



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume: 8    Issue: XI    Month of publication: November 2020**

**DOI: <https://doi.org/10.22214/ijraset.2020.32135>**

**[www.ijraset.com](http://www.ijraset.com)**

**Call:  08813907089**

**E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)**

# Generating Instrumental Cover from Songs

Ali Abbas Rizvi<sup>1</sup>, Srishti Prasad<sup>2</sup>, Bhavesh Singh<sup>3</sup>, Purva Raut<sup>4</sup>

<sup>1, 2, 3</sup>Student, Information Technology, <sup>4</sup>Associate HOD, Information Technology, Dwarkadas J. Sanghvi College of Engineering - Mumbai

**Abstract:** In today's digital world and the fast-growing musical instrument technology, the ability to adapt the songs to any particular instrumental sound, becomes a priceless possession for musicians. With numerous existing tools which cater to remixing original music with instruments, a tool to convert human lyrics into instrumental rhythms is a need of the hour, or music composers of background song tracks. Often, music students start with listening to popular songs, and then attempt to recreate them through instruments. This paper aims at filtering out the human vocals of songs using a library Spleeter [14], and turning them into instrumental versions, creating an ideal platform for young talented music enthusiasts to explore and learn to play various instruments. The vocals are converted to covers using the Differentiable Digital Signal Processing (DDSP) library [13] and pre-trained Convolutional Neural Networks (CNNs) - Resnet101 [15], which enables direct integration of signal processing elements with deep learning techniques. Currently in the music industry, after recording original songs with lyrics, they are re-recorded for covers. The proposed model can also help music professionals and artists in turning existing songs to a particular instrumental cover to suit their need and also allows them to release various versions of their songs, without the burden of extra effort and money.

**Keywords:** Differentiable Digital Signal Processing, Spleeter, Librosa, Vocals, Convolutional Neural Networks, Resnet101, Lyrics, Instrumental Rhythms.

## I. INTRODUCTION

As an art form, music plays an intrinsic part in people's lives. Not only for entertainment, but there are numerous instances where learning music can affect the intelligence in children. And there are many more benefits apart from this. Music as a discipline is available in schools to help young students develop psychological strengths, cognitive attributes like expression, and soft skills simultaneously. This is highly noticeable during the instrumental music sessions in classrooms.

Recently, the Modern Era has seen a tempestuous period with change in taste and style of songs. Majority of the modern "art-music" composers have traversed the path of unconventional sounds, and have located their music in terms of tone and texture of the instrumental beats, skipping the more regular and established features of harmony and melody in music.

The 20th Century music is an aesthetic stance underlying the period of change and development, demanding introduction of some advanced characteristics, which were not always present, which includes fewer lyrical melodies than other periods, and fun instrumental rhythms.

Unfortunately, the majority of the music students who think of taking up an instrument and learning, often give-up midway, because of the scarce availability of sources which allow them to hear and understand the beats and thus, to play instruments. There are lots of rules to learn, requiring an ample amount of time and no matter how young you start, there's always someone more prodigious than you.

Thus, there is a need to teach beginners to play a particular instrument in order to master it independently, irrespective of whether one excels in understanding notes or not. Due to the lack of instrument learning platforms and the parochial availability of song covers for certain types of instruments, it becomes necessary to develop an inexpensive system which would allow the young students to acquire diverse instrumental covers with ease.

The current system of creating covers involves professional musicians to manually understand the beats and play them on the various instruments. This process is more of a trial and error process where musicians continue to improve their covers until they sound perfect.

This is a very time-consuming process and also requires the musician to possess the instruments to play them and thus, to avoid this hard work there is a need for a low-cost, easily accessible system, that can automatically generate the covers not just for a single instrument, but for all possible instruments.

