

# Musical elements in the discrete-time representation of sound

RENATO FABBRI, University of São Paulo

VILSON VIEIRA DA SILVA JUNIOR, Cod.ai

ANTÔNIO CARLOS SILVANO PESSOTTI, Universidade Metodista de Piracicaba

DÉBORA CRISTINA CORRÊA, University of Western Australia

OSVALDO N. OLIVEIRA JR., University of São Paulo

---

The representation of basic elements of music in terms of discrete audio signals is often used in software for musical creation and design. Nevertheless, there is no unified approach that relates these elements to the discrete samples of digitized sound. In this article, each musical element is related by equations and algorithms to the discrete-time samples of sounds, and each of these relations are implemented in scripts within a software toolbox, referred to as MASS (Music and Audio in Sample Sequences). The fundamental element, the musical note with duration, volume, pitch and timbre, is related quantitatively to characteristics of the digital signal. Internal variations of a note, such as tremolos, vibratos and spectral fluctuations, are also considered, which enables the synthesis of notes inspired by real instruments and new sonorities. With this representation of notes, resources are provided for the generation of higher level musical structures, such as rhythmic meter, pitch intervals and cycles. This framework enables precise and trustful scientific experiments, data sonification and is useful for education and art. The efficacy of MASS is confirmed by the synthesis of small musical pieces using basic notes, elaborated notes and notes in music, which reflects the organization of the toolbox and thus of this article. It is possible to synthesize whole albums through collage of the scripts and settings specified by the user. With the open source paradigm, the toolbox can be promptly scrutinized, expanded in co-authorship processes and used with freedom by musicians, engineers and other interested parties. In fact, MASS has already been employed for diverse purposes which include music production, artistic presentations, psychoacoustic experiments and computer language diffusion where the appeal of audiovisual artifacts is exploited for education.

CCS Concepts: •**Applied computing** →**Sound and music computing**; •**Computing methodologies** →*Modeling methodologies*; •**General and reference** →*Surveys and overviews*; *Reference works*;

Additional Key Words and Phrases: music, acoustics, psychophysics, digital audio, signal processing

## ACM Reference format:

Renato Fabbri, Vilson Vieira da Silva Junior, Antônio Carlos Silvano Pessotti, Débora Cristina Corrêa, and Osvaldo N. Oliveira Jr.. 0. Musical elements in the discrete-time representation of sound. *ACM Comput. Surv.* 0, 0, Article 0 ( 0), 56 pages.

DOI: 0000001.0000001

---

## 1 INTRODUCTION

Music is usually defined as the art whose medium is sound. The definition might also state that the medium includes silences and temporal organization of structures, or that music is also a cultural

---

This work is supported by FAPESP and CNPq.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2017 Copyright held by the owner/author(s). 0360-0300/0/0-ART0 \$0.00

DOI: 0000001.0000001

activity or product. In physics and in this document, sounds are longitudinal waves of mechanical pressure. The human auditory system perceives sounds in the frequency bandwidth between  $20\text{Hz}$  and  $20\text{kHz}$ , with the actual boundaries depending on the person, climate conditions and the sonic characteristics themselves. Since the speed of sound is  $\approx 343.2\text{m/s}$ , such frequency limits corresponds to wavelengths of  $\frac{343.2}{20} \approx 17.16\text{ m}$  and  $\frac{343.2}{20000} \approx 17.16\text{ mm}$ . Hearing involves stimuli in bones, stomach, ears, transfer functions of head and torso, and processing by the nervous system. The ear is a dedicated organ for the appreciation of these waves, which decomposes them into their sinusoidal spectra and delivers to the nervous system. The sinusoidal components are crucial to musical phenomena, as one can recognize in the constitution of sounds of musical interest (such as harmonic sounds and noises, discussed in Sections 2 and 3), and higher level musical structures (such as tunings, scales and chords, in Section 4). [57]

