of \underline{M} the number of steps needed for finding the best combination of the instrument samples may require huge computational power. If we consider the playmode space to be continuous, it is not even possible to iterate through all the combinations. Finding an algorithm faster than brute force iteration, however, is out of the scope of this article.

7 Synthetic tests

The quality of the separation system described so far was evaluated by using a number of synthetic tests. This section features different test scenarios that illustrate the quality of the given method.

The test setup used for evaluating the performance of the separation system is shown in Figure 2. The test system was based upon the instrument sample collection of the University of Iowa, [27]. The waveforms were normalized and converted to mono, with a sampling rate of 44000Hz, 16 bits. Their DC offset was also corrected (shifted to zero) where it was necessary. The waveforms were then divided into samples containing only one instrument note using thresholding. Samples shorter than 500 ms were dropped, while samples longer than 2 seconds were cropped to 2 seconds. The above mentioned process resulted in 3841 waveforms of separate instrument notes of harmonic instruments. The instrument database contained samples of the following instruments: flute, saxophone, bass, clarinet, bassoon, trombone, cello, horn, oboe, piano, pizzicato, trombone, trumpet, tuba, viola and violin.

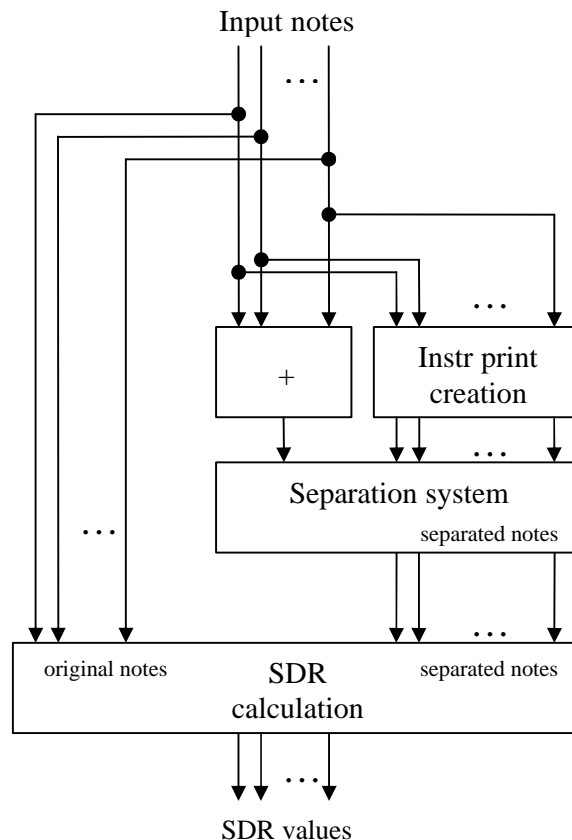


Figure 2: Block diagram of the automatic synthetic test system

In each of our tests a random set of instrument note waveforms were selected. The waveforms were converted to instrument prints using the technique described in section 4. The selected samples were then mixed together and fed to the separation system as the input recording.

Three testing series were performed, each one focusing on a specific aspect of note separation. In each series the performance of the system was tested for polyphony levels of 2 to 10. At each level a total of 50 individual tests were carried out.

In the first series the input waveforms were simply mixed together as in

$$\tilde{c}(n) = \sum_{i=1}^I \tilde{s}_i^{orig}(n), \quad (16)$$

where \tilde{c} and \tilde{s}_i^{orig} are the waveform of the mixed signal, and the note signals, respectively. This method ensures that all the notes share the same onset time. As the onset/ending times and note base frequencies were all known, user input was not necessary. This part was algorithmically fed to the separation system. Then the separation step was performed with the input mixed signal, the score data and the pre-sampled instruments.

The output note waveforms of the system were then compared to the original waveforms. Since the algorithm preserves the phase information of the original (mixed) input signal, an error between the original and the output waveforms can be obtained simply by subtracting the output signal from the original ones in time domain.

