

Jointly Detecting and Separating Singing Voice: A Multi-Task Approach

Daniel Stoller¹, Sebastian Ewert^{2*}, Simon Dixon¹

¹Centre for Digital Music
Queen Mary University of London

²Spotify
London

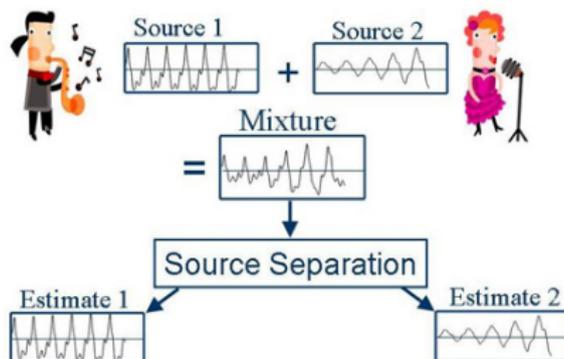
LVA ICA
05.07.2018

*Work was conducted at Queen Mary University of London

Vocal separation

Introduction

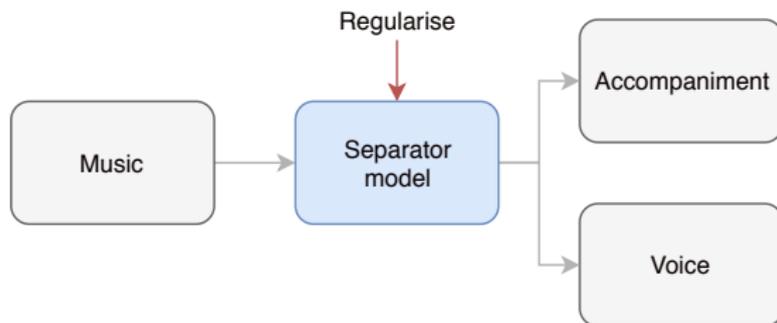
- Main task: Separate vocals from music pieces
- Applications: Karaoke generation, singer identification, voice analysis...



Vocal separation

Challenges

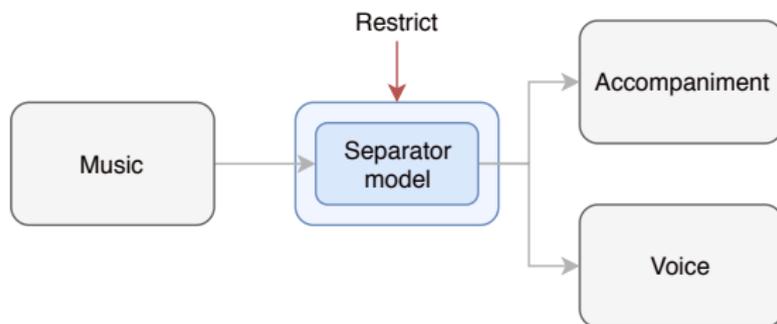
- Difficult task, small multi-track datasets \Rightarrow Overfitting
- Give model more knowledge:
 - Regularization (e.g. weight decay)



Vocal separation

Challenges

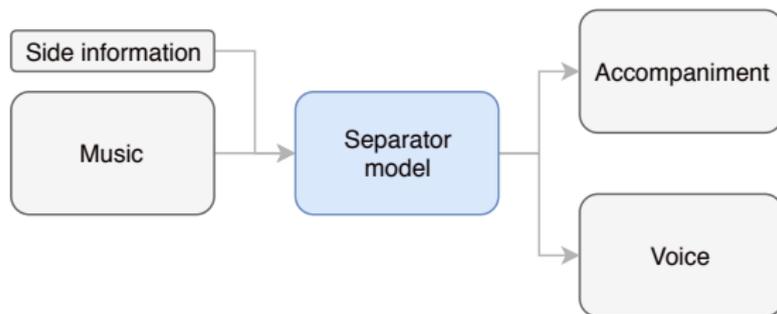
- Difficult task, small multi-track datasets \Rightarrow Overfitting
- Give model more knowledge:
 - Knowledge-driven (e.g. KAM [4])



Vocal separation

Challenges

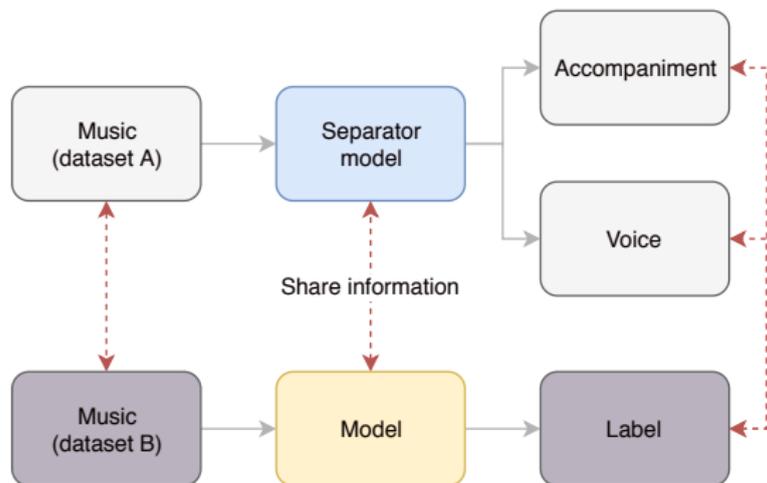
- Difficult task, small multi-track datasets \Rightarrow Overfitting
- Give model more knowledge:
 - Informed source separation [2]



