

Music demixing with the sliCQ transform

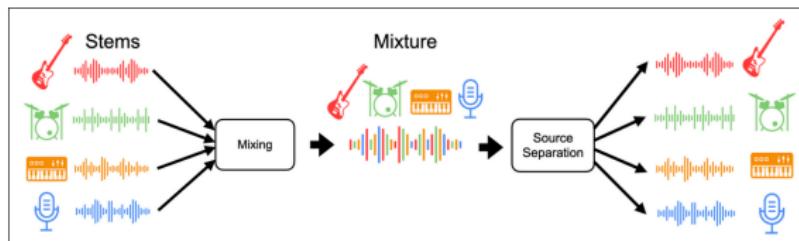
Sevag Hanssian

December 15, 2021

Music source separation, demixing, and unmixing

Music source separation: extract an estimate of an isolated source (or target) from mixed musical audio (e.g., harmonic/percussive, vocals, drums, bass, piano)

Music demixing (or unmixing): estimate multiple sources (vocals, drums, bass, other¹) that can be summed back to the original mix. Multiple MSS subproblems, reversing the linear mixing of stems in the recording studio (stem datasets can be used for mixing and demixing)

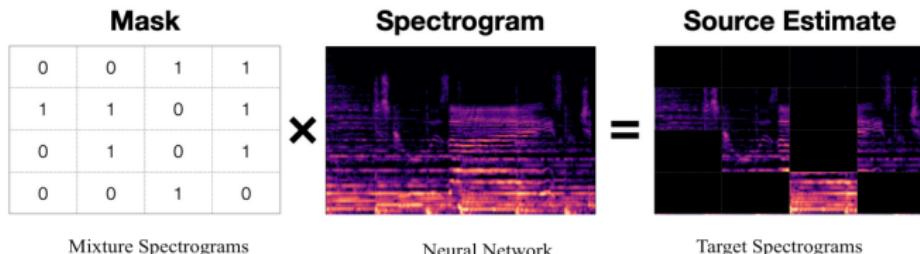
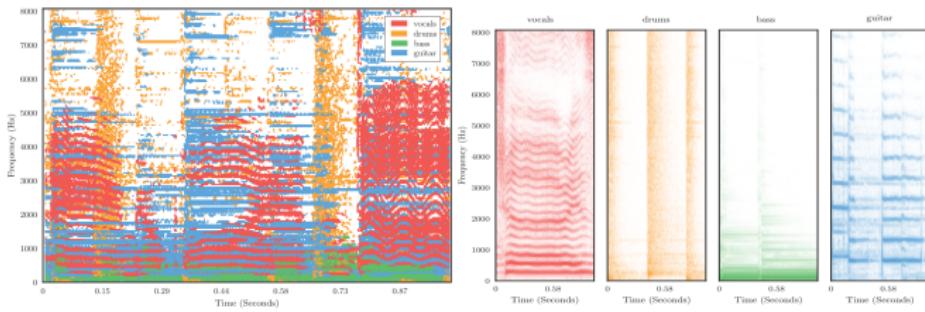


Popular approach: musical source models, which are “model-based approaches that attempt to capture the spectral characteristics of the target source”² with **time-frequency masks**

¹ Zafar Rafii et al. 2019.

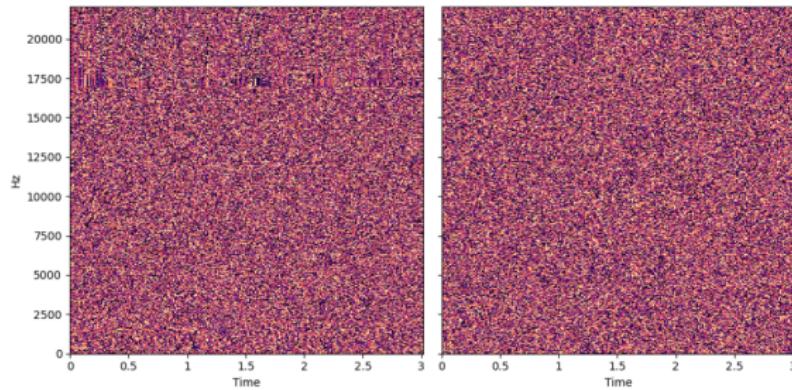
² Estefanía Cano et al. (2018), p. 36.

