

Real-time monophonic guitar pitch estimation

The feasibility of Fourier transform based methods

Luc de Jonckheere

October 22, 2021

Abstract

Short summary.

Contents

1	Introduction	1
2	Related work	2
3	Preliminaries	2
3.1	Audio processing	2
3.2	Fourier transform	2
3.3	Real-time	4
3.4	Music theory and notation	5
3.5	Envelope and transients	6
3.6	Fundamental frequency, overtones and timbre	6
4	Thesis content	6
4.1	Real-time constraint	7
4.2	Software amplification	7
5	Implementation	7
6	Experiments	7
7	Conclusions	7
8	Future work	7
A	Measuring the latency of the AXON AX 100 mkII	7
B	Effect of different window functions	7
	Bibliography	7

1 Introduction

Pitch estimation, which is also referred to as f_0 estimation, is an important subtask within the field of Automatic Music Transcription (AMT). The goal of

pitch estimation is to estimate the pitch or fundamental frequency f_0 of a given signal. In the context of AMT, pitch estimation is used to determine what note is played in a given signal.

Real-time pitch estimation is a subproblem where we want to estimate the note associated with the measured pitch while the musician is playing it with minimal latency. This entails we have to use the latest received signal. In contrast to non-real-time methods, we have no knowledge of what may happen ahead of time and signal corresponding to previous notes is irrelevant. This limits the methods we can use to solve this problem.

If pitch estimation can accurately be performed in real-time, it can be used to create a digital (MIDI) instrument from an acoustic instrument. This digital instrument can then be used as an input for audio synthesizers, allowing musicians to produce sounds from a wide variety of instruments. Furthermore, accurate real-time pitch estimation can be used to automatically correct detuned instruments by pitch shifting the original signal to the closest harmonious note.

The Fourier transform is often used to decompose a signal into the frequencies that make up the signal. Predominant frequencies in the signal show up as spectral peaks in the frequency domain. These peaks are important to human perception of melody [9]. Other popular methods used for pitch estimation include non-negative matrix factorization, autocorrelation, statistical model based estimation and hidden Markov model based estimation.

Our research focusses on monophonic pitch estimation. Here, we assume that the signal contains at most one note. It is much easier to perform monophonic pitch estimation compared to polyphonic pitch estimation [10], especially when using Fourier transform based methods, as fundamental limits of the Fourier transform inhibit our ability to discern two low pitched notes [4]. Furthermore, hexaphonic guitar pickups are becoming more widespread, which allows us to view the guitar as six monophonic in-

struments instead of one six-way polyphonic instrument. Commercial guitar synthesizer solutions from companies with big research departments such as Roland also use these hexaphonic pick-ups, which indicates the infeasibility of accurate and responsive polyphonic pitch estimation of a guitar.

This thesis builds upon a preliminary research project [7]. In our research project, we found that Fourier transform based pitch estimation methods might not be well suited for real-time use due to fundamental limitations of the Fourier transform. In this work, we will further research if Fourier transform based methods are viable, as real-time transcription research often relies Fourier transform based methods.

The goal of this thesis is to research the limits of Fourier transform based real-time pitch estimation. To correctly assess the limits, we develop a pitch estimation framework. This framework will focus on extensibility and the ability to perform automated tests. This is important, as much work in this field does not provide its associated source code. This limits the ability to build on other’s work and hinders direct comparisons between different methods. Our framework can provide a common ground for the different methods and algorithms to be implemented and compared in. The framework is available at www.github.com/lucmans/digistring.

2 Related work

Much research has been performed on Fourier transform based real-time pitch estimation. Fourier based methods are deemed infeasible by some due to low frequency resolution [5]. This is especially problematic when adhering to a real-time constraint, as extra short signal frames have to be used. Some papers circumvent this problem by choosing a very high real-time constraint [12, 11], however, this inhibits the use for real-world applications. Furthermore, conventional operating systems also have a latency when delivering audio samples to your program due to how audio drivers work [20].

