**Combined Audio and Video Analysis for Guitar Chord Identification**

A Thesis

Submitted to the Faculty

of

Drexel University

by

Alex Hrybyk

in partial fulfillment of the

requirements for the degree

of

MS in Electrical Engineering

June, 2010

## Acknowledgements

I'd like to thank my parents for the endless amount of love and devotion they have shown me throughout my life. They have always been and will forever be the source of my passion for music and science. I'd also like to thank my advisor for his constant support of my research and for all the valuable knowledge he has taught me.

## Table of Contents

## List of Figures

# List of Tables

## Abstract
### Combined Audio and Video Analysis for Guitar Chord Identification
#### Alex Hrybyk
#### Advisor: Youngmoo Kim, PhD

This thesis presents a multi-modal approach to automatically identifying guitar chords using audio and video of the performer. Chord identification is typically performed by analyzing the audio, using a chroma based feature to extract pitch class information, then identifying the chord with the appropriate label. Even if this method proves perfectly accurate, stringed instruments add extra ambiguity as a single chord or melody may be played in different positions on the fretboard. Preserving this information is important, because it signifies the original fingering, and implied "easiest" way to perform the selection. This chord identification system combines analysis of audio to determine the general *chord scale* (i.e. A major, G minor), and video of the guitarist to determine *chord voicing* (i.e. open, barred, inversion), to accurately identify the guitar chord.

## 1. Introduction

Music has oftentimes been referred to as a "language of emotion" [9]. In spoken language, feelings are conveyed by combining words to make sentences, sentences to form paragraphs, and so on. Likewise, in music, single notes or tones are combined to form chords, chords to form songs, and so on. Not surprisingly, Western music is heavily based upon chord progressions and the emotions that chords (and combinations thereof) convey. The guitar is a unique instrument in the way it can produce various chords. This research is devoted to designing a system which can automatically identify various types of guitar chords.

### 1.1  Notes and Chords

In Western music, the smallest difference in pitch formally recognized when playing notes is called the *half-step*. The half-step is easily observed by looking at two adjacent white keys on the piano that are not separated by a black key, such as E and F in Figure 1.1 [3]. The frequencies of these adjacent notes are related by a ratio of $2^{\frac{1}{12}}$. Furthermore the frequency of the $n^{\text{th}}$ half-step, relative to some starting frequency $f_0$, is related by

$$f_n = f_0 \cdot 2^{\frac{n}{12}} \quad n = \ldots -2, -1, 0, 1, 2\ldots \tag{1.1}$$

It is easy to see that when $n = 12$, the frequency $f_{12} = 2f_0$. This relationship is called the *octave*, where the note "repeats" but is doubled in pitch. More octaves are also found by looking at integer multiples of 12 notes higher or lower from a starting note. The notes in between the octave ($n = [1, 11]$) define the notes of the chromatic scale, which are used to make melodies, harmonies, chords, etc.

A musical chord is a combination of three or more notes sounding simultaneously [1]. Western music theory has grouped chords together based on the number of half-steps (called *intervals*) between successive notes in the chord. For instance, a *major* chord contains a 4 half-step, followed by a 3 half-step interval, whereas a *minor* chord contains a 3 half-step, then a 4 half-step interval. Some examples of these chords are shown in Figure 1.1.

(a) C Major



(b) C Minor



(c) C Diminished

Figure 1.1: Piano key and staff notation of various C chords.

## 1.2 Guitar Chords

The ability of an instrument to produce chords is a crucial element of that instrument's musical versatility. Trying to identify a chord, by ear or using signal processing, can be a drastically different problem depending on the instrument. The guitar is unique in the way that the strings are tuned. The range of notes which can be played on each string overlap, allowing the same note to be played in multiple positions along the neck of the fretboard (Figure 1.2). Extending this concept to multiple notes, chords can be played in various positions along the neck (Figure 1.3). On an instrument such



Figure 1.2: Guitar neck diagram, showing three places where an F3 note exists.

Figure 1.3: Guitar neck diagram, showing three places where an F major chord exists.

as the piano, this problem does not exist, as there is only one possible piano key to generate each note.

A common representation for written guitar music is known as *tablature*. Tablature uses horizontal lines to represent the strings of the guitar, and numbers placed on those lines to indicate the which fret to play on that string. Figure 1.4 depicts a C major scale in staff notation, followed by three representations in tablature form. All of these tablature notations are valid transcriptions, in that they produce the correct fundamental frequencies as the staff notation when performed. However, only one of these positions may correspond to the original, perhaps easiest fingering. This research seeks to disambiguate which version, or *voicing* of the chord, scale, or note was played.



Figure 1.4: Three voicings of a C major scale in staff and tablature notation, shown in various positions along the guitar fretboard.

## 1.3 Motivations

Detecting, separating, and transcribing the individual notes played by a musician is a widely researched area in audio signal processing. Even if the recorded audio contains only one source and is free of noise, using signal processing techniques to derive note information from the frequency spectrum is still a daunting task. Multiple notes, and natural overtones belonging to each note can crowd the spectrum and make fundamental frequencies difficult to locate. The cluttered state of the spectrum naturally leads us to use other data sources, besides audio, to aid in deriving information about the notes. Digital video of instrumentalists performing music is a widely available data source (now very abundant on the internet), which can contain relevant information pertaining to the notes.

