

Segmentation of electronic dance music

Tim Scarfe, Wouter M. Koolen and Yuri Kalnishkan

Computer Learning Research Centre and Department of Computer Science,
Royal Holloway, University of London, Egham, Surrey, TW20 0EX, United Kingdom
{tim,wouter,yura}@cs.rhul.ac.uk

September 26, 2014

Abstract

We consider the problem of segmenting DJ-mixed dance music recordings (pod-casts, radio shows, live events) into their original sequence of contiguous tracks. We present an algorithm to reconstruct a fixed number of segment boundaries as close as possible to what a human domain expert would create in respect of the same task. As the number of segments is known in advance we do not have to rely on local points-of-change heuristics prevalent in common segmentation algorithms. We also adapt the method to estimate the number of tracks in a recording and compare our method with others. The goal of DJ-mixing is to render track boundaries invisible from human perception.

Segmentation is performed on a self-similarity matrix which is derived from normalized cosines of various cost matrices which have themselves been derived from a time-series of Fourier based spectral features. The cost matrices introduced in this paper introduce notions of self-similarity, symmetry, contiguity and evolution in respect of time.

Our work incorporates self-similarity over a soft time horizon with a fixed upper boundary and is quantitatively assessed on a large corpus of radio show recordings which have been hand-labelled by a domain expert. The method could be used on other segmentation tasks and other domains.

Keywords. music, segmentation, DJ, mix, dynamic programming

1 Introduction

Electronic Dance Music tracks are usually mixed by a disc jockey (DJ). For this reason EDM music streams are unique compared to other genres of music. Mixing is the *modus operandi* in electronic music. We first transform the audio file into a time series of features discretized into adjacent tiles and transform them into a domain where pairs from the same track would be distinguishable by their cosine. Our features are based on a Fourier transformation with convolution filtering to accentuate prominent instruments and self-similarity within tracks. We create a similarity matrix from these cosines and then derive cost matrices showing the costs of fitting a track at a given time with a given length. We use Dynamic Programming to create the cost matrices and again to perform the most economical segmentation of the cost matrices to fit a predetermined number of tracks. Dynamic programming means solutions to a problem are described in terms of overlapping sub-solutions to achieve a significant improvement in time complexity and therefore execution time.

Contiguous-segmentation differs from the standard clustering problem in that the clusters arrive sequentially and are contiguous (AAABBBCCDDDD, not AAABBBCCBBBB). In the literature this may also be known as *time-dependent* clustering [1]. For brevity we will use the term *segmentation* from now on to describe this methodology.

Our work mainly considers a segmentation that is globally-optimised and based upon the number and

self-similarity of segments, avoiding transient points-of-change methods. Music and mixes of music have the property that they are made up of recursively repeating self-similar regions. Our method does not strictly require any training or tenuous heuristics to perform well. The distinguishing feature of our problem domain is that the number of segments is known a priori but the segmentation boundaries are not known, or ambiguous and subjective. However, computing the best solution is desirable.

The intended purpose of the algorithm is to reconstruct optimal boundaries given a fixed number of tracks known a priori. The track listing is usually published by the DJ so the number of tracks is known. This is relevant when one has recorded a show (perhaps automatically), downloaded a track list and needs to reconstruct the indices given that track list. The order of the indices reconstructed is critical so that we can align the correct track names with the reconstructed indices.

We extend it to estimate the number of tracks and compare to other methods.

One of the interesting features of audio is that you *cannot scrub through it, and get an overview in the same way you can with video*. Audio has a reduced *contextual continuum* when the user skips through it, perhaps due to the lack of redundant, persistent scene-setting information or indeed a psychological reason. Even in video applications, discovery, context and scrubbing are an active area of research [2]. Time index meta-data would allow click through monetisation, and allow improved use-case scenarios (for example publishing track names to social networks, information discovery and retrieval). Capturing meta-data in audio is a time consuming and error-prone process. Tzanetakis [3] found that it took users on average 2 hours to segment 10 minutes of audio using standard tools. While not directly relevant we might glean from those findings that there is a strong motivation to automate this process.

DJs always match the speed or beats per minute (BPM) of each adjacent track during a transition and align the major percussive elements in the time domain. This is the central concept of removing any dissonance from overlapping tracks. Tracks can overlap by any amount. DJs increase adjacent track compati-

bility further by selecting adjacent pairs that are harmonically compatible and by applying spectral transformations (EQ).

1.1 Literature Review

Audio segmentation in the literature is colloquially implemented in the context of structural analysis. Music structure denotes the organization of a composition by its melody, harmony, timbre and rhythm. Repetitions, transformations and evolutions of music structure also contribute to its identity and it is this semantic information that structural analysis algorithms aim to extract from music. An example structure for a song might be “Common” → “Verse” → “Chorus” → “Common”. Speaker diarization is another example of structural analysis

Segmentation in the context of structural analysis has been thus far been concerned with creating a novelty function to find points-of-change using distance-based metrics, rather than trying to find a fixed number of segments. Heuristics with hard decision boundaries have been used to find the best change points, for example Tzanetakis [3] used first-order derivatives of a time series of audio features.

The use of a similarity matrix to visualize and analyse local time dependencies (then called “Recurrence Plots”) was first proposed by Eckmann[4].

J. Foote et al [5, 6, 7, 8, 9] were the first to use self-similarity matrices to visualise and exploit time dependencies in music data. Foote evaluated a Gaussian ‘tapered’ checker board kernel along the diagonal of a music self similarity matrix to create a 1-dimensional novelty function that had the notion of self-similarity over a fixed time horizon. The kernel was ‘tapered’ down to zero on the top right and bottom left edges to reduce edge effects. The dynamic program in this paper allows any self-similarity time horizon up to a fixed limit (Foote’s work had a fixed kernel size).

Goodwin et al. used a dynamic program for segmentation [10]. Their intriguing supervised approach was to perform Linear Discriminant Analysis (LDA) on the features to transform them into a domain where segmentation boundaries would be emphasised and the feature weights normalized. Afterwards,

Goodwin formulated the problem into one of finding the globally minimum cost path through a state graph (‘cluster space trajectory’) modelling local and transition costs. Goodwin already demonstrated in [11] that novelty peaks often exist within segments, not only on the boundary of segments and took the approach of modelling all possible sequential transitions between all possible clusters.

A potential drawback of the approach by Goodwin and all other approaches in scene analysis segmentation that we are aware of; is that they are somewhat local methods that focus on points-of-change rather than optimizing for the best possible results for a fixed number of segments.

All contiguous segmentation algorithms that we looked at do not know a priori how many segments to find. In almost all applications, the number of known segments derives from knowledge of where boundaries are. This is the reason why the algorithm presented in this paper is unique.

[?]

A distinguishing feature of our approach is that we evaluate how well we are doing compared to humans in respect of the same task. We compare our reconstructed indices to the ones created by a human domain expert and the algorithm itself is optimised for the domain of mixed music.

In the coming sections we describe the corpus (see Section 2), human accuracy (see Section 3), the evaluation criteria (see Section 9), how we pre-process the data (see Section 4.1), how we perform feature extraction (see Section 4.2.1), designing the cost matrices using observed phenomena in the domain (see Section 4.2.2), computing the best segmentation (see Section 5), discussion of confidence intervals (see Section 6), our methodology (see Section ??), materials used (see Section 13), results (see Section 10) and finally the summary (see Section 11).

2 Corpus

We have been supplied with several broadcasts from three popular radio shows. These are: Magic Island, by Roger Shah (106 shows); A State of Trance with Armin Van Buuren (109 shows); and Trance Around

The World with Above and Beyond (88 shows). The show genres are a mix of Progressive Trance, Uplifting Trance and Tech-Trance.

There are no silent gaps after the introduction. The shows come in 44100 samples per second, 16 bit stereo MP3 files sampled at 192Kbs. We resampled these to 4000Hz 16 bit mono (left+right channel) WAV files to allow us to process them faster. We have used the “Sound eXchange”¹ program to do this. These shows are all 2 hours long. The overall average track length is 5 and a half minutes (slightly less for Magic Island (see Figure 1)) and normally distributed. The average number of tracks is 23 for ASOT and TATW, 19 for Magic Island. There is a guest mix on the second half of each show. The guest mix DJs show off their skills with technically convoluted mixing, so it is fair to say that the boundary “complexity” increases during the guest mix and is at least not fixed throughout the shows.