The mean-square level of each error signal was computed over the whole signal. A signal-to-distortion ratio could then be obtained by comparing these levels to the original:

$$SDR_i \text{ [dB]} = 10 \log_{10} \frac{\sum_n \tilde{s}_i^{orig}(n)^2}{\sum_n [\tilde{s}_i(n) - \tilde{s}_i^{orig}(n)]^2}, \quad (17)$$

where \tilde{s}_i is the waveform of the separated signal for note i .

Within one polyphony level the average and standard deviation of the individual SDR_i values were calculated, which is shown in Figure 3.

When listening to the separated output channels manually, it was found that the SDR gave false negative values in some cases. Although the quality of certain separated notes was satisfactory, the calculated SDR value still indicated very low separation quality for these notes. After investigating the issue we found that in most of

these cases the bad SDR scores can be explained by phase deviation between the original signals and the mixture that was generated from them. When two or more of the original signals contain energy on common frequencies then the phase of the resulting mixture may deviate from the phase of the original signals. This difference is not audible to human listeners, however it highly affects the SDR value.

To overcome the above mentioned issue we must realize that under normal circumstances a human listener is unable to sense any difference between two signals that differ by phase, but otherwise have the same energy content (except for very rare cases or some synthetic tests). Therefore we introduce here another, perceptual measure. Instead of comparing the time-domain waveforms to the time-domain residual, we recalculate the SDR error levels using the frequency-domain representation of the signals. The error can now be expressed as:

$$SDR_i^F \text{ [dB]} = 10 \log_{10} \frac{\sum_{\forall rt} \sum_{k=0}^K s_{i,k}^{orig}(rt)^2}{\sum_{\forall rt} \sum_{k=0}^K [s_{i,k}(rt) - s_{i,k}^{orig}(rt)]^2}, \quad (18)$$

Although the new measure represents the perceptual separation quality more efficiently, in this paper we still include the original SDR measurements for easier comparison with other works. We must note, that there is no definite relation between the two.

The second test series was very similar to the first one. In this case, however, the note onset times were dilated by a Δt_i random number, as in

$$\tilde{c}(rt) = \sum_{i=1}^I \tilde{s}_i^{orig}(rt - \Delta t_i) \quad (19)$$

where

$$0 \text{ sec} < \Delta t_i < 1 \text{ sec}. \quad (20)$$

The first test provides greater overlap of the notes in time, while in the second test it is easier to observe artifacts in one separated note caused by the onset of another note. However, the two tests returned similar SDR values. Measurement results are shown on Figure 4.

