

Time-Frequency Representations for Music Source Separation

Final project presentation

Sevag Hanssian

MUMT 622, Winter 2021

April 22, 2021

TF representations for music separation

2021-04-16

Time-Frequency Representations for Music Source Separation
Final project presentation

Sevag Hanssian
MUMT 622, Winter 2021
April 22, 2021

Music Source Separation

Musical sources are often categorized as either predominantly harmonic, predominantly percussive, or as singing voice.¹ In this project, I consider both cases (HPSS and harmonic/percussive/vocal)

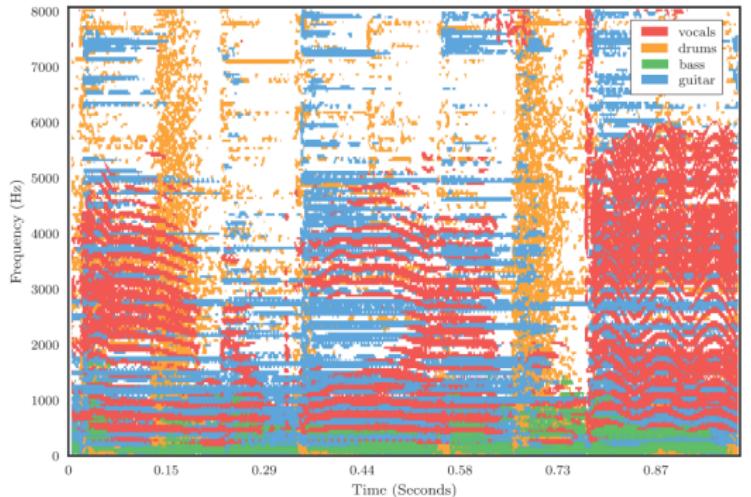


Figure: Typical music sources in a mix

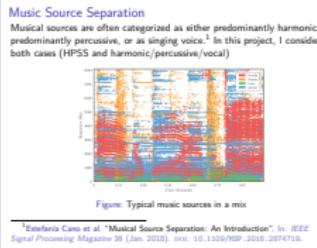
¹ Estefanía Cano et al. "Musical Source Separation: An Introduction". In: *IEEE Signal Processing Magazine* 36 (Jan. 2018). DOI: 10.1109/MSP.2018.2874719.

TF representations for music separation

2021-04-16

└ Music Source Separation

- steady-state/transient
- tonal/transient in Itfat world



Music Source Separation

A notable property of musical sources is that they are typically *sparse* in the sense that for the majority of points in time and frequency, the sources have very little energy present²

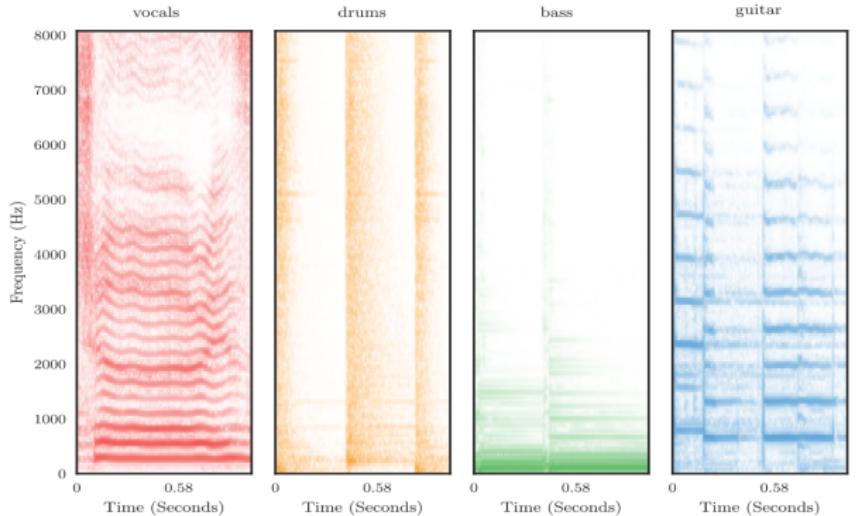


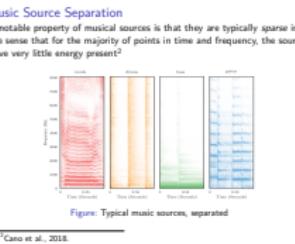
Figure: Typical music sources, separated

²Cano et al., 2018.

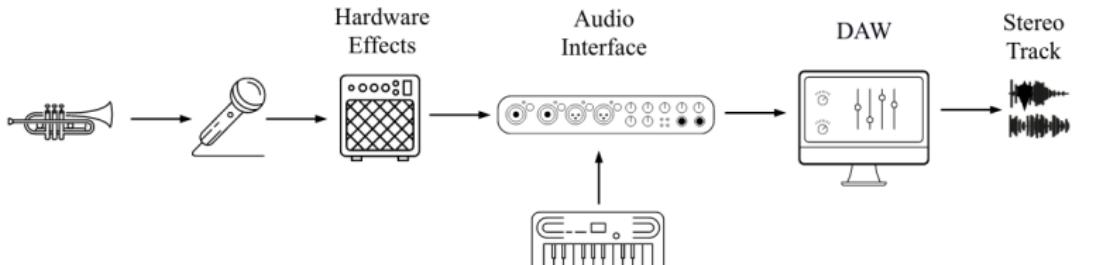
TF representations for music separation

2021-04-16

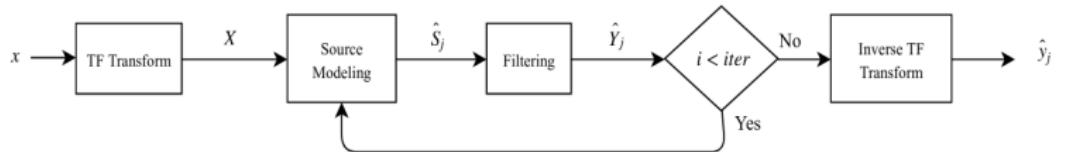
└ Music Source Separation



Music Source Separation



(a) Mixing block diagram



(b) Unmixing block diagram

Figure: Typical block diagrams for source mixing and source separation (aka “unmixing”)³

³Cano et al., 2018.