An additional dataset of 36 radio shows have been mixed by and annotated by Mikael Lindgren (the so called *lindmik* dataset). These shows are extremely useful because the DJ is the same person that created the ground truth time indices. Also there is less noise, for example voice-overs, guest mixes, radio show sounds, introductions etc. These shows also vary significantly in length from 1 hour to nearly 5 hours. There are 339 shows in total.

We believe this corpus is the largest of its kind used in the literature going on the comparative table of segmentation corpora listed by Peiszer et al in their literature review of audio segmentation [12]. More recently Badawy et al [13] used a corpus of 61 hours. The corpus we are using is longer than 680 hours in length.

There is already a large community of people interested in getting track/time meta-data for DJ sets. “CueNation”² is an example of this. CueNation is a website allowing people to submit *cue-sheets* for popular DJ Mixes and radio shows. A cue-sheet is a text file containing time meta-data (indices) for a media file.

The three main radio shows in the corpus were

¹<http://sox.sourceforge.net>

²<http://www.cuenation.com>

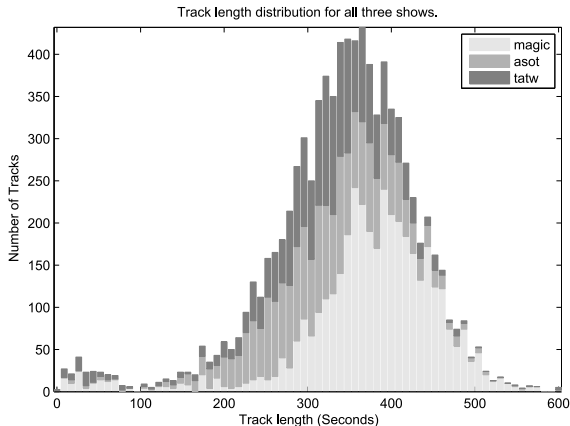


Figure 1: The track length distribution for all three radio shows. The *bump* of short tracks (less than 3 minutes) is often introductions or otherwise extraneous.

hand captured by Denis Goncharov; a domain expert and one of the principal contributors to *CueNation*. One of the significant problems with this task is that there is a random error variable associated with the human captured indices themselves. On some tracks, it is unclear where to place the optimal index and when analysing our results, we have noticed some obvious human errors. Regrettably, there is no clear way of quantifying this. Many of the cue-sheet authors themselves reject the idea of automating the task, citing the poor precision of any such result (they often place indices on the exact MP3 frame). However this sentiment seems misplaced given that they frequently make mistakes or that it is a matter of opinion where to place the track and some consistent method may be preferential. A potential outcome of our method could be an assistance mechanism to help them with initial placements. Our results demonstrate that it is indeed possible to automate this task and that while there is some uncertainty attached to the optimal placement, it is still largely predictable. Indeed on the majority of track indices the uncertainty is ostensibly small.

Denis Goncharov provided us with the following description of how he captures the indices. To quote from a personal email exchange with Denis:

Trance music is made in slices of 8 bars (1 bar is 4 beats. At 135 beats per minute, 8 bars is $\frac{60}{135}(4 \cdot 8) \approx 14.8$ sec). Trance music tends to be around 130-135 BPM. It is a matter of personal preference which point of the transition to call the index. My preference is to consider the index to be the point at which the second track becomes the focus of attention and the first track is sent to the background. Most of the time the index is the point at which the bass line (400Hz and lower) of the previous track is cut and the bass line of the second track is introduced. If the DJ decides to exchange the adjacent tracks gradually over the time instead of mixing them abruptly then it is up to the cue-sheet maker to listen further into the second track noting the musical qualities of both tracks and then go back and choose at which point the second track actually becomes the focus of attention.

The most obvious and pervasive element in dance music is the percussion (the beats). We believe on balance that ignoring the percussive information is advantageous, because DJs use percussion primarily to blur boundaries between tracks. We tried to capture percussive based features and found that the transitions between tracks and indeed groups of tracks appeared as stronger self-similar regions than the actual tracks. The percussive feature extractor transformed the autocorrelation of the audio samples in the time domain tiles, and compared the cosine of their absolute values. It was reasonably clear from that research that track boundaries are revealed with less uncertainty between instruments and harmonic content. However. We do not rule out looking at percussive features again in future research because we are currently ignoring potentially useful features.

Some DJs “mix harmonically” (by matching instruments as opposed to percussion) but this preys on human hearing and perception. An algorithm capturing the harmonic information would most likely be able to distinguish two harmonically compatible tracks.

3 Human Accuracy

We did some analysis on how accurate the humans themselves are at creating indices. In the absence of a perfect data set our analysis instead hinged on the amount to which the humans disagreed with each other aggregated over a large amount of historical data. Mikael Lindgren was kind enough to send us a dump of his cuesheet database to experiment with. As ASOT is such a popular show there were many independently captured cuesheets to compare against for all of the historical shows. We selected all the shows having at least 3 distinct cuesheets (not copies or shifted/misaligned copies of each other) and such that all the cuesheets had the same number of tracks. The first track was ignored (as it was always 0 seconds). We ended up with 115 shows with 3 authors, 65 shows with 4 authors and 30 shows with 5 authors. We generated a histogram of distances from the median time for each track, for each cuesheet and assumed values greater than 100 seconds or less than -100 seconds were outliers. The standard deviation of the ‘human disagreement’ variable is 9.13 seconds. See Figure 2 for an illustration. So at this stage it does not seem feasible for us to achieve a higher accuracy when we are evaluating against a method which is intrinsically error prone. An important caveat here is that ASOT turned out to be the most error-prone show to segment out of our corpus. The standard deviation of the bumps could be reduced if we normalized the times by the BPM of each transition.

4 Data Handling

4.1 Preprocessing

The corpus had some outliers that may have slightly distorted the analysis of our method. Many of the “tracks” in our data set (of indices) were in fact not tracks at all but rather introductions or voice-overs. Almost all of these outlier tracks were short in length. These are quite clearly visible on the distribution of track lengths on Figure 1. To ameliorate the situation we removed any tracks that were shorter than 180 seconds. We also removed any end tracks that were

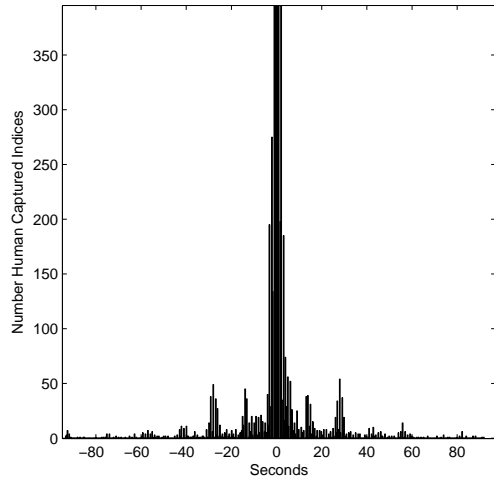


Figure 2: Illustration of the ‘human disagreement’ random variable (zoomed in at the bottom), standard deviation 9.13 seconds. Peaks are visible at intervals of 8 bars (≈ 14.8 seconds) which corroborates the analysis from Denis Goncharov in Section 2.

shorter than 240 seconds as very often the end tracks on a radio show contain strange elements (for example voice-overs, interviews, show-related ‘jingles’). This required some manipulation of the cue-sheets and audio files. The undesired segments of the audio files were chopped out, and the cue-sheets were re-flowed so that the time indices point to the correct location in the file.

The algorithm still performs similarly when removing just these indices and leaving the audio intact underneath, so it would not significantly affect any real-world implementation.

For those wishing to use this algorithm in practice with pre-recorded shows; the introductions at the start of the shows are often fixed length or at least predictable so error would be small on average.

The lindmik dataset which was noise-free did not require any preprocessing.

4.2 Feature Extraction

4.2.1 Music

We used SoX (see Sect. 2) to downsample the shows to 4000Hz. We are not particularly interested in frequencies above around and above 2000Hz because instrument harmonics become less visible in the spectrum as the frequency increases. The Nyquist theorem [15] states that the highest representable frequency is half the sampling rate, so this explains our reason to use 4000Hz. We will refer to the sample rate as R . Let L be the length of the show in samples.