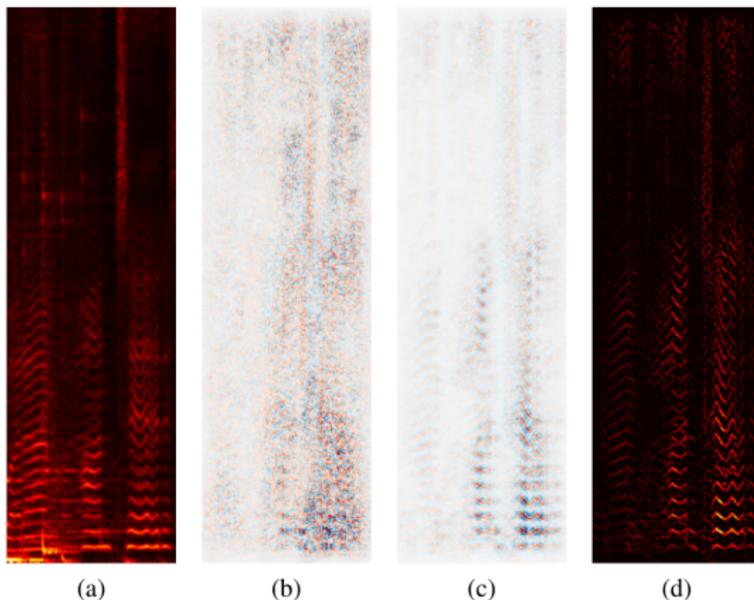
Vocal separation

Challenges

- Difficult task, small multi-track datasets \Rightarrow Overfitting
- Give model more knowledge:
 - **Integrate information from related tasks/datasets**



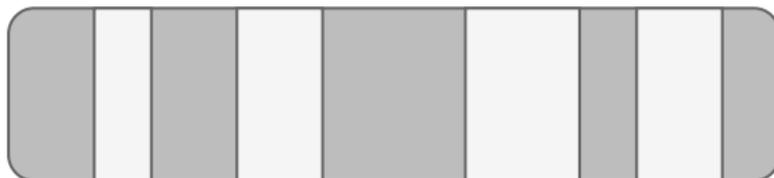
- Which other tasks could help?
- Vocal activity detection is promising:
 - Knowing vocal activity improves vocal separation [1]
 - Vocal detection networks learn a form of separation: [5]



Initial approach

Using additional non-vocal sections

- U-Net adaptation [3] as separator, MSE loss
 - **Sample instrumental sections also from SVD databases**
- ⇒ Diversifies instrumental training data