Guitar lessons are more accessible now than ever with the rise of streaming internet video and live interactive lessons. The research presented in this paper has direct applications to these multimedia sources. A system which can automatically identify and transcribe chord diagrams from audio and video lessons between student and teacher would be an invaluable tool to aid in the learning process.

For any guitarist, it is essential to know how a chord or melody maps to particular frets and strings of the guitar. Having to translate music in staff notation (Figure 1.4) to all possible frets/strings in order to determine which is most effective and convenient for the fingers is a tedious task (for advanced to beginning guitarists). Therefore, guitarists will benefit greatly by adding a component to a chord identification system which can detect the unique guitar voicing used.

Automatic chord identification algorithms have traditionally used an audio feature called *chroma*, introduced by Fujishima [4]. The chroma based approach (explained later in section 2.2), though intuitive and easily implemented, presents many problems due to the existence of overtones in the signal. This research avoids this problem by using a polyphonic pitch estimation method named *Specmurt Analysis* which filters out the overtones in the log-frequency spectrum to yield only a chord's fundamental frequencies [11]. The chroma feature also does not resolve the ambiguity of chord voicing - this research fixes this problem using video for chord voicing identification.

Visual approaches to guitar chord and melody transcription have been attempted. Most of these methods, while accurate, are obtrusive to the guitarist; cameras must be mounted to the guitar [2], or the guitarist must wear colored fingertips to be tracked [5]. The method presented here uses brightly colored dots placed at various points along the guitar's fretboard to be tracked by the camera. These dots, which are unobtrusive to the guitarist, are used as reference points to isolate the fretboard within the image, so that principal component analysis may be used to identify the

guitarist's particular voicing of that chord.

The multi-modal guitar chord identification algorithm presented in this research is as follows: first, using Specmurt Analysis, fundamental frequency information will be retrieved and the general *chord scale* identified (i.e. G major, A♯ minor, etc.). Next, using video analysis, the guitarist's particular *chord voicing* (i.e. open, barred, inversion, etc.) will be identified using principal components analysis (PCA) of the guitarist's fretting hand.

## 2. Related Work

Automatic chord identification is a widely researched area of audio and signal processing, as it has a wide array of applications to automating various music related tasks, which could otherwise be tedious. Automatic transcription of chords to written music has direct application to performance and educational purposes. Also, identifying a particular song's chord structure can be used for automatic music segmentation, detecting song similarity, as well as audio thumb-nailing.

### 2.1 Chord Identification vs. Transcription

Before explaining current approaches to chord identification, it is important to note that this research is a subset of a larger area of research called automatic music transcription. Transcription seeks to identify each individual note, onset time, and duration, in order to exactly replicate a piece of music into written notation. Chord identification seeks to analyze the occurrences of groups of notes, and identify a song's higher level structure.

### 2.2 The Chroma Feature

An audio feature called *chroma* has been used as the starting point for most chord recognition systems. The chroma feature reduces the audible frequency spectrum into a 12-dimensional vector, which represents the relative intensity of the twelve tones in the chromatic scale (i.e. C, C♯, D, etc.), irrespective of octave [7]. A flow diagram of the computation of chroma, and chroma over time (the *chromagram*) is shown in Figure 2.1. The chroma vector is computed to see how much of each note is contained in the signal, hopefully creating a more general, chord-like frequency representation of the signal. Fujishima first demonstrated that decomposing the discrete Fourier transform (DFT) of a signal into the 12 pitch classes and then using template matching of various known chords produces an accurate representation of a song's chord structure [4]. Pauws experimented with this method, looking for the maximum correlation between 24 chord templates, and receiving an accuracy of 75.1% - counting only exact matches as correct [8]. However, many problems exist with using the chroma feature to identify chords.

Figure 2.1: Calculation of the chroma feature, and chromagram (chroma over time).

## 2.2.1 Overtones

When using chroma to identify chords, the chroma feature vector for an audio signal is compared against templates for known chords. For example, a D Major chord would have an ideal, normalized chroma vector of $[1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0]$. However, when a single note is played, instruments naturally generate overtones located at integer multiples of the fundamental frequency. These frequencies do not always correspond to octaves of the fundamental frequency, causing the ideal 0's and 1's to fluctuate and can, sometimes, create false chord identifications (Figure 2.2).

Figure 2.2: Chroma of a D Major chord, showing noise from overtones in the signal.

Modified versions of chroma, such as the Enhanced Pitch Class Profile by Lee have been introduced to ease the effects of overtones in the signal [6]. This method computes the chroma vector from the harmonic product spectrum rather than the DFT, suppressing higher harmonics making the chroma vector more like the ideal binary template. However, this method fails to identify the unique guitar voicing of the chord.

### 2.2.2 Voicing Independency

As explained before in Section 1.1, in guitar chord identification, it is important to preserve the position or voicing information associated with a chord along with its scale. If two chords in different voicings contain the same notes, the chroma feature by design will only show the similarity in notes, leaving the ambiguity of voicing unresolved. This can be seen in Figure 2.3.

### 2.3 Visual Methods for Chord Identification

Burns et al. developed a visual system for left-hand finger position tracking with respect to a string/fret grid [2]. Their method relies on the circular shape of fingertips, using a circular Hough transform on an image of the left-hand to detect fingertip locations with respect to the underlying fretboard. However, this method requires mounting a camera on the headstock of the guitar, which poses many problems: it can be obtrusive to the guitar player's natural method of playing, and also only captures information about the first five frets of the guitar.