TF representations for music separation

2021-04-16

Music Source Separation

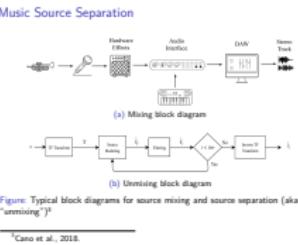
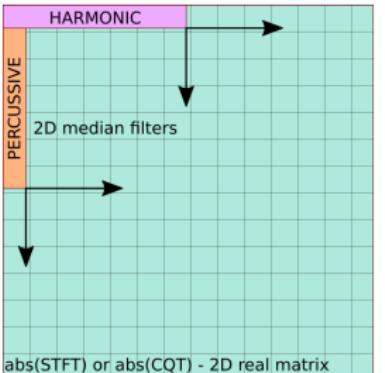


Figure: Typical block diagrams for source mixing and source separation (aka “unmixing”)³

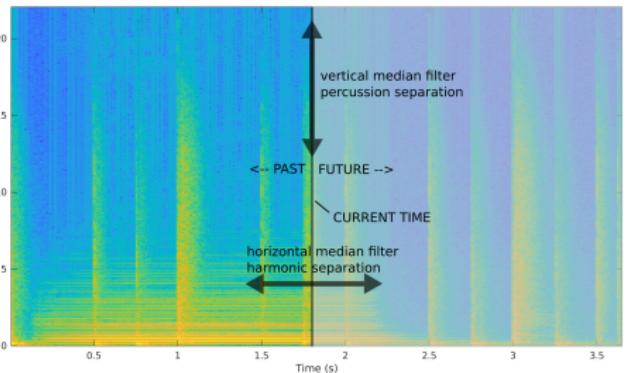
³Cano et al., 2018.

Median filtering HPSS

Form of Kernel Additive Model – describe harmonic sounds as horizontal, percussive sounds as vertical, and apply median filters on magnitude spectrogram to estimate each⁴



(a) Anticausal/offline



(b) Causal/realtimse

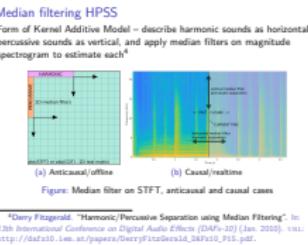
Figure: Median filter on STFT, anticausal and causal cases

⁴Derry Fitzgerald. "Harmonic/Percussive Separation using Median Filtering". In: *13th International Conference on Digital Audio Effects (DAFx-10)* (Jan. 2010). URL: http://dafx10.iem.at/papers/DerryFitzGerald_DAFx10_P15.pdf.

TF representations for music separation

2021-04-16

Median filtering HPSS



- sliding STFT, perform causal median filtering, then invert
- keep several columns of STFT in memory, perform windowing + overlap-add – presented this in 501

Median filtering HPSS

Use harmonic and percussive magnitude spectrogram estimates to compute soft masks⁵ or hard masks:⁶

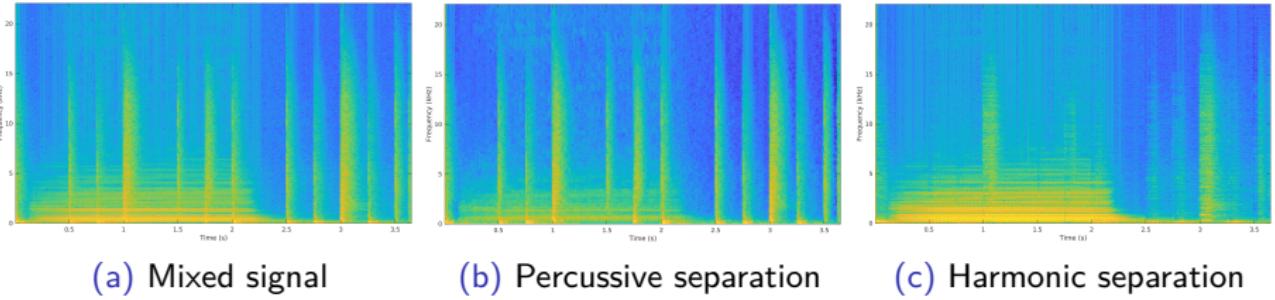


Figure: Example of median filtering HPSS

Originally based on STFT spectrogram. CQT works fine too.

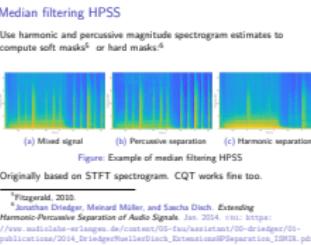
⁵Fitzgerald, 2010.

⁶Jonathan Driedger, Meinard Müller, and Sascha Disch. *Extending Harmonic-Percussive Separation of Audio Signals*. Jan. 2014. URL: https://www.audiolabs-erlangen.de/content/05-fau/assistant/00-driedger/01-publications/2014_DriedgerMuellerDisch_ExtensionsHPSeparation_ISMIR.pdf.

TF representations for music separation

2021-04-16

└ Median filtering HPSS



Median filtering HPSS

Soft mask, Wiener filter:⁷

$$M_{\text{harmonic}} = \frac{|\hat{S}_{\text{harmonic}}|^2}{|\hat{S}_{\text{harmonic}}|^2 + |\hat{S}_{\text{percussive}}|^2}$$

Hard mask, binary mask:⁸

$$M_{\text{harmonic}} = \frac{|\hat{S}_{\text{percussive}}|}{|\hat{S}_{\text{harmonic}}| + \epsilon} > \beta$$

Setting $\beta > 1.0$ leads to a third residual component:

$$M_{\text{residual}} = 1 - (M_{\text{harmonic}} + M_{\text{percussive}})$$

Soft mask gives higher audio quality⁹

⁷Fitzgerald, 2010.

⁸Driedger, Müller, and Disch, 2014.

⁹Gerkmann and Vincent, 2018.

