

# Transfer Learning with Jukebox for Music Source Separation

Wadhah Zai El Amri, Oliver Tautz, Helge Ritter, and Andrew Melnik

Bielefeld University, Germany

wadhah.zai.papers@gmail.com, oliver.tautz.papers@gmail.com,  
helge@techfak.uni-bielefeld.de, andrew.melnik.papers@gmail.com

**Abstract.** In this work, we demonstrate how a publicly available, pre-trained *Jukebox* model can be adapted for the problem of audio source separation from a single mixed audio channel. Our neural network architecture, which is using transfer learning, is quick to train and the results demonstrate performance comparable to other state-of-the-art approaches that require a lot more compute resources, training data, and time. We provide an open-source code implementation of our architecture (<https://github.com/wzaielamri/unmix>)

## 1 Introduction and Related Work

Source separation is an important issue in many fields such as audio processing, image processing [9], EEG [10,8], etc. Music source separation from mixed audio is a challenging problem, especially if the source itself should be learned from a dataset of examples. Additionally, such models are expensive to train from scratch. We tested our model on the MUSDB18-HQ [14] dataset which supplies full songs with ground truth stems of 'bass', 'drums', 'vocals' and 'other', which includes instruments such as guitars, synths, etc. The task is to separate a mixed audio channel into the separately recorded instruments, called stems here. Most baseline models in the Music Demixing Challenge 2021 [11] used masking of input transformed to the frequency domain by short-time Fourier transformation such as UMX [18]. This older technique transforms the waveform into a three dimensional input that can be viewed as a picture. The model then computes a 'mask' which subtracts parts of the audio in frequency domain before transforming back to listenable audio. *Demucs* [1] on the other hand showed a successful approach that works in waveform domain, where an autoencoder, based on a bidirectional long short-term memory network, is used.

This success of such a solution (using waveforms) encouraged us to adapt *Jukebox* [3], a powerful, generative model using multiple, deep Vector Quantized-Variational Autoencoders (VQ-VAE) [12] to automatically generate real sounding music, and using its publicly available pre-trained weights for the task.

Transfer learning helped deep learning models reach new heights in many domains, such as natural language processing [2,13] and computer vision [5,4]. Although relatively unexplored for the audio domain, [7] proved feature representation learned on speech data could be used to classify sound events. Their

results verify that cross-acoustic transfer learning performs significantly better than a baseline trained from scratch. TRILL [16] showed great results of pre-training deep learning models with an unsupervised task on a big dataset of speech samples. Its learned representations exceeded SOTA performance on several downstream tasks with datasets of limited size.

We take a similar approach that is heavily based on *Jukebox* [3]. It uses multiple VQ-VAEs to compress raw audio into discrete codes. They are trained self-supervised, on a large dataset of about 1.2 million songs, needing the compute power of 256 V100 to train in an acceptable time. Our experiments show that *Jukebox's* learned representations can be used for the task of source separation.

