

A New Approach to Music/Voice Separation Using Resonance-Based Signal Decomposition

Roozbeh Soleymani, New York University, Polytechnic School of Engineering



Abstract

Separation of the voice from background music has been the subject of many recent research studies. As voice and music are mixed in both time and frequency domains, conventional filtering and signal separation methods are not applicable to this problem.

State of the art solutions to music/voice separation take advantage of complex techniques such as pitch detection, texture detection, repeating pattern, temporal and spectral continuity, trend estimation, sparsity and wavelet and voice recognition. However many of the techniques operate based on assumptions and simplifications that are not always valid, thus resulting in high separation error rates.

This poster presents a new approach to music/voice separation based on the use of tunable Q wavelet transform and sparsity enabled method for resonance-based signal decomposition. Resonance-based signal decomposition method performs well when applied to the separation of high resonant and low resonant components of the audio signals.

Stage 1

The fundamental difference between music and voice in this case is that music has a more consistent behaviour when it comes to high and low resonant components and it mostly contains either high or low resonant components rather than a mixture of two (Since there is no exact definition for "Low" or "High" in this case if we carefully tune the tunable Q wavelet transforms we can have most of the music energy as high resonant component most of the time).

As opposed to music, the human voice does not have such a consistency and most of the time contains a mixture of high and low resonant components at the same time. In this method, this feature has been used to recognize the sections which are pure music from the sections which are the mixture of voice and music by defining a cost function. So we can simply classify the segments with a high percentage of high resonant component as "Music Only" or "Pure Music" parts and the segments which have a considerable amount of low resonant content (compared to previous case) as "Mixture of Voice/music" segments.