TF representations for music separation

2021-04-16

└ Median filtering HPSS

- define how different harmonic and percussive should be

Median filtering HPSS	
Soft mask, Wiener filter: ⁷	
$M_{\text{harmonic}} = \frac{ \hat{S}_{\text{harmonic}} ^2}{ \hat{S}_{\text{harmonic}} ^2 + \hat{S}_{\text{percussive}} ^2}$	$M_{\text{harmonic}} = \frac{ \hat{S}_{\text{harmonic}} ^2}{ \hat{S}_{\text{harmonic}} ^2 + \hat{S}_{\text{percussive}} ^2}$
Hard mask, binary mask: ⁸	$M_{\text{harmonic}} = \frac{ \hat{S}_{\text{percussive}} }{ \hat{S}_{\text{harmonic}} + \epsilon} > \beta$
Setting $\beta > 1.0$ leads to a third residual component:	
$M_{\text{residual}} = 1 - (M_{\text{harmonic}} + M_{\text{percussive}})$	
Soft mask gives higher audio quality ⁹	
<small>⁷Fitzgerald, 2010. ⁸Driedger, Müller, and Disch, 2014. ⁹Gerkmann and Vincent, 2018.</small>	

Iterative median filtering HPSS

2-pass HPSS,¹⁰ , harmonic/percussive/vocal:¹¹

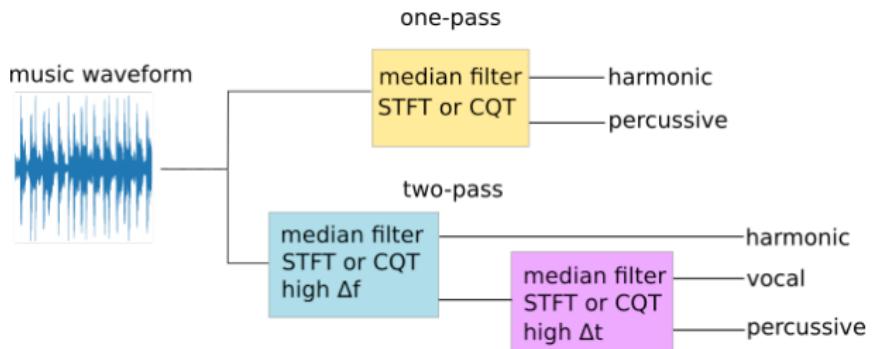


Figure: One- or two-pass median filter HPSS

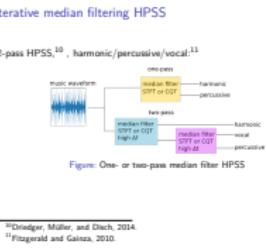
¹⁰Driedger, Müller, and Disch, 2014.

¹¹Fitzgerald and Gainza, 2010.

TF representations for music separation

2021-04-16

Iterative median filtering HPSS

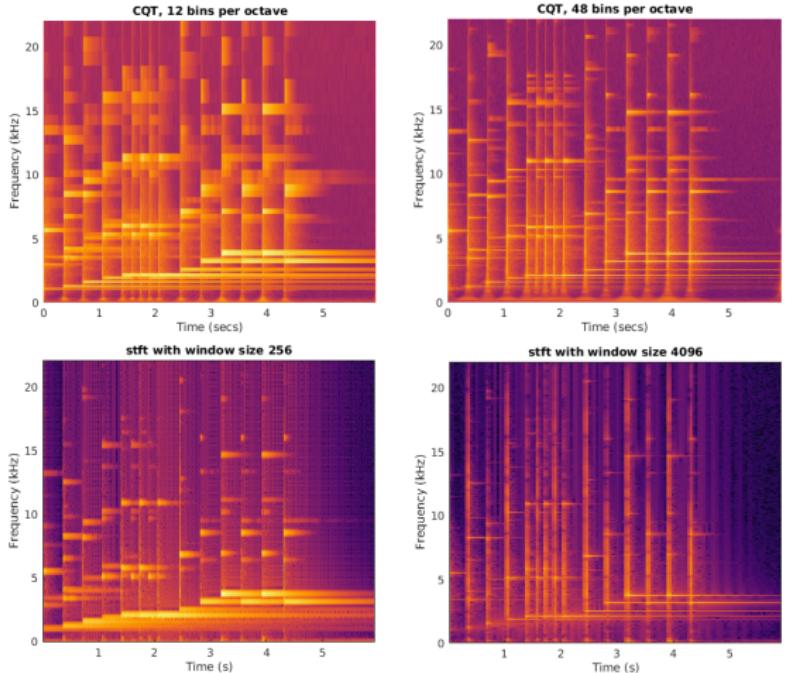


¹⁰Driedger, Müller, and Disch, 2014.
¹¹Fitzgerald and Gainza, 2010.

Figure: One- or two-pass median filter HPSS

STFT, CQT, and TF resolution

STFT vs. CQT¹² (based on NSGT¹³):



¹²<https://www.mathworks.com/help/wavelet/ref/cqt.html>

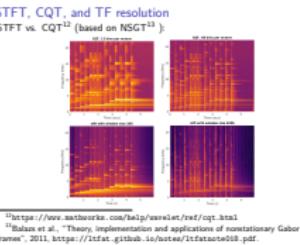
¹³Balazs et al., "Theory, implementation and applications of nonstationary Gabor frames", 2011, <https://ltfat.github.io/notes/ltfatnote018.pdf>.

TF representations for music separation

2021-04-16

└ STFT, CQT, and TF resolution

- glockenspiel = default file of ltfat, the canonical sound for tonal/transient
- loaded by gspi function
- look how good the cqt is – original motivation for this project, can we take the simple median filter algorithm and make it better with a better tf representation



¹²<https://www.mathworks.com/help/wavelet/ref/cqt.html>
¹³Balazs et al., "Theory, implementation and applications of nonstationary Gabor frames", 2011, <https://ltfat.github.io/notes/ltfatnote018.pdf>

Sparsity, entropy, and two window sizes

Sparsity and entropy are complementary concepts:¹⁴

- *Sparsity* is the property of concentrating most of the energy of x in few coefficients of w
- *Entropy* is the property of not concentrating most of the probability mass in few atoms of p – in other words, the entropy of a random variable is a concept of information theory that characterizes the unpredictability inherent in its outcomes

In every algorithm shown, the common element is a 2 dictionary wide + narrow window analysis, to represent tonal and transient parts of the input signal sparsely – or, to represent tonal and transient parts of the input signal with low entropy/high significance

¹⁴ Paul Honeine. *Entropy of Overcomplete Kernel Dictionaries*. 2014. arXiv: 1411.0161 [cs.IT]. URL: <https://arxiv.org/pdf/1411.0161.pdf>; Giancarlo Pastor et al. *Mathematics of Sparsity and Entropy: Axioms, Core Functions and Sparse Recovery*. 2015. arXiv: 1501.05126 [cs.IT].

