

Diploma Thesis

Estimating Frequency and Amplitude of Sinusoids in Harmonic Signals

A Survey and the Use of Shifted Fourier Transforms

Konrad Hofbauer

Graz University of Technology
Graz University of Music and Dramatic Arts

April 2004

Supervisors

Philippe DEPALLE
B.Sc., M.A., D.E.A., Ph.D.
McGill University
Music Technology Area
Montréal, Canada

Robert HÖLDRICH
o.Univ.-Prof. Mag. DI Dr.
Graz University of Music and Dramatic Arts
Institute of Electronic Music and Acoustics (IEM)
Graz, Austria



UNIVERSITY
OF MUSIC AND
DRAMATIC ARTS
GRAZ - AUSTRIA



McGill



Abstract

Additive synthesis is a meaningful sound model for analysis and processing of sound that represents a signal as a sum of individually controllable sinusoids. The accuracy of twelve frame-based parameter estimators which determine frequency and amplitude of the partials in harmonic sounds are tested with sinusoidal, static, time-varying and natural-like test sounds. The estimation errors are determined by the known frequency content of the test signals and the output of the estimators. The errors are smallest for the method proposed here which directly fits the analytical shape of the window main lobe into the frequency spectrum. Among the commonly known estimators, derivative algorithm, phase vocoder and spectrum reassignment give the best numerical results. The Matlab implementation provides a platform to evaluate new approaches for frame-based parameter estimators and readily shows the effectiveness compared to the standard algorithms.

Parameter estimation directly from MDCT coefficients has limitations and further research is necessary in this area. The ODFT, a frequency-shifted discrete Fourier transform, is shown as the natural link between the DFT and the MDCT. A parameter estimator based on ODFT is presented. It provides accurate results and becomes useful in MDCT based audio coding. Furthermore it is exemplarily demonstrated that a combined analysis with DFT and ODFT can reduce the estimation error at low frequencies.

Title: Estimating Frequency and Amplitude of Sinusoids in Harmonic Signals – a Survey and the Use of Shifted Fourier Transforms

Keywords: Additive synthesis, analysis of sound, parameter estimation, frequency estimation, shifted discrete Fourier transform, modified discrete Fourier transform, SDFT, ODFT, MDCT

Zusammenfassung

In dieser Arbeit werden die gebräuchlichsten Verfahren zur Parameterbestimmung für die Additive Synthese numerisch verglichen. Harmonische Klänge werden dabei in ihre Teiltöne (zeitvariable Sinuskomponenten) zerlegt und deren Frequenz und Amplitude möglichst genau bestimmt. Die erzielte Genauigkeit wird mittels geeigneter Testsignale mit exakt definierten Partialtönen ermittelt. Unter den herkömmlichen Methoden ist der Schätzfehler am geringsten, wenn neben dem Amplituden- auch das Phasenspektrum berücksichtigt wird (Phase Vocoder, Spectrum Reassignment und Derivative Algorithm). Die kleinsten Fehler werden jedoch durch das direkte Einpassen der kontinuierlichen Fourier-Transformierten des Analysefensters in das diskrete Frequenz-Spektrum des Signals erzielt.

Die Bestimmung der Teiltonparameter direkt aus MDCT-Koeffizienten funktioniert nur eingeschränkt und bedarf weiterer Entwicklung. Die ODFT, eine diskrete Fourier-Transformation, deren Abtastpunkte im Frequenzbereich verschoben sind, stellt die Verbindung zwischen MDCT und DFT dar. Ein Parameterschätzer basierend auf der ODFT-Transformation liefert gute Ergebnisse und kann im Bereich MDCT-basierter Audio-Kodieralgorithmen zu einer Effizienzsteigerung führen. Des Weiteren konnte gezeigt werden, dass sich durch eine kombinierte Analyse mit DFT und ODFT die Schätzfehler im Bereich niedriger Frequenzen verringern.

Titel: Frequenz- und Amplitudenbestimmung von Partialtönen harmonischer Signale: Vergleich und Entwurf von ODFT- und MDCT-basierten Verfahren

Schlagwörter: Additive Synthese, Klanganalyse, Parameterschätzer, Frequenzbestimmung, Verallgemeinerte diskrete Fourier-Transformation, SDFT, ODFT, MDCT

Gewidmet meinen Eltern.

For my parents.

Acknowledgements

First and foremost I want to thank my parents Konrad and Annemarie Hofbauer. They gave me incredible love and support throughout the course of my education and encouraged me in every decision I made. Without them my studies would have never been possible.

Enormous thanks also go to my advisor Professor Philippe Depalle at McGill University, where I conducted my research. I am deeply indebted to him. He gave my studies a meaningful direction and has been a constant inspiration and guidance. He highly supported my work and my understanding to the subject wherever possible. His hospitality made my stay a very pleasant experience.

I am equally in debt to my academic supervisor Professor Robert Höldrich, who encouraged me in my research plans from the very first. My gratitude as well to Dr. Alois Sontacchi, who kindly proof-read this document.

Many thanks to François Thibault, Mark Zadel and all the other fellows in the Music Technology Area at McGill for the interesting discussions and the friendly and supporting atmosphere in the SPCL lab. I hope we will get a chance to collaborate again some time in the future. Great support also came from Julien Boissinot, Richard McKenzie and Professor Marcelo Wanderley. They kindly supplied me with a superb infrastructure and thanks to them I have never been without a friendly ear for discussions.

Finally, deepest thanks to my girl-friend Agata: She willingly provided spiritual assistance whenever I needed it and supported me in my ups and downs during my thesis with inexhaustible patience.

Thank you all very much!

Contents

Abstract	2
Zusammenfassung	3
Acknowledgements	4
1. Introduction	9
2. Sound Models and Additive Synthesis	10
2.1. Sound Models	10
2.2. Additive Synthesis Model	10
2.3. Analysis for Additive Synthesis	12
3. Frequency Representations based on Fourier Transforms	15
3.1. Basic Fourier Transforms	15
3.2. Shifted Fourier Transforms	16
3.3. Discrete Cosine Transforms	19
3.4. Modified Discrete Cosine Transform	20
3.5. Complex Modified Discrete Cosine Transform	23
4. Parameter Estimators based on the Discrete Fourier Transform	25
4.1. Estimation Methods based on Interpolation	25
4.1.1. Zero Order Interpolation	25
4.1.2. First Order Interpolation	26
4.1.3. Second Order Interpolation	27
4.2. Estimation Methods based on Linear Phase Evolution	31
4.2.1. Derivative Algorithm	31
4.2.2. Spectral Reassignment	33
4.2.3. Phase Vocoder	34
5. Frequency Estimators based on Shifted Fourier Transforms	36
5.1. Motivation	36
5.1.1. MDCT in Perceptual Audio Coding	36
5.1.2. MDCT Domain Processing	38
5.2. Frequency Response of the Sinusoidal Analysis Window	38
5.3. ODFT Analysis of a Sinusoidal Input	42
5.4. MDCT Analysis of a Sinusoidal Input	45

6. Systematic Performance Comparison	55
6.1. Analysis Methods	55
6.2. Test Signals	57
6.2.1. Single Stationary Sinusoids	57
6.2.2. Two Stationary Sinusoids	59
6.2.3. Single Non-Stationary Sinusoids	60
6.2.4. Musical Instrument Sounds	61
6.3. Peak Filtering	62
6.4. Characteristic Touchstones	63
6.5. Comparison Results and Interpretation	64
6.5.1. Single Stationary Sinusoids	65
6.5.2. Single Stationary Sinusoids with Noise	66
6.5.3. Two Stationary Sinusoids	67
6.5.4. Single Non-Stationary Sinusoids	69
6.5.5. Musical Instrument Sounds	70
6.5.6. Hann Window for all Methods	71
6.5.7. Sinusoidal Window for all Methods	72
6.5.8. Conclusion	72
7. Improvements	75
7.1. Parabolic Interpolation Between ODFT Coefficients	75
7.2. Direct Fitting of the Analysis Window Main Lobe	79
7.3. Removal of Mirror Frequencies	84
8. Summary	87
A. Appendix	88
A.1. Software Implementation	88
A.2. Simulation Results	89
Bibliography	108

List of Figures

2.1.	Sound processing system based on additive synthesis	12
2.2.	Analysis stage for additive synthesis	13
3.1.	DTFT and DFT of a windowed sinusoid	16
3.2.	DFT and ODFT of an arbitrary input signal	19
3.3.	MDCT, Overlap-Add and Time Domain Alias Cancellation	22
4.1.	Fitting of a triangle into a sequence of FFT data points	26
4.2.	Analysis window with a triangular spectral shape	28
4.3.	DFT and DTFT of a sinusoid with a triangular spectral shape window	28
4.4.	Fitting of a parabola through three points of a DFT spectrum	29
4.5.	Phase vocoder used for frequency estimation	35
5.1.	MDCT based audio coder	37
5.2.	MDCT based audio coder using ODFT	37
5.3.	Frequency response of the normalized sinusoidal window	41
5.4.	ODFT coefficients of a windowed sinusoid	42
5.5.	Proportions of ODFT magnitude coefficients	43
5.6.	MDCT coefficients of a sine and a cosine signal	48
5.7.	MDCT spectra of different sinusoids	49
6.1.	Histogram of frequency of single stationary sinusoid test signals	58
6.2.	Frequency of sinusoids in the two stationary sinusoids test signals	60
6.3.	Detected peaks for pairs of closely spaced partials	68
6.4.	Detected peaks for pairs of closely spaced partials (with sinusoidal window)	73
7.1.	Frequency error of <code>ParInt</code> for stationary sinusoids	76
7.2.	Frequency error of <code>OdftParInt</code> for stationary sinusoids	77
7.3.	Frequency error of <code>DualParInt</code> for stationary sinusoids	78
7.4.	Frequency error of <code>ParInt</code> for stationary sinusoids (with Hann window)	79
7.5.	DFT spectrum of a windowed sinusoid	80
7.6.	Detected Peaks for pairs of closely spaced partials (new methods)	82
7.7.	Frequency error of <code>DirWinFitAnal</code> for stationary sinusoids	83
7.8.	Output of a FFT based Hilbert transformer	85
7.9.	Magnitude response of a one-sided FIR Hilbert transformer	86

List of Tables

3.1. Standard types of the Discrete Cosine Transform (DCT-I to DCT-IV).	20
5.1. Determination of the sinusoid's initial phase quadrant in the MDCT domain.	51
6.1. Tremolo and Vibrato test sound parameters.	61
6.2. Parameters of the natural-like test sounds.	62
6.3. Amplitude and frequency errors with noiseless single stationary sinusoids. . .	65
6.4. Amplitude and frequency errors with noisy single stationary sinusoids.	67
6.5. Amplitude and frequency errors with single non-stationary sinusoids.	70
6.6. Amplitude and frequency errors with musical instrument sounds.	71
A.1. Full simulation results with best suitable windows.	90
A.2. Full simulation results with Hann window.	96
A.3. Full simulation results with sine window.	102

Chapter 1.

Introduction

Computers offer new possibilities concerning sound processing. The ongoing advancement in computer performance reinforces the use of additive synthesis as sound model. Additive synthesis is a spectrum modelling technique that creates sound by a sum of a large number of independently controllable sinusoidal partials.

For applications in computer music systems and electronic music, natural existing sounds are analysed in order to modify their partials. It is therefore necessary to extract the sinusoids from a signal, which is often harmonic. The so-called ‘parameter estimators’ then try to accurately determine the sinusoids’ parameters, namely frequency, amplitude and initial phase.

Numerous research on DFT-based parameter estimation has been undertaken so far. Many different procedures are available, quite a few of them claiming outstanding performance. In order to provide a basis for comparisons and further improvements, this thesis objectively evaluates the accuracy of the amplitude and frequency estimation of the most common parameter estimators. They are therefore implemented and tested with different categories of reference signals.

Driven by the increasing usage of shifted transforms (in particular the MDCT in advanced audio coding), parameter estimators that apply shifted Fourier transforms are closely investigated. Their results are compared with the conventional methods and improvements suggested.

This thesis report concerns itself with the determination of the following objectives: Chapter 2 gives an introductory overview on sound models and shows in particular the functional units of an additive synthesis system. A mathematical description of the used sound model is provided. As there are many transformations that are closely related to the Fourier transform (e.g. SDFT, ODFT, MDCT, ...), chapter 3 gives their definitions and shows some interconnections between them.

Then chapter 4 presents an overview of the most commonly used state-of-the-art DFT-based parameter estimators and explains their functionality. The succeeding chapter moves on to parameter estimators based on shifted Fourier transforms. The motivation for the usage is demonstrated and their underlying functionality is mathematically derived.

The above mentioned survey on the accuracy of the analysis methods is presented in chapter 6. The used methodology is addressed in detail and the resulting simulation results are provided and commented. Driven by the observations obtained in the comparisons some improvements are presented in the last chapter. This includes mechanisms for suppressing the influence of negative mirror frequencies, optimum interpolation techniques and ways of combined application of shifted and non-shifted Fourier transforms.

Chapter 2.

Sound Models and Additive Synthesis

In order to successfully synthesize or transform digital sounds on a computer, an appropriate model that describes the sound is necessary. It provides the vital link between the natural and the mathematical representation of sound. A good sound model should be flexible in terms of signals which it can represent. For musical applications it should as well provide parameters that allow musically expressive and meaningful transformations. For more detailed information please refer to [24, 29, 42].

2.1. Sound Models

Mainly four families of sound models are used nowadays for musical sound generation:

Temporal Representations are the most obvious way to describe an audio signal, as they are closely related to the *physical* nature of sound. Simply the amplitude of the sound wave is given as a function over time. A variety of signal processing techniques can be applied to the signal. Unfortunately only a very small number of these transformations are related to musical parameters.

Physical Models try to put the natural acoustical formation of a sound in its source into mathematical equations. For the case of e. g. an organ pipe, the shape and the material of the pipe are considered, and the physical behaviour of the enclosed air – which finally dispreads the sound – is mathematically expressed.

Abstract Models give the rules to create a sound in purely mathematically inspired equations. Without any analysis stage, these models are well suited to create *new* sounds, but fail in reproducing or manipulating natural sounds. An example for these models is the well known Frequency-Modulation-Synthesis (FM-Synthesis).

Spectral Models represent sound with a spectrum in the frequency domain. This is driven by a perceptual perspective, as the human ear makes on the basilar membrane as well some sort of time-frequency transformation [53]. Transformation and feature extraction is possible in a way which is closer related to perception.

2.2. Additive Synthesis Model

This work addresses issues in the last category of sound models, spectrum modelling. The main advantage of this type is the existence of analysis methods to extract the synthesis

parameters out of real sounds. This allows reproduction and modification of existing sounds using the extracted parameters.

Our particular approach is based on modelling sounds with sinusoids. A pure continuous-time sine wave $x(t)$ is given by

$$x(t) = A \sin(\Omega t + \Phi) \quad \text{with} \quad \Omega = 2\pi f \quad (2.1)$$

All variables are real numbers and

- A : amplitude ($A \geq 0$)
- Ω : radian frequency in rad/s
- f : frequency in Hz
- t : time in s
- Φ : initial phase in rad.

Sampling this signal $x(t)$ every T_s seconds gives a discrete-time signal which is mathematically represented by a sequence of numbers $x[n] := x(nT_s)$, where $x(nT_s)$ is the sample taken at time $t = nT_s$. So with

$$x(nT_s) = A \sin(\Omega T_s n + \Phi)$$

and $\omega := \Omega T_s = \frac{\Omega}{f_s}$ equation 2.1 becomes to its sampled discrete-time equivalent

$$x[n] = A \sin(\omega n + \Phi) \quad (2.2)$$

with

- n : integer sample number
- T_s : sampling period in s
- $f_s = \frac{1}{T_s}$: sampling frequency in Hz
- ω : normalized radian frequency in rad¹
- [...]: used for independent variable of a discrete-time sequence
- (...): used for independent variable of a continuous-time function.

The *additive synthesis model* represents a sound as a sum of 'quasi-stable' sinusoids (sinusoids with slowly time-varying amplitude and frequency). These sinusoidal components of the modelled sound are usually called '*partials*'. When using the term 'harmonics' for them, it expresses that the partials' frequencies are harmonically related among each other.

The sound $a(t)$ is modelled in the additive synthesis model by

$$a(t) = \sum_{p=1}^P A_p(t) \sin(\phi_p(t)) \quad (2.3)$$

¹Instead of the frequency Ω its normalized version $\omega = \frac{\Omega}{f_s}$ is often used implicitly, which is equivalent to setting $T_s = 1$ s and $f_s = 1$ Hz.

with the instantaneous phase of the p^{th} partial

$$\phi_p(t) = \Phi_p + 2\pi \int_{u=0}^t f_p(u) du. \quad (2.4)$$

where $A_p(t)$ and $f_p(t)$ are the time-varying amplitude and frequency of the p^{th} partial, Φ_p the initial phase and P the total number of partials.

2.3. Analysis for Additive Synthesis

The analysis stage of an additive synthesis system examines the time-varying spectral characteristics of a sound. It maps these characteristics to the parameters of the additive synthesis model described above, namely to the slowly time-varying amplitude and frequency (and sometimes initial phase) of each sine wave.

The block diagram of an additive synthesis system basically consists of three blocks (figure 2.1): The analysis stage extracts the parameters, which are modified in a transformation stage according to the desired musical application. The synthesis part then uses these modified parameters to generate the output sound according to equation 2.3. Smith and Serra give in [44] a concise report about the implementation of an additive synthesis system.

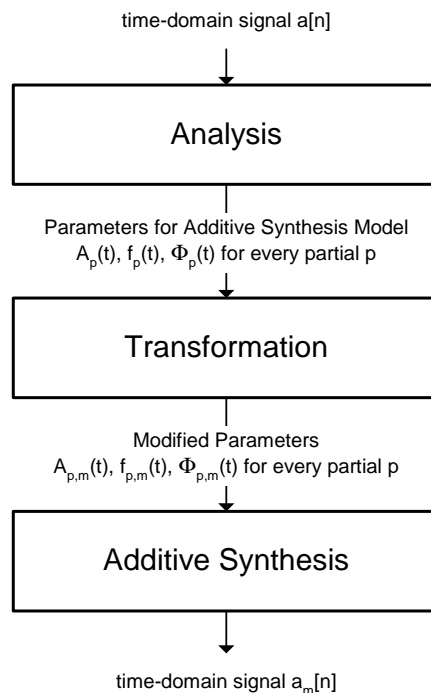


Figure 2.1.: Typical structure of a sound processing system based on additive synthesis.

We focus on the analysis stage, whereof figure 2.2 on the following page shows its general structure.

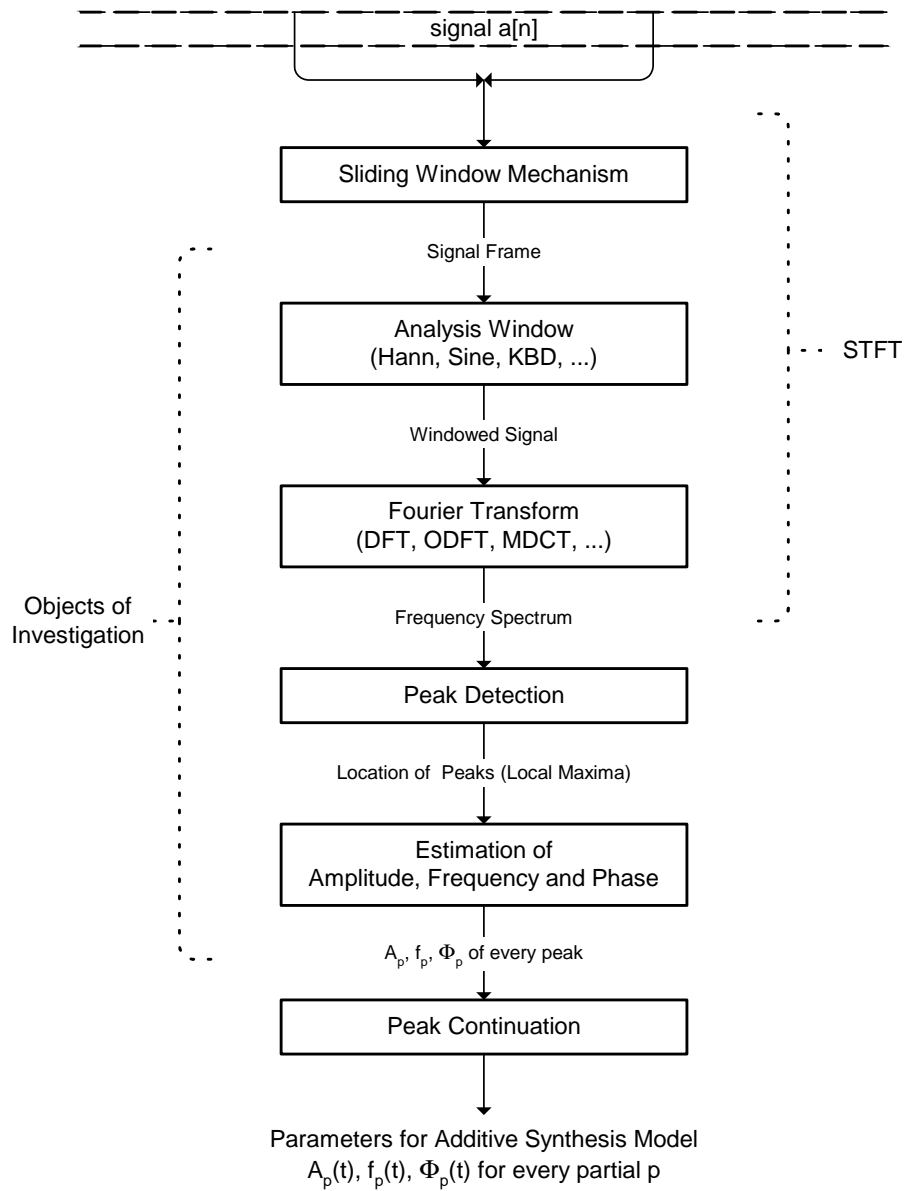


Figure 2.2.: Typical structure of the analysis stage for additive synthesis.

In computer music applications for the purpose of spectral analysis the *short-time Fourier transform* (STFT) is widely used [2, 44]. The STFT is often as well referred to as time dependent Fourier transform.

In the STFT first a sliding window mechanism is applied on the input signal $x[n]$. That is, the signal is structured in frames of length M whose start points are separated by the hop size R . As a consequence the frames overlap by $M - R$ samples. An ascending index r is assigned to the frames.

Every frame x_r of the signal is then multiplied by an analysis window function w of the same length to smooth out the impact of discontinuities of the time signal occurring at the frame borders. The choice of the analysis window is a well-developed topic and influences the spectral resolution of the analysis. A discussion on windows is beyond the scope of this document. An introduction is given by Harris in [16] with supplementations by Nuttall in [28].

The windowed frames are then possibly extended with zeros on one or both sides to a length N . The zero-padding and interpolation factor is given by the ratio $\frac{N}{M}$. Zero-padding in the time domain corresponds to interpolation in the frequency domain. On the zero-padded and windowed frames a DFT is applied.

With a zero-padding factor of $\frac{N}{M} = 1$ the STFT of a signal $x[n]$ is given by:

$$X_{STFT}(r, k) = \sum_{n=-\infty}^{\infty} w[n - rR]x[n]e^{-j\frac{2\pi}{N}nk} \quad (2.5)$$

The next step is to detect the prominent peaks in the spectrum of the current frame. The magnitude spectrum is therefore searched for local maxima, and every local maximum above a certain threshold is considered as a potential peak. As $X_{STFT}(r, k)$ is discrete in frequency, the frequency of the detected peak is only accurate within one frequency bin or $\frac{f_s}{N}$ (see chapter 3).

The accuracy can be highly increase by so-called ‘parameter estimators’. With different methods they interpolate between the discrete spectrum values or make use of the phase spectrum of the signal to estimate the frequency, amplitude and sometimes phase of the peak very accurately. The different methods available for this purpose are the main scope of this report.

These parameter estimators are usually much more efficient than interpolating by zero-padding, as the latter requires a much longer FFT length. Nevertheless a combination of both can be successfully applied. [42]

The last step in the analysis stage is to assign the detected peaks to the different partials of the additive synthesis model. This process of tracking peaks from frame to frame is as well called ‘partial tracking’. Therefore the peaks in adjacent frames are connected to each other by peak trajectories. The algorithm tries to find continuations of peaks in subsequent frames e.g. with a “birth death concept” and connection probabilities [24] or by the application of hidden Markov models [9]. A peak that was carried forward for several frames is then considered as partial of the sound.

Chapter 3.

Frequency Representations based on Fourier Transforms

As shown in chapter 2 we study a sound model that represents sound as a mix of sine waves, which is similar to the Fourier series representation of periodic signals. As the parameters of the model's sinusoids slowly evolve in time, the signals are strictly speaking no more periodic. We therefore extend the Fourier series approach and give an overview of Fourier transforms and related transformations. For more information see [43, 29, 32, 11, 24].

3.1. Basic Fourier Transforms

The continuous-time Fourier transform (FT) of a signal $x(t)$ with $-\infty \leq t \leq \infty$ is defined as

$$X(\omega) = \int_{t=-\infty}^{\infty} x(t)e^{-j\omega t} dt \quad (3.1)$$

with $j = \sqrt{-1}$ being the imaginary unit.

The discrete-time Fourier transform (DTFT) of a sequence $x[n]$ with $-\infty \leq n \leq \infty$ is similarly defined as

$$X_{DTFT}(\omega) = \sum_{n=-\infty}^{\infty} x[n]e^{-j\omega n} \quad (3.2)$$

which is often simply called the Fourier transform of a sequence $x[n]$ and is continuous in frequency.

If we consider only a finite length signal $x[n]$ with length N and non-zero values only between 0 and $N-1$, the DTFT becomes:

$$X_{DTFT}(\omega) = \sum_{n=0}^{N-1} x[n]e^{-j\omega n} \quad (3.3)$$

The discrete Fourier transform (DFT) of a finite length sequence $x[n]$ with $0 \leq n \leq N-1$ and length N is linked to the Fourier transform by sampling the DTFT at distinct uniformly spaced frequencies ω_k within the range $0 \leq \omega_k < 2\pi$ [51].

$$X_{DFT}(k) = \sum_{n=0}^{N-1} x[n]e^{-j\omega_k n} \quad \text{with} \quad \omega_k = \frac{2\pi}{N}k \quad (3.4)$$

where

$$\begin{aligned} X_{DFT}(k) &: \text{DFT of } x[n] \text{ or 'spectrum'} \\ k &: \text{DFT bin number with } 0 \leq k \leq N - 1 \\ \omega_k &: \text{normalized radian frequency at bin } k \end{aligned}$$

Although $x[n]$ and $X_{DFT}(k)$ are defined in the range $0 \leq n, k \leq N - 1$ only, in operations involving the DFT they take on an implicit periodicity with period N [38]. The inverse DFT results in a periodic extension of the signal $x[n]$, which is an infinite duration N -periodic signal created by concatenating the original finite length signal $x[n]$.

The so-called fast Fourier transform (FFT) is a computationally efficient implementation of the DFT. The choice of the transform length N is restricted, but the computational complexity is significantly reduced from $O(N^2)$ to $O(N \log(N))$. In practice we usually set the length N as a power of two.

Figure 3.1 illustrates the connection between the DFT and the DTFT: A pure sinusoid with frequency $l = 100.3$ bins is windowed by a Hann window. Its discrete-time Fourier transform DTFT is plotted with a continuous line. The DFT of the same windowed sinusoid is represented by equally spaced samples of the DTFT, marked in the figure with small circles.

The distance between the frequency samples is one frequency bin and equals sampling-rate f_s divided by the length N of the analysis window. E. g. with the used sampling rate of $f_s = 44,1$ kHz and a window length of $N = 1024$ samples the distance between two frequency bins is approximately 43 Hz.

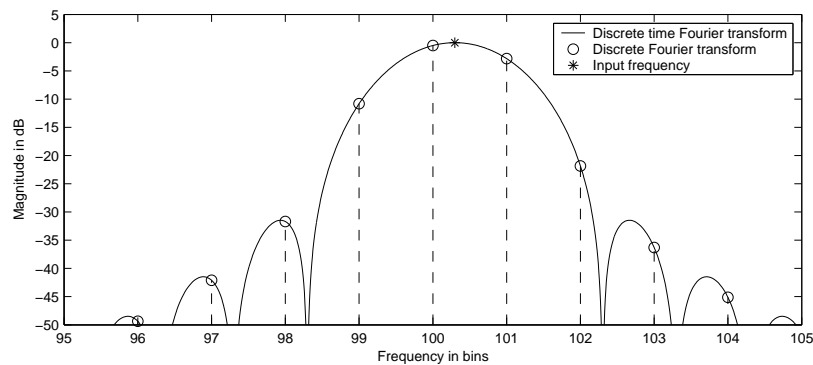


Figure 3.1.: DTFT and DFT of a windowed sinusoid ($N = 1024$, $l = 100.3$, Hann window).

3.2. Shifted Fourier Transforms

The shifted discrete Fourier transform (SDFT) is defined in its general form as:

$$X_{SDFT,a,b}(k) = \sum_{n=0}^{N-1} x[n] e^{-j \frac{2\pi}{N} (n+a)(k+b)} \quad (3.5)$$

with

a : arbitrary time domain shift

b : arbitrary frequency domain shift

The SDFT can be seen as a generalization of the DFT, which allows us to shift the samples in time as well as in the frequency domain with respect to the signal and its spectrum coordinate systems (see [52] and [48]).

In the conventional DFT, the periodicity is given simply by a periodic extension. Whereas in the generalized case, every further period of the sequence is rotated once more by an angle proportional to the frequency or time shift. The cyclicity can be expressed by:

$$x[n + Np] = x[n]e^{-j2\pi pb} \quad (3.6)$$

$$X_{SDFT,a,b}(k + Np) = X_{SDFT,a,b}(k)e^{j2\pi pa} \quad (3.7)$$

where p is an integer number.

As shown in [52], a *shifted sequence* rotates the SDFT coefficients by an angle proportional to the sequence shift *and* the transform shift parameters:

$$SDFT \{x[n + n_0]\} = X_{SDFT,a,b}(k)e^{-j2\pi n_0(k+b)/N} \quad (3.8)$$

$$X_{SDFT,a,b}(k + k_0) = SDFT \left\{ x[n]e^{-j2\pi k_0(n+a)/N} \right\} \quad (3.9)$$

Ordinary DFT

Several special cases of the SDFT can be considered, depending on the specific values for a and b .

With no time shift ($a = 0$) and no frequency shift ($b = 0$) the SDFT becomes the usual DFT as already given in equation 3.4 on page 15:

$$X_{SDFT,0,0}(k) = \sum_{n=0}^{N-1} x[n]e^{-j\frac{2\pi}{N}nk} = X_{DFT}(k) \quad (3.10)$$

This DFT is periodic both in time and frequency domain:

$$\begin{aligned} X_{DFT}(k) &= X_{DFT}(k + Np) && \text{(result of sampling the signal)} \\ x[n] &= x[n + Np] && \text{(result of sampling the spectrum)} \end{aligned}$$

where p is an integer number.

Time-Shifted DFT

With a time shift of half a sample ($a = \frac{1}{2}$) and no frequency shift ($b = 0$) the SDFT becomes the odd-time DFT:

$$X_{SDFT,\frac{1}{2},0}(k) = \sum_{n=0}^{N-1} x[n]e^{-j\frac{2\pi}{N}(n+\frac{1}{2})k} \quad (3.11)$$

It is periodic in time and ‘anti-periodic’ in the frequency domain:

$$\begin{aligned} X_{SDFT,\frac{1}{2},0}(k) &= (-1)^p X_{SDFT,\frac{1}{2},0}(k + Np) \\ x[n] &= x[n + Np] \end{aligned}$$

Frequency-Shifted DFT

With no time shift ($a = 0$) and a frequency shift of half a DFT bin ($b = \frac{1}{2}$) the SDFT becomes the odd-frequency DFT, or e. g. in [14] as well just called Odd-DFT (ODFT):

$$X_{ODFT}(k) = X_{SDFT,0,\frac{1}{2}}(k) = \sum_{n=0}^{N-1} x[n] e^{-j\frac{2\pi}{N}n(k+\frac{1}{2})} \quad (3.12)$$

It is periodic in the frequency and anti-periodic in the time domain:

$$\begin{aligned} X_{ODFT}(k) &= X_{ODFT}(k + Np) \\ x[n] &= (-1)^p x[n + Np] \end{aligned}$$

For real valued input signals $x[n] \in \mathbb{R}$ another symmetry property can be easily shown:

$$X_{ODFT}(k) = X_{ODFT}(N - k - 1)^*$$

where z^* denotes the complex conjugate of the complex number z .