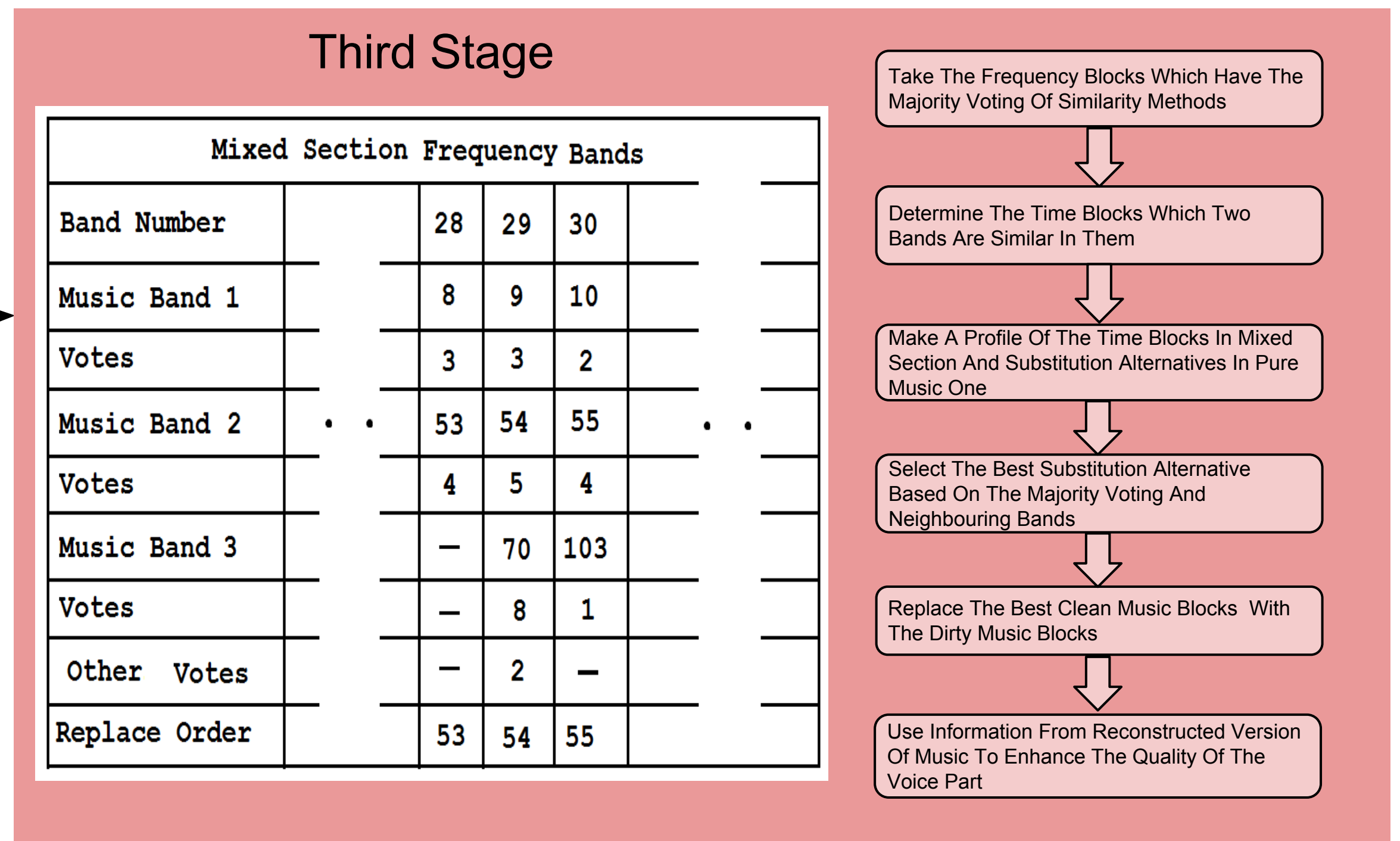
