

# Chapter Number

## Monaural Audio Separation Using Spectral Template and Isolated Note Information

Anil Lal and Wenwu Wang  
*Department of Electronic Engineering, University of Surrey,  
United Kingdom*

### 1. Introduction

Musical sound separation systems attempt to separate individual musical sources from sound mixtures. The human auditory system gives us the extraordinary capability of identifying instruments being played (pitched and non-pitched) from a piece of music and also hearing the rhythm/melody of the individual instrument being played. This task appears 'automatic' to us but has proved to be very difficult to replicate in computational systems. Many methods have been developed recently for addressing this challenging source separation problem. They can be broadly classified into two categories, respectively, statistical learning based techniques such as independent component analysis (ICA) and non-negative matrix/tensor factorization (NMF/NTF), and computational auditory scene analysis (CASA) based techniques.

One of the popular methods for source separation was based on ICA [1-10], where the underlying unknown sources are assumed to be statistically independent, so that a criterion for measuring the statistical distance between the distribution of the sources can be formed and optimised either adaptively [5] [10] [11], or collectively (in block or batch processing mode) [2], given mixtures as the input signals. Both high-order statistics (HOS) [2] [4] [5] and second order statistics (SOS) [12] have been used for this purpose. The ICA techniques have been developed extensively since the pioneering contributions in early 1990s, made for example by Jutten [1], Comon [3], and Cardoso [2]. The early work of ICA concentrates on the instantaneous model which was soon found to be limited for real audio applications such as in a cocktail party environment, where the sound sources reach listeners (microphones) through multi-path propagations (with surface reflections). Convolutional ICA [13] was then proposed to deal with such situations (see [21] for a comprehensive survey). Using Fourier transform, the convolutional ICA problem can be converted to multiple instantaneous but complex valued ICA problems in the frequency domain [14-17] thanks to its computational efficiency, and the sources can be separated after permutation correction for all the frequency bins [18-20]. Most of the aforementioned methods consider (over-)determined cases where the number of sources is assumed to be no greater than the number of observed signals. In practical situations, however, an underdetermined separation problem is usually encountered. A widely used method for tackling this problem is based on sparse signal representations [22-29], where the sources are assumed to be sparse in either the time domain or a transform domain such that the overlap between the sources at

1 each time instant (or time-frequency point) is minimal. Audio signals (such as music and  
2 speech) become sparser when transformed into the time-frequency domain, therefore, using  
3 such a representation, each source within the mixture can be identified based on the  
4 probability of each time-frequency point of the mixture that is dominated by a particular  
5 source, using either sparse coding [26], or time-frequency masking [30] [31] [33], based on  
6 the evaluation of various cues from the mixtures, including e.g. statistical cues [20], and  
7 binaural cues [30] [32]. Other methods for source separation include, for instance, the non-  
8 negative ICA method [34], the independent vector analysis (IVA) [35], and NMF/NTF [37-  
9 43]. Comprehensive review of ICA (and other statistical learning) methods is out of the  
10 scope of this chapter, and for more references, we refer the interested readers to the recent  
11 handbook on ICA edited by Comon and Jutten [36], and a book on NMF by Cichocki [44].

12 Many ICA methods discussed above can be broadly applied to different types of signals.  
13 In contrast, CASA is another important technique dealing specifically with audio signals,  
14 which is based on the principles of Auditory Scene Analysis (ASA). In [60], Bregman  
15 attempts to explain ASA principles by illustrating the ability of the human auditory  
16 system to identify and perceptually isolate several sources from acoustic mixtures by  
17 separating the sources into individual (perceptual) acoustic streams for each source,  
18 which suggests that the auditory system operates in two main stages, segmentation and  
19 grouping. The segmentation stage separates the mixture into (time-frequency)  
20 components that would relate to an individual source. The grouping stage then groups  
21 the components that are likely to be from the same source e.g. using information such as  
22 simultaneous onset/offset of particular frequency amplitudes or relationships of  
23 particular frequencies to source pitch [45-50]. It is well-known that the ICA technique is  
24 not effective in separating the underdetermined mixtures, for which, as mentioned above,  
25 one has to turn to, e.g. the technique of sparse representations, by sparsifying the  
26 underdetermined mixtures into a transform domain, and to reconstruct the sources using  
27 sparse recovery algorithms [51-53]. In contrast, CASA technique evaluates the temporal  
28 and frequency information of the sources directly from the mixtures, and therefore it has  
29 the advantage in dealing with underdetermined source separation problem, without  
30 having to assume explicitly the system to be (over-) determined, or the sources to be  
31 sparse. This is especially useful for addressing the monaural (single-channel) audio source  
32 separation problem, which is an extreme case of the underdetermined source separation  
33 problem. The task of computationally isolating acoustic sources from a mixture is  
34 extremely challenging, and recent efforts attempt to isolate speech/singing sources from  
35 monaural musical pieces or to isolate an individual's speech from a speech mixture [45]  
36 [54-59] [61] [62], and have achieved reasonable success. However, the task of separating  
37 musical sources from a monaural mixture has been, thus far, less successful in  
38 comparison.

39 The ability to isolate/extract individual musical components within an acoustic mixture  
40 would give an enormous amount of control over the sound. Musical pieces could be un-  
41 mixed and remixed for better musical fidelity. Signal processing, e.g. equalisation or  
42 compression, could be applied to individual instruments. Instruments could be removed  
43 from a mixture, possibly for musical students to accompany pieces of music for practice.  
44 Control over source location could be achieved in 3-D audio applications by placing the  
45 source in different locations within a 3D auditory scene.

1 Musical sources (instruments) have features in the frequency spectrum that are highly  
2 predictable due to the fact that they are typically constrained to specific notes (A to G# on  
3 the 12-tone musical scale) and so, frequencies are typically constrained to particular values.  
4 As such, harmonic frequencies are predictable as they can be derived from multiples of the  
5 fundamental frequency. If reliable pitch information for each source is available, harmonic  
6 frequencies for each source can be determined. With this information in hand, frequencies  
7 where harmonics from each source would overlap can be calculated. Non-overlapped  
8 harmonic frequencies in each source can therefore also be determined and non-overlapped  
9 and overlapped harmonic frequency regions in the mixture can be found, along with which  
10 particular source each non-overlapped harmonic would belong to. Existing systems [63-65]  
11 are successful in using this pitch information to identify non-overlapped harmonics and the  
12 source to which it belongs.

13 Polyphonic musical pieces typically have notes that complement each other (i.e. perfect 3rd,  
14 perfect 5th, minor 7th etc., explained by music theory) and so, result in a high, and regular,  
15 number of harmonics that overlap. For this reason, musical acoustic mixtures contain a  
16 larger number of overlapping harmonics in comparison to speech mixtures. Existing sound  
17 separation systems do not completely address the problem of resolving overlapping  
18 harmonics i.e. determining the contribution of each source to an overlapped harmonic. And  
19 so, because of typically higher numbers of overlapping harmonics in musical passages,  
20 musical sound separation is a difficult task and performance of existing source separation  
21 techniques has been limited. Therefore, the major challenge in musical sound separation is  
22 to effectively deal with overlapping harmonics.

23 A system proposed by Every and Szymanski [64] attempts to resolve overlapping harmonics  
24 by using adjacent non-overlapped harmonics to interpolate an estimate of the overlapped  
25 harmonic and so, 'fills out' the 'missing' harmonics for the spectrum of non-overlapped  
26 harmonics of each source. Nevertheless, this method relies heavily on the assumption that  
27 spectral envelopes are smooth and that amplitudes of any harmonic will have a 'middle  
28 value' of the amplitudes of the adjacent harmonics. In practice, however, spectral envelopes  
29 of real instruments rarely are smooth so this method produces varied results.

30 Hu [66] proposes a method of sound separation that uses onset/offset information (i.e.  
31 where performed notes start and end). Transient information in the amplitude envelope is  
32 used to determine onset/offset time by half-wave rectifying and low pass filtering the  
33 signals to obtain the amplitude envelope and the first order differential of the envelope  
34 highlights the time of sudden change in the envelope. This is a powerful cue as regions of  
35 isolated note performances can be determined. Li and Wang [63] also incorporate  
36 onset/offset information to separate sounds. However, the Li-Wang system uses the  
37 predetermined pitch information to find the onset/offset time; the time points where pitches  
38 change by at least a semi-tone are labelled appropriately as onset or offset times.

39 The Li-Woodruff-Wang system [67] incorporates a method utilizing common amplitude  
40 modulation (CAM) information to resolve overlapping harmonics. CAM suggests that all  
41 harmonics from a particular source have similar amplitude envelopes. The system uses the  
42 change in amplitude from the current time frame to the next of the strongest non-  
43 overlapped harmonic (in terms of a ratio), and the observed change in phase of the  
44 overlapped harmonic from the mixture to resolve the overlapped harmonic by means of  
45 least-squares estimation.

1 The focus of this chapter is to investigate the musical sound separation performance using  
2 pitch information and CAM principles described by Li-Woodruff-Wang [67] and proposing  
3 methods for the improvements of the system performance. The methods outlined by the  
4 pitch and CAM separation system have shown promising results, but only a small amount  
5 of research has been carried out that uses both pitch and CAM techniques together [67].  
6 Preliminary work reveals that the pitch and CAM based system produces good results for  
7 mixtures containing long notes with considerable sustained portions e.g. a violin holding a  
8 note, but produces poor quality results for attack sections of notes, i.e. mixtures containing  
9 instruments with smaller, or no sustain sections (just attack and decay sections), e.g. a piano.  
10 Modern music typically has a high number of non-sustained note performances so the pitch  
11 and CAM method would fail with a vast number of musical pieces. In addition, the pitch  
12 and CAM method has difficulty in dealing with the overlapping harmonics, in particular,  
13 for audio sources playing similar notes.

