



**HAL**  
open science

## Hidden Markov Modeling for humpback whale (*Megaptera novaeangliae*) call classification

Federica Pace, Paul White, Olivier Adam

► **To cite this version:**

Federica Pace, Paul White, Olivier Adam. Hidden Markov Modeling for humpback whale (*Megaptera novaeangliae*) call classification. Acoustics 2012, Apr 2012, Nantes, France. hal-00810807

**HAL Id: hal-00810807**

**<https://hal.science/hal-00810807>**

Submitted on 23 Apr 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# ACOUSTICS 2012

## Hidden Markov Modeling for humpback whale (*Megaptera novaeangliae*) call classification

F. Pace<sup>a,b</sup>, P. R. White<sup>a</sup> and O. Adam<sup>c,d</sup>

<sup>a</sup>ISVR, University of Southampton, Southampton, UK, SO17 1BJ Southampton, UK

<sup>b</sup>Baker Consultants, Cromford Station, DE4 5JJ Matlock Derbyshire, UK

<sup>c</sup>Institut Jean le Rond d'Alembert - Université Pierre et Marie Curie-Paris VI, 11 rue de  
Lourmel 75015 Paris

<sup>d</sup>CNRS, CNRS UMR 8195, Université Paris Sud, Bâtiment 446, 15, rue Georges Clemenceau,  
91405 Orsay Cedex, France

fp@isvr.soton.ac.uk

This study proposes a new approach for the classification of the calls detected in the songs with the use of Hidden Markov Models (HMMs) based on the concept of subunits as building blocks. HMMs have been used once before for such task but in an unsupervised algorithm with promising results, and they are used extensively in speech recognition and in few bioacoustics studies. Their flexibility suggests that they may be suitable for the analysis of the varied repertoire of humpback whale (*Megaptera novaeangliae*) calls because they cope well with variations in the call durations, which is a common feature in humpback whale vocalizations. Another attractive characteristic of HMMs is that highly developed tool-set is widely available as a consequence of the widespread use of their employment for human speech analysis.

We describe the HMM classification method and show that a high level of performance can be achieved with modest requirements both in terms of computational load and storage. Training stage requires minimal manual input and once trained the recognition process is fully automated. We will present how the classification performance is affected by different amount of training.

## 1 Introduction

Songs of humpback whales (*Megaptera novaeangliae*) have been extensively studied for the past four decades, since they were defined as such by Payne and McVay (1971), who noticed that the sounds emitted by humpback whales recorded in Hawaii followed a patterned sequence. Since then, the songs produced by these baleen whales have been studied across the World (Winn *et al.*, 1981; Helweg *et al.*, 1998; Noad *et al.*, 2000; Cerchio *et al.*, 2001; Razafindrakoto, 2001; Suzuki *et al.*, 2006; Whitlow *et al.*, 2006) because humpback whales breeding grounds, where songs are typically heard, are widespread. As a consequence of the vast amount of data collected through recording humpback whale songs in the wild, the need has arisen for the development of appropriate tools for the automatic classification of the song components. This would allow large scale comparisons of songs across populations and from year to year, which is necessary to understand how songs are culturally transmitted and learn about their migratory patterns (Noad *et al.*, 2000; Mercado III *et al.*, 2005; Oviedo *et al.*, 2008; Garland *et al.*, 2011). The task of song classification is still largely carried out manually or with the use of algorithms that require substantial human supervision, which is extremely time-consuming and not easily replicable when comparing songs analysed across research groups.

Our new approach for song classification of humpback whales using Hidden Markov Models (HMMs) showed high level of classification of songs recorded in Madagascar between 2007 and 2009 (Pace *et al.*, 2010; Pace *et al.*, 2011). HMMs have been used once before to model humpback whale calls, but with the implementation of an unsupervised algorithm (Rickwood and Taylor, 2008). The power of HMMs derives from their ability to model non-stationary random processes, specifically, they are particularly appropriate when modelling signals, such as humpback whale vocalisations, whose durations are stochastic. HMMs have become the basis of most modern algorithms for the classification of human speech (speech recognition) (Deller *et al.*, 1993). This central role in speech recognition (and their wider use in the field of speech analysis) has meant that considerable research effort has been dedicated to the study of HMMs, one consequence of which is a highly developed, and widely available, tool-set. This makes them attractive tools for application in a wide range of fields, including bioacoustics (Brown and Smaragdis, 2009; Ren *et al.*, 2009).

