

Entropy-based segmentation of birdcalls using Fourier transform phase

Author Name¹, Co-author Name²

¹Author Affiliation

²Co-author Affiliation

author@university.edu, coauthor@company.com

Abstract

In this paper we describe an entropy-based algorithm for the segmentation of birdcalls from recordings. The entropy of time-frequency blocks are estimated from the phase of the Fourier transform. To overcome difficulties in processing the phase, the group delay function from an all-pole filter is utilised. The group delay function has good frequency resolution properties, and hence provides reliable estimates of the entropy. Furthermore, spectral whitening is performed to smooth the entropy estimate and the extremities are determined. A threshold is applied on the difference to distinguish the call periods from the background. The algorithm is evaluated on two different datasets, one of which is recorded in more challenging field conditions. When compared to entropy estimated from the power spectrum, the entropy from the group delay function provides better detection accuracy at almost all operating points. The choice of model order of the all-pole filter for different bird species is also briefly investigated.

Index Terms: bioacoustics, birdcall segmentation, Fourier transform phase

1. Introduction

With the advent of automated recording devices, the collection of large amounts of bioacoustic data has become relatively easy. By analysing birdcalls collected in this manner, it is possible to perform tasks such as the tracking of migrant species or examining the avian biodiversity of a given region. Typically, the collected data is processed offline. In this process, the first step is usually to determine regions of interest in the recording. An entropy-based bird phrase segmentation technique was developed in [1]. In this paper, we propose a modified version of that technique, by using information from the phase of the short-term Fourier transform (STFT.) Most techniques for processing speech and audio signals have utilised the magnitude spectrum of the STFT. Although the phase spectrum of the STFT has useful information, its processing has remained difficult. A popular technique for exploiting information from the phase has been through group delay functions. In this work, we utilise information from group delay functions using parametric models, and apply it to segment birdcalls into active and inactive regions.

The group delay function has good frequency resolution properties, which enable it to be useful in tasks such as speech recognition and speaker recognition [2] [3] [4]. The same property is beneficial in the processing of bird vocalizations. In [1], the entropy within a sliding time-frequency window over the spectrogram has been effectively used for distinguishing active and inactive regions. The essential idea is that birdcalls have more structure (for eg. harmonics may be present), and thus have lower entropy when compared to background sounds,

which have higher entropy. This difference in entropy levels enable effective distinction between birdcalls and the background. The high resolution property of group delay functions enable accurate tracking of time-frequency information [3]. In this work, the entropy of a sliding time-frequency window is estimated from the group delay representation. Spectral whitening is applied to smooth the entropy estimates and thresholding on differences of extrema is applied to separate the birdcalls from the background.

Most of the bioacoustic studies have used manually segmented bird calls [5] [6] [7]. Time domain segmentation using energy has been used in many studies [8] [9] [10]. The energy based segmentation method is highly influenced by background noise and will where bird calls have low energy in comparison to the background. A KL-divergence based segmentation method is proposed in [11]. KL-divergence between normalized power spectral density of a frame and uniform distribution is computed. Local minima of KL divergence act as change points for bird vocalizations. In [12], time-frequency domain based segmentation using random forest classifier is proposed to segment the syllables from noisy audio signal.

2. Utilising Fourier transform phase

Commonly used features for processing speech and audio signals are based on the magnitude spectrum of the short-term Fourier transform. The phase spectrum has received relatively lesser attention due to signal processing difficulties, one of them being the need to unwrap the phase spectrum. The unwrapping problem can be bypassed by utilising the group delay function, which is the negative derivative of the phase spectrum. The group delay function can be computed using properties of the Fourier transform, and hence avoids the need for explicit computation of the phase spectrum [13]. However, this method can produce artifacts in the form of spurious peaks at spectral nulls. These nulls correspond to zeros close to the unit circle when the vocal tract transfer function is represented in the Z domain. Several methods have been proposed in the literature to overcome the effects due to these artifacts [4, 14]. Another technique to overcome this difficulty is to model the vocal tract as an all-pole filter, hence avoiding the nulls altogether. Such a technique derived using linear prediction analysis was used in the detection of formants in human speech [15].

There is strong evidence that birds use their vocal tract as a selective filter to modify the final sound [16]. Given this, the source-filter model developed for analysing human speech can be applied to bird vocalizations as well. Linear predictive (LP) analysis of human speech signals models the vocal tract spectrum as an all-pole filter [17] excited by a single source. When applied to birdcalls, this is a simplification of the ‘two-voice’ theory of avian vocalization [16], in that there is assumed to

be only one source, rather than two. Nevertheless, this reasonable assumption is followed in this work. A similar assumption has been made in [18], where LP analysis has been applied in analysing the song of the greater racket-tailed drongo.