14 This study aims to investigate more reliable methods of resolving harmonics for the pitch  
15 and CAM based technique of music separation which improves results, particularly for  
16 attack sections of note performances and overlapping harmonics. A method of using  
17 isolated (or relatively isolated) sections of performances in mixtures by obtaining  
18 onset/offset information is used to provide more reliable information to resolve harmonics.  
19 Such information is also used to generate a spectral template which is further used to  
20 improve the separation performance of overlapping spectral regions in the mixtures, based  
21 on the reliable information from non-overlapping regions. Implementation of the proposed  
22 methods is then attempted using a baseline pitch and CAM source separation algorithm,  
23 and system performance is evaluated.

## 24 **2. Pitch and CAM system and its performance analysis**

25 In general, the pitch and CAM system shows good performance for separating audio  
26 sources from single channel mixtures. However, according to our experimental evaluations  
27 briefly discussed below, its separation performance is limited for attack sections of notes  
28 and regions of same note performances.

29 We first evaluate the performance of the pitch and CAM system for separating the attack  
30 sections of music notes. To this end, we take the baseline pitch and CAM algorithm  
31 implemented in Matlab to test its performance. We use a sample database of real  
32 instrument recordings (available within ProTools music production software) to generate  
33 test files, so that the system performance on separating attack sections of notes could be  
34 evaluated. The audio file generated is a (monaural) single-channel mixture containing a  
35 melody played on a cello, and a different but complimentary melody played on a piano.  
36 The purpose of combining complimentary melodies from different sources is to generate a  
37 realistic amount of overlapping harmonics between sources, as would be found in typical  
38 musical pieces. Qualitative results show that the cello, which had long sustained portions  
39 of notes, is separated considerably well, while the attack sections of piano notes are in  
40 some cases lost as a result of the limited analysis time frame resolution. The piano has  
41 shorter notes with no sustain sections, only attacks and decays, but still contains  
42 considerable amount of harmonic content. As a result, the system performs less effectively  
43 in separating the piano source which highlights the difficulty the separation system has in  
44 isolating instruments playing short notes that are made up of regions of attack. Another

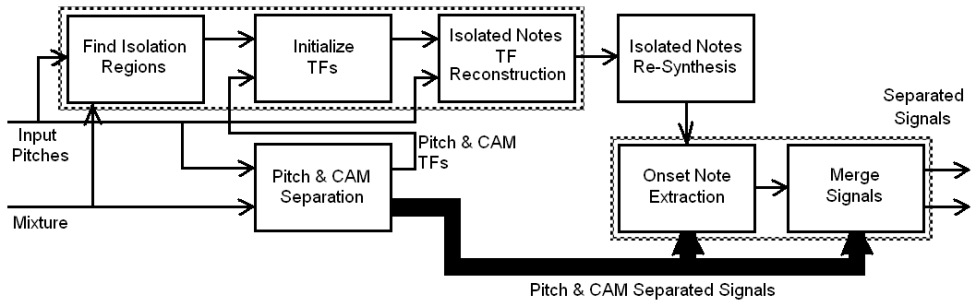
1 experiment on the mixture of audio sources played by clarinet and cello again confirms  
2 that the pitch and CAM system has difficulty in separating the soft attack sections of the  
3 notes played by the clarinet.

4 We then evaluate the system performance for the regions of same notes in the mixture.  
5 We generated a mixture containing a piano and a cello performing the same note (C4).  
6 Using the pitch and CAM system, the cello was separated from the mixture but with some  
7 artefacts and distortions. However, the system was unsuccessful in separating the piano  
8 source, and only a low level signal could be heard that did not resemble the original piano  
9 signal. In another experiment, we generated a mixture with a cello playing the note C4  
10 and a piano playing all notes in sequence from C4 to C5 (C4, C#4, D4, D#4... etc.) and  
11 ended on the note C4. The cello was separated well from the mixture as were all notes  
12 played by the piano except the notes C4 and C5 at the both ends of the sequence. Due to  
13 the slow attack of the cello, the C4 note played by the piano at the beginning of the piece  
14 was better separated than the C4 note at the end of the sequence, as the C4 note at the  
15 beginning is more isolated. In addition, we have examined the performance of the system  
16 for mixtures with the same note and varying octaves. To this end, we generated another  
17 mixture with a cello playing the note C3 and a piano playing notes C1 to C6 in sequence  
18 and then ending on note C2. The results again show that the cello was separated well but  
19 with high distortions in sections where the piano attacks occur. The piano notes C1 and  
20 C2 were separated with some distortions but notes C3 through to C6 were almost not  
21 separated at all.

22 In summary, the pitch and CAM system does not perform well for recovering the sharp  
23 transients of the amplitude envelope from mixtures due to the limited time frame  
24 resolution, and it also has difficulty in separating notes with same fundamental frequencies  
25 and harmonics, caused by insufficient data for resolving the overlapping harmonics and for  
26 extracting the CAM information. For example, if one source has a pitch frequency of 50 Hz,  
27 its harmonics would occur at 100 Hz, 150 Hz, etc. If the pitch frequency of a second source is  
28 an octave higher, i.e. 100 Hz, its harmonics would occur at 200 Hz, 300 Hz, etc. As a result,  
29 the harmonics of the second source will be overlapped with those of the first source. To  
30 address these problems, we suggest two methods to improve the pitch and CAM system,  
31 respectively, isolated note and spectral template methods, which attempt to better resolve  
32 the overlapping harmonics when the information used by the pitch and CAM system is  
33 considered to be unreliable, as described next in detail.

### 34 **3. Isolated note method**

35 The proposed isolated note system, shown in Figure 1, uses note onset/offset information to  
36 determine periods of isolated performance of an instrument so that the reliable spectral  
37 information from the isolated regions can be used to resolve overlapping harmonics in the  
38 remaining note performance regions. The proposed system is based on the pitch and CAM  
39 algorithm [67], with the addition of new processing stages shown in dotted lines in Figure 1.  
40 Same to the pitch and CAM system, the inputs to the proposed system are mixture signals  
41 and pitch information supplied by a pitch tracker. The details of each block in Figure 1 are  
42 explained below.



1  
2 Fig. 1. Diagram of Isolated Note System.

3 The first processing stage is the *Pitch and CAM Separation* stage. The mixture signal is  
4 separated using the method described in [67] and by using the pitch information provided.  
5 The separated signals are used later in *Onset Note Extraction* and *Merge Signals* stages by the  
6 isolated note system. When the pitch and CAM separation is carried out the time-frequency  
7 (TF) representations of the mixture signal and the separated signals are generated which are  
8 then utilized later by the *Initialize TFs* processing stage.

9 The next processing stage is the *Find Isolated Regions* stage. Using input pitch information,  
10 we attempt to find time frames for each source where isolated performances of notes occur.  
11 Each time frame of each source is evaluated to determine if other sources contain pitch  
12 information (i.e. if other notes are performing during the same time frame). A list of time  
13 frames for each source is created and a flag is raised (the time frame is set to 1) if the note for  
14 the current frame and current source is isolated. Each occurrence of an isolated region  
15 (indicated by the flag) in each source is then numbered so that each region can be identified  
16 and processed independently at a later stage (achieved by simply searching through time  
17 frames and incrementing the region number at each encounter of a transition from 0 to 1 in  
18 the list of flagged time frames).

19 Next, we determine the non-isolated regions for the notes that contain a region of isolated  
20 note performance. For each numbered isolated region we find the corresponding non-  
21 isolated note performance and generate a new list where time frames for the non-isolated  
22 regions are numbered with the number relating to the corresponding isolated region.  
23 Note that we do not number the isolated time frames themselves in the newly generated  
24 list.

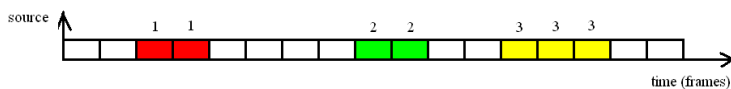
25 The new list is generated by searching back (previous frames) from the relevant isolated  
26 region and numbering all frames appropriately, and we then repeat by searching forward  
27 from the isolated region. Searches are terminated at endpoints of the note or at occurrences  
28 of another isolated region. Each isolated region that generates a set of corresponding non-  
29 isolated frames is saved in a new list separately, the list is then collapsed to form a final list  
30 where time frames for which we have non-isolated regions relating to two isolated regions  
31 are split halfway.

32 This is better illustrated by Fig. 2. Fig. 2(a) shows an occurrence of a note with three  
33 isolated regions for which information of time frames with isolated performance is  
34 determined. Fig. 2(b) illustrates that the non-isolated regions relating to each isolated

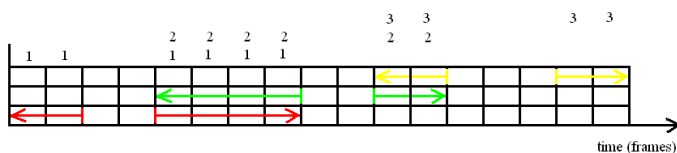
1 region, are found by searching forwards and backwards and terminating at endpoints of  
 2 notes or an occurrence of another isolated region. Each region is stored individually. Fig.  
 3 2(c) shows the final set of regions where time frames 'belonging' to two sets are split  
 4 halfway.