Regarding figure 3.1 on page 16 again, the ODFT is very similar to the DFT. Both are samples of the DTFT, the samples just taken at different places (or frequencies). In the case of the ODFT, its sample locations are positioned directly between two DFT sample bins. The sample bins are shifted by $b = \frac{1}{2}$ bins, which corresponds to a shift of $\omega = \frac{1}{2N}$ in normalized frequency. The connection to the DTFT definition in equation 3.3 on page 15 is given by:

$$\omega_k = \frac{2\pi}{N}(k + \frac{1}{2}) \quad (3.13)$$

The same values could be achieved by zero padding the input signal from length N to length $2N$, performing a DFT of length $2N$ and then only considering the odd values. This is because zero padding just provides a finer sampling of the spectrum and zero padding to double the length gives exactly one additional value each in between the original DFT bins.

Although the total number of frequency bins N is of course the same for DFT and ODFT, in the ODFT we need one spectral sample less to represent a real signal: As the ODFT does not sample the frequencies $f = 0$ and $f = \frac{f_s}{2}$, we have exactly $\frac{N}{2}$ bins below the Nyquist frequency, and the same number of mirror frequency bins above it. In contrast the DFT gives $\frac{N}{2} + 1$ frequency bins within $0 \leq f \leq \frac{f_s}{2}$. Figure 3.2 on the next page illustrates this fact by a DFT and ODFT transform of an artificially designed input signal. This somehow nicer symmetry of the ODFT coefficients could be useful in some applications, especially in terms of window design issues.

Time- and Frequency-Shifted DFT

With a time shift of half a sample ($a = \frac{1}{2}$) and a frequency shift of half a DFT bin ($b = \frac{1}{2}$) the SDFT becomes the odd-time odd-frequency DFT (O²-DFT):

$$X_{SDFT,\frac{1}{2},\frac{1}{2}}(k) = \sum_{n=0}^{N-1} x[n] e^{-j\frac{2\pi}{N}(n+\frac{1}{2})(k+\frac{1}{2})} \quad (3.14)$$

It is anti-periodic in the time and frequency domain:

$$\begin{aligned} X_{SDFT,\frac{1}{2},\frac{1}{2}}(k) &= (-1)^p X_{SDFT,\frac{1}{2},\frac{1}{2}}(k + Np) \\ x[n] &= (-1)^p x[n + Np] \end{aligned}$$

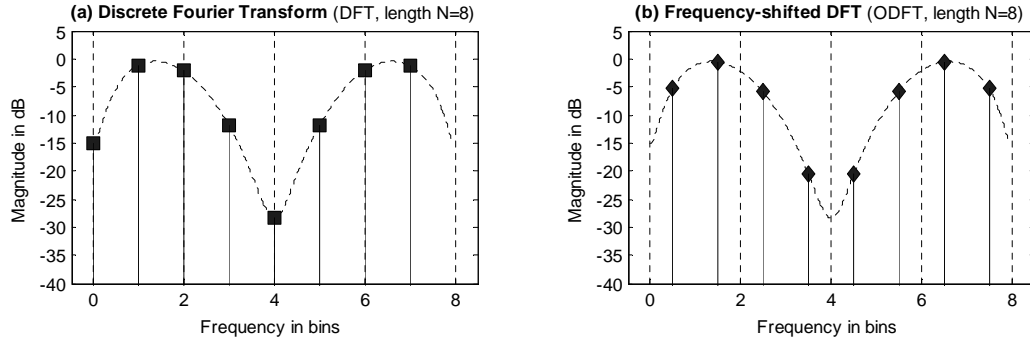


Figure 3.2.: 8-point DFT and ODFT of the same arbitrary input signal.

3.3. Discrete Cosine Transforms

The discrete cosine transform (DCT) and its advantages such as the close relation to the Karhunen-Loève transform where first shown in [1]. It is a frequency transform similar to the shifted discrete Fourier transform (SDFT) [50]:

$$X_{DCT,a,b}(k) = \sum_{n=0}^{N-1} x[n] \cos\left(\frac{\pi}{N}(n+a)(k+b)\right) \quad (3.15)$$

It is defined only for input signals $x[n]$ which are real valued and with even symmetry (i. e. $x(n) = x(-n)$). Parameters a and b are 0 or $1/2$ again.

Since the Fourier transform of a real and even function is respectively even and real again, the DCT is closely related to the shifted DFT. If we assume real and even input and take $X_{SDFT,a,b}(k)$ as defined in equation 3.5 on page 16, we can follow:

$$\begin{aligned} X_{SDFT,a,b}(k) &= \Re[X_{SDFT,a,b}(k)] \\ &= \Re\left[\sum_{n=0}^{N-1} x[n] e^{-j\frac{2\pi}{N}(n+a)(k+b)}\right] \\ &= \sum_{n=0}^{N-1} x[n] \cos\left(2\frac{\pi}{N}(n+a)(k+b)\right) \end{aligned} \quad (3.16)$$

This result equals almost the DCT definition in equation 3.15, except the factor of two in the cosine argument. Mathematically expressed, this simply means that the basis functions of the new signal space oscillate twice as fast.

In literature (e. g. [40] and [45]), mainly four different types of the DCT are defined, depending on the values of the parameters a and b (table 3.1 on the following page).

Let us take a look at the DCT-II, which is so defined by

$$X_{DCT,\frac{1}{2},0}(k) = \sum_{n=0}^{N-1} x[n] \cos\left(\frac{\pi}{N}(n+1/2)k\right) \quad (3.17)$$

in order to show its relation to the the SDFT with the same parameters $a = \frac{1}{2}$ and $b = 0$.

DCT type	a	b
DCT-I	0	0
DCT-II	1/2	0
DCT-III	0	1/2
DCT-IV	1/2	1/2

Table 3.1.: Standard types of the Discrete Cosine Transform (DCT-I to DCT-IV).

We therefore apply the time-shifted DFT on a double length real even sequence defined by:

$$x[n] = x[2N - 1 - n] \quad (3.18)$$

$$x[n] = x[n + 2Np] \quad \text{with } p \text{ integer and } 0 \leq n \leq 2N - 1 \quad (3.19)$$

$$X_{SDFT, \frac{1}{2}, 0}(k) = \sum_{n=0}^{2N-1} x[n] e^{-j \frac{2\pi}{2N} (n + \frac{1}{2})k} \quad (3.20)$$

Due to equation 3.18 $x[n]$ is even around time $n = N - \frac{1}{2}$, with $x[N - 1] = x[N]$. Additionally, by inspection it shows up that the exponential term is anti-symmetric around the same point $n = N - \frac{1}{2}$ so that $f(n) = (f(2N - 1 - n))^*$ with $f(n) = e^{-j \frac{2\pi}{2N} (n + \frac{1}{2})k}$. We so group together the two terms with the same value for $x[n]$, which is e.g. $n = 0$ and $n = 2N - 1$, $n = 1$ and $n = 2N - 2$, ... and so end up with half as much terms with $0 \leq n \leq N - 1$:

$$\begin{aligned} X_{SDFT, \frac{1}{2}, 0}(k) &= \sum_{n=0}^{N-1} x[n] \left[e^{-j \frac{2\pi}{2N} (n + \frac{1}{2})k} + e^{j \frac{2\pi}{2N} (n + \frac{1}{2})k} \right] \\ &= 2 \sum_{n=0}^{N-1} x[n] \cos \left(\frac{\pi}{N} (n + 1/2)k \right) \end{aligned} \quad (3.21)$$

which is (except a scaling factor of two) identical to the DCT-II of $x[n]$ with length N .

3.4. Modified Discrete Cosine Transform

The modified discrete cosine transform (MDCT) has become the most predominant time-frequency decomposition method for high-quality audio coding and compression.

The MDCT is used on blocks or frames of the input signal. Thereby usually a 50% time-domain window overlap is used. Putting some restrictions on the used window (namely the Princen-Bradley conditions [35, 36]), it is an interesting property of the MDCT, that the time-domain aliasing introduced by the transformation cancels out by the overlap and add. We so get perfect reconstruction of the input signal again. This process is called time domain aliasing cancellation (TDAC) and explained in detail in [35, 12].

The MDCT is usually defined as

$$X_{MDCT}(k) = \sum_{n=0}^{N-1} h[n] x[n] \cos \left(\frac{2\pi}{N} \left(n + \frac{1}{2} + \frac{N}{4} \right) \left(k + \frac{1}{2} \right) \right) \quad (3.22)$$

where $h[n]$ is the window function. A commonly used window is given by

$$h[n] = \begin{cases} \sin\left(\frac{\pi}{N}\left(n + \frac{1}{2}\right)\right) & \text{if } 0 \leq n \leq N - 1 \\ 0 & \text{else} \end{cases} \quad (3.23)$$

This sine window is the square root of a shifted Hann window, which can be easily shown applying the identity $\sin^2(x) = \frac{1}{2} - \frac{1}{2} \cos(2x)$. It fulfils the Princen-Bradley conditions, because:

1. The window overlap is not more than 50 %.

Proof: We use a step size of $\frac{N}{2}$ and a window length of N , so an overlap of 50 %.

2. $h[n]$ is symmetric with $h[n] = h[N - 1 - n]$.

Proof: $h[N - 1 - n] = \sin\left(\frac{\pi}{N}\left(N - 1 - n + \frac{1}{2}\right)\right) = \sin\left(\pi - \frac{\pi}{N}\left(n + \frac{1}{2}\right)\right) = h[n]$ (which additionally identifies $h[n]$ as linear phase filter).

3. $h[n]$ fulfils $h^2[n] + h^2\left[n + \frac{N}{2}\right] = 1$.

Proof: $h^2[n] + h^2\left[n + \frac{N}{2}\right] = \sin^2\left(\frac{\pi}{N}\left(n + \frac{1}{2}\right)\right) + \sin^2\left(\frac{\pi}{N}\left(n + \frac{N}{2} + \frac{1}{2}\right)\right) = \sin^2(a) + \sin^2\left(a + \frac{\pi}{2}\right) = \sin^2(a) + \cos^2(a) = 1$

In combination with the above given MDCT definition $h[n]$ so fulfils the requirements for perfect reconstruction.

The MDCT is usually applied with a 50 % window overlap. The length of the analysis is two times the (integer) window hop size and we can easily restrict N to be an even number. Then, as shown in [48], the MDCT transform coefficients show some symmetric properties, namely:

$$X_{MDCT}(N - 1 - k) = (-1)^{\frac{N}{2} + 1} X_{MDCT}(k) \quad (3.24)$$

For $N/2$ being an even number, as it is usually the case in audio applications, the MDCT coefficients show odd symmetry around $n = N/2$. The number of unique and independent MDCT output coefficients is so only half the number of input samples and we regard from now on only MDCT coefficients

$$X_{MDCT}(k) \quad \text{with} \quad 0 \leq k \leq \frac{N}{2} - 1. \quad (3.25)$$

Before showing in the next chapter the connection to the DFT, we here give first an illustrative example of the MDCT in conjunction with the TDAC concept in figure 3.3 on the following page [48, 49].

A 54 samples long arbitrary artificial time domain signal was created and is plotted in figure 3.3 on the next page, subplot (a). It is windowed by length N sine windows which overlap by 50 % (dashed line). An MDCT transform is then applied on the first window and the resulting coefficients plotted in subplot (b). As mentioned above, we can see the occurring subsampling: the N time domain samples result in only $N/2$ independent frequency domain coefficients. So obviously alias is introduced, which can be seen in subplot (c). This is the inverse MDCT of (b), the alias is highlighted by markers on the line. Subplots (d) and (e) are the MDCT and inverse MDCT of the second frame. The resulting time-domain signals (d) and (e) are then added together in the region where they overlap (between the time-points B and C). As can be seen in subplot (f), the original time-domain signal is perfectly

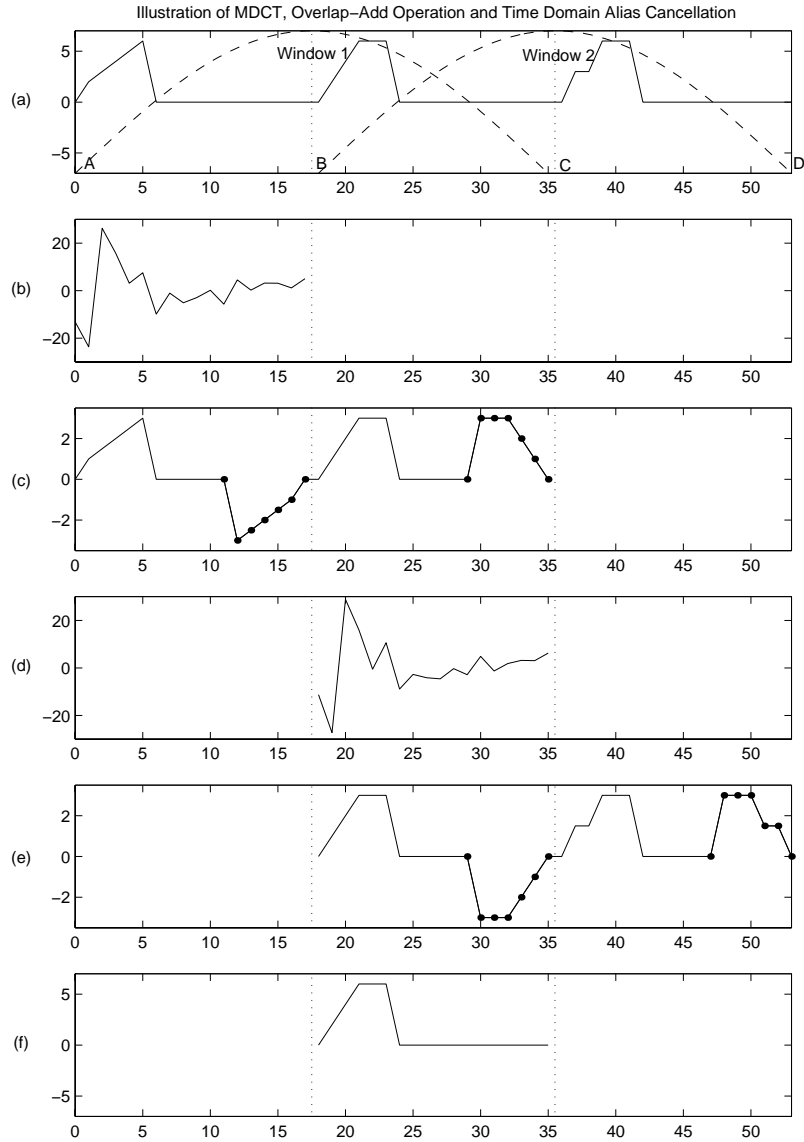


Figure 3.3.: Illustration of MDCT, Overlap-Add and Time Domain Alias Cancellation [48]:

- (a) Arbitrary artificial time signal. Dashed lines indicate overlapping windows.
- (b) MDCT of signal in frame one (window one).
- (c) IMDCT of Frame1-MDCT. Markers highlight introduced alias.
- (d) MDCT of signal in frame two (window two).
- (e) IMDCT of Frame2-MDCT. Markers highlight introduced alias.
- (f) Reconstructed signal between B and C by overlapping (c) and (e).

reconstructed in the overlapped part. As there is no overlapping available for the first half of the first window and the last half of the last window, perfect reconstruction cannot be achieved for these areas.

3.5. Complex Modified Discrete Cosine Transform

Similar to the MDCT the modified discrete sine transform (MDST) is defined by:

$$X_{MDST}(k) = \sum_{n=0}^{N-1} h[n]x[n] \sin\left(\frac{2\pi}{N}\left(n + \frac{1}{2} + \frac{N}{4}\right)\left(k + \frac{1}{2}\right)\right) \quad (3.26)$$

With the MDCT and the MDST of a windowed input signal $x[n]$ with $0 \leq k \leq \frac{N}{2} - 1$ we define the complex modified discrete cosine transform (CMDCT) $X_{CMDCT}(k)$:

$$X_{CMDCT}(k) = X_{MDCT}(k) + jX_{MDST}(k) \quad (3.27)$$

We want to show a connection of the MDCT to the DFT and evaluate therefore the complex conjugate of the CMDCT [23]:

$$\begin{aligned} X_{CMDCT}(k)^* &= X_{MDCT}(k) - jX_{MDST}(k) \\ &= \sum_{n=0}^{N-1} x[n] \cos\left(\frac{2\pi}{N}\left(n + \frac{1}{2} + \frac{N}{4}\right)\left(k + \frac{1}{2}\right)\right) \\ &\quad - j \sum_{n=0}^{N-1} x[n] \sin\left(\frac{2\pi}{N}\left(n + \frac{1}{2} + \frac{N}{4}\right)\left(k + \frac{1}{2}\right)\right) \\ &= \sum_{n=0}^{N-1} x[n] \cos\left(\frac{2\pi}{N}\left(k + \frac{1}{2}\right)n + \Phi_k\right) \\ &\quad - j \sum_{n=0}^{N-1} x[n] \sin\left(\frac{2\pi}{N}\left(k + \frac{1}{2}\right)n + \Phi_k\right) \end{aligned} \quad (3.28)$$

$$\text{with } \Phi_k = \left(\frac{2\pi}{N}\left(\frac{1}{2} + \frac{N}{4}\right)\left(k + \frac{1}{2}\right)\right)$$

Rearranging terms and using Euler's identity $e^{-j\phi} = \cos(\phi) - j \sin(\phi)$ we get

$$\begin{aligned} X_{CMDCT}(k)^* &= \sum_{n=0}^{N-1} x[n] e^{j\left(\frac{2\pi}{N}\left(k + \frac{1}{2}\right)n + \Phi_k\right)} \\ &= e^{-j\Phi_k} \sum_{n=0}^{N-1} x[n] e^{j\frac{2\pi}{N}\left(k + \frac{1}{2}\right)n} \\ &= e^{-j\Phi_k} X_{ODFT}(k) \end{aligned} \quad (3.29)$$

This shows, that the complex conjugate of the CMDCT is the frequency shifted DFT as given in equation 3.12 on page 18, phase-shifted by a constant angle Φ_k , which only depends

on the number k of the frequency bin. Furthermore this gives the connection between MDCT and SDFT, as:

$$\begin{aligned} X_{MDCT}(k) &= \Re[X_{CMDCT}(k)^*] = \Re[e^{-j\Phi_k} X_{ODFT}(k)] = \\ &= \Re[X_{ODFT}(k)] \cos(\Phi_k) + \Im[X_{ODFT}(k)] \sin(\Phi_k) \quad (3.30) \end{aligned}$$

So the MDCT coefficients can as well be obtained by computing the complex-valued ODFT, phase-shifting the coefficients by Φ_k as defined above and taking the real part of the resulting spectrum.

Chapter 4.

Parameter Estimators based on the Discrete Fourier Transform

This chapter presents an overview on the most widespread methods for estimating frequency, amplitude and phase of sinusoids in musical sounds. Their basic principles are described briefly, the choice of parameters used in the implementations is explained and pointers to detailed descriptions are provided. We focus on the methods selected in [17] and therefore restrict ourself on DFT-based estimators that extract parameters directly from a single signal frame.

For all methods the same basic analysis structure is used: an input signal $x[n]$ is split into (possibly overlapping) blocks of length M . These blocks are weighted by a temporal window function $w[n]$ and then an N -point DFT transforms these blocks into the frequency domain. The DFT length N equals the window length M , no oversampling or zero-padding is performed.

4.1. Estimation Methods based on Interpolation

As illustrated already in figure 3.1 on page 16 it seems possible to increase the accuracy of the parameter estimation by carefully interpolating between the DFT coefficients.

4.1.1. Zero Order Interpolation (Plain FFT without further Processing)

Zero order interpolation in this context means, that no interpolation at all is done in between the DFT samples. This method is included in the implementations for generic comparison only.

As analysis window $w[n]$ the ‘periodic’ Hann window is applied.¹ We compute the FFT magnitude spectrum and set all values below a certain threshold to zero. For the implementation a threshold value of -80 dB in terms of power is used. Therefore the amplitude gain of the Fourier transform has to be considered.

The local maxima k_m of this modified FFT magnitude spectrum are considered as direct representation of the input signal’s partials. The estimation is then given by the frequency corresponding to the FFT bin.

$$f_m = \frac{k_m f_s}{N} \quad (4.1)$$

The amplitude is scaled to one half of the maximum of the window’s Fourier transform $|W(0)|$, which can be expressed as the sum of the time domain window samples:

¹We use periodic windows (in contrast to symmetric ones) as defined in [16].

$$|W(0)| = \sum_{n=0}^{N-1} w[n] \quad (4.2)$$

So the amplitude estimation for the local maximum becomes to:

$$a_m = \frac{2|X_{DFT}(k_m)|}{\sum_{n=0}^{N-1} w[n]} \quad (4.3)$$

4.1.2. First Order Interpolation (Linear Interpolation)

The next logical step is to interpolate linearly between the DFT values, which means connecting the sample points with straight lines. This provides spectrum values in between the sample points, but does not improve parameter estimation, because we do not gain any more information about the position of the apex of the modulated window Fourier transform in the general case.

Triangle Analysis Algorithm

Keiler and Zölzer present in [18] and [3] an analysis algorithm which successfully applies linear approximation with using a special analysis window. The technique does not directly implement what is usually considered as linear interpolation between sample points, but extrapolates linearly from a set of sample points by a straight line to the desired location.

As explained before, a sinusoidal input modulates the window Fourier transform to the frequency of the input sine (and its negative mirror frequency). The idea is to use a window which magnitude frequency response has a triangular shape. So two straight lines can be fit into the DFT sample points around a local maximum, which correspond to the slopes of the triangle. As illustrated in figure 4.1 the exact frequency is then determined by the point of intersection of the two straight lines. For a good and more in-depth description of the algorithm please refer to [18].

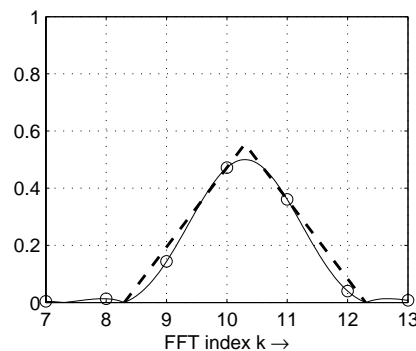


Figure 4.1.: Fitting of the triangle into a sequence of FFT data points in a least mean square error sense. [17]

Design of Analysis Windows with Triangular Spectral Shape

The window is designed in the frequency domain. It is defined by the sampling of a spectral triangle shape whose main parameter is the bandwidth of one triangle slope in frequency bins. In order to get the time domain window, an inverse DFT of length N is applied. This results in a temporal window, which is centred around $n = 0$ and is so non-causal. As a consequence it is delayed by $N/2$ samples to get a causal window.

The discrete-time Fourier transform of windows designed in this way do only approximately model the ideal triangle shape we started from. The resulting frequency response is a sum of scaled and frequency- and phase-shifted $\text{sinc}_N(\Omega)$ -functions. Figure 4.2 on the next page depicts the resulting window in time and frequency domain for a triangle slope bandwidth of $S = 4$ bins and a window length of $N = 1024$. Please note, that the DFT samples lie exactly on the predefined spots of the ideal triangle. Just the values of the continuous discrete-time Fourier transform in between them differ from the ideal straight line (which of course will have an impact on the accuracy of the parameter estimation).

The choice of the parameter S directly influences the window bandwidth. The length of the triangle base, which corresponds to the width of the main lobe of conventional windows, is double the length of S . For narrower triangles the frequency resolution increases. Wider triangles (bigger S) make the estimation more stable in the case of presence of noise, as a higher number of values is considered. Figure 4.3 on the following page shows the magnitude spectrum of a sinusoid on which a window of length $S = 4$ was applied.

In the implementation a slope bandwidth of $S = 2$ was used, in order to maintain a sufficient frequency resolution. The above described window design procedure then works as follows: Only three samples of the (normalized) discrete Fourier transform are given: the first is 1, the second 0.5 and the third is 0. To get a periodic window of length N , the N_{th} value is set to 0.5, too. All other values in between are zero. On this sequence an inverse DFT is applied. Only the real part is considered and for causality shifted by $\frac{N}{2}$, which results in the time domain analysis window.

One important information is neither given in [18], nor in [17], which uses $S = 2$, too: For this special case of $S = 2$ this design process results in the same sequence of samples than for the well known Hann window.

We can show this as well with the analytic expression of the resulting time domain window $w_S[n]$, which is given in [18, eq.7f]. With $S = 2$ this equation resolves to a scaled Hann window:

$$w_2[n] = \frac{1}{N} \left\{ 1 + 2 \sum_{k=1}^{S-1} \left(1 - \frac{k}{S} \right) \cos \left[\frac{2\pi}{N} k(n - N/2) \right] \right\} \Big|_{S=2} \quad (4.4)$$

$$= \frac{1}{N} \left[1 - \cos \left(\frac{2\pi}{N} n \right) \right] \quad (4.5)$$

4.1.3. Second Order Interpolation (Parabolic Interpolation)

Parabolic interpolation is a very widely used approach. It is based on the assumption that the main lobe of the discrete-time Fourier transform of the analysis window in a logarithmic dB-scale can be approximated by a parabola.

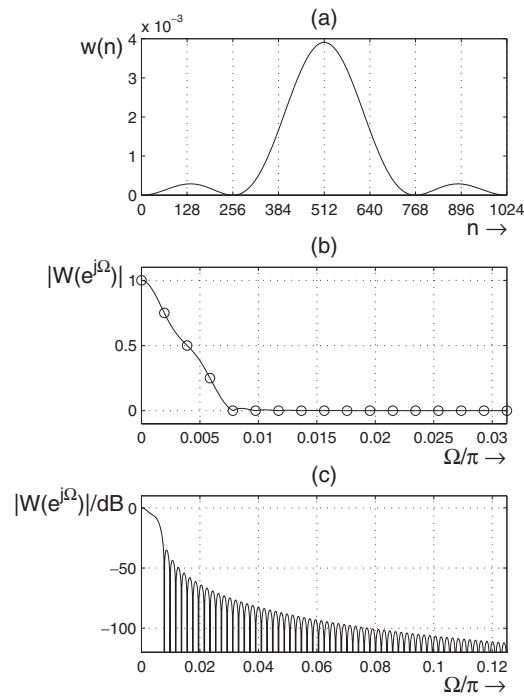


Figure 4.2.: Analysis window with a triangular spectral shape for $N = 1024$ and $S = 4$. (a) window function in time domain, (b) linear scale magnitude response and (c) logarithmic scale magnitude response. [18]

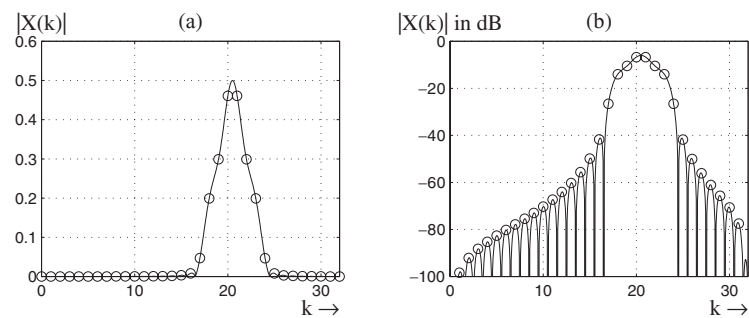


Figure 4.3.: Spectrum of frequency modulated window with triangular spectral shape, input signal $x[n] = \cos(\frac{2\pi}{N} 20.5n)$, $N = 64$ and $S = 4$. (a) Linear and (b) logarithmic scale. [18]

Again, a threshold-limited FFT magnitude spectrum of a windowed input signal is searched for local maxima. Then, a parabola (a second order polynomial) is fit into the DFT points around the maxima and the parabola's apex determines the frequency and amplitude estimation.

In our implementation we only use the local maximum and its right and left hand side direct neighbours. As long as these three points are not all positioned along a straight line, it is always possible to exactly fit a parabola through them. We can derive an analytic expression to get the position of the apex of the parabola with only the three magnitude values given. As mentioned before, we use the magnitude spectrum in a dB scale, so $X_{dB}(k) = 20 \log_{10}(|X_{DFT}(k)|)$, and we call k_0 the frequency index of the local maximum.

We start with the generic equation for a parabola which is parallel to the y-axis and open to the bottom:

$$(x - x_s)^2 = -2p(y - y_s) \quad (4.6)$$

with x_s, y_s being the coordinates of the apex A_0 and $p > 0$ being the parabola's parameter. The point A_0 is the apex of the parabola and A_{0+1} is the point on the parabola at the frequency of the apex plus one bin. We name the three DFT parameters which determine the parabola according to:

$$A_1 = X_{dB}(k_0 - 1) \quad A_2 = X_{dB}(k_0) \quad A_3 = X_{dB}(k_0 + 1) \quad (4.7)$$

The position of all these points is depicted in figure 4.4.

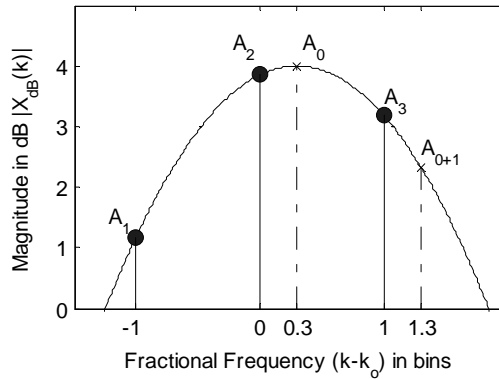


Figure 4.4.: Fitting of a parabola through three points of a DFT spectrum

We put the given values A_1, A_2, A_3 into the parabola equation which results in three equations for the three unknowns x_s, y_s and p :

$$(-1 - x_s)^2 = -2p(A_1 - y_s) \quad (4.8)$$

$$(0 - x_s)^2 = -2p(A_2 - y_s) \quad (4.9)$$

$$(1 - x_s)^2 = -2p(A_3 - y_s) \quad (4.10)$$

Solving this equation system leads to:

$$x_s = \frac{A_1 - A_3}{2(A_1 - 2A_2 + A_3)} \quad (4.11)$$

$$y_s = A_2 - \frac{(A_1 - A_3)^2}{8(A_1 - 2A_2 + A_3)} = A_2 - \frac{x_s}{4}(A_1 - A_3) \quad (4.12)$$

$$p = \frac{1}{-A_1 + 2A_2 - A_3} \quad (4.13)$$

The estimation of the partial's frequency is so given by $k_0 + x_s$ in DFT bins and its amplitude by $10^{(y_s/20)}$ times the amplitude scaling factor introduced by transformation and analysis window. Better amplitude estimation results are achieved by using the actual shape of the window main lobe. This technique is explained in detail in section 4.2.1, equation 4.31 on page 33.

If the local maximum is actually induced by a sinusoid, the resulting parabola should be approximately as wide as the window's main lobe. As a characteristic parameter we use the amplitude difference in dB between the apex A_0 and the amplitude A_{0+1} at a distance of one frequency bin away. This difference is compared to the corresponding value for the used analysis window (e.g. approximately -3 dB for a Hann window). If they differ too much, the corresponding peak is discarded.

Obviously this method is highly dependent on the used analysis window. The approximation error is made by using a parabola instead of the actual window shape. So the requirement on a "good" window for this method is, that its discrete-time Fourier transform main lobe in dB-scale is as close to a parabola as possible. The Hann window or a truncated Gaussian window could be used here, but best results are achieved with a special window described in [17]. It is again designed in the frequency domain and with inverse Fourier transform as described in section 4.1.2 on page 27, but starting from a parabolic shape in the frequency domain, instead of a triangular one.

The bandwidth S of the slope of the parabola in the frequency domain in bins is usually defined to $S = 2$ or $S = 3$. In the lower half of the frequency spectrum (the upper half is just mirrored) the first S values are set unequal to zero. With A_{0+1} being the desired amplitude of the magnitude spectrum in dB at bin $k = 1$ in terms of power of the magnitude at bin $k = 0$ (e.g. $A_{0+1} = -4.35$ dB), these S coefficients are given by:

$$X_{DFT}(k) = 10^{\left(\frac{A_{0+1}k^2}{20}\right)} \quad \text{for} \quad 0 \leq k \leq S - 1 \quad (4.14)$$

For a window with $S = 2$, in [17] a value of $A_{0+1} = -4.35$ dB is used.

Again, after mirroring the lower half around $\frac{N}{2} + 1$ to get a window which is periodic in frequency domain, an inverse Fourier transform of length N is performed. The resulting time domain window is shifted by $\frac{N}{2}$ to make it a causal and finally normalized to a maximum amplitude of one.

In order to make the algorithm more stable against the influence of noise, more than three points could be used. One can either consider not only the direct neighbours of the maximum but as well more adjacent points, what reduces the frequency resolution for closely spaced partials. Or one can oversample the signal by zero-padding. The parabola is then overdetermined and a numerical curve fitting approach is used which can consider a higher

number of points. An error function for the curve fitting error is set up and numerically minimized. An algorithm which proved to work well in practice is referred to as the Brent-method (or ‘Golden Section Search’), and its implementation is demonstrated in [34].

