

---

# BARNYARD KARAOKE: RECREATING MELODIES WITH ANIMAL SOUNDS

---

**Patricio Ovalle**  
Universitat Pompeu Fabra  
Barcelona, Spain  
patricio.ovalle01@estudiant.upf.edu

## 1 Introduction

Audio mosaicing refers to the technique of reconstructing a target audio using only material from a collection of source audio. This project explores using audio mosaicing as a means to orchestrate new arrangements of popular melodies using animal sounds. The goal is to produce audio that evokes a medley of barn animals howling in the yard for a carefree evening of karaoke.

The idea takes inspiration from *Meowify*<sup>1</sup>, a project that replaces singing voice with cats meowing. Their method consists of two steps: first, a source separation model is used to isolate the vocals from the rest of the mix; then, a timbre transfer model is used to morph the singing voice into a cat's meow before recombining the modified vocals with the rest of the mix. This project attempts to achieve a similar result via the technique of audio mosaicing. The code for this project is openly available online<sup>2</sup>.

## 2 Methodology

Audio mosaicing consists of two distinct components: a target audio and a set of source audio. To decide how the source audio should be split and recombined to reconstruct the target, several key decisions must be made:

- the characteristics to emulate in a target sound
- the selection of a target sound
- the selection of a set of source audio
- a method for selecting source audio frames
- a method for recombining selected frames

The following sections address each of these concerns and describe the strategies used at each step.

### 2.1 Target characteristics

The first step to define the characteristics of the target sound we want to preserve. Our aim is to reconstruct a melody, so the pitches and pitch onset times of the melody in the target audio are the most essential features to extract.

To analyze the pitch contour of the predominant melody from the audio, we use the MELODIA algorithm introduced by Salamon and Gómez [2012] and implemented in the Essentia audio analysis toolkit [Bogdanov et al., 2013]. To subsequently determine the pitch onset times of the melody, we use the pitch contour segmentation algorithm in Essentia introduced by McNab et al. [1995]. This segmentation algorithm is especially convenient for our needs as it takes a pitch contour as input, which is the output of the MELODIA algorithm. A downside of this algorithm is that it requires a fine-tuning of several parameters in order to function well. For its "hopSize" and "minDuration" parameters we choose

---

<sup>1</sup><https://github.com/gulnazaki/meowify>

<sup>2</sup><https://github.com/p3zo/barnyard-karaoke>

relatively small values to ensure it captures the rhythmic nuance of the melodies. Specifically, we derive the minimum note duration using the tempo of the target tracks together with the shortest note values used in their melodies. For a track at 90 beats per minute (bpm), the duration in seconds for an eighth note is given by:

$$\frac{60 \text{ seconds}}{90 \text{ bpm}} * \frac{1}{2} \text{ beat} = .33 \text{ seconds}$$

For each track we set the "minDuration" parameter to just below this value to account for some variation in tempo. We set the "rmsThreshold" parameter to be low as well because we deem over-detecting onsets to be better than under-detecting them for this application. The result of extraneous onsets is that more than one source frame will be used to reconstruct a given melodic pitch in the target. On the other hand the result of a missing onset is silence, which makes for a less exciting output. An example of an estimated pitch contour and its detected onsets are shown in Figures 1 and 2.

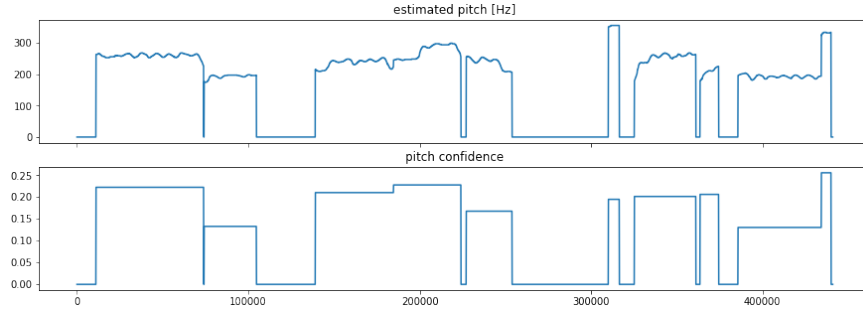


Figure 1: A pitch contour estimated using MELODIA



Figure 2: Original audio with estimated pitch onset time markers

## 2.2 Target sound selection

The next concern is to select the target sound for which we want to emulate the pitch contours. Our aim is to recreate melodies, so we choose musical tracks with well-known melodies as targets. In particular, we choose two tracks with simple vocal melodies that are easy to sing along to. The tracks are listed in Table 1 along with the timestamps used from each.

An important consideration is to select a sound with not too much complexity in the mixture to increase the effectiveness of a general pitch-tracking algorithm for melody extraction. We selected only target sounds in which the melodies were predominant and in a distinct pitch range from other instruments in the arrangement.

## 2.3 Source sound collection

The source collection in audio mosaicing is the set of audio from which frames can be sliced and used as material to recreate the melody of the target. Our aim is to create a set of animal sounds that are well-pitched and clearly recognizable.