5 The TF representation of each source is formed for the isolated notes in the *Initialize TFs*  
 6 stage. We initialize the TF representation by starting with an empty set of frequency  
 7 information for each time frame and then by searching through the list of isolated regions.  
 8 For time frames that are identified as an isolated performance of a note (from the list), we  
 9 copy all frequency information for those frames directly from the mixture to the  
 10 corresponding TF representation of the sources. This is shown in Fig. 3 where the time  
 11 frames for the isolated performances of the note C4 (in Fig. 3(a)) are copied directly to  
 12 initialize the TF representation. Fig. 3(b) shows that all of the harmonic information is  
 13 copied directly from the mixture; hence all harmonics are correctly present in the initialized  
 14 isolated note TF representation.

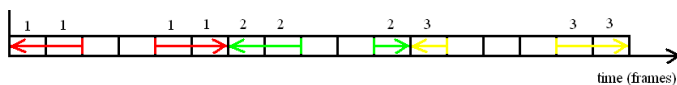
15  
16  
17  
18  
19  
20  
21



(a) Time Frames with Numbered Isolated Regions.



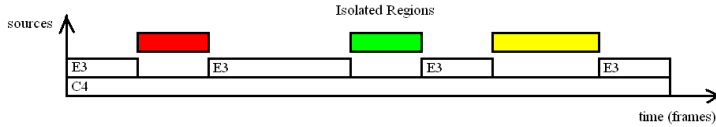
(b) Non-Isolated Regions Corresponding to Each Isolated Region.



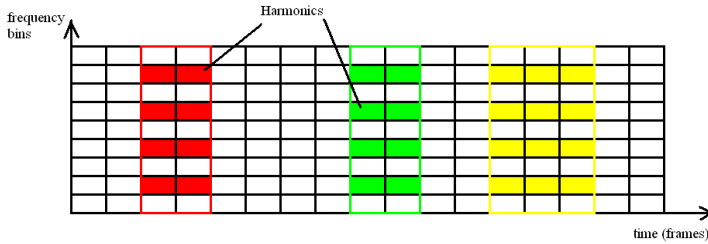
(c) Time Frames of Non-Isolated Regions Associated with Each Isolated Region.

28  
29  
30  
31  
32  
33  
34

Fig. 2. Method Used to Determine Non-Isolated Regions of Isolated Notes.



(a) Note Performance and Isolated Note Regions.

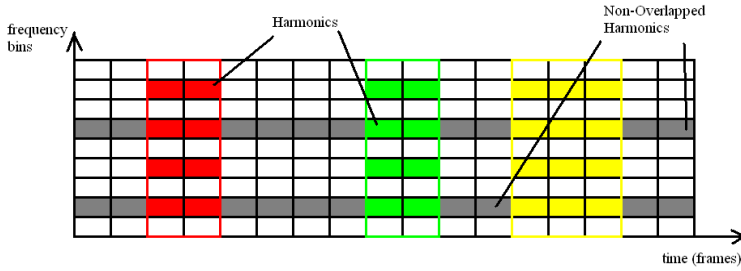


(b) Initialized TF Representations

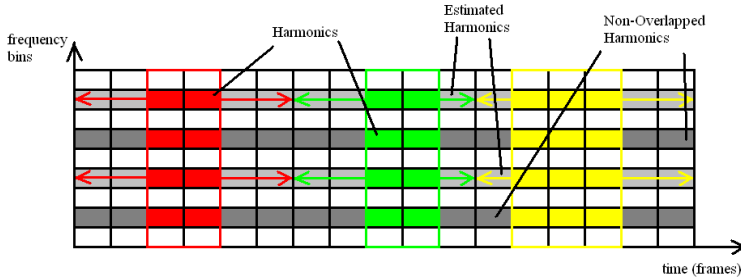
Fig. 3. Method Used to Initialize TFs.

After the TF initialization, the *Isolated Notes TF Reconstruction* stage extends these isolated performance regions for the remaining parts of the note performances that contain isolated performance sections. Each region is evaluated in turn using information from the list of time frames for note performances which contain regions of isolated performance. The note for each time frame in the current region, and notes performed by other sources in the same time frame are determined so that a binary harmonic mask can be generated. This is then used to extract the non-overlapped harmonics for the note during sections of non-isolated performance (shown in Fig. 4(a)), which are then passed to the TF representation for relevant time frames.





(a) TFs with Non-Overlapped Harmonics Added



(b) TFs with Overlapped Harmonics Estimated Using Harmonic Information in Isolated Regions

Fig. 4. Method Used to Reconstruct TFs.

Having used non-overlapped harmonic information to update the isolated note TF representation, we can begin to estimate the overlapped harmonics for the relevant time frames. By using harmonic information available in isolated regions (for which information on all harmonics are available), amplitudes of overlapped harmonics can be estimated. Phase information for the overlapping harmonics is obtained from the corresponding harmonics in the separated TF representations found from the *Pitch and CAM Separation* stage.

As detailed earlier, each set of time frames for each source, relating to non-isolated notes containing an isolated section, are derived from time frames of corresponding isolated regions. Based on the boundary time frames, i.e. the first and last time frames of the isolated regions, we can estimate overlapped harmonic amplitudes (shown in figure 4(b)) by using the spectral information in these frames as templates. We use the first time frame frequency information in an isolated region to process previous time frames, and use the last time frame in the isolated region to process subsequent time frames. According to the CAM principle, amplitude envelopes are assumed to be the same for all harmonics. Hence, by following harmonic envelopes for the subsequent or previous time frames, we can determine the amplitude ratio  $r_{t_0 \rightarrow t}$  between the template time frame  $t_0$  and the time frame currently being processed  $B_t^h$  associated with harmonic  $h$  in time frame  $t$

$$B_t^h = r_{t_0 \rightarrow t} B_{t_0}^h \quad (1)$$

1 Hence, by multiplying bins associated with an overlapped harmonic from the template  
2 frame with the ratio between frames, the amplitude for the corresponding bins in frame  $t$   
3 can be found.

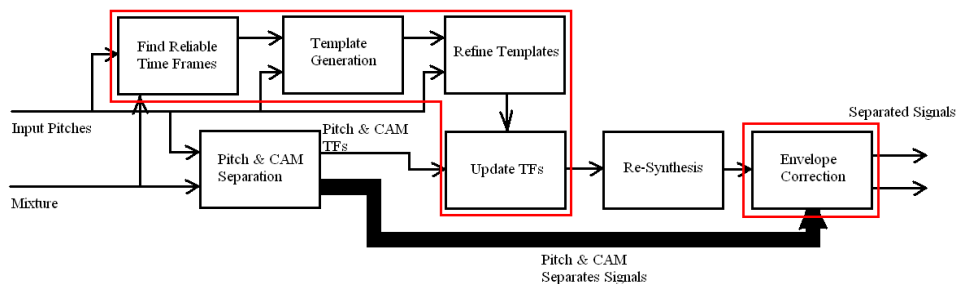
4 Once the TF information for notes with isolated performance regions has been constructed,  
5 it can be converted to the time domain as time-amplitude representation by the *Isolated Note*  
6 *Re-Synthesis* stage. The method is an adapted method used in [67]. Full frequency spectra are  
7 recreated from the half frequency spectra used in the TF representations, and the  
8 overlapped-add method is used to reconstruct the time-amplitude signals for each source.  
9 The system is designed to separate mixture signals comprising of two sources. Therefore,  
10 time domain signals of notes with isolated performance regions can be removed from the  
11 mixture signal to reveal the separated signal for the remaining source. We can simply  
12 subtract isolated note signal sample values from corresponding mixture signal sample  
13 values to generate the 'extracted' signals (performed by the *Onset Note Extraction* stage).

14 Finally, the *Merge Signals* stage uses isolated note signals, and the 'extracted' remaining  
15 signal, to update the separated signals obtained using the baseline pitch and CAM method.  
16 When isolated note information is available (determined by checking for a non-zero sample  
17 value) the final signal is updated with the corresponding sample in the isolated note signal  
18 for the current source being processed. The corresponding sample value is used to update  
19 the signal for the 'other' source, i.e. with the extracted signal. When isolated note  
20 information is unavailable (if a sample value of zero is encountered) the corresponding  
21 sample in the pitch and CAM separated signal, from the respective source, is used to update  
22 the final signal.