TF representations for music separation

2021-04-16

└ Sparsity, entropy, and two window sizes

Sparsity, entropy, and two window sizes

Sparsity and entropy are complementary concepts:¹⁴

- Sparsity is the property of concentrating most of the energy of x in few coefficients of w
- Entropy is the property of not concentrating most of the probability mass in few atoms of p – in other words, the entropy of a random variable is a concept of information theory that characterizes the unpredictability inherent in its outcomes

In every algorithm shown, the common element is a 2 dictionary wide + narrow window analysis, to represent tonal and transient parts of the input signal sparsely – or, to represent tonal and transient parts of the input signal with low entropy/high significance

¹⁴ Paul Honeine. *Entropy of Overcomplete Kernel Dictionaries*. 2014. arXiv: 1411.0161 [cs.IT]. URL: <https://arxiv.org/pdf/1411.0161.pdf>; Giancarlo Pastor et al. *Mathematics of Sparsity and Entropy: Axioms, Core Functions and Sparse Recovery*. 2015. arXiv: 1501.05126 [cs.IT].

Structured sparsity

Group-LASSO:¹⁵ Lasso shrinkage (aka linear least squares regression^{16 17}) to the transform coefficients in the time and frequency dimensions

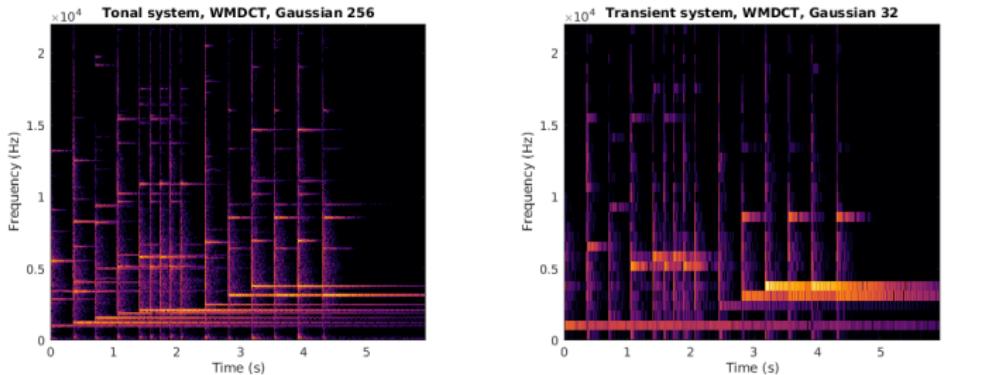


Figure: WMDCT frames for structured sparsity tonal/transient separation

¹⁵ Matthieu Kowalski and Bruno Torrésani. "Sparsity and persistence: Mixed norms provide simple signal models with dependent coefficients". In: *Signal Image and Video Processing* 3 (Sept. 2009). DOI: [10.1007/s11760-008-0076-1](https://doi.org/10.1007/s11760-008-0076-1).

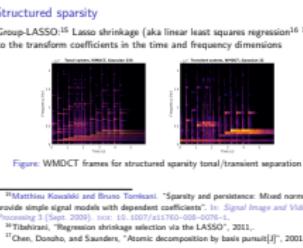
¹⁶ Tibshirani, "Regression shrinkage selection via the LASSO", 2011,.

¹⁷ Chen, Donoho, and Saunders, "Atomic decomposition by basis pursuit[J]", 2001,.

TF representations for music separation

2021-04-16

└ Structured sparsity



Structured sparsity

WMDCT¹⁸ + Group-LASSO – “audioshrink”¹⁹

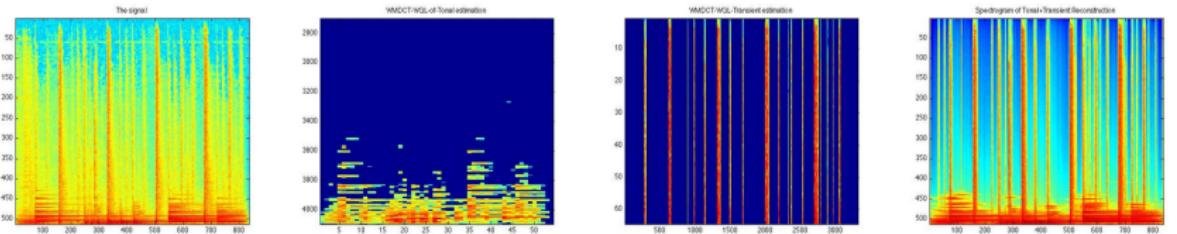


Figure: Audioshrink for tonal/transient separation in jazz music

Use 2 WMDCT transforms (wide + narrow window) + Group-LASSO to shrink input signal into significant coefficients in “time” and “frequency” groups

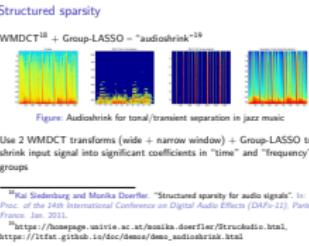
¹⁸Kai Siedenburg and Monika Doerfler. “Structured sparsity for audio signals”. In: *Proc. of the 14th International Conference on Digital Audio Effects (DAFx-11)*, Paris, France. Jan. 2011.

¹⁹<https://homepage.univie.ac.at/monika.doerfler/StrucAudio.html>, https://ltfat.github.io/doc/demos/demo_audioshrink.html

TF representations for music separation

2021-04-16

└ Structured sparsity



- Using Lasso-like methods (ℓ_1, ℓ_2 norm in linear least squares regression) assumes that the underlying systems are linear Gaussian, or approximately so²⁰
- If the system is not linear-Gaussian, linear least squares regression can lead to suboptimal results
- Rényi entropy, which is an adaptation of Shannon entropy²¹ is a measure of information in signal processing, and can be used as an optimization target (i.e., loss function) without imposing restrictions on the system being optimized

²⁰Ed Beadle et al. "An Overview of Renyi Entropy and Some Potential Applications". In: Nov. 2008, pp. 1698–1704. doi: [10.1109/ACSSC.2008.5074715](https://doi.org/10.1109/ACSSC.2008.5074715).