Fourier analysis allows one to represent a time domain process as a set of integer oscillations of trigonometric functions. We transform the tiles into the frequency domain using the discrete Fourier transform

$$F(x_k) = X_k = \sum_{n=0}^{N-1} x_n \cdot e^{-i2\pi \frac{k}{N}n}$$

which transforms a sequence of complex numbers x_0, \dots, x_N into another sequence of complex numbers X_0, \dots, X_N where

$$e^{-i2\pi \frac{k}{N}n}$$

are points on the complex unit circle. Note that the fftw algorithm [16] that we used to perform this computation operates significantly faster when N is a power of 2 so we zero pad the input to the next power of 2. We denote the tile width by M in seconds (an algorithm parameter). Note that

$$N = \frac{L}{M}$$

denotes the tile size in samples (length of show in samples over the tile size). Let

$$T = \left\lfloor \frac{L}{\tilde{M}} \right\rfloor$$

be the total number of tiles, and

$$\tilde{M} = \frac{L}{T}$$

the tile width in samples. Because we are passing real values into the $F(x_k)$, the second half of the result is a rotational copy of the first half.

Show samples are collated into a time series Q_i^y ($T \times N$) of contiguous, non-overlapping, adjacent *tiles* of equal size where $i = 1, 2, \dots, T$. Samples at the end of the show that do not fill a complete tile get discarded. The affect of this is increasingly negligible with decreasing tile size. Since we zero-pad N to the next power of two, this also decreases the affect.

As we are not always interested in the entire range of the spectrum, we use l to represent a low pass filter (in Hz) and h the high pass filter (in Hz). So we will capture the range from h to l on the first half of the result of F . Let $\hat{h} = \lceil h \cdot \frac{N}{R} \rceil + 1$ be the position of h in the spectrum, and $\hat{l} = \lceil l \cdot \frac{N}{R} \rceil + 1$ be the position of l in the spectrum.

Let D_e^y ($T \times \hat{l} - \hat{h} + 1$) denote the feature matrix.

For each tile $\bar{i} = 1, 2, \dots, T$ we assign

$$D_{\bar{i}}^{1, \dots, \hat{l} - \hat{h} + 1} = \left| F(Q_{\bar{i}}^{1, \dots, \tilde{M}})_{\hat{h}, \hat{h}+1, \dots, \hat{l}} \right|$$

selecting the part of the spectrum between the high and low pass filters h and l . We take the absolute values of the complex result of $F(x_k)$ (defined as its distance in the complex plane from the origin using the Pythagorean theorem).

To accentuate instrument harmonics we perform convolution filtering on the feature vectors in D , using a Gaussian first derivative filter. This works like an edge detection/transient filter but also expands the width of the transients (instrument harmonics) to ensure that feature vectors from the same song appear similar because their harmonics are aligned on any distance measure (we use the cosines). This is an issue because of the extremely high frequency resolution from having such large inputs into $F(t_i)$. For example with a tile size of 10 seconds and a sample rate of 4000 we have a frequency resolution of $\frac{1}{2} \cdot 10 \cdot 4000 = 20\text{KHz}$.

Typically a ‘short-time discrete Fourier transform’ is used which has smaller sized inputs (windows) into $F(t_i)$ which are usually overlapping and are multiplied by a window function, attenuating the tails to reduce spectral leakage. Usually these window functions look similar to a Gaussian, for example;

$$\text{Hann}_i = 0.5 - 0.5 \cos \frac{2\pi i}{n-1} w(i)$$

where n is the window size (see [17] for an example). The short-time Fourier transform is relevant when increased time precision is needed as there is a frequency-time resolution trade-off with respect of the input size to $F(t_i)$. This is not a concern in this particular application as our time resolution is never required to be better than 1 second which would still produce adequate frequency resolution.

The Gaussian first derivative filter is defined as

$$-\frac{2\hat{\lambda}}{v^2}e^{-\frac{\lambda^2}{v^2}}$$

where

$$\hat{\lambda} = \{ -\lfloor 2v \rfloor, \lfloor -2v + 1 \rfloor, \dots, \lfloor 2v \rfloor \},$$

and

$$v = b \frac{N}{R}.$$

b is the bandwidth of the filter in Hz and this is a parameter of the algorithm. After the convolution filter is applied to each feature vector in D , we take the absolute values and normalize on the vector lengths.

Because the application domain is well defined in this setting, we can design features that look specifically for what we are interested in (musical instruments). Typically in the literature; algorithms use an amalgam of general purpose feature extractors. For example; spectral centroid, spectral moments, pitch, harmonicity [3]. We construct a dissimilarity matrix of cosines as is common in the literature for similar applications [5]. The cosines are computable easily because they are the inner products of the respective features (the features have been normalized to unit length).

Let

$$S_{ij} = 1 - \langle D_i, D_j \rangle,$$

define the dissimilarity matrix.

Then we apply some normalizing transformations. First we center S around 0.5 by raising each element to the power $2s$, where $s = \frac{1}{T^2} \sum_{i,j=1}^T S_{ij}$. Since for $x \in [0, 1]$ and $y > 0$ we have $x^y \leq x$ if $y \geq 1$ and $x^y \geq x$ if $y \leq 1$, the transformation $S_{ij} \rightarrow S_{ij}^{2s}$ increases the values S_{ij} whenever the mean value $s < 0.5$ and decreases them whenever $s > 0.5$. Note that the transformation keeps the values S_{ij} in the interval

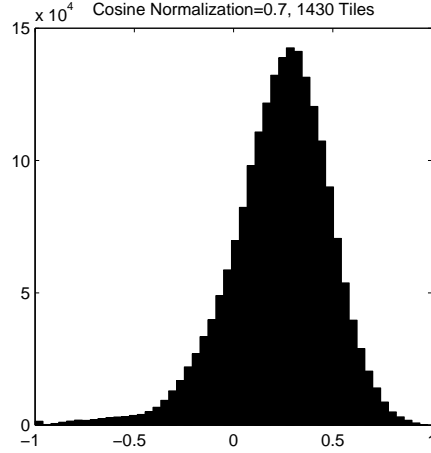


Figure 3: Illustration of the effect of normalization parameter $\hat{c} = 0.7$ on the values in S on radio show Magic Island 110. The small raised section on the left correspond to the tracks.

$[0, 1]$. We find this a convenient and gentle way to rescale S .

Secondly we raise each value S_{ij} to a power $\hat{c} \in [0.5, 1.5]$ and then rescale and translate them to $[0, 1]$ using $S_{i,j} \rightarrow 2S_{ij} - 1$. The parameter \hat{c} is tuned so as to achieve the right balance between negative *incentives* and positive *disincentives* for meaningful track placement. The distribution of values in S after the transformations will have a raised tail on left. This will become relevant when we discuss cost matrices as some of them depend on the sign of the value in S .

See Figure 4 for an illustration of S and Figure 3 for an illustration of the normalization.

4.2.2 Self Similarity

We now have a similarity matrix S_{ij} as described in Section 4.2.1.

Let w and W denote the minimum and maximum track length in seconds, these will be parameters of the algorithm that improve the time complexity while not significantly harming the results.

We proceed to constructing a cost matrix $C(f, t)$

that describes the cost of placing a track starting at f and finishing at t (and having length $t - f + 1$). After analysing the data set, we have created 7 cost matrices that exploit observed phenomena in S . We also provide an additional cost matrix which is just a Gaussian random function centred around the mean track length for all times which can be used to regularise the other 7 matrices or used on its own as a comparator to a more naive method of placement.

The cost matrices exploit themes such as contiguity, symmetry, evolution and change as well as simple summation of S as was presented in our last paper [14]. In our previous work S was on the interval $[0, 1]$ and the summation method could only consider disincentives. The new cost matrices have a parameter to shift the consideration of incentive versus disincentive and values on the interval $[-1, 1]$.

On the whole, a significant number of tile pairs within one track are similar to each other. Pairs of tiles that do not belong to the same track are expected to be dissimilar, most of the time. However, tracks have contiguous regions within them that are dissimilar to each other. Transitions between songs may appear as a self-similar region but usually also similar to each adjacent track to varying degrees.

Summation The most obvious strategy of all is to sum up all relevant tiles in S for each candidate track from tile f through tile t . We define $C(f, t)$, the cost of a candidate track from tile f through tile t , to be the sum of the similarities between all pairs of tiles inside it

$$\tilde{C}(f, t, \Omega) = \sum_{i,j=f}^t \hat{S}_{ij}$$

where

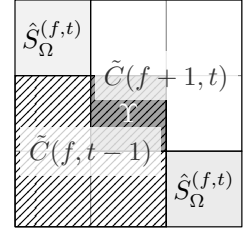
$$\hat{S}_{ij} = \begin{cases} \Omega S_{ij}, & \text{if } S_{ij} > 0 \\ (1 - \Omega) S_{ij}, & \text{otherwise} \end{cases}$$

for all $i, j \in S$ (the parameter Ω control the balance of positive and negative values). Direct computation using the definition takes $O(TW^3)$ time. We can improve this to $O(TW)$ by using the following recursion

(assume that $f + 1 \leq t - 1$):

$$\begin{aligned} \tilde{C}(f, t) &= \tilde{C}(f + 1, t) + \tilde{C}(f, t - 1) \\ &\quad - \tilde{C}(f + 1, t - 1) + \hat{S}_{ft} + \hat{S}_{tf}. \end{aligned}$$

The recursion implies that the cost of a track of length $L = t - f + 1$ can be calculated from the costs of shorter tracks using a constant number of operations. The following picture (where the middle cell corresponds to $\Upsilon = \tilde{C}(f + 1, t - 1)$) provides an illustration:



Symmetry A common feature on dance music tracks is partial mirror-symmetry. We build a cost matrix to capture that.

