

A STATISTICAL MULTIDIMENSIONAL HUMMING TRANSCRIPTION USING PHONE LEVEL HIDDEN MARKOV MODELS FOR QUERY BY HUMMING SYSTEMS

Hsuan-Huei Shih, Shrikanth S. Narayanan and C.-C. Jay Kuo

Integrated Media Systems Center and Department of Electrical Engineering
University of Southern California, Los Angeles, CA 90089-2564
E-mail: {hshih,shri,cckuo}@sipi.usc.edu

ABSTRACT

A new phone level hidden Markov model approach applied to human humming transcription is proposed in this research. A music note has two important attributes, *i.e.* pitch and duration. The proposed system generates multidimensional humming transcriptions, which contain both pitch and duration information. Query by humming provides a natural means for content-based retrieval from music databases, and this research provides a robust front-end for such an application. The segment of a note in the humming waveform is modeled by phone level hidden Markov models (HMM). The duration of the note segment is then labeled by a duration model. The pitch of the note is modeled by a pitch model using a Gaussian mixture model. Preliminary real-time recognition experiments are carried out with models trained by data obtained from eight human objects, and an overall correct recognition rate of around 84% is demonstrated.

1. INTRODUCTION

Querying music databases using human humming as the query key has recently gained attention as a viable option [1]. This requires signal processing for automatically mapping human humming waveforms to symbol strings representing the underlying melody and duration contours. This paper focuses on automatic humming recognition and transcription.

Most approaches to humming analysis, mainly developed for query by humming systems, use non-statistical signal processing. Ghias *et al.* [1], and Jang *et al.* [2] used autocorrelation to calculate pitch periods. McNab *et al.* [3, 4] adopted the Gold-Rabiner algorithm. A major problem with the non-statistical approach is its robustness to inter-speaker variability and other signal distortions. Another problem with the proposed non-statistical approaches is the lack of real-time processing capability. Most of these methods rely on full utterance level feature measurements that demand buffering of long humming data, thereby limiting real-time processing.

To overcome the above two problems, we proposed an HMM-based humming recognition system in [5]. The HMM-based humming recognition system transcribed a humming piece into a melody contour. Some similar work was performed by Raphael [6, 7] and Durey *et al.* [8, 9]. Raphael [6, 7] attempted to solve the segmentation problem of automatic musical accompaniment using HMMs. Durey *et al.* [8, 9] used HMMs to spot melody in music pieces. Both of them were based on instrumental or MIDI generated music pieces. These music pieces were accurate in tune and rhythm. Our proposed system aims at dealing with human hum-

ming which is much less accurate than instrumental music. The proposed system detects not only when a note is hummed for how long but also which note is hummed.

To enhance our previously proposed system in [5], we focus on recognizing both the melody contour and the duration contour of a piece of humming in this research. These two contours are combined to form a multidimensional humming transcription. The multidimensional transcription is assumed to be a sequence of notes in a piece of humming defined by their detected duration changes and pitch intervals. A note model is defined by phone level hidden Markov models with features modeled by Gaussian mixture models. A pitch model of a note is defined by pitch features and modeled by Gaussian mixture models. During the training phase, note and pitch models are trained with humming data obtained from real people. During recognition, the incoming piece of the humming waveform is decoded with trained note models for note segmentation, and then pitch values of segmented notes are detected with trained pitch models. The detection process is done in real-time.

The rest of the paper is organized as follows. In Sec. 2, the proposed algorithm is described in detail. Experimental results are given in Sec. 3 and conclusions and future work are presented in Sec. 4.

2. PROPOSED ALGORITHM

Our approach to multidimensional humming transcription is summarized in Fig. 1. Similar to any data-driven pattern recognition approach, models are derived from data representing the underlying classes for recognition. Details of database preparation are given in Sec. 2.1. The proposed algorithm can be divided into two stages. At the first stage, a humming piece is first passed into the note decoder for note segmentation, and a duration label of the segmented note is given at this stage. At the second stage, a segmented note of the humming piece is then passed to the pitch detector for pitch tracking. The phone level statistical models, selected features, training process, decoding process, and labeling process of the note segmentation stage are addressed in Sec. 2.2. Pitch feature selection, pitch analysis, and pitch model generation of the pitch tracking stage are described in Sec. 2.3. Finally, the generation of a multidimensional humming transcription is given in Sec. 2.4.