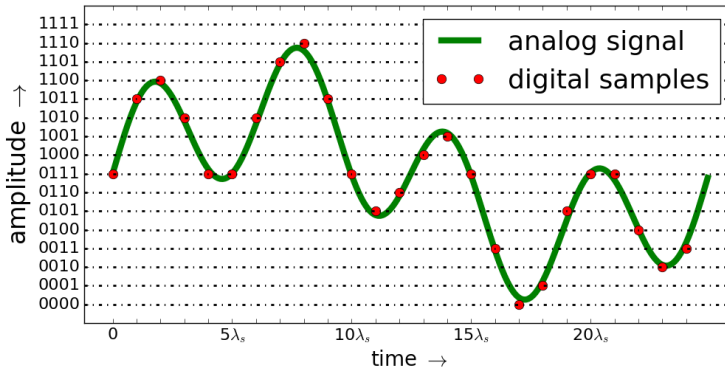


Fig. 1. Example of PCM audio: a sonic wave is represented by 25 samples equally spaced in time, where each sample has an amplitude specified with 4 bits.

The representation of sound can take many forms, from musical scores and texts in a phonetic language to electric analog signals and binary data. It includes sets of features such as wavelet or sinusoidal components. Although the terms 'audio' and 'sound' are often used without distinction and 'audio' has many definitions which depend on the context and the author, audio most often means a representation of the amplitude through time. In this sense, audio expresses sonic waves yield by synthesis or input by microphones, although these sources are not always neatly distinguishable e.g. as captured sounds are processed to generate new sonorities. Digital audio protocols often imply in quality loss (to achieve smaller files, ease storage and transfer) and are called *lossy* [49]. This is the case e.g. of MP3 and Ogg Vorbis. Non-lossy representations of digital audio, called *lossless* protocols or formats, on the other hand, assure perfect reconstruction of the analog wave within any convenient precision. The standard paradigm of lossless audio consists of representing the sound with samples equally spaced by a duration  $\delta_s$ , and specifying the amplitude of each sample by a fixed number of bits. This is the linear Pulse Code Modulation (LPCM) representation of sound, herein referred to as PCM. A PCM audio format has two essential attributes: a sampling frequency  $f_s = \frac{1}{\delta_s}$  (also called e.g. sampling rate or sample rate), which is the number of samples used for representing a second of sound; and a bit depth, which is the number of bits used for specifying the amplitude of each sample. Figure 1 shows 25 samples of a PCM audio with a bit depth of 4, which enables  $2^4 = 16$  possible values for the amplitude of each sample and a total of  $4 \times 25 = 100$  bits for representing the whole sound.

The fixed sampling frequency and bit depth yield the quantization error or quantization noise. This noise diminishes as the bit depth increases while greater sampling frequency allows higher frequencies to be represented. The Nyquist theorem asserts that the sampling frequency is twice the maximum frequency that the represented signal can contain [51]. Thus, for general musical purposes, it is suitable to use a sample rate of at least twice the highest frequency heard by humans, that is,  $f_s \geq 2 \times 20\text{kHz} = 40\text{kHz}$ . This is the basic reason for the adoption of sampling frequencies such as  $44.1\text{kHz}$  and  $48\text{kHz}$ , which are standards in Compact Disks (CD) and broadcast systems (radio and television), respectively.

Within this framework for representing sounds, musical notes can be characterized. The note often stands as the 'fundamental unit' of musical structures (such as atoms in matter or cells in macroscopic organisms) and, in practice, it can unfold into sounds that uphold other approaches to music. This is of capital importance because science and scholastic artists widened the traditional comprehension of music in the twentieth century to encompass discourse without explicit rhythm, melody or harmony. This is evident e.g. in the concrete, electronic, electroacoustic, and spectral musical styles. In the 1990s, it became evident that popular (commercial) music had also incorporated sound amalgams and abstract discursive arcs<sup>1</sup>. Notes are also convenient for another reason: the average listener – and a considerable part of the specialists – presupposes rhythmic and pitch organization (made explicit in Section 4) as fundamental musical properties, and these are developed in traditional musical theory in terms of notes. Thereafter, in this article we describe musical notes in PCM audio through equations and then indicate mechanisms for deriving higher level musical structures. We understand that this is not the unique approach to mathematically express music in digital audio, but musical theory and practice suggest that this is a proper framework for understanding and making computer music, as should become patent in the remainder of this text and is verifiable by usage of the MASS toolbox. Hopefully, the interested reader or programmer will be able to use this framework to synthesize music beyond traditional conceptualizations when intended.