Let $\Lambda(f, t, d)$ be the diagonal parallel to the minor diagonal of S and at the ‘distance’ d from it. We represent it as an ordered set

$$\begin{aligned} \Lambda(f, t, d) &= \langle S_{f+d,f}, S_{f+d+1,f+1}, S_{f+d+2,f+2}, \dots, S_{t,t-d} \rangle. \end{aligned}$$

For each such diagonal in one triangle/half of S we want to compare each element against its mirror counterpart. For an ordered set Λ we define its cost as

$$\bar{C}(f, t, \Omega)(\Lambda) = \sum_{i=1}^{|\Lambda|} \delta(\Lambda_i, \Lambda_{|\Lambda|-i+1}, \Omega)$$

where

$$\begin{aligned} \delta(p, q, \Omega) &= \begin{cases} 0, & \text{if } \text{sign}(p) \neq \text{sign}(q), \\ \Omega pq, & \text{if } \text{sign}(p) \geq 0 \text{ and } \text{sign}(q) \geq 0, \\ (1 - \Omega) pq, & \text{if } \text{sign}(p) < 0 \text{ and } \text{sign}(q) < 0, \end{cases} \end{aligned}$$

i.e., ‘symmetric’ pairs that have the same sign make positive contributions and pairs that have a different

sign contribute 0 to the cost. We define the cost matrix as

$$\tilde{C}(f, t, \Omega) = \sum_{d=1}^{t-f+1} \bar{C}(\Lambda(f, t, d, \Omega))$$

Clearly, one can reuse the cost for shorter intervals to calculate the cost of longer ones, namely, $\tilde{C}_{f+1, t-1}$ can be used to calculate \tilde{C}_{ft} this saving computation time.

Static Contiguity Horizontal contiguous traces in S indicate that the track is self-similar (negative values) or self-dissimilar (positive values) due to repetition. If a given tile is the same as a set of contiguous tiles following it, then there is some static contiguous region in the show. The word *static* denotes that the music is not evolving in respect of time.

Let $\Gamma(f, t, h)$ be the horizontal segment in the matrix S showing the similarity of tile $f + h - 1$ to ‘future’ tiles $f + h, f + h + 1, \dots, t$. We represent it as an ordered set

$$\begin{aligned} \Gamma(f, t, h) &= \langle S_{f+h-1, f+h-1}, S_{f+h-1, f+h}, S_{f+h-1, f+h+1}, \dots, \\ &\quad S_{f+h-1, t} \rangle. \end{aligned}$$

We define the cost of an ordered set $\Gamma = \langle \Gamma_1, \Gamma_2, \dots, \Gamma_{|\Gamma|} \rangle$ as

$$\bar{C}(\Gamma, \Omega) = \sum_{i=\rho}^{|\Gamma|} \hat{C}(\Gamma^i, \Omega),$$

where ρ is a parameter indicating the minimum number of contiguous tiles required, $\Gamma^i = \langle \Gamma_1, \Gamma_2, \dots, \Gamma_i \rangle$ and

$$\hat{C}(\Gamma^i, \Omega) = \begin{cases} \frac{1}{i} \sum_{j=1}^i \tilde{f}(\Gamma_j, \Omega), & \text{if the numbers} \\ & \text{sign } G_1, \dots, \text{sign } G_i \\ & \text{are the same,} \\ 0, & \text{otherwise,} \end{cases}$$

where

$$\tilde{f}(v, \Omega) = \begin{cases} \Omega v, & \text{if sign } v = 1 \\ (1 - \Omega)v, & \text{otherwise.} \end{cases}$$

The cost of an interval is defined as

$$\tilde{C}(f, t, \Omega) = \sum_{h=1}^{t-f} \bar{C}(f, t, \Omega) .$$

Again it is possible to reuse values $\tilde{C}(f, t, \Omega)$ from shorter intervals to save calculation time for longer intervals.

Evolutionary Contiguity Any diagonal traces in S that are parallel to the main diagonal are partial copies of the track in the future which evolve in respect of time. Evolutionary Contiguity is a cost matrix which compares all adjacent pairs on the diagonals in \hat{t} , using the comparator $\delta(p, q, \Omega)$ from the standard symmetry cost function and multiplies those values by the time horizon. All the values are summed and normalized by the track width squared.

Let the evolutionary cost matrix

$$\begin{aligned} C(f, t, \Omega) &= \hat{N}_\Omega \left(\frac{\sum_{d=1}^{t-f+1} \sum_{i=2}^{|\Lambda(f, t, d)|} \delta(\kappa_i, \kappa_{i-1}, \Omega) \cdot d}{(t-f+1)^2} \right) \end{aligned}$$

where $\kappa_i = \Lambda(f, t, i)$. As before \tilde{C} is normalized by track width and incentive bias.

Gaussian Let

$$G(\varpi, N)_{tw} = e^{-\frac{1}{2} \frac{\varpi n}{\frac{1}{2} W}^2}$$

for all $n = 1, 2, \dots, W$ denote the Gaussian matrix cost function of $N \times W$. $G(\varpi, N)$ is time-independent and every row is the same. We will use this cost function for regularising the others and for use on its own for comparison against a ‘naive’ competitor. Increasing values of ϖ will tighten up the Gaussian although after experimentation we observed that 1 was always the best value and stuck with that.

Normalization All cost matrices are normalized in the following manner. First we divide $\tilde{C}(f, t, \Omega)$ by the track length,

$$\tilde{C}(f, t, \Omega) \leftarrow \frac{\tilde{C}(f, t, \Omega)}{t-f+1}.$$

Then we scale and shift C according to the value of Ω . If $\Omega = 0$, we want the resulting values to fill the interval $[0, 1]$, if $\Omega = 0.5$ we want the resulting values to fill the interval $[-1, 1]$, and if $\Omega = 1$ we want the resulting values to fill the interval $[-1, 0]$. For intermediate values of Ω we want a linear combination.

This is achieved by applying the normalization function to each element:

$$\hat{N}(x, \Omega) = \frac{x - \min_{ft} C(f, t, \Omega)}{\max_{ft} C(f, t, \Omega) - \min_{ft} C(f, t, \Omega) - s_\Omega} \hat{h}_\Omega$$

where $\hat{h}_w = 2 - 2|0.5 - \Omega|$ and

$$s_\Omega = \begin{cases} 1 - 2(0.5 - \Omega) & \text{if } \Omega < 0.5, \\ 1 & \text{otherwise.} \end{cases}$$

We let

$$C(f, t, \Omega) = N(\tilde{C}(f, t, \Omega), \Omega)$$

for all f and t .

Mixing Cost Functions We mix cost matrices together by adding them. In our experiments we will have a parameter for each cost matrix $\in [0, 1]$ to show its contribution to the mixture. The cost matrices will be multiplied by this number before being mixed.

See Figure ?? for an illustration of a single cost matrix and a mixture.

5 Computing Best Segmentation

We obtain the cost of a full segmentation by summing the costs of its tracks. The goal is now to efficiently compute the segmentation of least cost.

We want to reconstruct m track boundaries ($m + 1$ tracks).

A sequence $\mathbf{t} = (t_1, \dots, t_{m+1})$ is called an m/T -segmentation if and only if

$$1 = t_1 < \dots < t_m < t_{m+1} = T + 1.$$

m is the number of tracks we are trying to find and is a parameter of the algorithm. We use the interpretation that track $i \in \{1, \dots, m\}$ comprises times

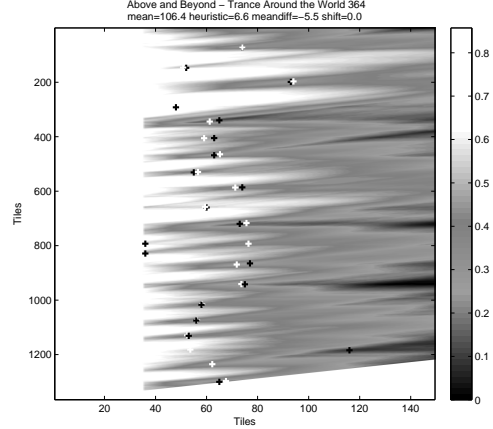


Figure 5: Summation cost matrices for Magic Island episode 110 with an incentive bias $\Omega = 1$ and therefore containing disincentives.