4.2. Estimation Methods based on Linear Phase Evolution

Although the interpolation methods often provide good results, they share a common insufficiency: they depend on a special analysis window. Therefore an alternative is to find methods which are no more based on interpolation of the spectrum but which try to exploit a special property that comes along with the discrete Fourier transform and the well defined class of signals we want to use the methods for: we can assume that there is only one sine wave which dominates the corresponding DFT channel at a local maximum. Furthermore we assume that the frequency of the partial changes only slowly in time and we can so consider that its frequency is constant within the analysis frame. As a consequence, the modelling of the phase evolution simplifies to

$$\phi_p(t) = \Phi_p + 2\pi f_p t \quad (4.15)$$

where f_p is the frequency of the p^{th} partial and Φ_p its initial phase. This describes that the phase of the sinusoid evolves linearly in time and results from equation 2.4 on page 12, with considering the assumptions stated above.

4.2.1. Derivative Algorithm

The derivative Algorithm was recently developed by Sylvain Marchand and is presented in detail in [24] and [10], with some improvements shown in [18]. The basic idea behind this algorithm is quite simple: the derivative of a sinusoid is a sinusoid again, multiplied by its frequency. E. g. for a continuous-time signal $x(t) = A \sin(\Omega t + \Phi)$ its derivative is given by $x'(t) = \Omega A \cos(\Omega t + \Phi)$.

Continuous-Time Signal

Considering more generally the expression of a single partial $a_p(t)$ as defined in equation 2.3 on page 11, its derivative $a_p'(t) = \frac{da_p}{dt}$ is given by:

$$a_p(t) = A_p(t) \sin(\phi(t)) \quad (4.16)$$

$$\frac{da_p(t)}{dt} = \frac{d}{dt}\{A_p(t)\} \sin(\phi(t)) + A_p(t) \frac{d}{dt}\{\sin(\phi(t))\} \quad (4.17)$$

As A_p is only slowly time varying in comparison to the second term, we assume $\frac{d}{dt}\{A_p(t)\} \approx 0$. We can so rewrite this equation to:

$$\frac{da_p(t)}{dt} \approx A_p(t) \cos(\phi(t)) \frac{d}{dt}\{\phi(t)\} \quad (4.18)$$

$$\text{with } \frac{d\phi(t)}{dt} = \frac{d}{dt} \left\{ \Phi_p + 2\pi \int_{u=0}^t f_p(u) du \right\} = 2\pi f_p(t) \quad (4.19)$$

$$\frac{da_p(t)}{dt} \approx 2\pi A_p(t) f_p(t) \sin\left(\phi(t) + \frac{\pi}{2}\right) \quad (4.20)$$

The frequency of the partial $f_p(t)$ is only slowly time varying, too. So we approximate it as being constant with f_p . We denote the discrete-time Fourier transform of $a_p(t)$ as $FT^0(f)$ and the Fourier transform of the derivative $\frac{da_p(t)}{dt}$ as $FT^1(f)$. For a given partial the magnitude spectra of both Fourier transforms have their maximum at the partial's frequency f_p . Putting this together with equations 4.16 and 4.18 we can conclude that:

$$f_p = \frac{FT^1(f_p)}{2\pi FT^0(f_p)} \quad (4.21)$$

Please note, that the effect of the analysis window is the same on $FT^0(f)$ and $FT^1(f)$ and is cancelled out by the above division.

Discrete-Time Signal

For a discrete-time signal $s[n]$ it is shown in [10] that

$$s'[n] = f_s \{s[n] - s[n-1]\} \quad (4.22)$$

is a valid approximation for the derivative of $s[n]$. The z -transform $S'(z)$ of this expression using the linearity and time shifting properties of the z -transform becomes to:

$$S'(z) = f_s (S(z) - z^{-1}S(z)) = H(z)S(z) \quad \text{with} \quad H(z) = f_s(1 - z^{-1}) \quad (4.23)$$

Evaluating the magnitude of $H(z)$ at the complex unit circle $z = e^{j\omega}$ results in the frequency response $H(\omega) := H(e^{j\omega})$:

$$|H(\omega)| = f_s \sqrt{2 - 2\cos(\omega)} = 2f_s \sin\left(\frac{\omega}{2}\right) \quad (4.24)$$

Let us now consider an input partial $a_p[n] = A_p \sin(\omega_p n + \Phi_p)$, which amplitude and frequency is constant within one analysis frame. It is windowed by a finite length window function $w[n]$, which Fourier transform is given by $W(\omega)$. Chapter 5.2 on page 38 will show, that the discrete-time Fourier transform $X_{DTFT}(\omega)$ of the windowed input partial is given by:

$$X_{DTFT}(\omega) = \frac{A}{2j} [e^{j\Phi_p} W(\omega - \omega_p) - e^{-j\Phi_p} W(\omega + \omega_p)] \quad (4.25)$$

The discrete-time Fourier transform $X'_{DTFT}(\omega)$ of the derivative of $a_p[n]$ according to equation 4.22 and using equation 4.23 is given by:

$$X'_{DTFT}(\omega) = H(\omega)X_{DTFT}(\omega) \quad (4.26)$$

For both ω_p and ω not being very low frequencies, we can assume that $W(\omega + \omega_p) \approx 0$. This is especially true for analysis windows with high side lobe attenuation. When we apply this approximation and look at the maxima of both DTFTs at ω_p , we can so compute a ratio $\gamma(\omega_p)$:

$$\gamma(\omega_p) = \frac{|X'_{DTFT}(\omega_p)|}{|X_{DTFT}(\omega_p)|} = \left| \frac{X_{DTFT}(\omega_p)H(\omega_p)}{X_{DTFT}(\omega_p)} \right| = 2f_s \sin\left(\frac{\omega_p}{2}\right) \quad (4.27)$$

We regroup this for ω_p :

$$\omega_p = 2 \arcsin\left(\frac{\gamma(\omega_p)}{2f_s}\right) \quad (4.28)$$

In order to use this for frequency estimation, we approximate $\gamma(\omega_p)$ by its value at the integer DFT frequency bin closest to ω_p : We compute the DFT spectra $X_{DFT}(k)$ and $X'_{DFT}(k)$ of the windowed partials $a_p[n]$ and $a'_p[n]$. This is equivalent to evaluating $X_{DFT}(\omega)$ and $X'_{DFT}(\omega)$ at the frequencies $\omega_k = \frac{2\pi k}{N}$, $0 \leq k \leq N - 1$. The maximum of $|X_{DFT}(k)|$ is at the frequency bin $k = k_0$, which is as well the position of the maximum of $|X'_{DFT}(k)|$. This location k_0 is the frequency bin closest to ω_p , i. e. $\left| \frac{\omega_p N}{2\pi} - k_0 \right| \leq 0.5$. With $\omega_0 = \frac{2\pi k_0}{N}$ we therefore get as frequency estimation

$$\omega_p \approx 2 \arcsin \left(\frac{\gamma(\omega_0)}{2f_s} \right) \quad (4.29)$$

or with f_p being the partial's frequency in Hz:

$$f_p \approx \frac{f_s}{\pi} \arcsin \left(\frac{|X'_{DFT}(k_0)|}{2f_s |X_{DFT}(k_0)|} \right) \quad (4.30)$$

The estimation error is consequently not made in the approximation of the derivative of the input signal, but by evaluating $\gamma(\omega)$ only at the DFT sample points instead of at the exact (but unknown) frequency ω_p .

For the amplitude estimation we again neglect the contributions of negative mirror frequencies. It is then based on the frequency estimation and the gain of the window at the computed fractional frequency, i. e.

$$A_p = 2 \frac{|X_{DFT}(k_0)|}{|W(\omega_p - \omega_0)|} \quad (4.31)$$

This requires the discrete-time Fourier transform of the window function, which is for most of the common windows given by analytical expressions. Otherwise a lookup table can be used, which is pre-computed by highly oversampling the analysis window (zero-padding the time window before applying the DFT).

4.2.2. Spectral Reassignment

Spectral Reassignment was initially designed by Kodera, Gendrin and Villedary to improve the readability of spectrograms [19]. Due to difficulties when implementing it in practice it was not widely used. Auger and Flandrin reformulated the algorithm [5, 4] and generalized it for all kinds of bilinear transforms. They presented efficient ways for implementing it and showed its practical use for audio analysis.

In classical short time Fourier analysis (STFT) the energy of a windowed signal frame is represented in single points, which are given by the centre of the analysis window (with a given length and overlap factor). The spectrogram is so not continuous, but discrete in time (determined by the hop size of the window) and discrete in frequency (determined by the window and DFT length). The idea of spectrum reassignment is, to move these points away from the time-frequency lattice to the centre of gravity of the energy included in the window. I. e. the point representing the frames energy is not any more in the centre of the window, but on the position where most of the energy is located in the window (as it is in the general case unevenly distributed). This leads to a reassignment of both the time and the frequency value. This is done by exploiting the phase information which is provided by the DFT but

not used in e.g. all the interpolation based methods. It helps to reduce the smearing of energy and to sharpen the amplitude and frequency estimates provided by the spectrogram [15]. An application to natural speech signals is presented in [33].

For a given time signal frame, the frequency reassignment is done by the following steps: First, the time signal $x[n]$ is multiplied with an arbitrary analysis window function $w[n]$ and its discrete Fourier transform $X_{DFT,w}(k)$ is computed. Second, the DFT $X_{DFT,w'}(k)$ of the same signal frame $x[n]$ is evaluated, but instead using the analysis window $w[n]$ its analytical derivative $w'[n]$ is applied.

For a given partial in the input signal we get a maximum in the magnitude spectrum $|X_{DFT,w}(k)|$ at $k = k_0$. The frequency reassignment f_p for this partial is then given by

$$f_p = \frac{k_0 f_s}{N} - \frac{f_s}{2\pi} \Im \left\{ \frac{X_{DFT,w'}(k_0)}{X_{DFT,w}(k_0)} \right\} \quad (4.32)$$

The theoretical background which leads to this expression is explained in detail in [5] and not repeated here.

This reassigned frequency is directly used as estimation for the partial. The amplitude is again estimated using the shape of the window main lobe, in the same way as shown in equation 4.31 on the previous page.

4.2.3. Phase Vocoder

Similar to the spectral reassignment method, the phase vocoder approach uses the phases of the complex DFT coefficients to improve the accuracy in frequency and amplitude estimation. Puckette and Brown give in [37] a thorough presentation on this. For implementation details please refer to [11, p. 337].

The basic idea is, that the fundamental frequency ω_p of a partial is mathematically given by the derivative of the instantaneous phase, i.e. $\omega_p = \frac{d\phi(n)}{dn}$. We approximate this derivative by the phase difference between two DFTs separated by the hop size H . To get a frequency estimation, this phase difference is evaluated at the local maximum of the first DFT and divided by the used hop size. So $\omega_p \approx \frac{\Delta\phi(k_0)}{H}$.

Figure 4.5 on the following page shows the basic framework: Similar to the STFT the input signal is framed with an analysis window of length N , starting at every H^{th} sample, with H being the hop size. For two consecutive windows (only shifted by the hop size) the corresponding DFTs $X_{DFT,1}(k)$ and $X_{DFT,2}(k)$ are computed. At the local maximum $k = k_0$ of $X_{DFT,1}(k)$ the phases of $X_{DFT,1}(k_0)$ and $X_{DFT,2}(k_0)$ are computed, unwrapped, and their difference divided by the hop size H . This results in the frequency estimation ω_p .

The amplitude is again estimated using the shape of the window main lobe, in the same way as shown in equation 4.31 on the previous page.

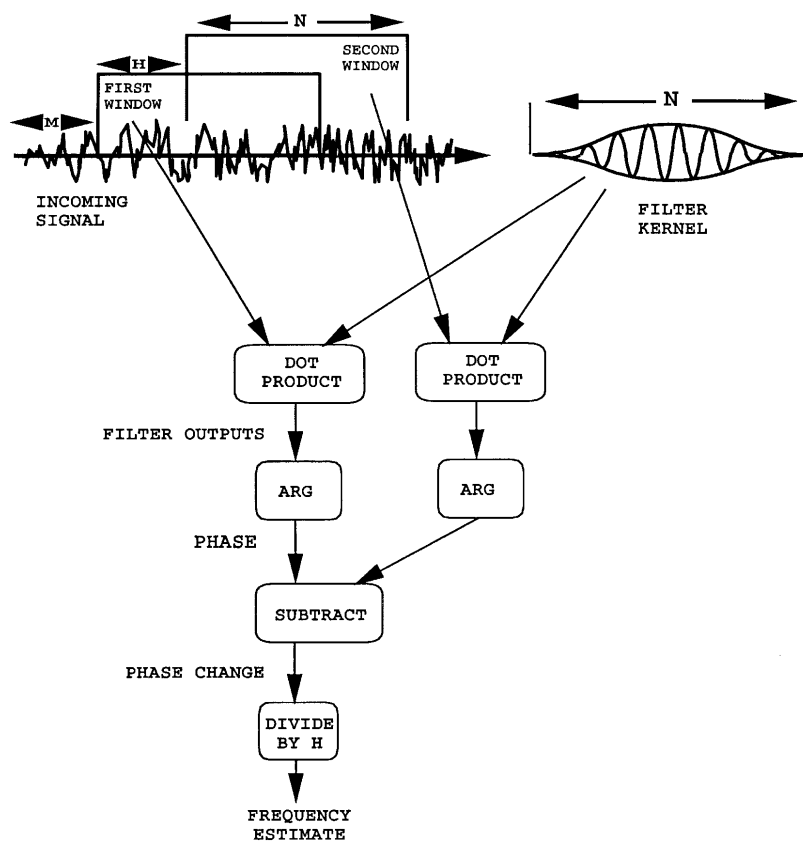


Figure 4.5.: Block diagram of the phase vocoder used for frequency estimation. [37]

Chapter 5.

Frequency Estimators based on Shifted Fourier Transforms

We will introduce in this chapter two estimation methods based on shifted Fourier transforms. They promise good results in terms of numerical exactness.

5.1. Motivation

In order to understand their usefulness and advantages, we take a look at the motivation behind these methods.

5.1.1. MDCT in Perceptual Audio Coding

The MDCT in combination with the sinusoidal window (square root of a shifted Hann window) is widely used in audio coding. It can be computed efficiently (e. g. [27]) and the energy compaction is known to be good (although the differences are marginal [47]). Furthermore it showed up that it serves as a very suitable filter bank for perceptual audio coders, applied e. g. in the MPEG audio coders MPEG-1 Layer 3 and MPEG-2 AAC. For detailed information on perceptual audio coders please refer to the excellent review and tutorial by Painter and Spanias [31].

Figure 5.1 on the following page shows a simplified block diagram of the MPEG-2 AAC coder, which is common to most MDCT based audio coders. A frequency domain representation of the audio signal is computed with the MDCT, usually using a sine window or a Kaiser-Bessel derived window (KBD window). By requantisation a quantiser reduces the number of bits used for coding every MDCT coefficient to a value, which is provided by the psychoacoustic model (quantisation noise shaping). The psychoacoustic model is mostly based on temporal and frequency masking and requires therefore the computation of an FFT to compute the masking thresholds and the maximum allowed quantization noise level. In the MPEG standard a Hann window is used in combination with the FFT.

For computational simplicity it would be desirable, to replace the two different frequency transforms by one transformation which is then common to the audio coding *and* the psychoacoustic model. We showed in chapter 3.5 on page 23 that the shifted DFTs and the MDCT are closely related to each other. We can express the MDCT by taking the real part of the complex ODFT coefficients (shifted by a constant phase term) and therefore replace in the audio coder model the DFT and the MDCT by the ODFT [14]. Figure 5.2 on the following page illustrates this.

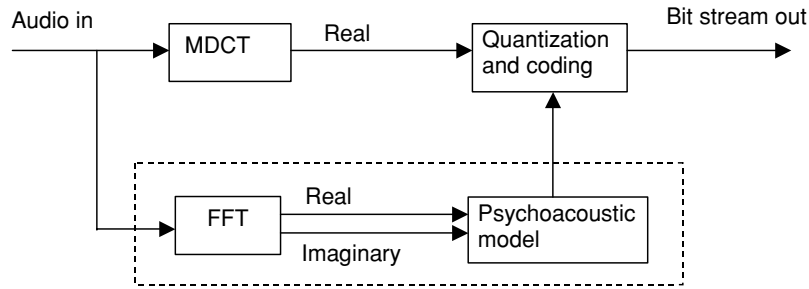


Figure 5.1.: Block diagram of a MDCT based audio coder (e. g. MPEG-2 AAC). [46]

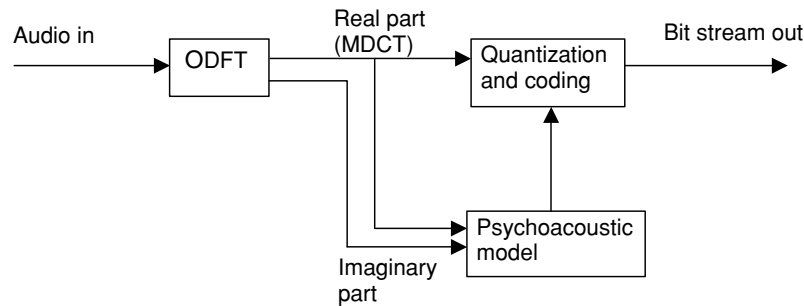


Figure 5.2.: Modified structure for a MDCT based audio coder using ODFT. [46]

In [46] instead of the ODFT another shifted DFT is used, which has a frequency shift $b = v = 0.5$ and a time shift of $a = u = \frac{N+1}{2}$. So the constant phase shift of the coefficients to get the MDCT values is not necessary there anymore.

Nevertheless we focus on the ODFT representation ($a = 0, b = \frac{1}{2}$), as the resulting spectrum is much closer to the DFT spectrum and the obtained MDCT coefficients are identical anyway.

As a consequence for both the signal coding part and the psychoacoustic model the same sine window is used then. The use of a e. g. Hann window is impossible, as perfect reconstruction in the signal coding part would not be maintained anymore. The impact of the use of the sine window instead of the Hann window in the psychoacoustic model is examined in detail in [21]. The result is, that the use of the sine or the KBD window does not make any difference to the final masking threshold compared to when the Hann window is used.

In [23] a similar approach for the MP3-encoder MPEG1 Layer3 is proposed. It is based on the same idea, uses the complex CMDCT and considers the additional filter bank that is present in the analysis stage.

From this perspective it is not astonishing anymore that we want to carry out parameter estimation based on the ODFT.

Daudet and Sandler presented a “pure MDCT” coder which approximates the FFT by the regularized S-spectrum [7]. They estimate the parameters directly from the MDCT coefficients with the method described below in section 5.4.

5.1.2. MDCT Domain Processing

Compressed domain processing is in [22, ch.6] defined as “... the ability to perform modifications on the audio while in its compressed form.” Its importance is showed in detail there. The goal is to do modifications on the sound or even only feature-extraction without having to perform a time-consuming decoding and furthermore quality-degrading re-encoding of the audio data. As shown above, many common audio coder outputs basically consist of re-quantised MDCT coefficients. For a large number of processing applications the first step is the extraction of parameters of the sound’s frequency content.

So a method to perform parameter estimation directly from the MDCT coefficients is desirable. The only algorithm we found that is capable of doing this is the MDCT estimator presented in section 5.4. In its original version, as published in [26], it is capable of detecting only the most dominant spectral peak per frame. We will generalize the method to detect as much partials of the input sound as possible.

5.2. Frequency Response of the Sinusoidal Analysis Window

In the same way as for most of the methods shown in chapter 4, for simplicity we restrict the theoretical description of the estimators on a single sinusoid which is stationary within the analysis frame. This is by far the biggest limitation. We so assume that the partials are only slowly changing over time. Furthermore they must be widely spaced in frequency, so that their spectral components do not interfere significantly. Nevertheless, as shown in chapter 6.5, the estimators still perform quite well when analysing musical sounds which do not fully hold the given constraints anymore.

So we assume a discrete-time sinusoidal input $x[n]$ as defined in equation 2.2 on page 11,

$$x[n] = A \sin \left[\frac{2\pi}{N} (l_0 + \Delta l)n + \Phi \right] \quad (5.1)$$

with $l = l_0 + \Delta l$

where l is the normalized frequency (or frequency in DFT bins), l_0 its integer part and Δl its fractional part (so $0 \leq \Delta l < 1$).

This input is windowed by a sinusoidal window $h[n]$ as defined in equation 3.23 on page 21.

$$h[n] = \begin{cases} \sin \left(\frac{\pi}{N} (n + \frac{1}{2}) \right) & \text{if } 0 \leq n \leq N - 1 \\ 0 & \text{else} \end{cases} \quad (5.2)$$

This sinusoidal window is as well a plain sinusoid as defined in equation 5.1, with

$$A = 1 \quad l = \frac{1}{2} \quad \omega_0 = \frac{\pi}{N} \quad \Phi = \frac{2\pi}{4N} \quad (5.3)$$

and which is multiplied by a rectangular window $r[n]$ of length N .

$$r[n] = \begin{cases} 1 & \text{if } 0 \leq n \leq N - 1 \\ 0 & \text{else.} \end{cases} \quad (5.4)$$

In the following steps we derive analytically the discrete-time Fourier transform (DTFT) of the analysis window $h[n]$. (We thereby omit the subscript ‘DTFT’ for the frequency representations $R_{DTFT}(\omega)$, $H_{DTFT}(\omega)$ and $X_{DTFT}(\omega)$.)

With $z \in \mathbb{Z}$ being an arbitrary complex number, the so-called geometric series and its closed form are given by the following formula [29]:

$$\sum_{n=N_1}^{N_2} z^n = \frac{z^{N_1} - z^{N_2+1}}{1 - z} \quad \text{with } N_1 \leq N_2 \quad (5.5)$$

which gives

$$\sum_{n=0}^{N-1} z^n = \frac{1 - z^N}{1 - z} \quad (5.6)$$

With the aid of that the DTFT $R(\omega)$ of the rectangular window $r[n]$ calculates to:

$$\begin{aligned} R(\omega) &= \sum_{n=-\infty}^{\infty} r[n]e^{-j\omega n} = \sum_{n=0}^{N-1} 1(e^{-j\omega})^n \\ &= \frac{1 - e^{-j\omega N}}{1 - e^{-j\omega}} = \frac{e^{-j\omega N/2}(e^{j\omega N/2} - e^{-j\omega N/2})\frac{1}{2j}}{e^{-j\omega/2}(e^{j\omega/2} - e^{-j\omega/2})\frac{1}{2j}} \\ &= e^{-j\omega/2(N-1)} \frac{\sin(\omega N/2)}{\sin(\omega/2)} \end{aligned} \quad (5.7)$$

With $\omega = \omega_0$ and with $\delta(l)$ being the Kronecker delta defined by

$$\delta(l) = \begin{cases} 1 & \text{for } l = 0 \\ 0 & \text{else} \end{cases} \quad (5.8)$$

one period of the DTFT of the sinusoid $x[n]$ is given by:

$$X(\omega) = \frac{A}{2j} [e^{j\Phi}\delta(\omega - \omega_0) - e^{-j\Phi}\delta(\omega + \omega_0)] \quad (5.9)$$

The longish proof for that is shown in [30] and not repeated here.

The window $h[n]$ is a sinusoid multiplied by a rectangular window $r[n]$. As a consequence, its DTFT $H(\omega)$ can be easily computed by convolution of $X(\omega)$ and $R(\omega)$:

$$\begin{aligned} H(\omega) &= X(\omega) * R(\omega) \\ &= \frac{A}{2j} [e^{j\Phi}R(\omega - \omega_0) - e^{-j\Phi}R(\omega + \omega_0)] \end{aligned} \quad (5.10)$$

With the values given in equation 5.3 on the preceding page and equation 5.7 $H(\omega)$ can be

expressed as:

$$\begin{aligned}
 H(\omega) &= \frac{1}{2} e^{j\frac{3\pi}{2}} \left[\begin{array}{l} e^{j\frac{2\pi}{4N}} e^{-j\frac{1}{2}(\omega - \frac{\pi}{N})(N-1)} \frac{\sin[(\omega - \frac{\pi}{N})\frac{N}{2}]}{\sin[(\omega - \frac{\pi}{N})\frac{1}{2}]} \\ -e^{-j\frac{2\pi}{4N}} e^{-j\frac{1}{2}(\omega + \frac{\pi}{N})(N-1)} \frac{\sin[(\omega + \frac{\pi}{N})\frac{N}{2}]}{\sin[(\omega + \frac{\pi}{N})\frac{1}{2}]} \end{array} \right] \\
 &= \frac{1}{2} \left[\begin{array}{l} e^{-j\frac{\omega}{2}(N-1)} \frac{\cos(\frac{\omega N}{2})}{\sin(\frac{\pi}{2N} - \frac{\omega}{2})} \\ + e^{-j\frac{\omega}{2}(N-1)} \frac{\cos(\frac{\omega N}{2})}{\sin(\frac{\pi}{2N} + \frac{\omega}{2})} \end{array} \right]
 \end{aligned}$$

This gives as final result:

$$H(\omega) = \frac{1}{2} e^{-j\frac{\omega}{2}(N-1)} \cos\left(\frac{\omega N}{2}\right) \left[\frac{1}{\sin(\frac{\pi}{2N} + \frac{\omega}{2})} + \frac{1}{\sin(\frac{\pi}{2N} - \frac{\omega}{2})} \right] \quad (5.11)$$

As $H(\omega)|_{\omega=0} = \frac{1}{\sin(\frac{\pi}{2N})}$, the normalized frequency response $\hat{H}(\omega)$ evaluates to

$$\hat{H}(\omega) = \frac{1}{2} e^{-j\frac{\omega}{2}(N-1)} \sin\left(\frac{\pi}{2N}\right) \cos\left(\frac{\omega N}{2}\right) \left[\frac{1}{\sin(\frac{\pi}{2N} + \frac{\omega}{2})} + \frac{1}{\sin(\frac{\pi}{2N} - \frac{\omega}{2})} \right] \quad (5.12)$$

One can see in equation 5.11 already, that the frequency response is a sum of two components H_1 and H_2 differing only in the $\frac{1}{\sin}$ term within the brackets. Figure 5.3 on the following page illustrates this. Its frequency axes are given in DFT-bins ($\omega = k \cdot \frac{2\pi}{N}$).

Subplots (a) and (b) (real and imaginary part) show, that the two components $H_1(\omega)$ and $H_2(\omega)$ are identical, except that one of them is mirrored around the y-axis. They add up to $H(\omega)$. In the magnitude plot (subplot (c)) it is not visible anymore that the frequency response of the window actually consists of two overlying spectral components. It has zeros at

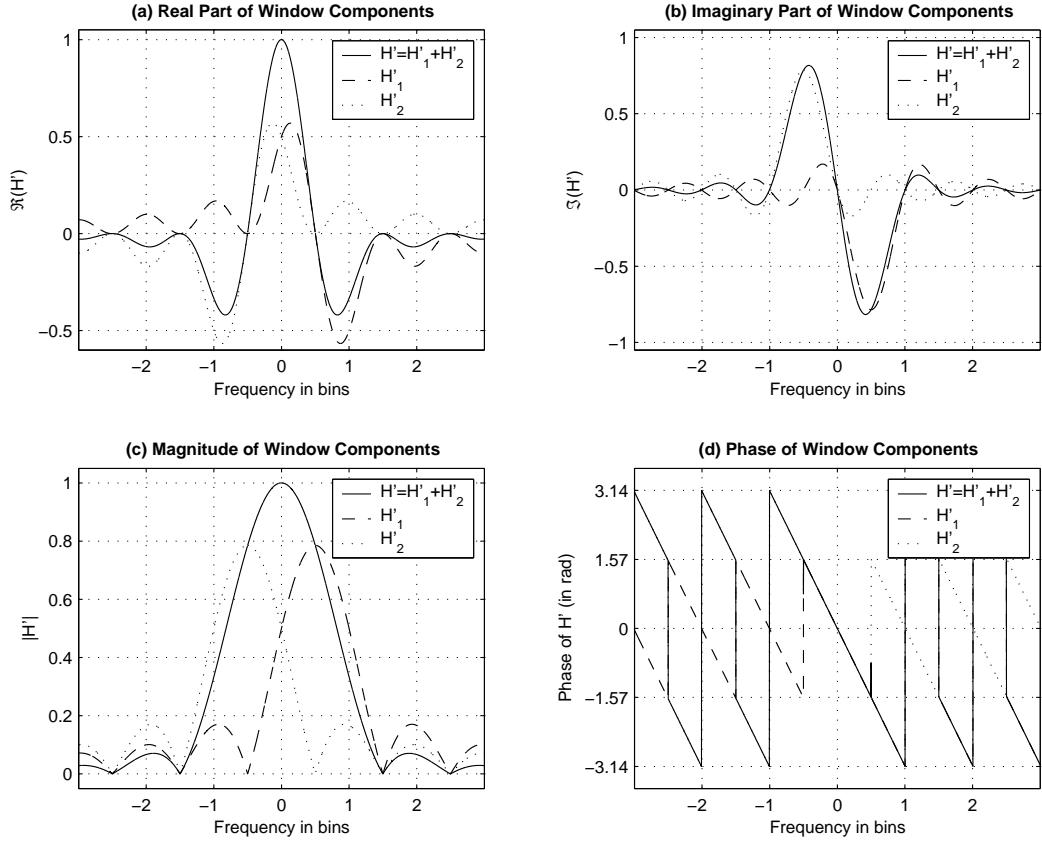
$$\omega = \pm \left(\frac{\pi}{N} + k \frac{2\pi}{N} \right) \quad \text{with } k \text{ integer, } k \geq 1. \quad (5.13)$$

The width of the main lobe is $3\frac{2\pi}{N}$ or three frequency bins. As a consequence, in a corresponding DFT spectrum with sinusoidal input no more than three bins are excited by the main lobe.

The phase plot in sub-figure (d) demonstrates the linear phase property of this special window. The linear phase can be shown as well directly with the definition $h[n]$ of the window, as it fulfils the condition $h[n] = h[N-1-n]$. Please notice that the phase jumps of π are positioned at the zeros of $H(\omega)$.

The width of the side lobes is one bin, the highest side lobe is at -23 dB and an approximation for the side lobe falloff for high frequencies is given in [13].

Furthermore subplot (c) gives some insight on how the input signal's power is distributed after a transformation on the different spectral peaks. It is twice split up to finally four different locations: on the one hand to the positive and to the negative mirror frequencies,


 Figure 5.3.: Frequency response of the normalized sinusoidal window $\hat{H}(\omega)$.

and on the other hand to the two spectral components of the window frequency response as shown in the plot. Including the factor N which will be introduced by the discretisation of the Fourier transform, in our equations the amplitude factor will show up with $\frac{|NA|}{4}$.

We therefore want to normalize the frequency response in a way, so that the maximum of $|\hat{H}_1(\omega)|$ at $\omega = \frac{\pi}{N}$ (which equals $|\hat{H}(\frac{\pi}{N})|$) equals one. $|\hat{H}_1(\omega)|$ has a zero at this point which is lifted by a pole at the same location. We therefore compute the limit by:

$$\begin{aligned}
 \lim_{\omega \rightarrow \frac{\pi}{N}} |\hat{H}_1(\omega)| &= \lim_{\omega \rightarrow \frac{\pi}{N}} \left[\frac{1}{2} \sin\left(\frac{\pi}{2N}\right) \frac{\cos\left(\frac{\omega N}{2}\right)}{\sin\left(\frac{\pi}{2N} - \frac{\omega}{2}\right)} \right] \\
 &= \lim_{z \rightarrow 0} \left[\frac{1}{2} \sin\left(\frac{\pi}{2N}\right) \frac{\sin\left(\frac{z}{2}\right)}{\sin\left(\frac{z}{2N}\right)} \right] && \text{with } z := \omega N - \pi \quad (5.14) \\
 &= \frac{N}{2} \sin\left(\frac{\pi}{2N}\right) && \text{with l'Hospital's rule} \\
 &\approx \frac{\pi}{4} && \text{for } N \gg
 \end{aligned}$$

As a consequence we rescale our normalized frequency response $|\hat{H}(\omega)|$ by the multiplicative factor $\frac{2}{N \sin(\frac{\pi}{2N})} \approx \frac{4}{\pi}$.

5.3. ODFT Analysis of a Sinusoidal Input

We stepwise derive this method based on the ODFT, which was recently published by Aníbal J. S. Ferreira in [14].

We assume the input signal $x[n]$ as given in equation 5.1 on page 38 with a given frequency l . The task is, to determine l_0 and Δl from the signal's ODFT spectrum.

Because in time domain $x[n]$ is windowed by $h[n]$, this means that their DTFT spectra are convoluted and the resulting spectrum $Y(\omega)$ is basically $H(\omega)$ frequency shifted (modulated) to the frequency l of the input sinusoid. In the strict sense, $H(\omega)$ is as well modulated to the negative mirror frequency $-l$, which gives additional small components in the positive frequency range.