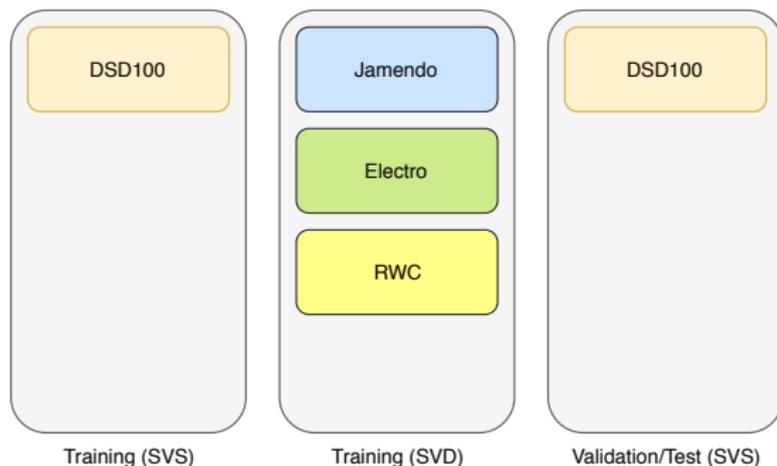
Song A
SVS Database



Song B
SVD Database

Initial approach

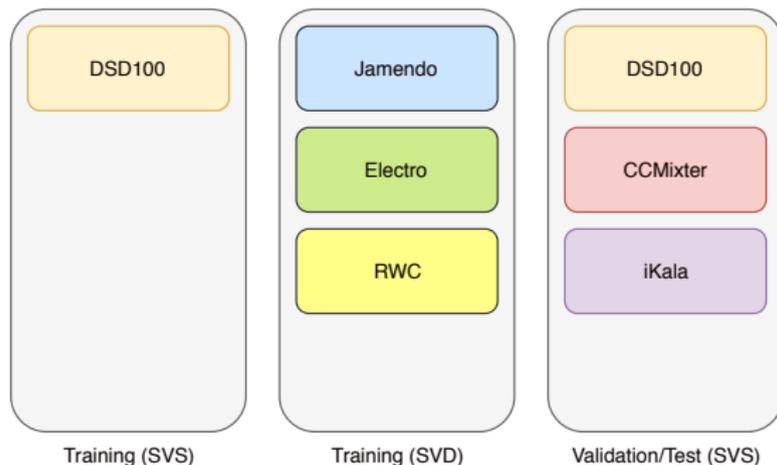
Results



Performance **decrease**

Initial approach

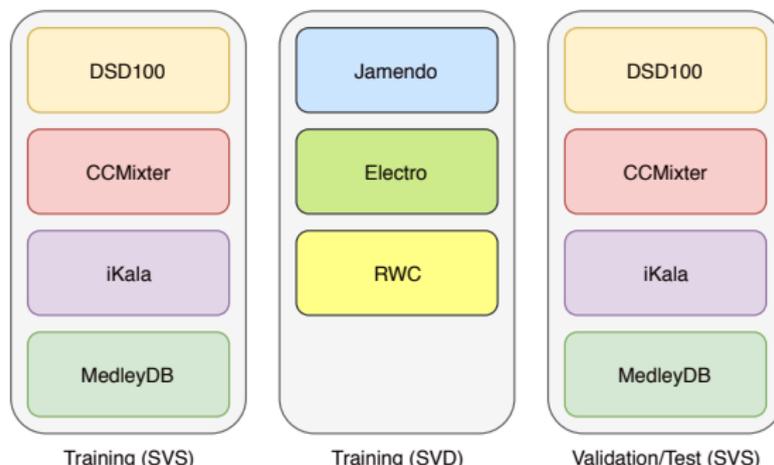
Results



Performance **increase**

Initial approach

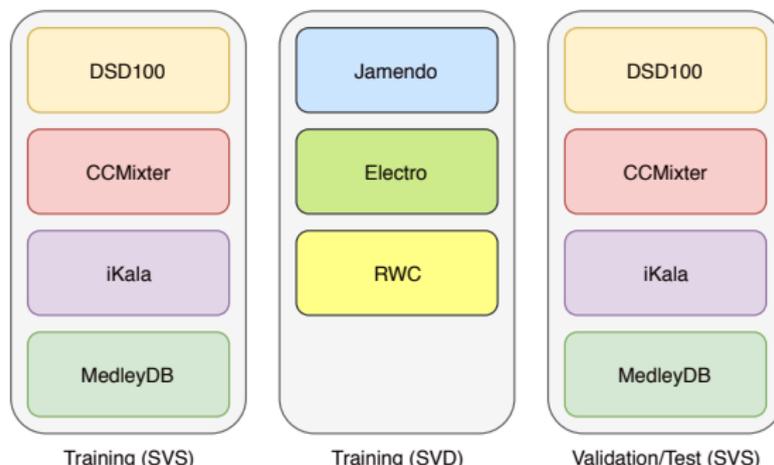
Results



Performance **decrease**

Initial approach

Results

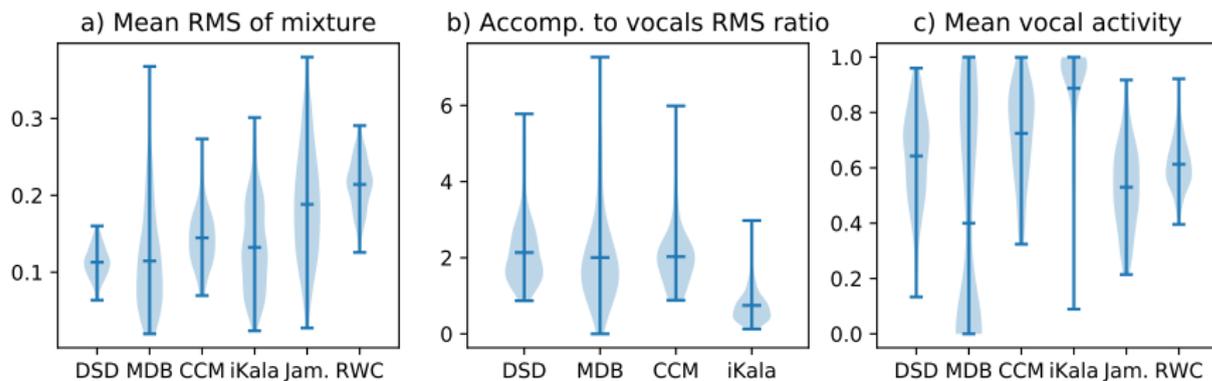


Performance **decrease**

Dataset bias?