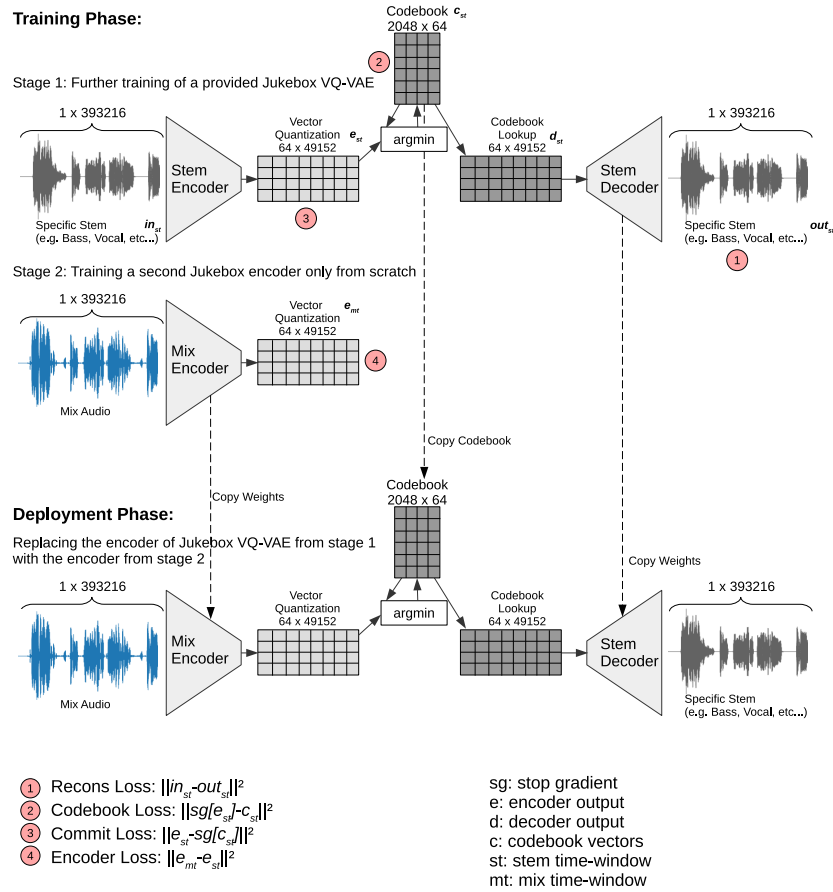


Fig. 1: Visualization of the proposed transfer learning model architecture.

## 2 Method

### 2.1 Architecture

Our architecture utilizes the default *Jukebox*'s [3] variant of the publicly available pre-trained VQ-VAE model. *Jukebox* uses three separated VQ-VAEs. We use only the smallest one with the strongest compression, due to the low resource usage. It employs dilated 1-D convolutions in multiple residual blocks to find a less complex sequence representation of music. An audio sequence  $x_t$  gets mapped by an encoder  $E_1$  to a latent space  $e_t = E_1(x_t)$  of 64 dimensions so that it can be mapped to the closest prototype vector in a collection  $C$  of vectors called *codebook*. These 2048 prototype vectors, denoted  $c_{st}$ , are learned in training and help to form a high-quality representation.

The rate of compression for a sequence is called the hop length, for which a value of 8 is used. It depends on the stride values of the convolutional layers. We set the stride value to 2 as well as the downsampling to 3. All other values remain as defined in [3]. After mapping to the codebook, a decoder  $D$  aims to reconstruct the original sequence. In summary, equation (1)

$$y_t = D(\operatorname{argmin}_c(\|E_1(x_t) - c\|)) \text{ for } c \in C \quad (1)$$

describes a full forward pass through the VQ-VAE, where  $y_t$  is the prediction for an input sequence  $x_t$  and  $\|\cdot\|$  is the euclidean norm. For further technical details on the used VQ-VAE architecture, refer to the paper of Dhariwal et al[3]. The model is fine-tuned on data for one stem, learning good representations for a single instrument. In addition, we train a second encoder  $E_2$ , identical to the one already mentioned, to project an input sequence of the mixture to the space already known by the codebook and decoder. For deployment, the encoder of the VQ-VAE is switched with the second one, effectively mapping from the full mixture to one stem.

### 2.2 Data

Our models are trained on the MUSDB18-HQ [14] dataset, also used in the music demixing challenge [11]. It consists of 150 full-length songs, sampled at 44KHz, providing the full audio mixture and four stems, 'vocals', 'bass', 'drums', and 'other' for each sample, which can be regarded as ground truth in the context of source separation. We train on the full train set composed of 100 songs, testing is done on the remaining 50.