As an approximation we neglect these aliasing terms. Usually not being a big restriction, this increases the estimation error for very low and very high frequencies. By doing this we do not take into account components that are caused by the second term of equation 5.9 on page 39 and approximate it by:

$$X(\omega) \approx \frac{A}{2j} e^{j\Phi} \delta(\omega - \omega_0) \quad \text{for } 0 \leq \omega \leq \pi. \quad (5.15)$$

As shown in chapter 3.2, the ODFT $Y_{ODFT}(k)$ is then obtained by taking samples of $Y(\omega)$ at the frequency points $\omega_k = \frac{2\pi}{N}(k + \frac{1}{2})$.

Figure 5.4 shows a typical situation of the ODFT coefficients around the input frequency l . In this case $l = 100.3$, so $l_0 = 100$ and $\Delta l = 0.3$. With the used window length $N=1024$ this corresponds to a normalized frequency of $\omega_0 = \frac{2\pi l}{N} = 0.61543$.

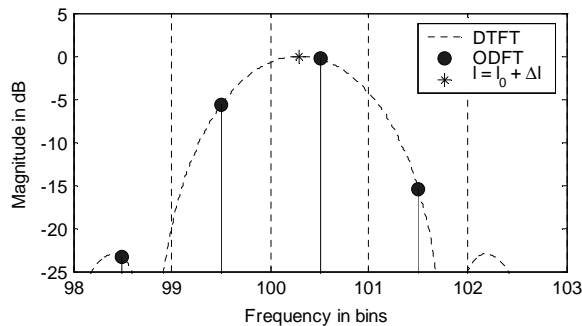


Figure 5.4.: ODFT coefficients of a windowed sinusoid with frequency $l = 100.3$.

Figure 5.5 on the following page gives a qualitative overview on possible magnitude proportions of the three ODFT bins that are excited by a sinusoid. Their relation is dependent on Δl only. $Y_{ODFT}(l_0)$ is a local maximum in the cases (b), (c) and (d). Therefore, l_0 can be easily determined by scanning $Y_{ODFT}(k)$ for its local maximum. In case (a), where only two bins at all show up and have the same magnitude, the higher number coefficient on the right gives the correct value.

Furthermore, obviously the relative magnitude of the ODFT coefficients $l_0 - 1$ and $l_0 + 1$ give information about Δl . We will exploit this fact to determine the input frequency of the sinusoid. As you can see on the figures, all coefficients are determined by the window main lobe, which we will fit into the local maximum l_0 and its two neighbours.

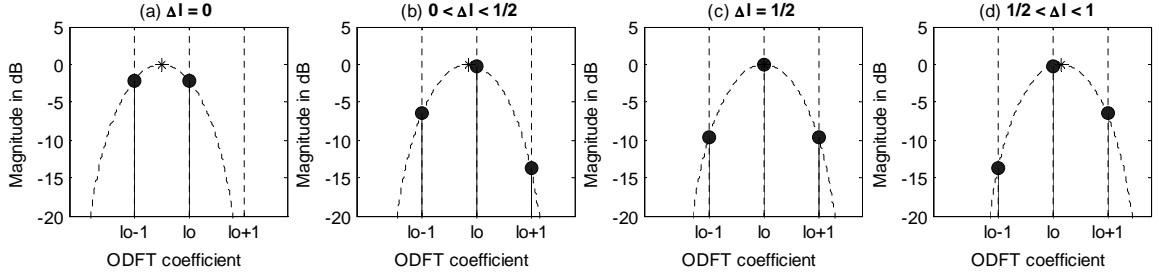


Figure 5.5.: Proportions of the ODFT magnitude coefficients $l_0 - 1$, l_0 and $l_0 + 1$ with an input sinusoid of frequency $\omega = \frac{2\pi}{N}(l_0 + \Delta l)$ (denoted by ‘*’).

Extraction of the Frequency

Let us define an amplitude scaling factor κ corresponding to equation 5.14 on page 41:

$$\kappa = \frac{|NA|}{4} \frac{2}{N \sin(\frac{\pi}{2N})} \quad (5.16)$$

From figure 5.4 on the preceding page we can easily derive the connection between the magnitudes of our three ODFT coefficients $Y_{ODFT}(k)$ and the window frequency response $\hat{H}(\omega)$:

$$|Y_{ODFT}(l_0 - 1)| = \kappa |\hat{H}(\frac{2\pi}{N}(0 - \Delta l - \frac{1}{2}))| = \kappa |\hat{H}(\frac{2\pi}{N}(\Delta l + \frac{1}{2}))| \quad (5.17)$$

$$|Y_{ODFT}(l_0)| = \kappa |\hat{H}(\frac{2\pi}{N}(0 - \Delta l + \frac{1}{2}))| = \kappa |\hat{H}(\frac{2\pi}{N}(\Delta l - \frac{1}{2}))| \quad (5.18)$$

$$|Y_{ODFT}(l_0 + 1)| = \kappa |\hat{H}(\frac{2\pi}{N}(0 - \Delta l + \frac{3}{2}))| = \kappa |\hat{H}(\frac{2\pi}{N}(\Delta l - \frac{3}{2}))| \quad (5.19)$$

The ratio of two of these signals is given by:

$$\frac{|Y_{ODFT}(l_0 - 1)|}{|Y_{ODFT}(l_0 + 1)|} = \frac{|\hat{H}(\frac{2\pi}{N}(\Delta l + \frac{1}{2}))|}{|\hat{H}(\frac{2\pi}{N}(\Delta l - \frac{3}{2}))|} \quad (5.20)$$

The left side of this equation consists of measured values from the ODFT spectrum, the right hand side is an analytical expression, in which the only unknown is Δl . As a consequence, it is – at least in theory – possible to extract Δl from this equation.

Unfortunately $|\hat{H}(\omega)|$ is an analytically complex expression, so that we cannot easily extract Δl from that. Because of this, we introduce a more simple function which models the main lobe of $|\hat{H}(\omega)|$. This is of course an approximation only. Therefore this method is somehow similar to e.g. the parabolic interpolation method, which models the window main lobe by a parabola.

As approximation we use the following function, which is traceable by the parameter G :

$$|\hat{H}(\omega)| \approx \left[\cos\left(\frac{N\omega}{6}\right) \right]^G \quad \text{with } |\omega| < \frac{3\pi}{N} \quad (5.21)$$

Parameter G is an arbitrary real constant. In [14] an optimization process on G was carried out. Its best value to minimize the maximum absolute frequency estimation error was found and given with $G = 1.37$.

We can now exchange $|\hat{H}(\omega)|$ in equation 5.20 on the previous page by the approximated $|\hat{H}(\omega)|$ and easily resolve this new equation for Δl , which ends up in:

$$\Delta l \approx \frac{3}{\pi} \arctan \left[\frac{\sqrt{3}}{2 \sqrt{\frac{|Y_{ODFT}(l_0-1)|}{|Y_{ODFT}(l_0+1)|} + 1}} \right] \quad (5.22)$$

This is the final estimation of the input frequency l , l_0 given through the local maximum and Δl approximated by equation 5.22.

Extraction of the Phase

The phase estimation is based on the value of Δl , so its estimation error mainly depends on an accurate frequency estimation.

The phase of a coefficient $Y_{ODFT}(l_0)$ is determined by three components:

1. Phase shift $e^{j\Phi}$ given by the initial phase Φ (equation 5.9 on page 39).
2. Phase shift $e^{-j\frac{\pi}{2}}$ given by the input transform $X(\omega)$ (equation 5.9).
3. Phase shift $e^{-j\frac{\omega}{2}(N-1)}$ given by the window frequency response $H(\omega)$ in equation 5.11, with $\omega = \frac{2\pi}{N}(\frac{1}{2} - \Delta l)$, which is derived from figure 5.4 again.

So the phase for the l_0 -th ODFT coefficient evaluates to:

$$\begin{aligned} \angle(Y_{ODFT}(l_0)) &= \Phi - \frac{\pi}{2} - \frac{\pi}{N}(N-1)(\frac{1}{2} - \Delta l) \\ &= \Phi + \Delta l \pi (1 - \frac{1}{N}) - \pi(1 - \frac{1}{2N}) \end{aligned} \quad (5.23)$$

Resolved for Φ and with Δl from equation 5.22 we get the estimation for the initial phase of the input sinusoid:

$$\Phi = \angle(Y_{ODFT}(l_0)) - \Delta l \pi (1 - \frac{1}{N}) + \pi(1 - \frac{1}{2N}) \quad (5.24)$$

As mentioned before, the estimation error of Δl propagates linearly into the phase estimation.

Extraction of the Amplitude

In the same way as the phase, the amplitude estimation is based on an accurate extraction of Δl . Furthermore, it is based on a main lobe model similar to equation 5.21 on the previous page, with the slight difference that only a smaller part of the main lobe is modelled.

$$|\hat{H}(\omega)| \approx \left[\cos \left(\frac{N\omega}{6} \right) \right]^F \quad \text{with } |\omega| \leq \frac{\pi}{N} \quad (5.25)$$

Again $\hat{H}(\omega)$ was optimized in [14] in the parameter F and the smallest relative magnitude estimation error showed up for $F = 1.65$. This value is different from G (as seen above), as not the whole main lobe has to be covered by the model.

To get the amplitude, we use the same relation as in equation 5.18. The rescaling is done similar to equation 5.14 on page 41, but using the approximation model to get the rescaling factor:

$$|\hat{H}(\frac{\pi}{N})| \approx \left[\cos \left(\frac{\pi}{6} \right) \right]^F = \left[\frac{\sqrt{3}}{2} \right]^F \quad (5.26)$$

Putting all this together results in:

$$\begin{aligned} |Y_{ODFT}(l_0)| &\approx \frac{|NA|}{4} \left[\frac{2}{\sqrt{3}} \right]^F |\hat{H}(\frac{2\pi}{N}(\Delta l - \frac{1}{2}))| \\ &\approx \frac{|NA|}{4} \left[\frac{2}{\sqrt{3}} \cos(\frac{\pi}{6}(2\Delta l - 1)) \right]^F \end{aligned} \quad (5.27)$$

So finally the estimation for the amplitude is given by resolving this equation for A :

$$A \approx \frac{4|Y_{ODFT}(l_0)|}{N} \left[\frac{\sqrt{3}}{2 \cos(\frac{\pi}{6}(2\Delta l - 1))} \right]^F \quad (5.28)$$

Extension to Sounds with Multiple Sinusoidal Components

In the same way as with the conventional methods described before, the ODFT magnitude spectrum is searched for local maxima. These are considered as the integer parts l_0 of the partials frequencies. Together with their left and right hand side neighbour bins the parameters for every peak are separately estimated.

5.4. MDCT Analysis of a Sinusoidal Input

In the section before, we got the frequency representation by the ODFT and extracted the signal parameters from that spectrum. If required, the MDCT coefficients are then obtained from equation 3.30 on page 24.

Another approach, recently published by Merdjani and Daudet in [26], estimates the parameters directly from the MDCT coefficients. The following section will now derive this method including the enhancement for multiple peaks.

We first want to get an explicit expression for the MDCT coefficients of a sinusoidal signal which is windowed by the sinusoidal window. For the derivation we start from the MDCT definition $X'_{MDCT}(k)$ as given by Daudet in [26] together with the corresponding window and input signal definitions. They slightly differ from the usually used representation as introduced in chapter 3.4. We will therefore implicitly give the connection between these two variants, by reshaping it to a form as given by equation 3.30 on page 24. Dashes after a variable, like X' , h' , etc., depict the functions or sequences as defined by Daudet in [26].

$$X'_{MDCT}(k) = \sum_{n=-\frac{L}{2}}^{\frac{3L}{2}-1} h'[n] x'[n] \sqrt{\frac{2}{L}} \cos \left[\frac{2\pi}{2L} \left(n + \frac{1}{2} \right) \left(k + \frac{1}{2} \right) \right] \quad (5.29)$$

We substitute $2L$ by N and the summation index n by $q = n + \frac{N}{4}$ and get:

$$X'_{MDCT}(k) = \sum_{q=0}^{N-1} h'[q - \frac{N}{4}] x'[q - \frac{N}{4}] \sqrt{\frac{4}{N}} \cos \left[\frac{2\pi}{N} \left(q + \frac{1}{2} - \frac{N}{4} \right) \left(k + \frac{1}{2} \right) \right] \quad (5.30)$$

Let us define $h[n]$ and $x[n]$ similar to equation 5.1 and 5.2 on page 38. In contrast to [26] we do not use the sliding window formalism here and therefore set the frame index $p = 0$.

$$h[q] := \begin{cases} h'[q - \frac{N}{4}] = \sin \left[\frac{\pi}{N} \left(q + \frac{1}{2} \right) \right] & \text{if } 0 \leq q \leq N - 1 \\ 0 & \text{else} \end{cases} \quad (5.31)$$

$$x[q] := x'[q - \frac{N}{4}] = A \sin \left[\frac{2\pi}{N} f(q - \frac{N}{4}) + \Phi' \right] \quad (5.32)$$

With substituting $l := f$ and setting the initial phase

$$\Phi := \Phi' - l \frac{\pi}{2} \quad (5.33)$$

we get:

$$x[q] = A \sin \left[\frac{2\pi}{N} l q + \Phi \right] \quad (5.34)$$

With $\cos(-x) = \cos(x)$ we can now rewrite equation 5.30 by:

$$\begin{aligned} X'_{MDCT}(k) &= \sum_{q=0}^{N-1} \sqrt{\frac{4}{N}} h[q] x[q] \cos \left[-\frac{2\pi}{N} \left(q + \frac{1}{2} - \frac{N}{4} \right) \left(k + \frac{1}{2} \right) \right] \\ &= \Re \left\{ \sum_{q=0}^{N-1} \sqrt{\frac{4}{N}} h[q] x[q] e^{-j \frac{2\pi}{N} (q + \frac{1}{2} - \frac{N}{4}) (k + \frac{1}{2})} \right\} \\ &= \Re \left\{ e^{-j \frac{2\pi}{N} (\frac{1}{2} - \frac{N}{4}) (k + \frac{1}{2})} \sqrt{\frac{4}{N}} \sum_{q=0}^{N-1} h[q] x[q] e^{-j \frac{2\pi}{N} (k + \frac{1}{2}) q} \right\} \end{aligned} \quad (5.35)$$

The equation is now of the form as shown in equation 3.30 on page 24 again. Its kernel is an ODFT transform of the input signal $x[n]$ as defined before, just with another initial phase term. Additionally, we have a slightly different constant phase shift Φ_k .

We express equation 5.35 by

$$X'_{MDCT}(k) = \Re \{ z_k \} = \Re \{ r_k \cdot e^{j\theta_k} \} = r_k \cos(\theta_k) \quad (5.36)$$

in order to compute magnitude r_k and phase $e^{j\theta_k}$ of the complex expression z_k in 5.35 separately.

As we have an ODFT transform again, we can completely reuse the knowledge we gained from sections 5.2 and 5.3.

So the following terms contribute to the phase of z_k :

1. Phase shift given by the modified initial phase Φ (equations 5.23 and 5.33):

$$e^{j\Phi} = e^{j\Phi'} e^{-j l \frac{\pi}{2}} \quad (5.37)$$

2. Phase shift given by the input transform $X(\omega)$ (equations 5.15 and 5.23):

$$e^{-j \frac{\pi}{2}} \quad (5.38)$$

3. Phase shift given by the window frequency response $H(\omega - \omega_0)$ with $\omega = \frac{2\pi}{N}(k + \frac{1}{2})$ and $\omega_0 = \frac{2\pi}{N}l$ (equations 5.11, 5.34 and 5.35):

$$e^{-j \frac{\omega - \omega_0}{2} (N-1)} = e^{-j \frac{\pi}{N} (k + \frac{1}{2} - l) (N-1)} \quad (5.39)$$

4. Phase shift given by the ODFT to MDCT transform (equation 5.35):

$$e^{-j\frac{2\pi}{N}(\frac{1}{2}-\frac{N}{4})(k+\frac{1}{2})} \quad (5.40)$$

The longish process of putting all four terms together, evaluating the products and regrouping gives:

$$e^{j\theta_k} = e^{j\{\Phi' - l\frac{\pi}{2} - \frac{\pi}{2} - \frac{\pi}{N}(k+\frac{1}{2}-l)(N-1) - \frac{2\pi}{N}(\frac{1}{2}-\frac{N}{4})(k+\frac{1}{2})\}} \quad (5.41)$$

$$= e^{j\{\Phi' - \frac{3}{4}\pi - \frac{k}{2}\pi + (\frac{N-1}{N} - \frac{1}{2})l\pi\}} \quad (5.42)$$

$$= (-1)^k e^{j\{\Phi' - \frac{3}{4}\pi + \frac{k}{2}\pi + \frac{N-1}{N}l\pi\}} \quad \text{with } e^{j\pi k} = (-1)^k, k \text{ integer} \quad (5.43)$$

We compute r_k , the magnitude of z_k in the same way:

1. Amplitude term given by the input transform $X(\omega)$ (equation 5.15): $\frac{A}{2}$

2. Amplitude term given by the scaling factor in equation 5.35: $\sqrt{\frac{4}{N}}$

3. Amplitude terms given by the window frequency response $H(\omega - \omega_0)$ with $\omega = \frac{2\pi}{N}(k + \frac{1}{2})$ and $\omega_0 = \frac{2\pi}{N}l$ (equations 5.11, 5.34 and 5.35):

$$\begin{aligned} & \frac{1}{2} \cos\left(\pi\left(k + \frac{1}{2} - l\right)\right) \left[\frac{1}{\sin\left(\frac{\pi}{2N} + \frac{\pi}{N}\left(k + \frac{1}{2} - l\right)\right)} + \frac{1}{\sin\left(\frac{\pi}{2N} - \frac{\pi}{N}\left(k + \frac{1}{2} - l\right)\right)} \right] \\ & = \frac{1}{2} (-1)^k \sin(\pi l) \left[\frac{1}{\sin\left(\frac{\pi}{N}(k - l + 1)\right)} + \frac{1}{\sin\left(-\frac{\pi}{N}(k - l)\right)} \right] \end{aligned} \quad (5.44)$$

Putting these three terms together results in the magnitude r_k :

$$r_k = \frac{A}{2\sqrt{N}} (-1)^k \sin(\pi l) \left[\frac{1}{\sin\left(\frac{\pi}{N}(k - l + 1)\right)} + \frac{1}{\sin\left(-\frac{\pi}{N}(k - l)\right)} \right] \quad (5.45)$$

This is an expression similar to the one given by equation 5.11 on page 40. Ferreira introduced as approximation model a simple cosine function as given in equation 5.21 on page 43 in order to overcome the analytical complexity of this equation. We here go another way, which makes the basic difference between these two methods.

As we are looking for an expression for the purpose of peak extraction, we will later on use this equation only in the vicinity of the input frequency l , so $k - l$ is always a relatively small number. As N is large, the arguments of the two sine functions within the square brackets, which are basically of the structure $\frac{k-l}{N}$, are small and we can use the approximation

$$\sin(x) \approx x \quad \text{for } x \ll 1 \quad (5.46)$$

The error that is made by this approximation gets smaller with increasing N .

By doing this, equation 5.45 becomes to:

$$r_k \approx (-1)^k A \frac{\sqrt{N}}{2\pi} \sin(\pi l) \left[\frac{1}{(k - l + 1)} + \frac{1}{l - k} \right] = (-1)^k \frac{\sqrt{N}}{2\pi} \frac{-A \sin(\pi l)}{(l - k)(l - k - 1)} \quad (5.47)$$

We can now put r_k from equation 5.47 and θ_k from equation 5.43 on the previous page back into equation 5.36 on page 46 to get our final explicit expression for the MDCT coefficients $X'_{MDCT}(k)$:

$$X'_{MDCT}(k) \approx -\frac{\sqrt{N}}{2\pi} \frac{A \sin(\pi l)}{(l-k)(l-k-1)} \cos \left[\Phi' - \frac{3}{4}\pi + \frac{k}{2}\pi + \frac{N-1}{N}l\pi \right] \quad (5.48)$$

Although starting from the same prerequisites, Daudet provides in [26] a different result for this equation. In his result an additional addend of $\frac{\pi}{2}$ shows up in the cosine term.

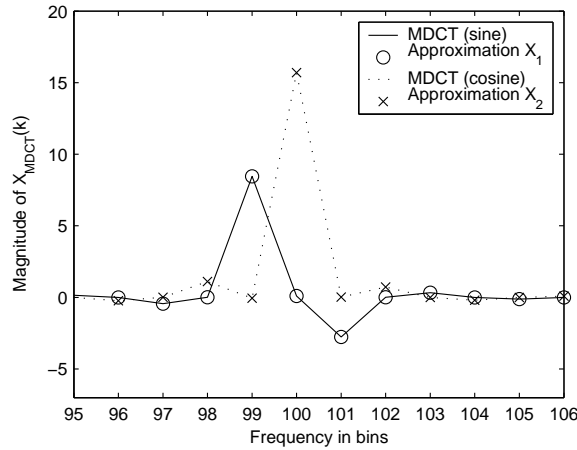


Figure 5.6.: MDCT coefficients of a sine and a cosine signal of frequency $l=100.3$ bins and its approximations ($N=1024$).

Figure 5.6 show a plot of the MDCT transform of a sine function with frequency $l = 100.3$, initial phase $\Phi' = \frac{\pi}{5}$ and amplitude $A = 0.8$ as defined in equations 5.1 and 5.29 and its corresponding approximation X_1 as given by 5.48. Furthermore it plots the approximation formula X_2 which is given in [26]. Although otherwise stated in the paper, this is obviously the MDCT transform of a *cosine* function with the same frequency and initial phase, which is plotted as well. This can also be easily derived from equation 5.48 by rewriting the input cosine by a sine function with an additional initial phase of $\frac{\pi}{2}$:

$$x[n] = A \cos \left[\frac{2\pi}{N}ln + \Phi \right] = A \sin \left[\frac{2\pi}{N}ln + \Phi + \frac{\pi}{2} \right] \quad (5.49)$$

This changes then equation 5.48 to the one published in [26].

Additionally please notice, that for this simple case of single sinusoids the initial phase shift is nothing else than a time shift of the input signal. Nevertheless this time shift has obviously (figure 5.6) a big impact on the resulting MDCT coefficients. Paper [7] gives more details on the non-invariance of the MDCT coefficients in respect to time shifts and its impact on audio coding.

Extraction of the Frequency

Extracting the integer part l_0 of the input frequency l is not as straightforward as before. Figure 5.7 on the following page illustrates this. In all three cases $l_0 = 100$, but nevertheless

coefficient $l_0 - 1$, l_0 or $l_0 + 1$ can be the local maximum, depending on the fractional frequency Δl and the initial phase Φ' .

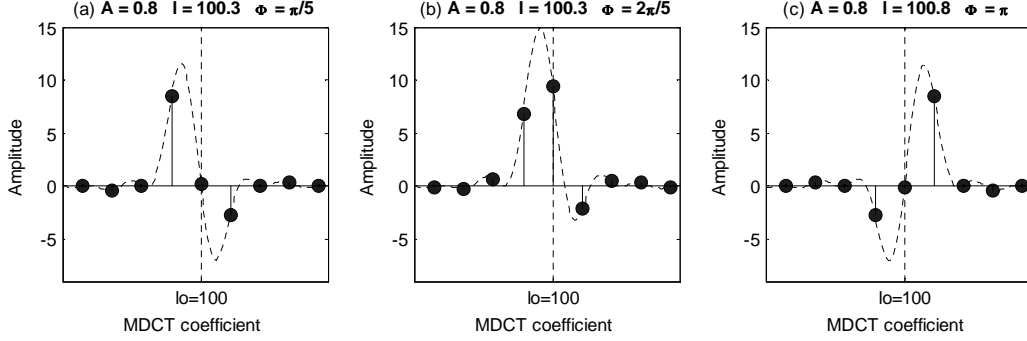


Figure 5.7.: MDCT coefficients of windowed sinusoids with different fractional frequency Δl and initial phase Φ' .

In order to circumvent this problem, we represent the MDCT coefficients in a regularized spectrum [7], which coefficients are defined by the corresponding MDCT coefficient and its direct neighbours:

$$S(k) = \sqrt{(X_{MDCT}(k))^2 + (X_{MDCT}(k+1) - X_{MDCT}(k-1))^2} \quad (5.50)$$

with $0 \leq k \leq \frac{N}{2} - 1$ and $X_{MDCT}(-1) = X_{MDCT}(\frac{N}{2}) = 0$.

If we neglect the aliasing terms from negative frequencies again, this so called S-spectrum has its absolute maximum always at the frequency bin l_0 [6, 26]. As a consequence, we obtain the integer part of the input frequency l_0 by computing the S-spectrum from the MDCT coefficients and scanning the spectrum for the absolute maximum.

In order to get the fractional part Δl , we replace l in equation 5.48 on the previous page by $l_0 + \Delta l$. Then the amplitude ratio α between the two neighbouring peaks is on the one hand given by the actual MDCT coefficients, and can be on the other hand expressed in terms of the approximated explicit expressions. This equation can then be resolved for Δl . So we define α as:

$$\alpha = -\frac{X_{MDCT}(l_0 - 1)}{X_{MDCT}(l_0 + 1)} \quad (5.51)$$

We put equation 5.48 into this equation, and after cancelling and regrouping α becomes to:

$$\alpha = \frac{(\Delta l - 1)(\Delta l - 2)}{(\Delta l + 1)\Delta l} \quad (5.52)$$

As Δl is by definition (eqn. 5.1) $0 \leq \Delta l < 1$, the two terms in the numerator are both always negative (or zero), and the denominator is always positive, so we could adhere that $\alpha \geq 0$ always. Solving this equation for Δl (a second order polynomial) gives exactly one solution which lies in the valid range of Δl :

$$\Delta l = \frac{3 + \alpha - \sqrt{\alpha^2 + 14\alpha + 1}}{2(1 - \alpha)} \quad \text{with } \alpha \neq 1 \quad (5.53)$$

For the pole at $\alpha = 1$, which means that the two coefficients have same amplitude but opposite sign, Δl evaluates to $\Delta l = 0.5$. This equation together with the local maximum of the S-spectrum at l_0 is the frequency estimation.

Nevertheless, as stated in [26], there are in practice cases, where the two coefficients at $l_0 - 1$ and $l_0 + 1$ are both very small and therefore α becomes fault-prone and leads to spurious results. In these cases then the next two neighbouring coefficients $l_0 - 2$ and $l_0 + 2$ are automatically large, and we better consider the ratio β between these two coefficients:

$$\beta = \frac{X_{MDCT}(l_0 - 2)}{X_{MDCT}(l_0 + 2)} \quad (5.54)$$

In the same way as before Δl computes to

$$\Delta l = \frac{5 + 3\beta - \sqrt{\beta^2 + 62\beta + 1}}{2(1 - \beta)} \quad \text{with } \beta \neq 1 \quad (5.55)$$

and $\Delta l = 0.5$ for $\beta = 1$.

The decision whether to use α or β to estimate Δl bases on a ratio λ defined by

$$\lambda = \frac{|X_{MDCT}(l_0)|}{S(l_0)} \quad (5.56)$$

Considering the definition for $S(k)$ in equation 5.50 on the preceding page one can see that $S(l_0) \geq |X_{MDCT}(l_0)|$ always, so $0 \leq \lambda \leq 1$. As more the value of $S(l_0)$ decreases down to its minimum $|X_{MDCT}(l_0)|$ (and so λ gets close to one), as more contribution to $S(l_0)$ is coming from the coefficient l_0 . As a consequence the neighbour coefficients at $l_0 - 1$ and $l_0 + 1$ are small. In [26] a threshold value of $\lambda_0 = 0.9685$ is numerically derived. So if λ is below this value, we use α to estimate Δl , and β otherwise.

Additionally, in practice we cannot assume that $\alpha, \beta \geq 0$ always, as equation 5.48, from which this assumption is derived from, is an approximation only. Therefore we have to restrict the valid range for α and β to get meaningful results, as $\Delta l \in [0, 1[$ only. Equation 5.53 resolves to $\Delta l = 1$ for $\alpha = 0$ and converges monotonously to $\Delta l = 0$ for increasing α . For negative α the result becomes complex (which is not meaningful here). Equation 5.55 resolves to $\Delta l = 1$ for $\beta = \frac{1}{3}$ and decreases monotonously to $\Delta l = 0$ at $\beta = 3$. For negative β the result becomes complex, again.

So if $\alpha \notin [0, \infty[$ use β and if $\beta \notin [1/3, 3]$ use α . If both parameters are not in their valid range, discard this peak, as the resulting out-of-range fractional frequency would result in a very spurious amplitude estimation.

Extraction of the Phase

Extracting the phase Φ of the input sinusoid is not as straight forward as proposed in [26]. The process again uses the ratio between two adjacent coefficients. This is again on the one hand given by the transform coefficients and on the other hand can be expressed by equation 5.48 on page 48. After cancelling down and regrouping we get:

$$\gamma = \frac{X_{MDCT}(l_0 - 1)}{X_{MDCT}(l_0)} = \frac{\Delta l - 1}{\Delta l + 1} \tan \left[\Phi' - \frac{3}{4}\pi + \frac{l_0}{2}\pi + \frac{\frac{N}{2} - 1}{N}(l_0 + \Delta l)\pi \right] \quad (5.57)$$

Please note, that in [26] another expression for this is given, which is probably only a type-setting error. Furthermore the difference in equation 5.48 as described above propagates into this equation. Resolving it for the initial phase Φ' gives:

$$\Phi' = \arctan \left[\gamma \frac{\Delta l + 1}{\Delta l - 1} \right] + \frac{3}{4}\pi - \frac{l_0}{2}\pi - \frac{\frac{N}{2} - 1}{N}(l_0 + \Delta l)\pi \quad (5.58)$$

As the tan and arctan functions are periodic in π instead of 2π , we always only get the initial phase reduced by $k\pi$ so that it is in a range between $-\pi/2$ and $\pi/2$, as:

$$a = \tan(\Phi) \quad \Leftrightarrow \quad \Phi = \arctan(a) + k\pi \quad (k \in \mathbb{Z}) \quad (5.59)$$

The “loss of information” about the phase occurs during the cancelling of terms in the quotient of the two expressions for the MDCT coefficients and the substitution of $\frac{\sin(x)}{\cos(x)} = \tan(x)$. We can determine the quadrant in which the phase angle has to be by looking at the single coefficients itself.

Equation 5.48 becomes for the two used coefficients and using w_1 , w_2 and Φ_t as abbreviations to:

$$X_{MDCT}(l_0) = -w_2 \sin[\pi(l_0 + \Delta l)] \cos[\Phi' + \Phi_t] \quad (5.60)$$

$$X_{MDCT}(l_0 - 1) = -w_1 \sin[\pi(l_0 + \Delta l)] \sin[\Phi' + \Phi_t] \quad (5.61)$$

with

$$w_1 = \frac{\sqrt{N}}{2\pi} \frac{A}{(\Delta l + 1)\Delta l} \geq 0 \quad w_2 = \frac{\sqrt{N}}{2\pi} \frac{A}{(\Delta l - 1)\Delta l} \leq 0 \quad (5.62)$$

$$\Phi_t = \left[-\frac{3}{4}\pi + \frac{l_0}{2}\pi + \frac{\frac{N}{2} - 1}{N}(l_0 + \Delta l)\pi \right] \quad (5.63)$$

Equations 5.60 and 5.61 both define a specific range in which for given parameters the sum of $\Phi' + \Phi_t$ lies. Table 5.1 shows four conditions. Two of them are always fulfilled and therefore the quadrant in which $\Phi' + \Phi_t$ lies is fully determined.

Condition	$\Phi' + \Phi_t$ is in quadrant			
	I	II	III	IV
$\frac{X_{MDCT}(l_0-1)}{-\sin[\pi(l_0+\Delta l)]} \geq 0$	X	X		
$\frac{X_{MDCT}(l_0-1)}{-\sin[\pi(l_0+\Delta l)]} < 0$			X	X
$\frac{X_{MDCT}(l_0)}{-\sin[\pi(l_0+\Delta l)]} \geq 0$		X	X	
$\frac{X_{MDCT}(l_0)}{-\sin[\pi(l_0+\Delta l)]} < 0$	X			X

Table 5.1.: Determination of the sinusoid’s initial phase quadrant in the MDCT domain.

We so compute Φ' with equation 5.58, check if $\Phi' + \Phi_t$ lies in the required quadrant and add or subtract $\pi/2$ or π accordingly.

Extraction of the Amplitude

As we already estimated frequency and phase of the input sinusoid, we could directly extract the amplitude using these two values. But in terms of error propagation we favourably use two coefficients again and only the frequency estimation.