A big problem with Fourier based pitch estimation is the occurrence of overtones [21]. Especially octaves are a problem, as the fundamental and all overtones of the higher note overlap with overtones of the lower note. This is referred to as the octave problem [23]. Overtones are periodic in nature, as they diminishingly repeat every multiple of the fundamental frequency. As a consequence, they could also be detected using a subsequent Fourier transform [17] on the frequency domain. However, this

does not solve the octave problem.

Many different transform have been researched for pitch estimation, however, Fourier transform remains popular as it is broadly studied and its behaviour is well known [19]. Lately, the CQT transform is gaining popularity as it may provide higher resolution in the frequency domain [26] at the cost of lower computational efficiency [24]. However, the limitations of the Fourier transform are mainly due to the low frequency resolution and since the CQT is often implemented using Fourier transforms [6], we are still left with the same problem. One big advantage is that the frequency bins can perfectly align with the notes of an instrument [8]. However, as described in Section 3.6, overtones are dissonant with respect to our notes and consequently, the CQT bins do not align with the overtones. If a note perfectly aligns with a Fourier bin, all overtones will also align. In order to cover every note, we could instead perform 12 Fourier transform in parallel. This difference is important when performing polyphonic transcription.

3 Preliminaries

Jargon required to understand this paper.

3.1 Audio processing

Audio is represented by samples. We get these samples from the OS audio driver. Driver latency.

3.2 Fourier transform

The Fourier transform is a mathematical transform which transforms a function of time to a complex valued function of frequency and phase. Here, the magnitude represents the amplitude and the argument represents the phase of the corresponding sine wave. The Fourier transform works on continuous functions and assumes an infinite time interval. Concepts such as continuous and infinite cannot be represented by a computer. Consequently, the discrete Fourier transform (DFT) has to be used for Fourier analyses on computers. It can be efficiently implemented using the fast Fourier transform (FFT) algorithm.

The DFT transforms a finite sequence of equally spaced samples, which we will refer to as a frame, into an equal number of complex values representing the amplitude and phase, which we refer to as bins. When working with audio, the samples are real valued, and the DFT output is symmetrical. Because of this, we can discard the second half of the output. In the rest of this thesis, we will only consider the first half of the output. Each bin corresponds to a

specific frequency. All other frequencies are spread out over multiple bins due to spectral leakage, which will be discussed later. Given a frame F , the number of samples in the frame is $n_F = |F|$. Using n_F and sample rate f_{SR} , we can calculate the distance between bins:

$$\Delta f_{bin} = \frac{f_{SR}}{n_F}$$

This is also referred to as the frequency resolution. Closely related to the frequency resolution is the frame length, which is calculated as follows:

$$t_F = \frac{n_F}{f_{SR}}$$

Given a bin number $i \in [0, \lfloor \frac{n_F}{2} \rfloor]$, the frequency of a bin can be calculated as:

$$f_{bin} = \Delta f_{bin} * i$$

The 0 Hz bin corresponds to the so called DC offset. This is the average amplitude of the signal. The frequency of the last bin is also called the Nyquist frequency, which is the highest frequency which can be sampled without aliasing. The Nyquist frequency is equal to half the sample rate.

The DFT assumes the frame to be periodic. In other words, the frame is regarded as infinitely repeating. This may lead to aliasing if the beginning and end of a frame do not align, see Figure 1. Here, we take a frame shown by the red lines. The frame is not aligned to a period of the sine wave and causes aliasing as seen in the second graph. This kind of aliasing is called spectral leakage and causes other frequencies which are not a multiple of Δf_{bin} to spread out over multiple bins. By applying a window function, the beginning and end of a frame are forced to align.

Spectral leakage can be controlled using window functions. Given signal $s(n)$ and window function $w(n)$, we get the resulting signal $res(n)$ using:

$$res(n) = s(n) * w(n)$$

Figure 2 shows the working of a window function on a frame graphically.

The different window functions trade off between a narrow center lobe and little overall leakage [13], see Figure 3 for some examples. Using the rectangular window could be considered as using no window function. It simply only takes n_F samples from the wave without any alteration. The Hann window is an all round performing window and as a result is often used. The Welch window is a window with a very narrow center lobe. The Dolph–Chebyshev window has little and very evenly spread overall leakage.