Dataset bias

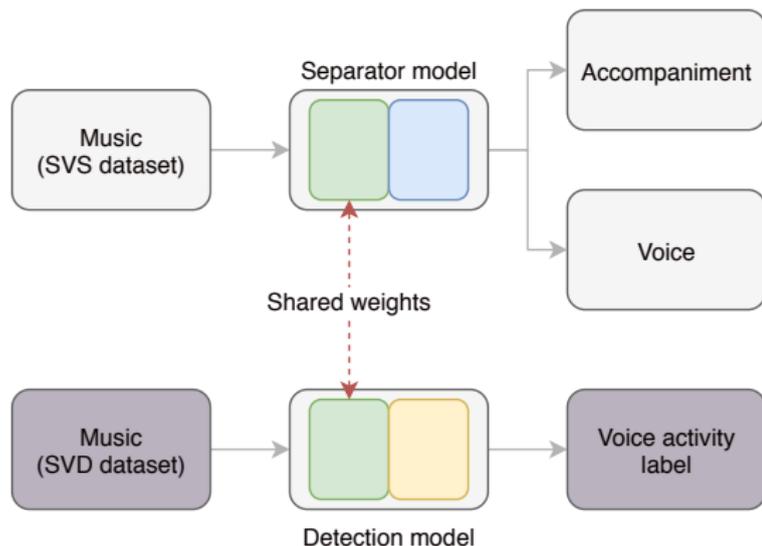
Analysis



Multi-task approach

Introduction and motivation

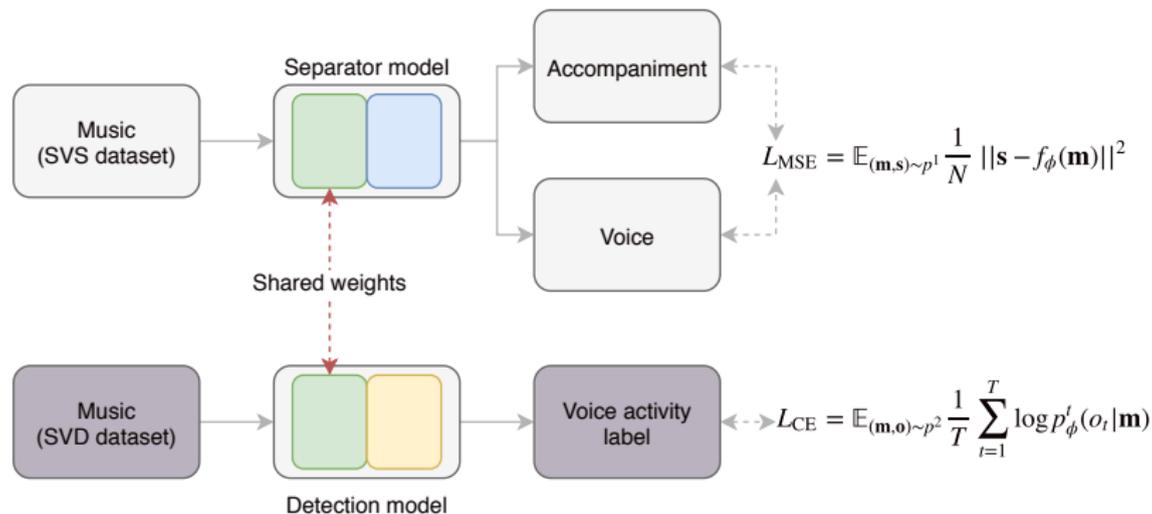
Key idea: Predict both audio and label



Multi-task approach

Introduction and motivation

Key idea: Predict both audio and label

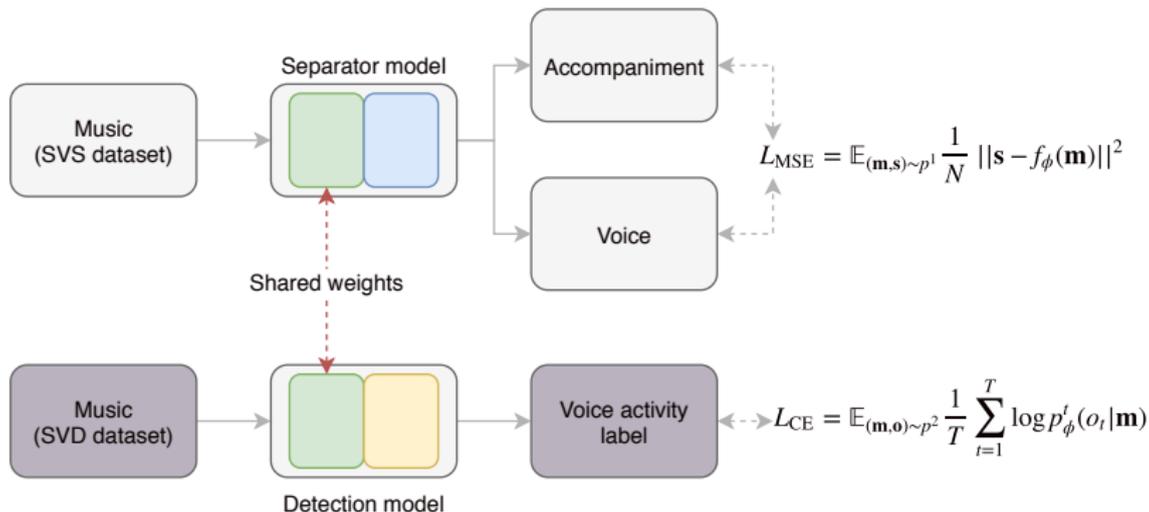


$$L_{\text{MTL}} = \alpha L_{\text{MSE}} + (1 - \alpha) L_{\text{CE}}$$

Multi-task approach

Introduction and motivation

Key idea: Predict both audio and label



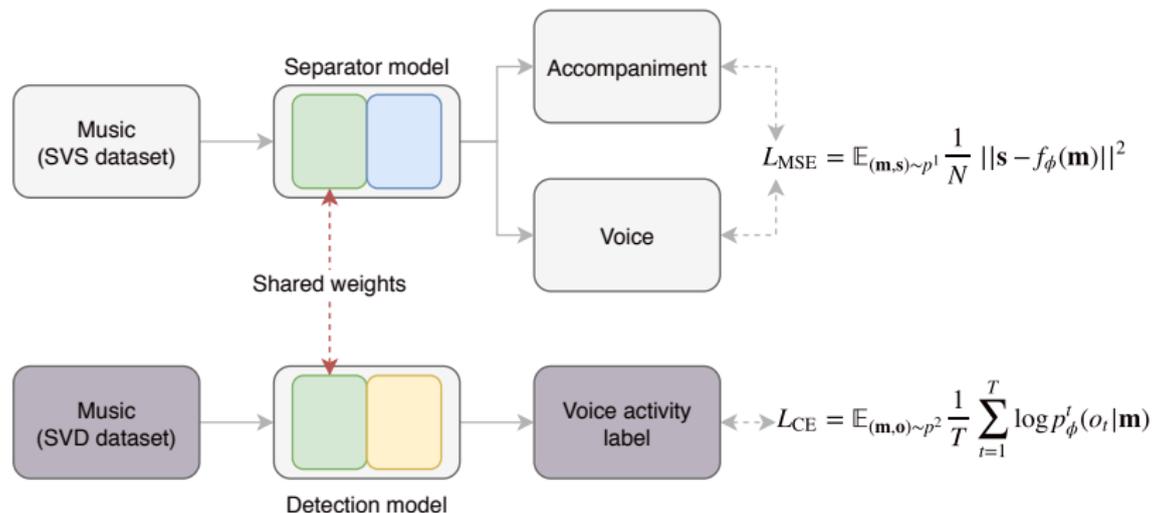
$$L_{\text{MTL}} = \alpha L_{\text{MSE}} + (1 - \alpha) L_{\text{CE}}$$