We first give two expressions which are the squares of two adjacent MDCT coefficients multiplied by compatible factors. The motivation for that is to eliminate the cosine term by applying the identity $\sin^2(x) + \cos^2(x) = 1$.

$$[\Delta l(\Delta l - 1)X_{MDCT}(l_0)]^2 = \left(\frac{\sqrt{N}}{2\pi}A \sin(\pi\Delta l)\right)^2 \cos^2 \left[\Phi' - \frac{3}{4}\pi + \frac{l_0}{2}\pi + \frac{\frac{N}{2} - 1}{N}l\pi \right] \quad (5.64)$$

$$[(\Delta l - 1)(\Delta l - 2)X_{MDCT}(l_0 + 1)]^2 = \left(\frac{\sqrt{N}}{2\pi}A \sin(\pi\Delta l)\right)^2 \sin^2 \left[\Phi' - \frac{3}{4}\pi + \frac{l_0}{2}\pi + \frac{\frac{N}{2} - 1}{N}l\pi \right] \quad (5.65)$$

These two equations added together result in:

$$\zeta_r := [(\Delta l - 1)\Delta l X_{MDCT}(l_0)]^2 + [(\Delta l - 1)(\Delta l - 2)X_{MDCT}(l_0 + 1)]^2 = \frac{N}{4\pi^2}A^2 \sin^2(\pi\Delta l) \quad (5.66)$$

Doing the same with the MDCT coefficients $X_{MDCT}(l_0)$ and $X_{MDCT}(l_0 - 1)$ results in:

$$\zeta_l := [(\Delta l - 1)\Delta l X_{MDCT}(l_0)]^2 + [\Delta l(\Delta l + 1)X_{MDCT}(l_0 - 1)]^2 = \frac{N}{4\pi^2}A^2 \sin^2(\pi\Delta l) \quad (5.67)$$

We can compute the left side of the equations (abbreviated by $\zeta_{l/r}$) from the given MDCT coefficients and the estimated frequency. On the right side of the equation the only unknown variables are the amplitudes, which so result in:

$$A_l = \sqrt{\frac{4\pi^2\zeta_l}{N \sin^2(\pi\Delta l)}} \quad A_r = \sqrt{\frac{4\pi^2\zeta_r}{N \sin^2(\pi\Delta l)}} \quad (5.68)$$

In the theoretical case A_l and A_r are identical. In practice this is hardly the case, as the MDCT coefficients are disturbed by negative mirror frequency components, the influence of other partials, and noise. We so use the normalized difference between the two amplitude measures as a criteria to determine if the actual peak is caused by a partial in the input signal or if we should better discard this peak.

No such test is performed in [26]. We use as final amplitude estimation the mean value of the two amplitude estimations $A = \frac{A_l + A_r}{2}$ and discard the peak, if the difference between the two values (normalized by the smaller one of the two values) is bigger than a certain percentage. So A_l and A_r have to hold the condition:

$$\frac{|A_r - A_l|}{\min[A_l, A_r]} \leq a_{dif} \quad (5.69)$$

A maximum ratio of $a_{dif} = 0.05$ gives satisfactory results.

Extension to Sounds with Multiple Sinusoidal Components

In order to detect more than one partial of the signal, a local peak picking has to be performed. The observation made above, that a partial of frequency $l_0 + \Delta l$ induces a maximum in the S-Spectrum at bin l_0 still holds for the case of multiple widely spaced partials and considering l_0 as only local maximum of the spectrum. Unfortunately, the rather complex structure of the MDCT and S-spectra are both by far not monotonously increasing or decreasing between even widely separated peaks. So a simple search for local maxima does not give satisfactory results. Too many local maxima are present in the spectrum which are caused by the transformation and not by an actual partial.

Basically by inspection of spectra of known signals, a peak picking algorithm was developed and is presented in this section.

In a first step, the MDCT spectrum is transformed to the S-Spectrum as described above. It is shown in [7], that the S-Spectrum roughly approximates the DFT spectrum of the input signal. In the same way as with the other parameter estimators, we set all values below a certain threshold to zero. In the simulations we use a threshold of -80 dB in terms of power of the DFT spectrum. Please note that the threshold value has to be multiplied with the additional scaling factor $\sqrt{\frac{4}{N}}$, which originates from the applied MDCT definition (and is not present in the DFT definition).

Beside the overall trend that a partial causes a local maximum, the actual values beside the maximum peak down to low amplitudes show a periodic oscillation around that overall trend. This results in unwanted local maxima at every second frequency bin. In order to circumvent this, we apply a moving average filter of length two on the S-spectrum. This results in an averaged S-spectrum, which we denote as S_a -spectrum.

$$S_a(k) = \frac{S(k) + S(k-1)}{2} \quad (5.70)$$

We scan the S_a -spectrum for local maxima. We detect the locations of the biggest S-Spectrum values in the direct vicinity (plus or minus one bin) of the S_a -maxima. These frequency bins form a first set of candidates for being peaks caused by partials. Still, even in the noiseless case and with a low number of partials, we have much more candidates than partials in the used signal.

In a further step we only keep those candidates k_i (with i denoting the i^{th} candidate) which hold the condition

$$S_a(k_i) > \frac{S_a(k_i-1) + S_a(k_i+1)}{2} \quad (5.71)$$

This translates with equation 5.70 into the equivalent condition:

$$S(k_i) > S(k_i-2) - S(k_i-1) + S(k_i+1) \quad (5.72)$$

Again, there are still more candidates than partials in the input signal. Examining these false candidates reveals, that at least one of the second neighbours of the peak candidates (namely $S(k_i-2)$ or $S(k_i+2)$) should have a by a certain factor a_{thr} smaller amplitude than the candidates S-spectrum value $S(k_i)$. In other words, one of the following conditions has to be fulfilled:

$$S(k_i) > \frac{S(k_i-2)}{a_{thr}} \quad \text{or} \quad S(k_i) > \frac{S(k_i+2)}{a_{thr}} \quad (5.73)$$

In practice, a threshold factor of $a_{thr} = 0.5$ turned out to work well (meaning that the candidate must be at least double as big as the smallest of its second neighbours).

Unfortunately this last step removes in rare cases candidates which are actually induced by a partial. To prevent this, we want to keep those peaks which have a very “symmetrical” shape and have almost equally strong first and second neighbours. We therefore define

$$d_0 = S(k_i) \quad (5.74)$$

$$d_1 = S(k_i-1) - S(k_i+1) \quad (5.75)$$

$$d_2 = S(k_i-2) - S(k_i+2) \quad (5.76)$$

and keep those which fulfil either the condition

$$\frac{|d_1| + |d_2|}{d_0} < a_{sym} \cdot \text{sign}(d_1 d_2) \quad (5.77)$$

or, in the special case that all five values $S(k_i-2)$ to $S(k_i+2)$ have the same sign, the condition:

$$\frac{|d_1| + |d_2|}{d_0} < a_{sym} \cdot \frac{1}{100} \quad (5.78)$$

In practice, setting parameter $a_{sym} = 0.17$ gives good results.

For every peak that is finally remaining after this filtering process, frequency and amplitude are computed as described above.

Chapter 6.

Systematic Performance Comparison

A main scope of this work is to present a concise overview on the performance of the different parameter estimators. Inspired from previous work done by Keiler and Marchand in [17] and in order to verify their results we implemented the parameter estimators shown in chapter 4 and 5. We test them with different test signals to determine their performance.

We restrict the survey on the accuracy of the amplitude and frequency estimation. In the context of analysis for additive synthesis the estimation of the initial phase is usually omitted. We therefore do not include phase estimation in the survey, although it might be beneficial to take this issue up again in a future extension.

Computational costs for the different methods highly depend on their implementation. The estimators were implemented using Matlab by Mathworks. The computation time of a method in Matlab does not provide a sufficient estimate for the numerical complexity of one or the other method or its computational demand for execution in real time on a digital signal processor.

In the following chapters we first take a look at the test setup and then describe the simulation results.

6.1. Analysis Methods

The survey includes all parameter estimation methods mentioned in [17]. The implementations were cross-checked with the source code of Keiler's and Marchand's implementations, that became available at a late state of the project. Additionally the survey includes two parameter estimators based on shifted Fourier transforms, as they represent an essential part of this work.

The following list provides an overview of the used parameters for the different analysis methods. For a detailed description and the meaning of the parameters please refer to the explanations of the estimators in the chapters 4 and 5.

1. Plain FFT with no Further Processing (**PlainFFT**)
 - *Analysis Window*: Periodic Hann Window
 - *FFT magnitude threshold*: -80 dB (Power)
2. Parabolic Interpolation (**ParInt**)
 - *Analysis Window*: IFFT of a Parabola
 - *FFT magnitude threshold*: -80 dB (Power)

- *Window shape error threshold*: 0.3 (multiple of reference value)
3. Triangle Algorithm (**TrglAlg**)
 - *Analysis Window*: IFFT of a Triangle ($S = 2$, Periodic Hann Window)
 - *FFT magnitude threshold*: -80 dB (Power)
 - *Slope bandwidth S* : 2 (frequency bins)
 - *Error energy threshold*: 0.05 (multiple of reference energy)
 - *Phase error threshold*: 0.8 (radians)
 - *Frequency difference threshold*: 1 (frequency bins)
 4. Spectrum Reassignment (**SpcReass**)
 - *Analysis Window*: Periodic Hann Window
 - *FFT magnitude threshold*: -80 dB (Power)
 5. Derivative Algorithm (**DrvAlg**)
 - *Analysis Window*: Periodic Hann Window
 - *FFT magnitude threshold*: -80 dB (Power)
 6. Phase Vocoder (**PhsVoc**)
 - *Analysis Window*: Periodic Hann Window
 - *FFT magnitude threshold*: -80 dB (Power)
 - *Hop size H between two analysis frames*: 1 (frequency bin)
 7. ODFT Analysis (**OdfAlg**)
 - *Analysis Window*: Sinusoidal Window (square root of a shifted Hann window)
 - *FFT magnitude threshold*: -80 dB (Power)
 - *Window main lobe approximation parameter G* : 1.37 (scalar)
 - *Window main lobe approximation parameter F* : 1.65 (scalar)
 8. MDCT Analysis (**MdctAlg**)
 - *Analysis Window*: Sinusoidal Window (square root of a shifted Hann window)
 - *FFT magnitude threshold*: -80 dB (Power)
 - *Coefficient selection threshold λ_0* : 0.9685 (ratio)
 - *Amplitude difference threshold a_{dif}* : 0.05 (ratio)
 - *Second neighbour magnitude threshold a_{thr}* : 0.5 (multiple of peak magnitude)
 - *Symmetry threshold a_{sym}* : 0.17 (ratio)

All methods use the same window length and transform length of $N = 1024$ and a sample rate of $f_s = 44100$ Hz. For the amplitude estimation using the shape of the window main lobe, a lookup table is used instead of the analytical expression. The spectrum of the window is oversampled by a factor of 2^{12} . After the FFT the first 2^{12} values form the lookup table. These values are equivalent to a frequency range of one DFT bin ($0 \leq \Delta l < 1$).

The most important parameter for all the estimators is the choice of the analysis window. The accuracy of several estimators varies highly with different windows. Some methods do not work with every window, some others are even dependent on a special analysis window. In the first sequence of tests every estimator is used in combination with that analysis window with which the estimator performs best. Of course, in a specific application, other windows than the ones optimum for the estimator might be preferable, depending on the input signal and the information that should be extracted. As a consequence we additionally test all methods with the periodic Hann window and the sinusoidal window (square root of a shifted Hann window).

6.2. Test Signals

It is important to choose the right kinds and a wide variety of test signals. On the one hand, they provide some insight on advantages and drawbacks of certain methods. On the other hand, they can show limitations that are common to all estimators and given just by the nature of the signal. Simple signals allow the study of the behaviour of the estimators. Nevertheless this gives only a very limited picture on how the estimator performs for more complex signals. This is especially the case when the assumptions made in the theoretic derivations do not hold any more. The design of the test signals is similar to the ones used in [17] in order to enable comparability and verifiability of results. Nevertheless there are subtle implementation differences which result in slightly different numerics. Additionally in [17] consistently cosine instead of sine functions are used to compute the signals.

6.2.1. Single Stationary Sinusoids

Single Stationary Sinusoids (SiStSi)

The first and simplest signal is the use of pure sine tones which are stationary within the analysis frame. That means the amplitude and the frequency of the tone is constant and does not change within the analysis window.

$$x[n] = A \sin \left[\frac{2\pi}{N} (l_0 + \Delta l)n + \Phi \right] \quad (6.1)$$

with $l = l_0 + \Delta l$

The tests are performed at different frequencies of the input sinusoid. We use an integer frequency bin range from $l_0 = 5$ to $l_0 = 100$ (corresponding to a frequency range from 215 Hz to 4321 Hz). The number of frequencies within one integer frequency bin in this range decreases logarithmically from 100 to 1 and results in a total number of 2090 sinusoids or frames. We so give more weight to lower frequencies, in order to conform with perception and the distribution of energy in most musical sounds. The fractional frequency Δl is chosen

randomly for every frame between 0 and 1 bin with a uniform distribution. For numerical reproducibility of the results it is important to reset the random number generator to a well defined state every time it is used. Figure 6.1 depicts the distribution of the chosen frequencies.

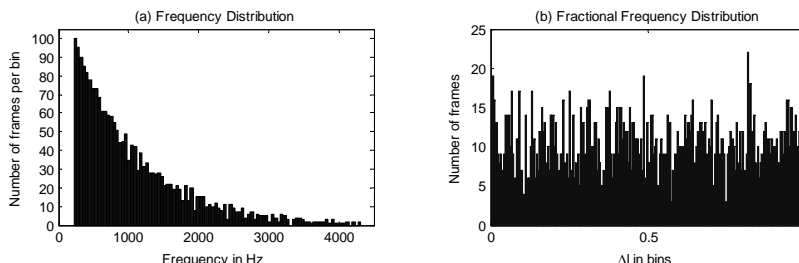


Figure 6.1.: Frequency of occurrence of tones with a certain frequency over (a) the tone's pitch and (b) the fractional frequency of the tone (histograms).

The initial phase Φ of each frame is set in a way so that there are no discontinuities between the end of one frame and the beginning of the next frame. The amplitude is set to $A = 1$ and constant for all frames.

Single Stationary Sinusoids at Higher Frequencies (SiStSiMF)

This test uses the same signals as in the SiStSi case, except that only sinusoids with a frequency at around one fourth of the sampling frequency are used. We can therefore assume, that even for the worst analysis windows used there are practically no more contributions of negative mirror frequencies in the observed spectrum. Beside analytical signals, which are described later, this is the 'cleanest' test signal possible.

Single Stationary Sinusoids with High Noise Level (SiStSiHN)

To the pure stationary sinusoids random noise is added. In general the presence of noise decreases the performance of the parameter estimators. They are differently stable against the influence of noise, and spurious peaks are detected.

We use the same frequency distribution as in the SiStSi case with the same amplitudes and initial phases. We add to this signal Gaussian white noise with a signal-to-noise ratio of $SNR = 12$ dB in terms of power.

The Gaussian white noise signal $x_n[n]$ is created using the Matlab `randn()` function, again with the random number generator reset to a defined state. This results in a set of random numbers with a mean value $\mu = 0$ and a variance and standard deviation of $\sigma^2 = \sigma = 1$ (for an infinite number of samples).

We compute the power P_{sig} of the pure sinusoidal signal and the power P_{noise} of the noise signal separately for every frame. [43, Signal Metrics]

$$P_{sig} = \frac{A^2}{2} \quad P_{noise} = \frac{1}{N} \sum_{n=0}^{N-1} |x_n[n]|^2 \quad (6.2)$$

The required noise power P'_{noise} is determined by the signal-to-noise ratio SNR and computes to

$$P'_{noise} = P_{sig} \cdot 10^{-0.1SNR} \quad (6.3)$$

We therefore rescale the amplitude of our noise signal by the multiplicative factor a_n :

$$a_n = \sqrt{\frac{P'_{noise}}{P_{noise}}} \quad (6.4)$$

Please note, that only for $N \rightarrow \infty$ this value converges to $a_n = A \cdot 10^{-0.05SNR}$. The final noisy test signal $x'[n]$ is so given by:

$$x'[n] = x[n] + a_n x_n[n] \quad (6.5)$$

Single Stationary Sinusoids with Low Noise Level (SiStSiLN)

The SiStSiHN signal presented above is suitable for testing the stability of the estimators. In practical applications one usually encounters a much lower noise level. We therefore use a second test signal set with a signal-to-noise ratio of $SNR = 40$ dB. All other parameters are set as in the SiStSiHN case.

6.2.2. Two Stationary Sinusoids (TwoStSi)

This set of test signals consists of a sum of two sinusoids with different frequencies, but which are both stationary within the analysis frame.

This test reveals the influence of the window side lobes of one sinusoid on the estimation accuracy of the other sinusoid, and how good or bad this is handled by the different estimators. Furthermore, as closer the frequencies of the two sinusoids become as more difficult it is for the estimators to distinct between the two sinusoids. This is closely related to the spectral resolution of the applied transform and the analysis window.

The first of the two sinusoids has constant frequency of 10.3 frequency bins or $f_1 \approx 444$ Hz over all frames. The frequency of the second sinusoid is as well constant within the analysis frame, but the frequency difference between the two sinusoids is different in every frame. Therefore the frequency of the second sinusoid increases stepwise (with a constant increment) from 5 frequency bins ($\cong 215$ Hz) to 16 frequency bins ($\cong 689$ Hz) within the 1099 frames. As the superposition of two sinusoids with exactly the same frequency would give misleading results, we omit for the second sinusoid a frequency range of one percent of a half-tone around the first sinusoid's frequency. That means we skip a frequency interval of approximately 0.5 Hz:

$$f_2 \notin [2 - 2^{(1/1200)}, 2^{(1/1200)}] f_1 \quad (6.6)$$

Figure 6.2 on the following page shows the frequencies of the two sinusoids for all frames. Both sinusoids have a constant amplitude of $A = 0.8$. The initial phase Φ of each sinusoid is set again in a way so that there are no discontinuities in the resulting time signal between the end of one frame and the beginning of the next frame.

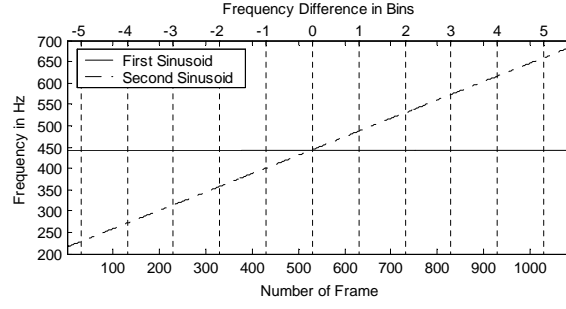


Figure 6.2.: Frequencies of the two stationary sinusoids in different analysis frames.

6.2.3. Single Non-Stationary Sinusoids

This test section contains signals which are single sinusoids but which are not anymore stationary within the analysis frame. Their amplitude or frequency is modulated by another function. All analysis methods assume constant conditions during one frame and are not able to resolve frequency or amplitude changes during one frame. This results in artefacts in the parameter estimations.

In the general case, a continuous-time time varying sinusoid $s(t)$ is given by:

$$s(t) = A(t) \sin(\phi(t)) \quad (6.7)$$

$$\text{with the instantaneous phase } \phi(t) = \Phi + 2\pi \int_{u=0}^t f(u) du \quad (6.8)$$

$$\text{and with the instantaneous frequency } f(t) = \frac{1}{2\pi} \frac{d\phi(t)}{dt}. \quad (6.9)$$

The total length of each signal is set to $\tau_s = 2s$. The continuous-time signals are sampled at the times $t = \frac{n}{f_s}$. A window overlapping is used for analysis. The hop size is 64 samples, which results in an overlap of 960 samples (93.75%) and a total number of 1363 frames. The arithmetic mean values of frequency and amplitude within the frames are considered as the reference values for the evaluation of the parameter estimation.

Linear Frequency Sweep (Sweep)

A linear frequency sweep from $f_0 = 440$ Hz to $f_1 = 880$ Hz is performed. We use a constant amplitude of $A_0 = 0.8$. The continuous-time signal $s(t)$ is then given by:

$$f(t) = f_0 + \frac{f_1 - f_0}{\tau_s} t \quad (6.10)$$

$$s(t) = A_0 \sin \left[2\pi \left(f_0 t + \frac{f_1 - f_0}{2\tau_s} t^2 \right) \right] \quad (6.11)$$

The discrete-time signal $s[n]$ with n integer and $0 \leq n \leq \tau_s f_s - 1$ then evaluates to:

$$s[n] = A_0 \sin \left[\frac{2\pi}{f_s} \left(f_0 + \frac{f_1 - f_0}{2\tau_s f_s} n \right) n \right] \quad (6.12)$$

Tremolo		Vibrato	
$A_0 = 0.8$	amplitude basis	$A_0 = 0.8$	amplitude (const.)
$f_0 = 440$	frequency (const.)	$f_0 = 440$	frequency basis
$A_1 = 0.15$	tremolo depth	$f_1 = 10$	vibrato depth
$f_t = 5$	tremolo frequency	$f_v = 10$	vibrato frequency

Table 6.1.: Tremolo and Vibrato test sound parameters.

Tremolo (Tremolo)

The tremolo is characterized by a constant frequency and an oscillating amplitude. Table 6.1 gives an overview about the parameters used for the tremolo and vibrato sounds. For the continuous-time signal $s(t)$ we therefore get:

$$A(t) = A_0 + A_1 \sin(2\pi f_t t) \quad (6.13)$$

$$s(t) = A(t) \sin(2\pi f_0 t) \quad (6.14)$$

Vibrato (Vibrato)

The vibrato is characterized by a constant amplitude and an oscillating frequency. Again with the parameters given in table 6.1 we get for the continuous-time signal $s(t)$:

$$f(t) = f_0 + f_1 \cos(2\pi f_v t) \quad (6.15)$$

$$s(t) = A_0 \sin \left[2\pi \left(f_0 t + \frac{f_1}{2\pi f_v} \sin(2\pi f_v t) \right) \right] \quad (6.16)$$

Combined Tremolo and Vibrato (TremVib)

The combined tremolo and vibrato is characterized by oscillating amplitude and frequency. Again with the parameters given in table 6.1 we get for the continuous-time signal $s(t)$:

$$A(t) = A_0 + A_1 \sin(2\pi f_t t) \quad (6.17)$$

$$f(t) = f_0 + f_1 \cos(2\pi f_v t) \quad (6.18)$$

$$s(t) = A(t) \sin \left[2\pi \left(f_0 t + \frac{f_1}{2\pi f_v} \sin(2\pi f_v t) \right) \right] \quad (6.19)$$

6.2.4. Musical Instrument Sounds (Guitar, Sax, Voice)

The final test is performed using samples of musical instruments. This provides insight in the ‘real world’ performance of the estimators and proves or disproves their practical usefulness for musical sounds.

We need reliable and accurate reference amplitudes and frequencies of the test sound’s partials. This information is not available for recorded sounds. We therefore use high quality synthesized sounds of musical instruments. Their partials change amplitude and frequency every 64 samples. We so have by far not constant partials within the analysis frame. Beside

their length in frames and seconds table 6.2 provides the average number of partials per frame included in every sound. We use a hop size of 64 samples for analysis, so we have again a window overlap of 960 samples or 93.75%.

The sounds are the same as used in [17]: a short phrase of a synthesized singing male voice, a saxophone tone with strong tremolo and vibrato and a plucked guitar sound with a long sustain phase.

	Guitar	Sax	Voice
Number of Frames	4722	1947	989
Length in seconds	6.87	2.85	1.45
Avg. Number of Partial	11.75	30.46	40.06

Table 6.2.: Parameters of the natural-like test sounds.

6.3. Peak Filtering

In order to get an accurate estimation of frequency and amplitude of a partial, as the first thing the estimator has to detect the partial at all. For the test signals described above with more than one sinusoidal component this is by far not always the case. On the other hand the estimators detect spurious peaks, which are actually not present in the test signal and do not represent a partial in the reference signal. They result from noise or from inter-products of closely spaced partials in the reference signal.

The algorithm has to check, if the detected peaks really correspond to a partial in the reference signal or not. There are peaks that are in the original signal and not in the estimation and vice versa.

For all frequency components (partials) present in the test signal, the algorithm looks for corresponding candidates within the *detected peaks* which fulfil the following prerequisites:

1. The frequency difference between the reference sinusoid and the candidate is smaller than one frequency bin. So with f_r being the frequency of the test signal partial, the frequency f_e of the estimated peak has to be within the interval:

$$f_e \in \left] f_r - \frac{f_s}{N}; f_r + \frac{f_s}{N} \right[\quad (6.20)$$

2. The amplitude difference between the reference sinusoid and the candidate is smaller than 3 dB in terms of power. So with A_r being the amplitude of the test signal partial, the amplitude A_e of the estimated peak has to be within the interval:

$$A_e \in \left] A_r 10^{(-3/20)}; A_r 10^{(3/20)} \right[\quad (6.21)$$

3. From the remaining candidates that fulfil criteria one and two (if there are any), the algorithm selects the one with the smallest amplitude difference. If there are more than one with equal (and smallest) amplitude difference, the one of these with the smallest frequency difference is selected. If there are still more than one candidates, i. e. two or

more candidates with smallest and equal amplitude and frequency difference, simply the first one with the lowest frequency is selected.

If there is no candidate left after this peak filtering process, the estimator failed to detect the partial. We call this partial a *missed peak* or *missed partial*. Otherwise the one and final candidate is assumed as the estimation for the corresponding reference sinusoid. The estimator successfully detected this partial, and we call this a *successfully detected partial*. Although this candidate could fulfil the criteria for several (closely spaced) reference partials, this candidate is removed from the list of potential candidates for subsequent reference sinusoids. Doing this or not is a matter of interpretation and not done in [17]. But consider for example, if we have only one maximum in the FFT spectrum (and so only one peak estimation), we cannot distinguish if this peak is caused by one or more input sinusoids. It so seems more natural, to count this as one detected peak only.

At the end of this process, the algorithm picked for every input sinusoid a corresponding peak estimation and removed it from the list of available peak estimations or candidates. Especially in the presence of noise there are still many candidates left which were not assigned to a reference partial in the process before: we call these falsely detected peaks *spurious peaks* or *phantom peaks*.

Please note that in [17] these spurious peaks are referred to as ‘misdetected peaks’ with the subscript ‘miss’. Missed reference partials are denoted by the subscript ‘missed’ there.

The number of successfully detected partials plus the number of missed partials per frame is equal to the total number of reference partials in that frame. In the same way the sum of phantom peaks and successfully detected peaks has to equal the total number of peaks the estimator detected.

6.4. Characteristic Touchstones

Different criteria are applied to numerically compare the performance of the estimation methods. It is not a goal of this project, to provide or find a single-number value which describes the overall performance of a particular method. This is probably not very useful, as the performance of a special method highly depends on the application it is used for. We therefore better look at different criteria which provide insight to the suitability of an estimator according to given requirements.

For the comparison we consider a number of values which give information about the performance of the estimator in the tests described above. The following parameters are taken into account:

1. Frequency error (Freq)

The frequency error of the estimation is the frequency difference between the reference sinusoid and its selected estimation candidate. We express this error in cents (‘percent of number of half-tones’). It is a logarithmic measure of intervals with which an equally tempered half tone corresponds to 100 cents. 1200 cents equal one octave. With k being the reference partial’s frequency in bins and k_{err} the absolute value of the frequency error in FFT bins, the frequency error factor γ_0 computes to:

$$\gamma_0 = \frac{k + k_{err}}{k} > 1 \quad (6.22)$$

With the half-tone factor $\gamma_{ht} = \sqrt[12]{2}$ the frequency error $f_e^{[ct]}$ in cents computes to [11]:

$$f_e^{[ct]} = 100 \frac{\ln(\gamma_0)}{\ln(\gamma_{ht})} \quad (6.23)$$

2. Amplitude error (**Amp1**)

The amplitude error of the estimation is the amplitude difference between the reference sinusoid and its selected estimation candidate. We express this error in dB in terms of power. With A being the amplitude of the reference partial and A' its estimation, the amplitude error $A_e^{[dB]}$ in dB computes to:

$$A_e^{[dB]} = \left| 20 \log_{10} \left(\frac{A'}{A} \right) \right| \quad (6.24)$$

3. Amplitude of missed partials (**Amp1Missed**)

The amplitudes A of partials that are not detected by the estimator are given in dB in terms of power: $A_{Missed} = 20 \log_{10}(A)$

4. Amplitude of phantom peaks (**Amp1Phantom**)

The amplitudes A of spurious peaks that do not correspond to a reference partial are given in dB in terms of power: $A_{Phantom} = 20 \log_{10}(A)$

5. Total number of peaks per frame detected by the analysis method (**nPeaksTot**)

6. Number of successfully detected partials per frame (**nPeaksCorr**)

7. Number of missed partials per frame (**nMissed**)

8. Number of phantom peaks (spurious peaks) per frame (**nPhantom**)

The detection rate of an estimator for a given signal is determined by:

$$\text{DetRate} = \frac{\text{nPeaksCorr}}{\text{nPeaksCorr} + \text{nMissed}} \quad (6.25)$$

For all values the mean value μ , the standard deviation σ and the minimum and maximum values over all frames of the according test signal are provided.

6.5. Comparison Results and Interpretation

The following chapter does not provide many numerical results but focuses on discussion and interpretation. The full tables with all simulation results are included in the appendix. Three sets of simulations with different analysis windows were carried out:

1. *Analysis window best suitable for the estimation method*

The windows and parameters listed in section 6.1 were used. Table A.1 on page 90 in the appendix shows the simulation results.

2. *Periodic Hann window*

All methods apply a periodic Hann window. The results are given in table A.2 on page 96.

3. Sinusoidal window (square root of a shifted Hann window)

Using the sinusoidal window as analysis window for all methods, leads to the results provided in table A.3 on page 102.

In the following discussion we foremost concentrate on the first case, the simulations with the best suitable analysis windows. For the analysis methods and test signals the abbreviations as introduced in section 6.1 and 6.2 are used.

The plain FFT with no further processing (**PlainFFT**) should be the worst case limit for all methods. If an estimator gives worse results than the plain FFT with no correction at all, the method is doing more bad than good and should not be used at all. This is not the case for all of the tested estimators.

6.5.1. Pure Single Stationary Sinusoids

With the **SiStSi** and **SiStSiMF** test signals, all test methods detect exactly one sinusoid for all frames. So there are no missed partials or spurious peaks. In these two tests the methods **SpcReass**, **DrvAlg** and **PhsVoc** produce almost the same (very accurate) results (see table 6.3). They have in common, that they use the same analysis window and are all based on the linear phase evolution model. Especially in the **SiStSiMF** case, all assumptions that are made in the theoretical derivations of these methods are met: The partials are widely spaced (there is only one), they are constant within the frame and there is no interference from negative mirror frequencies. All the conventional methods show the same amplitude error in the **SiStSiMF** case. The amplitude estimation is based on the frequency estimation. Although the frequency error tends towards zero, the amplitude error reaches a lower limit. This limit is given by the resolution of the main lobe shape lookup table which is used for the amplitude estimation. With the analytical expression of the main lobe shape instead of the lookup table, the amplitude error tends towards zero, too.

SiStSi			SiStSiMF		
	Freq	Ampl		Freq	Ampl
	μ , in cent	μ , in dB		μ , in cent	μ , in dB
ParInt	0.057545	0.001267	ParInt	0.000232	0.000174
TrglAlg	0.006238	0.000185	TrglAlg	0.000000	0.000175
SpcReass	0.052656	0.000924	SpcReass	0.000010	0.000175
DrvAlg	0.052561	0.000926	DrvAlg	0.000002	0.000175
PhsVoc	0.052622	0.000921	PhsVoc	0.000001	0.000175
OdfAlg	0.709787	0.045253	OdfAlg	0.044944	0.045254
MdctAlg	0.040463	0.001610	MdctAlg	0.000003	0.000018

Table 6.3.: Amplitude and frequency errors with noiseless single stationary sinusoids.

The **TrglAlg** outperforms the other conventional methods in the **SiStSiMF** case, its frequency error is basically zero. The **ParInt** does not provide as good results here, the frequency error reaches a lower bound. This is probably caused by the non-ideal parabolic shape of the analysis window main lobe, which is the theoretical restriction of this method.

For all conventional methods the results in the **SiStSi** case with lower frequencies are not as superior as at $f \approx \frac{f_s}{4}$. This results from the contributions of the negative mirror frequencies. The side lobes of the negative mirror partial unfavourably influence the parameter estimation. This effect is stronger for the **ParInt**, as the analysis window that is used with this method has a much lower side lobe attenuation.

Again, the `TrglAlg` outperforms the other conventional methods, approximately by a factor of 10. Although it uses a Hann window like the other methods, it seems like it is less susceptible to the disturbance of the negative mirror partial. The author states in [18] a maximum frequency error of 0.1 Hz, which translates at $l_0 = 5$ to approximately 0.8 cents. Our simulations show a much smaller maximum frequency error. Even with using signals starting from $l_0 = 4$ as stated in the paper, the maximum frequency error occurring is only 0.3 cents.