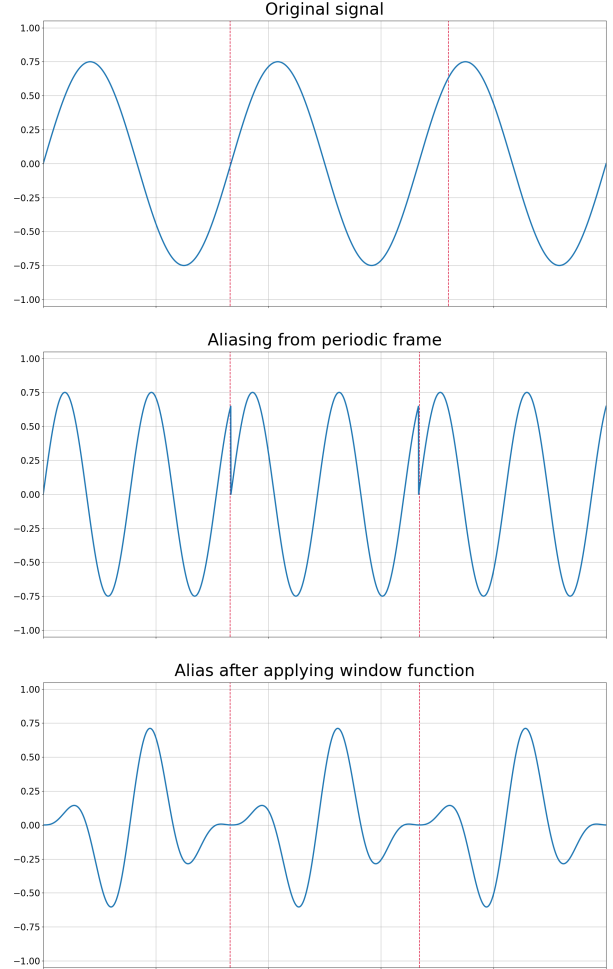


Figure 1: Aliasing that occurs when performing the discrete Fourier transform

Zero-padding can be used to increase the number of output bins of the DFT. It is important to note that it does not increase the resolution of the DFT, it merely interpolates the coarse spectrum to become more smooth [2]. Two frequencies closer than Δf_{bin} together will still fall in the same bin [1]. This form of interpolation is relatively compute intensive [27].

A less compute intensive method of interpolation is quadratic interpolation. It interpolates the actual location of a spectral peak between bins by fitting a quadratic function through the bins on a decibel scale [16, 14]. We calculate a value $p \in [-\frac{1}{2}, \frac{1}{2}]$, which is the offset in bins of the interpolated peak with respect to the peak bin. Using the magnitude of the peak $y(0)$ and the magnitude of the neighbouring

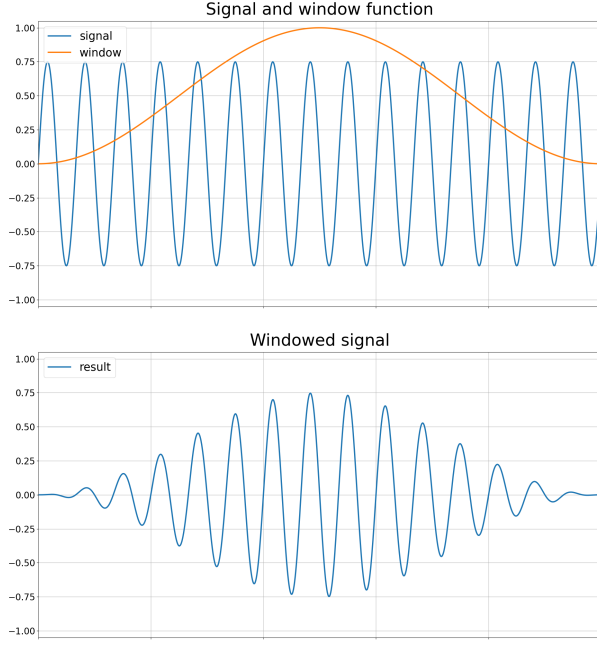


Figure 2: An example of using a window function on a signal

Figure 3: Examples of the spectral leakage of some window functions

bins $y(-1)$ and $y(1)$, we define:

$$\begin{aligned}\alpha &= 20 * 10 \log y(-1) \\ \beta &= 20 * 10 \log y(0) \\ \gamma &= 20 * 10 \log y(1)\end{aligned}$$