Speech models are well suited to describing the mechanisms of our vocal apparatus and our hearing, and through the extensive research that has been carried out on the matter, sentences and words can be characterised and classified with a high level of accuracy. The fact that we aim to build a classifier for humpback whale song that mimics the accuracy of a trained human listener justifies the adoption of processing tools that have been developed based on the human perception of speech. The underlying idea is that the model is tuned to the way humans perceive whale vocalisations, given that we can classify their songs accurately in a biological significant way.

Models borrowed from speech production have already been implemented for the classification of humpback whale songs (Mercado III and Kuh, 1998; Mercado III *et al.*, 2010). One such instance is the source-filter model where vocalisations are modelled a system of voiced sounds and unvoiced sounds that are produced by a source, i.e. the voice-box. The sound then is filtered using an all-pole infinite impulse response filter (IIR) (Deller *et al.*, 1993). This model is widely accepted to describe the mechanisms of speech production but it accounts only for the flow of air from the source to one filtering chamber and its propagation out to the environment. In the case of sound production in baleen whales, a more complex model is required because more mechanisms are likely to be involved. Although the specific pathways of sound production and propagation in baleen whales remain to be understood, it is recognised that air must be recycled within the vocal tract to allow continuous production of sound underwater and considering the lack of bubble emission during sound generation (Reindenberg and Laitman, 2007; Mercado III *et al.*, 2010).

In this study, we assess the performance of HMMs for the classification of humpback whale songs using different levels of algorithm training. Indeed, the goal is to maximise the recognition performance of the individual calls present in a song sequence, whilst minimising the amount of training data needed, so that human input is reduced as are time consumption and computational load.

## 2 Methods

### 2.1 Data collection and pre-processing

Humpback whale songs were recorded in the Ste. Marie Island Channel which is located between the Island of Ste. Marie and the North East Coast of Madagascar (Indian Ocean). Whales are present in this area during the winter months, i.e. June to October, and come from Antarctica for

purposes of breeding. Data were recorded from a small boat using a single CO.L.MAR Italia GP280 hydrophone connected to a TASCAM HD-P2 recorder which recorded data at a sampling frequency of 44.1 kHz and digitalized using 16 bits. The hydrophone was located at a depth of approximately 20 m for all recordings and the bathymetry ranged between 28 to 40 m.

Prior to classification, the song analysed in this study was segmented into its component units, i.e. continuous sounds between two silences, as defined by Payne and McVay (1971). This segmentation was performed automatically using an energy detector with a double threshold, and then refined manually to ensure that the limits of the start and the end of each call were accurate. Although the classification might as well work with partially detected calls where the start and end thresholds do not correspond precisely to the start and the end of the call, we preferred to adjust the markers manually to ensure that the performance of the classifier was not affected by incorrect segmentation of certain calls, and hence maximise the performance outcome.

## 2.2 Feature extraction

The individual vocalisations detected during the song segmentation stage were not directly inputted in the classification algorithm; instead they were reduced to a series of coefficients that described the essential characteristics of the call. The efficiency of three feature sets that are commonly adopted to represent bioacoustics signals was tested using our Madagascar recordings in a previous study on humpback whale call classification (Pace, *et al.*, 2009). These included Linear Prediction Coefficients (LPCs), coepstrum, and Mel-frequency coepstrum coefficients (MFCCs). The results showed that MFCCs performed better than the other two feature sets for nearly all call types despite the fact that there are based on an anthropomorphic perception of sound. Hence, they were chosen for the feature extraction stage of our classification algorithm.

Each call was represented through a series of 24 features; specifically these are 12 MFCCs and 12 corresponding delta coefficients (i.e.  $\Delta$ MFCCs). The number of coefficients was selected to optimise performance, as described in (Pace *et al.*, 2011). These coefficients were calculated using the standard function included in the HTK toolkit (Young *et al.*, 2000) which was used for the HMM implementation. This toolkit was chosen because it has been highly developed by Cambridge University and applied for several bioacoustics studies, as well as speech recognition tasks (Ren *et al.*, 2009).