The `MdctAlg` method performs comparable to the ones based on the linear phase evolution model. With the `SiStSiMF` signals the estimation error tends towards zero. In the other case with the `SiStSi` frequencies, the method gives slightly better frequency estimations and slightly worse amplitude estimations than the other conventional methods (except `TrglAlg`). Please note, that only by using the parameters for the `MdctAlg` algorithm as provided in section 6.1 exactly one peak is detected. Otherwise a small number of partials is missed and spurious peaks with low amplitudes can occur.

Far behind with these test signals is the `OdftAlg` method. Even in the `SiStSiMF` case the frequency and amplitude estimation errors are relatively high. This means, that even in optimal conditions the algorithm does not converge to the true value. In the lower frequency signal the frequency error is approximately ten times larger than with the other methods. The amplitude error is approximately 300 times bigger, which confirms the superiority of the amplitude correction using the window main lobe. The author of the method declares in [14] that the maximum frequency estimation error is approximately 1% of the bin width, which translates to 0.43 Hz, and is rather constant over the whole frequency spectrum. At the frequency bin $l_0 = 5$, which is the lowest component of the test signal, this 0.43 Hz error translates into an error of 3.46 cents. This complies with the maximum frequency error of the simulations of 3.3 cents. All other methods have an approximately three times smaller maximum frequency error.

Summarizing, the `TrglAlg` achieves best results here. The `SpcReass`, `DrvAlg`, `PhsVoc` and `MdctAlg` give approximately equal results, closely followed by `ParInt`. The `OdftAlg` performs comparatively unfavourable here.

6.5.2. Single Stationary Sinusoids with Noise

With the test signals `SiStSiHN` with $SNR = 12$ dB and `SiStSiLN` with $SNR = 40$ dB every method (except the `MdctAlg`) clearly detects the input signal's partial. Additionally all methods detect a large number of spurious peaks, induced by the presence of noise.

In the low noise case the maximum amplitude of these spurious peaks is between -44 dB and -55 dB, the mean at around -60 dB, so well below the noise power level of -40 dB. Same in the high level noise case: with a mean at around -35 dB and the maxima at around -21 dB the amplitude of the phantom peaks is far below the noise level of -12 dB. This shows, that the amplitude estimations of all the methods are stable even in the case of a pure noise signal: with a random signal as basis, the amplitude estimation does *not* give random values. This would lead to very disturbing spurious peaks with high amplitude. Instead of that all values are below the noise level.

The number of detected peaks should be as low as possible. It is basically constant for both cases and is around 130 peaks per frame for `SpcReass`, `DrvAlg`, `PhsVoc` and `OdftAlg`. Only around 50 peaks per frame are found with `ParInt`, `MdctAlg` and `TrglAlg`. What makes the

difference here, is that the later methods include in their algorithms some sort of detection: They check if the found peak is actually caused by a partial in the input signal or not. The shape of the window main lobe should be present in the spectrum in the case of a real peak.

The `MdctAlg` fails to detect the signal's partial in two frames of the `SiStSiLN` input signal and in approximately 5% of the frames of the `SiStSiHN` signal. The algorithm is lacking robustness in these situations: both the coefficients α and β are out of range and so the peak has to be discarded. Otherwise this would result in a spurious peak with a very high amplitude (see chapter 5.4 on page 48).

The low noise test signals confirm the picture of the noiseless case: As table 6.4 shows, again the `TrglAlg` achieves best results in amplitude and frequency estimation. `SpcReass`, `DrvAlg` and `PhsVoc` perform identically, closely followed by the `ParInt`. The `MdctAlg` degrades in the presence of noise. The `OdftAlg` is again endmost.

In the presence of high noise the `TrglAlg` loses its top most position and performs slightly worse than `SpcReass`, `DrvAlg` and `PhsVoc`, which again give equal results. `ParInt` follows closely, the `OdftAlg` only with some distance. The high noise impairs the `MdctAlg` more than all the other algorithms. The mean frequency error is approximately four times as high as with `SpcReass`.

SiStSiLN			SiStSiHN		
	Freq			Freq	
	μ , in cent	μ , in dB		μ , in cent	μ , in dB
ParInt	0.077415	0.003117	ParInt	0.838156	0.073368
TrglAlg	0.029341	0.002771	TrglAlg	0.693989	0.069515
SpcReass	0.059189	0.003021	SpcReass	0.591390	0.064072
DrvAlg	0.059246	0.003024	DrvAlg	0.593009	0.064101
PhsVoc	0.059156	0.003016	PhsVoc	0.591396	0.064066
OdftAlg	0.719328	0.046593	OdftAlg	1.230023	0.084639
MdctAlg	0.159518	0.008283	MdctAlg	2.774101	0.144906

Table 6.4.: Amplitude and frequency errors with noisy single stationary sinusoids.

To put it in a nutshell, all methods degrade in the presence of noise. The ranking between the methods stays approximately the same as in the noiseless case. Only the `MdctAlg` shows less robustness against noise compared to the other methods.

6.5.3. Two Stationary Sinusoids

The `TwoStSi` signal shows the spectral resolution of the analysis methods. Figure 6.3 on the next page gives an overview of the peaks that are detected by every method. As soon as the reference partials get too close together, the analysis algorithms are not capable any more to distinct between them. This leads to an unstable behaviour: Either the one or the other partial is detected, or the power of the partials is accumulated into one peak. This leads to spurious peaks with up to double the amplitude of the input partials. The maximum amplitude of the phantom peaks is around 6 dB.

The frequent changes of the number of detected peaks in the transition range result from the initial phase of the two sinusoids. The initial phase of the two partials, and so the phase difference between the two partials, is different in every frame. It is set in a way so that there are no time signal discontinuities at the analysis frame borders.

The `OdftAlg` shows the best spectral resolution of approximately 1.8 bins. The frequency

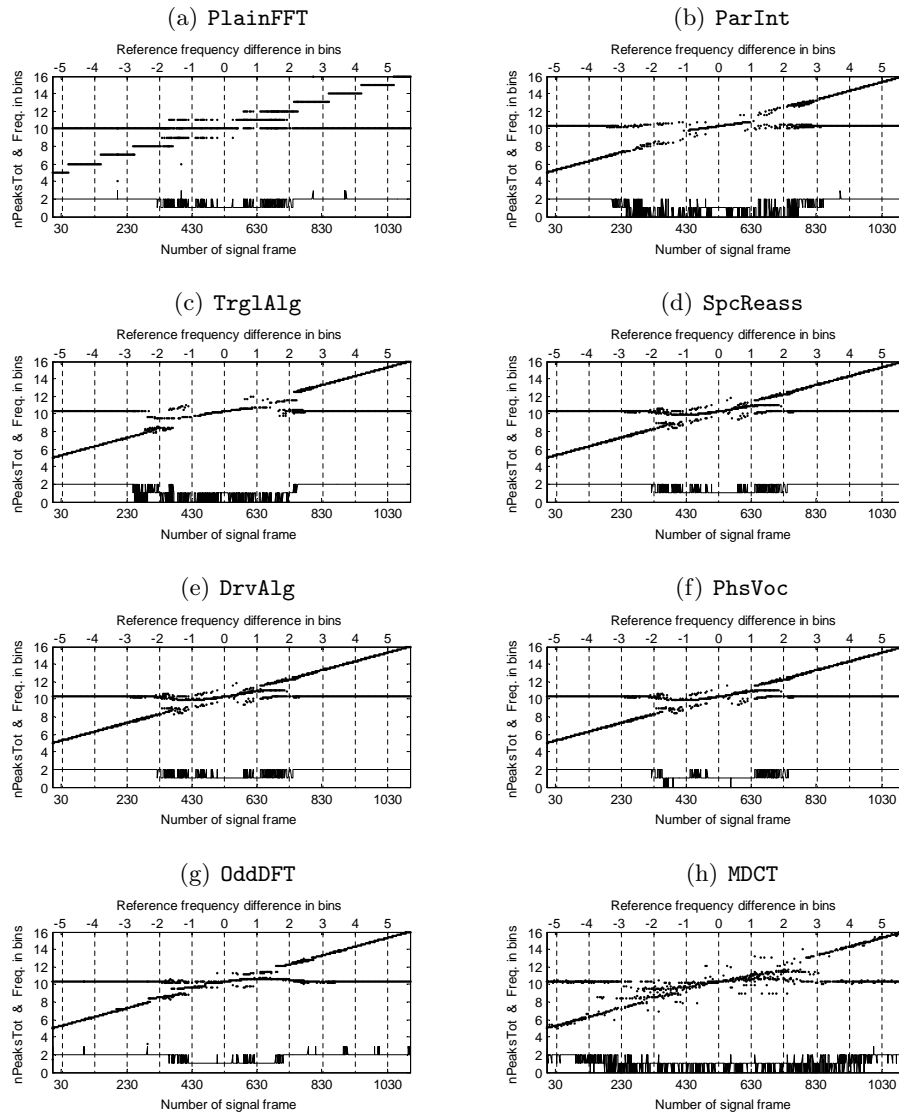


Figure 6.3.: Detected peaks for pairs of closely spaced stationary partials in dependence of their frequency difference. Bottom lines show the total number of detected peaks per frame.

resolution of the `ParInt` algorithm evaluates to approximately 3.3 bins, the one of `TrglAlg` to 3 bins. `SpcReass`, `DrvAlg`, `PlainFFT` and `PhsVoc` all show the same resolution of approximately 2.1 bins.

One actually observes here the spectral resolution of the underlying analysis window. The later ones all use a Hann window with a main lobe width of approximately of 4 bins, so 2 bins towards every side of the peak. `OdfAlg` uses the sinusoidal window which has a significantly smaller main lobe (see figure 5.3 on page 41). The artificially designed window for `ParInt` has a wider window main lobe of around 6 bins. In the `TrglAlg` case on the right side of the transition zone, the algorithm restarts to detect two partials again as soon as they have a frequency distance of two bins, which is according to the used Hann window. For some reason on the left side the detection fails as soon as the frequency distance is lower than three bins.

The `MdctAlg` is an exception here, too: although it uses the same narrow sinusoidal window as the `OdfAlg`, the frequency resolution is low. Obviously the crossproducts between the partials strongly degrade the peak detection capability of the algorithm. The `MdctAlg` does not provide stable results for closely spaced partials any more. The frequency resolution is around 5 bins only.

The spectral resolution expresses itself as well in the mean value of the number of correctly detected partials. As better the resolution, as more often the frequency and amplitude estimation is successful. The number of correctly detected partials has to be considered carefully when looking at the estimation accuracy, as only correctly detected peaks are considered for calculating the estimation errors. As higher the number of correctly detected peaks, as closer spaced partial pairs are considered. At those the impact from one partial lobe onto the other partial is bigger, and so the estimation error increases.

`SpcReass`, `DrvAlg` and `PhsVoc` again give the same results. The `OdfAlg` has slightly bigger estimation errors, but has a higher detection rate, whereas with `ParInt` and `TrglAlg` it is just the other way round: their accuracy is slightly higher, therefore less partials were correctly detected. The `MdctAlg` generates estimation errors comparable to the `OdfAlg`, but the detection rate is the worst of all methods.

Please note, that the maximum amplitude error has an upper bound of 3 dB. If the amplitude differs more than 3 dB from the reference partial, the detected peak is considered as phantom peak.

As a conclusion this test mainly shows the dependency between the spectral resolution of the estimator and the used analysis window.

6.5.4. Single Non-Stationary Sinusoids

Again, with the non-stationary test signals `Sweep`, `Tremolo`, `Vibrato` and `TremVib` all methods except the `MdctAlg` clearly detect exactly one partial.¹

The `MdctAlg` again sometimes fails to detect a partial at all, as it is unable to extract valid parameters from the signal. This happens particularly often (for up to 10% of the analysis frames) when the frequency of the partial is changing within the frame. Additionally the frequency estimation error of the `MdctAlg` is very sensitive to frequency modulations. The performance degrades much more than for the other methods, so that the frequency

¹Only in very rare cases with a vibrato test signal, the `OdfAlg` sometimes detects a second partial. But its amplitude is almost as small as the threshold value for zeroing the transform coefficients and so negligible.

error is approximately five times as high. The method is more resistant against amplitude modulations.

As expected, the accuracy of all methods degrades in the presence of time variations (see table 6.5). Again `SpcReass`, `DrvAlg` and `PhsVoc` perform identically well. The amplitude and frequency errors of `ParInt` and `TrglAlg` are sometimes a little bit better, sometimes a little bit worse than the first three methods, but basically in the same range. In contrast to that, the `OdftAlg` produces clearly worse results, but still better than the `MdctAlg`.

Sweep			Tremolo		
	Freq	Ampl		Freq	Ampl
	μ , in cent	μ , in dB		μ , in cent	μ , in dB
<code>ParInt</code>	0.042790	0.001397	<code>ParInt</code>	0.046026	0.016799
<code>TrglAlg</code>	0.010170	0.000667	<code>TrglAlg</code>	0.078389	0.015752
<code>SpcReass</code>	0.024623	0.001309	<code>SpcReass</code>	0.088084	0.013586
<code>DrvAlg</code>	0.028425	0.001006	<code>DrvAlg</code>	0.114295	0.013591
<code>PhsVoc</code>	0.027385	0.001331	<code>PhsVoc</code>	0.088171	0.013592
<code>OdftAlg</code>	0.899852	0.056877	<code>OdftAlg</code>	1.108961	0.047567
<code>MdctAlg</code>	1.425574	0.098121	<code>MdctAlg</code>	1.953542	0.065469

Vibrato			TremVib		
	Freq	Ampl		Freq	Ampl
	μ , in cent	μ , in dB		μ , in cent	μ , in dB
<code>ParInt</code>	1.653547	0.003081	<code>ParInt</code>	1.655514	0.016680
<code>TrglAlg</code>	1.540901	0.002679	<code>TrglAlg</code>	1.548640	0.015060
<code>SpcReass</code>	1.335832	0.005106	<code>SpcReass</code>	1.332840	0.014120
<code>DrvAlg</code>	1.337586	0.005088	<code>DrvAlg</code>	1.337831	0.014421
<code>PhsVoc</code>	1.336262	0.005109	<code>PhsVoc</code>	1.333235	0.014112
<code>OdftAlg</code>	1.563133	0.084149	<code>OdftAlg</code>	1.680363	0.087315
<code>MdctAlg</code>	5.403369	0.178947	<code>MdctAlg</code>	5.653225	0.177485

Table 6.5.: Amplitude and frequency errors with single non-stationary sinusoids.

6.5.5. Musical Instrument Sounds

The test with the `Sax`, `Guitar` and `Voice` test signals gives a very mixed picture: the overall performance of all the estimators together highly depends on the used test signal.

The following remarks apply for all analysis methods except the `MdctAlg` (which performs less good in many points):

In the `Guitar` sound basically all partials are detected correctly. There are hardly any missed peaks or spurious peaks. The methods perform equally well, with some better results in the frequency estimation of `ParInt` and `TrglAlg`. The mean frequency error is less than 2% and the mean amplitude error is below 0.33 dB for all the methods, which can be considered as good results for a natural-like sound.

The opposite is the case in the `Voice` sound: The analysis methods miss approximately 20% of the partials. A large number of these missed partials has very low amplitudes. Most of the other ones are not detected due to the restricted spectral resolution of the analysis. The mean amplitude of the missed partials is around -70 dB with maxima up to -13 dB. The power of partials that were missed due to lack of spectral resolution often leads to phantom peaks, which show up with amplitudes up to -16 dB and a mean value of -60 dB.

The mean frequency error is with 4.5 cents approximately the same for all the methods, which compares to 12.5 cents with `PlainFFT` without any processing. The mean amplitude error is for all methods almost as big as when no amplitude correction at all is conducted.

The results perceived with the **Sax** sound are basically in between the two extreme cases **Guitar** and **Voice**. The reason for the big differences are quite obvious: where the guitar sound has a low number of partials and changes after the attack only very slowly in time, the voice sound has much more partials and is not stationary at all. Table 6.6 summarizes the results.

As mentioned before, the above comments do not include the **MdctAlg**: its detection rate and frequency estimation accuracy is noticeably lower than with all the other methods.

Guitar				Sax			
	Freq	Ampl	DetRate		Freq	Ampl	DetRate
	μ , in cent	μ , in dB	μ		μ , in cent	μ , in dB	μ
PlainFFT	27.946719	0.659374	0.992	PlainFFT	7.480726	0.916054	0.918
ParInt	0.824619	0.207893	0.936	ParInt	1.270960	0.715987	0.864
TrglAlg	0.756137	0.215648	0.991	TrglAlg	1.249471	0.733356	0.924
SpcReass	1.715053	0.227606	0.992	SpcReass	1.258186	0.734249	0.928
DrvAlg	1.726644	0.227958	0.992	DrvAlg	1.371690	0.719767	0.927
PhsVoc	1.703284	0.226826	0.991	PhsVoc	1.259803	0.734434	0.928
OdftAlg	1.886747	0.321352	0.991	OdftAlg	1.306614	0.739361	0.927
MdctAlg	3.117508	0.340208	0.774	MdctAlg	1.515316	0.726742	0.817

Voice			
	Freq	Ampl	DetRate
	μ , in cent	μ , in dB	μ
PlainFFT	12.509649	0.861879	0.793
ParInt	4.421437	0.749896	0.747
TrglAlg	4.321297	0.746076	0.789
SpcReass	4.591920	0.767413	0.806
DrvAlg	4.581522	0.768416	0.806
PhsVoc	4.592742	0.766632	0.804
OdftAlg	4.619156	0.796307	0.801
MdctAlg	5.126278	0.825206	0.534

Table 6.6.: Amplitude and frequency errors and detection rate with musical instrument sounds.

6.5.6. Hann Window for all Methods

Instead of their designated windows given in 6.1, for **ParInt**, **OdftAlg** and **MdctAlg** a periodic Hann window is used for signal analysis. All other methods anyway used the Hann window already in the simulations considered beforehand. Table A.2 on page 96 in the appendix shows the full results.

With the Hann window, the performance of **OdftAlg** and **MdctAlg** decreases dramatically. The mean frequency errors increase up to 20 cents and 30 cents respectively. So changing the window makes these methods basically unusable. They heavily depend on the use of the special sinusoidal window.

The results are not as bad with **ParInt**. Nevertheless its frequency and amplitude errors are higher than for all the other conventional methods then. The detection rate for closely spaced sinusoids is quite low, too. Only for the musical instrument sounds (**Sax**, **Guitar**, **Voice**) the method still performs comparable to the other estimators, although the detection rate is getting worse there as well.

6.5.7. Sinusoidal Window for all Methods

The performance of the estimators in the case of that for given circumstances the sinusoidal window (square root of a shifted Hann window) should be used is verified, too. This window is commonly applied in conjunction with the modified discrete cosine transform. The full results are provided in table A.3 on page 102 in the appendix.

The `ParInt` and `TrglAlg` estimators do not work well anymore with this analysis window. Even with the simple `SiStSi` and `SiStSiMF` test signals the peak detection with the `ParInt` method is not reliable anymore. In a high number of frames the partial is missed. With the `TrglAlg` at least all partials are detected, but the frequency estimation errors are unacceptably high.

Again, the group of the three conventional estimators `SpcReass`, `DrvAlg` and `PhsVoc` gives almost identical results. They all degrade with the use of this window. The accuracy compared to the other two methods based on shifted Fourier transforms depends on the test signal:

For the noiseless single stationary sinusoids the `MdctAlg` provides more accurate estimations than `SpcReass`, `DrvAlg` and `PhsVoc`. The results of the `OdftAlg` are again not convincing with these test signals.

As more noisy the signal gets, as more the `MdctAlg` degrades and finally gives in the `SiStSiHN` tests the worst results.

The spectral resolution of `SpcReass`, `DrvAlg` and `PhsVoc` improves according to the increased resolution of the window. The results are as good as with the `OdftAlg`. Again the spectral resolution of the `MdctAlg` is quite bad and results in a low detection rate. Figure 6.4 on the next page gives an overview.

For the time-varying test signals the group of `SpcReass`, `DrvAlg` and `PhsVoc` gives throughout the best results, closely followed by the `OdftAlg`. The `MdctAlg` shows significantly bigger estimation errors.

The test with the the synthesized musical instrument sounds `Sax`, `Guitar` and `Voice` gives a surprise: Although the `OdftAlg` did not perform particularly good with the synthetic signals, it here provides by far the best peak detection rate, combined with the smallest frequency and amplitude errors. The `SpcReass`, `DrvAlg` and `PhsVoc` show lower accuracies and detection rates. The `MdctAlg` again highly degrades with increasing complexity of the sound.

6.5.8. Conclusion

The following list gives a short and slightly simplified overview on the achieved results:

PlainFFT The plain FFT without any post-processing is included for reference only. Do not seriously consider this as a parameter estimator, except a very low accuracy is sufficient.

ParInt The parabolic interpolation only gives good results when the specially designed window (the inverse FFT of a parabola) is used. Its performance degrades with the use of different windows. Other estimation methods become more favourable then.

TrglAlg When using the Hann window, the triangular algorithm outperforms the other methods. With other analysis windows by far not as good results are achieved.

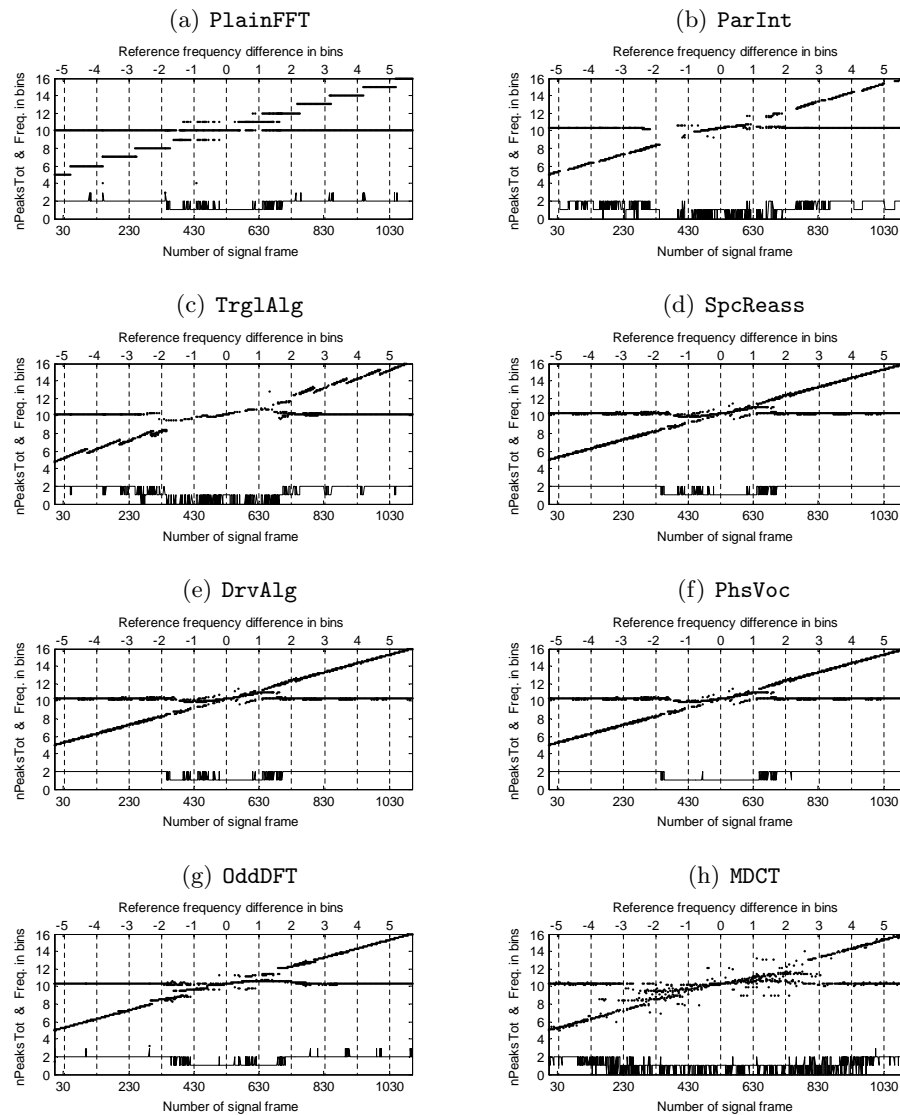


Figure 6.4.: Detected Peaks for pairs of closely spaced stationary partials in dependence of their frequency difference. Bottom lines show the total number of detected peaks per frame. The sinusoidal window is used for all analysis methods.

SpcReass, DrvAlg, PhsVoc The three linear phase evolution based estimators provide the best overall quality. Although they do not provide the most accurate results for every test, they combine a high flexibility with a low computational effort and perform comparably well for all test methods and, what is even more important, with all tested analysis windows.

Although their functional principles are very different, these three methods perform in *all* tests identically well. One can arbitrary select one of them. The actual choice then depends on which implementation fits best in the given framework. The computational effort is comparable, as they all three require a second FFT of the same length. If a sliding window approach with a small hop size is used anyway, the **PhsVoc** method might be computationally advantageous.

OdfAlg The use of this estimator is only beneficial if for given circumstances the sinusoidal window (square root of a shifted Hann window) has to be applied. The analysis of musical instrument sounds shows the best results then.

MdctAlg This method is only recommendable if following prerequisites are fulfilled:

1. Due to external circumstances you want to use the sinusoidal window (square root of a shifted Hann window).
2. You have a very clean and static signal.

The performance degrades very quickly in the presence of noise and time-variations. Additionally the spectral resolution is dissatisfactory.

Chapter 7.

Improvements

Chapter 6 confirmed the high practicality and accuracy of the well-established DFT based estimators like the spectrum reassignment method, the derivative algorithm and the phase vocoder method. The implemented ODFT and MDCT estimation algorithms showed significant limitations.

7.1. Parabolic Interpolation Between ODFT Coefficients

We showed in chapter 5.1 the need for an estimator based on ODFT coefficients. The survey revealed one limitation of the `OdftAlg`: the frequency estimation does not converge to the true value in the case of very ‘clean’ test signals. As the ODFT spectrum is very similar to the DFT spectrum, we try a parameter estimation based on the ODFT and using parabolic interpolation. Similar to the naming introduced before we call this estimation method `OdftParInt`.

The method follows exactly the same procedure as the conventional DFT based parabolic interpolation `ParInt`, which is described in section 4.1.3 on page 27.

Instead of the DFT the ODFT of the input signal is applied. We can compute the ODFT by using a fast FFT implementation. The ODFT of an input signal $x[n]$ which is windowed by a window function $w[n]$ is given by equation 3.12 and repeated here. The following equation shows a computationally effective way to compute the ODFT by transforming it to a DFT with a modified window:

$$X_{ODFT}(k) = \sum_{n=0}^{N-1} w[n]x[n]e^{-j\frac{2\pi}{N}n(k+\frac{1}{2})} = \sum_{n=0}^{N-1} w[n]e^{-j\frac{\pi}{N}n}x[n]e^{-j\frac{2\pi}{N}nk} \quad (7.1)$$

$$= \sum_{n=0}^{N-1} w'[n]x[n]e^{-j\frac{2\pi}{N}nk} \quad \text{with} \quad w'[n] = w[n]e^{-j\frac{\pi}{N}n} \quad (7.2)$$

This is the DFT of the input signal $x[n]$ windowed by a modified window $w'[n]$. The complex valued $w'[n]$ is a phase shifted version of the original window $w[n]$ and can be computed once beforehand. The DFT of the then complex valued sequence $w'[n]x[n]$ is finally computed using a standard FFT algorithm.

For the parabolic interpolation the same equations as before are used. Only the final frequency estimation has to be increased by 0.5 frequency bins.

The `OdftParInt` method is included in the simulations and its results can be found again in the tables A.1, A.2 and A.3 in the appendix. The same parameters as for the `ParInt` method were used.

As expected, `OdfParInt` provides basically the same results as the conventional `ParInt` throughout all the tests. Using the specially designed parabolic window the `OdfParInt` provides much better results than the `OdfAlg`.

Unfortunately this does not help for the setting described in 5.1: for the intended application in MDCT based audio coders, a sinusoidal window or a KBD window has to be used. But the parabolic interpolation methods highly degrade with the use of these windows. It might be helpful to adopt one of the `SpcReass`, `DrvAlg` or `PhsVoc` methods to the use of the ODFT transform. But even in the case they would then give as good results as with the conventional DFT, they would not outperform the `OdfAlg` (see section 6.5.7 on page 72 and table A.3).

Let us nevertheless take a closer look at the frequency estimation errors of the two parabolic interpolation methods with DFT and ODFT. As before we denote the input partial's frequency $l = l_0 + \Delta l$ by its integer part l_0 and its fractional part Δl . We expect a systematic error of the frequency estimation. The accuracy should depend on the position Δl of the input partial between two frequency bins.

Figure 7.1 plots the absolute value of the the frequency error of the estimation of single stationary sinusoids. The integer frequency of the input sinusoid is varied in 50 steps in between $l_0 = 12$ and $l_0 = 502$ bins (with $N = 1024$). The fractional frequency Δl is varied in 200 steps from $\Delta l = 0$ to $\Delta l = 1$. This results in a total number of 8000 analysed sinusoids, which were analysed by parabolic interpolation between the DFT coefficients (`ParInt`) using the specially designed parabola IFFT window.

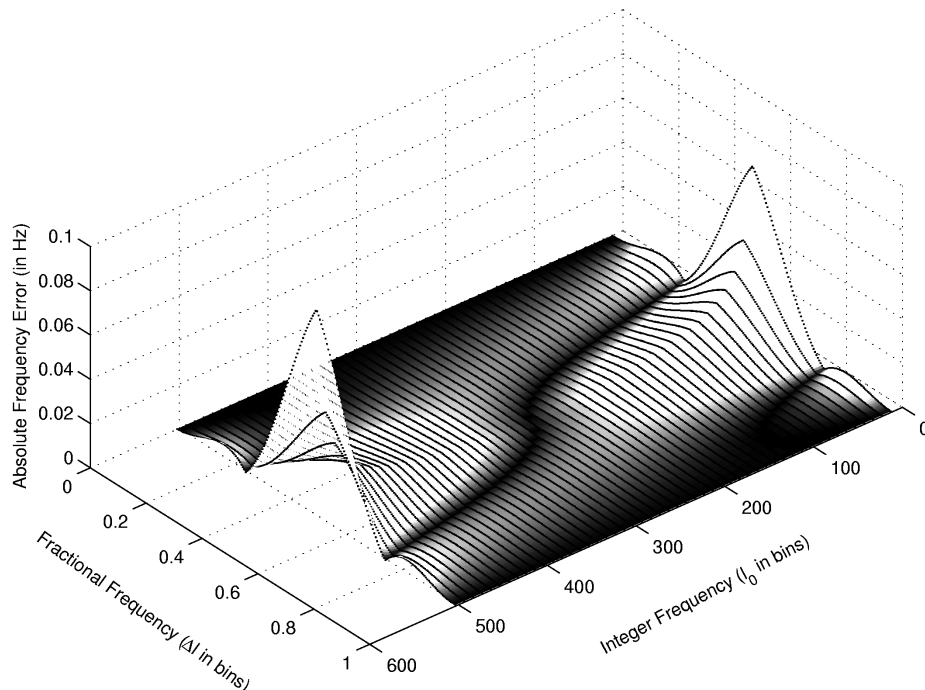


Figure 7.1.: Frequency error of `ParInt` for stationary sinusoids ($N = 1024$, $l_{0,min} = 12$).

The plot clearly shows a dependency between the estimation error and the fractional frequency. The estimation error is in general biggest in-between two frequency bins ($\Delta l = 0.5$).

Varying the integer frequency, the estimation error increases towards the boundaries of the frequency spectrum. This is caused by the increasing influence of the side lobes of the negative mirror peaks, which disturb the estimation. The slightly asymmetrical shape along the $\Delta l = 0.5$ axis results as well from the negative mirror peaks: at the high integer frequencies, the partials with $\Delta l > 0.5$ are closer to their mirror peaks than those with $\Delta l < 0.5$. For partials with low integer frequencies it is just the other way round. It is more than noteworthy, that the biggest estimation errors occur in the interval $0.25 \leq \Delta l \leq 0.75$.

The same plot is given in figure 7.2 using parabolic interpolation between the ODFT coefficients (`OdftParInt`). The same observations as before can be made, except that everything is shifted by 0.5 frequency bins. The biggest estimation errors now occur in the intervals $0 \leq \Delta l \leq 0.25$ and $0.75 \leq \Delta l \leq 1$ with the maximum errors at $\Delta l = 0$, so just in the areas where the DFT method has low frequency errors.