Robust to dataset bias and label accuracy

Multi-task approach

Introduction and motivation

Key idea: Predict both audio and label



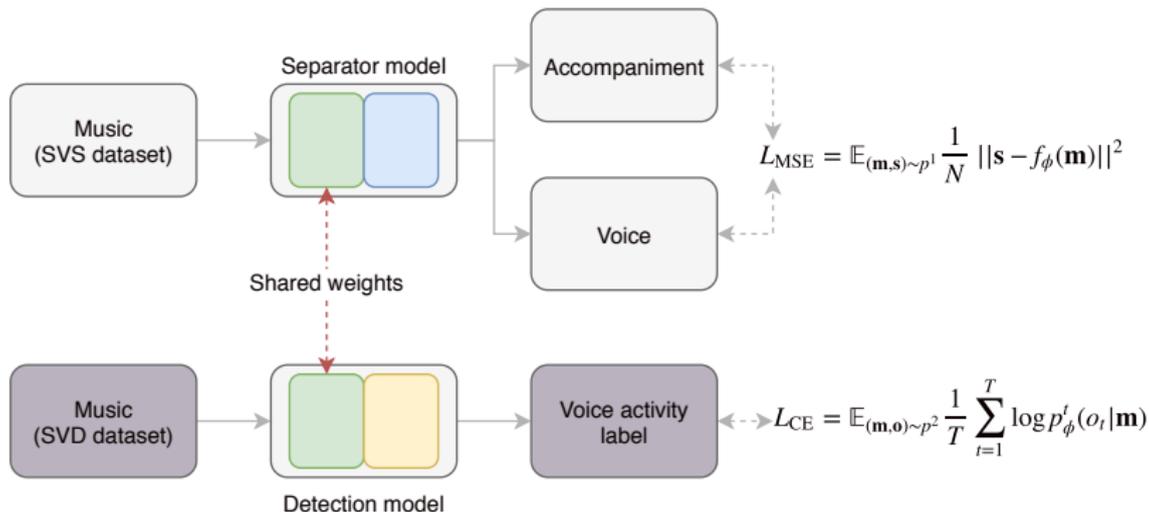
$$L_{\text{MTL}} = \alpha L_{\text{MSE}} + (1 - \alpha) L_{\text{CE}}$$

Can train with vocal sections from SVD data

Multi-task approach

Introduction and motivation

Key idea: Predict both audio and label



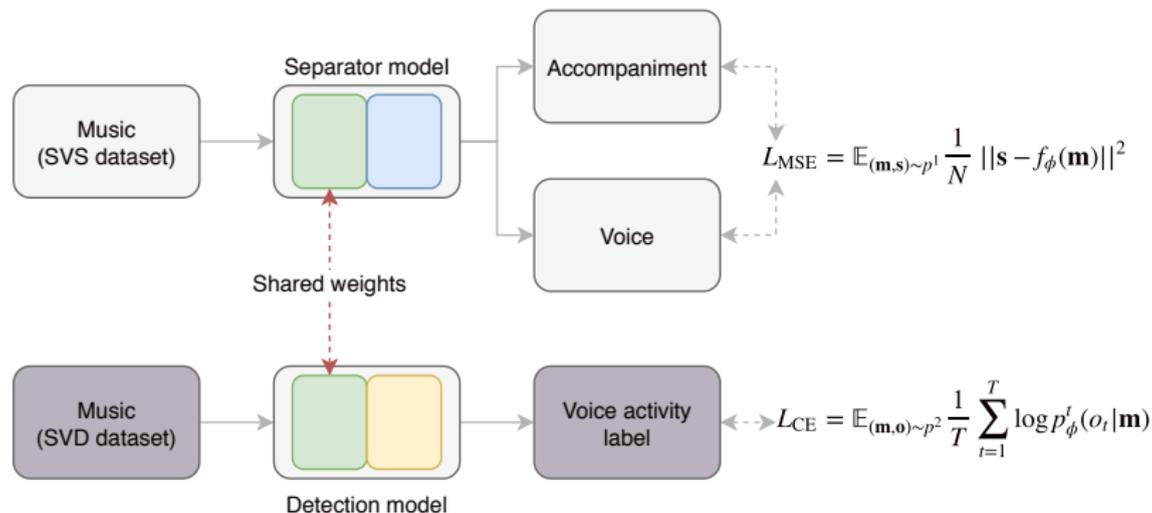
$$L_{\text{MTL}} = \alpha L_{\text{MSE}} + (1 - \alpha) L_{\text{CE}}$$

