

# Zero-shot Crate Digging → DJ Tool retrieval using Speech Activity, Music Structure and CLAP embeddings

Iroro Orife

## What are DJ Tools?

In genres like Hip-Hop, RnB, Reggae, Dancehall and every Electronic/Dance/Club style, DJ tools are a special set of audio files of short, simplified musical phrases retrieved, curated from existing music, or composed specifically to heighten the DJ's musical performance and creative mixing choices.

- ▶ Acapella loops, vocal chants, spoken word, one shot vocal samples
- ▶ Purely instrumental (no drums) piano, guitar, strings, horns
- ▶ Drum loops, solo percussion, drum breaks
- ▶ Sound effects: risers, sweeps, sirens, vinyl fx,
- ▶ Scratch and Battle loops

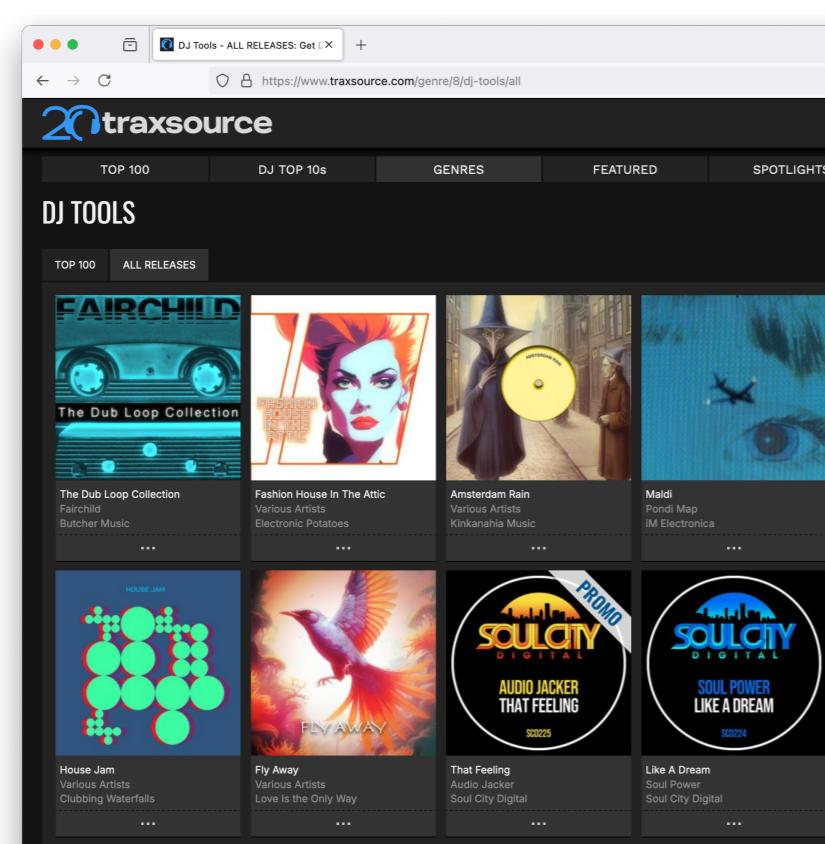


Figure 1: Traxsource DJ Tools shop

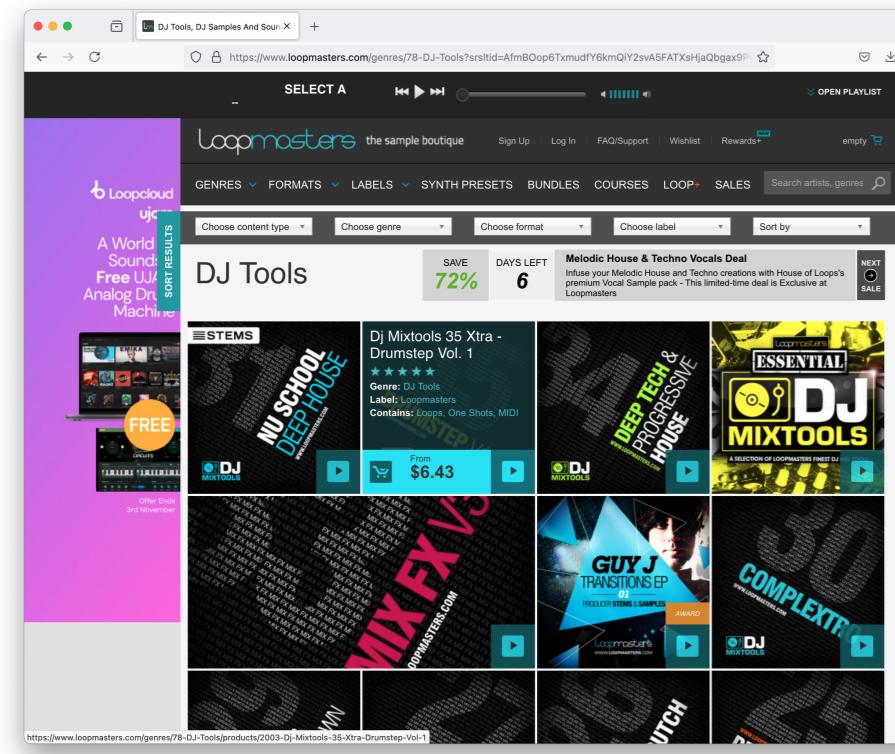