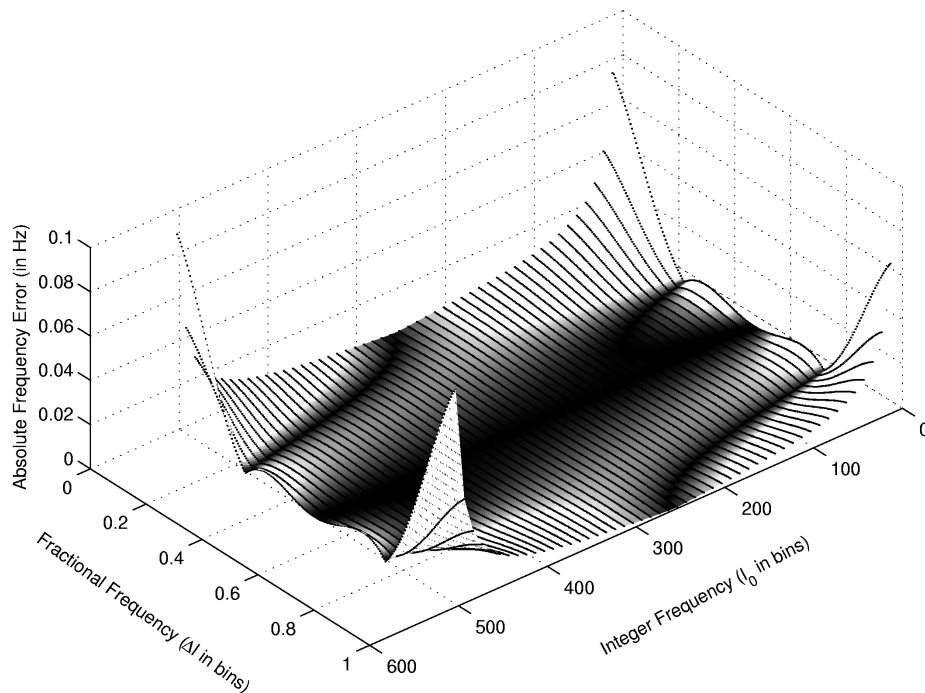


Figure 7.2.: Frequency error of `OdftParInt` for stationary sinusoids ($N = 1024$, $l_{0,min} = 12$).

We therefore propose an improvement to the widely used parabolic interpolation that tries to suppress the influence of the negative mirror peaks by combining the two ODFT and DFT based methods. We will denote this method by `DualParInt`.

The new method basically performs the parabolic interpolation algorithm in parallel for the DFT and ODFT coefficients. Then the results of one or the other method are used, depending on the just estimated fractional frequency Δl .

So after applying the parabola IFFT window on the input signal, its DFT and ODFT coefficients are computed. The resulting spectra are searched for local maxima. We only keep peaks that are present in both spectra. Then frequency and amplitude of every peak is computed independently from both spectra with parabolic interpolation as described before. We arbitrary choose the DFT based fractional frequency estimation Δl_{dft} as selection criterion

for the final estimation value: if Δl_{dft} is within the interval $0.25 \leq \Delta l_{dft} \leq 0.75$, the frequency and amplitude estimation values obtained from the ODFT spectrum are used for that peak. Otherwise the estimation values obtained from the DFT spectrum are assigned.

Figure 7.3 shows the resulting frequency estimation error. The same test signals as described above are used. Please note the different scaling of the frequency error axis. The maximum error for these test signals is approximately ten times smaller than in figure 7.1 and 7.2.

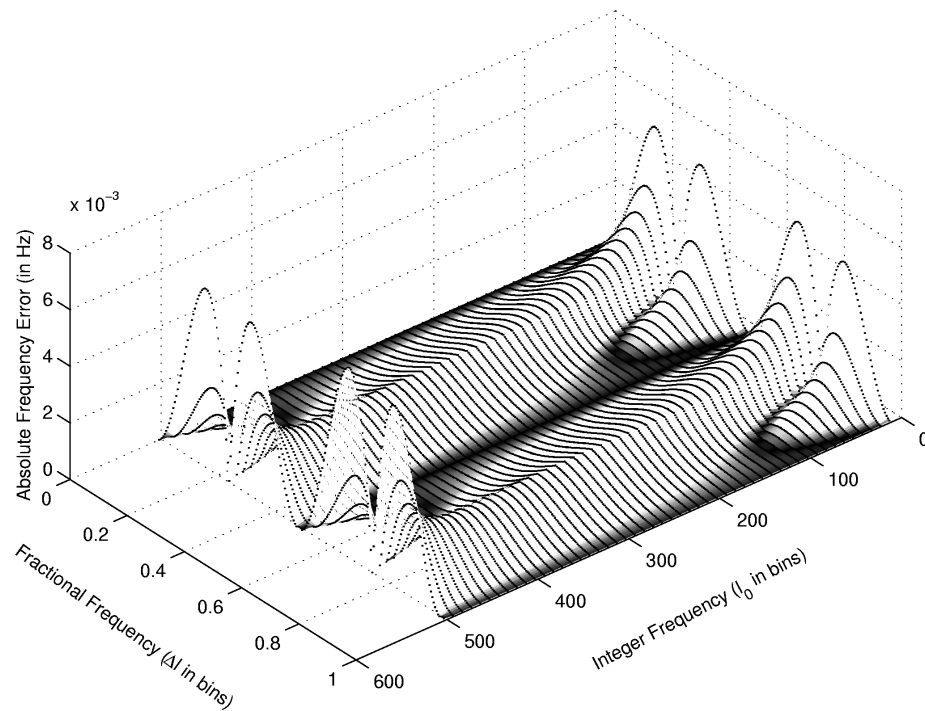


Figure 7.3.: Frequency error of `DualParInt` for stationary sinusoids ($N = 1024$, $l_{0,min} = 12$).

This `DualParInt` estimation method is as well included in our simulation results presented in tables A.1, A.2 and A.3 in the appendix. With the specially designed IFFT parabola window and the stationary low-noise signals the results are significantly better than with `ParInt` and `OdftParInt`. With increasing complexity of the sound (time variations, high noise, ...) other disturbances become dominant. They then overrule the influence of negative mirror peaks in the estimation, which is the issue we addressed here, and the difference to the `ParInt` becomes negligible.

In the same way, another source of error that becomes dominant is the use of a different window than the IFFT of a parabola. We use now a periodic Hann window for analysis of the same signals as before. The frequency estimation error of the `ParInt` method is plotted in figure 7.4 on the following page.

Compared to figure 7.1 the absolute frequency error is now much bigger and constant over the entire frequency range. The integer frequency has basically no more influence. The error only depends on the position of the partial between two frequency bins. This traces back to the fact that the error is caused by the interpolation itself: the used Hann window is by far

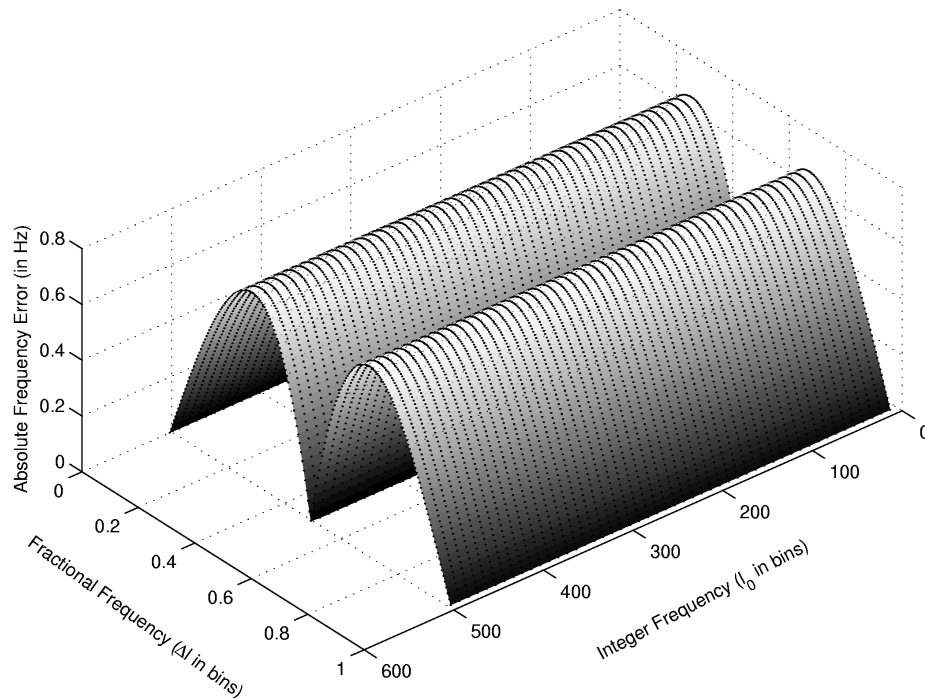


Figure 7.4.: Frequency error of `ParInt` for stationary sinusoids ($N = 1024$, $l_{0,min} = 12$, Hann window).

not parabolic anymore.

As the plot shows, the frequency error is systematically predictable. It so might be possible to improve the estimation with the help of an error correction function that depends on the fractional frequency only.

The frequency error is very low for $\Delta l = 0$ and $\Delta l = 0.5$. This is because as long as the window shape is symmetric around the centre, the parabolic interpolation always gives correct values there by definition. The frequency errors do not become zero at these fractional frequencies though. They only converge to the values shown in figure 7.1 for same Δl , because the influence of the mirror peaks is of course still present.

7.2. Direct Fitting of the Analysis Window Main Lobe

The common idea of all interpolation methods is, that the spectrum basically consists of a sum of modulated and scaled analysis window Fourier transforms. For single stationary sinusoids the spectrum of the windowed signal consists of the shape of the analysis window in the frequency domain, shifted to the frequency and scaled to the amplitude of the input partial.

For a given analysis window the classic interpolation based methods theoretically derive an analytic expression that describes the spectrum as a function of the input partial's frequency and amplitude. Then the given spectrum values are put into the equation and frequency and amplitude of the input partial remain the only unknown. Usually it is too difficult to resolve this expression for the unknown quantities. Therefore the equation for the window main lobe

is approximated by an analytically less complex function.

We want to circumvent this approximation and therefore approach the problem the other way round: We again use the analytic expression of the shape of the windowed sinusoid in the frequency domain. We then modify the two parameters amplitude A and frequency l of the expression to make it match the actual spectrum around a local maximum as closely as possible.

This process could be visualized by moving the window main lobe along the x-axis and rescaling it in a way that the spectrum values around a peak all lie on the window main lobe. With a single sinusoid only and neglecting negative mirror frequencies this can be perfectly achieved. The apex of the so positioned window main lobe then determines the amplitude and frequency estimation. In the non-ideal situation the lobe is positioned in a way to minimize the amplitude differences between the spectrum coefficients and the according main lobe values in a least squared error sense. Again, to maintain comparability to the estimation methods introduced before, we consider three consecutive samples of the frequency spectrum, although considering more samples would improve the stability of the estimation.

Let us again assume an input sinusoid $x[n]$ with frequency $l = l_0 + \Delta l$ and amplitude A as given in equation 5.1. The input signal is windowed by a window $w[n]$, whose magnitude response we denote by $|W(\omega)|$ with $\omega = \frac{2\pi}{N}(k + \Delta k)$. It is an even function, so $|W(\omega)| = |W(-\omega)|$. The DFT spectrum $X_{DFT}(k)$ of the windowed sinusoid has a local maximum at $k = l_0$. Similar to equation 4.3 we define an amplitude scaling factor $\alpha = \frac{A}{2}$. Figure 7.5 illustrates the considered setup.

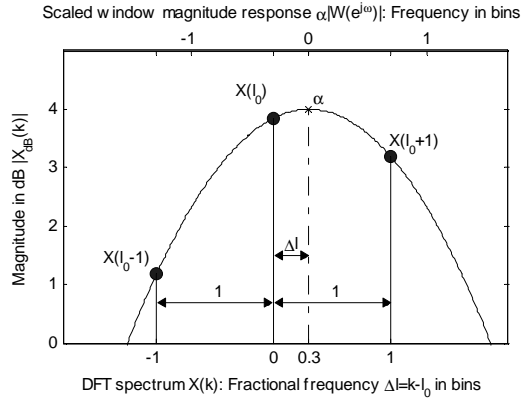


Figure 7.5.: DFT spectrum of a windowed sinusoid with fractional frequency $\Delta l = 0.3$ and the scaled window magnitude response.

Obtained from the figure the amplitudes $|\hat{X}_{DFT}(k)|$ of the DFT coefficients are given by:

$$|\hat{X}_{DFT}(l_0 - 1)| = \alpha |W(\frac{2\pi}{N}(\Delta l + 1))| \quad (7.3)$$

$$|\hat{X}_{DFT}(l_0)| = \alpha |W(\frac{2\pi}{N}\Delta l)| \quad (7.4)$$

$$|\hat{X}_{DFT}(l_0 + 1)| = \alpha |W(\frac{2\pi}{N}(\Delta l - 1))| \quad (7.5)$$

We now want to numerically fit the analytic expression of the window main lobe into the three DFT samples points. It is equivalent to determining the amplitude scaling factor α and the frequency shift Δl . In practice, the window shape will not exactly fit into the given

points. We then want to determine Δl and α in a way to minimize the squared error of the curve fitting. We could define the error as a function of two variables. This could then be minimized by one of the mathematically highly developed non-linear least-squares (LS) optimization methods like the Gauss-Newton or the Levenberg-Marquardt method [8].

Δl has a non-linear influence on the error, as $|W(\omega)|$ is a non-linear function. Fortunately α is a linear term which can be removed from the error function by considering ratios instead of the absolute values. The optimization in two variables is therefore reduced to a minimization in one variable, which is computationally much more inexpensive.

We therefore consider the ratios γ_l and γ_r between the maximum peak and its neighbours. With the theoretical amplitudes $|\hat{X}_{DFT}(k)|$ the left and right-hand side window main lobe ratios are given by:

$$\frac{|\hat{X}_{DFT}(l_0)|}{|\hat{X}_{DFT}(l_0 - 1)|} = \frac{|W(\frac{2\pi}{N}\Delta l)|}{|W(\frac{2\pi}{N}(\Delta l + 1))|} = \gamma_l \quad (7.6)$$

$$\frac{|\hat{X}_{DFT}(l_0)|}{|\hat{X}_{DFT}(l_0 + 1)|} = \frac{|W(\frac{2\pi}{N}\Delta l)|}{|W(\frac{2\pi}{N}(\Delta l - 1))|} = \gamma_r \quad (7.7)$$

The quadratic error Q_e between the spectrum values and the scaled and shifted Fourier representation of the window as function of only $\Delta k = k - l_0$ and with the DFT spectrum values $|X_{DFT}(k)|$ evaluates to:

$$Q_e(\Delta k) = \left| \frac{|X_{DFT}(l_0)|}{|X_{DFT}(l_0 - 1)|} - \frac{|W(\frac{2\pi}{N}\Delta k)|}{|W(\frac{2\pi}{N}(\Delta k + 1))|} \right|^2 + \left| \frac{|X_{DFT}(l_0)|}{|X_{DFT}(l_0 + 1)|} - \frac{|W(\frac{2\pi}{N}\Delta k)|}{|W(\frac{2\pi}{N}(\Delta k - 1))|} \right|^2 \quad (7.8)$$

$Q_e(\Delta k)$ has a minimum at $\Delta k = \Delta k_m$. Neglecting negative mirror frequencies again, the local minimum is at $\Delta k_m = \Delta l$, the fractional frequency of the input sinusoid. So the final frequency estimation is given by $l_0 + \Delta k_m$ with l_0 being the position of the local maximum. The minimum of Q_e is numerically determined, using an algorithm that is based on Golden Section search and parabolic interpolation [25]. The search interval for the minimization of $Q_e(\Delta k)$ can be restricted to $-0.5 \leq \Delta k \leq 0.5$.

The remaining error $Q_e(\Delta k_m)$ is a good estimate, how well the window main lobe fits into the spectrum. As a consequence it is used as a criteria to decide if the detected peak is considered as a partial or if it is induced by e. g. noise. So the least square error $Q_e(\Delta k_m)$ has to be smaller than a certain threshold a_{thr} . Otherwise the peak is discarded. In practice a value of $a_{thr} = 0.01$ gives satisfactory results.

As with all the other methods the amplitude estimation is based on the frequency estimation and calculated using equation 4.31 on page 33.

In order to improve performance, instead of the analytic expression for $|W(\omega)|$ a lookup table for the Fourier representation of the window can be used. It is so possible to compute the ratios γ_l and γ_r for equation 7.8 beforehand. Then two ways are possible:

1. $Q_e(\Delta k)$ is computed for all discrete values of Δk_i and the Δk_i is selected that corresponds to the smallest quadratic error.
2. The minimum of $Q_e(\Delta k_i)$ is numerically determined, using a minimization algorithm that considers the discrete nature of Δk_i .

In this case the resolution of the lookup table is the limiting factor for the accuracy with clean signals.

These methods are included in the comparisons. Using the direct window fitting with the analytical expression for the main lobe, the method is referred to as `DirWinFitAnal`. In the case of the table lookup, the name is `DirWinFitTbl`. For the lookup table, the Fourier representation of the window is oversampled by a factor of 2^{12} again. Tables A.1, A.2 and A.3 in the appendix show the full results. For the first test run with the windows best suitable to the algorithms, the periodic Hann window is used.

Both methods perform comparably well in all tests. Especially with the `SiStSi` test signals they clearly outperform the other methods. Only the `TrglAlg` can keep up with them, but it degrades heavily with other analysis windows. The `DirWinFit` algorithms perform good with the sinusoidal window, too. In the `SiStSiMF` case, the amplitude error is practically zero with `DirWinFitAnal`, as it uses the analytical shape for the amplitude correction. All the other methods use a lookup table here. This test signal as well reveals the limited resolution of the `DirWinFitTbl` method, as the frequency error does not converge to zero.

With the `SiStSiHN` signals with high noise level, the methods `SpcReass`, `DrvAlg` and `PhsVoc` perform equally well or slightly better than the `DirWinFit` algorithms. The former seem to be slightly more stable against the influence of noise.

In the test for the spectral resolution `TwoStSi` the detection rate of the `DirWinFit` algorithm is similar to the other ones that include a shape criteria. The accuracy is significantly better here, though. Figure 7.6 gives an overview of the detected peaks of the two methods with both the sinusoidal and the Hann window.

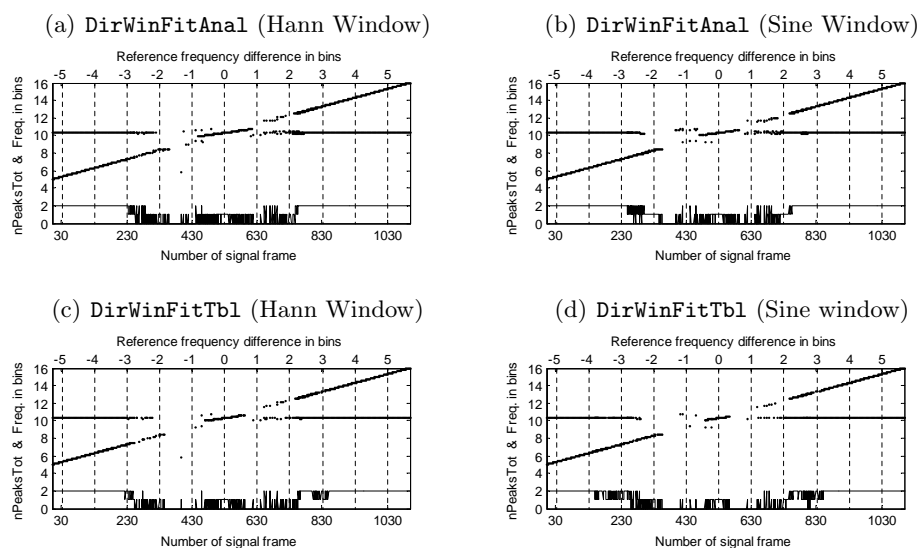


Figure 7.6.: Detected Peaks for pairs of closely spaced stationary partials in dependence of their frequency difference. Bottom lines show the total number of detected peaks per frame.

With the `Sweep` and `Tremolo` signals the performance is best with the `DirWinFit` methods. With `Vibrato`, `TremVib` and the musical instrument sounds, the performance is comparable to the best of the other methods tested before.

With the use of the sinusoidal window (square root of a shifted Hann window) for all

methods the `OdfitAlg` came out ahead for the musical instrument sounds before. This is not the case anymore compared to the `DirWinFit` methods. They provide highest accuracy combined with good detection rates. As well with all the other test signals the `DirWinFit` methods provide comparable or significantly better results.

The degradation of the `DirWinFirtbl` method against `DirWinFitAnal` is noticeable but not big, and depends on the used test signal. The level of impairment is traceable by the resolution of the lookup table. It is a direct trade-off between accuracy and computational effort and can so be easily scaled according to the application.

Similarly to figure 7.1 on page 76 we plot in figure 7.7 the absolute value of the the frequency error of the estimation of single stationary sinusoids. The integer frequency of the input sinusoid is varied in 50 steps in between $l_0 = 3$ and $l_0 = 52$ bins (with $N = 1024$). The fractional frequency Δl is varied in 200 steps from $\Delta l = 0$ to $\Delta l = 1$. This results in a total number of 8000 analysed sinusoids, which were analysed by `DirWinFitAnal` using a periodic Hann window.

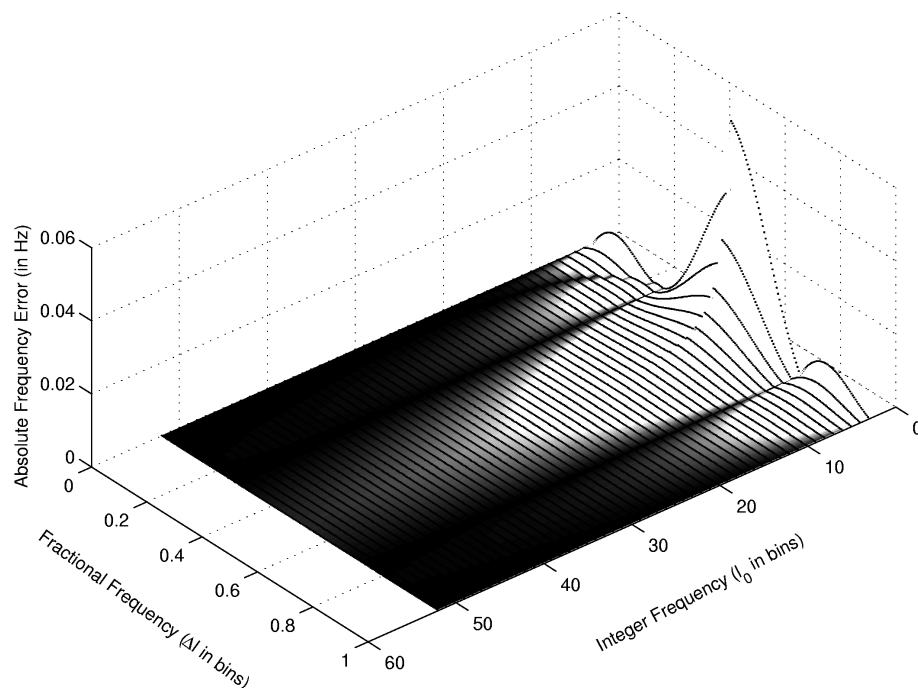


Figure 7.7.: Frequency error of `DirWinFitAnal` for stationary sinusoids ($N = 1024$, $l_{0,min} = 3$, Hann window).

When comparing this figure to the ones before, it is important to take into consideration that the frequency error is plotted for much lower frequencies, starting from $l_0 = 3$ instead of $l_0 = 12$. For the same frequencies the error is much smaller than before. The discontinuity at the fractional frequency $\Delta l = 0.5$ results from the change of the position of the local maximum from one bin to the next higher one.

Although the error is on a much lower level, we have basically the same situation again as shown in section 7.1: The estimation error is biggest in-between two frequency bins ($\Delta l =$

0.5) and is caused by the influence of the negative mirror peaks. We therefore could again simultaneously perform an ODFT analysis with direct fitting of the window and use its results for the values with a fractional frequency within $0.25 \leq \Delta l \leq 0.75$. Due to time restrictions this experiment was not conducted anymore, but an improvement in the frequency error estimation for low frequencies similar to the one shown with figure 7.3 is expected.

7.3. Removal of Mirror Frequencies

In section 7.1 we showed an effective way to reduce the impairment of the estimation induced by the negative mirror frequencies. Therefore we applied a combined analysis based on both the DFT and the ODFT.

We here want to discuss another approach which uses the concept of *analytic signals*: A signal is called ‘analytic’ if its Fourier transform is zero for negative frequencies. For a real signal $s(t)$ let us denote its corresponding analytic signal $s_a(t)$ by:

$$s_a(t) = s(t) + js_h(t) \quad (7.9)$$

This is a complex signal. Its real part is the original signal $s(t)$. The imaginary part $s_h(t)$ is $s(t)$ phase shifted by $\frac{\pi}{2}$. For e. g. a cosine signal $s(t) = \cos(\omega t)$ the corresponding analytic signal $s_a(t)$ is given by Euler’s equation:

$$s_a(t) = e^{j\omega t} = \cos(\omega t) + j \sin(\omega t) \quad (7.10)$$

In order to verify that the frequency error in figure 7.7 on the preceding page is indeed caused by negative mirror frequencies, we perform the same test but using the analytic signals of the same frequencies and amplitudes as complex input signals. As they have by definition no negative frequencies, the frequency error shown in that figure should vanish. According to equation 7.10 we therefore replace the $a \sin(\dots)$ terms in the test signal computations by $\frac{a}{2} e^{j(\dots)}$.

The resulting frequency error fully confirms the assumption: it is very small and constant over the integer frequency l_0 . The absolute value of the frequency error is around 10^{-7} Hz only.

Of course we can easily create analytic signals for test signals with *known* frequency content. But the computation of the analytic signal of an arbitrary real input signal is not as straight forward.

The required phase-shifted signal $s_h(t)$ is in theory obtained by applying the *Hilbert transform* on $s(t)$. A concise treatment of the Hilbert transform is given in [29]. The continuous-time transform equation is given by [39]:

$$s_h(t) = \frac{1}{\pi} \int_{\tau=-\infty}^{\infty} \frac{s(\tau)}{t - \tau} d\tau \quad (7.11)$$

We performed tests with the most common practical discrete-time Hilbert transformers. They did not improve the estimation performance. Numerous research was undertaken in the past about the computation of the analytic signal of a given real-valued signal [39, 41, 20].

Mainly two approaches are common for practical Hilbert transformers, based on the discrete Fourier transform or based on filter design techniques.

To obtain the analytic signal $s_a(n)$ of a signal $s(n)$, the discrete Fourier transform based methods compute the FFT of $s(n)$. In the resulting spectrum the coefficients which belong to negative frequencies are replaced with zeros. The result is subject to an inverse FFT, which gives the analytic signal $s_a(n)$.

But this is an approximation only. Zeroing the negative coefficients removes the signal components with negative frequencies in the corresponding negative frequency range. But it does not remove the contributions of the negative frequencies in the positive frequency range. They are caused by the side lobes of the analysis window when performing the DFT. These contributions become bigger for lower frequencies and they are exactly the disturbing factor that we initially wanted to address. As a consequence it is not surprising that with this kind of Hilbert transformers no improvement for the parameter estimation at low frequencies is achieved.

The MATLAB `hilbert()` function is based on this principle. The analytic signal $s_a(n)$ of a discrete-time sinusoid $s(n) = \sin(2\pi 0.0022n)$ is computed. The real part of $s_a(n)$ is the signal $s(n)$ itself. The imaginary part and its computed approximation are plotted in figure 7.8.

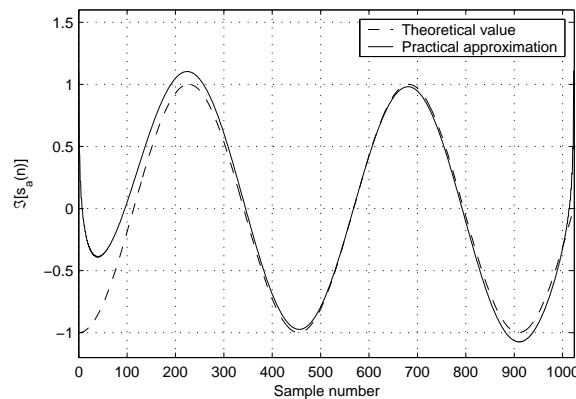


Figure 7.8.: Output of a FFT based Hilbert transformer: Theoretical value and approximation of the imaginary part of the analytic signal of a pure sinusoid.

The second approach regards the Hilbert transformer as a 90-degree phase shifter and applies standard filter design techniques to obtain an appropriate filter. The demands are similar to those of an ideal low-pass filter or an ideal band-limited differentiator. As known from digital filter theory, in practice only approximations are possible.

The most commonly used method is the Parks-McClellan algorithm. It provides FIR Hilbert transformer approximations which have a phase shift of exactly 90 degrees and an equiripple magnitude approximation. IIR filter realizations, which have beside a magnitude response error as well a phase response error, often consist of two all-pass systems. The phase responses of the all-pass filters are designed so that they differ by 90 degrees. Such a system is often referred to as a ‘phase splitter’. [29]

Figure 7.9 on the following page shows a typical example of a Hilbert transformer magnitude plot. There is a certain transition zone between the stop-band and the pass-band. Furthermore the magnitude response shows some ripple in the pass-band. The figure results from a FIR Hilbert transformer designed with the complex Remez algorithm as described in

[20]. The length of the complex impulse response is 22, the passband ripple is 0.74 dB and the stopband attenuation is 21 dB.

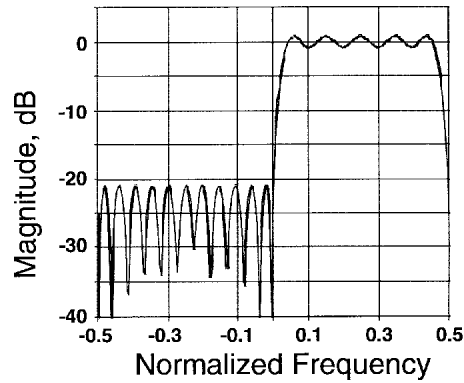


Figure 7.9.: Magnitude response of a one-sided FIR Hilbert transformer [20].

The figure reveals, why as well these kinds of Hilbert transform approximations do not improve parameter estimation. The undesired effect of negative mirror frequencies mainly shows up at very low (or very high) frequencies. But just there is as well the transition zone of the Hilbert transform filter, which prevents accurate estimation results in this frequency range. Furthermore the pass-band magnitude ripple impairs the estimation over the whole frequency range.

In conclusion the theoretical concept of the Hilbert transform looks very appealing. In practice only approximations of the analytic signal can be obtained. Unfortunately the specific nature of the approximation error prevents an improvement of the frequency or amplitude estimation.

Chapter 8.

Summary

High accuracy when extracting sinusoid parameters from harmonic sounds is an important key element in the analysis of sounds for additive synthesis. This report gives an comprehensive overview on the most common DFT based parameter estimators and evaluates their estimation results by means of synthesized test signals with known frequency content. The best and very similar results achieve the estimators which take in one or the other way into account the linear phase evolution of a stationary sinusoid, namely the derivative algorithm, the spectrum reassignment and the phase vocoder methods.

To what extent the numerical differences in the estimation errors have impact on the perceived quality, was not part of this thesis. On the one hand this effect depends not only on the choice of the analysis parameters, but as well highly on the other elements of the analysis stage. This includes in particular the peak tracking algorithm and its compatibility to the parameter estimation stage. Additionally this would imply the accomplishment of listening tests.

Nevertheless the software implementation provides a good device to compare new frame-based parameter estimators to the common ones. New methods that are in development can be easily included into the frame work, and they will get tested with the same signals. By this the effectiveness of new approaches for common test signals becomes immediately visible.

Based on this framework we investigated methods based on shifted transforms. In order to understand the functional principles, we showed the close relations between the DFT, ODFT and MDCT. The ODFT estimation algorithm provides good results and might be beneficial for the psycho-acoustic model in MDCT based audio coding. The MDCT-based algorithm was extended to detect multiple sinusoids and provides as the only method the possibility to extract partials directly from MDCT coefficients. Furthermore we were able to improve the parabolic interpolation method by exploiting different distributions of the estimation error when using DFT and ODFT.

In this report additionally a new approach was developed, which aims to pursue optimum DFT magnitude interpolation without loosing the possibility to freely choose the analysis window. This was achieved by directly fitting the analytic expression of the analysis window main lobe into the DFT spectrum. The approach was confirmed by including it into the above framework: it outperforms the other estimators. At the moment the high computational costs of this method prevent its use in standard real-time environments. But there are many applications which do not require analysis in real-time and so can benefit from the outstanding accuracy of this approach (e.g. for the extraction of reference data for additive synthesizers).

Further research might include initial phase estimation into the survey, improve the stability of the MDCT algorithm against noise and closely spaced neighbouring peaks, and generalize the combined DFT/ODFT approach for use with other estimators.

A. Appendix

A.1. Software Implementation

All methods and tests as described in chapters 6 and 7 were implemented in Mathworks' Matlab. The code was written in Matlab 6.5 R13 and exemplarily tested on Windows XP, Redhat Linux and Apple OS-X platforms.