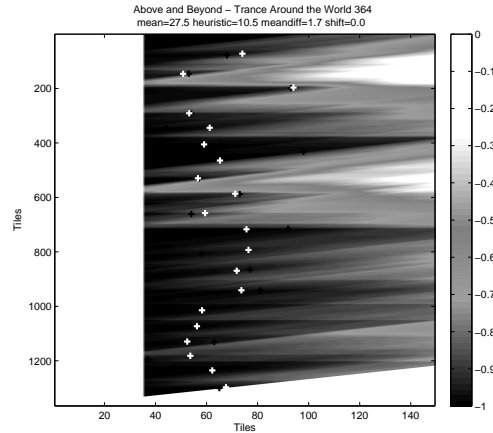


Figure 6: Summation cost matrices for Magic Island episode 110 with an incentive bias $\Omega = 0$ and therefore containing incentives.

$\{t_i, \dots, t_{i+1} - 1\}$. Let \mathbb{S}_m^T be the set of all m/T -segmentations. Note that there are a very large number of possible segmentations

$$|\mathbb{S}_m^T| = \binom{T-1}{m-1} = \frac{(T-1)!}{(m-1)!(T-m)!} = \frac{(T-1)(T-2) \cdots (T-m+1)}{(m-1)!} \geq \left(\frac{T}{m}\right)^{m-1}.$$

For large values of T , considering all possible segmentations using brute force is infeasible. For example, a two hour long show with 25 tracks would have more than

$$\left(\frac{60^2 \times 2}{25}\right)^{24} \approx 1.06 \cdot 10^{59}$$

possible segmentations.

We can reduce this number slightly by imposing upper and lower bounds on the song length. Recall that W is the upper bound (in seconds) of the song length, w the lower bound (in seconds) and m the number of tracks. With the track length restriction in place, the number of possible segmentations is still massive. A number now on the order of 10^{56} for a two hour show with 25 tracks, $w = 190$ and $W = 60 \cdot 15$.

Let $N(T, W, w, m)$ be the number of segmentations with time T (in tiles),

We can write the recursive relation

$$N(T, W, w, m) = \sum N(t_m - 1, W, w, m - 1),$$

where the sum is taken over t_m such that

$$\begin{aligned} t_m &\leq T - w + 1 & t_m &\geq T - W + 1 \\ t_m &\geq (m-1)w + 1 & t_m &\leq (m-1)W + 1 \end{aligned}$$

The first two inequalities mean that the length of the last track is within an acceptable boundary between w and W . The last two inequalities mean that the lengths of the first $m-1$ tracks are within the same boundaries.

We calculated the value of $N(7000, 60 \cdot 15, 190, 25)$ and got $5.20 \cdot 10^{56}$ which is still infeasible to compute with brute force.

Our solution to this problem is to find a dynamic programming recursion.

The loss of an m/T -segmentation \mathbf{t} is

$$\ell(\mathbf{t}) = \sum_{i=1}^m C(t_i, t_{i+1} - 1)$$

We want to compute

$$\mathcal{V}_m^T = \min_{\mathbf{t} \in \mathbb{S}_m^T} \ell(\mathbf{t})$$

To this end, we write the recurrence

$$\mathcal{V}_1^t = C(1, t)$$

and for $i \geq 2$

$$\begin{aligned} \mathcal{V}_i^t &= \min_{\mathbf{t} \in \mathbb{S}_i^t} \ell(\mathbf{t}) \\ &= \min_{t_i} \min_{\mathbf{t} \in \mathbb{S}_{i-1}^{t_i-1}} \ell(\mathbf{t}) + C(t_i, t) \\ &= \min_{t_i} C(t_i, t) + \min_{\mathbf{t} \in \mathbb{S}_{i-1}^{t_i-1}} \ell(\mathbf{t}) \\ &= \min_{t_i} C(t_i, t) + \mathcal{V}_{i-1}^{t_i-1} \end{aligned}$$

In this formula t_i ranges from $t - W$ to $t - w$. We have $T \times m$ values of \mathcal{V}_m^T and calculating each takes at most $O(W)$ steps. The total time complexity is $O(TWm)$.

6 Confidence Intervals

It may be useful for some applications to build a framework to allow confidence intervals for our predicted indices. This may also be useful for meaningful comparison of cost matrices.

6.1 Posterior Marginal of Song Boundary

Fix a learning rate η , and fix T and m . Let

$$P(j, s) = \frac{\sum_{\mathbf{t} \in \mathbb{S}_m^T: t_j = s} e^{-\eta \ell(\mathbf{t})}}{\sum_{\mathbf{t} \in \mathbb{S}_m^T} e^{-\eta \ell(\mathbf{t})}}$$

That is, $P(j, s)$ is the “posterior probability” that song j starts at time s .

To compute $P(j, s)$, we need an extended notion of segmentation. We call \mathbf{t} a $m/F : T$ segmentation if

$$F = t_1 < \dots < t_m < t_{m+1} = T + 1.$$

Let $\mathbb{S}_m^{F:T}$ be the set of all $m/F - T$ -segmentations. We have

$$\begin{aligned} \sum_{\mathbf{t} \in \mathbb{S}_m^{F:T} : t_j = s} e^{-\eta \ell(\mathbf{t})} &= \sum_{\substack{\mathbf{t} \in \mathbb{S}_{j-1}^{s-1}, \\ \mathbf{t}' \in \mathbb{S}_{m-j+1}^{s:T}}} e^{-\eta(\ell(\mathbf{t}) + \ell(\mathbf{t}'))} = \\ &\left(\sum_{\mathbf{t} \in \mathbb{S}_{j-1}^{s-1}} e^{-\eta \ell(\mathbf{t})} \right) \left(\sum_{\mathbf{t}' \in \mathbb{S}_{m-j+1}^{s:T}} e^{-\eta \ell(\mathbf{t}')} \right) \end{aligned}$$

which upon abbreviating

$$\mathcal{H}_m^t = \sum_{\mathbf{t} \in \mathbb{S}_m^t} e^{-\eta \ell(\mathbf{t})} \quad \mathcal{T}_m^f = \sum_{\mathbf{t} \in \mathbb{S}_m^{f,T}} e^{-\eta \ell(\mathbf{t})}$$

means that we can write

$$P(j, s) = \frac{\mathcal{H}_{j-1}^{s-1} \cdot \mathcal{T}_{m-j+1}^s}{\mathcal{H}_m^T}.$$

So it suffices to compute \mathcal{H}_m^t and \mathcal{T}_m^t for all relevant t and m . We use

$$\mathcal{H}_1^t = e^{-\eta C(1,t)} \quad \mathcal{T}_1^f = e^{-\eta C(f,T-f+1)}$$

and for $m \geq 2$

$$\begin{aligned} \mathcal{H}_m^t &= \sum_{t_m} \sum_{\mathbf{t} \in \mathbb{S}_{m-1}^{t_m-1}} e^{-\eta(\ell(\mathbf{t}) + C(t_m, t - t_m + 1))} \\ &= \sum_{t_m} e^{-\eta C(t_m, t - t_m + 1)} \sum_{\mathbf{t} \in \mathbb{S}_{m-1}^{t_m-1}} e^{-\eta \ell(\mathbf{t})} \\ &= \sum_{t_m} e^{-\eta C(t_m, t - t_m + 1)} \mathcal{H}_{m-1}^{t_m-1} \\ \mathcal{T}_m^f &= \sum_{t_2} \sum_{\mathbf{t} \in \mathbb{S}_{m-1}^{t_2:T}} e^{-\eta(C(f, t_2 - f) + \ell(\mathbf{t}))} \\ &= \sum_{t_2} e^{-\eta C(f, t_2 - f)} \sum_{\mathbf{t} \in \mathbb{S}_{m-1}^{t_2:T}} e^{-\eta \ell(\mathbf{t})} \\ &= \sum_{t_2} e^{-\eta C(f, t_2 - f)} \mathcal{T}_{m-1}^{t_2} \end{aligned}$$