### 2.3 Training

For each stem  $i=1..4$ , we train a model in two phases (see Fig. 1). In the first phase, the model is trained on data that present the chosen stem in isolation (i.e. not embedded in a mixture). This produces a VQ-VAE with a "single stem encoder" ( $SE_i$ ) that can map a single stem into a good latent representation,

followed by a "stem decoder" ( $SD_i$ ) tuned to reconstruct the input after the "discretization bottleneck" as faithfully as possible. Training of each VQ-VAE is based on the same three losses as chosen in the original Jukebox paper [3] :  $L = L_{recons} + L_{codebook} + \beta L_{commit}$ . However, our final goal is to process each stem when it is part of a mixture of all four stems. Such embedding will introduce distortion of each stem, requiring to replace each single stem encoder  $SE_i$  from phase 1 by a corresponding "mixture stem encoder" ( $ME_i$ ) that is trained in phase 2, to map its changed (mixture embedded) stem  $i$  input onto the representation (stem- $i$  codebook prototypes) created in phase by the  $SE_i$ - $SD_i$  encoder-decoder pair. So, for each stem  $i$  (now omitting index  $i$  in the following) for each training sample ( $x_{mt}$ : the sequence of the mixed audio,  $x_{st}$  : the sequence of stem audio), we feed  $x_{st}$ , to the already trained encoder SE, producing  $e_{st}$ . Separately, the full mixture  $x_{mt}$  is passed through the new encoder ME, yielding  $e_{mt}$ . Now, we can backpropagate through ME using MSE loss  $\|e_{st} - e_{mt}\|^2$  (keeping SE fixed throughout phase 2). It would also be possible to train E2 end-to-end, similar to training phase 1. This could be done in two ways, freezing the weights of the rest of the VQ-VAE or fine-tuning further. With frozen weights, the best reachable performance would be identical to our current approach, but training time would be increased. Continuing to fine-tune would most probably lead to more source bleeding because the embeddings no longer represent only the stem they were trained on in training phase 1, so we decided against training end-to-end.

Finally, we obtain our mixture-adapted final VQ-VAE by concatenating the trained mixture stem encoder ME with the stem decoder SD. Note that this procedure will be carried out for each of the four stems, yielding four correspondingly optimized "stem mixture encoder- decoders" that together provide our decomposition of the mixture input into its stem constituents. On a more technical note, in both training phases and deployment, the data is processed chunk-wise, with a size of about 9 seconds.

For a clear overview of the content of this section, refer to figure 1. All conducted experiments that will be defined in the next section were computed on two Tesla GPUs with 16Gb each. The length of each input sequence is equal to 393216 data points as used by *Jukebox*. The batch size is equal to 4.

$$SDR_{stem} = 10 \log_{10} \frac{\sum_n \|s(n)\|^2 + \epsilon}{\sum_n \|s(n) - \hat{s}(n)\|^2 + \epsilon} \quad (2)$$

To benchmark the conducted experiments, signal-to-distortion ratio (SDR) metric is used, which is a common metric in other SOTA papers[1][18][6][15][17]. It is computed by equation (2), as stated in [11], where  $s(n)$  is the values of the ground truth and  $\hat{s}(n)$  depicts the values of the prediction. A small value  $\epsilon = 10^{-7}$  is added to avoid division by zero. 'Total' SDR is the mean SDR for all stems.

SDR Values					
Method	Drum	Bass	Other	Vocal	Total
Our Approach (i)	4.925	4.073	2.695	5.060	4.188
Our Approach trained from scratch (ii)	-0.002	-0.087	-0.026	0.00	-0.028
Our Approach without finetuning(iii)	1.06	-1.072	0.79	0.33	0.279

Table 1: Comparison of SDR values per stem and in total for three different versions of our approach.

SDR Values					
Method	Drum	Bass	Other	Vocal	Total
DEMUCS	6.509	6.470	4.018	6.496	5.873
Our Approach	4.925	4.073	2.695	5.060	4.188
Wave-U-Net	4.22	3.21	2.25	3.25	3.23
ScaledMixturePredictor	0.578	0.745	1.136	1.090	0.887

Table 2: Comparison of SDR values per stem and in total. Our approach outperforms both the *ScaledMixturePredictor*, the basic baseline in the Music Demixing Challenge [11] and *Wave-U-Net* [17], a classic approach of source separation in the waveform domain while *Demucs* [1] achieves current SOTA performance on the Dataset.