Kerdvibulvech et al. proposed to track the fingering positions of a guitarist relative to the guitar's

Figure 2.3: Similar chromagrams of two different voicings of a C Major chord.

position in 3D space [5] . This is done by using two cameras to form a 3D model of the fretboard. Finger position was tracked using color detection of bright caps placed on each of the guitarist's fingertips. Again, this can hinder the physical capabilities and creative expression of the guitarist, which should not happen in the transcription process.

An alternative method for inferring guitarist fingering position and/or tablature from a musical staff notation was proposed by Radicioni et al. [10]. Their method first builds a grid of all finger/fret/string combinations, and then uses a least-cost path algorithm to compute the most likely position and fingering to use. Their method, however, does not account for chords, which are a critical element of guitar playing.

## 3. Audio Analysis

When playing a single note, instruments produce natural harmonics (overtones) in addition to the note's fundamental frequency. Therefore, when playing multiple notes, the frequency spectrum of the audio appears cluttered, making detection of the fundamental frequencies (the actual notes) hard to locate. Using *Specmurt* analysis, the notes of a guitar chord can be extracted from the audio signal [11].

### 3.1  Specmurt Analysis

Multiple fundamental frequency estimation using Specmurt analysis is performed by inverse filtering the log-scale frequency domain with a common harmonic structure of that instrument [11]. The resulting log-frequency spectrum contains only impulses located at the log-fundamental frequencies.

### Log-Scaled Frequency Domain

Harmonics of fundamental frequencies normally occur at integer multiples of that fundamental. Furthermore, if the fundamental frequency changes by some $\Delta f$, the change in frequency of its respective higher harmonics will also be integer multiples of that $\Delta f$.

$$
\begin{aligned}
f_n &= n f_0 \qquad n = 1, 2, 3... \\
(f + \Delta f)_n &= n(f_0 + \Delta f) \\
&= n\Delta f + n f_0 \\
&= n\Delta f + f_n
\end{aligned}
$$

By resampling the frequency domain to have a log-scaled axis, this allows changes in harmonics to simply be a sum of the log fundemental and $\Delta f$.

$$\begin{aligned}
\hat{f}_n &= \log f_n \\
&= \log(nf_0) \qquad n = 1, 2, 3... \\
&= \log(n) + \log(f_0)
\end{aligned}$$

This allows the harmonics of a frequency to be consistantly spaced by $\log(n) + \log(f_0)$, independent of fundamental frequency.

$$\hat{f} = \log f \tag{3.1}$$

### 3.1.1   Common Harmonic Structure

Using the log-scale frequency axis, we can assume that the harmonic frequencies are located at $\hat{f} + \log 2, \hat{f} + \log 3, ..., \hat{f} + \log n$. When a chord is played on an instrument, each note will presumably contain these same harmonic frequencies, beginning at different $\hat{f}$'s. Therefore, we can assume that the log-scaled multipitch spectrum, $c(\hat{f})$, is a combination of these harmonic structures, shifted and weighted differently per note. Specifically, the resulting log-scale frequency spectrum, $c(\hat{f})$, is equal to the convolution of a common harmonic structure, $h(\hat{f})$, with a fundamental frequency distribution, $g(\hat{f})$.

$$c(\hat{f}) = h(\hat{f}) * g(\hat{f}) \tag{3.2}$$

The harmonic structure can be written in terms of its log-frequency axis spacing, $\hat{f}_{0n}$, and its harmonic weights, $W_n$, where $n = 1, 2...N$ harmonics.

$$h(\hat{f}, W) = \sum_{n=1}^{N} W_n \delta(\hat{f} - \hat{f}_{0n}) \tag{3.3}$$

The harmonic weights will initially be a guess, which will be refined later using an iterative process to minimize the overall error of Specmurt analysis. An example initial harmonic structure, $h_i(\hat{f})$ could be of the form

$$h_i(\hat{f}) = \sum_{n=1}^{N} \frac{1}{n} \delta(\hat{f} - \hat{f}_{0n}) \tag{3.4}$$

Figure 3.1: Log-spaced frequency domain $c(\hat{f})$ as a convolution of common harmonic structure $h(\hat{f})$ with fundamental frequency distribution $g(\hat{f})$.

### 3.1.2 Specmurt Domain

In order to determine the desired fundamental frequency distribution, $g(\hat{f})$, one can solve (3.2) by deconvolving the log-spectrum with the common harmonic structure. An easier way of obtaining $g(\hat{f})$ would utilize the duality of the time/frequency-convolution/multiplication relationship (shown in Figure 3.1). Therefore, taking the inverse Fourier transform would yield the relationship

$$
\mathcal{F}^{-1}\{c(\hat{f})\} = \mathcal{F}^{-1}\{h(\hat{f}) * g(\hat{f})\} \tag{3.5}
$$

$$
C(\hat{s}) = H(\hat{s})G(\hat{s}) \tag{3.6}
$$

where $\hat{s}$ is a temporary Specmurt domain variable. Simple algebra followed by a Fourier tranform of $G(\hat{s})$ will yield the resulting fundamental frequency spectrum.