Then, we can calculate p as follows:

$$p = \frac{1}{2} \cdot \frac{\alpha - \gamma}{\alpha - 2\beta + \gamma}$$

The amplitude a_i corresponding to the interpolated peak is:

$$a_i = \beta - \frac{(\alpha - \gamma) * p}{4}$$

Given the bin number b of the spectral peak location, the frequency f_i corresponding to the interpolated peak is:

$$f_i = \Delta f_{bin} * (b + p)$$

Quadratic interpolation can be combined with zero-padding for better results [27].

3.3 Real-time

Real-time is a difficult concept, as it has many definitions. Here are a few examples of relevant definitions:

1. Being synced with actual clock time (or wall time). This is for instance relevant when playing media such as audio and video. When such media is played at an incorrect speed, it could be considered distorted. The hardware which keeps track of the clock time is called a real-time clock.
2. A system must response within a specified time constraint, which is called the real-time constraint or deadline. This constraint is usually a relatively short time. This definition comes from real-time computing and is relevant when making car airbags or airplane control systems. Failing to response within the real-time constraint leads to failures of the system. Real-time systems are often classified into hard, firm and soft real-time based on the impact of missing the deadline [18]. Note that system that have to respond seemingly instantaneous, such as graphical user interfaces or instant chatting/calling, are not real-time. There are no hard deadlines which the system has to respond within and the system does not fail if some delay does occur. Only user experience is slightly impacted.
3. **A system that can process data faster than it acquires data.** This is technically not real-time, however, it is often used as such in academic literature. It is important for real-time system to process data faster than it acquires it so it does not start to lag behind, however, this is an implicit deadline. Not having this deadline explicit may lead to non-sensible expectations of the system.

Even though the first definition is very relevant when working with audio, it is not relevant for us. The audio drivers of operating systems handle all timing for us. We simply have to wait for samples to be recorded and made available to our program. We only have to keep the sampling rate in mind when working directly with the samples.

In order to allow guitarists to use their guitar as a MIDI instrument, our system has to respond within a small time frame. On top of that, if the system fails to respond quickly enough, the usefulness of the result degrades, as timing is very important when playing an instrument. These restrictions would classify our system as a soft real-time system. We choose a real-time of constraint of TODO milliseconds. We elaborate on this choice in section 4.1.

Other work in real-time pitch estimation often uses the third definition of real-time. This is a problem when using Fourier transform based methods, as many papers choose large frame sizes to get a high

resolution in the frequency domain. For instance, in order to discern the two lowest notes on a guitar which are 4.9 Hz apart, we would need a frame length of 204 milliseconds. This is well over our real-time constraint and would be unplayable for any musician. Other papers which do explicitly set a real-time constraint, choose very high constraints from 140 ms [12] up to 360 ms [11]. We have found no papers which choose a sensible real-time constraint.

3.4 Music theory and notation

In modern western music, we use the twelve-tone equal temperament (12-TET) music system. This system divides an octave, which is the interval between a pitch and another pitch with double the frequency, into twelve equally spaced semitones on the logarithmic scale. The logarithmic scale is used such that the perceived interval between two adjacent notes is constant [22]. As a result, the ratio between two frequencies in an n -semitone interval is $\sqrt[12]{2}^n$ or $2^{\frac{n}{12}}$, invariant to pitch. A semitone further be divided into 100 logarithmically scales cents. Cents are often used to measure dissonance.

Using scientific pitch notation, every note can be uniquely identified by combining the traditional note names A to G (with accidentals such as \sharp and \flat) with an octave number (e.g. E_3^\flat). An octave starts at C, which means the octave number increases between B and C. This counter intuitively implies that A_3 is higher than C_3 . Note that in 12-TET, C^\sharp and D^\flat are enharmonically equivalent. In this thesis, we will always refer to the sharp (\sharp) note instead of the enharmonically equivalent flat note (\flat). The range of a typical guitar in standard tuning is from E_2 up to E_6 .