## II. LITERATURE SURVEY

Separating vocals from songs with the help of a Convolutional Neural Network [8] is a process where Andrew J. R. Simpson et al. trained a convolutional Deep Neural Network (DNN) to output probabilistic approximations of the absolute binary mask for extraction of vocals. The audio signals from every song were categorized as vocal or non-vocal. All these signals are sampled at 44 kHz which they transform to spectrograms using Short-Time Fourier Transform (STFT) using certain parameters. A binary mask is calculated with the help of spectrograms of the sound signal, where each component of the mask is found by contrasting the magnitudes of the corresponding component. This is done by allocating the mask a '0' if the vocal spectrogram has lesser magnitude and '1' otherwise. Now these spectrograms are divided into various windows and hence they have a huge dataset which they send to feed-forward deep neural networks. The biased-sigmoid activation function is used throughout the DNN with no bias for the output layer. The DNN is prepared by applying some parameter  $k$  of iterations of stochastic gradient descent [10]. The model is then utilized as a feed-forward probabilistic system after the training process.

Another system is the music/voice separation using a similarity matrix [9] which is an improvisation over the traditional approach of extracting vocals from music. In the conventional approach the principle was to distinguish the musical component from the vocal component in a musical mix, by simply separating the primary recurrent structure which assumed that the background had certain repeated patterns over specific patterns. Hence, Zafar Rafii et al. have proposed a generalised approach for doing so considering redundancies occur irregularly or without a fixed period, consequently permitting the handling of music pieces with quick differing recurring structures and remote repeating components. In this approach they have used a similarity matrix which is a 2-D representation in which each point is used to find the dissimilarity amongst any component pairs of a sequence. As the recurrent patterns are responsible for creating the music structure, a similarity matrix derived from a sound signal might assist in revealing the musical structure that is present within it, which is later used to produce spectrograms to identify repeating elements.

R. Hennequin et al. have managed to produce a remarkable library Spleeter [1][14] which is an efficient tool for music source separation from a given soundtrack. It contains pretrained models which are based on TensorFlow and can separate audio files into 2, 4 or 5 stems. They have also claimed that Spleeter is extremely fast and have backed their statements with statistical data where they state that it can separate a mix audio file into 4 stems about 100 times faster than real-time on a single GPU using the pre-trained 4-stems model. With the help of this open source library published by R. Hennequin et al., vocal extraction can be easily carried out in this use case.

Human understanding is affected by general patterns as well as in-depth waveform coherence. Due to this, efficient audio synthesis becomes a fundamentally complex machine learning task. Autoregressive models, like WaveNet [4][5][6], replicate local structure although they present slow repetitive sampling and show a paucity of global hidden structure [2]. Global hidden conditioning and systematic parallel sampling, on the other hand, are present in Generative Adversarial Networks (GAN), however, GANs scuffle in producing locally-coherent sound waves. Through replication of log magnitudes and expeditious frequencies with adequate frequency resolution in the spectral domain, Jesse Engel et al. in their paper demonstrate that GANs can produce audio which has a high degree of accuracy. After thorough tests carried out on the NSynth dataset, they portray that GANs can perform far better than strong WaveNet [4] [5] [6] standards on automatic as well as human assessment criteria, and produce sound multiple times quicker.

Jesse Engel et al. have also proposed a novel system named Differentiable Digital Signal Processing (DDSP) which helps to combine traditional signal processing elements with deep learning methods [3][12][13]. In simpler terms, they have put forward an easily understandable and modular approach to generative modelling without losing out on the advantages of deep learning. They have proposed a complete system which links Digital Signal Processing to deep learning, retaining the benefits of strong inductive biases without sacrificing the expressive power of neural networks. Their library can be used to accommodate deep learning techniques for producing the instrumental covers from the vocals extracted using Spleeter.

## III. PROPOSED METHODOLOGY

The Proposed System can be broadly divided into two modules, the first being the Vocal Extraction and the second Vocal Conversion. In this section, these modules are described in detail.

### A. Vocal Extraction

In this step, the aim is to extract the human vocals from the input song with a minimal amount of noise and background notes. The system uses Spleeter [1][14], which is a deep learning based library for extracting various components of songs separately. These components include the vocals, different musical notes and noises.