The vocal tract is represented in the LP model as

$$H(\omega) = \frac{G}{1 - \sum_{k=1}^P a(k)e^{-j\omega k}}, \quad (1)$$

where the predictor model order is P , G represents the gain and $a(k)$ are the predictor coefficients [17]. The filter represented by $H(\omega)$ is an all-pole filter, and its group delay function does not suffer from the artifacts mentioned earlier. The group delay function computed in this manner is termed as all-pole group delay function (APGDF.) Figure 1 shows the magnitude spectrum, LP spectrum and APGDF derived from a 20 ms call of Cassins vireo (*Vireo cassinii*.) As can be seen, the APGDF emphasises the formants, as compared to the DFT magnitude spectrum or the LP magnitude spectrum. Two peaks which are merged in the magnitude spectra around the 100th frequency bin appear distinctly in the APGDF.

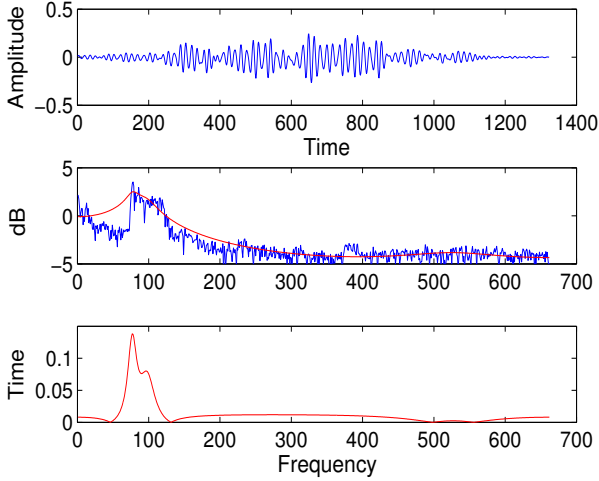


Figure 1: based on all pole model (in the bottom panel).

Recently, feature vectors derived from such a representation were used in speech [19] and speaker recognition [3].

3. Entropy-based segmentation of birdsong

Computing the APGDF for every frame enables good time-frequency resolution of the audio recording. Unlike the power spectrum, the APGDF can be negative. Since here we are interested only in the magnitudes and locations of the frequency components, the sign of the APGDF is ignored by taking the absolute value. Henceforth, APGDF means positive APGDF. A sliding time-frequency window of width w frames and frequency range $fmin$ to $fmax$ is considered over the APGDF vectors. As in [1], the entropy of this window is estimated using the expression

$$h_k = - \sum_{n=kT+1}^{kT+w} \sum_{f=fmin}^{fmax} \tau_N(n, f) \ln \tau_N(n, f), \quad (2)$$

where T is the time-frequency window shift and $\tau_N(n, f)$ is the normalized APGDF. $\tau_N(n, f)$ is computed as

$$\tau_N(n, f) = \frac{\tau(n, f)}{\sum_{n=kT+1}^{kT+w} \sum_{f=fmin}^{fmax} \tau(n, f)}, \quad (3)$$

where $\tau(n, f)$ is the APGDF at frame n and frequency f .

The entropy of the time-frequency windows containing a bird call is lower than the one containing only the background. *What happens at transistions?* The entropy calculated from time-frequency window is less susceptible to sudden changes in background as compared to the entropy calculated at each time instance [1]. *check this*

3.1. Whitening APGDF before Entropy Calculation

The drop in entropy during the presence of bird vocalizations makes it possible to detect call periods. Due to the presence of various background sounds e.g. rain, thunder, other animals etc., entropy of the background can vary rapidly making it difficult to distinguish it from the entropy during a call period. To mitigate this problem to some extent, following [1], the APGDF is whitened before calculating the entropy. This makes the entropy of the background relatively constant and dips enough to mark the presence of bird vocalizations. Due to this, small change in entropy can be detected quite reliably.

To whiten the APGDF ($\tau(n, f)$), the covariance matrix C of the mean-subtracted APGDF is estimated. The eigenvalue matrix S and the eigenvector matrix U of C are determined. The whitened APGDF ($\tau_W(n, f)$) is calculated as

$$\tau_W(n, f) = \text{diag} \left(\frac{1}{\sqrt{\text{diag}(S) + \epsilon}} \right) * U' * \tau(n, f) \quad (4)$$

check the formula above

Therefore, we use the whitened APGDF in our experiments. The entropy from the orinal APGDF and the whitened APGDF are shown in figure 2.