2.1. Database Collection and Preparation

There is no publicly available database for human humming. A small humming database was created in our previous work [5] for

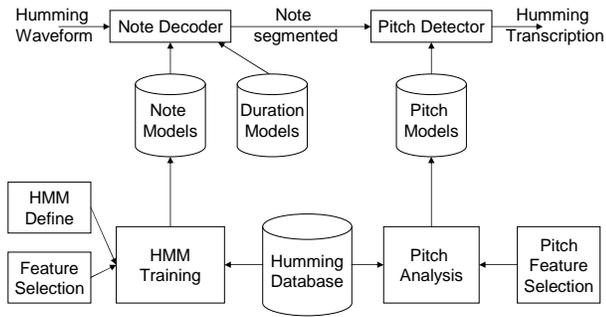


Fig. 1. The functional block diagram of the humming recognition system.

system design and preliminary experiments. The same database is used here and briefly described below. Humming pieces were recorded and sampled at 16kHz using high-quality close talking microphone. A humming piece contains a sequence of notes, which were labeled by human listeners. Manual segmentation of notes was included, since it provided the boundary information for pitch model training. A note has two main attributes, namely, the pitch (measured by the fundamental frequency of voicing) and the duration. Since the goal of this work is to map the humming signal to notes, a more robust and general method that focuses on relative changes in pitch values of a melody contour is adopted. In other words, relative pitch values were used to label a humming piece instead of absolute pitch values in our work.

The labeling convention is based on the rationale that a human is sensitive to the pitch value of adjacent notes rather than the first note. According to this convention, the humming piece for *do-re-mi-fa* will be labeled as “R-U2-U2-U1”, and a humming piece corresponding to *do-ti-la-sol* is labeled as “R-D1-D2-D2” where “R” is the reference note or no change in the pitch value, “U2” denotes a pitch value higher than the reference by two half-steps and “D1” denotes a pitch value lower than the reference by one half-step. The numbers following “D” or “U” are variable, depending on the humming data. Transcriptions were saved in separate files and used during supervised training of note and pitch models and to provide the reference transcriptions to evaluate recognition results.

2.2. Note Segmentation

The first stage of the proposed algorithm is note segmentation, where the process of segmenting notes of a humming piece is conducted. First, a feature set which can characterize a note is chosen. Next, the HMM definition is chosen before training. During the training phase, notes’ phone level HMMs are trained using the selected feature set. The trained note models are then used by the note decoder for note segmentation. Finally, the duration of a segmented note is labeled according to its relative duration change.

A. Feature Selection

The choice of good features is the key to good humming recognition performance. Since human humming production is similar to speech, features used to characterize a phoneme in automatic speech recognition (ASR) are considered for modeling notes in humming recognition. Features used in our base feature set include

mel-frequency cepstral coefficients (MFCC), energy measures and their first- and second-derivatives.

Mel-Frequency Cepstral Coefficients (MFCCs) is used to characterize the acoustical shape of humming notes. For our analysis, 26 filterbank channels are chosen, and the first 12 MFCCs are selected as features.

Energy is an important feature in humming recognition especially to provide temporal segmentation of notes. Typically, a distinct variation in energy will occur during the transition from one note to another. This effect is especially enhanced since users are asked to hum using basic sounds that are a combination of a stop consonant and a vowel (e.g., “da”, “la”).

The 39-element feature vector contains 12 MFCCs, 1 energy measure, and their first and second derivatives.

B. Phone Level Hidden Markov Model

Phone level hidden Markov models (HMM) with Gaussian mixture models (GMM) for observations corresponding to each state of an HMM are used to characterize a note model. Each note is modeled by 3-state left-to-right phone level HMMs. The use of HMM provides the ability to model temporal aspects of a note especially in dealing with time elasticity. The features corresponding to each state occupation in an HMM are modeled by a mixture of 2 Gaussians.