Music source separation with DNN and spectrograms



Phase performance ceiling

- ① Simplifying assumption: estimate magnitude spectrograms, use the phase of the original mixed audio. Called “noisy phase”³. Done by Open-Unmix (UMX), CrossNet-Open-Unmix (X-UMX)⁴, and many other popular & near-SOTA models
- ② Why? Phase is hard to model!⁵



³ Gordon Wichern et al. *arXiv preprint arXiv:1907.01160* (2019).

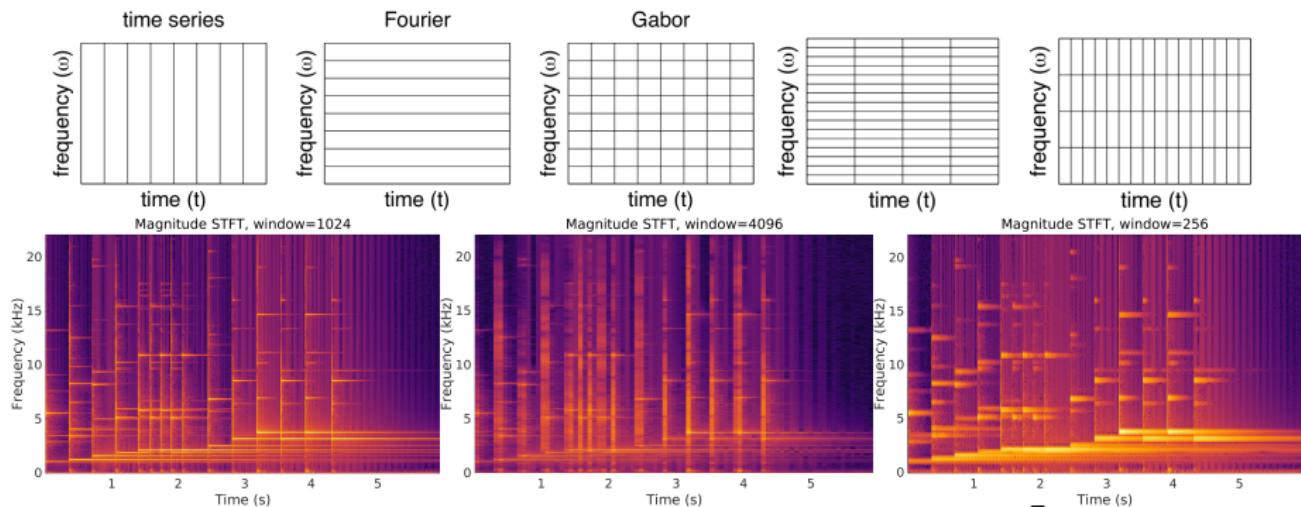
⁴ Fabian-Robert Stöter et al. *Journal of Open Source Software* (2019);

Ryosuke Sawata et al. *arXiv preprint arXiv:2010.04228* (2021).

⁵<https://source-separation.github.io/tutorial/basics/phase.html#why-we-don-t-model-phase>

Time-frequency tradeoff in the STFT and MDX

Joint time-frequency analysis is important for signals whose frequencies change with time.⁶ Take the Fourier transform of local windows of the signal, i.e., the STFT. Change window size to trade off time and frequency:



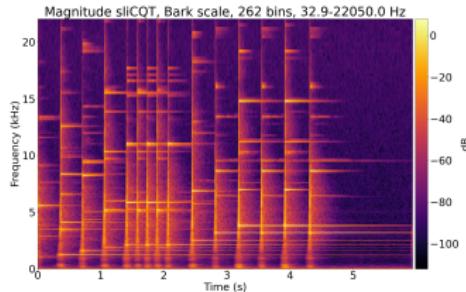
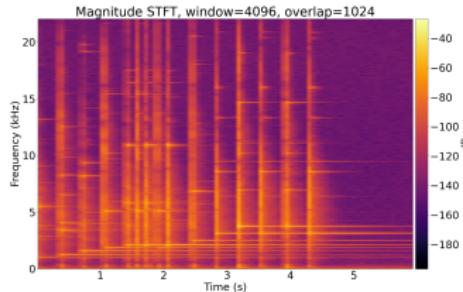
In music source separation, window size matters per-target.⁷ Short-window for percussion, long-window for harmonic

⁶ Dennis Gabor. *Journal of Institution of Electrical Engineers* (1946).