This document provides a fundamental description of musical structures in discrete-time audio. The results include mathematical relations, usually in terms of musical characteristics and PCM samples, concise musical theory considerations, and their implementation as software routines both as very raw and straightforward algorithms and in the context of rendering musical pieces. Despite the general interests involved, there are only a few books and computer implementations that tackle the subject directly. These mainly focus on computer implementations and ways to mimic traditional instruments, with scattered mathematical formalisms for the basic notions. Articles on the topic appear to be lacking, to the best of our knowledge, which contrasts with the advanced and specialized developments often reported. A compilation of such works and their contributions is in the Appendix G of [23]. Although current music software uses the analytical descriptions presented here, there is no concise mathematical description of them, and it is far from trivial to achieve the equations by analyzing the available software implementations.

Accordingly, the objectives of this paper are to:

- (1) Present a concise set of mathematical and algorithmic relations between basic musical elements and sequences of PCM audio samples.
- (2) Introduce a framework for sound and musical synthesis with control at sample level which entails potential uses in psychoacoustic experiments, data sonification and synthesis with extreme precision (recap in Section 5).

<sup>1</sup>There are well known incidences of such characteristics in ethnic music, such as in Pygmy music, but western theory assimilated them only in the last century [76].

- (3) Provide a powerful theoretical framework which can be used to synthesize musical pieces and albums.
- (4) Provide approachability to the developed framework<sup>2</sup>.
- (5) Provide a didactic presentation of the content, which is highly multidisciplinary, involving signal processing, music, psychoacoustics and programming.

The reminder of this article is organized as follows: Section 2 characterizes the basic musical note; Section 3 develops internal dynamics of musical notes; Section 4 tackles the organization of musical notes into higher level musical structures [15, 43, 44, 56, 64, 74, 76, 78]. As these descriptions require knowledge on topics such as psychoacoustics, cultural traditions, and mathematical formalisms, the text points to external complements as needed and presents methods, results and discussions altogether. Section 5 is dedicated to final considerations and further work.

### 1.1 Additional material

One Supporting Information document [28] holds commented listings of all the equations, figures, tables and sections in this document and the scripts in the MASS toolbox. Another Supporting Information document [29] is a PDF version of the code that implements the equations and concepts in each section<sup>3</sup>. The Git repository [27] holds all the PDF documents and Python scripts. The rendered musical pieces are referenced when convenient and linked directly through URLs, and constitute another component of the framework. They are not very traditional, which facilitates the understanding of specific techniques and the extrapolation of the note concept. There are MASS-based software packages [24, 26] and further musical pieces that are linked in the Git repository.

### 1.2 Synonymy, polysemy and theoretical frames (disclaimer)

Given that the main topic of this article (the expression of musical elements in PCM audio) is multidisciplinary and involves art, the reader should be aware that much of the vocabulary admits different choices of terms and definitions. More specifically, it is often the case where many words can express the same concept and where one word can carry different meanings. This is a very deep issue which might receive a dedicated manuscript. The reader might need to read the rest of this document to understand this small selection of synonymy and polysemy in the literature, but it is important to illustrate the point before the more dense sections:

- a “note” can mean a pitch or an abstract construct with pitch and duration or a sound emitted from a musical instrument or a specific note in a score or a music.
- The sampling rate (discussed above) is also called the sampling frequency or sample rate.

<sup>2</sup> All the analytic relations presented in this article are implemented as small scripts in public domain. They constitute the MASS toolbox, available in an open source Git repository [10]. These routines are written in Python and make use of Numpy, which performs numerical routines efficiently (e.g. through LAPACK), but the language and packages are by no means mandatory. Part of the scripts has been ported to JavaScript (which favors their use in Web browsers such as Firefox and Chromium) and native Python [50, 58, 72]. These are all open technologies, published using licenses that grant permission for copying, distributing, modifying and usage in research, development, art and education. Hence, the work presented here aims at being compliant with recommended practices for availability and validation and should ease co-authorship processes [45, 54].

<sup>3</sup> The toolbox contains a collection of Python scripts which:

- implement each of the equations;
- render music and illustrate the concepts;
- render each of the figures used in this article.

The documentation of the toolbox consists of this article, the Supporting Information documents and the scripts themselves.