#### 23 **4. Spectral template method**

24 This method aims to generate a database of spectral envelope templates of the sources from  
25 the mixtures, and then use the templates to resolve the overlapped harmonics when the  
26 pitch and CAM information is known to be unreliable. In this method, we generate a  
27 spectral envelope template for each note, using the information from the mixture.  
28 Eventually, it builds a database of spectral envelopes for all notes that are performed for  
29 each source, e.g. spectral envelopes for notes C4, E5, D# etc. The note information occurring  
30 in the mixture can be determined from the supplied pitch information. In particular, we use  
31 the non-overlapped harmonics from the most reliable sections of the mixture to fill in the  
32 spectral template for each note that appears in the mixture, where the most reliable section  
33 is regarded as the time section having the most non-overlapped harmonics for a particular  
34 instance of a note occurrence. The number of non-overlapped harmonics can vary,  
35 depending on the other notes being played simultaneously. Within this most reliable time  
36 section, the frequency spectrum at the time frame in which the largest harmonic occurs is  
37 used to train the template. Other occurrences of the note within the mixture are used to  
38 update the template for remaining unknown harmonics by analysing the ratio to adjacent  
39 non-overlapped harmonics (CAM information), based on the extraction of the 'exposed'  
40 non-overlapping harmonics. For example, when the note C5 from a source is being played  
41 together with another note G6, the 'exposed' non-overlapped harmonics of C5 can be used  
42 to train the C5 note template. Other occurrences of C5 from the same source, whilst the note  
43 A7 from the other source is being played, would 'expose' a different set of non-overlapping  
44 harmonics. These non-overlapped harmonics can be used to update the spectral template in

1 order to ‘fill out’ the unknown harmonics by using the relative amplitudes of the harmonics.  
 2 This provides a ‘backup’ set of information for the estimation of the overlapped harmonics  
 3 and also enables us to better handle situations where other information for resolving  
 4 overlapping harmonics is limited or unreliable e.g. concurrent same note occurrence. Figure  
 5 5 shows the diagram of the proposed system that uses the spectral envelope model, which  
 6 uses the *Pitch and CAM Separation* algorithm (developed by Li, Woodruff and Wang [67]) as  
 7 a basis and adds several components shown in large red blocks (implemented in Matlab),  
 8 including *Find Reliable Time Frames*, *Template Generation*, *Refine Templates*, *Update TFs*, and  
 9 *Envelope Correction*, discussed next.



10

11 Fig.5. Diagram of the spectral template method for audio mixture separation.

12 The proposed spectral template system has two inputs: the mixture signal and pitch  
 13 information. The input signals are the audio mixtures we attempt to separate, which can be  
 14 a time-domain representation. Pitch information of each source can be extracted from the  
 15 time-frequency representation of the signals, using a pitch estimator or a pitch tracker, as  
 16 done in the pitch and CAM system [67]. In our proposed system, however, we use the  
 17 supplied pitch information as inputs, and this essentially eliminates the influence of pitch  
 18 estimation process on the separation performance. The pitch information is needed by the  
 19 pitch and CAM algorithm (shown in the *Pitch and CAM Separation* stage in Figure 5) for  
 20 producing an initial estimate of the sources from the TF representations of the mixtures. It  
 21 is also used in *Find Reliable Time Frames* stage to determine the time frames within the TF  
 22 representations that would convey the most reliable harmonic information. These time  
 23 frames are then passed onto the *Template Generation* stage, and the harmonic information  
 24 from these frames is used to initialize the template. In the *Refine Templates* stage, the missing  
 25 harmonics of each template are estimated from the templates of other notes, when limited  
 26 information is available in the mixture. The *Update TFs* stage then uses the templates at time-  
 27 frames with non-overlapped harmonics to resolve the overlapped harmonics (by the *Pitch*  
 28 *and CAM Separation* stage). These modified TF representations are passed onto the *Re-*  
 29 *Synthesis* stage for the reconstruction of the time domain signals of each source. The *Envelope*  
 30 *Correction* stage obtains envelope information by subtracting all but the current source from  
 31 the mixture, and then use it to correct the envelope for time regions of the sources where the  
 32 template was used.

33 In the *Pitch and CAM Separation* stage, we use the baseline algorithm developed by Li,  
 34 Woodruff and Wang [67] to separate the audio mixtures, using the additional pitch contour  
 35 information. More specifically, the audio mixture is transformed to the TF domain using  
 36 short-time Fourier transform (STFT) with overlaps between adjacent time frames. TF

1 representations are generated for each separated source by the *pitch and CAM separation*  
 2 algorithm, and are updated in later processing stages with improved information for time  
 3 frames of unreliable information before being finally transformed back to the time domain.  
 4 The separated time domain signals are also used in the *Envelope Correction* stage to obtain  
 5 envelope information for the refinement of the separated signal.

6 In the *Find Reliable Time Frames* stage, we first find the time frames of the mixture that are  
 7 most likely to yield the best set of harmonics, and we then use them to generate the  
 8 spectral templates. Notes played by different instruments may have different harmonic  
 9 structures, and many of them contain unreliable harmonic content. This is especially true  
 10 for attack sections of many notes, due to the sharp transients and the noise content in the  
 11 attack. For example, when a string is struck, its initial oscillations caused by the initial  
 12 displacement will be non-periodic, and it takes a short amount of time for the string to  
 13 settle into stable resonances of the instrument, and hence provide more reliable harmonic  
 14 information. Some instruments may have long, slow and weak attack section, and in such  
 15 a case, the harmonic content only becomes reliable at some time after the onset of the note.  
 16 A similar problem also happens for notes of a short duration. In order to provide reliable  
 17 frequency information for updating the note templates, we generate a list of time frames  
 18 that does not include time frames containing short note performances and attack sections  
 19 of note performances.

20 The pitch information is supplied to the *Find Reliable Time Frames* stage of the system in the  
 21 form of fundamental frequencies for each source and for all time frames. The fundamental  
 22 frequencies of the notes are converted into numbers representing the closest note in the 12  
 23 tone scale, i.e. C0 is 1, C#0 is 2, B0 is 12, C1 is 13 and so on up to 96 representing note B7. To  
 24 find the corresponding note numbers for each frequency in the input pitch information we  
 25 first determine which octave range the frequency is in by selecting an integer  $m$  such that

26  $f_{\min} < \frac{f}{2^m} \leq f_{\max}$ , where the lower and upper frequency limits of the first octave (C0 to B0)

27 are  $f_{\min}$  and  $f_{\max}$  respectively and  $f$  is the (fundamental) frequency value that we wish  
 28 to convert to a note number. In practice,  $f_{\min}$  is selected as the frequency value between C0  
 29 and the note that is one semi-tone lower (in theory, note B1), and  $f_{\max}$  is selected as the  
 30 frequency value between B0 and C1. The integer  $m$  can be determined by repeatedly  
 31 halving the frequency until it falls within the first octave range. Once the octave range has  
 32 been found, the note from A to G# on the 12-tone scale can then be found by further  
 33 narrowing the searching range in terms of multiples of  $f_{\min}$ . In other words, we choose the  
 34 integer  $n$  that satisfies the following inequality

$$35 \quad \left(\sqrt[12]{2}\right)^{n-1} f_{\min} < \frac{f}{2^m} \leq \left(\sqrt[12]{2}\right)^n f_{\min} \quad (2)$$

36 where  $m$  is the octave range value found previously. Once the octave range  $m$  and the note  
 37  $n$  are found, the list of note numbers at each time frame for each source can be easily  
 38 calculated. From the list of notes selected, we further remove the invalid notes if they are  
 39 from the attack sections of the notes, or their duration is too short. In our case, any notes  
 40 whose duration is shorter than six time frames will be set to zero.

1 In the *Template Generation* stage, we update the spectral templates for each note with spectral  
 2 information from the list of time frames that contains the valid notes obtained above. We  
 3 search over each time frame in the list (also for each source in turn), and ignore the invalid  
 4 time frames (with values of zero). We then determine the note performed by the current  
 5 source and the notes by all other sources for the current time frame. Using such a particular  
 6 note combination, we can generate binary harmonic masks to extract the non-overlapping  
 7 harmonics from the TF representation of the mixtures for each of the frames. More  
 8 specifically, the note performed by the current source is used to determine the frequencies of  
 9 all harmonics. Notes performed simultaneously by all other sources are used to determine  
 10 which of the current source harmonics are overlapped by all other sources thus, indicating  
 11 the ‘exposed’ non-overlapping harmonics for the current note. Using such information, we  
 12 can set frequency bins that are associated with non-overlapped harmonics to 1 and all other  
 13 bins to 0. Firstly, the frequency of the note value for the current source must be found.  
 14 According to the international standard (ISO 16 [68]), the note frequency for A4 is 440Hz,  
 15 and note C0 is 57 semi-tones below A4. Hence, the frequency of C0 can be used as a basis to  
 16 find the fundamental frequency of other notes such as A4 using  $f = f_{C0} \left( \sqrt[12]{2} \right)^{p-1}$ , where  $p$   
 17 is the note value and  $f_{C0}$  is the fundamental frequency of note C0. We then associate  
 18 frequency bin  $b$  to harmonic  $h_i$  for the current source  $i$  using a similar method to that in  
 19 [67], if it satisfies  $|bf_a - h_i f| < \theta_1$ , where  $\theta_1$  is a threshold and  $f_a$  is the frequency resolution  
 20 of the TF representation (both  $\theta_1$  and  $f_a$  are determined previously in the *Pitch and CAM*  
 21 *Separation* stage). We use a second threshold  $\theta_2$  to define the range in which the current  
 22 source harmonic  $h_i$  is overlapped with any other source harmonic  $h_j$ , i.e.  $|f_{h_j} - h_i f| < \theta_2$ ,  
 23 where  $f_{h_j}$  is the frequency of harmonic  $h_j$ . Again, this is a similar method to that in [67],  
 24 hence  $\theta_2$  can be chosen in the same way as used in the *Pitch and CAM Separation* stage. As a  
 25 result, we can define a TF mask  $M_b$  which takes 1 if  $|bf_a - h_i f| < \theta_1$ , otherwise 0. This binary  
 26 mask is then used to extract all non-overlapped harmonics for all time frames from the TF  
 27 representation of the mixture. All the harmonic sets for the current note combination are  
 28 evaluated to find the set that contains the largest amplitude harmonic, which is then used to  
 29 update the note template (or simply stored if the template is empty and has not yet been  
 30 initialized). We continue to go through the whole list of valid note regions, and when a new  
 31 note combination is encountered, we update the note templates based on the new harmonic  
 32 mask generated using the new set of ‘exposed’ non-overlapped harmonics. After all note  
 33 combinations have been evaluated, the note templates may contain several sets of harmonics  
 34 for each note combination. If this happens, we merge them to create a final set for the note  
 35 template. Note that one may wish to apply scaling to each set of harmonic templates to  
 36 ensure harmonics are of correct magnitude when merging the template.