### 3 Experiments and Results

The main key point of this paper consists of demonstrating that it is possible to get decent audio quality by using transfer learning. For this, we did three different experiments (i), (ii), and (iii) on the four audio stems. In experiment (i) we trained each audio stem’s stem encoder SE from scratch without using any pretraining values. For the second experiment (ii) we trained the SE’s with initial weights chosen as the pre-trained weights of Jukebox. Experiment (iii) is conducted to show the impact of fine-tuning the SE. We skip fine-tuning in training phase 1 and use the weights of Jukebox as is and train a ME like in the other phases.

For all these VQ-VAE, we pick the checkpoint 80K and train the corresponding mixture encoders ME in phase 2. For these, and in both experiments (i) (ii) and (iii), we initialized their weights randomly. For the first experiment, we found out that all results are bad, and no good audio quality is reached for the extracted stems. The SDR values are equal to or near 0 for all four stems. For the second experiment, the model converges after 32 hours of training in total on two Tesla GPU units with 16GB of VRAM each. (iii) shows that fine-tuning in phase 1 is important.

Figure 2 demonstrates decent SDR values for networks trained with pre-trained weights in comparison to others trained with randomly initialized weights from scratch. It can also be deduced that in order to get fairly good SDR val-

ues, it is enough to train until early checkpoint values, such as 20K. Then, the checkpoint 20K is reached after 16 hours for each of the two models on two Tesla GPUs. Table 2 gives a comparison of different approaches for audio signal separation. Here, our approach achieves comparable results when benchmarked against other state-of-the-art networks.

In terms of deployment, we need 0.73 seconds of CPU processing time for an audio chunk of 8.91 seconds per stem, which translates into an RTF (Real-Time Factor) of 0.082. That seems quick, but the current implementation of our model takes an audio chunk of 8.91 seconds. As with most approaches to audio source separation, direct online processing is not possible. However, a single code vector of the VQ-VAE represents a sliding window of an input vector of 49152 timesteps, which corresponds to about 1s of audio with a sampling rate of 44.1k. So, we think that an adapted architecture could be more flexible.

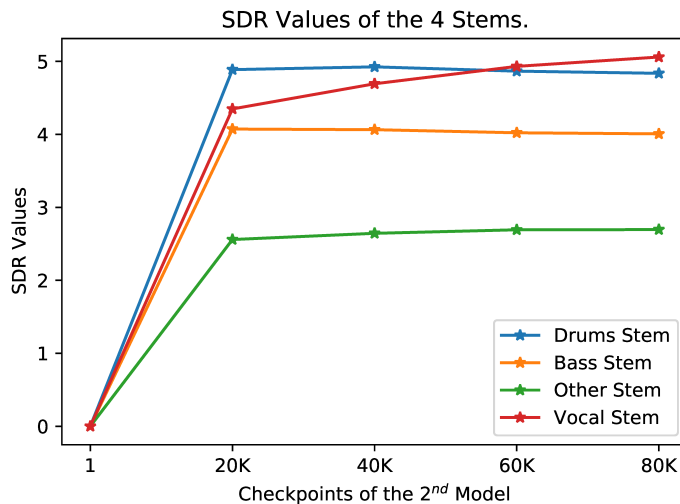


Fig. 2: SDR results of the 4 audio signal stems for the second experiment.

## 4 Conclusion

We demonstrate how transfer learning can be used for a problem of audio signal processing, in particular for the separation of an audio signal from a single mixed audio channel into four different stems: 'drums', 'bass', 'vocals', and 'other'. We show that it is possible to be successful with a small data set and relatively short training time on just two GPUs by fine-tuning pre-trained weights of *Jukebox* [3]. Similar results could not be achieved in a comparable timeframe when training from scratch or without fine-tuning, showing the potential for reduced training

times and improving results through the use of transfer learning in the audio domain.