²¹A. Rényi. "On Measures of Entropy and Information". In: *Proceedings IV Berkeley Symposium on Mathematical Statistics and Probability, Berkeley, 20-30 June 1961* 1 (1961), pp. 547–561; C. E. Shannon. "A mathematical theory of communication". In: *The Bell System Technical Journal* 27.3 (1948), pp. 379–423. doi: [10.1002/j.1538-7305.1948.tb01338.x](https://doi.org/10.1002/j.1538-7305.1948.tb01338.x).

TF representations for music separation

2021-04-16

└ Rényi entropy vs. Lasso

Rényi entropy vs. Lasso

- Using Lasso-like methods (ℓ_1, ℓ_2 norm in linear least squares regression) assumes that the underlying systems are linear Gaussian, or approximately so²⁰
- If the system is not linear-Gaussian, linear least squares regression can lead to suboptimal results
- Rényi entropy, which is an adaptation of Shannon entropy²¹ is a measure of information in signal processing, and can be used as an optimization target (i.e., loss function) without imposing restrictions on the system being optimized

²⁰Ed Beadle et al. "An Overview of Renyi Entropy and Some Potential Applications". In: Nov. 2008, pp. 1698–1704. doi: [10.1109/ACSSC.2008.5074715](https://doi.org/10.1109/ACSSC.2008.5074715).

²¹A. Rényi. "On Measures of Entropy and Information". In: *Proceedings IV Berkeley Symposium on Mathematical Statistics and Probability, Berkeley, 20-30 June 1961* 1 (1961), pp. 547–561; C. E. Shannon. "A mathematical theory of communication". In: *The Bell System Technical Journal* 27.3 (1948), pp. 379–423. doi: [10.1002/j.1538-7305.1948.tb01338.x](https://doi.org/10.1002/j.1538-7305.1948.tb01338.x).

Time-Frequency Jigsaw Puzzle

- ① Create time-frequency “super-tiles” by superimposing a large window + small window Gabor analysis
- ② Use Rényi entropy to set coefficients to zero where sound has more entropy than random white noise

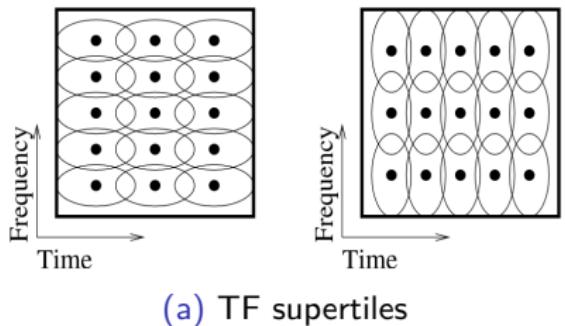


Figure: TF Jigsaw Puzzle tonal/transient separation²²

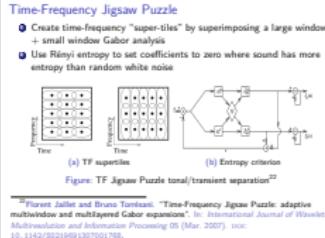
²²Florent Jaillet and Bruno Torrésani. "Time-Frequency Jigsaw Puzzle: adaptive multiwindow and multilayered Gabor expansions". In: *International Journal of Wavelets, Multiresolution and Information Processing* 05 (Mar. 2007). DOI: 10.1142/S0219691307001768.

TF representations for music separation

2021-04-16

Time-Frequency Jigsaw Puzzle

- i.e. good tonal/good transient
- high entropy = indicating sound is poorly represented



TFJigsaw

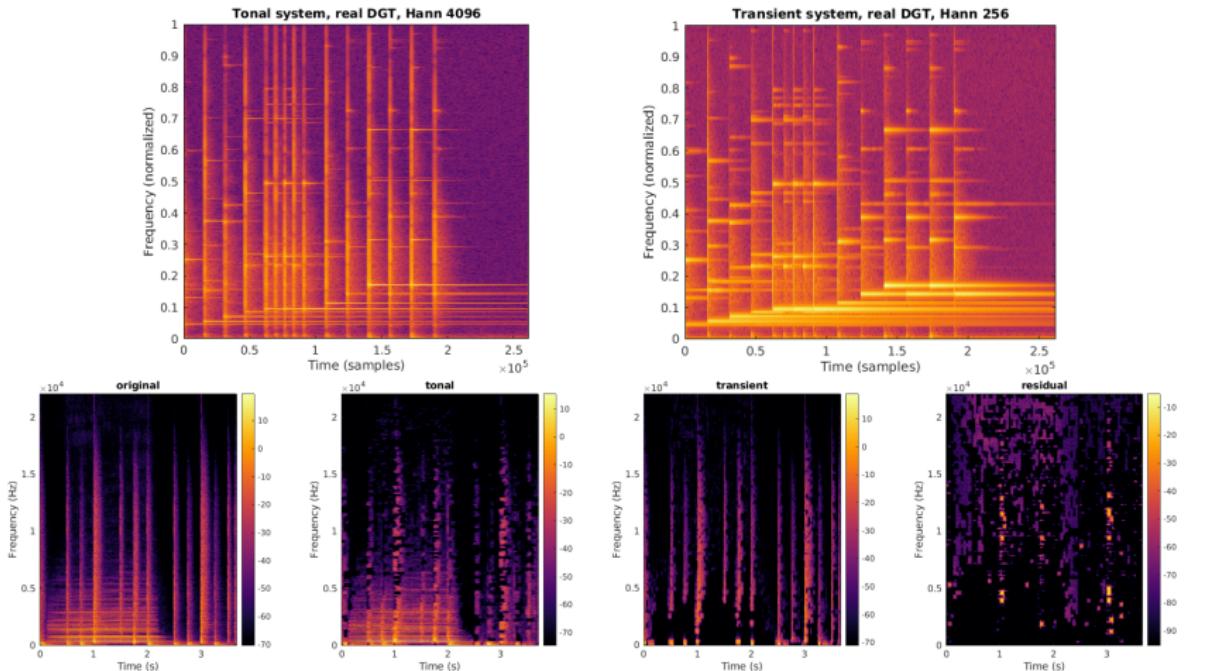


Figure: TF Jigsaw Puzzle tonal/transient separation²³

²³<https://ltfat.github.io/doc/sigproc/tfjigsawsep.html>,
https://github.com/ltfat/ltfat/blob/master/demos/demo_tfjigsawsep.m