The idea of using phone level HMM for a humming note is very similar to speech recognition. The word level note models initially proposed in [5] limited the ability to better characterizing a humming note (consider a humming using the syllable “da”). Since a stop consonant and a vowel have quite different acoustical characteristics, two distinctive phone level HMMs are defined for “d”, and “a”. The HMM of “d” is used to model the stop consonant of a humming note. The HMM of “a” is used to model the vowel of a humming note. A humming note is then represented by the HMMs of “d” followed by “a”.

A new silence model is designed to improve the robustness. Back ground noise and other distortion may cause unwanted segmentation of notes. In the new silence model, an extra transition from state 1 to 3 and then from state 3 to 1 is added in the original 3-state left-to-right HMM model. By doing so, the silence model can allow each model to absorb the various impulsive noise. When an impulsive noise comes in, the backward skip provides a mechanism to absorb the noise without exiting the silence model. At this point, a 1 state short pause “sp” model is created. This is called the “tee-model”, which has a direct transition from the entry to the exit node. The emitting state is tied with the center state (state 2) of the new silence model. The topology of the new silence model is shown in Fig. 2. A “Rest” in a melody is then represented by the HMM of “Silence”.

C. Training Process

The parameters of HMMs are estimated during a supervised training process using the maximum likelihood approach with Baum-Welch re-estimation. An initial 3-state left-to-right HMM silence model is used at the first two Baum-Welch iterations to initialize the silence model. The tee-model (“sp” model) extracted from the silence model and a backward 3-to-1 state transition are added after the second Baum-Welch iteration.

D. Note Decoding and Duration Labeling

The Viterbi decoding algorithm is used in the decoding process. The recognition problem is to find a state sequence of a model which is most likely to have been generated by the data.

After a note is segmented, the duration of the segment is labeled by a duration label according to the duration model. Instead

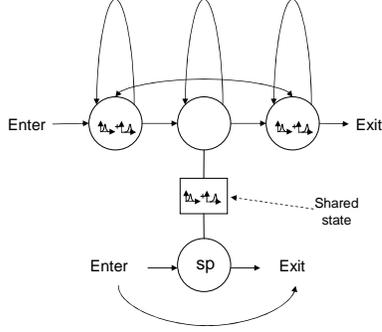


Fig. 2. The new silence model with a one state short pause “sp” model tied to the center state (*i.e.* state 2).

of directly using the absolute length of a duration, the relative duration change is used in the labeling process. The relative duration change of a note is based on its previous note. The relative duration is calculated as

$$\text{relative duration} = \log_2\left(\frac{\text{current duration}}{\text{previous duration}}\right) \quad (1)$$

The system assumes that the shortest note of a humming piece is the thirty second note. Total of 11 duration models which are -5, -4, -3, -2, -1, 0, 1, 2, 3, 4, 5 cover possible differences from a whole note to a thirty second note. The duration label of a note segment is then represented by the duration model which is closest to calculated relative duration value. The first note’s duration label is labeled as “0”, since no previous reference note exists.

2.3. Pitch tracking

After a note is segmented from a humming piece, it is passed to the second stage to decide its pitch. The pitch detector decides the pitch of a segmented note based on the statistical information of pitch models. The statistical information of pitch models is obtained from the humming database off-line. The detailed implementation of each component of the pitch detector is given below.

A. Pitch Analysis

Short-time autocorrelation is chosen for pitch analysis. The main advantage of using short-time autocorrelation is its relative low computational cost in comparison with other existing pitch detection algorithms. A frame size of 20 msec with 10 msec frame overlap was adopted throughout our experiments. The frame-based analysis is performed on a note segment, which usually has several frames. Multiple frames of a segmented note are used for pitch model analysis. After applying autocorrelation to those frames, pitch features are extracted. The selected pitch features include: the first harmonic, the pitch median, and the pitch log standard deviation.

The first harmonic is also known as the fundamental frequency or the pitch. The first harmonic provides the most important pitch information. Because of noise, some frames’ pitch values are very different when compared with other frames’ pitch values in the same note segment. Taking their average is not a good choice, since distant pitch values move the mean to the location where it is away from the target value. The median pitch value of a note segment proves to be a better choice in our experiments.