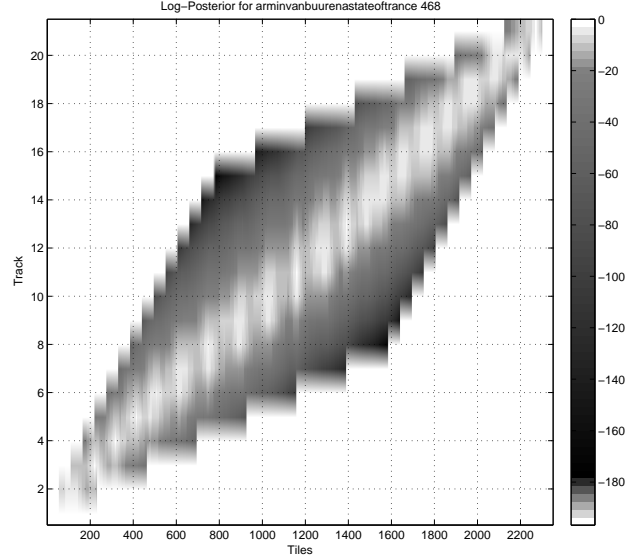
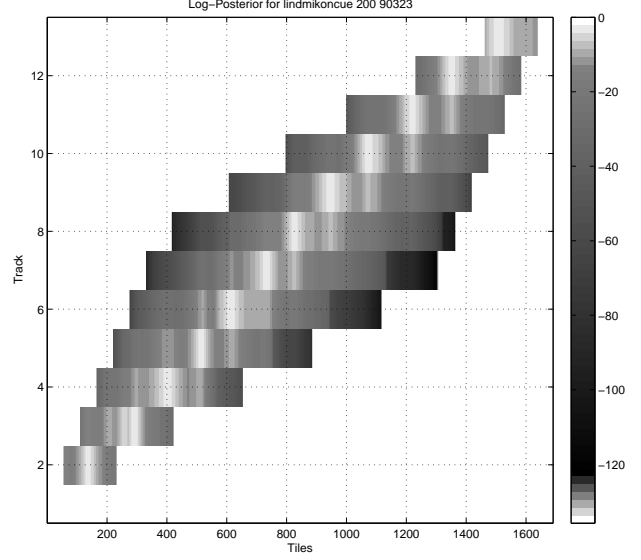


Figure 7: A visualization of $\log(P(j, s))$ ($\eta = 10$) for one of the shows in the test set using the cost matrix parameters from experiment 10. The actual tracks are overlaid as white crosses.

See Figure ?? for an example of the posterior for a radio show.

6.2 Posterior Marginal of Song Position

Fix a learning rate η , and fix T and m . Let

$$P(j, s, f) = \frac{\sum_{t \in \mathbb{S}_m^T: t_j = s \wedge t_{j+1} - 1 = f} e^{-\eta \ell(t)}}{\sum_{t \in \mathbb{S}_m^T} e^{-\eta \ell(t)}}$$

That is, $P(j, s, f)$ is the “posterior probability” that song j starts at time s and finishes at time f . In the same vein as the last section, we now get

$$P(j, s, f) = \frac{\mathcal{H}_{j-1}^{s-1} \cdot e^{-\eta C(s, f-s+1)} \cdot \mathcal{T}_{m-j}^{f+1}}{\mathcal{H}_m^T}.$$

6.3 Track Index Confidence

We can use the posterior marginal of song boundary to give estimates of confidence on track index placement and time accuracy.

To estimate the uncertainty of correct track alignment, we select the next highest probability of other track placements at the same time of the optimal placement from $P(j, s)$ and normalize them by the probability of the optimal placement.

Let track index confidence

$$I(j) = 1 - \frac{P(\{1, \dots, M\} \setminus j, \text{SortInd}(P(j, 1, \dots, T)_1)_2)}{\max(P(j, 1, \dots, T)_2)}$$

where $\text{SortInd}(l)$ will return the original indices corresponding to the sorted list of l .

6.4 Track Time Confidence

Track time uncertainty is estimated by normalizing the value of the next most significant peak in $P(j, \forall s)$ by the probability of actual track placement (which is the most significant peak). Let

$$\tilde{I}(j) = 1 - \frac{\text{Peaks}(P(j, 1, \dots, T))_2}{\text{Peaks}(P(j, 1, \dots, T))_1}$$

where $\text{Peaks}(s_i)$ is a peak finding algorithm that returns the peaks in descending order of magnitude (we used the `findpeaks` function in MatLab).

7 Experiments

7.1 Training Set

We selected 6 shows at random (two of each show type) to create a training set, which we will refer to as the *GitHub training set*. See Table 1 to see the shows we selected.

7.2 Finding The Best Cost Matrix

We used the GitHub training set to develop the cost matrices from first principles and to find robust parameters using a genetic algorithm search.

In our experiments we decided to fix the tile size at 5 seconds for the sake of speed. A lower tile size does increase accuracy but only marginally. Higher tile sizes can perform more robustly (fewer catastrophic misalignments) but progressively lose out on accuracy.

8 Estimating Segment Count

The original goal of our work was providing the best possible segmentation given a fixed number of tracks m , rather than estimating m . The problem domain is reasonably unique; the number of segments is known a priori but segmentation itself is not.

In the most basic setting. The number of tracks could be estimated as the variable of track lengths is Gaussian (see Figure 1).

We propose the following method of adapting our framework to estimate the number of contiguous segments in a data stream.

9 Evaluation

It is challenging to quantify the performance of our method because if we misplace any tracks, it may have a cascade effect. For example if we place one

Table 1: The shows randomly selected for inclusion in the *GitHub training set*.

#	Show Name	Artist	Date Broadcast
1	A State Of Trance 453	Armin Van Buuren	April 2010
2	A State Of Trance 462	Armin Van Buuren	June 2010
3	Magic Island 098	Roger Shah	March 2010
4	Magic Island 112	Roger Shah	July 2010
5	Trance Around The World 364	Above & Beyond	March 2011
6	Trance Around The World 372	Above & Beyond	May 2011

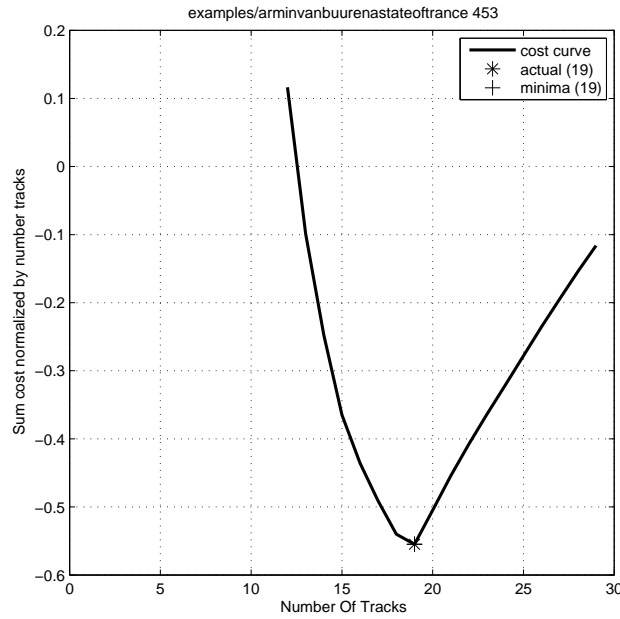


Figure 8: Number of tracks estimated correctly a show in the GitHub training set after a genetics algorithm was run to select a new set of configuration parameters.

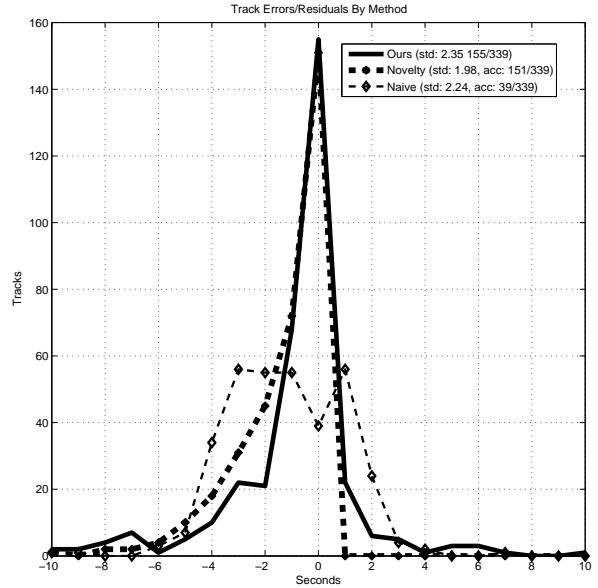


Figure 9:

track too many early on in a show, many of the subsequent tracks may be correctly detected but placed out of alignment.