Once the human vocals are extracted using Spleeter, the next task is to remove the zero frequency amplitudes from the extracted vocals. This is important because the songs may or may not consist of vocals throughout its duration. Thus, for efficient computation, the vocals at zero frequency amplitudes need to be trimmed. This can be carried out by the use of Notch Filter (band-stop filter). In signal processing, a band-stop filter or band-rejection filter is a filter that passes most frequencies unaltered, but attenuates those in a specific range to very low levels. Thus, these filters are used to overcome the problems which provide a better computational efficiency. Once the vocal is extracted it is divided into various sub parts, and spectrograms are generated which are passed to the next module.

### B. Vocal Conversion

This module comprises two different pre-trained Convolutional Neural Network (CNN) architectures used for converting the song into its instrumental cover [8]. CNNs are a special subset of neural networks which are specifically used and designed for extracting important features from images. A Convolution Neural Network is a deep neural network which has multiple layers and takes in an input image, extracts the important features from it which can either be used for classifying the image or for some other purpose. In the convolution layer, various features maps are mapped on the input image and convolution operation is applied. Feature maps are just filter matrices which are responsible for extracting different features from an input image. These filter maps are learned automatically by the network to extract important features from the image. The Pooling layers section would reduce the number of parameters when the images are too large but retain important information. Pooling can be of different types such as Average, Max or Sum Pooling. After this step, the two-dimensional image vector is flattened into a single large vector which is passed through dense layers. Dense layers are a series of linear neural layers. A Deep Convolution Neural is the one which consists of many layers of convolution, pooling together one after the other followed by dense layers.

The first CNN model is a Resnet101 model [15] with pre-trained weights. The weights are further modified by training it on the Nsynth dataset [16] of a violin. Once the system works for the violin, it can be generalized for other instruments as well. These violin notes are divided into subparts and spectrograms are then created. Spectrograms are just image representations of these notes [11]. These images of spectrograms are passed through the resnet101 model which has to learn the task of classifying the instrument as a violin. This step is carried out beforehand so that this custom pre-trained model can be accommodated in the system.

At the time of training, the system receives a band of spectrograms for a given vocal song. This band of spectrograms are iterated over continuously until they convert into the desired instrumental spectrogram. Predefined sine and cosine noises are added alternatively to these spectrograms which are then passed to a CNN model for extracting its image representations. These image representations are passed through the custom pretrained resnet101 model to classify it as a violin. If it doesn't classify it as a violin spectrogram, then some noise is added and the same procedure is iterated repeatedly until the pretrained model classifies it as a violin.

After all the entire vocal is converted into instrumental spectrograms, the spectrograms are converted back into the song thereby yielding the desired outputs. Figure 1 shows a complete flow architecture of the proposed System.

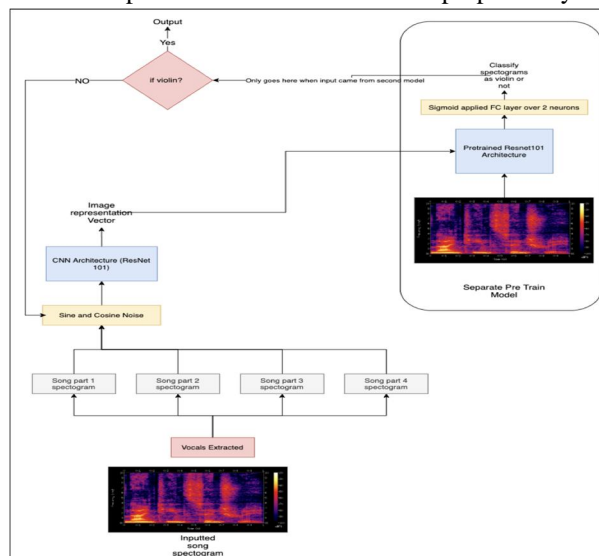


Figure 1. System Workflow

#### IV. CONCLUSION

The voice signals of the song are easily distinguishable from the music track using Spleeter [1][14] and hence generating the cover from these vocals is possible through DDSP [3][12][13] and pretrained models. Thus, on completion, this system will act as an instrumental cover generator which will readily produce covers for a profusion of instruments, making covers easily accessible to those who rely on it for various purposes. It will have the following effects:

- A. This system will remove the dependency of the instrument learning process of callow students on musicians.
- B. This will enable game developers, YouTubers and yogis to more readily use instrumental covers in their games, videos and yoga sessions respectively.
- C. Thus, the system will not only boost the availability of instrument covers but also provide a platform for inexperienced learners to continue growing independently.

#### REFERENCES

- [1] L. Pr  t, R. Hennequin, J. Royo-Letelier and A. Vaglio, "Singing Voice Separation: A Study on Training Data," ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Brighton, United Kingdom, 2019, pp. 506-510.
- [2] Jesse Engel and Kumar Krishna Agrawal and Shuo Chen and Ishaan Gulrajani and Chris Donahue and Adam Roberts, "GANSynth: Adversarial Neural Audio Synthesis".
- [3] Jesse Engel and Lamtharn Hantrakul and Chenjie Gu and Adam Roberts, "DDSP: Differentiable Digital Signal Processing".
- [4] Aaron van den Oord and Sander Dieleman and Heiga Zen and Karen Simonyan and Oriol Vinyals and Alex Graves and Nal Kalchbrenner and Andrew Senior and Koray Kavukcuoglu, "WaveNet: A Generative Model for Raw Audio".
- [5] Comparison metrics taken from DeepMind's Website. {last updated: Sep 8, 2016, url : "https://deepmind.com/blog/article/wavenet-generative-model-raw-audio".}
- [6] Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Douglas Eck, Karen Simonyan, and Mohammad Norouzi. "Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders." 2017.
- [7] Book by Ian Goodfellow and Yoshua Bengio and Aaron Courville, "Deep Learning", MIT Press 2016. @book {Goodfellow-et-al-2016, title=Deep Learning, author=Ian Goodfellow and Yoshua Bengio and Aaron Courville, publisher=MIT Press, note=url : http://www.deeplearningbook.org, year=2016}
- [8] Andrew J. R. Simpson and Gerard Roma and Mark D. Plumbley, 2015 "Deep Karaoke: Extracting Vocals from Musical Mixtures Using a Convolutional Deep Neural Network"
- [9] Zafar Rafii and Bryan Pardo, "Music/Voice separation using the Similarity Matrix". In 13th International Society for Music Information Retrieval Conference (ISMIR 2012).
- [10] Referred from {last updated: Dec 21, 2017, url : "https://towardsdatascience.com/gradient-descent-algorithm-and-its-variants-10f652806a3".}
- [11] R. Kasantikul and W. Kusakunniran, "Improving Supervised Microaneurysm Segmentation using Autoencoder-Regularized Neural Network," 2018 Digital Image Computing: Techniques and Applications (DICTA), Canberra, Australia, 2018, pp. 1-7.
- [12] Jay K. Patel and E. S. Gopi, "Musical Notes Identification using Digital Signal Processing". In 3rd International Conference on Recent Trends in Computing 2015 (ICRTC-2015) Published by Elsevier B.V.
- [13] Architecture Referred from the official magenta/ddsp website {last updated: Jan 15, 2020, url : "https://magenta.tensorflow.org/ddsp".}
- [14] AI tool Spleeter {last updated : Nov 5, 2019, url : "https://www.theverge.com/2019/11/5/20949338/vocal-isolation-ai-machine-learning-deezer-spleeter-automated-open-source-tensorflow".}
- [15] He, K., Zhang, X., Ren, S., and Sun, J., "Deep Residual Learning for Image Recognition", <i>arXiv e-prints</i>, 2015.
- [16] Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Douglas Eck, Karen Simonyan, and Mohammad Norouzi. "Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders." 2017.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)