$$G(\hat{s}) = \frac{C(\hat{s})}{H(\hat{s})} \tag{3.7}$$

$$\mathcal{F}\{G(\hat{s})\} = g(\hat{f}) \tag{3.8}$$

### 3.1.3 Thresholding

When it is known that the harmonics of the signals should be of lesser magnitude than the fundamental frequencies, the minor peaks in $g(\hat{f})$ can be attenuated using a non-linear thresholding. This is used to compensate for remaining error, after minimization in Section 3.1.4. The thresholded $g(\hat{f})$ function, $\bar{g}(\hat{f})$, is given by

$$\bar{g}(\hat{f}) = \frac{1}{1 + \exp\left\{-\alpha\left(\frac{g(\hat{f})}{g_{max}} - \beta\right)\right\}} g(\hat{f}) \tag{3.9}$$

The thresholding function is controlled by two parameters, $\beta$ and $\alpha$. $\beta$ corresponds to the "cut-off" value under which frequency components are assumed to be unwanted, and $\alpha$ represents the degree of strictness of the thresholding function around $\beta$.

### 3.1.4 Error Minimization

The squared error after performing Specmurt analysis can be defined as

$$E(W_n) = \int_{-\infty}^{+\infty} \left\{c(\hat{f}) - h(\hat{f}, W_n) * g(\hat{f})\right\}^2 d\hat{f} \tag{3.10}$$

Minimizing the error of Specmurt is done by refining the harmonic weights, $W_n$, of the harmonic structure. This is done by setting the error's $N$ partial derivatives to zero, and solving the system of equations for $W_n$. For $N$ harmonics, this yields a system of equations which is equivalent to the following matrix equation with coefficients $a_{i,j}$ and $b_j$.

$$\frac{\partial E}{\partial W_n} = 0 \qquad n = 1...N \tag{3.11}$$

$$\begin{pmatrix} a_{1,1} & \cdots & a_{1,n} & \cdots & a_{1,N} \\ \vdots & & \vdots & & \vdots \\ a_{n,1} & \cdots & a_{n,n} & \cdots & a_{n,N} \\ \vdots & & \vdots & & \vdots \\ a_{N,1} & \cdots & a_{N,n} & \cdots & a_{N,N} \end{pmatrix} \begin{pmatrix} W_1 \\ \vdots \\ W_n \\ \vdots \\ W_N \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ b_n \\ \vdots \\ b_N \end{pmatrix} \tag{3.12}$$

where

$$a_{i,j} = \int_{-\infty}^{+\infty} \bar{g}(\hat{f} - W_i)\bar{g}(\hat{f} - W_j) \, d\hat{f} \tag{3.13}$$

$$b_j = \int_{-\infty}^{+\infty} c(\hat{f})\bar{g}(\hat{f} - W_j) \, d\hat{f} \tag{3.14}$$

The original Specmurt formulation assumed that the first weight, $W_1 = 1$, of the normalized common harmonic structure. After experimentation with various guitar signals, the higher harmonics were sometimes of larger magnitude than the fundamental frequency. By allowing the first harmonic's magnitude to vary, the algorithm was able to better identify fundamental frequencies.

## 4. Video Analysis

In order to visually identify the performing guitarist's chord voicing, the guitar fretboard must first be located and isolated within the image. However, the guitar can be held in many different orientations relative to the camera, making it difficult to find the location or coordinates of the fretboard in the image plane.

The frets of a guitar are logarithmically spaced to produce the 12 tones of the western scale. The coordinates in the $(x, y)$ plane are plotted in Figure 4.1, where the $x_i$ coordinates are related by

$$x_i = \sum_{k=0}^{i} x_0 \times 2^{\frac{k}{12}} \tag{4.1}$$

### 4.1 Homography

Homography is the process of applying a projective linear transformation to a scene (a 2D image or 3D space), to describe how perceived positions of observed objects change when the point of view of the observer (a camera) changes. Homography will be used to determine the correct perspective transformation, i.e. rectify or warp the original image to fit the ideal fretboard spacing in Figure 4.1. This will make it easy to isolate the fretboard in the image for analysis. The general homography matrix equation

$$w\mathbf{p}' = \mathbf{H}\mathbf{p} \tag{4.2}$$

states that points in the image, $\mathbf{p}'$ can be expressed as a warping of ideal points $\mathbf{p}$ with a homography matrix $\mathbf{H}$. The homography matrix is a transformation matrix between the two images, based on which a one-to-one relationship between the features points $\mathbf{p}'$ and $\mathbf{p}$ [13]. Specifically, the points will have two dimensions, $x$ and $y$, and will be expressed in terms of a 3x3 homography matrix with elements $h_{ij}$.

$$\begin{bmatrix} x_i' \\ y_i' \\ 1 \end{bmatrix} \approx \begin{bmatrix} h_{00} & h_{10} & h_{20} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} \tag{4.3}$$