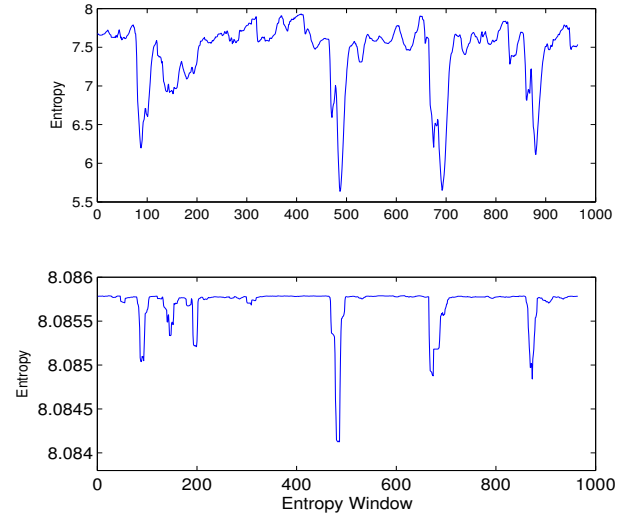


Figure 2: Entropy calculated from phase spectrum and whitened phase spectrum respectively.

3.2. Detecting change points using thresholding

To detect the change points, extrema-based thresholding is used. Local minima and maxima are estimated on the entropy. A threshold is applied on the difference of adjacent maxima and minima to determine the change point. Two contiguous change points correspond to the start and end of a bird vocalization. Using the XXX and YYY, these change points can be tracked back to get the start and end time of the vocalization in sound recording.

4. Experimentation and Performance Analysis

The proposed algorithm is evaluated on two datasets. The first consists of recordings of a single species, the Cassin's vireo, and second is a subset of the MLSP 2013 bird classification challenge, consisting of recordings from several species. In the Cassin's vireo dataset [20], the total duration of recordings is about 45 minutes, out of which about 5 minutes correspond to the calls. These recordings are fairly clean and have less background noise. Phrase annotations are provided with the dataset, the start and end of which are used as the ground truth.

The MLSP 2013 dataset [21] is much noisier and includes wind and rain in the background. Twenty five files of ten seconds each

The second dataset used is a subset of MLSP bird classification challenge 2013 dataset [21]. The dataset is collected in the H. J. Andrews (HJA) long-term experimental research forest in Oregon. The audio recordings also include rain and wind in background. The recordings are done for over two years at 13 different locations. Twenty Five files of ten seconds length are taken from the data to evaluate the proposed method.

Manual annotations are done to mark the presence of bird vocalizations in these audio files.

The proposed algorithm is compared with a modified version of technique proposed in [1]. The entropy estimated from the whitened power spectrum is used to determine the change points. Similar to the proposed algorithm, extrema-based thresholding is utilised to determine the change points.

The front-end processing utilises a frame length of 20ms and frame shift of 5ms. A time-frequency window of length $w=138.8$ ms is used along with increment of $T = 15$ ms to estimate entropy (see equation XXX.) The frequency range of the block is from XXX 1.5 kHz to 7 kHz as in [1].

Receiver operating characteristic (ROC) curves are used to analyze the performance of both the techniques. True positive rate and false alarm rate are calculated as

$$TPR(\%) = \frac{\text{Correctly Classified Call Frames}}{\text{Total Call Frames}} \times 100 \quad (5)$$

$$FAR(\%) = \frac{\text{Wrongly Classified As Call Frames}}{\text{Total Background Frames}} \times 100 \quad (6)$$

Figure 3 depicts ROC curves comparing the performance of the methods based on entropy calculated from power spectrum and the entropy calculated from APGDF. It is clear from (figure, not Figure XXX) Figure 3 that the APGDF method outperforms the power spectrum method. On what dataset is this???

Figure 4 shows ROC curves comparing performance of methods based on entropy calculated from APDGF and power

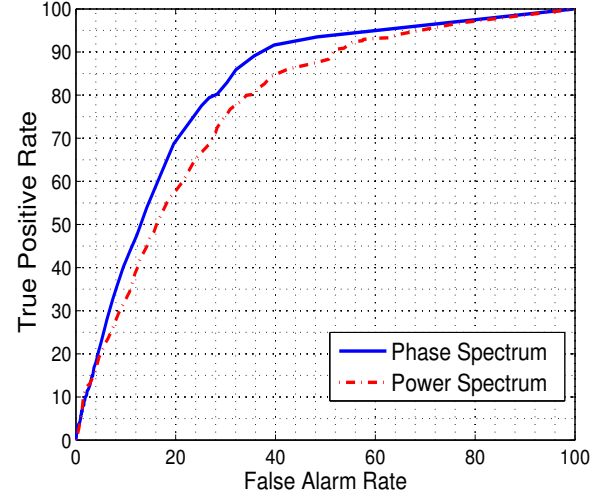


Figure 3: ROC curves comparing performance of methods based on white phase spectrum and white power spectrum on Cassin's Vireo dataset.

spectrum on dataset 2. Here too, APGDF based method outperforms power spectrum based method for almost all the operating points.