37 As done in the *Refine Templates* stage, the spectral templates generated, are further refined  
 38 and improved by using information from all the templates. The reason that the spectral  
 39 templates need to be refined is because for some notes, there may be only a limited set of  
 40 non-overlapped harmonics, as some harmonics may not be available in the mixture. To  
 41 improve the templates, harmonic information from other note templates that are available  
 42 within a specified range of notes is used. Spectra of other note templates are pitch shifted to

1 match the note we intend to improve, so that information for correlating harmonics can be  
 2 obtained (after harmonics are aligned). However, spectral quality tends to deteriorate as the  
 3 degree of pitch shifting increases. Therefore we first use the templates of notes that are  
 4 closest in frequency to the note template for which we wish to improve, and then continue  
 5 with templates of decreasing quality. In addition, lower frequency note templates yield  
 6 higher quality spectra when the pitch is shifted up to match the frequency of the note  
 7 template we wish to improve, and vice versa. Hence, we limit the range of notes and also  
 8 the number of note templates to be used for improving note templates. This essentially  
 9 excludes note templates that have been excessively pitch shifted, and also improves  
 10 computational efficiency of the proposed system.

11 In the *Update TFs* stage, we update the TF representations of the separated sources from the  
 12 *Pitch and CAM Separation* stage, using the note templates. Pitch information is used to  
 13 determine, for each source, the time frames where reliable non-overlapped harmonics are  
 14 unavailable for separation. As already mentioned, if a source is playing a note which is one  
 15 octave lower than a note played by another source, the former one would have every other  
 16 harmonic overlapped whereas the harmonics of the latter one would be totally overlapped  
 17 by those of the former one. As a consequence, no reliable information is available to resolve  
 18 the overlapped harmonics of the latter source. However, there are many other note  
 19 combinations, leading to unavailable non-overlapping harmonics to be used to resolve the  
 20 overlapped harmonics, e.g. when one source performs a note 7 semi-tones higher (perfect  
 21 fifth interval) than the other it would result in every third harmonic of the latter source  
 22 being overlapped by the former one. Of course, it would be exhaustive to find all possible  
 23 combinations of notes that result in all of the source harmonics being overlapped. Using  
 24 pitch information is an efficient way to calculate the resulting number of overlapped  
 25 harmonics at each time frame for each source. The number of overlapping harmonics  $\varphi_i(t)$   
 26 for source  $i$  at time frame  $t$  can be determined by finding the number of harmonics in a  
 27 complete set  $H_{N_i}(t)$  that is not in the set of non-overlapped harmonics  $\tilde{H}_{N_i}(t)$  based on the  
 28 pitch information of note  $N_i(t)$ .

29 We use the same method discussed above to generate binary masks, using current note  
 30 information and information on all other notes that are performed simultaneously. We also  
 31 create a binary mask with a complete set of harmonics from which the mask with non-  
 32 overlapped harmonics is subtracted. This gives a mask containing harmonics that are  
 33 overlapped. Evaluating the magnitude at bins closest to the expected harmonic frequencies  
 34 allows the number of overlapped harmonics present to be determined. For all  $t$  where  
 35  $\varphi_i(t)=0$  i.e. time frames for source  $i$  that have no reliable information for source  
 36 separation, frequency spectra for the respective note templates are used to replace the  
 37 frequency spectra in the TF representation of the separated source.

38 The *Re-Synthesis* stage, adapted from [67], involves the reconstruction of the time domain  
 39 signals from the TF representations for each source. Specifically, symmetric frequency  
 40 spectra are created from the half spectra used in the TF representations and the overlap-add  
 41 method is used to generate the time domain signals.

42 No amplitude envelope information has been conveyed in the note templates for refining  
 43 the separated sources. Hence, in the *Envelope Correction* stage, for the time regions with

1 unresolved overlapped harmonics, the amplitude envelopes of the separated sources will be  
 2 corrected. All sources that have been separated (in the *Pitch and CAM Separation* stage)  
 3 except the current source, for which the envelope is being corrected, are removed from the  
 4 original mixture signal. The remaining signal would then be a crude representation of the  
 5 source we are attempting to correct as most of the high energy components from all other  
 6 sources are removed. The envelope of the remaining signal is found by finding peaks of  
 7 absolute amplitude values. We detect peaks at time instances where the first order  
 8 derivative of the absolute time-amplitude signal is zero. The envelopes of the separated  
 9 sources are then adjusted by applying a certain amount of scaling determined by the desired  
 10 envelope obtained above.

## 11 **5. System evaluation**

### 12 **5.1 Evaluation method**

13 The system is evaluated using test signals specifically designed to highlight differences  
 14 between the proposed systems and the original pitch and CAM separation system. The  
 15 proposed systems aim to address the weak points of the pitch and CAM system, i.e. the lack of  
 16 time domain detail arising from poor separation of attack regions of notes, and its difficulty in  
 17 resolving the overlapping harmonics due to similar note performances. Hence, tests were  
 18 designed to evaluate differences in these particular points between the proposed systems and  
 19 the original system, rather than an evaluation of overall performance of the system.

20 For the proposed isolated note system, test signals which were generated using real  
 21 instrument recordings with different musical scores, contain isolated performances of notes  
 22 in order to show the effectiveness of the proposed system. The isolated note system aims to  
 23 better resolve attack sections of notes for which the pitch and CAM system performs poorly.  
 24 Hence, instruments with fast attacks and relatively higher energy in the higher frequency  
 25 range (of the attacks), e.g. instruments that are struck, or particular instruments that are  
 26 plucked were selected for the test signals. Two test signals meeting these criteria were  
 27 generated; the first signal (test signal 1) was a two-source mixture containing a cello and a  
 28 piano performance, the cello was played throughout the signal and the piano had sections of  
 29 performance interspersed with sections of silence giving the cello regions of isolated  
 30 performance. The second signal (test signal 2) was also a two-source mixture containing a  
 31 string section and a guitar performance, again, the string section was played throughout the  
 32 test signal and the guitar had interspersed sections of silence. Both test mixtures were  
 33 created by mixing clean source signals (16-bit, 44100Hz sample rate).

34 For the spectral template system, two test signals with the same musical score are generated  
 35 containing sections with the same note performance and also sections with sufficient  
 36 information to train the templates. The first piece was a two source mixture of a cello and a  
 37 piano, the second piece was a two source mixture of a cello and a clarinet (both pieces  
 38 approximately four seconds long at 16 bit, 44100 kHz sampling rate). All the test signals  
 39 were created using ProTools music production software and instruments were selected to  
 40 avoid synthesized replications to achieve performances as realistic as possible (this avoids  
 41 signals being created with stable frequency spectra for note performances). A database of  
 42 real recordings of instruments within the music production software was used to generate  
 43 the test signals. Pitch and CAM separation was performed with default values.

1 System performance is evaluated by calculating the SNR for the pitch and CAM system and  
 2 the proposed system with each test signal using

$$3 \quad SNR(dB) = 10 \log_{10} \frac{\sum_n (x[n])^2}{\sum_n (x[n] - \hat{x}[n])^2} \quad (3)$$

4 where  $x[n]$  is the original signal and  $\hat{x}[n]$  is the separated signal ( $n$  is the sample index).  
 5 This allows us to quantify the sample-wise resemblance between the clean source signals  
 6 and the separated signals generated by each of the systems.

7 For the evaluation of the isolated system, a direct comparison of SNR values for both  
 8 systems would reveal the gains made by the isolated note system. However, differences in  
 9 the attack sections only are difficult to quantify when evaluating the entire signal as they  
 10 make up only a small proportion of the test signal. Hence, we expect the differences in  
 11 perceptual quality to be more significant (i.e. differences would be heard, but are not  
 12 represented as well in comparison using SNR measurements). Therefore, a listening test was  
 13 also performed to observe the perceptual difference between the separated signals obtained  
 14 using the pitch and CAM and the isolated note methods. Test signals for the listening test  
 15 were generated by including the original clean source signal, followed by a one second  
 16 silence, and then the separated signal allowing for a direct comparison to be made between  
 17 the clean source and separated signals. 26 participants were asked to score the signals from  
 18 0 to 5, with 0 being extremely poor and 5 being perceptually transparent (with reference to  
 19 the original signal). Scores were based on the details of attack sections as well as overall  
 20 separation performance between the two systems (i.e. which system 'sounds better') all test  
 21 signals were presented in a random order for each participant.

22 For the evaluation of the spectral template system, the separated signals are modified to  
 23 remove the pitch and CAM sections so that the signals contain only same note  
 24 performances, and the influence of the pitch and CAM results is ignored. Test signals are  
 25 created by including the original signal at the start, followed by a one second silence, and  
 26 then followed by the separated signal; this allows the listener to hear the original before  
 27 hearing the separated signal so a direct comparison can be made. Test signals were  
 28 generated for both pitch and CAM and note template systems. All test signals were played  
 29 in a random order so that identification of each system remains unknown and cannot be  
 30 anticipated. Signals were allowed to be repeated as many times as needed to assess signal  
 31 quality.