$x_i$ are determined from (4.1) and $y_i$ are determined as an arbitrary guitar neck width (from the ideal, rectangular fretboard). The corresponding reference points $(x_i', y_i')$ in the image now need to
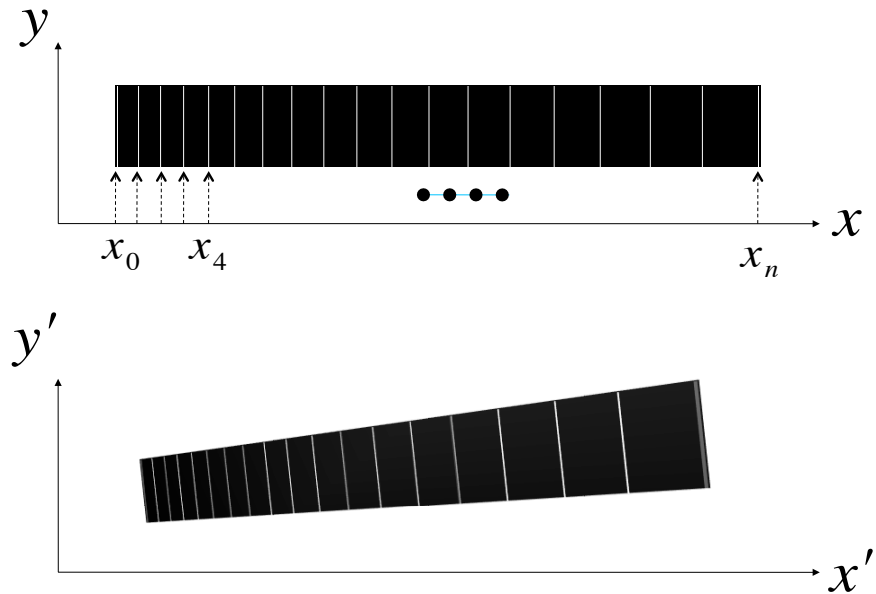
Figure 4.1: Ideal fretboard (top) with logarithmic $x$ spacing of $n$ frets, and arbitrary neck width in $y$ direction, and seen image (bottom) with warped spacing.

be established, to compute the homography matrix, **H**.

## 4.2 Reference Point Tracking

In order to perform the homography rectification concepts in 4.1, the correct reference points in the image must be determined. Attempts were made at using an iterative non-linear error minimization method, which proved initially unsuccessful (see later section 6). Instead, distinct bright colored stickers were placed at various fret locations on the neck of the guitar. The coordinates of these points $(x_i', y_i')$ were tracked in each frame of video using a simple color masking followed by a k-means clustering algorithm. The small stickers were placed on the neck of the guitar on either side of the metal frets, so as not to interfere with the guitarist's playing and the timbre of the instrument.

A set of $(x_i, y_i)$ and $(x_i', y_i')$ now exist, corresponding to the frets of the guitar. The homography matrix is determined by minimizing the mean square error of (4.3) using these points. Equation 4.3 can be rewritten as
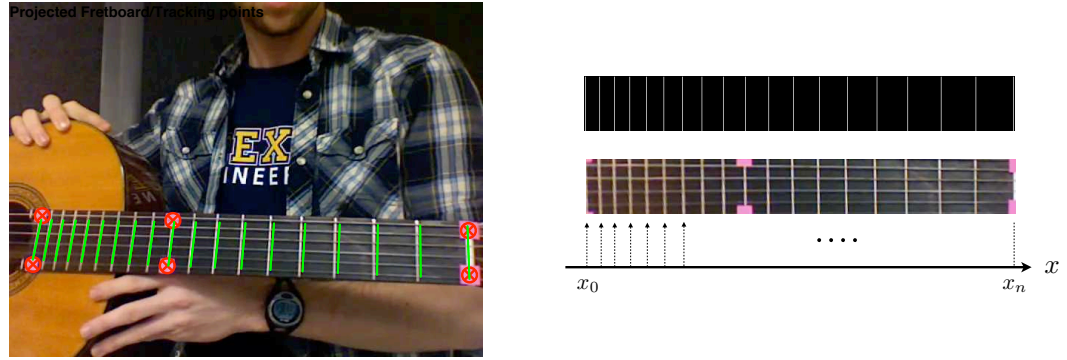
Figure 4.2: (left) Original image showing tracking points (in red), projected frets (in green) using the homography matrix. (right) Ideal fretboard, and subsection of original image after applying homography matrix to each coordinate.

$$
\begin{bmatrix}
x_1 & y_1 & 1 & 0 & 0 & 0 & -x_1'x_1 & -x_1'y_1 & -x_1 \\
0 & 0 & 0 & x_1 & y_1 & 1 & -y_1'x_1 & -y_1'y_1 & -y_1 \\
& & & & \vdots & & & & \\
x_n & y_n & 1 & 0 & 0 & 0 & -x_n'x_n & -x_n'y_n & -x_n \\
0 & 0 & 0 & x_n & y_n & 1 & -y_n'x_n & -y_n'y_n & -y_n
\end{bmatrix}
\begin{bmatrix}
h_{00} \\ h_{10} \\ h_{20} \\ h_{10} \\ h_{11} \\ h_{12} \\ h_{20} \\ h_{21} \\ h_{22}
\end{bmatrix}
=
\begin{bmatrix}
0 \\ 0 \\ \vdots \\ 0
\end{bmatrix}
\tag{4.4}
$$

$$
\begin{matrix}
\mathbf{A} & \mathbf{h} \\
2n \times 9 & 9 \times 1
\end{matrix}
=
\begin{matrix}
\mathbf{0} \\
2n \times 1
\end{matrix}
\tag{4.5}
$$

An estimate for $\mathbf{h}$ can be found by taking the eigenvector of $\mathbf{A^T A}$ with the smallest eigenvalue.

Applying the inverse transformation, $\mathbf{H^{-1}}$, to the ideal grid in Figure 4.1 yields frets that overlay perfectly with the frets in the image (Figure 4.2). Applying $\mathbf{H}$ to the original image and taking the subsection of coordinates yields the rectified fretboard (Figure 4.2), whose fret spacings are known from Equation (4.1). The rectified fretboard is now isolated and in a usable form for PCA.