Figure 2: Loopmasters online shop. Note the specification of Loops, One-shots

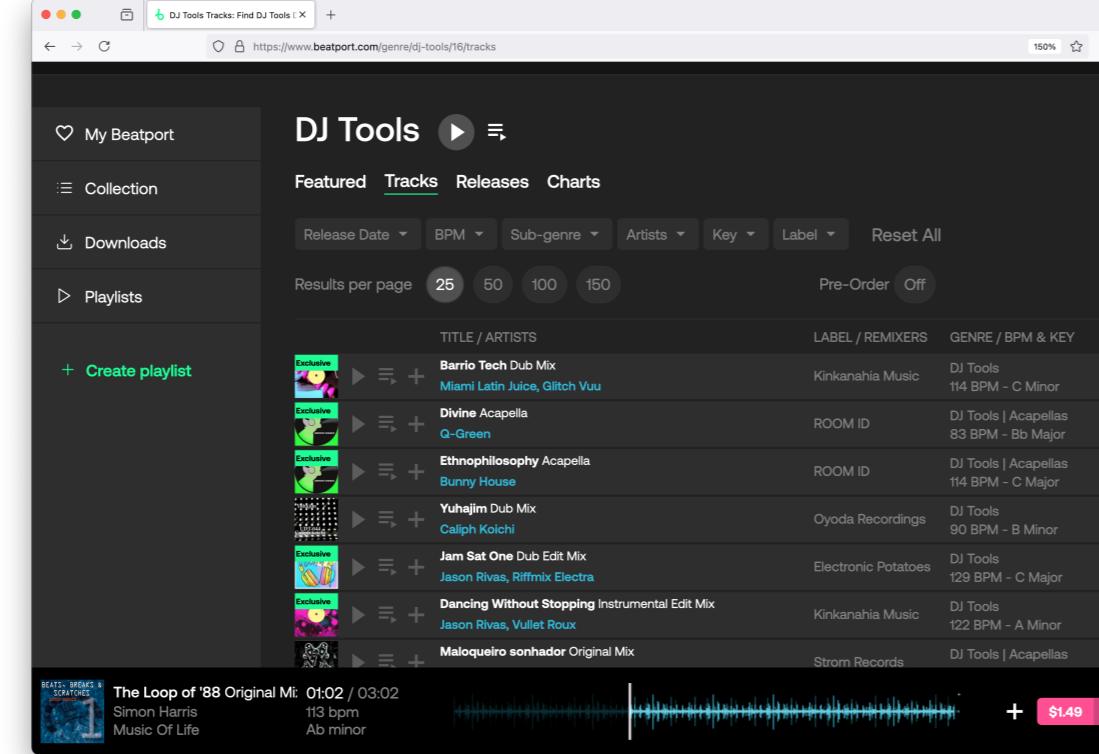


Figure 3: Beatport DJ Tools selection. Note the designated Genre, Key and BPM metadata

## Software Demo Implementation

The DJ-tool classification system and its dependencies are fully open-source software (OSS) written in Python. Firstly, we run MSAF<sup>a</sup> and SMAD<sup>b</sup> to generate raw activities and structural boundaries, saved as CSV files[2, 3]. Next, these are post-processed to yield lists of segment time ranges. Figure 5 shows the relationship between the boundaries and the speech activity for one song.

For zero-shot classification, we utilize the LAION version of CLAP<sup>c</sup> hosted on Huggingface[1, 4]. Finally, we manually create a list of descriptive text strings for each class, as listed in Table 1.