- A harmonic in a sound is most often a sinusoidal component which is in the harmonic series of the fundamental frequency. Many times, however, the terms harmonic and component are not distinguished. A harmonic can also be a note performed in an instrument by preventing certain overtones (components).
- Harmony can refer to chords or to note sets related to chords or even to “harmony” in a more general sense, as a kind of balance and consistency.
- A “tremolo” can mean different things: e.g. in a piano score, a tremolo is a fast alternation of two notes (pitches) while in computer music theory it is (most often) an oscillation of loudness.

We strived to avoid nomenclature clashes and the use of more terms than needed. Also, there are many theoretical standpoints for understanding musical phenomena, which is an evidence that most often there is not a single way to express or characterize musical structures. Therefore, in this article, adjectives such as “often”, “commonly” and “frequently” are abundant and they would probably be even more numerous if we wanted to be pedantically precise. Some of these issues are exposed when the context is convenient, such as in the first considerations of timbre.

## 2 CHARACTERIZATION OF THE MUSICAL NOTE IN DISCRETE-TIME AUDIO

In diverse artistic and theoretical contexts, music is conceived as constituted by fundamental units referred to as notes, “atoms” that constitute music itself [46, 74, 76]. In a cognitive perspective, notes are understood as discernible elements that facilitate and enrich the transmission of information through music [43, 57]. Canonically, the basic characteristics of a musical note are duration, loudness, pitch and timbre [43]. All relations described in this section are implemented in the file `src/sections/eqs2.1.py`. The musical pieces *5 sonic portraits* and *reduced-fi* are also available online to corroborate and illustrate the concepts.

### 2.1 Duration

The sample frequency  $f_s$  is defined as the number of samples in each second of the discrete-time signal. Let  $T = \{t_i\}$  be an ordered set of real samples separated by  $\delta_s = 1/f_s$  seconds ( $f_s = 44.1kHz \Rightarrow \delta_s = 1/44100 \approx 0.023ms$ ). A musical note of duration  $\Delta$  seconds can be expressed as a sequence  $T^\Delta$  with  $\Lambda = \lfloor \Delta \cdot f_s \rfloor$  samples. That is, the integer part of the multiplication is considered, and an error of at most  $\delta_s$  missing seconds is admitted, which is usually fine for musical purposes. Thus:

$$T^\Delta = \{t_i\}_{i=0}^{\lfloor \Delta \cdot f_s \rfloor - 1} = \{t_i\}_0^{\Lambda-1} \quad (1)$$

### 2.2 Loudness

Loudness<sup>4</sup> is a perception of sonic intensity that depends on reverberation, spectrum and other characteristics described in Section 3 [12]. One can achieve loudness variations through the power of the wave [12]:

$$pow(T) = \frac{\sum_{i=0}^{\Lambda-1} t_i^2}{\Lambda} \quad (2)$$

<sup>4</sup>Loudness and “volume” are often used indistinctly. In technical contexts, loudness is used for the subjective perception of sound intensity while volume might be used for some measurement of loudness or to a change in the intensity of the signal by equipment. Accordingly, one can perceive a sound as loud or soft and change the volume by turning a knob. We will use the term loudness and avoid the more ambiguous term volume.

The final loudness is dependent on the amplification of the signal by the speakers. Thus, what matters is the relative power of a note in relation to the others around it, or the power of a musical section in relation to the rest. Differences in loudness are the result of complex psychophysical phenomena but can often be reasoned about in terms of decibels, calculated directly from the amplitudes through energy or power:

$$V_{dB} = 10 \log_{10} \frac{\text{pow}(T')}{\text{pow}(T)} \quad (3)$$

The quantity  $V_{dB}$  has the decibel unit ( $dB$ ). By standard, a “doubled loudness” is associated to a gain of  $10dB$  (10 violins yield double the loudness of a violin). A handy reference is  $10dB$  for each step in the musical intensity scale: *pianissimo*, *piano*, *mezzoforte*, *forte* and *fortissimo*. Other useful references are  $dB$  values related to double amplitude or power:

$$t'_i = 2t_i \Rightarrow \text{pow}(T') = 4\text{pow}(T) \Rightarrow V'_{dB} = 10 \log_{10} 4 \approx 6dB \quad (4)$$