TF representations for music separation

2021-04-16

└ TFJigsaw

- Hann DGT = stft, basically

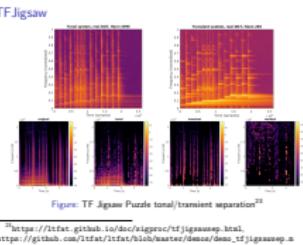


Figure: TF Jigsaw Puzzle tonal/transient separation²³
<https://ltfat.github.io/doc/sigproc/tfjigsawsep.html>,
https://github.com/ltfat/ltfat/blob/master/demos/demo_tfjigsawsep.m

Evaluation testbench

Inspired by SigSep²⁴, SISEC (Signal Separation Evaluation Campaign):

- BSS²⁵ (BSSv4 variant²⁶) and PEASS²⁷ (MATLAB toolkit²⁸). BSS vs. PEASS?²⁹ BSSv4 is used widely in modern literature,³⁰ but perceptual measures are important! Use PEASS
- Testing files: MUSDB18-HQ³¹
- MATLAB/Python testbench using file system + JSON interchanges
- Open-Unmix³² as a reference, open, near-SOTA neural solution
- Compare different configurations of each algorithm in “group stages,” winners move to next stage and may be combined in hybrid algorithms

²⁴<https://sigsep.github.io/>

²⁵Vincent, Gribonval, and Févotte, 2006.

²⁶<https://github.com/sigsep/bsseval>

²⁷Emiya et al., 2011.

²⁸<http://bass-db.gforge.inria.fr/peass/>

²⁹Ward et al., 2018.

³⁰Stöter, Liutkus, and Ito, 2018.

³¹Rafii et al., 2019.

³²Stöter et al., 2019.

TF representations for music separation

2021-04-16

└ Evaluation testbench

- global score omitted. target, interference, artifact

Evaluation testbench
Inspired by SigSep²⁴, SISEC (Signal Separation Evaluation Campaign)
• BSS²⁵ (BSSv4 variant²⁶) and PEASS²⁷ (MATLAB toolkit²⁸). BSS vs. PEASS?²⁹ BSSv4 is used widely in modern literature,³⁰ but perceptual measures are important! Use PEASS
• Testing files: MUSDB18-HQ³¹
• MATLAB/Python testbench using file system + JSON interchanges
• Open-Unmix³² as a reference, open, near-SOTA neural solution
• Compare different configurations of each algorithm in “group stages,” winners move to next stage and may be combined in hybrid algorithms
²⁴<https://sigsep.github.io/>
²⁵Vincent, Gribonval, and Févotte, 2006.
²⁶<https://github.com/sigsep/bsseval>
²⁷Emiya et al., 2011.
²⁸<http://bass-db.gforge.inria.fr/peass/>
²⁹Ward et al., 2018.
³⁰Stöter, Liutkus, and Ito, 2018.
³¹Rafii et al., 2019.
³²Stöter et al., 2019.