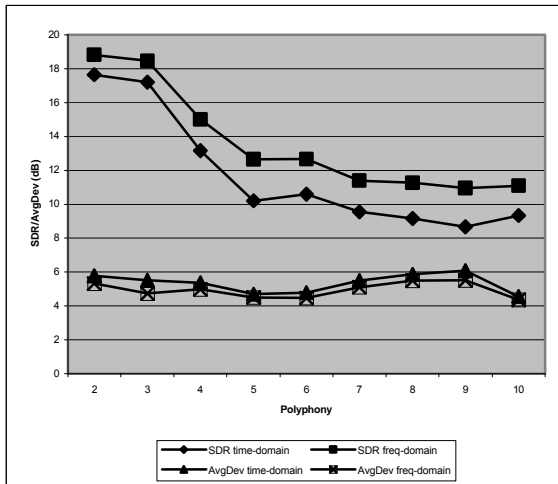


Figure 3: Signal to Distortion Ratio for test 1

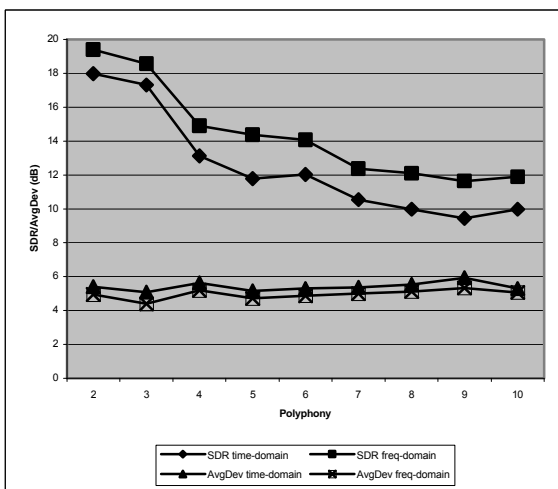


Figure 4: Signal to Distortion Ratio for test 2

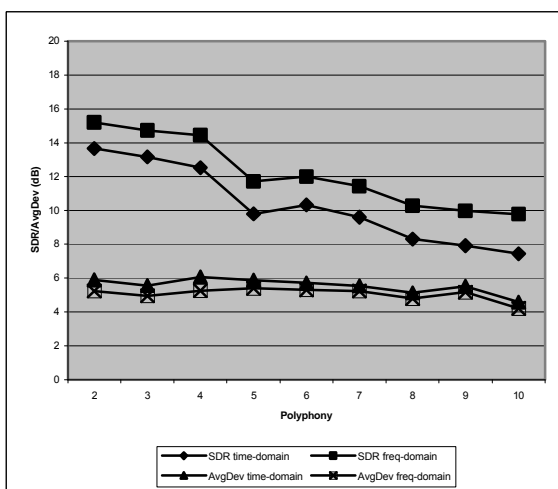


Figure 5: Signal to Distortion Ratio for test 3

The third test covered cases with notes that are in overtone relation with each other. For this reason Test 1 was repeated with input channels meeting the following condition:

$$\exists i \neq j \ni \frac{f_{base,i}}{f_{base,j}} \approx \frac{m}{n} \quad m, n \in \{1, 2, 3, 4, 5\} \quad (21)$$

where i and j represent the input instrument notes. Condition (21) ensures that any two notes that are mixed together for testing purposes should be in close overtone relation with at least one input note in the set.

Figure 5 presents measurement results for Test 3. As can be expected, the SDR values in this test were lower than in the previous series. This can be attributed to the fact that the vibrations generated by the notes cancelled each other in most of the test cases. Such cases inevitably introduce the beating effect in the generated mixture. The approach we took, by applying simplification (10) to the system, does not handle the beating, causing some artifacts to appear in the output channels. However, despite the fact that we cannot get back the original tracks, it was found that the separated channels still resemble real-life instruments even in higher polyphony levels.

8 Implementation and real-life test scenario

This section shows an example scenario of a difficult problem in the area of sound separation. The results illustrate the quality of the method proposed for isolating notes in practice.

In the test case two instruments – a piano and a saxophone – played the tune in Figure 6. As shown, the two instruments played on the same frequency, which is one of the hardest cases of sound separation. The instruments were recorded separately, one after the other. Their signals were then mixed together and fed to the separation algorithm.



Figure 6: musical score of the analyzed fragment

Instrument prints were also built from these two instruments by sampling all halftones in the frequency range in question at three different playmodes ('soft', 'neutral' and 'hard').