For the task of computing the best cost segmentation when the number of tracks are known a priori, we can use simple statistical descriptions of the track residuals $|P_{st} - A_{st}|$ where P is the predicted track, A the actual track, for show s and track t (for all s and t in the corpus). The mean average will give a good indication of the accuracy combined with the amount of misplacements. The median of the residuals will indicate the actual track accuracy invariant to any catastrophic misplacements. The standard deviation of the residuals will indicate the spread of error.

10 Results

Please see Table ?? for the main table of results, Figure 10 for a histogram of predicted versus actual differences for experiment 10.

321	321		dsa		
312	312	dsa	dsa		
			d	dasd	
			asdas	dsa	

On our previous work we were using a disincen- tive only summation matrix, and found that normal- izing it on the square root of the width produced the best result. This would have been necessary to encourage placement of longer tracks as no incentive was present. So experiment 11 is roughly comparable and indeed produces the same overall mean average to that previous experiment ($\approx 20S$). Note that we no longer discard any shows from evaluation which makes the result stronger.

11 Summary

We believe our algorithm would be useful for seg- menting DJ-mixed audio streams in batch mode. It would be excellent if Sound Cloud³ for example started to do something similar. Sound Cloud is an

³<http://www.soundcloud.com>

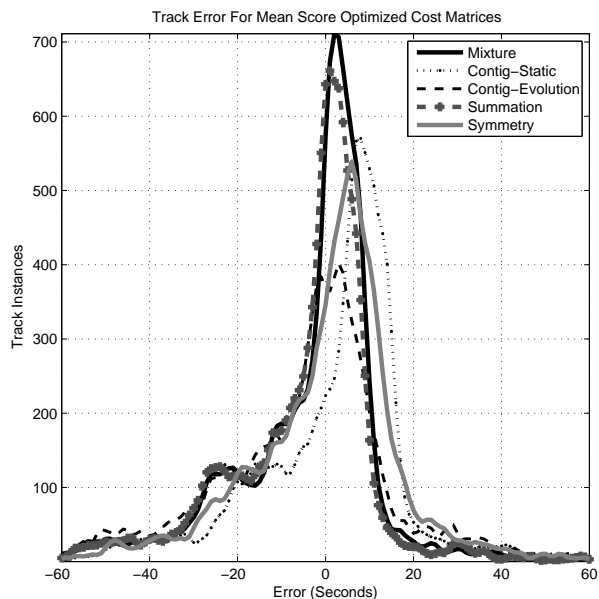


Figure 10: Histogram of the residuals (errors) between reconstructed and human captured time indices. Apart from obvious noise there appears to be a tendency for the algorithm to place an index slightly earlier. Contig-static cost matrix apparently contained parameters selected through optimization that shifted it to remove some of this effect.

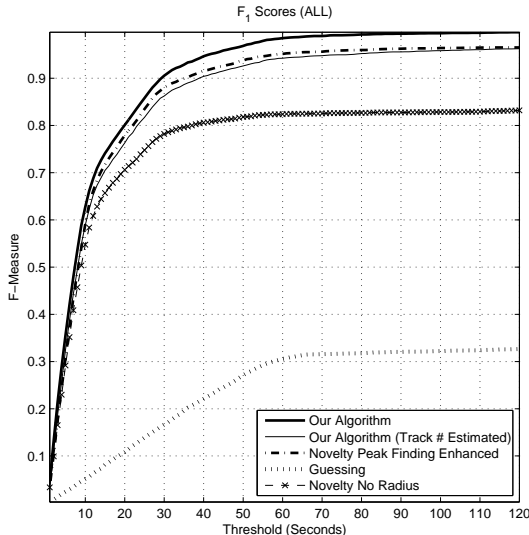


Figure 12: Comparison between our algorithm and the Foote novelty peak finding approach on all of the datasets.

on-line music service with many electronic dance music radio shows with the track listing in text. This method would allow them to reliably segment the shows, and they could display an interactive segmentation in the music player.

The new cost matrices in combination improve robustness significantly over single matrices or regularised single matrices (as in our last paper). We are seeing about a 50% improvement in overall mean accuracy over single cost matrices that are correctly normalized. The new cost matrices improve on many drawbacks of our previous work (mainly that it was vulnerable to dissimilar regions within tracks).

Our method still has one key drawback that we are aware of. This is the rare instance where there are head or tail segments to a track that seem independent from the rest of the track. When these are small they usually get absorbed without any problems but they can cause misplacements. In spite of this issue we suspect that our predictions are more accurate and more consistent than the human equivalents while not being as precise in situations when

our index agrees with theirs.

A more precise corpus where the DJ was also the cuesheet author would allow us to tune the parameters more succinctly and also the generation of an artificial dataset to test against. This work is forthcoming.

We would also like to implement some of the methods in the literature (which were mostly designed for scene analysis) to see if we outperform them. It would be tricky to get an exact comparison because we could not find an unsupervised deterministic algorithm which finds a fixed number of strictly contiguous clusters. We could however adapt existing algorithms to get a like for like comparison.

12 Acknowledgements

We would like to thank Mikael Lindgren and Denis Goncharov from *cuenation*⁴ for their help explaining how they make cue-sheets and for providing the data set to test the algorithm on.

13 Materials

All the code presented in this paper with the training set is available on GitHub⁵. The large data set ($\approx 130\text{GB}$) we received from Denis Goncharov and Mikael Lindgren can easily be made available on request (it is in a cloud storage account).

References

- [1] R. Curticapean, “Clustering-based audio segmentation with applications to music structure analysis,”
- [2] J. Matejka, T. Grossman, and G. Fitzmaurice, “Swifter: Improved online video scrubbing,” in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI ’13, (New York, NY, USA), pp. 1159–1168, ACM, 2013.

⁴<http://www.cuenation.com>

⁵github.com/ecsplendid/DanceMusicSegmentation

- [3] G. Tzanetakis and F. Cook, "A framework for audio analysis based on classification and temporal segmentation," in *EUROMICRO Conference, 1999. Proceedings. 25th*, vol. 2, pp. 61–67, IEEE, 1999.
- [4] J.-P. Eckmann, S. O. Kamphorst, and D. Ruelle, "Recurrence plots of dynamical systems," *EPL (Europhysics Letters)*, vol. 4, no. 9, p. 973, 1987.
- [5] J. Foote, "Visualizing music and audio using self-similarity," in *Proceedings of the seventh ACM international conference on Multimedia (Part 1)*, pp. 77–80, ACM, 1999.
- [6] J. Foote, "A similarity measure for automatic audio classification," in *Proc. AAAI 1997 Spring Symposium on Intelligent Integration and Use of Text, Image, Video, and Audio Corpora*, 1997.
- [7] J. Foote, "Automatic audio segmentation using a measure of audio novelty," in *Multimedia and Expo, 2000. ICME 2000. 2000 IEEE International Conference on*, vol. 1, pp. 452–455, IEEE, 2000.
- [8] J. T. Foote and M. L. Cooper, "Media segmentation using self-similarity decomposition," in *Electronic Imaging 2003*, pp. 167–175, International Society for Optics and Photonics, 2003.
- [9] J. Foote and M. Cooper, "Visualizing musical structure and rhythm via self-similarity," in *Proceedings of the 2001 International Computer Music Conference*, pp. 419–422, 2001.
- [10] M. M. Goodwin and J. Laroche, "A dynamic programming approach to audio segmentation and speech/music discrimination," in *Acoustics, Speech, and Signal Processing, 2004. Proceedings.(ICASSP'04). IEEE International Conference on*, vol. 4, pp. iv–309, IEEE, 2004.
- [11] M. M. Goodwin and J. Laroche, "Audio segmentation by feature-space clustering using linear discriminant analysis and dynamic programming," in *Applications of Signal Processing to Audio and Acoustics, 2003 IEEE Workshop on*, pp. 131–134, IEEE, 2003.
- [12] E. Peiszer, T. Lidy, and A. Rauber, "Automatic audio segmentation: Segment boundary and structure detection in popular music," *Proc. of LSAS*, 2008.
- [13] D. El Badawy, P. Marmaroli, and H. Lissek, "Audio novelty-based segmentation of music concerts,"
- [14] T. Scarfe, W. M. Koolen, and Y. Kalnishkan, "A long-range self-similarity approach to segmenting dj mixed music streams," in *Artificial Intelligence Applications and Innovations*, pp. 235–244, Springer, 2013.
- [15] H. Nyquist, "Certain topics in telegraph transmission theory," *American Institute of Electrical Engineers, Transactions of the*, vol. 47, no. 2, pp. 617–644, 1928.
- [16] M. Frigo and S. G. Johnson, "The fftw web page," 2004.
- [17] G. Tzanetakis and P. Cook, "Multifeature audio segmentation for browsing and annotation," in *Applications of Signal Processing to Audio and Acoustics, 1999 IEEE Workshop on*, pp. 103–106, IEEE, 1999.