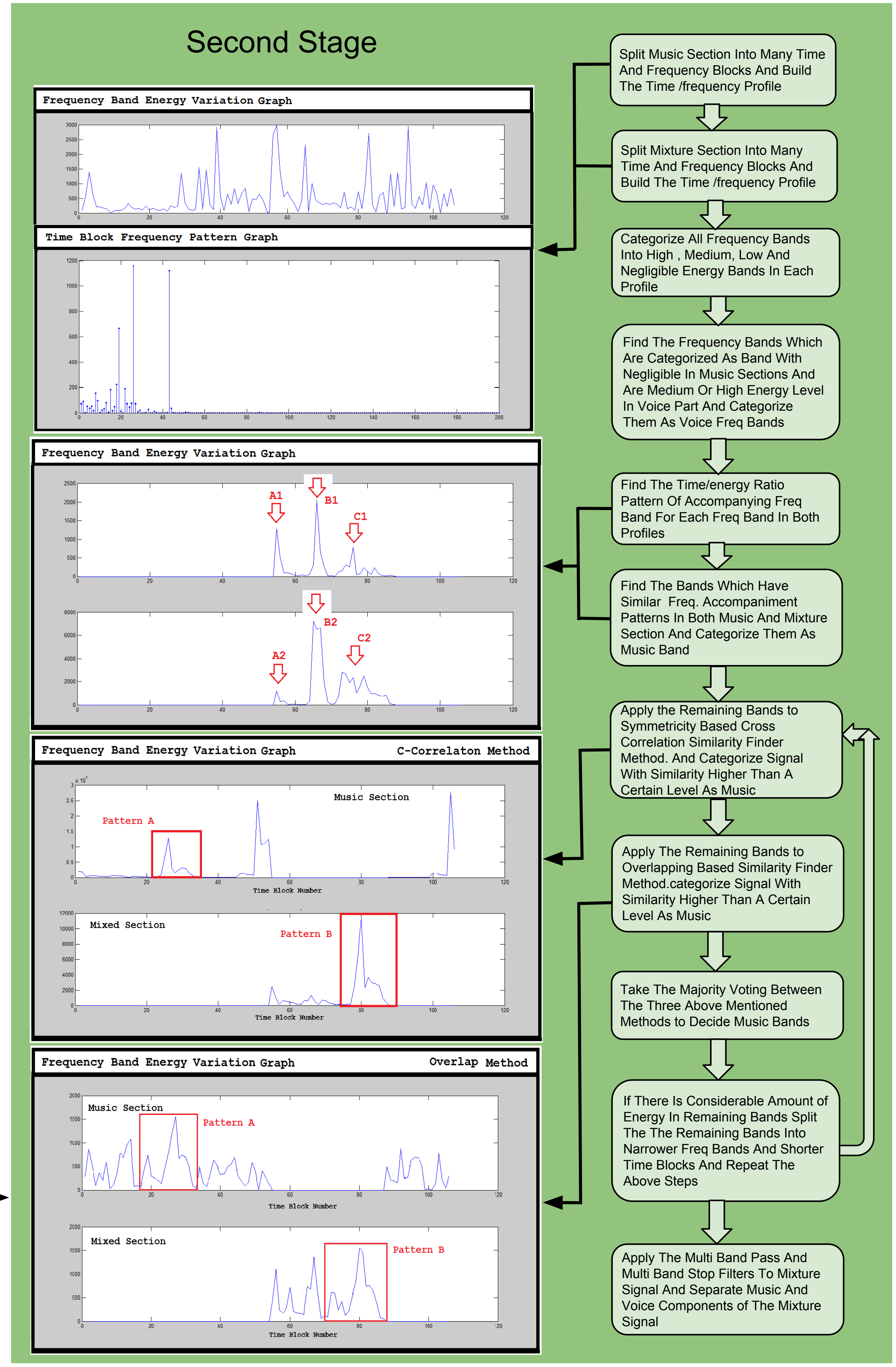
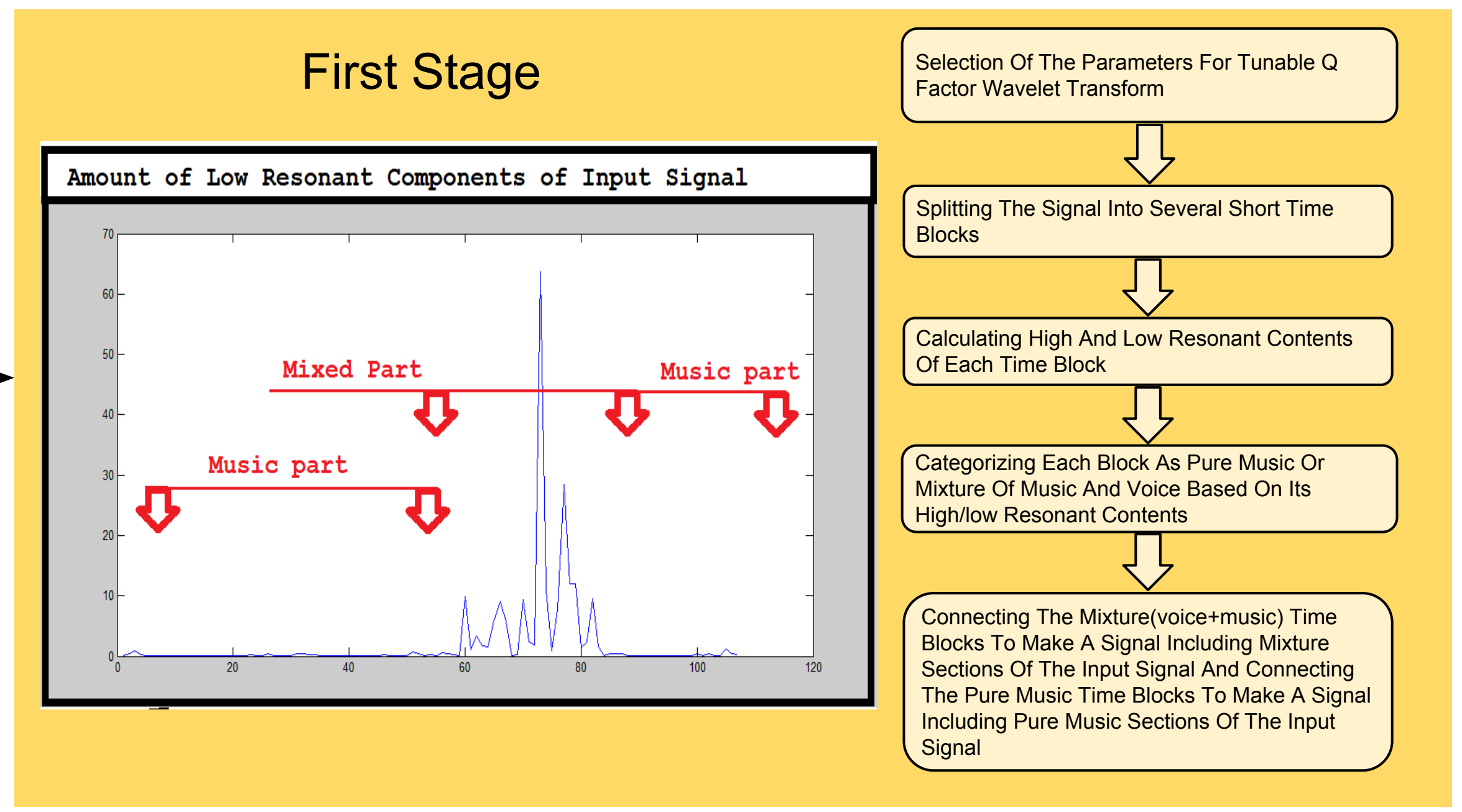