## 32 **5.2 Results**

33 The results of the isolated note system are shown in Tables 1 and 2. When comparing results  
 34 for test signal 1, source 1 (cello), we observe a reduction of -3.75 dB in SNR between the two  
 35 systems. Nevertheless, this source contains sections of isolated performance which we use to  
 36 better separate attack sections of source 2 (for which this study concerns). As can be seen for  
 37 source 2 (piano), SNR of the proposed system is 15.08 dB higher than the pitch and CAM  
 38 system, so a significant gain in separation performance is achieved. Looking at SNR results  
 39 for test signal 2 source 1 (string section) we see a marginal increase of 0.34 dB in separation  
 40 performance from the isolated note system, again, this source contains the isolated region of



1 performance which is used to improve separation of source 2. For source 2 (guitar), we see  
2 a significant improvement in separation performance by the isolated note system with a SNR  
3 8.44 dB higher than the pitch and CAM system.  
4

Test Signal	Source	Pitch and CAM System	Isolated Note System
1	1	19.04	15.29
	2	5.87	20.95
2	1	15.78	16.09
	2	3.63	12.07

5 Table 1. SNR (dB) results for Isolated Note System as compared with the pitch and CAM  
6 system.

Test Signal	Source	Pitch and CAM Mean Score	Isolated Note Mean Score
1	1	4.88	4.73
	2	2.50	4.50
2	1	3.85	3.69
	2	1.65	3.54

7  
8 Table 2. Listening test results for Isolated Note System as compared with the pitch and CAM  
9 system.

10 For test signal 1 we can see similar mean opinion scores for separation of source 1 by both  
11 systems suggesting a similar level of separation performance between the two systems.  
12 However, listening test results suggest a significant improvement of separation performance  
13 by the isolated note system for source 2. For test signal 1 and source 2, the pitch and CAM  
14 system achieved a mean score of 2.50 and the isolated note system achieved a mean score of  
15 4.50. Again, the isolated note system achieved similar separation performance compared to the  
16 pitch and CAM system for test signal 2, source 1, while giving a significant improvement for  
17 source 2. The pitch and CAM system achieved a mean score of 1.65 whereas the isolated note  
18 achieved a higher mean score of 3.54. Both SNR and listening test results indicate that the note  
19 isolation separation system achieves better separation performance. We can see significant  
20 quantitative gains from the SNR results for signals with fast attacks (source 2 in both test  
21 signals 1 and 2). Qualitative results from the listening test also show significant perceptual  
22 gains obtained in the separation of attack sections in addition to overall separation.

23 The results of the spectral template system are summarised in Tables 3 and 4. Table 3 shows  
24 SNR results for the proposed note template separation system compared with the pitch and  
25 CAM separation system. We can see that for both test signals, we have the same separation  
26 performance for source 1 (cello). Sufficient harmonic information is available for source 1 to  
27 resolve overlapping harmonics so the note template system also uses the pitch and CAM  
28 method to separate the signal which is why the same performance result can be observed.  
29 However, for source 2 (piano), SNR results appear to be poor. For test signal 1 we see that the  
30 pitch and CAM system has a SNR of 0.79 dB whereas the note template system has a SNR of -  
31 2.35 dB, suggesting that the level of noise introduced by the system is greater than the level of  
32 input signal. Likewise, test signal 2 shows poor SNR results for source 2, the pitch and CAM  
33 system has a SNR of 2.90 dB while the note template system has a SNR of -3.65 dB.

Test Signal	Source	Pitch and CAM System	Note Template System
1	1	2.62	2.62
	2	0.79	-2.35
2	1	7.79	7.79
	2	2.90	-3.65

Table 3. SNR (dB) results for Note Template System as compared with the pitch and CAM system.

Test Signal	Source	Pitch and CAM System	Note Template System
1	1	4.08	3.77
	2	1.96	0.92
2	1	4.77	4.81
	2	1.65	0.92

Table 4. Listening test results for Note Template System as compared with the pitch and CAM system.

Table 4 shows average results for the listening test for the pitch and CAM separation system and the note template separation system. For test signal 1, source 1, we see a mean score of 4.08 for the pitch and CAM separation system and a mean score of 3.77 for the note template system despite the same pitch and CAM separated signal being used by both systems, as explained earlier. For test signal 1, source 2, we see a mean score of 4.77 for the pitch and CAM system. We see a reduction of the score for the note template system, with a mean score of 0.92. Comparing scores for test signal 2, similar scores for source 1 can be seen for both systems, with the pitch and CAM system scoring a mean of 4.77 and the note template system scoring a mean of 4.81. Again, both systems use the pitch and CAM separated signals for source 1, as explained earlier. The score for the note template system is lower than the score for the pitch and CAM system, for test signal 2, source 2; we see a mean score of 1.65 for the pitch and CAM system and a mean score of 0.92 for the note template system. The spectral template system does not work as promising as we would have expected, due to the following possible reasons. The templates trained from mixtures may not be accurate enough to represent the sources, because of the limited number of non-overlapped harmonics and isolated notes within the mixture. Using clean music source data (instead of the monaural mixture) to train the templates may be able to mitigate this problem and further to improve the results. Also, in the proposed template systems, pitch shifting which was used to fill up the missing notes that are not available in the mixture, apparently introduces errors in harmonic estimation. These are interesting points for future investigation.

## 6. Conclusions

We have presented two new methods for music source separation from monaural mixture using the isolated note information and note spectral template, both evaluated from the sound mixture. The proposed methods were designed to improve the separation performance of the baseline pitch and CAM system especially for the separation of attack sections of notes, and overlapping time-frequency regions. In the

1 pitch and CAM system, the fast attack sections are almost completely lost in the  
2 separated signals, resulting in poor separation results for the transient part of the signal.  
3 In the proposed isolate note system, accurate harmonic information available in the  
4 isolated regions is used to reconstruct harmonic content for the entire note performance,  
5 and so, the harmonic content can be removed from the mixture to reveal the remaining  
6 note performance (in a two-source case). The isolated note system has been shown to be  
7 successful in improving the separation performance of attack sections of notes, offering a  
8 large improvement in separation quality over the baseline system. In the proposed note  
9 template system, the overlapping time-frequency regions of the mixtures are resolved  
10 using the reliable information from the non-overlapping regions of the sources, based on  
11 the spectral template matching. Preliminary results show that the spectral templates  
12 evaluated from the mixtures can be noisy and may degrade the results. Using spectral  
13 template generated directly from clean training data (i.e. containing single signals,  
14 instead of mixtures) has the potential to improve the system performance which will be  
15 our future study.

## 16 7. Future directions

17 We have studied the potentials of using spectral template and isolated note information  
18 for music sound separation. A major challenge is however to identify the regions from  
19 which the note information can be regarded as reliable and thereby used to estimate the  
20 note information for the unreliable and overlapped regions. Under noisy and multiple  
21 source conditions, more ambiguous regions may be identified, and using such  
22 information may further distort the separation results. Pitch information is relatively  
23 reliable under noisy conditions and can be used to improve the system performance [81].  
24 Another potential direction is to use the property of the sources and noise/interferences,  
25 such as sparseness, to facilitate the identification of the reliable regions within the mixture  
26 that can be used to estimate the sources [74-77]. This is mainly due to the following three  
27 reasons. Firstly, as mentioned earlier, music audio can be made sparser if it is transformed  
28 into another domain, such as the TF domain, using an analytically pre-defined dictionary  
29 such as discrete Fourier transform (DFT) or discrete cosine transform (DCT) [69] [70].  
30 Recent studies show that signal dictionaries directly adapted from training data using  
31 machine learning techniques, based on some optimisation criterion (such as the  
32 reconstruction error regularised by a sparsity constraint), can offer better performance  
33 than the pre-defined dictionary [71] [72]. Secondly, the sparse techniques using learned  
34 dictionary have been shown to possess certain denoising capability for corrupted signals  
35 [72]. Thirdly, identification of reliable regions from sound mixtures, and estimation of the  
36 probability of each TF point dominated by a source can be potentially cast as an audio-  
37 inpainting [73] or matrix completion problem. This naturally links the two important  
38 areas: source separation and sparse coding. Hence, the emerging algorithms developed in  
39 the sparse coding area could be potentially used for the CASA based monaural separation  
40 system. Separating music sources from mixtures with uncertainties [78] [79], such as  
41 under the condition of unknown number of sources, is also a promising direction for  
42 future research, as required in many practical applications. In addition, online  
43 optimisation will be necessary when the separation algorithms operate on resource  
44 limited platforms [80].