Needs only mixture at test time

Multi-task approach

Introduction and motivation

Key idea: Predict both audio and label

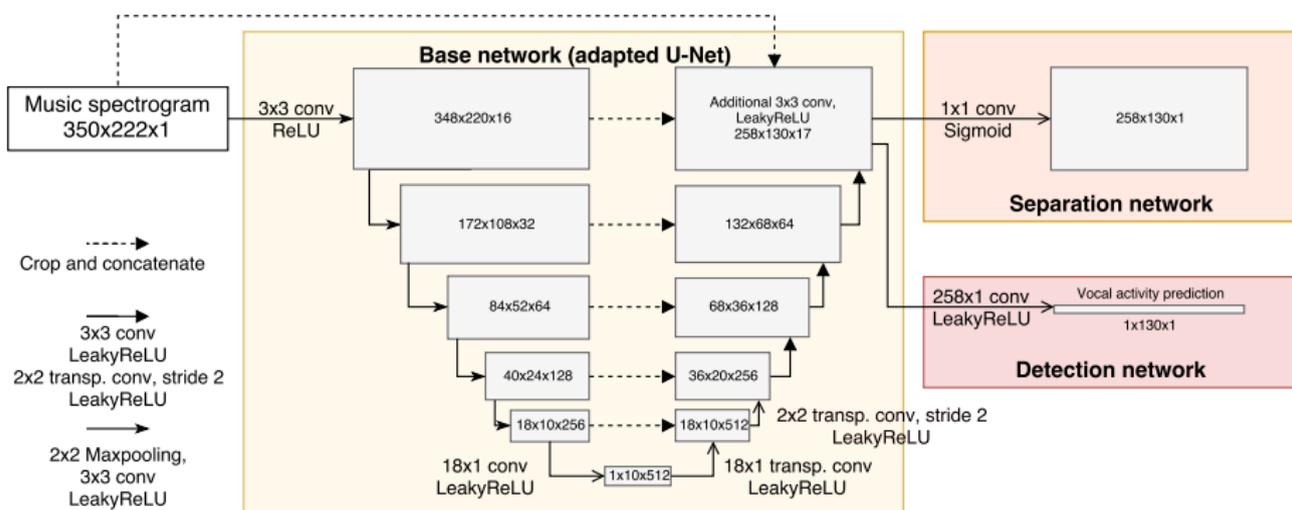


$$L_{\text{MTL}} = \alpha L_{\text{MSE}} + (1 - \alpha) L_{\text{CE}}$$

Solves two tasks at once

Experimental setup

Model architecture and dataset



- DSD100 as SVS, Jamendo as SVD training data

Experimental setup

Evaluation metrics: AU-ROC, MSE, SDR

- AU-ROC for SVD
 - MSE and SDR/SIR/SAR for separation
 - SDR gives $\log(0)$ for non-vocal sections ($\approx 10\%$)
- ⇒ Also measure RMS of vocal estimates for non-vocal sections

Results

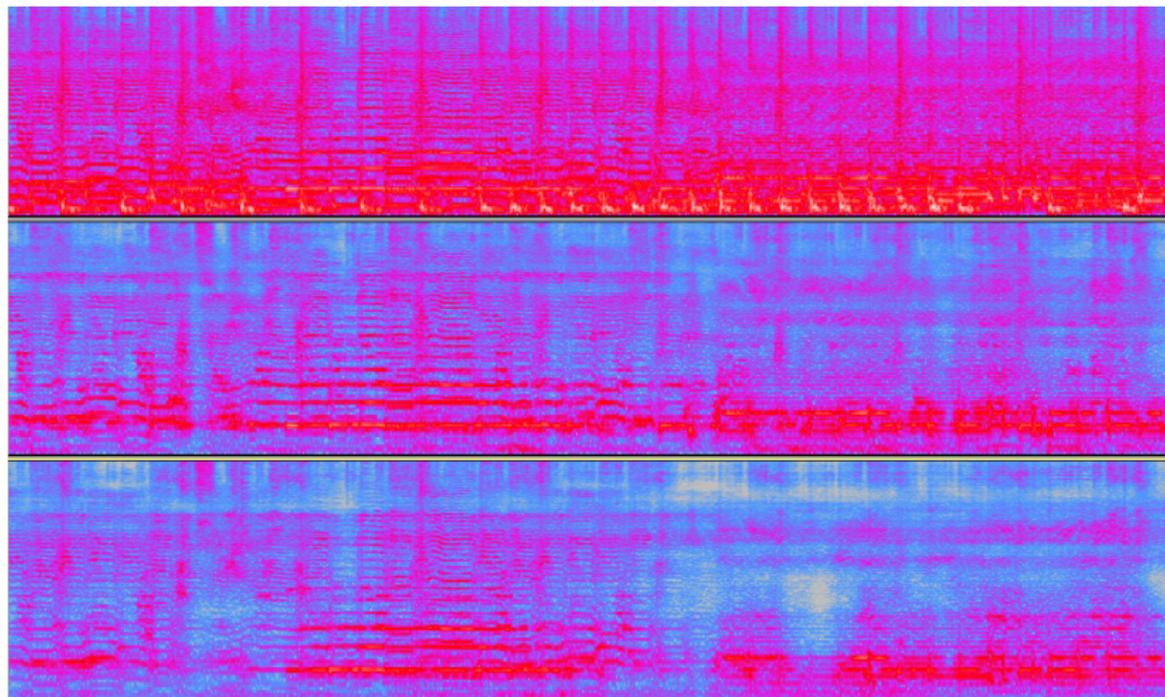
Single-task vs. multi-task model

		Metric									
		AU-ROC	MSE	Non-voc. RMS	Vocals			Accompaniment			
					SDR	SIR	SAR	SDR	SIR	SAR	
Model		SVD	0.9239	-	-	-	-	-	-	-	
		SVS	-	0.01865	0.0194	2.83	5.27	6.88	6.71	14.75	13.25
		Ours	0.9250	0.01755	0.0155	2.86	5.56	6.23	6.69	13.24	14.11

Table: Comparing SVS and SVD baseline with our approach

Results

Qualitative comparison



- Current SotA methods only use multi-track data
- Our approach also uses SVD databases
- Improved separation and detection performance
- Future work: Larger datasets, more related tasks



T.-S. Chan, T.-C. Yeh, Z.-C. Fan, H.-W. Chen, L. Su, Y.-H. Yang, and R. Jang.

Vocal activity informed singing voice separation with the ikala dataset. In *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*, pages 718–722. IEEE, 2015.



S. Ewert, B. Pardo, M. Müller, and M. D. Plumbley.

Score-informed source separation for musical audio recordings: An overview.

IEEE Signal Processing Magazine, 31(3):116–124, 2014.



A. Jansson, E. J. Humphrey, N. Montecchio, R. Bittner, A. Kumar, and T. Weyde.

Singing voice separation with deep U-Net convolutional networks. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, pages 323–332, 2017.



A. Liutkus, D. Fitzgerald, and Z. Rafii.

Scalable audio separation with light kernel additive modelling.

In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 76–80. IEEE, 2015.



J. Schlüter.

Learning to pinpoint singing voice from weakly labeled examples.
In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, pages 44–50, 2016.