## 2.3 HMM implementation

Each call class is represented by one left to right HMM with one state if the call frequency is stable throughout its duration or two states if the frequency is varying, e.g. in the case of an up-sweep or down-sweep, plus two “non-emitting” states at the start and at the end of each model. A definition file was therefore created for each HMM which includes a general description of the features, i.e. type and vector size of the feature sets, the number of states and the transition probabilities matrix. The number of states of each HMM determines the size of the transition matrix: a three state HMM will have a  $3 \times 3$  transition matrix.

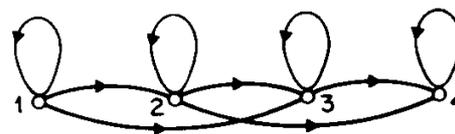


Figure 1: Bakis diagram of left to right HMM with four transition states. States 1 and 4 represent the start and end of each call and are termed ‘non-emitting’ whilst states 2 and 3 are termed ‘emitting states’, and represent the transitions within states of a single call. The HMM allows for transitions skipping one state as well as visiting all of the states in succession.

The transition between observations was modelled by a Gaussian mixture where, by definition, each probability was a real number from 0 to 1 and sum to unity. Given that each recording segment containing a call could include a silent part at the start and/or at the end of the sound clip, one HMM was created with the same characteristics described above to model the silences.

During the training phase a database of labelled (manually classified) data is employed. The manual classification is performed by an experienced observer, through visual inspection of the spectrograms and aurally.

The training stage deals with calculating the maximum likelihood estimates (MLE) of the transition probabilities matrix of the states. In practice, this means that starting from a prototype HMM after the training process one obtains a model whose mean, variance and transition probabilities are calculated based on the statistical properties of the data present in the training set. This is achieved in two steps:

i) The Viterbi algorithm (Forney, 1978) is used to find the most likely state sequence corresponding to each training sample;

ii) A Baum-Welch (Baum *et al.*, 1970) re-estimation is performed to find the probability of being in each state at each time frame using the *Forward-Backward* algorithm. This probability is then used to form weighted averages for the HMM parameters. A thorough review of the use of HMMs is provided by Rabiner (1989).

For the recognition stage, a Viterbi alignment (Viterbi, 1967) was performed to match each call of the testing dataset the best matching HMM. The output of the HMM recognition was then compared to the manual classification and the correction classification rate computed as a percentage.

In this study, three scenarios of training were analysed: initially the recognition performance was tested by training the HMM on 50% of the calls for each call class, up to a maximum of 12 training calls per class. Then we tested the HMM with 25% of data training, and lastly 10% data training was used. In the latter case, the amount of training data per class was sometimes slightly above 10% because the minimum number of samples required to train each HMM is 3 calls. A table showing the number of calls trained in each class and the recognition results is presented in the results section.

### 3 Results

The performance 2008 based on the recognition of individual calls of the HMM model for the classification of a humpback whale song recorded in Madagascar is presented in this section.

The calls identified during the segmentation were manually classified as described in the methods and named alphabetically, according to the sequence in which they were encountered in the song. Hence the first unit type was termed a, the second was termed b, and so on.

The song recorded in Madagascar in 2008 was segmented into 334 units, which were divided into 16 classes. Three classes were omitted from the analysis because they contained fewer than 6 calls, i.e. the minimum number to be able to run both the training and testing stages of the algorithm. The call types identified are presented in the table below, as well as the number of calls used to train each category for each training scenario (Table 1).

Table 1: table showing the call types identified in the recording analysed, as well as the number of calls used during the training stage for each of the training scenario. The classification performance for each of the scenarios is presented as a percentage of the total number of calls tested for each call type. The training set number denoted by a ‘\*’ mean that the actual number of calls used for the training stage should have been less than three if we calculated the appropriate percentage of calls for the training scenario; however, we had to train the HMM with 3 calls which is the minimum number required for running the algorithm. Also note that the number of calls used for the training was rounded to the nearest integer.