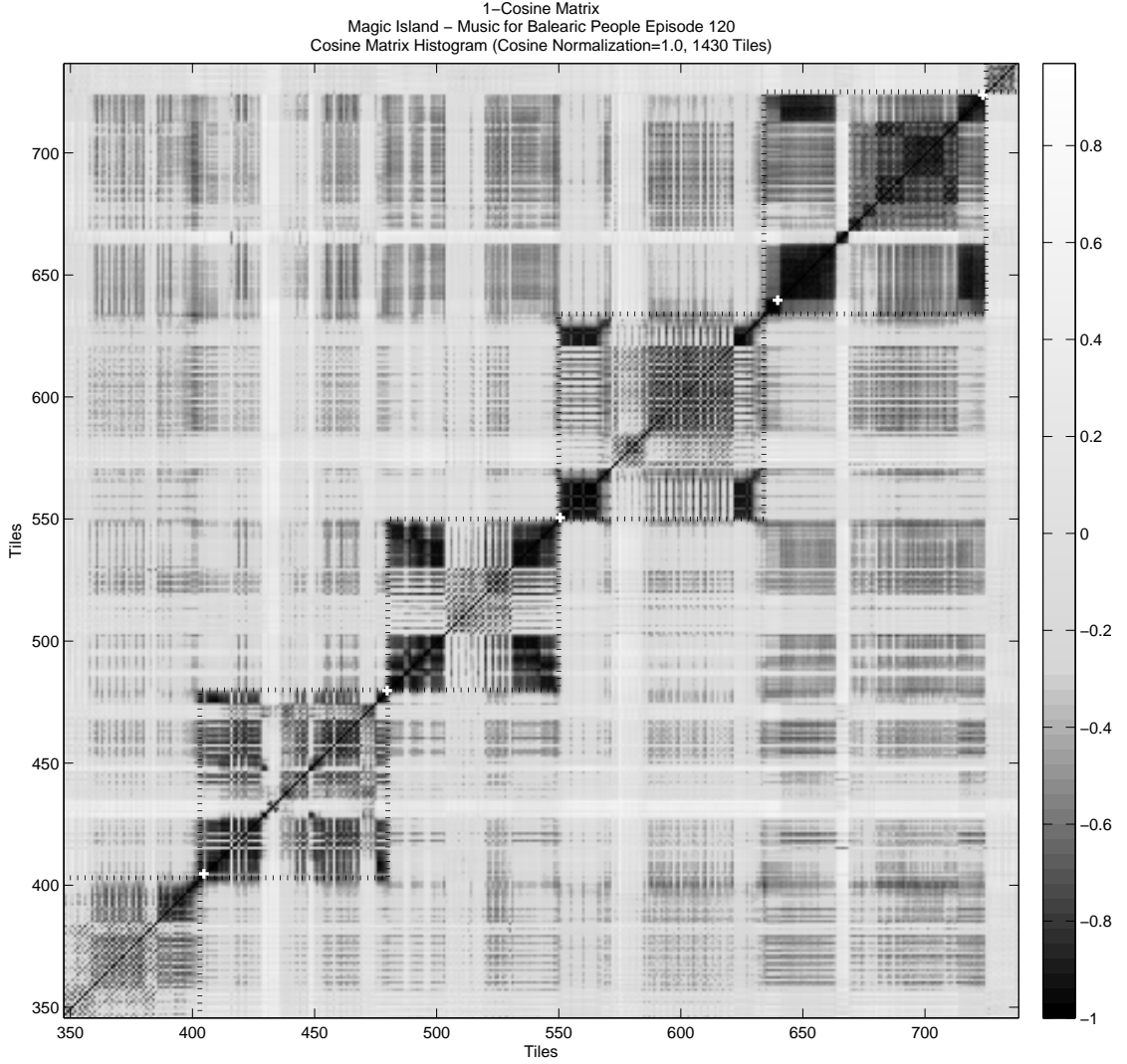


Figure 4: An illustration of the similarity matrix S with the actual indices drawn on with white crosses, and our reconstructed indices indicated with the black dotted lines. There are examples here of evolutionary repetition ($t = 500, \dots, 550$), static contiguity everywhere where there is solid black, and symmetry on the middle two tracks.

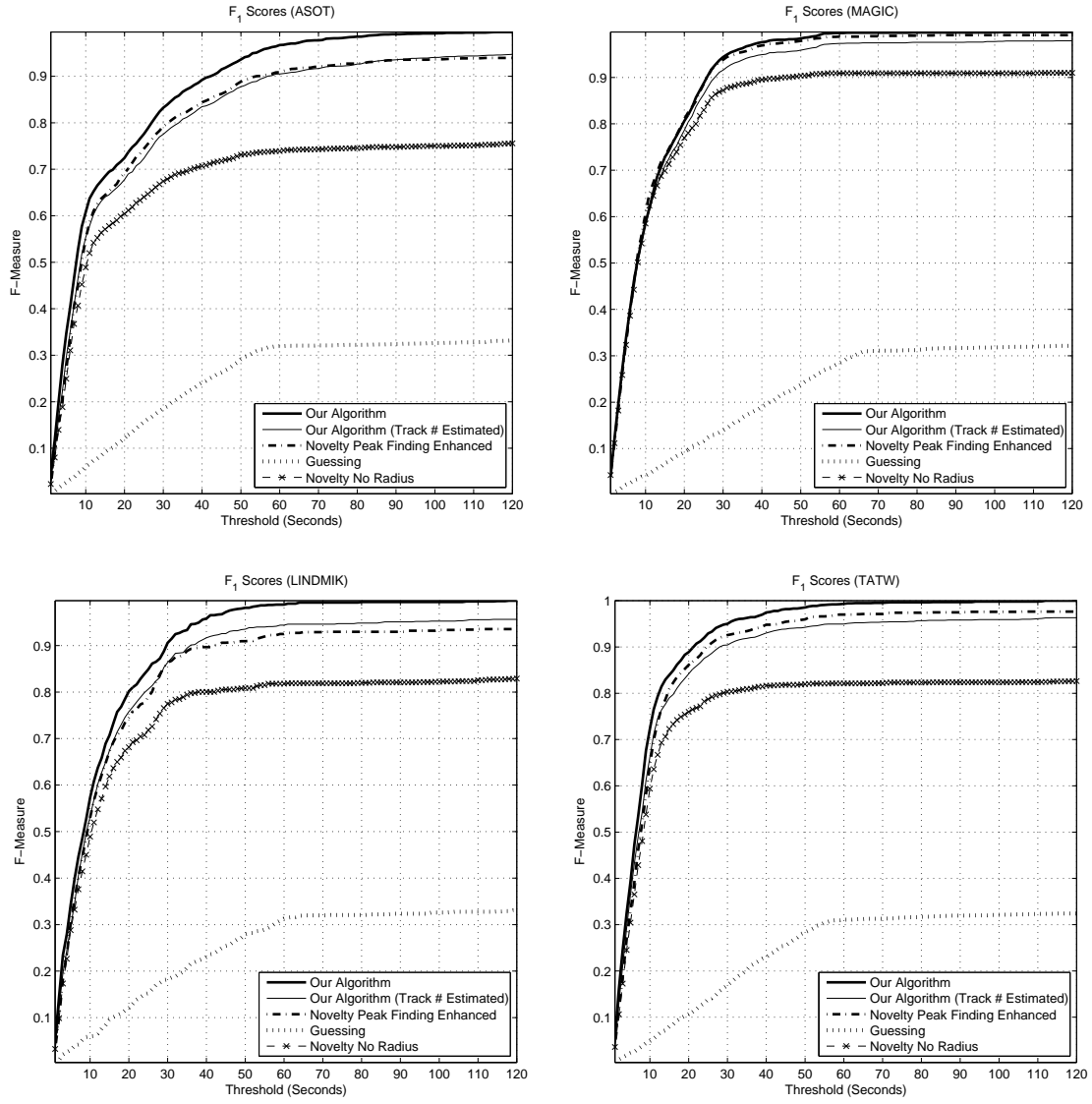


Figure 11: lalala