HPSS – PEASS results

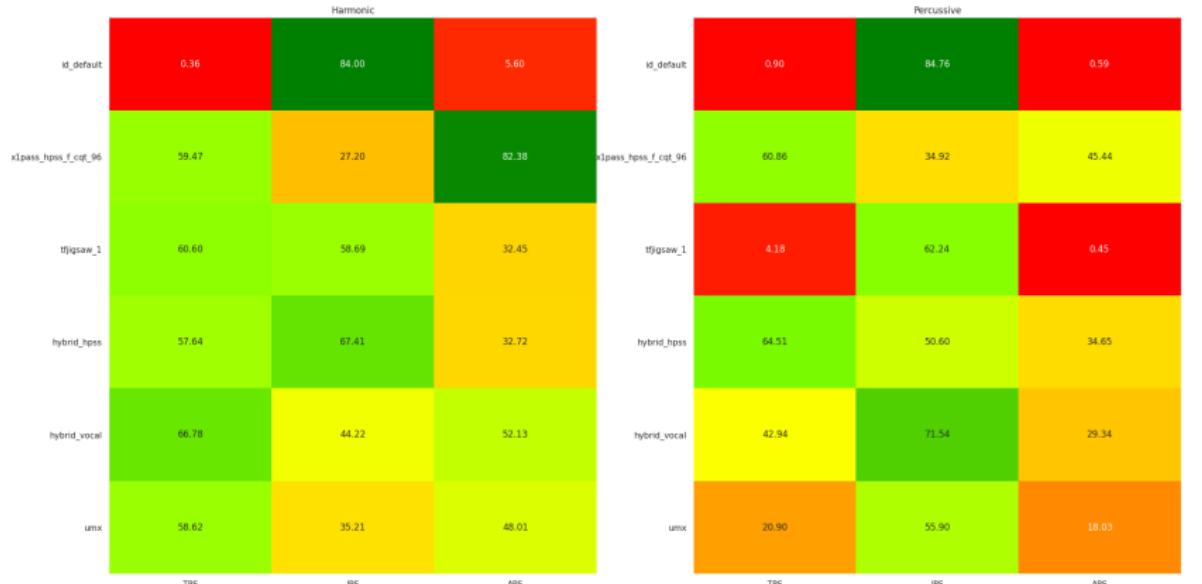


Figure: HPSS algorithms – final heatmap, PEASS scores

TF representations for music separation

2021-04-16

└ HPSS – PEASS results

HPSS – PEASS results

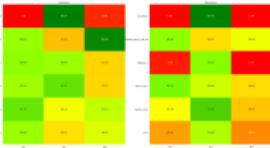


Figure: HPSS algorithms – final heatmap, PEASS scores

HPSS – BSSv4 results

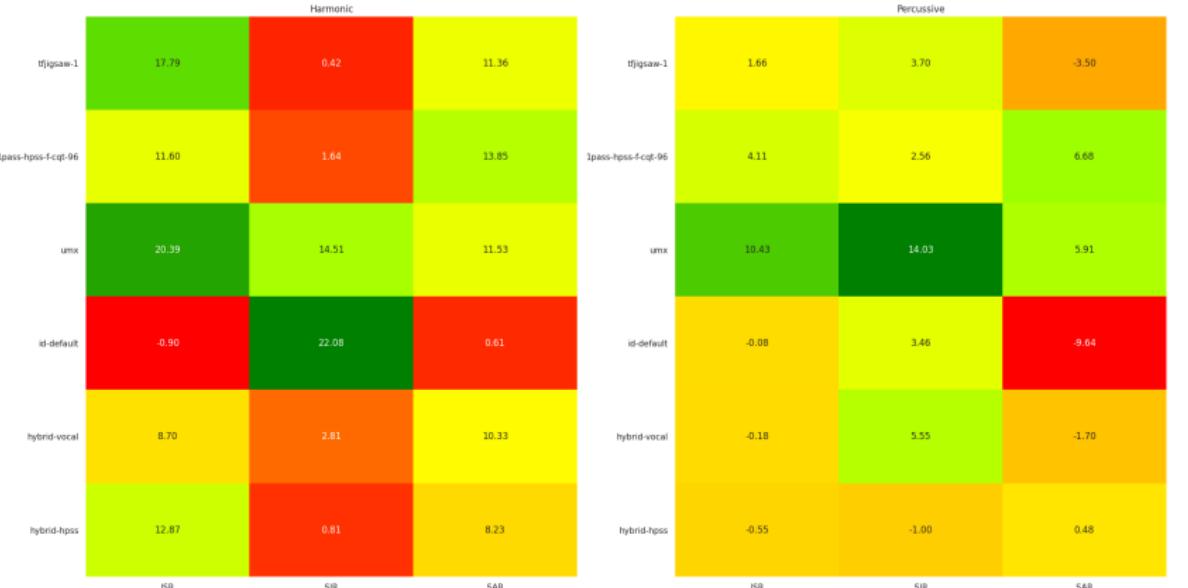


Figure: HPSS algorithms – final heatmap, BSSv4 scores

TF representations for music separation

2021-04-16

└ HPSS – BSSv4 results



HPSS + vocal – PEASS results

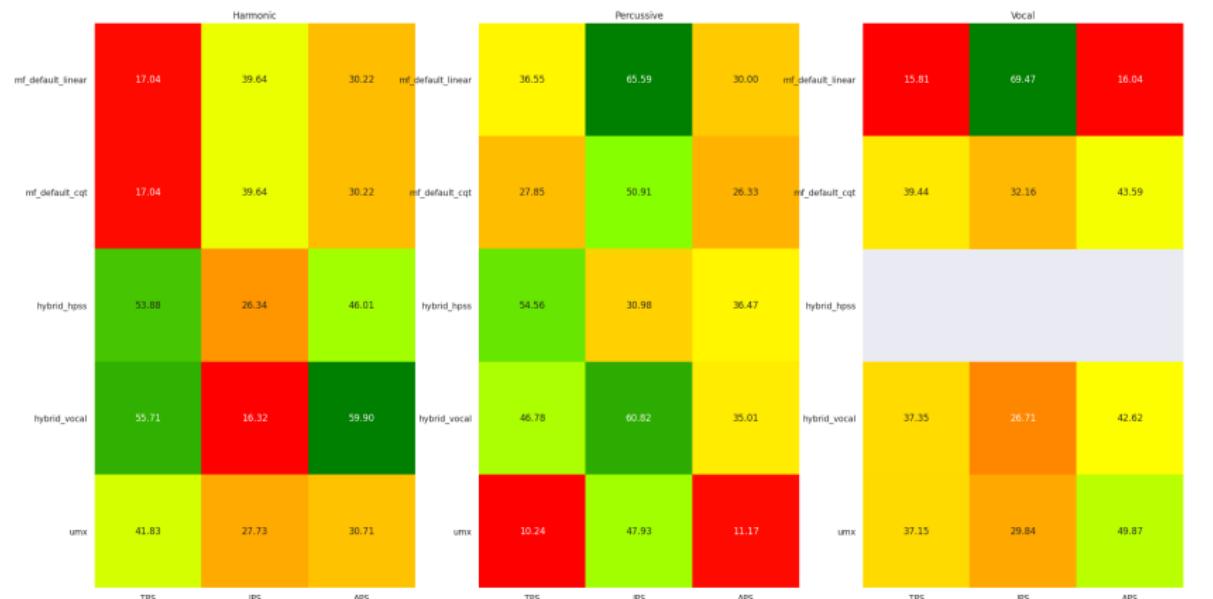
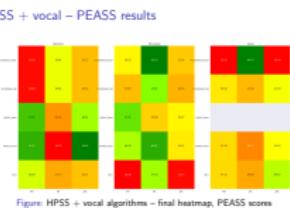


Figure: HPSS + vocal algorithms – final heatmap, PEASS scores

TF representations for music separation

2021-04-16

└ HPSS + vocal – PEASS results



HPSS + vocal – BSSv4 results



Figure: HPSS + vocal algorithms – final heatmap, BSSv4 scores

TF representations for music separation

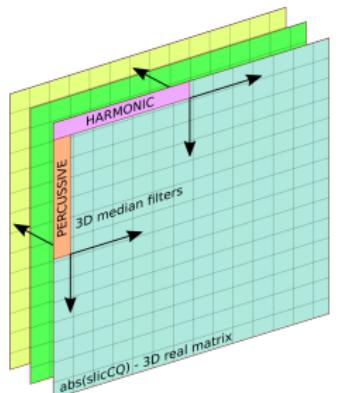
2021-04-16

└ HPSS + vocal – BSSv4 results



Best 1-pass HPSS – CQT-96, sliCQ median filter

Best performing 1-pass offline/anticausal algorithm: CQT with 96 bins-per-octave + median filter + **soft** mask. Realtime/causal: sliCQ³³ with 12 bins-per-octave + median filter + **hard** mask



Realtime STFT:³⁴ input stream = hop size, $2 \times$ hop ringbuffer (window).
sliCQ transform: input stream = trlen, $4 \times$ trlen ringbuffer (sllen)

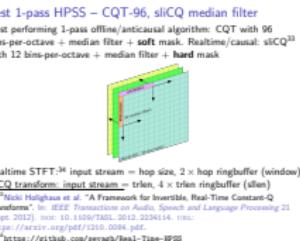
³³ Nicki Holighaus et al. "A Framework for Invertible, Real-Time Constant-Q Transforms". In: *IEEE Transactions on Audio, Speech and Language Processing* 21 (Sept. 2012). DOI: 10.1109/TASL.2012.2234114. URL: <https://arxiv.org/pdf/1210.0084.pdf>.