## 8. References

- [1] Jutten, C., & Herault, J. (1991). Blind Separation of Sources, Part I: An Adaptive Algorithm Based on Neuromimetic Architecture, *Signal Processing*, vol. 24, pp. 1-10.
- [2] Cardoso, J.-F., & Souloumiac, A. (1993). Blind Beamforming for Non Gaussian Signals, *IEE Proc. F, Radar Signal Processing*, vol. 140, no. 6, pp. 362-370.
- [3] Comon, P. (1994). Independent Component Analysis: a New Concept?, *Signal Processing*, vol. 36, no. 3, pp. 287-314.
- [4] Bell, A. J., & Sejnowski, T. J. (1995). An Information Maximization Approach to Blind Separation and Blind Deconvolution, *Neural Computation*, vol. 7, no. 6, pp. 1129-1159.
- [5] Amari, S.-I., Cichocki, A., & Yang, H. (1996). A New Learning Algorithm for Blind Signal Separation, *Advances Neural Information Processing System*, vol. 8, pp. 757-763.
- [6] Cardoso, J.-F. (1998). Blind Signal Separation: Statistical Principles, *Proceedings of the IEEE*, vol. 9, no. 10, pp. 2009-2025.
- [7] Lee, T.-W. (1998). *Independent Component Analysis: Theory and Applications*. Boston, MA: Kluwer Academic.
- [8] Haykin, S. (2000). *Unsupervised Adaptive Filtering, Volume 1, Blind Source Separation*. New York: Wiley.
- [9] Hyvärinen, A., Karhunen, J., & Oja, E. (2001). *Independent Component Analysis*. New York: Wiley.
- [10] Cichocki, A., & Amari, S.-I. (2002). *Adaptive Blind Signal and Image Processing, Learning Algorithm and Applications*. New York: Wiley.
- [11] Cardoso, J. F., & Laheld, B. (1996). Equivariant Adaptive Source Separation, *IEEE Transactions Signal Processing*, vol. 44, no. 12, pp. 3017-3030.
- [12] Belouchrani, A., Abed-Meraim, K., Cardoso, J., & Moulines, E. (1997). A Blind Source Separation Technique Using Second-Order Statistics, *IEEE Transactions Signal Processing*, vol. 45, no. 2, pp. 434-444.
- [13] Thi, H., & Jutten, C. (1995). Blind Source Separation for Convolutional Mixtures, *Signal Processing*, vol. 45, pp. 209-229.
- [14] Smaragdakis, P. (1998). Blind Separation of Convolved Mixtures in the Frequency Domain, *Neurocomputing*, vol. 22, pp. 21-34.
- [15] Parra, L., & Spence, C. (2000). Convolutional Blind Source Separation of Nonstationary Sources, *IEEE Transactions Speech Audio Processing*, vol. 8, no. 3, pp. 320-327.
- [16] Rahbar, K., & Reilly, J. (2001). Blind Source Separation of Convolved Sources by Joint Approximate Diagonalization of Cross-Spectral Density Matrices, in *Proc. IEEE International Conference Acoustics, Speech and Signal Processing*, Utah, USA.
- [17] Davies, M. (2002). Audio Source Separation, in *Mathematics in Signal Separation V*. Oxford, U.K.: Oxford Univ. Press.
- [18] Sawada, H., Mukai, R., Araki, S., & Makino, S. (2004). A Robust and Precise Method for Solving the Permutation Problem of Frequency-Domain Blind Source Separation, *IEEE Transactions Speech and Audio Processing*, vol.12, no. 5, pp. 530-538.
- [19] Wang, W., Sanei, S., & Chambers, J.A. (2005). Penalty Function Based Joint Diagonalization Approach for Convolutional Blind Separation of Nonstationary Sources, *IEEE Transactions on Signal Processing*, vol. 53, no. 5, pp. 1654-1669.
- [20] Sawada, H., Araki, S., & Makino, S. (2010). Underdetermined Convolutional Blind Source Separation via Frequency Bin-Wise Clustering and Permutation Alignment, *IEEE Transactions on Audio Speech and Language Processing*, vol. 18, pp. 516-527.

- 1 [21] Pedersen, M., Larsen, J., Kjems, U., & Parra, L. (2007). A Survey on Convolutional Blind  
2 Source Separation Methods, *Handbook on Speech Processing and Speech Communication*,  
3 Springer.
- 4 [22] Belouchrani, A., & Amin, M. G. (1998). Blind Source Separation Based on Time-  
5 Frequency Signal Representations, *IEEE Transactions Signal Processing*, vol. 46, no. 11,  
6 pp. 2888-2897, 1998.
- 7 [23] Chen, S., Donoho, D. L., & Saunders, M. A. (1998). Atomic Decomposition by Basis  
8 Pursuit, *SIAM Journal Scientific Computing*, vol. 20, no. 1, pp. 33-61.
- 9 [24] Lee, T., Lewicki, M., Girolami, M., & Sejnowski, T. (1998), Blind Source Separation of  
10 More Sources Than Mixtures Using Overcomplete Representations, *IEEE Signal*  
11 *Processing Letters*, vol. 6, no. 4, pp. 87-90.
- 12 [25] Lewicki, M. S., & Sejnowski, T. J. (1998). Learning Overcomplete Representations,  
13 *Neural Computation*, vol. 12, no. 2, pp. 337-365.
- 14 [26] Bofill, P., & Zibulevsky, M. (2001). Underdetermined Blind Source Separation Using  
15 Sparse Representation, *Signal Processing*, vol. 81, pp. 2253-2362.
- 16 [27] Zibulevsky, M., & Pearlmutter, B. A. (2001). Blind Source Separation by Sparse  
17 Decomposition in a Signal Dictionary, *Neural Computation*, vol. 13, no. 4, pp. 863-882.
- 18 [28] Li, Y., Amari, S., Cichocki, A., Ho, D. W. C., & Xie, S. (2006). Underdetermined Blind  
19 Source Separation Based on Sparse Representation. *IEEE Transactions on Signal*  
20 *Processing*, vol. 54, no. 2, pp. 423-437.
- 21 [29] He, Z., Cichocki, A., Li, Y., Xie, S., & Sanei, S. (2009). K-Hyperline Clustering Learning  
22 for Sparse Component Analysis. *Signal Processing*, vol. 89, no. 6, pp. 1011-1022.
- 23 [30] Yilmaz, O., & Richard, S. (2004). Blind Separation of Speech Mixtures via Time-  
24 Frequency Masking, *IEEE Transactions on Signal Processing*, vol. 52, no. 7, pp. 1830-1847.
- 25 [31] Wang, D.L. (2005). On Ideal Binary Mask as the Computational Goal of Auditory Scene  
26 Analysis. In Divenyi P. (ed.), *Speech Separation by Humans and Machines*, pp. 181-197,  
27 Kluwer Academic, Norwell MA.
- 28 [32] Mandel, M. I., Weiss, R. J., & Ellis, D. P. W. (2010). Model-based Expectation  
29 Maximisation Source Separation and Localisation, *IEEE Transactions on Audio Speech and*  
30 *Language Processing*, vol. 18, pp. 382-394.
- 31 [33] Duong, N. Q. K., Vicent, E., & Gribonval, R. (2010). Under-determined Reverberant  
32 Audio Source Separation Using a Full-Rank Spatial Covariance Model, *IEEE*  
33 *Transactions on Audio Speech and Language Processing*, vol. 18, pp. 1830-1840.
- 34 [34] Plumbley, M. D. (2003). Algorithms for Nonnegative Independent Component  
35 Analysis, *IEEE Transactions on Neural Networks*, vol. 14, no. 2, pp. 534-543.
- 36 [35] Kim, T., Attias, H., & Lee, T.-W. (2007). Blind Source Separation Exploiting Higher-  
37 Order Frequency Dependencies, *IEEE Transactions on Audio Speech and Language*  
38 *Processing*, vol. 15, pp. 70-79.
- 39 [36] Comon, P., & Jutten, C., eds. (2010). *Handbook of Blind Source Separation, Independent*  
40 *Component Analysis and Applications*, Academic Press.
- 41 [37] Smaragdis, P. (2004). Non-Negative Matrix Factor Deconvolution; Extraction of  
42 Multiple Sound Sources from Monophonic Inputs. In *Proceedings of the 5th International*  
43 *Conference on Independent Component Analysis and Blind Signal Separation*, Grenada, Spain.
- 44 [38] Schmidt, M. N., & Mørup, M. (2006). Nonnegative Matrix Factor 2-D Deconvolution for  
45 Blind Single Channel Source Separation, in *Proceedings of the International Conference on*  
46 *Independent Component Analysis and Signal Separation*, Charleston, USA.