$$t'_i = \sqrt{2}t_i \Rightarrow \text{pow}(T') = 2\text{pow}(T) \Rightarrow V'_{dB} = 10 \log_{10} 2 \approx 3dB \quad (5)$$

and the amplitude gain for a sequence whose loudness has been doubled ( $10dB$ ):

$$\begin{aligned} 10 \log_{10} \frac{\text{pot}(T')}{\text{pot}(T)} &= 10 \Rightarrow \\ \Rightarrow \sum_{i=0}^{\lfloor \Delta \cdot f_s \rfloor - 1} t_i'^2 &= 10 \sum_{i=0}^{\Lambda-1} t_i^2 = \sum_{i=0}^{\Lambda-1} (\sqrt{10} \cdot t_i)^2 \\ \therefore t'_i &= \sqrt{10} t_i \Rightarrow t'_i \approx 3.16 t_i \end{aligned} \quad (6)$$

Thus, an amplitude increase by a factor slightly above 3 is required for achieving a doubled loudness. These values are guides for increasing or decreasing the absolute values in sample sequences. The conversion from decibels to amplitude gain (or attenuation) is straightforward:

$$A = 10^{\frac{V_{dB}}{20}} \quad (7)$$

where  $A$  is the multiplicative factor that relates the amplitudes before and after amplification.

### 2.3 Pitch

The perception of sounds as ‘higher’ or ‘lower’ is usually thought in terms of pitch. An exponential progression of frequency ( $f_i = f \cdot X^i, \forall X > 0, i \geq 1$ ) yields a linear variation of the pitch, a fact that will be further exploited in Sections 3 and 4. Accordingly, a pitch is specified by a (fundamental) frequency  $f$  whose cycle has duration  $\delta = 1/f$ . This duration, multiplied by the sampling frequency  $f_s$ , yields the number of samples per cycle:  $\lambda = f_s \cdot \delta = f_s / f$ . For didactic reasons, let  $f$  divide  $f_s$  and result  $\lambda$  integer. If  $T^f$  is a sonic sequence with fundamental frequency  $f$ , then:

$$T^f = \{t_i^f\} = \{t_{i+\lambda}^f\} = \left\{t_{i+\frac{f_s}{f}}^f\right\} \quad (8)$$

In the next section, frequencies  $f$  that do not divide  $f_s$  will be considered. This restriction does not imply a loss of the generality of this current section’s content.

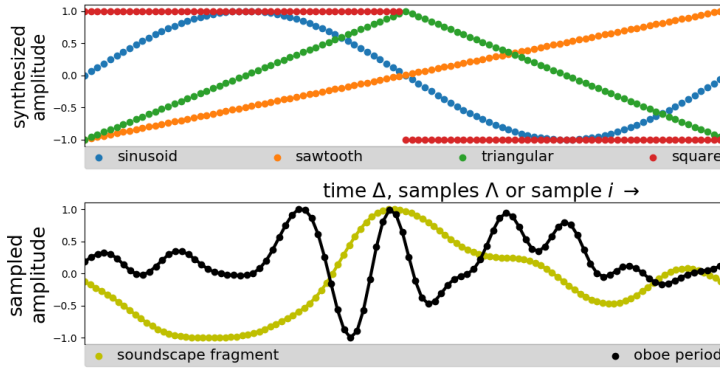


Fig. 2. Basic musical waveforms: (a) the basic synthetic waveforms given by the Equations 10, 11, 12 and 13; (b) real waveforms. Because of the period with  $\approx 100$  samples ( $\lambda_f \approx 100$ ), if  $f_s = 44.1\text{kHz}$  the basic and oboe waves have a fundamental frequency of  $f = \frac{f_s}{\lambda_f} \approx \frac{44100}{100} = 441\text{ Hz}$ , whatever the waveform is.