### Algorithm 1: Zero-shot Crate Digging

```
1 Create text prompts for  $M$  classes  $\{X_1^t, \dots, X_M^t\}$ 
2  $\{E_1^t, \dots, E_M^t\} \leftarrow \text{TextEncoder}(\{X_1^t, \dots, X_M^t\})$ 
3 for song  $s_i$  in the library  $S_N$  do
4    $W^s, W^m \leftarrow \text{SMAD}(s_i)$ 
5    $B \leftarrow \text{MSAF}(s_i)$ 
6   Adjust boundaries  $B$  to nearest  $W^s$ 
7   Use  $B$  to cut  $s_i$  to audio files  $\{X_1^a, \dots, X_N^a\}$ 
8   for  $j := 1$  to  $N$  do
9      $E_j^a \leftarrow \text{AudioEncoder}(X_j^a)$ 
10     $D \leftarrow \text{Similarity}(E_j^a, \{E_1^t, \dots, E_M^t\})$ 
11     $\vec{z} \leftarrow \text{softmax}(\vec{D})$ 
12    Store predicted class, the arg max of  $\vec{z}$ 
```

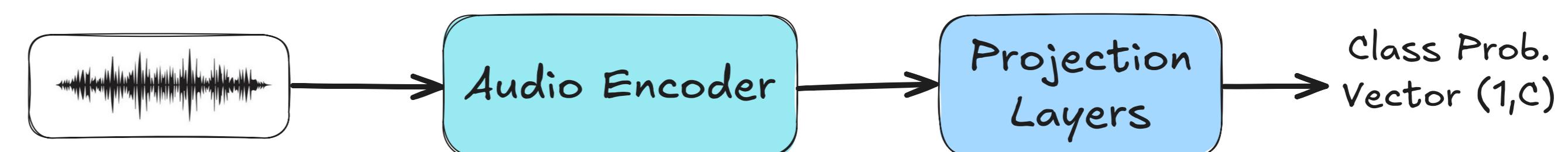
<sup>a</sup><https://github.com/urunieto/msaf>