The 12-TET music system only describes the relation between two notes in an interval. In order to play with other musicians in harmony, an arbitrary note has to be tuned to a specific frequency. Per ISO 16, the standard tuning frequency of the A_4 is 440 Hz within an accuracy of 0.5 Hz [3].

Using the above information, we can translate frequencies into scientific note names and vice versa. In order to numerically work with note names, we assign a value to each note as shown in Table 1.

name	number	name	number
C	0	F \sharp	6
C \sharp	1	G	7
D	2	G \sharp	8
D \sharp	3	A	9
E	4	A \sharp	10
F	5	B	11

Table 1: The number corresponding to each note name

In order to make calculations easier, we use C_0 as a tuning note instead of A_4 . We can calculate the frequency of C_0 using the fact that C_0 is 57 semitones lower than A_4 :

$$f_{C_0} = f_{A_4} * 2^{\frac{-57}{12}} = 440 * 2^{\frac{-57}{12}} = 16.352 \text{ Hz}$$

We assume that $a \bmod b$ always return a number c for which $0 \leq c < b$. Some programming languages allow the modulo operator to return a value c for which $-b < c < b$, resulting in $-13 \bmod 10 = -3$ instead of $-13 \bmod 10 = 7$.

We can calculate the frequency f_{N_O} , where N is the note name which is represented by a numerical value given by Table 1 and O is the octave number using:

$$f_{N_O} = f_{C_0} * 2^O * 2^{\frac{N}{12}} = f_{C_0} * 2^{O + \frac{N}{12}}$$

To calculate the closest note N_O corresponding to a frequency f , we first calculate the number of semitones n_s between the tuning note f_{C_0} and f :

$$n_s = \left\lceil 12 * {}^2\log \frac{f}{f_{C_0}} \right\rceil$$

Here, $\lceil \dots \rceil$ denotes rounding to the nearest integer. Rounding this number causes us to find the closest note to f . Now we can calculate N and O as follows:

$$N = n_s \bmod 12$$

$$O = \left\lfloor \frac{n_s}{12} \right\rfloor$$

In order to calculate the error e in cents of the given frequency to the closest tuned note, we first calculate the tuned frequency f_t of the closest note:

$$f_t = f_{C_0} * 2^{\frac{n_s}{12}}$$

Then the error e can be calculated using:

$$e = 1200 * {}^2\log \frac{f}{f_t}$$

3.5 Envelope and transients

The perceived loudness of a note over time can be described using an ADSR envelope. The ADSR envelope of a played note is the convex hull of the waveform of the signal, see Figure 4 for an example. This convex hull can be divided in four parts: Attack, Decay, Sustain, Release. When a note is strummed on the guitar, a percussive sound is generated which causes a loud and sharp attack along with the note. This percussive sound quickly decays and only the actually fretted note will sustain. Finally, when the note is released, it dies out quickly.

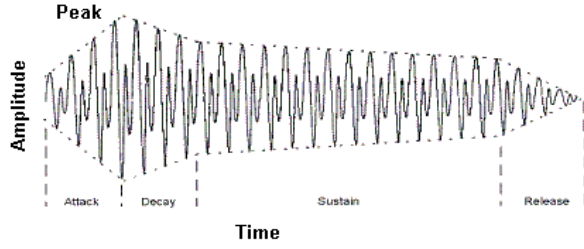


Figure 4: Example of an ADSR envelope (TODO: Better image)

As mentioned before, when strumming a note, a percussive sound is generated. This sound is called a transient and contains a high degree of non-periodic components. Transients appear very chaotically in the frequency domain and thus is often considered as noise. Since a transient is of high amplitude, it overshadows the note that will eventually sustain. Because of this, we cannot use the samples from a transient for pitch estimation. This in turn increases our minimum latency, as we have to wait for samples which do not contain the transient any more.