## 4.3 Determination of Chord Voicing

The next goal is to determine which chord voicing, given the subset of voicings that exist for a particular chord. PCA is used to decompose the rectified fretboard in its "eigen-chord" components, and determine the correct chord voicing.

Let the training set of fretboard images be $F_1, F_2...F_M$ which are vectors of length $LW$ for an image with dimensions $L$ by $W$. An example training set of fretboard images is shown in Figure 4.3. The average image is $A = \frac{1}{M} \sum_{i=1}^{M} F_i$, and each image with subtracted mean is $\bar{F}_i = F_i - A$. PCA seeks to find the eigenvectors and eigenvalues of the covariance matrix

$$\mathbf{C} \quad = \quad \frac{1}{M} \sum_{i=1}^{M} \bar{F}_i \bar{F}_i^T \tag{4.6}$$

$$= \quad \mathbf{SS}^T \tag{4.7}$$

where $\mathbf{S} = [\bar{F}_1, \bar{F}_2...\bar{F}_M]$ is a set of training of images in matrix form. However, $\mathbf{C}$ is of dimension $LW$; the images used in this experiment are of size 80x640 pixels, and computing 51200 eigenvectors and eigenvalues is computationally intractable. Turk et. al presented a method for solving for the $LW$ eigenvectors by first solving for the eigenvectors of an $MxM$ matrix $\mathbf{S}^T\mathbf{S}$ [12]. The $M$ eigenvectors $v_l$ are used to form the eigenvectors $u_l$ of $\mathbf{C}$.

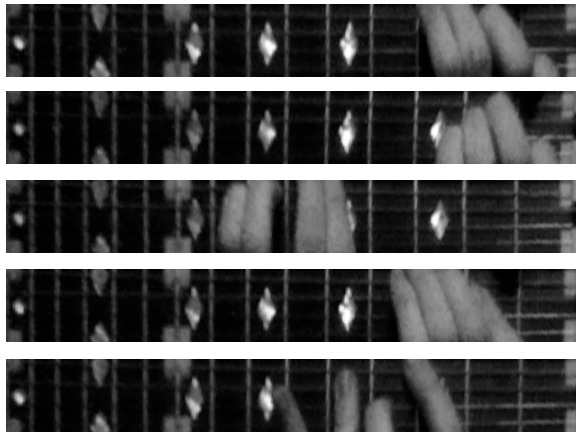$$u_l = \sum_{i=1}^{M} v_l \bar{F}_i \quad l = 1...M \tag{4.8}$$



Figure 4.3: Example fretboard images used for training.

A new image $\tilde{F}$ can be reduced to its eigen-chord components, $c_k$, using the $M'$ eigenvectors $(M' < M)$ which correspond to the larger eigenvalues of $\mathbf{S}^T\mathbf{S}$.

$$c_k = u_k(\tilde{F} - A) \quad k = 1...M' \tag{4.9}$$

## 5. Experimental Results

Three guitarists were asked to perform a sequence of chords from chord diagrams. The chords were a selection drawn from eight scales (major and minor), each in three voicing-dependent positions: open (traditional open stringed), barred, and a 1st inversion, totaling 24 chords all together. The system was evaluated using various combinations of features derived from audio only, video only, and combinations thereof. All experiments were performed using leave-one-out training of audio and video when using PCA.

### 5.1  Audio Only

First, the chroma feature was used to identify chord scale and voicing. The chord scale was identified using chroma template matching of the chords to be detected. The with the maximum correlation was selected to be the correct chord. This method showed an accuracy of 15.28%; this is better than random chance ($\frac{1}{24} \rightarrow 4.2\%$), but clearly can be improved upon.

The output of Specmurt analysis is a piano-roll vector of size 48, each element corresponding to the energy of a chromatic note from C2 to B5 (4 octaves, 12 notes per octave). An example of a piano-roll vector over multiple time frames is shown in Figure 5.1.

Two methods were used to calculate the correctness of the chord scale and voicing using this vector. It is known what notes make up each major and minor scale. Therefore, the chord scale was evaluated by summing the energy over all octaves of the notes belonging to that chord - similar to chroma analysis. The chord scale with the highest energy was assumed to be correct, yielding an accuracy of 98.6%.

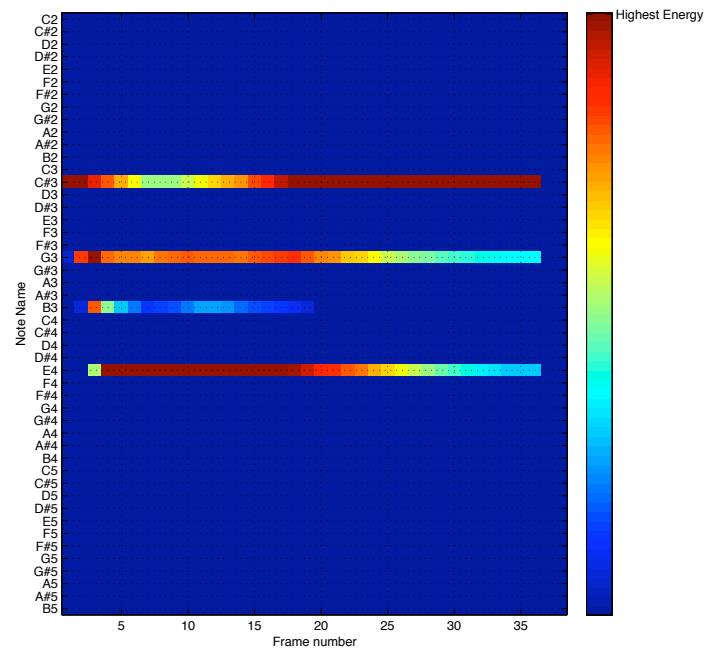It is not deterministic, however, as to which