Wave-U-Net

A Multi-Scale Neural Network for
End-to-End Audio Source Separation

DANIEL STOLLER¹, SEBASTIAN EWERT², SIMON DIXON¹

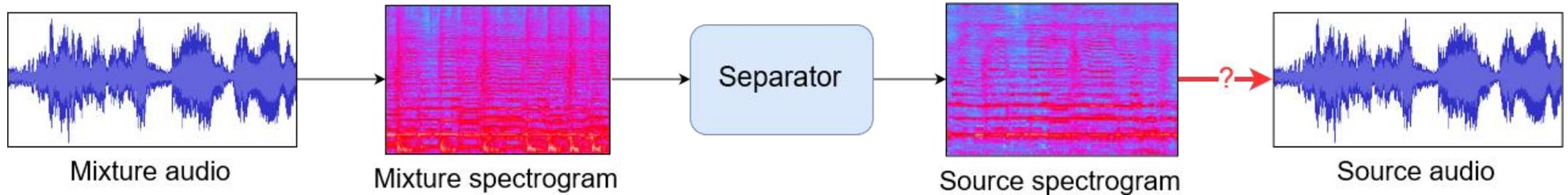
¹ QUEEN MARY UNIVERSITY OF LONDON

² SPOTIFY

Previous work

Mostly spectrogram-based [1,2,3]

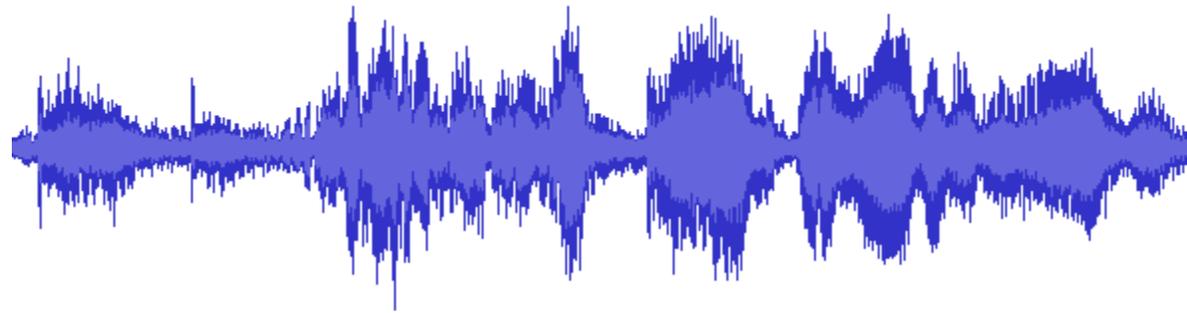
- Problem: Reconstruct source signal from its spectrogram estimates
- Result: Output artifacts



Previous work

Recently: Few time-domain approaches [4,5]

- Problem: Model long-term dependencies in raw audio
- Result: Context-deprived [4] or slow [5] models