During the implementation importance was attached to the basic principles of software engineering and coding style. All computations are highly vectorized. The source code is well documented, the input and output variables are clearly specified, and many additional explanations are given where appropriate. The following naming scheme was consistently used:

- **Functions** are named and written like `'FunctionName(...).'`
- **Files** are named according to the contained function `'FunctionName.m'`.
- **Constants** are written like `'CONSTANT'`.
- **Scalars** are written like `'scalar_variable'`.
- **Vectors and matrices** are written like `'VectorOrMatrix'`.

The entire program design is highly modular and hierarchically partitioned: All units and analysis methods are implemented in separate functions, which are themselves reusable in other programs. Furthermore high attention was paid to the fact that it is easy to integrate further analysis methods, test signals and analysis windows. Only a very small number of files have to be modified therefore. Simple step-by-step instructions on how to do this are enclosed with the source code.

Program Control

The file `'PeakExtraction.m'` is the central function which controls the simulations and includes all parameter settings. No parameters are passed to this function, it is supposed to be directly and manually edited. This function controls the complete data flow, calls the subsidiary functions and passes the appropriate parameters to the analysis methods.

The file starts with a section for enabling or disabling certain analysis methods or test signals for the comparisons by commenting out the corresponding name of the method or test signal. As some methods require rather long computation time, it is recommendable practice not to compute all methods and test signals at a time as long as it is not necessary. Furthermore the plotting of some illustrative figures can be enabled or disabled.

It follows a long list of parameters for the analysis methods. Beside the parameters N and f_s , which are common to all methods, the parameters are saved in the structure `'APrm'`. Its

field designators are named by the scheme ‘APrm.pf_AM’, with ‘pf’ being the functionality of the parameter and ‘AM’ being the analysis method the parameter belongs to. The names that are used for the analysis methods, test signals and analysis windows are defined in the files ‘ErrorResNames.m’ and ‘WindowNames.m’, which are applied as ‘C’-style include-files in most of the source files.

The consecutive section of the main function computes the test signals. Each signal ‘TestSigs’ is passed to the sub-function ‘DoAnalysis()’. In there the test signals are passed on to the parameter estimators. These are enclosed in functions (files) that are named by the parameter estimator name with a prefix ‘PE_’. A wrapper between the matrix- and vector-based signal representation is demonstrated in the ‘PE_Trg1Alg.m’ file, as the implementation for the triangle algorithm was taken over from a vector-based framework. The detected peaks ‘EstdPeaks’ are reviewed according to section 6.3.

From the remaining peaks ‘FiltPeaks’ the estimation errors are measured and stored in the multi-dimensional array ‘ErrorRes’. In this array all information about the estimation results for all analysis methods and test signals is collected. After all tests are carried out, the ‘ErrorRes’ array is converted to a human-readable format and saved in the sub-folder ‘results’ into the plain-text file ‘ErrorRes.txt’. The file is overwritten every time the program is executed.

As many Matlab functions as e.g. ‘fft()’ operate not only on vectors but as well column-wise on matrices, the signal representation is chosen accordingly: the framed signal is stored in a two-dimensional matrix. Every column of the matrix represents one signal frame. As a consequence the number of rows of the matrix equals the frame length, and the number of columns equals the number of frames the signal consists of. All signal and estimation data is passed between the functions in similarly structured matrices.

A.2. Simulation Results

The following tables show the full simulation results as discussed in chapter 6. The amplitude error is given in dB in terms of power, the frequency error in cents (see section 6.4).

In table A.1 on the next page for every method the window which gives best results in terms of estimation accuracy is used. The periodic Hann window was applied for table A.2 on page 96. With the sine window used for every method, the results are given in table A.3 on page 102.

TestSignal: TremVib		Freq		My		Sigma		Maxx		Minn		Amp1		Sigma		Maxx		Minn		AmplMissed		nPhantom		AmplPhantom		AmplMissed		nPhantom		AmplPhantom						
		My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn			
PlainFFT	:	35.470587	24.299136	69.777842	0.801543	0.412070	0.397427	1.112925	0.000105																											
ParInt	:	1.656514	0.809778	2.746516	0.001766	0.016680	0.009863	0.033984	0.000038																											
TrglAlg	:	1.548640	0.776981	2.794543	0.009823	0.015060	0.008725	0.029251	0.000025																											
SpkReass	:	1.332840	0.670569	2.359588	0.001397	0.014120	0.009568	0.027343	0.000008																											
DrvAlg	:	1.337831	0.674472	2.466520	0.002297	0.014421	0.009261	0.027211	0.000007																											
PhvVoc	:	1.333235	0.670749	2.360519	0.001981	0.014112	0.009597	0.027452	0.000012																											
OdfAlg	:	1.680363	1.568248	7.724613	0.006531	0.087315	0.082375	0.319208	0.000042																											
MdctAlg	:	5.653225	5.970915	50.049164	0.006368	0.177485	0.181478	1.673836	0.000393																											
OdfParInt	:	1.633689	0.799419	2.846584	0.006686	0.017459	0.010606	0.035206	0.000091																											
DirWinFitTbl	:	1.593986	0.781487	2.636870	0.005897	0.017857	0.009393	0.033348	0.000007																											
DirWinFitAhal	:	1.593952	0.781887	2.621177	0.000220	0.014258	0.009216	0.029006	0.000026																											
DualParInt	:	1.702649	0.831839	2.746516	0.001766	0.016525	0.010344	0.033984	0.000009																											
nPeaksTot		My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn			
PlainFFT	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00			
ParInt	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00			
TrglAlg	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
SpkReass	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
DrvAlg	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
PhvVoc	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
OdfAlg	:	1.02	0.12	2.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
MdctAlg	:	0.91	0.30	2.00	0.00	0.90	0.30	1.00	0.00	0.10	0.30	1.00	0.00	0.00	0.01	0.08	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
OdfParInt	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
DirWinFitTbl	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
DirWinFitAhal	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
DualParInt	:	1.00	0.00	1.00	1.00	1.00	1.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
nPeaksCorr		My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn			
PlainFFT	:	12.509649	22.244322	209.493584	0.000651	0.861879	0.691733	2.999928	0.000027																											
ParInt	:	4.421437	5.405019	65.790663	0.000126	0.749896	0.672668	2.999938	0.000091																											
TrglAlg	:	4.321297	5.275603	68.453666	0.000085	0.746076	0.667470	2.999974	0.000064																											
SpkReass	:	4.591920	5.458968	73.242667	0.000281	0.767413	0.682621	2.999488	0.000016																											
DrvAlg	:	4.581522	5.437991	74.303472	0.000088	0.768416	0.683233	2.999375	0.000091																											
PhvVoc	:	4.592742	5.464721	73.354253	0.000762	0.768632	0.682107	2.999793	0.000071																											
OdfAlg	:	4.619156	5.611506	76.657661	0.000023	0.796307	0.694033	2.997290	0.000005																											
MdctAlg	:	5.126278	7.653191	185.667823	0.000098	0.825206	0.707105	2.998839	0.000021																											
OdfParInt	:	4.413025	5.400914	70.832716	0.000129	0.748597	0.670251	2.999648	0.000061																											
DirWinFitTbl	:	4.273951	5.243380	65.690501	0.000031	0.748855	0.670982	2.997902	0.000016																											
DirWinFitAhal	:	4.277313	5.239265	65.691437	0.000116	0.747108	0.671051	2.999650	0.000012																											
DualParInt	:	4.384800	5.387922	66.597193	0.000129	0.753079	0.671457	2.999938	0.000006																											
nPeaksTot		My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn	My	Sigma	Maxx	Minn			
PlainFFT	:	38.39	9.59	67.00	9.00	31.77	11.71	50.00	0.00	8.30	10.71	46.00	0.00	6.63	10.38	58.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
ParInt	:	33.27	8.94	49.00	6.00	29.91	11.07	49.00	0.00	10.16	10.67	48																								

TestSignal: TwoSt5i												*-----*												
												* Amplitude errors in dB												
												* Frequency errors in cent												
												* n... in total number per frame												
												* Amplitudes in dB												

My	Sigma	Maxx	Minn	My	Ampl	Minn	My	Sigma	Maxx	Minn	My	Ampl	Minn	My	Sigma	Maxx	Minn	My	Ampl	Minn	My	Sigma	Maxx	Minn
49.002819	26.198536	180.614694	0.000000	0.782077	0.543133	2.982229	0.000502	0.43133	2.982229	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301
8.600256	11.035216	143.186898	0.004534	0.266673	0.363064	2.988229	0.000301	0.266673	2.988229	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301
22.783207	19.019082	175.621279	0.025939	0.386205	0.285397	2.942853	0.000833	0.285397	2.942853	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833	0.000833
11.742244	23.280282	175.176710	0.003310	0.302424	0.517921	2.999717	0.000301	0.302424	2.999717	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301	0.000301
11.786515	23.489536	177.242927	0.009661	0.302312	0.517704	2.983283	0.000201	0.302312	2.983283	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201	0.000201
10.615044	20.043123	175.286759	0.001694	0.287246	0.491989	2.955161	0.000156	0.287246	2.955161	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156	0.000156
12.710696	24.570374	163.935601	0.000307	0.310189	0.502392	2.995161	0.000541	0.310189	2.995161	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541	0.000541
10.257814	24.361901	187.787355	0.007251	0.350283	0.580596	2.992469	0.000691	0.350283	2.992469	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691	0.000691
8.068892	11.644908	175.897566	0.039197	0.219562	0.328416	2.985445	0.000162	0.219562	2.985445	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162
2.811143	8.773305	146.196331	0.000067	0.082070	0.307669	2.904509	0.000001	0.082070	2.904509	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001
2.777992	8.122833	146.182471	0.002511	0.102626	0.321136	2.959338	0.000051	0.102626	2.959338	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051	0.000051
10.135109	16.711933	139.955776	0.021842	0.226678	0.369280	2.988545	0.000162	0.226678	2.988545	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162	0.000162
My	Sigma	Maxx	Minn	My	nPeaks	Maxx	Minn	My	nMissed	Maxx	Minn	My	nPhantom	Maxx	Minn	My	nMissed	Maxx	Minn	My	nPhantom	Maxx	Minn	
1.74	0.47	3.00	1.00	1.61	0.67	2.00	0.00	0.39	0.67	2.00	0.00	0.13	0.37	2.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.20	0.76	2.00	0.00	1.12	0.82	2.00	0.00	0.88	0.82	2.00	0.00	0.07	0.26	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.30	0.81	2.00	0.00	1.23	0.88	2.00	0.00	0.77	0.88	2.00	0.00	0.07	0.26	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.73	0.44	2.00	1.00	1.62	0.66	2.00	0.00	0.38	0.66	2.00	0.00	0.11	0.32	2.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.73	0.44	2.00	1.00	1.62	0.66	2.00	0.00	0.38	0.66	2.00	0.00	0.11	0.32	2.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.71	0.45	2.00	1.00	1.60	0.67	2.00	0.00	0.40	0.67	2.00	0.00	0.10	0.31	2.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.76	0.46	3.00	1.00	1.61	0.67	2.00	0.00	0.39	0.67	2.00	0.00	0.13	0.36	2.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.11	0.70	3.00	0.00	0.82	0.86	2.00	0.00	1.18	0.86	2.00	0.00	0.29	0.46	2.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.26	0.75	2.00	0.00	1.17	0.83	2.00	0.00	0.83	0.83	2.00	0.00	0.09	0.29	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.36	0.82	2.00	0.00	1.30	0.88	2.00	0.00	0.70	0.88	2.00	0.00	0.06	0.24	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.36	0.82	2.00	0.00	1.30	0.88	2.00	0.00	0.70	0.88	2.00	0.00	0.06	0.24	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.18	0.76	2.00	0.00	1.11	0.82	2.00	0.00	0.89	0.82	2.00	0.00	0.07	0.26	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
My	Sigma	Maxx	Minn	My	nPeaks	Maxx	Minn	My	nMissed	Maxx	Minn	My	nPhantom	Maxx	Minn	My	nMissed	Maxx	Minn	My	nPhantom	Maxx	Minn	
29.311666	18.327565	80.302633	0.011005	0.693323	0.619807	2.098337	0.000621	0.693323	2.098337	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	0.000621	
3.481647	1.784246	7.351360	0.000610	0.120622	0.091789	0.264282	0.000040	0.120622	0.264282	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	
15.507237	8.238485	38.523935	1.848960	0.286359	0.120776	0.507945	0.001007	0.286359	0.507945	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	0.001007	
0.424165	0.345383	1.762461	0.000191	0.011456	0.068584	0.038017	0.000022	0.011456	0.068584	0.038017	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	0.000022	
0.429879	0.358003	1.875320	0.001393	0.011604	0.068846	0.040561	0.000001	0.011604	0.068846	0.040561	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	
0.438189	0.352846	1.872825	0.000943	0.012358	0.069164	0.041831	0.000001	0.012358	0.069164	0.041831	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	
0.899852	0.733463	5.313133	0.000189	0.056877	0.047199	0.299715	0.000116	0.056877	0.299715	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	0.000116	
1.425574	2.295992	29.386172	0.000071	0.098121	0.103621	1.513920	0.000012	0.098121	1.513920	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	0.000012	
3.476204	1.733798	7.010678	0.004592	0.118774	0.090716	0.263754	0.000027	0.118774	0.263754	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	0.000027	
0.028800	0.034624	0.265532	0.000061	0.000347	0.000214	0.001386	0.000001	0.000347	0.000214	0.001386	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	0.000001	
0.027198	0.034964	0.263969	0.000017	0.002926	0.002398	0.015393	0.000000	0.002926	0.015393	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
3.233192	1.589132	6.986249	0.000610	0.090089	0.064538	0.223297	0.000040	0.090089	0.223297	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	0.000040	
My	Sigma	Maxx	Minn	My	nPeaks	Maxx	Minn	My	nMissed	Maxx	Minn	My	nPhantom	Maxx	Minn	My	nMissed	Maxx	Minn	My	nPhantom	Maxx	Minn	
1.00	0.00	1.00	1.00	1.00	0.00	1.00	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
0.72	0.45	1.00	0.00	0.72	0.45	1.00	0.00	0.28	0.45	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
0.92	0.27	1.00	0.00	0.92	0.27	1.00	0.00	0.08	0.27	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.00	0.00	1.00	1.00	1.00	0.00	1.00	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
1.00	0.00	1.00	1.00	1.00	0.00	1.00	0.00	1.00	0.00	1.00	0.00	0.00	0.00											

TestSignal: TremVib		Freq		Amp1		Sigma		Minn		My		nPhantom		Amplitude errors in dB		
My	Sigma	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	
PlainFFT	24.299136	69.777842	0.801543	0.595428	0.801351	0.000135	0.580840	1.618743	0.000329	0.000135	0.580840	1.618743	0.000329	0.000135	0.580840	
ParInt	2.715504	8.544285	0.003602	0.089198	0.003329	0.000329	0.087196	0.258491	0.000329	0.000329	0.087196	0.258491	0.000329	0.000329	0.087196	
TrglAlg	14.103799	45.579011	2.823835	0.311101	0.311101	0.000875	0.126454	0.522267	0.000875	0.000875	0.126454	0.522267	0.000875	0.000875	0.126454	
SpkReass	1.371748	3.659593	0.000464	0.019977	0.000201	0.000201	0.014101	0.067075	0.000201	0.000201	0.014101	0.067075	0.000201	0.000201	0.014101	
DrvAlg	1.394376	3.755944	0.012557	0.021019	0.015306	0.000167	0.015115	0.064232	0.000167	0.000167	0.015115	0.064232	0.000167	0.000167	0.015115	
PhvVoc	1.388448	3.655570	0.000120	0.020896	0.000046	0.000046	0.015115	0.064232	0.000046	0.000046	0.015115	0.064232	0.000046	0.000046	0.015115	
OdfAlg	1.680363	1.568248	0.000531	0.087315	0.082375	0.000042	0.087315	0.319208	0.000042	0.000042	0.087315	0.319208	0.000042	0.000042	0.087315	
MdctAlg	5.653225	5.970915	0.049164	0.177485	0.167383	0.000393	0.181478	1.673836	0.000393	0.000393	0.181478	1.673836	0.000393	0.000393	0.181478	
OdfParInt	2.715196	7.614202	0.002957	0.101188	0.082910	0.004451	0.082910	0.259966	0.004451	0.004451	0.082910	0.259966	0.004451	0.004451	0.082910	
DirWinFitTbl	1.405722	0.712649	2.447119	0.001491	0.015092	0.000008	0.007498	0.029078	0.000008	0.000008	0.007498	0.029078	0.000008	0.000008	0.007498	
DirWinFitAhal	1.404955	0.713215	2.432359	0.001041	0.015121	0.000007	0.009079	0.031821	0.000007	0.000007	0.009079	0.031821	0.000007	0.000007	0.009079	
DualParInt	3.321926	1.904018	7.221786	0.003602	0.061074	0.000329	0.054066	0.212807	0.000329	0.000329	0.054066	0.212807	0.000329	0.000329	0.054066	
nPeaksTot		nPeaksCorr		nMissed		nPhantom		Amplitude errors in dB		Amplitude errors in dB		Amplitude errors in dB		Amplitude errors in dB		
My	Sigma	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	
1.00	0.00	1.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.74	0.44	1.00	0.74	0.44	1.00	0.00	0.26	0.44	1.00	0.00	0.26	0.44	1.00	0.00	0.26	
0.95	0.22	1.00	0.95	0.22	1.00	0.00	0.05	0.22	1.00	0.00	0.05	0.22	1.00	0.00	0.05	
1.00	0.00	1.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
1.00	0.00	1.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
1.00	0.00	1.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
1.02	0.12	2.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.91	0.30	2.00	0.90	0.30	1.00	0.00	0.10	0.30	1.00	0.00	0.10	0.30	1.00	0.00	0.10	
0.64	0.48	1.00	0.64	0.48	1.00	0.00	0.36	0.48	1.00	0.00	0.36	0.48	1.00	0.00	0.36	
1.00	0.00	1.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
1.00	0.00	1.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
0.74	0.44	1.00	0.74	0.44	1.00	0.00	0.26	0.44	1.00	0.00	0.26	0.44	1.00	0.00	0.26	
TestSignal: Voice		Freq		Amp1		Sigma		Minn		My		nPhantom		Amplitude errors in dB		
My	Sigma	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	
12.574771	22.327116	205.968019	0.000651	0.974277	0.733619	0.000043	2.999824	2.999824	0.000043	0.000043	2.999824	2.999824	0.000043	0.000043	2.999824	
4.942994	5.639846	63.395330	0.000436	0.758696	0.656813	0.000143	2.999049	2.999049	0.000143	0.000143	2.999049	2.999049	0.000143	0.000143	2.999049	
9.156681	14.615254	169.409113	0.000228	0.760022	0.616022	0.000021	2.988490	2.988490	0.000021	0.000021	2.988490	2.988490	0.000021	0.000021	2.988490	
6.904614	8.336189	133.467688	0.000347	0.908627	0.747771	0.000116	2.999708	2.999708	0.000116	0.000116	2.999708	2.999708	0.000116	0.000116	2.999708	
6.999642	8.408620	135.224321	0.000541	0.910353	0.749858	0.000029	2.999809	2.999809	0.000029	0.000029	2.999809	2.999809	0.000029	0.000029	2.999809	
6.890819	8.374957	134.331125	0.000336	0.874695	0.723886	0.000083	2.999544	2.999544	0.000083	0.000083	2.999544	2.999544	0.000083	0.000083	2.999544	
4.619156	5.611506	76.657661	0.000023	0.796307	0.694033	0.000005	2.997290	2.997290	0.000005	0.000005	2.997290	2.997290	0.000005	0.000005	2.997290	
5.126278	7.653191	185.667823	0.000098	0.825206	0.707105	0.000002	2.988839	2.988839	0.000002	0.000002	2.988839	2.988839	0.000002	0.000002	2.988839	
4.947122	5.785691	75.597126	0.000481	0.761322	0.658104	0.000001	2.998235	2.998235	0.000001	0.000001	2.998235	2.998235	0.000001	0.000001	2.998235	
3.872924	4.681709	64.852113	0.000143	0.702397	0.638818	0.000045	2.999529	2.999529	0.000045	0.000045	2.999529	2.999529	0.000045	0.000045	2.999529	
3.874161	4.670306	64.893887	0.000103	0.702109	0.635526	0.000013	2.999248	2.999248	0.000013	0.000013	2.999248	2.999248	0.000013	0.000013	2.999248	
4.909277	5.709538	65.660401	0.000436	0.756309	0.656035	0.000143	2.998235	2.998235	0.000143	0.000143	2.998235	2.998235	0.000143	0.000143	2.998235	
nPeaksTot		nPeaksCorr		nMissed		nPhantom		Amplitude errors in dB		Amplitude errors in dB		Amplitude errors in dB		Amplitude errors in dB		
My	Sigma	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	My	Maxx	Minn	
38.80	10.19	72.00	9.00	31.45	11.64	49.00	0.00	10.72	45.00	0.00	7.35	11.62	64.00	0.00	-69.03	
24.91	8.93	44.00	1.00	23.18	9.72	42.00	0.00	16.88	9.05	48.00	1.00	1.73	2.90	17.00	0.00	-61.76
23.26	8.98	43.00	0.00	22.21	9.31	43.00	0.00	17.85	11.19	50.00	1.00	1.04	2.10	15.00	0.00	-60.83
36.51	9.71	67.00	8.00	28.34	10.31	48.00	0.00	11.72	9.96	45.00	0.00	8.16	9.77	58.00	0.00	-67.10
35.99	9.42	65.00	8.00	27.70	9.89	47.00	0.00	12.36	9.81	45.00	0.00	8.29	9.25	56.00	0.00	-66.86
32.63	9.05	64.00	7.00	27.09	9.93	46.00	0.00	12.97	10.04	47.00	0.00	5.64	9.11	53.00	0.00	-66.02
38.82	10.20	70.00	9.00	32.09	11.69	50.00	0.00	7.97	10.48	45.00	0.00	6.73	11.40	60.00	0.00	-70.47
24.97	9.78	46.00	1.00	21.39	10.18	46.00	0.00	18.68	9.37	49.00	3.00	3.68	5.51	30.00	0.00	-59.29
24.52	9.16	45.00	1.00	22.83	9.78	44.00	0.00	17.23	8.60	51.00	4.00	1.69	2.79	18.00	0.00	-60.68
28.35	11.50	49.00	0.00	27.12	12.06	49.00	0.00	12.95	12.23	50.00	0.00	1.23	1.87	15.00	0.00	-64.35
28.35	11.50	49.00	0.00	27.14	12.07	49.00	0.00	12.25	12.25	50.00	0.00	1.21	1.85	15.00	0.00	-64.35
24.84	8.99	44.00	1.00	23.15	9.74	42.00	0.00	16.92	9.11	49.00	1.00	1.70	2.77	17.00	0.00	-61.75

Table A.3.: Full simulation results with sine window. (cont.)

Bibliography

- [1] N. Ahmed, T. Natarajan, and K. R. Rao. Discrete cosine transform. *IEEE Transactions on Computers*, January 1974.
- [2] Jont B. Allen and Lawrence R. Rabiner. A unified approach to short-time fourier analysis and synthesis. *Proceedings of the IEEE*, 65(11):1558–1564, November 1977.
- [3] Rasmus Althoff, Florian Keiler, and Udo Zölzer. Extracting sinusoids from harmonic signals. In *Proceedings of the 2nd COST G-6 Workshop on Digital Audio Effects (DAFx-99)*. DAFx-99, December 1999.
- [4] François Auger and Patrick Flandrin. The why and how of time-frequency reassignment. In *Proceedings of the IEEE-SP International Symposium on Time-Frequency and Time-Scale Analysis 1994*, pages 197–200. IEEE, October 1994.
- [5] François Auger and Patrick Flandrin. Improving the readability of time-frequency and time-scale representations by the reassignment method. *IEEE Transactions on Signal Processing*, 43(5):1068–1089, May 1995.
- [6] Laurent Daudet and Mark Sandler. MDCT analysis of sinusoids: Exact results and applications to coding artifacts reduction. *Submitted to IEEE Transactions on Acoustics, Speech, and Signal Processing*, 2002.
- [7] Laurent Daudet and Mark Sandler. MDCT analysis of sinusoids and applications to coding artifacts reduction. In *Preprints of the 114th Convention of the AES*. AES, March 2003.
- [8] J. E. Dennis. Nonlinear least-squares. In D. Jacobs, editor, *State of the Art in Numerical Analysis*, pages 269–312. Academic Press, 1977.
- [9] Philippe Depalle, Guillermo Garcia, and Xavier Rodet. Analysis of sound for additive synthesis: Tracking of partials using hidden markov models. In *Proceedings of the International Computer Music Conference (ICMC 93)*. ICMA, 1993.
- [10] Myriam Desainte-Catherine and Sylvain Marchand. High-precision fourier analysis of sounds using signal derivatives. *Journal of the Audio Engineering Society*, 48(7/8):654–667, July/August 2000.
- [11] Udo Zölzer et al. *DAFX - Digital Audio Effects*. John Wiley and Sons, 2003.
- [12] Aníbal J. S. Ferreira. *Spectral Coding and Post-Processing of High Quality Audio*. PhD thesis, University of Porto, 1998.

-
- [13] Aníbal J. S. Ferreira. An Odd-DFT based approach to time-scale expansion of audio signals. *IEEE Transactions on Speech and Audio Processing*, 7(4):441–453, July 1999.
- [14] Aníbal J. S. Ferreira. Accurate estimation in the ODFT domain of the frequency, phase and magnitude of stationary sinusoids. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pages 47–50. IEEE, October 2001.
- [15] Stephen W. Hainsworth and Patrick J. Wolfe. Time-frequency reassignment for music analysis. In *Proceedings of the International Computer Music Conference (ICMC 2001)*, pages 14–17. ICMA, 2001.
- [16] Frederic J. Harris. On the use of windows for harmonic analysis with the discrete fourier transform. *Proceedings of the IEEE*, 66(1):51–83, January 1978.
- [17] Florian Keiler and Sylvain Marchand. Survey on extraction of sinusoids in stationary sounds. In *Proceedings of the 5th International Conference on Digital Audio Effects (DAFx-02)*. DAFx-02, September 2002.
- [18] Florian Keiler and Udo Zölzer. Extracting sinusoids from harmonic signals. *Journal of New Music Research*, 30(3):243–258, 2001.
- [19] Kunihiro Kodera, Roger Gendrin, and Claude de Villedary. Analysis of time-varying signals with small BT values. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 26(1):64–76, February 1978.
- [20] Michael Z. Komodromos, Steve F. Russell, and Ping Tak Peter Tang. Design of FIR hilbert transformers and differentiators in the complex domain. *IEEE Transactions on Circuits and Systems–I: Fundamental Theory and Applications*, 45(1):64–67, January 1998.
- [21] E. Kurniawati, J. Absar, S. George, C. T. Lau, and B. Premkumar. Single transform perceptual audio encoder. In *Proceedings of the 14th International Conference on Digital Signal Processing (DSP 2002)*, volume 2 of 1-3, pages 599–602. IEEE, July 2002.
- [22] Scott Nathan Levine. *Audio Representations for Data Compression and Compressed Domain Processing*. Dissertation, Stanford University, 1998.
- [23] Mathew Manu, Bhat Vashuda, Thomas M. Shine, and Yim Changhoon. Modified MP3 encoder using complex modified discrete cosine transform. In *Proceedings of the International Conference on Multimedia & Expo (ICME 2003)*. IEEE, July 2003.
- [24] Sylvain Marchand. *Sound Models for Computer Music - Analysis, Transformation, Synthesis of Musical Sounds*. Phd-thesis, Université Bordeaux 1, December 2000.
- [25] The MathWorks. *Using MATLAB (Version 6.5 R13)*, August 2002.
- [26] S. Merdjani and Laurent Daudet. Direct estimation of frequency from MDCT-encoded files. In *Proceedings of the 6th International Conference on Digital Audio Effects (DAFx-03)*. DAFx-03, September 2003.

-
- [27] Vladimir Nikolajevic and Gerhard Fettweis. Computation of forward and inverse MDCT using clenshaw's recurrence formula. *IEEE Transactions on Signal Processing*, 51(5), May 2003.
- [28] Albert H. Nuttall. Some windows with very good sidelobe behaviour. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, ASSP-29(1):84–91, February 1981.
- [29] Alan V. Oppenheim and Ronald W. Schaffer. *Discrete-Time Signal Processing*. Prentice-Hall International, 1989.
- [30] Alan V. Oppenheim, Alan. S. Willsky, and Ian T. Young. *Signals and Systems*. Prentice-Hall International, 1983.
- [31] Ted Painter and Andreas Spanias. Perceptual coding of digital audio. *Proceedings of the IEEE*, 88(4):449–515, April 2000.
- [32] Alan Peevers. Real time audio signal processing on RISC workstations via short time fourier transforms. In *Preprints of the 97th Convention of the AES*. AES, November 1994.
- [33] F. Plante, G. Meyer, and W. A. Ainsworth. Improvement of speech spectrogram accuracy by the method of reassignment. *IEEE Transactions on Speech and Audio Processing*, 6(3):282–286, May 1998.
- [34] William H. Press, Brian P. Flannery, Saul A. Teukolsky, and William T. Vetterling. *Numerical Recipes in C: The Art Of Scientific Computing*. Cambridge University Press, 2nd edition, 1992.
- [35] John P. Princen and Alan Bernard Bradley. Analysis/synthesis filter bank design based on time domain aliasing cancellation. *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-34(5), October 1986.
- [36] John P. Princen, A. W. Johnson, and Alan Bernard Bradley. Subband/transform coding using filter bank designs based on time domain aliasing cancellation. In *Proceedings of the International Conference on Acoustics Speech and Signal Processing (ICASSP 87)*. IEEE, April 1987.
- [37] Miller S. Puckette and Judith C. Brown. Accuracy of frequency estimates using the phase vocoder. *IEEE Transactions on Speech and Audio Processing*, 6(2):166–176, March 1998.
- [38] Thomas F. Quatieri. *Discrete-Time Speech Signal Processing: Principles and Practice*. Prentice Hall PTR, 2002.
- [39] Ashwin Rao and Ramdas Kumaresan. A parametric modeling approach to hilbert transformation. *IEEE Signal Processing Letters*, 5(1):15–17, January 1998.
- [40] K. R. Rao and P. Yip. *Discrete Cosine Transform - Algorithms, Advantages, Applications*. Academic Press, 1990.
- [41] Andrew Reilly, Gordon Frazer, and Boualem Boashash. Analytic signal generation – tips and traps. *IEEE Transactions on Signal Processing*, 42(11):3241–3245, November 1994.

-
- [42] Xavier Serra. Musical sound modeling with sinusoids plus noise. In C. Roads, S. Pope, A. Picialli, and G. de Poli, editors, *Musical Signal Processing*. Swets & Zeitlinger Publishers, 1997.
- [43] Julius O. Smith. *Draft: Mathematics of the Discrete Fourier Transform (DFT)*. <http://www-ccrma.stanford.edu/~jos/mdft/>, January 2002.
- [44] Julius O. Smith and Xavier Serra. PARSHL: an analysis/synthesis program for non-harmonic sounds based on a sinusoidal representation. Superset of ICMC-87 conference paper, Stanford University, 1994.
- [45] Gilbert Strang. The discrete cosine transform. *Society for Industrial and Applied Mathematics SIAM Review*, 41(1):135–147, 1999.
- [46] Ye Wang and Miikka Vilermo. Modified discrete cosine transform – its implications for audio coding and error concealment. *Journal of the Audio Engineering Society*, 51(1/2):52–61, January/February 2003.
- [47] Ye Wang, Miikka Vilermo, and Leonid Yaroslavsky. Energy compaction property of the MDCT in comparison with other transforms. In *Preprints of the 109th Convention of the AES*. AES, September 2000.
- [48] Ye Wang, Leonid Yaroslavsky, and Miikka Vilermo. On the relationship between MDCT, SDFT and DFT. In *Proceedings of the 5th International Conference on Signal Processing (ICSP 2000)*. IEEE, 2000.
- [49] Ye Wang, Leonid Yaroslavsky, Miikka Vilermo, and Mauri Väänänen. Some peculiar properties of the MDCT. In *Proceedings of the 5th International Conference on Signal Processing (ICSP 2000)*. IEEE, 2000.
- [50] Wikipedia: The Free Encyclopedia. Discrete cosine transform (06:57 utc, 16 aug 2003). Retrieved from http://en.wikipedia.org/wiki/Discrete_cosine_transform, August 2003.
- [51] Wikipedia: The Free Encyclopedia. Discrete fourier transform (01:40 utc, 12 sep 2003). Retrieved from http://en.wikipedia.org/wiki/Discrete_Fourier_transform, September 2003.
- [52] Leonid Yaroslavsky and Murray Eden. *Fundamentals of Digital Optics: Digital Signal Processing in Optics and Holography*. Birkhäuser, 1996.
- [53] Eberhard Zwicker and H. Fastl. *Psychoacoustics - Facts and Models*. Springer Verlag, Heidelberg, 1990.