Figure 5.1: Specmurt piano-roll output of a C♯m7♭5 jazz chord.

chord voicing created a particular set of notes,

or chord. For example, both the G major open

chord and G major barred chord contain six

notes total, all of the same fundamental frequencies, but the notes are rearranged on different

strings, and hence use different fingerings. Therefore, a training set using the piano-roll energy

vector was developed for each chord scale. Using PCA to identify chord voicing from the piano-roll

vector shows some accuracy (62%) but is understandably low, as the difference in note energies may

be very fine and inseparable for different voicings with similar notes.

## 5.2 Video Only

A training set of 240 images was used to build the eigen-chord space for each test. Frames of

video were then projected into the chord-space using three eigen-chords of the training set using

(4.9), and its closest centroid was assumed to be the correct chord.

Chord scale identification using only video performed extremely poorly (34%). This is expected,

as the chord scale centroids in the projected chord-space after PCA are somewhat meaningless. For
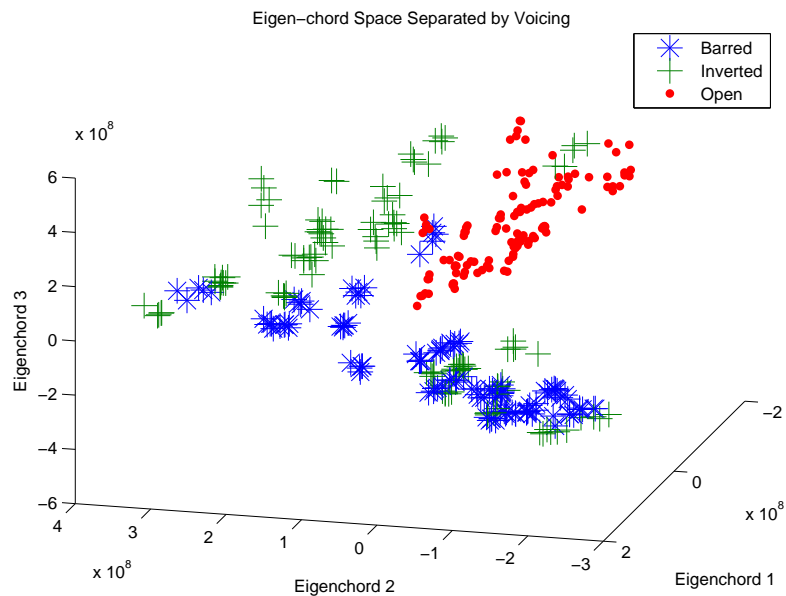


Figure 5.2: Three voicings from A minor, G major, and C major, after being projected into the chord-space. Various colors and symbols show the how the voicing of chords remain grouped after dimensionality reduction.

a particular chord scale, many different voicings exist at various points on the fretboard, which is what we hope to separate by using PCA.

For chord voicing however, very high accuracy was achieved (94.4%). Figure 5.2 shows how various voicings of chords, irrespective of scale, tend to group together due to the similar hand shapes used by the guitarist.

## 5.3 Combined System

The system which performs the best in terms of correctly identifying the overall chord (scale and voicing) utilizes the strong points of scale and voicing identification within the audio and video results. Since Specmurt analysis yielded extremely high accuracy for determining scale, it was used as a preprocessing step to voicing identification via video.

|         | Audio only | Video only | Combined System |
|---------|:----------:|:----------:|:---------------:|
| Scale   | **98.6**   | 34.8       | **98.6**        |
| Voicing | 62.0       | **94.4**   | **94.4**        |
| Both    | 61.1       | 32.8       | **93.1**        |

Table 5.1: Accuracy results (%) for various combinations of modes of information. Combined accuracy results using Specmurt for scale identification, and video for voicing identification showing highest accuracy.

## 6. Conclusion

This paper has presented an alternate approach to automatic guitar chord identification using both audio and video of the performer. The accuracy of chord identification increases from 61.1% to 93.1% when using audio for scale identification, and video for voicing. The "eigen-chord" decomposition of fretboard images proved extremely successful in distinguishing between a given chords voicings (normal, barred, inverted) if the chord scale is known (94.4%). The video and audio components of this guitar chord identification system also have the potential to be expanded upon.

### 6.1   Automatic Fretboard Registration

Placing colored tracking points along the neck of the guitar presents additional constraints on how the guitar fretboard can be rectified: all the tracking points must be visible in the frame of video, and nothing else in the frame may have similar color. Ideally, we would like to locate the fretboard without these points. By looking at the edge-detected image of a guitar, this produces a fairly accurate representation of where the frets are - the color of the metal frets contrasts heavily with that of the wooden neck, providing edges at frets (Figure 6.1).