## References

1. Défossez, A., Usunier, N., Bottou, L., Bach, F.R.: Demucs: Deep extractor for music sources with extra unlabeled data remixed. CoRR **abs/1909.01174** (2019), <http://arxiv.org/abs/1909.01174>
2. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. CoRR **abs/1810.04805** (2018), <http://arxiv.org/abs/1810.04805>
3. Dhariwal, P., Jun, H., Payne, C., Kim, J.W., Radford, A., Sutskever, I.: Jukebox: A generative model for music (2020)
4. Gao, Y., Mosalam, K.M.: Deep transfer learning for image-based structural damage recognition. *Computer-Aided Civil and Infrastructure Engineering* **33**(9), 748–768 (2018). <https://doi.org/https://doi.org/10.1111/mice.12363>, <https://onlinelibrary.wiley.com/doi/abs/10.1111/mice.12363>
5. Han, D., Liu, Q., Fan, W.: A new image classification method using cnn transfer learning and web data augmentation. *Expert Systems with Applications* **95**, 43–56 (2018). <https://doi.org/https://doi.org/10.1016/j.eswa.2017.11.028>, <https://www.sciencedirect.com/science/article/pii/S0957417417307844>
6. Hennequin, R., Khlif, A., Voituret, F., Moussallam, M.: Spleeter: a fast and efficient music source separation tool with pre-trained models. *Journal of Open Source Software* **5**(50), 2154 (2020). <https://doi.org/10.21105/joss.02154>, <https://doi.org/10.21105/joss.02154>
7. Lim, H., Kim, M.J., Kim, H.: Cross-acoustic transfer learning for sound event classification. In: 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 2504–2508 (2016). <https://doi.org/10.1109/ICASSP.2016.7472128>
8. Melnik, A., Hairston, W.D., Ferris, D.P., König, P.: Eeg correlates of sensorimotor processing: independent components involved in sensory and motor processing. *Scientific Reports* **7**(1), 1–15 (2017)
9. Melnik, A., Harter, A., Limberg, C., Rana, K., Sünderhauf, N., Ritter, H.: Critic guided segmentation of rewarding objects in first-person views. In: German Conference on Artificial Intelligence (Künstliche Intelligenz). pp. 338–348. Springer (2021)
10. Melnik, A., Legkov, P., Izdebski, K., Kärcher, S.M., Hairston, W.D., Ferris, D.P., König, P.: Systems, subjects, sessions: to what extent do these factors influence eeg data? *Frontiers in human neuroscience* **11**, 150 (2017)
11. Mitsufuji, Y., Fabbro, G., Uhlich, S., Stöter, F.R.: Music demixing challenge 2021 (2021). <https://doi.org/10.48550/ARXIV.2108.13559>, <https://arxiv.org/abs/2108.13559>
12. van den Oord, A., Vinyals, O., Kavukcuoglu, K.: Neural discrete representation learning. CoRR **abs/1711.00937** (2017), <http://arxiv.org/abs/1711.00937>
13. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. CoRR **abs/1910.10683** (2019), <http://arxiv.org/abs/1910.10683>
14. Rafii, Z., Liutkus, A., Stöter, F.R., Mimitakis, S.I., Bitner, R.: MUSDB18-HQ - an uncompressed version of

- musdb18 (Dec 2019). <https://doi.org/10.5281/zenodo.3338373>, <https://doi.org/10.5281/zenodo.3338373>
15. Sawata, R., Uhlich, S., Takahashi, S., Mitsufuji, Y.: All for one and one for all: Improving music separation by bridging networks (2021)
  16. Shor, J., Jansen, A., Maor, R., Lang, O., Tuval, O., Quitry, F.d.C., Tagliasacchi, M., Shavitt, I., Emanuel, D., Haviv, Y.: Towards learning a universal non-semantic representation of speech. *Interspeech 2020* (Oct 2020). <https://doi.org/10.21437/interspeech.2020-1242>, <http://dx.doi.org/10.21437/Interspeech.2020-1242>
  17. Stoller, D., Ewert, S., Dixon, S.: Wave-u-net: A multi-scale neural network for end-to-end audio source separation (2018)
  18. Stöter, F.R., Uhlich, S., Liutkus, A., Mitsufuji, Y.: Open-unmix - a reference implementation for music source separation. *Journal of Open Source Software* **4**(41), 1667 (2019). <https://doi.org/10.21105/joss.01667>, <https://doi.org/10.21105/joss.01667>