The outlying pitch values also impact the standard deviation of a note segment. To overcome this problem, these outlying pitch

values should be moved back to the range where most pitch values belong to. Since the smallest interval of two different notes is a semitone, we claim that pitch values different from the median value by more than one semitone have a significant drift. The pitch values drifted by more than a semitone are moved back to the median. Then, the standard deviation is calculated. Pitch values of notes are not linear in the frequency domain. In fact, they are linearly distributed in the log frequency domain. Therefore, calculating the standard deviation in the log scale makes more sense. The log pitch mean and the log standard deviation of a note segment are calculated.

B. Pitch Model

Pitch models are used for different pitch intervals which are defined to be the difference in semitones of two adjacent notes

$$\text{pitch interval} = \frac{\log(\text{current pitch}) - \log(\text{previous pitch})}{\log \sqrt[12]{2}} \quad (2)$$

The pitch models cover 2 octaves of pitch intervals, which are from D12 semitones to U12 semitones. A pitch model has two attributes: the length of the interval (in terms of the number of semitones) and the pitch log standard deviation in the interval. The two attributes are modeled by the Gaussian mixture model. The boundary information and the ground truth of a pitch interval were obtained from manual transcription. The calculated pitch intervals and log standard deviations, which correspond to the ground truth pitch interval, are collected. A 2-D Gaussian mixture model is generated based on the collected information.

C. Pitch Detector

The pitch detector decides the pitch change of a segmented note with respect to its previous note. The first note of a humming piece is always marked as the reference note, and its detecting is in principle not needed. However, the first note’s pitch is still calculated as reference in our experiments. The later notes of the humming piece are detected by the pitch detector. The pitch intervals and the pitch log standard deviations are calculated. They are used to select the best model that gives the maximum likelihood value as the detected result.

2.4. Transcription Generation

After the note segmentation stage and the pitch detection stage, a humming piece has all the information required for transcription. The transcription of the humming piece results in an sequence of length N with two attributes per symbol, where N is the number of notes. The two attributes are the duration change (or relative duration) of a note and the pitch change (or pitch interval) of a note. The “Rest” note is labeled as “Rest” in the pitch interval attribute, since they do not have a pitch value.

3. EXPERIMENTAL RESULTS

The proposed algorithm consists of two stages: note segmentation and pitch detection. For the note segmentation stage, the 39-element feature vector consists of 12 MFCCs, 1 energy measure and their first and second derivatives. The frame size was chosen to be 20 msec, and the frame skip 3 msec (which means two consecutive frames have an overlap of 17 msec.). The detail of choosing the frame size and the frame rate was given in [5]. A 3-state left-to-right HMM with 2 GMMs was used to model notes, and each model was trained 10 times. For the pitch detection stage, the frame size was chosen to be 20 msec and the frame skip 10 msec.

The number of Gaussian mixtures was set to one since the data set was quite limited.

A Graphical User Interface (GUI) program, called HTKedit, was written based on the Hidden Markov Model Toolkit (HTK). The program can be used to train note models and pitch models. It has a pitch analysis tool for the pitch model study. It can also take the trained note and pitch models to perform both off-line and real-time transcription. It can further convert a N-by-2 transcription into a music score. Fig. 3 is the screenshot of the HTKedit program and demo can be found at URL:<http://sail.usc.edu/music>.

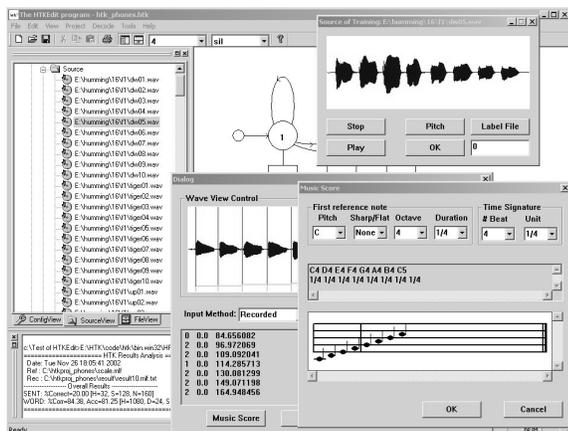


Fig. 3. The screenshot of HTKedit with a humming recognizer and a music score editor.