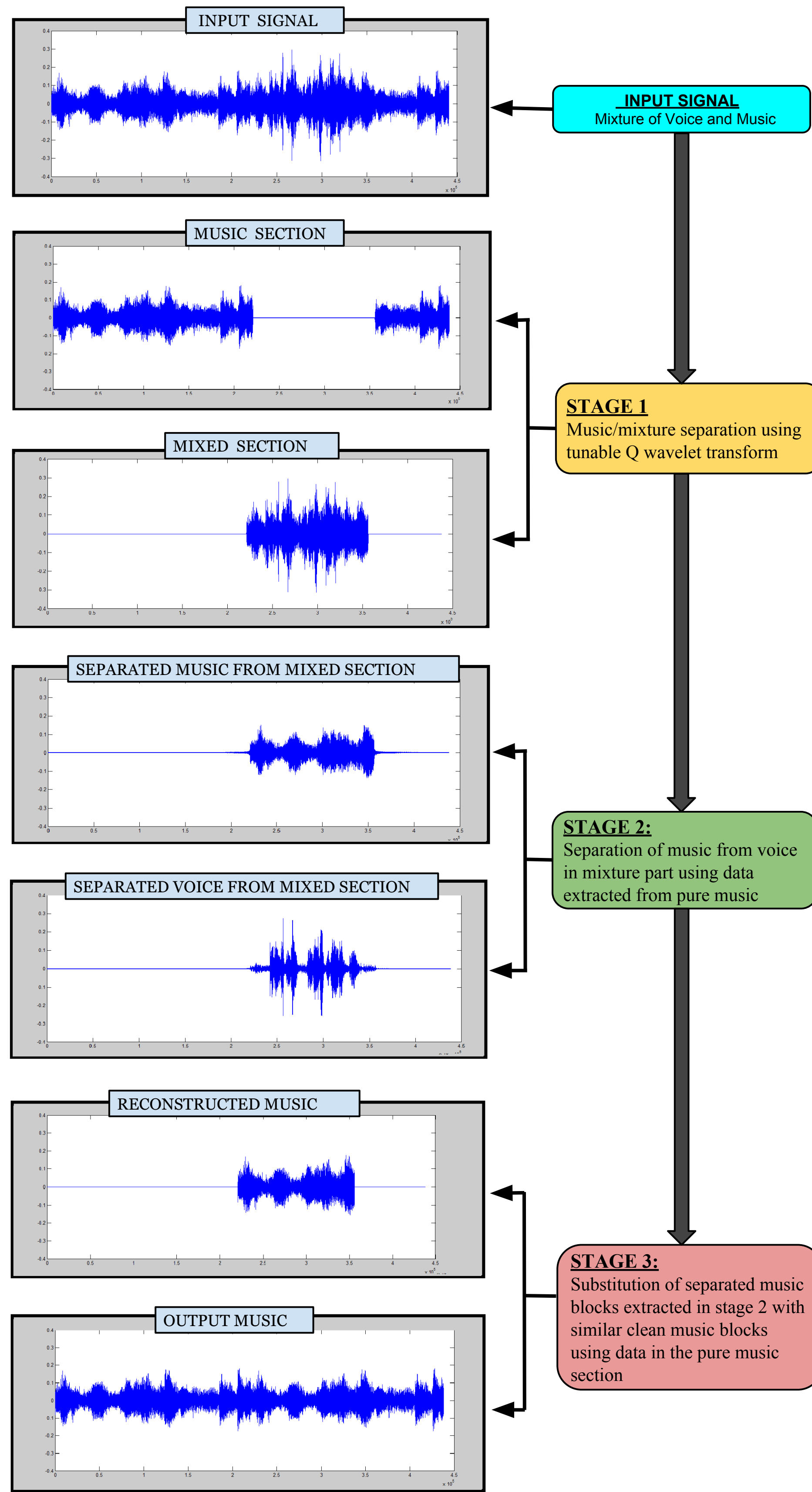
Stage 2