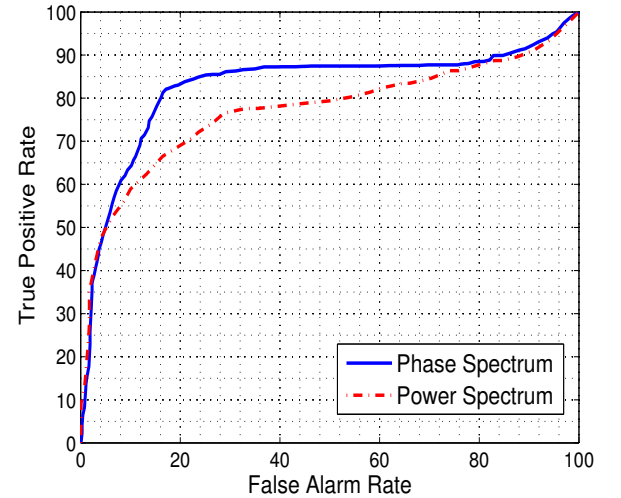


Figure 4: ROC curves comparing performance of methods based on white phase spectrum and white power spectrum on MLSP 2013 single species dataset

4.1. Model Order vs AIC

To establish the optimum model order, Akaike Information criteria (AIC) [17] is used. Figure 5 depicts the AIC for different model orders for Cassin's Vireo phrases.

Figure 6 depicts the AIC for different model orders for three different bird species i.e. Cuckoo, Great Barbet and Laughing Thrush.

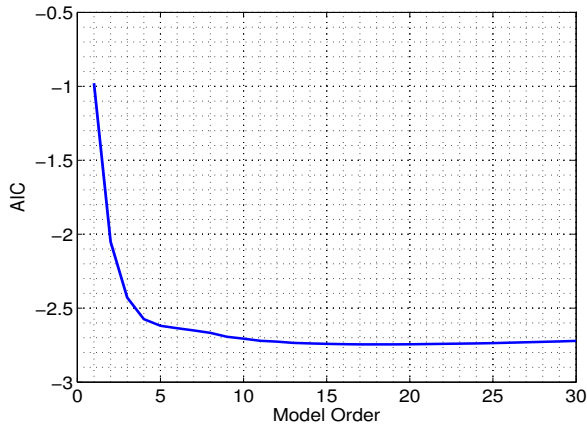


Figure 5: AIC vs model order for Cassin's Vireo

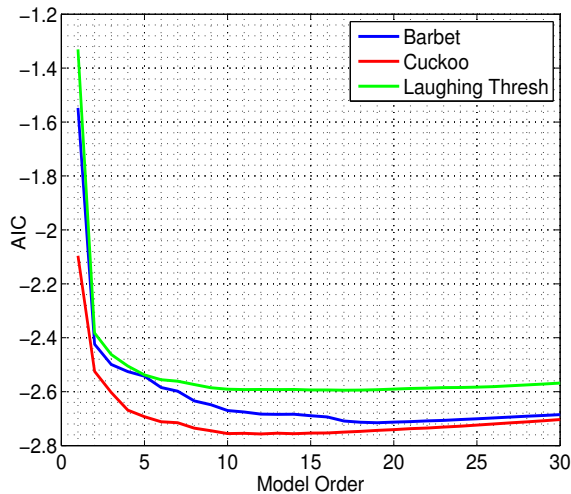


Figure 6: AIC vs model order for three different species

5. Conclusion

We propose an entropy based bird vocalization segmentation method where entropy is calculated from Group Delay phase spectrum. It is also established that whitening the power or phase spectrum before entropy calculation improves the performance of entropy based segmentation. From experimentation, it is clear that the proposed group delay based method outperforms the power spectrum based method.

