

Better Music Demixing with the sliCQ Transform

Sevag Hanssian¹

¹Independent researcher (sevag.xyz), Montréal, Canada
sevagh@protonmail.com

Abstract

Music source separation, or music demixing, is the task of decomposing a song into its constituent sources, which are typically isolated instruments (e.g., drums, bass, and vocals). Open-Unmix (UMX), and the improved variant CrossNet-Open-Unmix (X-UMX), are high-performing models that use Short-Time Fourier Transform (STFT) as the representation of music signals, and apply masks to the magnitude STFT to separate mixed music into four sources: vocals, drums, bass, and other.

The time-frequency uncertainty principle states that the STFT of a signal cannot be maximally precise in both time and frequency. The tradeoff in time-frequency resolution can significantly affect music demixing results. For the Cadenza Challenge in 2023, we submitted a model, xumx-sliCQ-V2,¹ which replaces the STFT with the sliCQT, a time-frequency transform with varying time-frequency resolution. Our system achieved an SDR score of 4.4 dB on the MUSDB18-HQ test set.

Index Terms: music source separation, music demixing, deep neural networks, time-frequency resolution, MUSDB18-HQ

1. Introduction

The STFT is computed by applying the Discrete Fourier Transform on fixed-size windows of the input signal. From both auditory and musical motivations, variable-size windows are preferred, with long windows in low-frequency regions to capture detailed harmonic information with a high frequency resolution, and short windows in high-frequency regions to capture transients with a high time resolution [1]. The sliCQ Transform (sliCQT) [2] is a time-frequency transform with complex Fourier coefficients and perfect inverse that uses varying windows to achieve nonlinear time or frequency resolution. An example application of the sliCQT is an invertible Constant-Q Transform (CQT) [3].

2. Methodology

From the guidelines of the Cadenza challenge and to ensure reproducibility, we only relied on the standard and widely-available MUSDB18-HQ dataset [4] for training and evaluation of xumx-sliCQ-V2.

In xumx-sliCQ-V2, we kept the same sliCQT parameters from the older variant, xumx-sliCQ [5]. The sliCQT parameters are 262 bins on the Bark scale between 32.9–22050 Hz, chosen in a random parameter

search to maximize the mix-phase or noisy-phase oracle [5]. STFT and sliCQT spectrograms of a glockenspiel signal are shown in Figure 1.

The STFT outputs a single time-frequency matrix where all of the frequency bins are spaced uniformly apart and have the same time resolution. The sliCQT groups frequency bins, which may be nonuniformly spaced, in a ragged list of time-frequency matrices, where each matrix contains frequency bins that share the same time resolution. In xumx-sliCQ-V2, convolutional layers were applied separately to each time-frequency matrix, shown in Figure 2.

We made three significant changes to the older system, xumx-sliCQ, which account for the improved performance of xumx-sliCQ-V2.

2.1. Improved overlap-add

The sliCQT subdivides the input signal into “slices” of length N that are “symmetrically zero-padded to length $2N$ ” [2, 10]. To create a spectrogram, adjacent slices need to be 50% overlap-added with each other, with no inverse operation. In xumx-sliCQ-V2, we incorporated the slice size in the kernel and stride of the first convolution layer and last transpose convolution layer to avoid the non-invertible overlap-add procedure, shown in Figure 3.

2.2. Differentiable Wiener filtering and complex MSE

In xumx-sliCQ, the mean-squared error (MSE) loss function is applied to the magnitude spectrogram estimates of the neural network. The post-processing Wiener filtering step is then used to further improve the separation results and create the estimated complex spectrograms. Danna-Sep [6], a winning system from MDX 21, incorporated the Wiener filtering step into the neural network to output complex spectrograms, and used the complex mean-squared error as the loss function. We did the same in xumx-sliCQ-V2.

2.3. Mask-sum loss

The final activation layer of xumx-sliCQ-V2 is the sigmoid function ($\in [0.0, 1.0]$), to apply as a soft mask (or ratio mask) to the magnitude spectrogram of the mix. An underlying simplifying assumption in music demixing is that the mix is a linear sum of the sources. Therefore, the sum of the four target masks should be exactly equal to one, shown in Equation (1). In xumx-sliCQ-V2, we introduce an additional loss term called the “mask-sum loss” which computes the MSE of the sum of the esti-

¹<https://github.com/sevag/xumx-sliCQ/tree/v2>