2 s, but
88200 samples

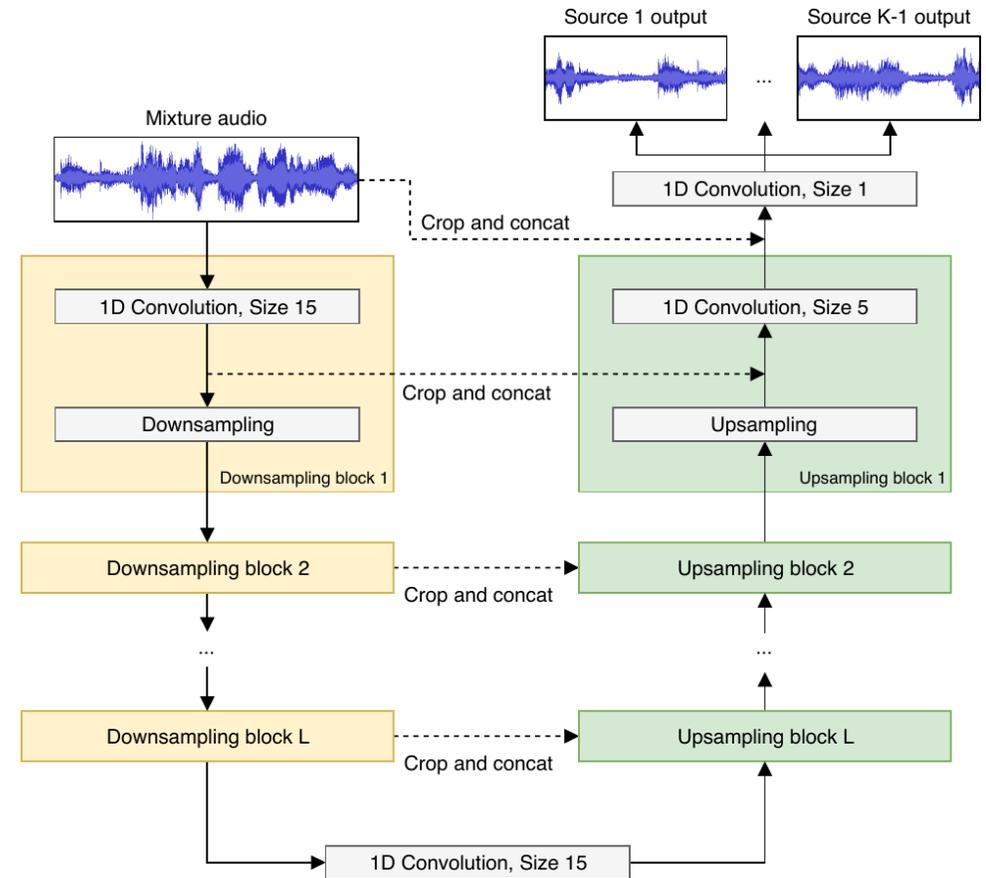
Our solution: Wave-U-Net

Inspired by U-Net [1,6] and Wavelets

Core idea: Feature hierarchy

- Features at different timescales
- Efficient long-term dependency modelling

Simple system



Upsampling

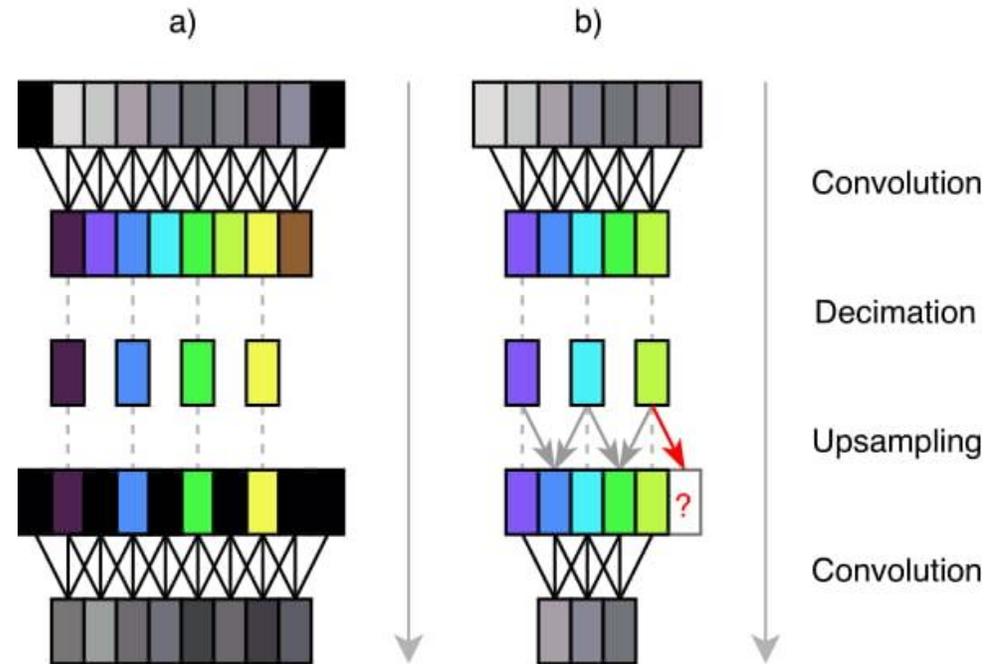
Commonly used:
Transposed convolutions

Introduces high-frequency noise

Solutions:

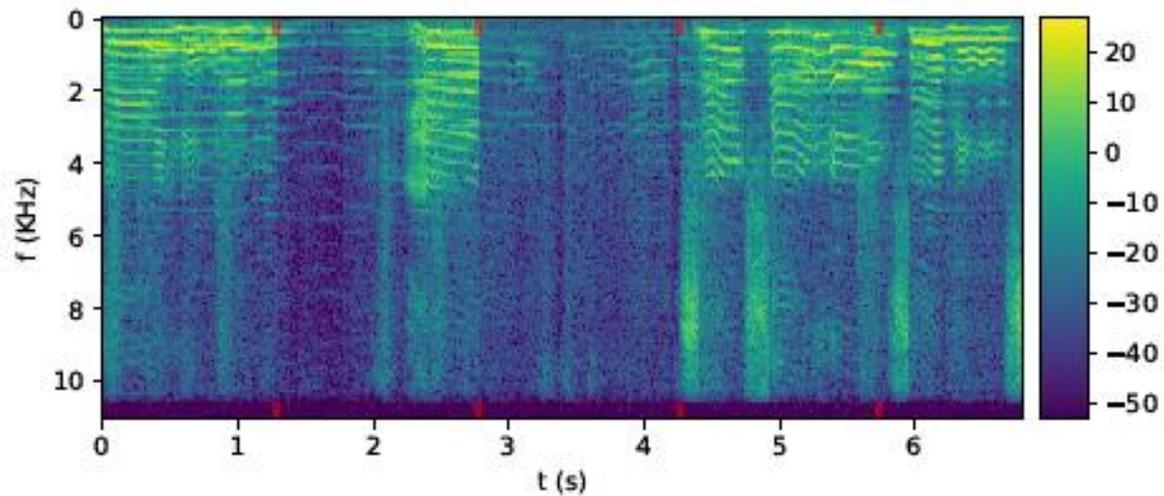
- Linear interpolation
- Learned upsampling:

$$f_{t+0.5} = \sigma(w) \odot f_t + (1 - \sigma(w)) \odot f_{t+1}$$



Context-aware prediction

Border artifacts with existing systems (equal no. of input & output timesteps):



Solution: No zero-padding for convolutions

=> Prediction of source only for centre piece of mixture input

Results

Improvements over spectrogram-based equivalent

Encouraging performance in SiSec challenge

Further improvements with

- Context-aware prediction
- Stereo handling



Code, trained model and audio examples:

<https://github.com/f90/Wave-U-Net>

References

- [1] Jansson, A.; Humphrey, E. J.; Montecchio, N.; Bittner, R.; Kumar, A. & Weyde, T.
Singing Voice Separation with Deep U-Net Convolutional Networks
Proceedings of the International Society for Music Information Retrieval Conference (ISMIR), **2017**, 323-332
- [2] Huang, P.-S.; Chen, S. D.; Smaragdis, P. & Hasegawa-Johnson, M.
Singing-voice separation from monaural recordings using robust principal component analysis
2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), **2012**, 57-60
- [3] Uhlich, S.; Giron, F. & Mitsufuji, Y.
Deep neural network based instrument extraction from music
2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), **2015**, 2135-2139
- [4] Grais, E. M.; Ward, D. & Plumbley, M. D.
Raw Multi-Channel Audio Source Separation using Multi-Resolution Convolutional Auto-Encoders
arXiv preprint arXiv:1803.00702, **2018**
- [5] Luo, Y. & Mesgarani, N.
TasNet: time-domain audio separation network for real-time, single-channel speech separation
CoRR, **2017**, *abs/1711.00541*
- [6] Ronneberger, O.; Fischer, P. & Brox, T.
U-net: Convolutional networks for biomedical image segmentation
International Conference on Medical Image Computing and Computer-Assisted Intervention, **2015**, 234-241