3.6 Fundamental frequency, overtones and timbre

When playing a note on an instrument, many sine waves are generated. The most notable frequency is called the fundamental frequency and determines what note is actually played. Integer multiples of the fundamental frequency can resonate and give rise to harmonic overtones [25]. In practice, these overtones are not exact integer multiples due to non-linear effect.

Furthermore, many other frequencies are generated along with the fundamental and its overtones. The instrument specific pattern of these frequencies, along with its envelope, is called the timbre of the instrument [15]. The timbre is what separates the

sound of the same note played on two different instruments [22]. Generally, the amplitude of the timbre frequencies is low compared to the fundamental frequency and can be disregarded as noise in the frequency domain.

In Section 2, we mentioned overtones are dissonant with respect to notes in 12-TET. This is true for all overtones, except for octaves, which are all overtones numbers equal to $2^n - 1$ for every n . In Table 2, we show an example for the overtones of C_4 . Note that the series of errors is the same, regardless of what the starting note is.

n	f_{overtone}	closest note	f_{note}	error
0	261.626	C_4	261.626	-
1	523.251	C_5	523.251	0
2	784.877	G_5	783.991	1.955
3	1046.502	C_6	1046.502	0
4	1308.128	E_6	1318.510	-13.686
5	1569.753	G_6	1567.982	1.955
6	1831.379	$A_6^\#$	1864.655	-31.174

Table 2: Example of overtone series from C_4 and the errors compared to the closest note

TODO: Overtone overlap and polyphonic difficulty.

n	$f_0^{C_3} * n$	n	$f_0^{E_3} * n$	n	$f_0^{G_3} * n$
1	130.813	1	164.814	1	195.998
2	261.626	2	329.628	2	391.995
3	392.438	3	494.441	3	587.993
4	523.251	4	659.255	4	783.991
5	654.064	5	824.069	5	979.989

Table 3: Overtones of C, E and G

Note	f	Δf
$f_2^{C_3}$	392.438	0.443
$f_1^{G_3}$	391.995	
$f_4^{C_3}$	654.064	5.191
$f_3^{E_3}$	659.255	

Table 4: Colliding overtones

4 Thesis content

Actual research etc. Summary what was done in the research project that we build on. Notes on future work of research project.

4.1 Real-time constraint

Start with the factors coming into play when choosing the real-time constraint for our system (latency, played notes per second etc). Note on measured latency in Appendix A. Empirically found bounds on latency using our latency program

4.2 Software amplification

(Software representation of samples and FFT). The FFT works on floating point numbers but most audio interfaces give up to 24 bit integers... We found empirically that when amplifying the input signal in software, peaks in the frequency domain are much easier to detect. However, it has to be done carefully to prevent distortion artifacts.

5 Implementation

Details about the program I've written. Usage instructions, code choices, code structure, screenshots, expandability

6 Experiments

Note that the parameters were empirically optimized with informal experiments. Datasets. Actual experiments.

7 Conclusions

What we did in this thesis. Reflection on the performance of the system. Final reference to the source code.

8 Future work

What could still be improved/further researched.

A Measuring the latency of the AXON AX 100 mkII

Lekker meten en weten.

B Effect of different window functions

Plots met effect van verschillende window functions.

References

- [1] Webpage on the limits of zero-padding. <https://dspillustrations.com/pages/posts/misc/spectral-leakage-zero-padding-and-frequency-resolution.html#Frequency-Resolution>
Last accessed on 20-10-2021.
- [2] Webpage with interactive tool that shows the effect of zero-padding. <https://jackschaedler.github.io/circles-sines-signals/zeropadding.html>
Last accessed on 20-10-2021.
- [3] Iso 16:1975. acoustics — standard tuning frequency (standard musical pitch). 1975.
- [4] Eric J. Anderson. Limitations of short-time fourier transforms in polyphonic pitch recognition. Ph.d. qualifying project report, University of Washington, Department of Computer Science and Engineering, 1997.
- [5] Eric J. Anderson. Limitations of short-time Fourier transforms in polyphonic pitch recognition. *Technical report, Department of Computer Science and Engineering, University of Washington*, 1997.
- [6] Judith Brown and Miller Puckette. An efficient algorithm for the calculation of a constant Q transform. *Journal of the Acoustical Society of America*, 92:2698–2701, 11 1992.
- [7] Luc de Jonckheere. Real-time guitar transcription using Fourier transform based methods; a pitch estimation framework and overtone sieve algorithm. 2021.
- [8] Robert Dobre and Cristian Negrescu. Automatic music transcription software based on constant Q transform. *8th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)*, pages 1–4, 2016.
- [9] Zhiyao Duan, Bryan Pardo, and Changshui Zhang. Multiple fundamental frequency estimation by modeling spectral peaks and non-peak regions. *IEEE Transactions on Audio, Speech, and Language Processing*, 18(8):2121–2133, 2010.
- [10] Paolo Nesi Fabrizio Argenti and Gianni Pantaleo. *Automatic Music Transcription: From Monophonic to Polyphonic*, pages 27–46. Springer Berlin Heidelberg, 2011.