Youtube ID	Title	Timestamp
V1bFr2SWP1I	Somewhere Over the Rainbow	1:05 - 1:15
_6HzoUcx3eo	Old Macdonald Had A Farm	0:15 - 0:35

Table 1: Target songs

We download our collection of audio samples from Freesound<sup>3</sup> using their API. The parameters used in the Freesound API queries are shown in Table 2. We select a set of animals that produce sounds in different registers overlapping with the range of the human voice. We also select animals with distinct timbres from one another in hopes of adding a bit of dynamism to the result. Each API query included a sort condition by descending user rating to get the highest-rated sounds first. In a manual exploration of downloaded results we noticed that shorter recordings were typically of higher quality than longer recordings, which often tended to be field recordings clouded with background noise and various non-animal sounds. Each collection was filtered to a short duration in order to increase the chance that the query would return single-shot sounds. The moos and bleats were allowed to be slightly longer because it seems that cows and sheep take longer to vocalize their thoughts.

For comparison, we also create a second source collection of tuned, harmonic sounds. Specifically, we download a set of violin notes from Freesound using the "ac\_single\_event" query parameter to specify that the audio should contain a single event, i.e. a single note played on the violin.

## 2.4 Frame selection strategy

When selecting which frames to use from the source to recreate the desired characteristics of the target audio, a similarity score is computed between each frame in the target audio and all of the frames in the source sound collection. Among the top 10 highest scoring source frames, we choose one randomly so as to add a bit more variation in the result.

We experimented with several sets of features for computing this similarity, namely:

- Pitch
- Pitch, loudness
- Pitch, loudness, and MFCCs

As previously mentioned, the pitch values were estimated from the MELODIA algorithm in Essentia. The loudness and MFCC features were extracted using built-in algorithms in Essentia. The goal with including loudness was to mimic the energy of the target melody, and the intent of using MFCC features was to capture the timbre of the instrument playing the melody. We computed similarity with and without MFCCs with the expectation that timbre would be less important than pitch in producing a convincing result. We also experimented with using pitch alone to see if omitting the other features would improve the general accuracy of pitch matching.

## 2.5 Recombination strategy

The final step in the audio mosaicing process is to define a strategy for how the selected source frames will be combined to reconstruct the target audio. We choose to recombine the selected frames at the pitch onset times determined by the segmentation algorithm. We elect to not allow the overlapping of source frames during recombination.

## 3 Results

Demos of both the animal and violin melody reconstructions can be heard at:

<https://github.com/p3zo/barnyard-karaoke/tree/main/demo>.

## 4 Discussion

Overall, the methodology described in this paper has a ways to go to produce the animal symphonies originally envisioned. The biggest improvement may come from assembling a well-curated sound collection. One glaring issue in

---

<sup>3</sup><https://freesound.org>

Count	Keyword	Max duration (seconds)
5	horse whinny	5
5	horse neighing	5
10	cat meow	5
20	dog bark	1
20	bleat	10
20	cow moo	10
20	bird chirp	10

Table 2: Query parameters used in Freesound API requests

the current animal sound collection is that many of the returned query results are actually sounds that are orthogonal to the query. For example, the query for "cat meow" returns results including the sound of pouring food in a cat bowl. Manual verification of the source sounds to remove non-animal sounds could go a long way. Additionally, there is a high degree of variability in the audio quality of the downloaded sounds. It would be helpful to have a few attributes related to audio quality to filter on when querying Freesound. For example a "signal-to-noise ratio" attribute would help filter out sounds with a lot of background noise. A "harmonic" attribute would be helpful in determining whether a sound is suitable to be used to match any pitch.

Another major improvement could come by using octave reduction. Some source frames match pitches in the target melodies but in a different octave. Labeling each source frame with its pitch class rather than its frequency would increase the coverage of pitch classes that can be matched exactly.

While the pitch detection algorithm used worked well on the selected target sounds, there are many complex mixtures for which it does not work well. For these cases, a source separation pre-processing step can be used to isolate the instrument playing the predominant melody, so as to simplify the input for pitch detection.

## References

- Justin Salamon and Emilia Gómez. Melody extraction from polyphonic music signals using pitch contour characteristics. *IEEE transactions on audio, speech, and language processing*, 20(6):1759–1770, 2012.
- Dmitry Bogdanov, Nicolas Wack, Emilia Gómez Gutiérrez, Sankalp Gulati, Herrera Boyer, Oscar Mayor, Gerard Roma Trepas, Justin Salamon, José Ricardo Zapata González, Xavier Serra, et al. Essentia: An audio analysis library for music information retrieval. In *Britto A, Gouyon F, Dixon S, editors. 14th Conference of the International Society for Music Information Retrieval (ISMIR); 2013 Nov 4-8; Curitiba, Brazil.[place unknown]: ISMIR; 2013. p. 493-8.* International Society for Music Information Retrieval (ISMIR), 2013.
- Rodger J McNab, Lloyd A Smith, and Ian H Witten. Signal processing for melody transcription. 1995.