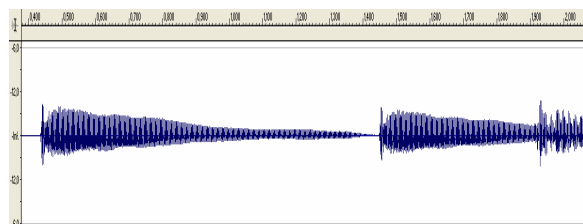


Figure 7: waveform of the original piano track

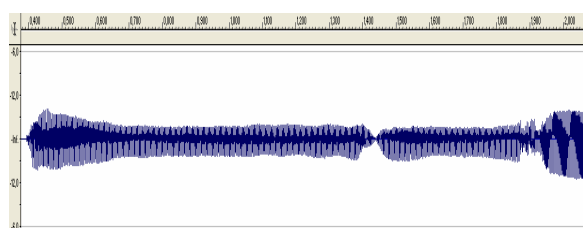


Figure 8: waveform of the original sax. track

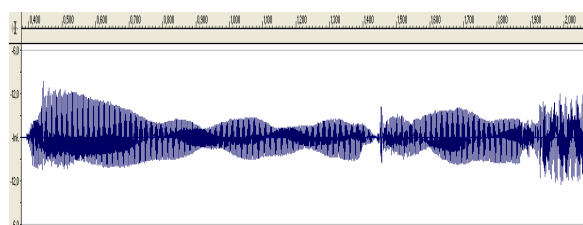


Figure 9: waveform of the mixed signal

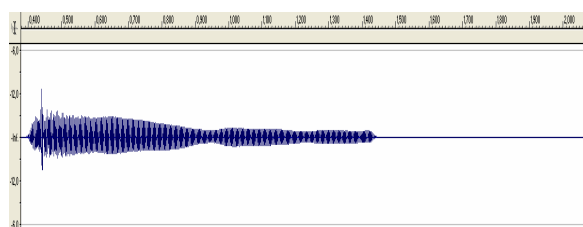


Figure 10: waveform of the first separated piano note

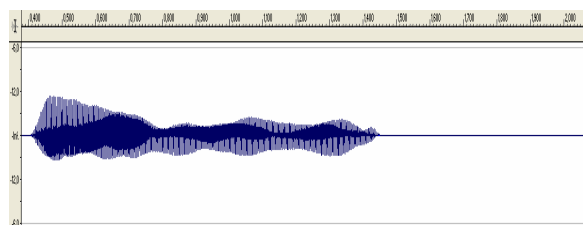


Figure 11: waveform of the first separated sax. note

Figure 7 to Figure 11 show the two input instruments, the mixed signal, and two separated notes, respectively. Human listeners confirmed that the separated notes did sound like real instruments, even if somewhat distorted.

9 Conclusion

This paper has shown a method for separating single instrument notes from a recording using pre-recorded instrument prints and the Simplified Energy Splitter algorithm. The novelty of the approach lies in the way the instrument prints are modeled, the algorithm for the SES and the system architecture. The results are quite promising. Simulation experiments on generated mixtures of pitched real-life musical instruments show that the proposed method allow average SDR above 18 dB for two simultaneous sources, and the quality decreases gradually as the level of polyphony increases.

Example waveforms from Section 7 and 8 along with a number of other separation samples can be downloaded from

<http://avalon.aut.bme.hu/~aczelkri/separation>.

For recordings that only contain harmonically unrelated notes the algorithm provides very clear results. In real life, however, consonant notes with overlapping overtones are usually favored over dissonant ones. Our test results show that even in cases where many notes are located on each other's base or overtone frequencies the separation provides reasonably good results.

Acknowledgements:

This work has been supported by the fund of the Hungarian Research Fund (grant number T68370)

References:

- [1] G. E. Poliner, D. P. W. Ellis, A. F. Ehmann, E. Gomez, S. Streich; B. Ong, Melody Transcription From Music Audio: Approaches and Evaluation, *IEEE Transactions on Audio, Speech and Language Processing*, Vol. 15., No. 4, 2007, pp. 1247 – 1256
- [2] H. Thornburg, R. J. Leistikow, J. Berger, Melody Extraction and Musical Onset Detection via Probabilistic Models of Framewise STFT Peak Data, *IEEE Transactions on Audio, Speech and Language Processing*, Vol. 15, No. 4, 2007, pp. 1257 – 1272