Let us define the following parameters: 1) N: the no. of correct notes, 2) D: the no. of deletion errors, 3) S: the no. of substitution errors, and 4) I: the no. of insertion errors. Two performance measures, i.e. the correct recognition rate (CRR) and the accuracy rate (AR), are adopted for comparison. They are defined as

$$CRR = (N - D - S)/N, \quad AR = (N - D - S - I)/N.$$

Both off-line and real-time recognition experiments were performed. The off-line recognition was conducted with the leave-one-out method (trained on 7 speakers and tested on 1). Among seven speakers, the average AR was 81.25% and the average CRR was 84.38%. For real-time recognition, 8 speakers's humming data was used in training and humming pieces of two participants, who were not in the humming database, were tested in real time. Among two participants, the average CRR was 83.23% and the average AR was 80.45%.

4. CONCLUSION AND FUTURE WORK

A new statistical approach to speaker-independent humming recognition was proposed in this work. Phone level hidden Markov models are used to better characterize humming notes. A robust silence (or "Rest") model is created to overcome unexpected note segments caused by back ground noise and some distortions. Features used in note modeling are extracted directly from the data. However, pitch values of hummed data are usually based on the previous note as a reference. Preliminary experimental results showed that our approach is a promising one for further refinement.

There are a few of issues to be investigated as future extension. First, the role of inter-note context should be investigated through context-dependent modeling of melody and tempo contours. The error made by the note decoder can be corrected by a tempo's context model before passing the note to the pitch detector. A pitch's context model can further improve detected results at the pitch level. Second, a more comprehensive database of human humming from a larger set of human subjects should be gathered to enable detailed modeling and evaluation of the recognition performance. This is being performed right now. Third, it is worthwhile to consider a joint note segmentation and pitch detection algorithm. Finally, the music retrieval performance based on the output of the humming recognizer should be investigated, especially from the viewpoint of recognition errors. A long term goal is to optimize the performance of the humming recognizer to maximize the overall retrieval accuracy.

5. ACKNOWLEDGEMENTS

The work is supported in part by National Science Foundation (NSF ITR 53-4533-2720, EEC-9529152) and in part by ALi Microelectronics Corp.

6. REFERENCES

- [1] A. Ghias, J. Logan, D. Chamberlin, and B. C. Smith, "Query by Humming: Musical information retrieval in an audio database," in *Proceedings of ACM Multimedia Conference '95*, San Francisco, California, November 1995.
- [2] Jyh-Shing Roger Jang, Hong-Ru Lee, and Ming-Yang Kao, "Content-based music retrieval using linear scaling and branch-and-bound tree search," in *2001 IEEE International Conference on Multimedia and Expo*, August 2001, pp. 405–408.
- [3] R. J. McNab, L. A. Smith, and Jan H. Witten, "Signal processing for melody transcription," in *the 19th Australasian Computer Science Conference*, 1996.
- [4] R. J. McNab, L. A. Smith, I. H. Witten, C. L. Henderson, and S. J. Cunningham, "Towards the digital music library: Tune retrieval from acoustic input," in *Digital Libraries Conference*, 1996.
- [5] H.-H. Shih, S. S. Narayanan, and C.-C. J. Kuo, "An hmm-based approach to humming transcription," in *2002 IEEE International Conference on Multimedia and Expo (ICME2002)*, August 2002.
- [6] Christopher Raphael, "Automatic segmentation of acoustic musical signals using hidden markov models," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1999, pp. 360–370.
- [7] Christopher Raphael, "Automated rhythm transcription," in *International Symposium on Music Information Retrieval (ISMIR 2001)*, 2001.
- [8] Adriane Swalm Durey and Mark A. Clements, "Melody spotting using hidden markov models," in *International Symposium on Music Information Retrieval (ISMIR 2001)*, 2001, pp. 109–117.
- [9] Adriane Swalm Durey and Mark A. Clements, "Features for spotting using hidden markov models," in *ICASSP 2002*, 2002.