- [11] Xander Fiss. Real-time software electric guitar audio transcription. Master's thesis, Rochester Institute of Technology, 2011.
- [12] T. A. Goodman and I. Batten. Real-time polyphonic pitch detection on acoustic musical signals. *2018 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT)*, pages 1–6, 2018.
- [13] F.J. Harris. On the use of windows for harmonic analysis with the discrete fourier transform. *Proceedings of the IEEE*, 66(1):51–83, 1978.
- [14] Julius Orion Smith III. Spectral audio signal processing. 2013.
https://www.dsprelated.com/freebooks/sasp/Quadratic_Interpolation_Spectral_Peaks.html
Last accessed on 20-10-2021.
- [15] Kristoffer Jensen. Timbre models of musical sounds. 2021. PhD thesis, University of Copenhagen.
- [16] Xavier Serra Julius O. Smith III. PARSHL: An analysis/synthesis program for non-harmonic sounds based on a sinusoidal representation. 1987.
- [17] S.S. Limaye K.A. Akant, R. Pande. Accurate monophonic pitch tracking algorithm for QBH and microtone research. *The Pacific Journal of Science and Technology*, 11(2):342–352, 2010.
- [18] Hermann Kopetz. *Real-Time Systems: Design Principles for Distributed Embedded Applications*. Springer US, 1997.
- [19] Tiago Fernandes Tavares, Jayme Garcia Arnal Barbedo, Romis Attux, Amauri Lopes. Survey on automatic transcription of music. *Journal of the Brazilian Computer Society*, 19:589–604, 2013.
- [20] Michael F. Zbyszynski Matthew Wright, Ryan J. Cassidy. Audio and gesture latency measurements on Linux and OSX. *In Proceedings of the ICMC*, pages 423–429, 2004.
- [21] James A. Moorer. On the transcription of musical sound by computer. *Computer Music Journal*, 1(4):32–38, 1977.
- [22] Michael J. O'Donnell. Digital sound modeling: Perceptual foundations of sound.
http://people.cs.uchicago.edu/~odonnell/Scholar/Work_in_progress/Digital_Sound_Modelling/lectnotes/node4.html
Last accessed on 04-10-2021.
- [23] A. Schutz and D. Slock. Periodic signal modeling for the octave problem in music transcription. *in Proceedings of the 16th International Conference on Digital Signal Processing (DSP'09)*, pages 1–6, 2009.
- [24] Christian Schörkhuber and Anssi Klapuri. Constant-q transform toolbox for music processing. *Proc. 7th Sound and Music Computing Conf.*, 2010.
- [25] R. S. Shankland and J. W. Coltman. The departure of the overtones of a vibrating wire from a true harmonic series. *The Journal of the Acoustical Society of America*, 10(3):161 ff, 1939.
- [26] Gino Angelo Velasco, Nicki Holighaus, Monika Doerfler, and Thomas Grill. Constructing an invertible constant-q transform with nonstationary gabor frames. *Proceedings of the 14th International Conference on Digital Audio Effects, DAFx 2011*, 2011.
- [27] Kurt James Werner. The XQIFFT: Increasing the accuracy of quadratic interpolation of spectral peaks via exponential magnitude spectrum weighting. *Proceedings of the International Computer Music Conference*, pages 326–333, 2015.