<sup>b</sup><https://github.com/biboamy/TVSM-dataset/tree/master>

<sup>c</sup><https://huggingface.co/laion/clap-htsat-unfused>

## What is Zero-shot Audio Classification?

### Supervised Audio Classification



### Zero-shot Audio Classification

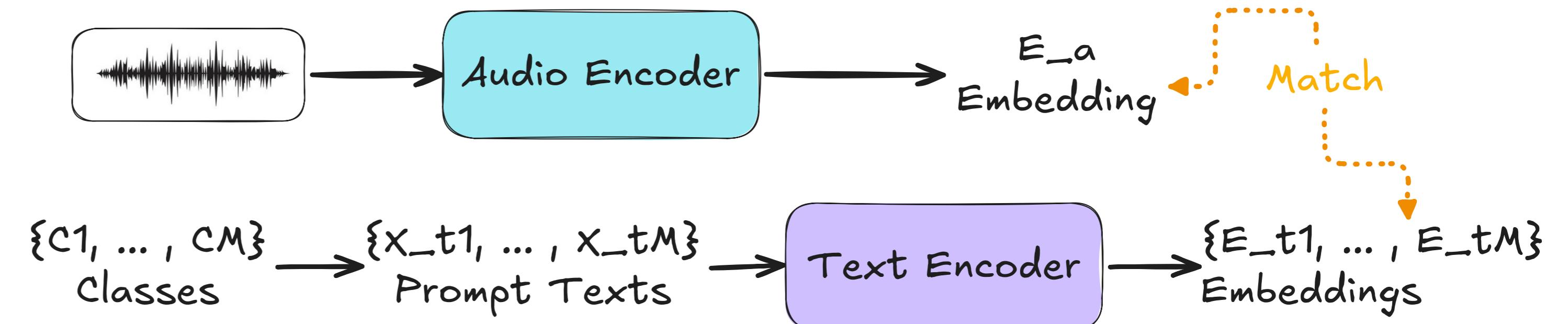


Figure 4

## Speech Activity & Music Structure Analysis

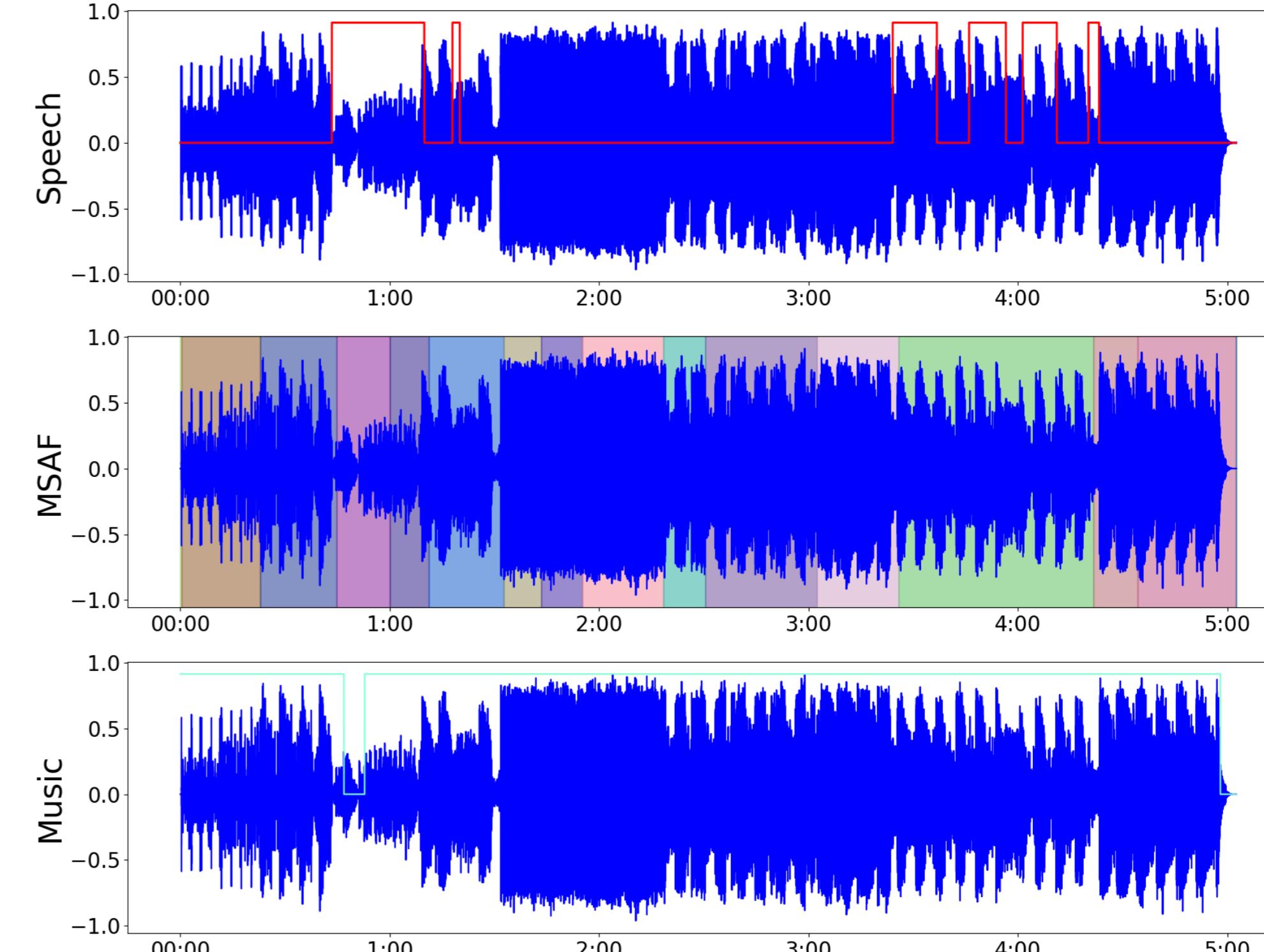


Figure 5: A 5 minute Ragga Jungle song overlaid with detected speech, music activity and music-structural boundaries

### DJ Tool class

DJ Tool class	Example text description
Acapella loops	"acapella, expressively sung vocal tracks, purely acoustic"
Sound effects	"siren, riser sound effects, whoosh, crash, synthetic"
Drums breaks	"drums, funky drummer, drum solo, Amen breakbeat"
Melodic hooks	"strings, solo guitar, piano melodies, synths instrumental"
Battle tracks	"vinyl scratch loop, turntablist DJ battle sounds"

Table 1: In practise, text descriptions are more tortuous

## Results: Making music with the system

We evaluated the system on 10+ albums from the author's DJ library →

- ▶ Vocal and drums tool classes perform better than short, burstier sound effects
- ▶ Performance is sensitive to *precise language*. Adding generic terms or entire genres (e.g. "hip-hop, drum and bass") to a class description led to markedly worse results overall.
- ▶ For best results, the class descriptions should not overlap
- ▶ One challenge was to determine optimal segmentation for CLAP classification, while slicing the song appropriately for use by DJs.
- ▶ Playing with the system and making music with the retrieved DJ tools led to unexpectedly delightful mixing sessions.

## References

- [1] B. Elizalde, S. Deshmukh, M. A. Ismail, and H. Wang. Clap: Learning audio concepts from natural language supervision, 2022.
- [2] Y.-N. Hung, C.-W. Wu, I. Orife, A. Hippel, W. Wolcott, and A. Lerch. A large tv dataset for speech and music activity detection. *EURASIP Journal on Audio, Speech, and Music Processing*, 2022(1):21, 2022.
- [3] O. Nieto and J. P. Bello. Systematic exploration of computational music structure research. In *ISMIR*, pages 547–553, 2016.
- [4] Y. Wu, K. Chen, T. Zhang, Y. Hui, T. Berg-Kirkpatrick, and S. Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2023.