⁷ Andrew Simpson. *arXiv preprint arXiv:1504.07372* (2015).

CQT, NSGT, sliCQT

- ① For musical and auditory reasons, we want high frequency resolution at low frequencies and high time resolution at high frequencies⁸
- ② CQT⁹ uses long windows in low frequencies and short windows in high frequencies for the 12-tone Western pitch scale
- ③ Nonstationary Gabor Transform (NSGT) and sliCQT¹⁰ are TF transforms with Fourier coefficients, perfect inverse, and varying windows to create a varying time-frequency resolution
- ④ sliCQT params chosen for max quality of the noisy-phase waveform



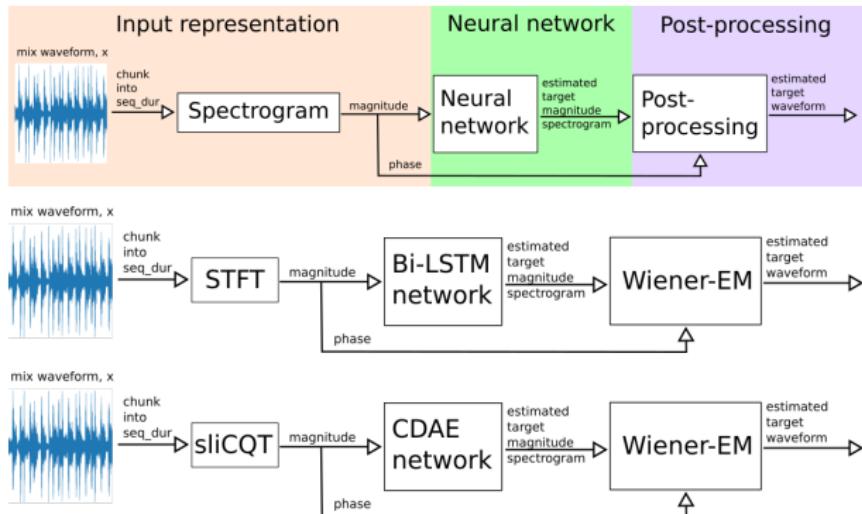
⁸ Monika Dörfler. PhD thesis. 2002.

⁹ Judith Brown. (1991).

¹⁰ Peter Balazs et al. (2011); Nicki Holighaus et al. (2013).

xumx-sliCQ

- ① My goal: improve Open-Unmix by replacing STFT with sliCQT
- ② My model submitted to the MDX21 challenge and workshop:
<https://github.com/sevagh/xumx-sliCQ>
- ③ Use Convolutional Denoising Autoencoder¹¹ neural architecture
- ④ Scored 3.6 dB vs. 4.6 dB (UMX) and 5.54 dB (X-UMX); there is still room for improvement



¹¹ Emad M. Grais et al. 2017; Emad M. Grais et al. 2021.

MDX 21 winners and current trends

- ① Previously, music demixing systems were submitted to and evaluated at SiSEC (Signal Separation Evaluation Campaign). This year: MDX (Music Demixing Challenge) ISMIR 2021 @ AICrowd, follow-up MDX21 workshop, satellite @ ISMIR 2021
- ② **ISMIR 2021:** Model that uses the complex spectrogram (i.e. includes phase) and uses complex masks¹²
- ③ **MDX21:** 1: Demucs++¹³ (waveforms + complex spectrogram), 2: KUIELAB-MDX-Net¹⁴ (waveforms + magnitude spectrogram), 3: Danna-Sep¹⁵ (waveform + magnitude spectrogram, use complex spectrogram in loss function)

Properties in common: blending networks, waveforms (implicitly includes phase), complex spectrograms/masks, mixing spectrogram and waveform models

¹² Qiuqiang Kong et al. 2021.

¹³ Alexandre Défossez. *arXiv preprint arXiv:2111.03600* (2021).

¹⁴ Minseok Kim et al. *arXiv preprint arXiv:2111.12203* (2021).

¹⁵ Chin-Yun Yu et al. *arXiv preprint arXiv:2112.03752* (2021).