- [3] J. P. Bello, L. Daudet; M. B. Sandler, Automatic Piano Transcription Using Frequency and Time-Domain Information, *IEEE Transactions on Audio, Speech and Language Processing*, Vol. 14, No. 6, 2006, pp. 2242 – 2251
- [4] A. T. Cemgil, H. J. Kappen; D. Barber, A generative model for music transcription, *IEEE Transactions on Audio, Speech and Language Processing*, Vol. 14, No. 2, 2006, pp. 679 – 694
- [5] S. A. Abdallah and M. D. Plumbley. Polyphonic transcription by nonnegative sparse coding of power spectra, *Proc of International Conference on Music Information Retrieval*, Barcelona, Spain, 2004, pp 318 – 325
- [6] S. A. Abdallah, M. D. Plumbley, Unsupervised analysis of polyphonic music by sparse coding, *IEEE Transactions on Neural Networks*, Vol. 17., No. 1., 2006, pp. 179 - 196
- [7] A. A. Wiczorkowska, A. Czyzewski, Rough Set Based Automatic Classification of Musical Instrument Sounds, *Electronic Notes in Theoretical Computer Science*, Vol. 82., No. 4., 2003, pp. 298 – 309
- [8] J. C. Brown, Computer identification of musical instruments using pattern recognition with cepstral coefficients as features, *J. Acoust. Soc. of Am.*, Vol. 105, 1999, pp. 1933 – 1941.
- [9] P. Herrera, X. Amatriain, E. Batlle, X. Serra, Towards instrument segmentation for music content description: a critical review of instrument classification techniques, *Proc. of International Symposium on Music Information Retrieval (ISMIR 2000)*, Plymouth, MA, 2000.
- [10] P. Herrera-Boyer, G. Peeters, S. Dubnov, 2003. Automatic classification of musical instrument sounds, *J. New Music Res.* Vol. 32, pp. 3 – 21.
- [11] D. Barry, R. Lawlor, E. Coyle, Sound Source Separation: Azimuth Discrimination and Resynthesis, *Proc. of 7th International Conference on Digital Audio Effects*, DAFX'04, Naples, Italy, 2004.
- [12] A. Klapuri, Multipitch estimation and sound source separation by the spectral smoothness principle, *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, Salt Lake City, USA, 2001.
- [13] T. Virtanen, A. Klapuri, Separation of Harmonic Sounds Using Multipitch Analysis and Iterative Parameter Estimation, *Proc. of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, New York, USA, 2001.
- [14] T. Virtanen, A. Klapuri, Separation of Harmonic Sound Sources Using Sinusoidal Modeling, *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, Istanbul, Turkey, 2000.
- [15] T. Virtanen, A. Klapuri, Separation of harmonic sounds using linear models for the overtone series, *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing*, Orlando, Fla, USA, 2002.
- [16] P. COMON, Independent Component Analysis, a new concept, *Signal Processing*, Elsevier, Vol. 36., No. 3., 1994, pp. 287 – 314
- [17] N. Mitianoudis, M. E. Davies, Using Beamforming in the audio source separation problem, *7th Int Symp on Signal Processing and its Applications*, Paris, 2003
- [18] Lee, D.D. and H.S. Seung, Learning the parts of objects by non-negative matrix factorization, *Nature* 401, 1999, pp 788 – 791
- [19] P. Smaragdis and J. C. Brown, Non-Negative Matrix Factorization for polyphonic music transcription, *Proc. of 2003 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, New York, USA, 2003
- [20] T. Virtanen, Sound Source Separation in Monoaural Music Signals, *PhD thesis, University of Kuopio*, 2006
- [21] R. Pintelon, J. Schoukens, *System Identification, A frequency domain approach*, ISBN 0-7803-6000-1, Wiley-IEEE Press, 2001, pp. 33 – 44
- [22] S. Gade, H. Herlufsen, Use Of Weighting Functions in DFT/FFT Analysis (Part I), *Brüel & Kjør Technical Review*, No. 3., 1987
- [23] K. Aczél, I. Vajk, Note separation of polyphonic music by energy split, *Proc. of WSEAS International Conference on Signal Processing, Robotics and Automation*, Cambridge, England, 2008
- [24] K. Aczél, Sz. Iváncsy, Sound separation of polyphonic music using instrument prints, *Proc of EUSIPCO 2007*, Poznan, Poland, 2007.
- [25] S. M. Bernsee, Pitch Shifting Using the Fourier Transform
<http://www.bernsee.com/dspdimension.com/html/pshiftstft.html> (04-04-2008)
- [26] P. Iverson; C. L. Krumhansl, Isolating the dynamic attributes of musical timbre, *The J. Acoust. Soc. of America*, Vol. 94, No. 5, 1993, pp.2595 – 2603
- [27] The University of Iowa Musical Instrument Samples Database.
<http://theremin.music.uiowa.edu>, 04-04-2008.