Call type	Training			Classification performance (% out of 181)		
	50%	25%	10%	50%	25%	10%
a	12	11	4	97	100	84
b	10	5	*	100	100	100
d	6	3	*	86	100	86
f	7	4	*	100	50	100
g	6	3	*	100	100	100
h	10	5	*	100	60	10
i	10	10	4	97	91	85
j	8	4	*	86	57	14
k	10	6	*	88	81	38
l	5	*	*	100	100	100
m	10	6	*	100	93	93
n	9	8	3	87	96	91
o	7	4	*	88	100	88
overall	110	72	41	94	90	78

The results show that the best performance overall was achieved using 50% of the data for training the HMMs and

the other 50% for testing the classifier; however, this was not true of all the call classes tested Figure 2.

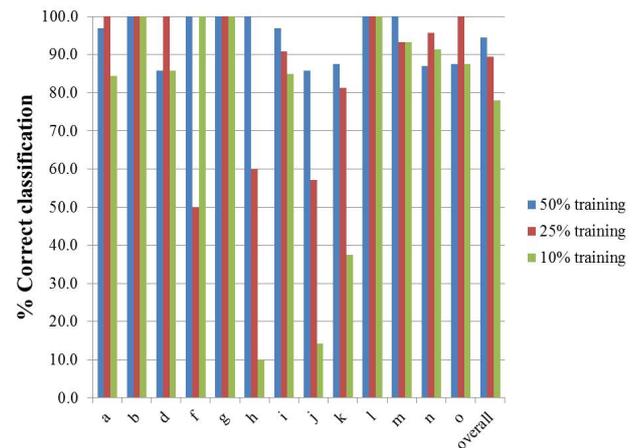


Figure 2: percentage correct classification of the Hidden Markov Modelling classification obtained for three different training scenarios for each call type (or unit type) and overall.

With a 25% percentage reduction in training data, the classification performance decreases only by 4% but the mistakes affect the various call types differentially. Indeed, in three instances, namely units ‘f’, ‘h’ and ‘j’, the classification accuracy halved (or nearly halved). On the other hand, there are 3 instances in which more units were correctly classified when there were 25% rather than 50% calls used for training.

In the last training scenario, when the HMMs were trained using only 10% of the data (or slightly more) the overall classification performance reduced to 78%. Again here some call types were more affected than others by the change in training set size. Specifically, units ‘h’, ‘j’ and ‘k’ were classified very poorly (<40% correct classification), whilst the classification of the other unit types was nearly equal to the one obtained with the other training scenarios.

### 4 Discussion

High levels of classification performance were obtained using Hidden Markov Modelling for classifying humpback whale calls, as demonstrated in our previous work which compared classification performance across songs of different years and emitted by a variety of singers (Pace *et al.*, 2010; Pace *et al.*, 2011). Whilst, in this study we did not compare songs from different years, we aimed at analysing in more detail how different amount of training affects the classification performance. For an automatic classifier to be efficient and widely used, the amount of training required to run the algorithm should be minimized to reduce the human input, which introduces subjectivity making it hard to replicate studies across research groups, and decrease the computational load so that the whole recognition task can be implemented quickly even with large datasets.

The results presented in this study show how three different training scenarios affect the performance of the automatic classification; the information that can be extrapolated can help choosing which scenario is better for the task that one needs to perform, considering the trade-off

between amount of time and human input required at the training stage and the performance outcome.

The data showed that the largest training set size led to a higher classification performance, with calls being correctly classified in more than 85% of the cases for all call types. Decreasing the training set led to a reduced classification performance, but this decrease was not linear and affected different call types differentially. The smallest training sample size resulted in an overall decrease of 16% in classification performance compared to the 50% training set scenario. Whilst this is not a huge decrease, it can be quite conspicuous when analysing large amounts of data, as is customary when dealing with humpback whale song classification tasks. In addition, the recording used for the analysis had quite a high signal to noise ratio (SNR), which is difficult to achieve for continuous recordings taken in the field. We expect the classification performance to be worse when the quality of the recordings is lower, and a larger amount of training being required in these cases so that the HMMs can recognize the characteristics of the original signal, rather than the artefacts of the noise that may be present.

