# Vocal Separation from Audio

## 1. Introduction

### 1.1 Background

Vocal separation, or source separation, is a critical problem in audio signal processing that involves isolating vocal tracks from mixed audio signals. This task is relevant in various applications, including music production, audio enhancement, and karaoke systems.

### 1.2 Objective

This report aims to elucidate the engineering mathematical principles used for vocal separation from audio signals, with a focus on the techniques implemented by vocalremover.org. It covers the theoretical foundations, mathematical formulations, and practical implementation of vocal separation algorithms.

## 2. Theoretical Background

### 2.1 Audio Signal Representation

Audio signals are typically represented as a time-series of amplitude values. In the frequency domain, they are analyzed using techniques like the Short-Time Fourier Transform (STFT) to understand the frequency components over time.

### 2.2 Vocal Separation Techniques

Several techniques have been developed for vocal separation, including:

- **Independent Component Analysis (ICA)**
- **Non-Negative Matrix Factorization (NMF)**
- **Deep Learning Approaches**

## 3. Mathematical Formulations

### 3.1 Short-Time Fourier Transform (STFT)

The STFT is a tool for analyzing non-stationary signals. It is defined as:

$$X(t,f) = \int_{-\infty}^{\infty} x(\tau)\omega(\tau - t)e^{-2\pi f\tau}d\tau$$

Where $x(\tau)$ is the audio signal, $\omega(\tau - t)$ is the window function, t is the time index, and f is the frequency index.

### 3.2 Independent Component Analysis (ICA)

ICA is used for separating mixed signals into statistically independent components. It is based on the assumption that the observed mixed signals x are linear mixtures of source signals s: x=As where A is the mixing matrix. ICA estimates the unmixing matrix W such that: s=Wx.

### 3.3 Non-Negative Matrix Factorization (NMF)

NMF decomposes a matrix into non-negative factors. For vocal separation, it is applied to the magnitude spectrogram V: V≈WH where W and H are non-negative matrices representing the basis and activation matrices, respectively.

### 3.4 Deep Learning Approaches

Recent advancements leverage deep learning for vocal separation. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) are commonly used for this purpose, leveraging large datasets to train models that can effectively separate vocals from music.

# 4. Implementation

### 4.1 Preprocessing

- **Resampling**: Audio signals are resampled to a standard rate to ensure consistency.
- **Normalization**: Signals are normalized to a common amplitude range.

### 4.2 Algorithm Implementation

- **STFT**: Apply STFT to obtain the time-frequency representation of the audio signal.
- **ICA/NMF**: Implement ICA or NMF algorithms to decompose the spectrogram into separate components.
- **Deep Learning Model**: Train and apply a deep learning model to predict vocal and accompaniment components.

### 4.3 Postprocessing

- **Reconstruction**: Combine the separated components back into time-domain signals using inverse STFT.
- **Smoothing**: Apply postprocessing filters to enhance the quality of the separated vocals.