³⁴ <https://github.com/sevagh/Real-Time-HPSS>

TF representations for music separation

2021-04-16

└ Best 1-pass HPSS – CQT-96, sliCQ median filter



- 3D sliCQ coefficients can be median filtered, analogous to RGB/RGBA images

Best 2-pass HPSS – hybrid

Hybrid HPSS: combine TFJigsaw and STFT median filtering

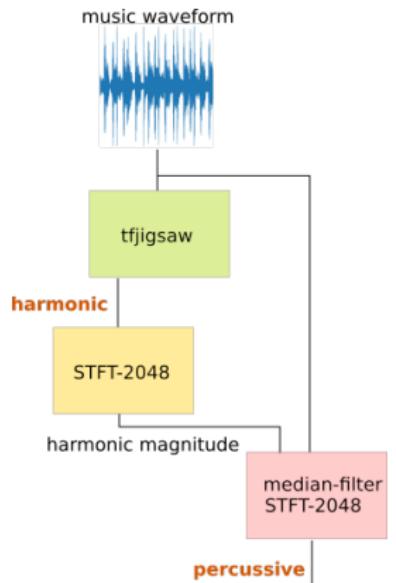


Figure: Block diagram of hybrid HPSS algorithm

TF representations for music separation

2021-04-16

└ Best 2-pass HPSS – hybrid

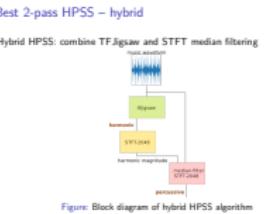
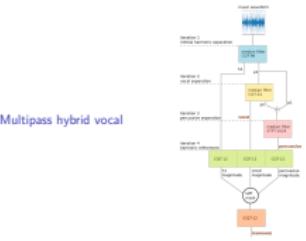


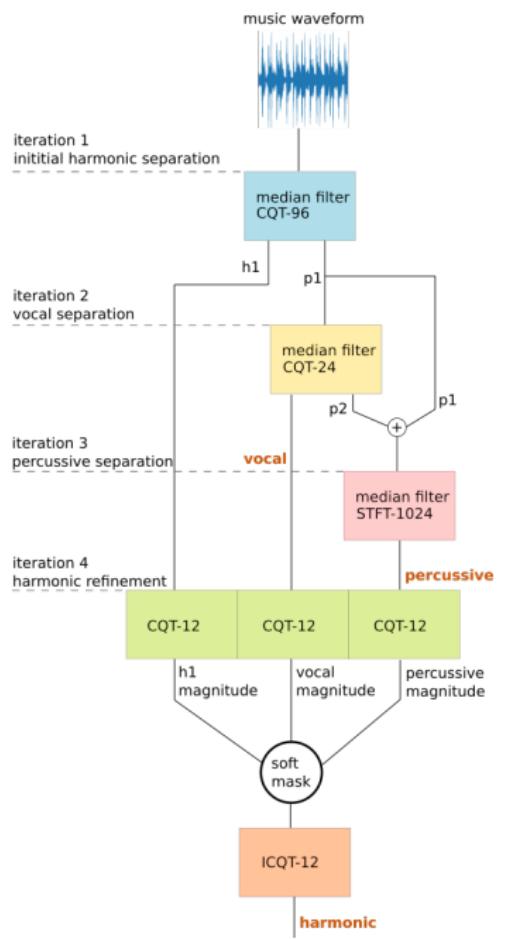
Figure: Block diagram of hybrid HPSS algorithm



TF representations for music separation

2021-04-16

Multipass hybrid vocal



HPSS – audio clips

Algorithm	Harmonic	Percussive
Reference mix	🔊 h	🔊 p
1-pass CQT 96 bins + soft mask ³⁵	🔊 h	🔊 p
Realtime sliCQ 12 bins + hard mask	🔊 h	🔊 p
Iterative Driedger ³⁶	🔊 h	🔊 p
Hybrid-HPSS	🔊 h	🔊 p
UMX	🔊 h	🔊 p

³⁵Fitzgerald, 2010.

³⁶Driedger, Müller, and Disch, 2014.

TF representations for music separation

2021-04-16

└ HPSS – audio clips

Algorithm	Harmonic	Percussive
Reference mix	🔊 h	🔊 p
1-pass CQT 96 bins + soft mask ³⁵	🔊 h	🔊 p
Realtime sliCQ 12 bins + hard mask	🔊 h	🔊 p
Iterative Driedger ³⁶	🔊 h	🔊 p
Hybrid-HPSS	🔊 h	🔊 p
UMX	🔊 h	🔊 p

³⁵Fitzgerald, 2010.
³⁶Driedger, Müller, and Disch, 2014.

HPSS + vocal – audio clips

Algorithm	Harmonic	Percussive	Vocal
Reference  mix	 h	 p	 v
Iterative Fitzgerald, vocal ³⁷	 h	 p	 v
Iterative Fitzgerald, percussive	 h	 p	 v
Hybrid-Vocal	 h	 p	 v
UMX	 h	 p	 v

³⁷Fitzgerald and Gainza, 2010.

TF representations for music separation

2021-04-16

└ HPSS + vocal – audio clips

HPSS + vocal – audio clips

Algorithm	Harmonic	Percussive	Vocal
Reference  mix	 h	 p	 v
Iterative Fitzgerald, vocal ³⁷	 h	 p	 v
Iterative Fitzgerald, percussive	 h	 p	 v
Hybrid-Vocal	 h	 p	 v
UMX	 h	 p	 v

³⁷Fitzgerald and Gainza, 2010.

- Real-world application of 622 concepts – sparsity, entropy, overcomplete dictionaries, pursuit
- Swap STFT for “better” TF representations in simple algorithms to improve source separation results
- Competitive PEASS separation results in hybrid algorithms based on advanced DSP/time-frequency analysis (not so good in BSSv4)
- Swap STFT for NSGT in both traditional DSP algorithms, and machine/deep learning networks³⁸ – lots of future potential

³⁸<https://github.com/sevagh/MiXiN>

TF representations for music separation

2021-04-16

└ Conclusions

Conclusions

- Real-world application of 622 concepts – sparsity, entropy, overcomplete dictionaries, pursuit
- Swap STFT for “better” TF representations in simple algorithms to improve source separation results
- Competitive PEASS separation results in hybrid algorithms based on advanced DSP/time-frequency analysis (not so good in BSSv4)
- Swap STFT for NSGT in both traditional DSP algorithms, and machine/deep learning networks³⁸ – lots of future potential

³⁸<https://github.com/sevagh/MiXiN>