The fact that some call types were more affected than others by the change in training set size suggests that training should be tuned to the type of call. Given that humpback whale songs are composed of units that vary considerably in characteristics, it is feasible that different types of calls may need different amount of training. Indeed, the repertoire of humpback whale includes tonal harmonic calls, broadband sounds, and fast up-sweeps and down-sweeps (Thompson *et al.*, 1977; Dunlop *et al.*, 2007). Considering that with the 50% training scenario, the classification performance was similar across call types, one can conclude that the differential performance is not due to the performance of the feature set used. This could have been a possibility given that MFCCs are based on the Fourier representation of the signals, and therefore are particularly suited for characterising harmonic sounds. The unit types that were most affected by changes in the amount of training data used were either broadband calls (units 'h' and 'j') or calls where sudden changes in frequency could be observed (units 'f' and 'k'). This is unsurprising considering that the few calls present in the training set might differ from one another and not give an accurate enough representation of the characteristics of the other calls that belong to the same class that were present in the test set.

Further work will consist in comparing training data sets sizes for larger test data sets, and to test songs emitted by different singers in one or more years to check if the results are consistent with these findings. We would expect a larger amount of training being required for correct classification of calls emitted by different singers to account for individual variability in the sound characteristics. In addition, the same study will be extended to the classification task based on the segmentation of songs into smaller building blocks, which we defined as subunits (Pace *et al.*, 2010; Pace *et al.*, 2011). Subunit segmentation was also proposed for killer whale calls (Shapiro *et al.*, 2011). We would expect the differential response across call types to be greatly reduced when classifying songs based on subunits because the calls identified in such categories are more stable in the frequency domain because where sudden frequency shifts are observed in a unit, this will lead to splitting it into two (or more) subunits.

## 6 Conclusion

Hidden Markov Models are well suited and easily adaptable for the classification of humpback whale calls. The classification performance with three training scenarios performed in this study suggests that, as expected, a larger training set leads to more accurate classification; however, given that halving the amount of training required, leads to only a 4% decrease in performance, one could choose to favour this scenario to reduce human input and effort considerably. Further analysis is required to consolidate the results and test the performance on larger test sets and songs emitted by different whales.

## Acknowledgments

The authors wish to thank the staff of CetaMada, and Princesse Bora Lodge and Spa in Madagascar for the financial and technical support during the field seasons.

## References

- [1] Baum, J. F., Petrie, T., Souler, G., and Weiss, N. (1970). "A maximisation technique occurring in the statistical analysis of probabilistic functions of Markov chains," *Am. Math. Statist.* **41**, 164-171.
- [2] Brown, J. C., and Smaragdis, P. (2009). "Hidden Markov and Gaussian mixture models for automatic call classification," *J. Acoust. Soc. Am.*, EL221-EL224.
- [3] Cerchio, S., Jacobsen, J. K., and Norris, T. F. (2001). "Temporal and geographical variation in songs of humpback whales, *Megaptera novaeangliae*: synchronous change in Hawaiian and Mexican breeding assemblages," *Animal behaviour* **62**, 313-329.
- [4] Deller, J. R., Proakis, J. G., and Hansen, J. H. L. (1993). *Discrete-time processing of speech signals* (Macmillian Publishing Company, New York).
- [5] Dunlop, R. A., Noad, M. J., Cato, D. H., and Stokes, D. (2007). "The social vocalization repertoire of east Australian migrating humpback whales (*Megaptera novaeangliae*)," *The Journal of the Acoustical Society of America* **122**, 2893-2905.
- [6] Forney, G. D. (1978). "The Viterbi Algorithm," *Proceedings of the IEEE* **61**.
- [7] Garland, Ellen C., Goldizen, Anne W., Rekdahl, Melinda L., Constantine, R., Garrigue, C., Hauser, Nan D., Poole, M. M., Robbins, J., and Noad, Michael J. (2011). "Dynamic Horizontal Cultural Transmission of Humpback Whale Song at the Ocean Basin Scale," *Current biology : CB* **21**, 687-691.
- [8] Helweg, D. A., Cato, D. H., Jenkins, P. F., Garrigue, C., and McCauley, R. D. (1998). "Geographic Variation in South Pacific Humpback Whale Songs," *Behaviour* **135**, 1-27.
- [9] Mercado III, E., Herman, L., and Pack, A. (2005). "Song copying by humpback whales: themes and variations," *Animal Cognition* **8**, 93-102.
- [10] Mercado III, E., and Kuh, A. (1998). "Classification of humpback whale vocalizations using a self-organizing neural network," in *IEEE World*