Using the homography concept in 4.1, the points denoted as edges, $\mathbf{p'}$, should be warped using $\mathbf{H^{-1}}$ to align with the ideal fret-grid points $\mathbf{p}$. This is equivalent to minimizing an error function defined as

$$E(\mathbf{H}) = ||\mathbf{p} - \mathbf{H^{-1}}\mathbf{p'}||^2 \tag{6.1}$$

$$\mathbf{H} = \underset{\mathbf{H}}{\operatorname{argmin}}(E(\mathbf{H})) \tag{6.2}$$

After experimentation, the error function $E(\mathbf{H})$ is noticeably non-convex, and contains local minima in $\mathbf{H}$. The two fret-grids "align" in alternate orientations which are incorrect, but still minimize the error function. This area of research is being continued with the motive of constraining (6.1) and (6.2), such that the error function will always be convex, and converge to a global minimum when the two images are correctly aligned.
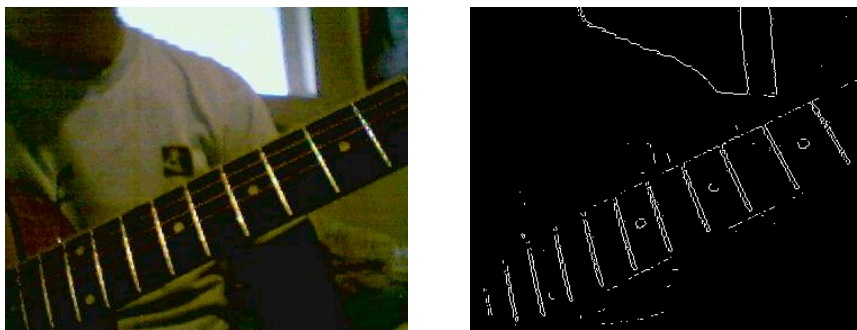
Figure 6.1: Guitar image (left) and edge-thresholded image (right).

## 6.2 Larger Training Sets

Very high accuracy of video voicing identification (94.4%) was achieved using image data from only three guitarists. A more robust classifier of chord voicings could be created by collecting more data, to account for players who use non-traditional finger orientations for chords. With more data, the accuracy of determining chord scale from video may increase (34.8%), as scales may then form more meaningful distributions in the eigen-chord space.

## 6.3 Additional Chord Types

This system is very extendable to detect different chord scales besides major and minor. Detection of diminished, augmented, 7th, and other jazz chords are easily implemented with the chroma-style analysis of Specmurt's output, and refined search using the eigen-chord decomposition of the fretboard image.

## 6.4 Fusing Audio/Video Data

Currently the system uses Specmurt analysis to determine a chord's scale as a pre-processing step to eigen-chord decomposition of the fretboard to determine voicing. This means that any error introduced by Specmurt propagates throughout the rest of the system. Therefore it is desired to jointly estimate the scale and voicing together using audio and video features simultaneously.

# Bibliography

[1] Bruce Benward and Marilyn Saker. *Music in Theory and Practice*, volume 1. McGraw-Hill Companies, Inc., 8th edition, 2009.

[2] Anne-Marie Burns and Marcelo M. Wanderley. Visual methods for the retrieval of guitarist fingering. In *NIME '06: Proceedings of the 2006 conference on New interfaces for musical expression*, pages 196–199, Paris, France, France, 2006. IRCAM — Centre Pompidou.

[3] William Christ, Richard DeLone, Vernon Kliewer, Lewis Rowell, and William Thomson. *Materials and Structure of Music*, volume 1. Prentice-Hall, Englewood Cliffs, N.J., 3rd edition, 1972.

[4] T. Fujishima. Realtime chord recognition of musical sound: A system using common lisp music. In *Proceedings of the International Computer Music Conference*, 1999.

[5] C. Kerdvibulvech and H. Saito. Vision-based guitarist fingering tracking using a bayesian classifier and particle filters. In *PSIVT07*, pages 625–638, 2007.

[6] Kyogu Lee. Automatic chord recognition from audio using enhanced pitch class profile. In *Proceedings of the International Computer Music Conference*, 2006.

[7] Kyogu Lee. *A System for Acoustic Chord Transcription and Key Extraction from Audio Using Hidden Markov Models Trained on Synthesized Audio*. PhD thesis, Stanford University, 2008.

[8] Steffen Pauws. Musical key extraction from audio. In *Proceedings of the International Conference on Music Information Retrieval*, 2004.

[9] Carroll C. Pratt. *Music as the language of emotion*. The Library of Congress, December 1950.

[10] D. P. Radicioni, L. Anselma, and V. Lombardo. A segmentation-based prototype to compute string instruments fingering. In R. Parncutt, A. Kessler, and F. Zimmer, editors, *Procs. of the 1st Conference on Interdisciplinary Musicology (CIM04)*, Graz, Austria, 2004.

[11] S. Saito, H. Kameoka, K. Takahashi, T. Nishimoto, and S. Sagayama. Specmurt analysis of polyphonic music signals. *Audio, Speech, and Language Processing, IEEE Transactions on*, 16(3):639–650, February 2008.

[12] M.A. Turk and A.P. Pentland. Face recognition using eigenfaces. In *Computer Vision and Pattern Recognition, 1991. Proceedings CVPR '91., IEEE Computer Society Conference on*, pages 586 –591, June 1991.

[13] Xiaohua Wang and Bing Yang. Automatic image registration based on natural characteristic points and global homography. In *Computer Science and Software Engineering, 2008 International Conference on*, volume 5, pages 1365 –1370, dec. 2008.