The second stage of this method attempts to restore the rest of the music which is mixed with singing voice. The key information which is very useful in this stage is the music features that are extracted from the segments that have been classified as pure music in the first stage. Both pure music and mixture part are divided into several narrow frequency bands in frequency domain and several short time blocks in time domain to make a profile which contains the information of the signal energy variation in each of these frequency bands over time also frequency arrangement in each time block. At this stage three independent classification methods are applied to each frequency subband to decide whether the subband in mixed section bear any resemblance to its counterpart in pure music section or not. If there is a considerable similarity we can assume that band as a music band as well. The three classification methods are based on frequency accompaniment pattern detection (An extended concept of harmonic detection), symmetry based cross correlation similarity finder, and overlap based similarity finder. The final decision is made based on majority voting from three classifiers.

Stage 3

The output of the second stage will separate music and voice signal in mixture part of the input. However both music and voice include low energy remnants of each other. Even though these components have a very low energy compared to the main signal, they still can be heard by human ear and therefore have a high negative impact on separation quality.

In order to remove these parasitic low energy components in third stage again using the information from the music segments (stage 1) and results of similarity finder methods (stage 2) the time blocks in the music signal which is separated from the voice (and contaminated with some voice remnants) will be substituted with a similar clean pure music time blocks obtained from pure music section. Here, it is assumed that every note which is played in the mixture section has been played at least once in the music section. Which seems to be a relatively reasonable assumption. (Normally singer does not start singing immediately after the music starts). This approach has been tested with several music/voice signals and whenever the input signal met assumption requirements the results demonstrated almost a perfect separation efficiency. However in cases which the input signal does not meet the above mentioned assumption the result is still very good.



References

- [1] Ivan W. Selesnick, Wavelet Transform With Tunable Q-Factor. IEEE Transactions on Signal Processing 59(8): 3560-3575 (2011).
- [2] Ivan W. Selesnick, Resonance-based signal decomposition: A new sparsity-enabled signal analysis method. Signal Processing 91(12): 2793-2809 (2011).
- [3] David Havelock, Sonoko Kuwano, Michael Vorländer. Handbook of Signal Processing in Acoustics. Springer. pp. 399-417. ISBN 978-0-387-77698-9.
- [4] Zafar Rafii, Bryan Pardo, Repeating Pattern Extraction Technique (REPET): A Simple Method for Music/Voice Separation. IEEE Trans. on Audio, Speech & Language Processing 21(1): 71-82 (2013).
- [5] Chao-Ling Hsu and Jyh-Shing Roger Jang. 2010. On the improvement of singing voice separation for monaural recordings using the MIR-1K dataset. Trans. Audio, Speech and Lang. Proc. 18, 2 (February 2010), 310-319.