- Congress on Computational Intelligence*, edited by N. N. Proceedings.
- [11] Mercado III, E., Schneider, J. N., Pack, A. A., and Herman, L. M. (2010). "Sound production by singing humpback whales," *Journ. Acous. Soc. Am.* **127**, 2678-2691.
- [12] Noad, M. J., Cato, D. H., Bryden, M. M., Jenner, M. N., and Jenner, K. C. S. (2000). "Cultural revolution in whale songs," *Nature* **408**, 537-537.
- [13] Oviedo, L., Guzman, H. M., Florez-Gonzalez, L., Alzueta, J. C., and Mair, J. M. (2008). "The Song of the Southeast Pacific Humpback Whale (*Megaptera novaeangliae*) off Las Perlas Arcipelago, Panama: Preliminary Characterization," *Aquatic Mammals* **34**, 458-463.
- [14] Pace, F., Benard, F., Glotin, H., Adam, O., and White, P. (2010). "Subunit definition and analysis for humpback whale call classification," *Applied Acoustics* **71**, 1107-1112.
- [15] Pace, F., White, P. R., and Adam, O. (2011). "Hidden Markov Models for classification of humpback whale songs collected in Ste Marie, Madagascar, over three years.," *J. Acoust. Soc. Am.* **In Review**.
- [16] Payne, R. S., and McVay, S. (1971). "Songs of Humpback Whales," *Science* **173**, 585-597.
- [17] Rabiner, L. R. (1989). "A tutorial on Hidden Markov Models and selected applications in speech recognition," in *IEEE*, pp. 257-286.
- [18] Razafindrakoto, Y. (2001). "First description of humpback whale song from Antongil Bay, madagascar," *Marine Mammal Science* **17**, 180-186.
- [19] Reindenberg, J. S., and Laitman, J. T. (2007). "Discovery of a low frequency sound source in Mysticeti (baleen whales): anatomical establishment of a vocal fold homolog," *The Anatomical Record* **290**, 745-759.
- [20] Ren, Y., Johnson, P. C., Darre, M., Suart Glaeser, S., Osiejuk, T. S., and Out-Nyarko, E. (2009). "A Framework for bioacoustic Vocalization Analysis Using Hidden Markov Models," *Algorithms*, 1410-1428.
- [21] Rickwood, P., and Taylor, A. (2008). "Methods for automatically analyzing humpback song units," *The Journal of the Acoustical Society of America* **123**, 1763-1772.
- [22] Shapiro, A. D., Tyack, P. L., and Seneff, S. (2011). "Comparing call-based versus subunit-based methods for categorizing Norwegian killer whale, *Orcinus orca*, vocalizations," *Animal behaviour* **81**, 377-386.
- [23] Suzuki, P., Buck, J. R., and Tyack, P. L. (2006). "Information entropy of humpback whale songs," *J. Acoust. Soc. Am.* **119**, 1849-1866.
- [24] Thompson, P. O., Cummings, W. C., and Kennison, S. J. (1977). "Sound production of humpback whales, *Megaptera novaeangliae*, in Alaskan waters," *J Acoust Soc Am* **62**, S89.
- [25] Viterbi, A. J. (1967). "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," *IEEE Trans on Inf Th* **13**, 260-269.
- [26] Whitlow, W. L. A., Adam, A. P., Marc, O. L., Louis, M. H., Mark, H. D., and Kim, A. (2006). "Acoustic properties of humpback whale songs," *The Journal of the Acoustical Society of America* **120**, 1103-1110.
- [27] Winn, H. E., Thompson, T. J., Cummings, W. C., Hain, J., Hudnall, J., Hays, H., and Steiner, W. W. (1981). "Song of the humpback whale — Population comparisons," *Behavioral Ecology and Sociobiology* **8**, 41-46.
- [28] Young, S., Kershaw, D., Odell, J., Ollason, D., Valtchev, V., and Woodland, P. (2000). "The HTK book," (Microsoft Corporation).
- [29]
- [30]