6. References

- [1] N. C. Wang, R. E. Hudson, L. N. Tan, C. E. Taylor, A. Alwan, and K. Yao, "Bird phrase segmentation by entropy-driven change point detection," in *Proc. Int. Conf. Acoust. Speech, Signal Process.*, 2013, pp. 773–777.
- [2] H. A. Murthy and V. R. R. Gadde, "The modified group delay function and its application to phoneme recognition," in *Proc. Int. Conf. Acoust., Speech, Signal Process.*, vol. 1, 2003, pp. 68–71.
- [3] P. Rajan, T. Kinnunen, C. Hanili, J. Pohjalainen, and P. Alku, "Using group delay functions from all-pole models for speaker recognition," in *Proc. Interspeech*, 2013, pp. 2489–2493.
- [4] R. M. Hegde, H. A. Murthy, and V. R. Gadde, "Significance of the modified group delay feature in speech recognition," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 15, pp. 190–202, Jan. 2007.
- [5] V. M. Trifa, A. N. G. Kirschel, C. E. Taylor, and E. E. Vallejo, "Automated species recognition of antbirds in a mexican rainforest using hidden markov models," *Jnl. Acoust. Soc. Amer.*, vol. 123, no. 4, pp. 2424–2431, Apr 2008.
- [6] C. H. Lee, C. C. Han, and C. C. Chuang, "Automatic classification of bird species from their sounds using two-dimensional cepstral coefficients," *IEEE Trans. Audio, Speech, Language Process.*, vol. 16, no. 8, pp. 1541–1550, Nov 2008.
- [7] K. Kaewtip, L. N. Tan, A. Alwan, and C. E. Taylor, "A robust automatic bird phrase classifier using dynamic time-warping with prominent region identification," in *Proc. Int. Conf. Acoust. Speech, Signal Process.*, 2013, pp. 768–772.
- [8] A. Harma and P. Somervuo, "Classification of the harmonic structure in bird vocalization," in *Proc. Int. Conf. Acoust. Speech, Signal Process.*, vol. 5, 2004, pp. 701–704.
- [9] P. Somervuo, A. Harma, and S. Fagerlund, "Parametric representations of bird sounds for automatic species recognition," *IEEE Trans. Audio, Speech, Language Process.*, vol. 14, no. 6, pp. 2252–2263, Nov 2006.
- [10] S. Fagerlund, "Bird species recognition using support vector machines," *EURASIP J. Appl. Signal Process.*, vol. 2007, no. 1, pp. 64–64, Jan. 2007.
- [11] B. Lakshminarayanan, R. Raich, and X. Fern, "A syllable-level probabilistic framework for bird species identification," in *Proc. Int. Conf. Mach. Learn. Applicat.*, 2009, pp. 53–59.
- [12] L. Neal, F. Briggs, R. Raich, and X. Z. Fern, "Time-frequency segmentation of bird song in noisy acoustic environments," in *Proc. Int. Conf. Acoust. Speech, Signal Process.*, 2011, pp. 2012–2015.
- [13] H. Banno, J. Lu, S. Nakamura, K. Shikano, and H. Kawahara, "Efficient representation of short-time phase based on group delay," in *Proc. Int. Conf. Acoust. Speech Signal Process.*, vol. 2, 1998, pp. 861–864.
- [14] D. Zhu and K. Paliwal, "Product of power spectrum and group delay function for speech recognition," in *Proc. Int. Conf. Acoust., Speech, Signal Process.*, 2004, pp. 125–128.
- [15] B. Yegnanarayana, "Formant extraction from linear-prediction phase spectra," *Jnl. Acoust. Soc. Amer.*, vol. 63, no. 5, pp. 1638–1640, 1978.
- [16] C. K. Catchpole and P. J. Slater, *Bird song: biological themes and variations*. Cambridge university press, 2003.
- [17] J. Makhoul, "Linear prediction: A tutorial review," *Proceedings of the IEEE*, vol. 63, no. 4, pp. 561–580, 1975.
- [18] S. Agnihotri, P. Sundeep, C. S. Seelamantula, and R. Balakrishnan, "Quantifying vocal mimicry in the greater racket-tailed drongo: a comparison of automated methods and human assessment," *PloS one*, vol. 9, no. 3, p. e89540, 2014.
- [19] E. Loweimi, S. M. Ahadi, and T. Drugman, "A new phase-based feature representation for robust speech recognition," in *Proc. Int. Conf. Acoust. Speech, Signal Process.*, 2013.
- [20] "Cassin's vireo recordings," <http://taylor0.biology.ucla.edu/al/bioacoustics/>, accessed: 2016-03-20.
- [21] "Mlsp bird classification challenge 2013," <https://www.kaggle.com/c/mlsp-2013-birds/data>, accessed: 2016-03-20.