- 1 [39] Wang, W., Cichocki, A., & Chambers, J. A. (2009). A Multiplicative Algorithm for  
2 Convolutional Non-negative Matrix Factorization Based on Squared Euclidean Distance,  
3 *IEEE Transactions on Signal Processing*, vol. 57, no. 7, pp. 2858-2864.
- 4 [40] Ozerov, A., & Févotte, C. (2010). Multichannel Nonnegative Matrix Factorization in  
5 Convolutional Mixtures for Audio Source Separation, *IEEE Transactions on Audio, Speech  
6 and Language Processing*.
- 7 [41] Mysore, G., Smaragdis, P., & Raj, B. (2010). Non-Negative Hidden Markov Modeling of  
8 Audio with Application to Source Separation. In *Proceedings of the 9th international  
9 conference on Latent Variable Analysis and Signal Separation (LCA/ICA)*. St. Malo, France.
- 10 [42] Ozerov, A., Févotte, C., Blouet, R., & Durrieu, J.L. (2011). Multichannel Nonnegative  
11 Tensor Factorization with Structured Constraints for User-guided Audio Source  
12 Separation, *Proceedings of the IEEE International Conference Acoustics, Speech and Signal  
13 Processing*, Prague, Czech Republic.
- 14 [43] Wang, W. & Mustafa, H. (2011). Single Channel Music Sound Separation Based on  
15 Spectrogram Decomposition and Note Classification, in *Computer Music Modelling and  
16 Retrieval*, Springer.
- 17 [44] Cichocki, A., Zdunek, R., Phan, A.H., & Amari, S. (2009). *Nonnegative Matrix and Tensor  
18 Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source  
19 Separation*. Wiley.
- 20 [45] Brown, G.J., & Cooke, M.P. (1994). Computational Auditory Scene Analysis, *Computer  
21 Speech and Language*, vol. 8, pp. 297-336.
- 22 [46] Wrigley, S. N., Brown, G. J., Renals, S., & Wan, V. (2005). Speech and Crosstalk  
23 Detection in Multi-Channel Audio, *IEEE Transactions on Speech and Audio Processing*, vol.  
24 13, no. 1, pp. 84-91.
- 25 [47] Palomäki, K. J., Brown, G. J., & Wang, D. L. (2004). A Binaural Processor for Missing  
26 Data Speech Recognition in the Presence of Noise and Small-Room  
27 Reverberation, *Speech Communication*, vol. 43, no. 4, pp. 361-378.
- 28 [48] Wang, D.L., & Brown, G. (2006). *Computational Auditory Scene Analysis: Principles,  
29 Algorithms, and Applications*, Wiley/IEEE.
- 30 [49] Shao, Y., & Wang, D.L. (2009). Sequential Organization of Speech in Computational  
31 Auditory Scene Analysis. *Speech Communication*, vol. 51, pp. 657-667.
- 32 [50] Hu, K., & Wang, D.L. (2011). Unvoiced Speech Segregation from Nonspeech  
33 Interference via CASA and Spectral Subtraction. *IEEE Transactions on Audio, Speech, and  
34 Language Processing*, vol. 19, pp. 1600-1609.
- 35 [51] Xu, T., & Wang, W. (2009). A Compressed Sensing Approach for Underdetermined  
36 Blind Audio Source Separation with Sparse Representations, in *Proceedings of the IEEE  
37 International Workshop on Statistical Signal Processing*, Cardiff, UK.
- 38 [52] Xu, T., & Wang, W. (2010). A Block-based Compressed Sensing Method for  
39 Underdetermined Blind Speech Separation Incorporating Binary Mask, in *Proceedings of  
40 the IEEE International Conference on Acoustics, Speech and Signal Processing*, Dallas, Texas,  
41 USA.
- 42 [53] Xu, T., & Wang, W. (2011). Methods for Learning Adaptive Dictionary for  
43 Underdetermined Speech Separation, in *Proceedings of the IEEE 21st International  
44 Workshop on Machine Learning for Signal Processing*, Beijing, China.

- 1 [54] Kim, M., & Choi, S. (2006). Monaural Music Source Separation: Nonnegativity,  
2 Sparseness, and Shift-Invariance, in *Proceedings of the IEEE International Conference on*  
3 *Independent Component Analysis and Blind Signal Separation*, Charleston, USA.
- 4 [55] Virtanen, T. (2006). *Sound Source Separation in Monaural Music Signals*, PhD Thesis,  
5 Tampere University of Technology.
- 6 [56] Virtanen, T. (2007). Monaural Sound Source Separation by Nonnegative Matrix  
7 Factorization With Temporal Continuity and Sparseness Criteria, *IEEE Transactions on*  
8 *Audio, Speech and Language Processing*, vol. 15, no. 3, pp. 1066-1074.
- 9 [57] Ozerov, A., Philippe, P., Bimbot, F., & Gribonval, R. (2007). Adaptation of Bayesian  
10 Models for Single-Channel Source Separation and its Application to Voice/Music  
11 Separation in Popular Songs, *Proceedings of the IEEE International Conference on Acoustics,*  
12 *Speech, and Signal Processing*, Hawaii, USA.
- 13 [58] Richard, G., & David, B. (2009). An Iterative Approach to Monaural Musical Mixture  
14 De-Soloing, in *Proceedings of the IEEE International Conference on Acoustics, Speech, and*  
15 *Signal Processing*, Taipei, Taiwan.
- 16 [59] Klapuri, A., Virtanen, T., & Heittola, T. (2010). Sound Source Separation in Monaural  
17 Music Signals Using Excitation-Filter Model and EM Algorithm, in *Proceedings of the*  
18 *IEEE International Conference on Acoustics, Speech, and Signal Processing*, Dallas, USA.
- 19 [60] Bregman, A. S. (1990). *Auditory Scene Analysis*, MIT Press.
- 20 [61] Li, Y. & Wang, D. L. (2007). Separation of Singing Voice From Music Accompaniment  
21 for Monaural Recordings, *IEEE Transactions on Audio, Speech and Language Processing*,  
22 vol. 15, no. 4, pp. 1475-1487.
- 23 [62] Parsons, T. W. (1976). Separation of Speech from Interfering Speech By Means of  
24 Harmonic Selection, *Journal of the Acoustical Society of America*, vol. 60, no. 4, pp. 911-918,  
25 1976.
- 26 [63] Li, Y. & Wang, D. L. (2009). Musical Sound Separation Based on Binary Time-Frequency  
27 Masking, *EURASIP Journal on Audio, Speech and Music Processing*, article ID 130567.
- 28 [64] Every, M. R. & Szymanski, J. E. (2006). Separation of Synchronous Pitched Notes by  
29 Spectral Filtering of Harmonics, *IEEE Transactions on Audio, Speech and Language*  
30 *Processing*, vol. 14, no. 5, pp. 1845-1856.
- 31 [65] Virtanen, T. & Klapuri, A. (2001). Separation of Harmonic Sounds Using Multipitch  
32 Analysis and Iterative Parameter Estimation, in *Proceedings of the IEEE Workshop on*  
33 *Applications of Signal Processing in Audio and Acoustics*, pp. 83-86.
- 34 [66] Hu, G. (2006). *Monaural Speech Organization and Segregation*, Ph.D. Thesis, The Ohio State  
35 University, USA.
- 36 [67] Li, Y., Woodruff, J. & Wang, D. L. (2009). Monaural Musical Sound Separation Based on  
37 Pitch and Common Amplitude Modulation, *IEEE Transactions on Audio, Speech and*  
38 *Language Processing*, vol. 17, no. 7, pp. 1361-1371.
- 39 [68] ISO. *Acoustics - Standard Tuning Frequency (Standard Musical Pitch)*, ISO 16:1975,  
40 International Organization for Standardization, Geneva, 1975.
- 41 [69] Nesbit, A., Jafari, M. G., Vincent, E., & Plumbley, M. D. (2010). Audio Source Separation  
42 Using Sparse Representations. In W. Wang (Ed), *Machine Audition: Principles, Algorithms*  
43 *and Systems*. Chapter 10, pp.246-264. IGI Global.
- 44 [70] Plumbley, M. D., Blumensath, T., Daudet, L., Gribonval, R., & Davies, M. E. (2010).  
45 Sparse Representations in Audio and Music: from Coding to Source Separation,  
46 *Proceedings of the IEEE*, vol. 98, pp. 995-1005.

- 1 [71] Dai, W., Xu, T. & Wang, W. (2012). Dictionary Learning and Update based on  
2 Simultaneous Codeword Optimisation (SIMCO), *Proceedings of the IEEE International*  
3 *Conference on Acoustics, Speech and Signal Processing*, Kyoto, Japan.
- 4 [72] Dai, W., Xu, T., & Wang, W. (2011). Simultaneous Codeword Optimisation (SimCO) for  
5 Dictionary Update and Learning, *arXiv:1109.5302*.
- 6 [73] Adler, A., Emiya V., Jafari, M.G., Elad, M., Gribonval, G., & Plumbley, M.D. (2012).  
7 Audio Inpainting, *IEEE Transactions on Audio, Speech and Language Processing*, vol. 20,  
8 pp. 922-932.
- 9 [74] Wang, W. (2011). *Machine Audition: Principles, Algorithms and Systems*, IGI Global Press.
- 10 [75] Jan, T. & Wang, W. (2011). Cocktail Party Problem: Source Separation Issues and  
11 Computational Methods, in W. Wang (ed), *Machine Audition: Principles, Algorithms and*  
12 *Systems*, IGI Global Press, pp. 61-79.
- 13 [76] Jan, T., Wang, W., & Wang, D.L. (2011). A Multistage Approach to Blind Separation of  
14 Convolutional Speech Mixtures. *Speech Communication*, vol. 53, pp. 524-539.
- 15 [77] Luo, Y., Wang, W., Chambers, J. A., Lambotharan, S., & Proudler, I. (2006). Exploitation  
16 of Source Non-stationarity for Underdetermined Blind Source Separation With  
17 Advanced Clustering Techniques, *IEEE Transactions on Signal Processing*, vol. 54, no. 6,  
18 pp. 2198-2212.
- 19 [78] Adiloglu, K. & Vincent, E. (2011). An Uncertainty Estimation Approach for the  
20 Extraction of Source Features in Multisource Recordings, in *Proceedings of the European*  
21 *Signal Processing Conference*, Barcelona, Spain.
- 22 [79] Adiloglu, K., & Vincent, E. (2012). A General Variational Bayesian Framework for  
23 Robust Feature Extraction in Multisource Recordings, in *Proceedings of the IEEE*  
24 *International Conference on Acoustics, Speech and Signal Processing*, Kyoto, Japan.
- 25 [80] Simon, L. S. R., & Vincent, E. (2012). A General Framework for Online Audio Source  
26 Separation, in *Proceedings of the International conference on Latent Variable Analysis and*  
27 *Signal Separation*, Tel-Aviv, Israel.
- 28 [81] Hsu, C.-L., Wang, D.L., Jang J.-S.R., & Hu, K. (2012). A Tandem Algorithm for Singing  
29 Pitch Extraction and Voice Separation from Music Accompaniment, *IEEE Transactions*  
30 *on Audio, Speech, and Language Processing*, vol. 20, pp. 1482-1491.