## 2.4 Timbre

A spectrum is said harmonic if all the (sinusoidal) frequencies  $f_n$  it contains are (whole number) multiples of a fundamental frequency  $f_0$  (lowest frequency):  $f_n = (n + 1)f_0$ . From a musical perspective, it is critical to internalize that energy in a component with frequency  $f$  is a sinusoidal oscillation in the constitution of the sound in that frequency  $f$ . This energy, specifically concentrated on  $f$ , is separated from other frequencies by the ear for further cognitive processes (this separation is performed by diverse living organisms by mechanisms similar to what is achieved by the human cochlea). The sinusoidal components are responsible for timbre<sup>5</sup> qualities (including pitch). If their frequencies do not relate by small integers, the sound is perceived as noisy or dissonant, in opposition to sonorities with an unequivocally established fundamental. Accordingly, the perception of absolute pitch relies on the similarity of the spectrum to the harmonic series. [57]

A sound with a harmonic spectrum has a wave period (wave cycle duration) which corresponds to the inverse of the fundamental frequency. The trajectory of the wave inside the period is the *waveform* and implies a specific combination of amplitudes and phases of the harmonic spectrum. Sonic spectra with minimal differences can result in timbres with crucial differences and, consequently, distinct timbres can be produced using different waveforms.

High curvatures in the waveform hint that there is energy in the high frequencies. Figure 2 depicts a wave, labeled as “soundscape fragment”. The same figure also displays a sampled period from an oboe note. One can notice from the curvatures: the oboe’s rich spectrum at high frequencies and the greater contribution of the lower frequencies in the spectrum of the soundscape fragment.

The sequence  $R = \{r_i\}_0^{\lambda_f - 1}$  of samples in a real sound (e.g. of Figure 2) can be taken as a basis for a sound  $T^f$  in the following way:

$$T^f = \{t_i^f\} = \left\{ r_{(i \% \lambda_f)} \right\} \quad (9)$$

<sup>5</sup>The timbre of a sound is a subjective and complex characteristic. The timbre can be considered by the temporal evolution of energy in the spectral components that are harmonic or noisy (and by deviations of the harmonics from the ideal harmonic spectrum). In addition, the word timbre is used to designate different things: one same note can have (be produced with) different timbres, an instrument has different timbres, two instruments of the same family have, at the same time, the same timbre that blends them into the same family, and different timbres as they are different instruments. Timbre is not only about spectrum: culture and context alter our perception of timbre. [57]



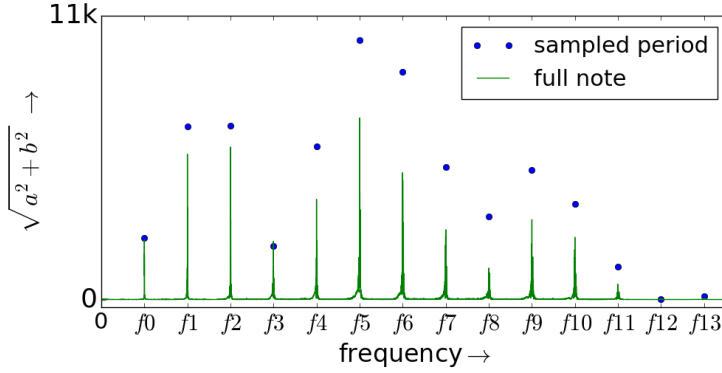


Fig. 3. Spectra of the sonic waves of a natural oboe note and obtained through a sampled period. The natural sound has fluctuations in the harmonics and in its noise, while the sampled period note has a perfectly harmonic (and static) spectrum.

The resulting sound has the spectrum of the original waveform. As a consequence of the identical repetitions, the spectrum is perfectly harmonic, without noise or variations of the components which are typical of natural phenomena. This can be observed in Figure 3, which shows the spectrum of the original oboe note and a note with the same duration, whose samples consist of the repetition of the cycle on Figure 2.

The simplest case is the spectrum with only one frequency, which is a sinusoid, often regarded as a “pure” oscillation (e.g. in terms of the *simple harmonic motion*). Let  $S^f$  be a sequence whose samples  $s_i^f$  describe a sinusoid with frequency  $f$ :

$$S^f = \{s_i^f\} = \left\{ \sin\left(2\pi \frac{i}{\lambda_f}\right) \right\} = \left\{ \sin\left(2\pi f \frac{i}{f_s}\right) \right\} \quad (10)$$

where  $\lambda_f = \frac{f_s}{f} = \frac{\delta_f}{\delta_s}$  is the number of samples in the period.

Other artificial waveforms are used in music for their spectral qualities and simplicity. While the sinusoid is an isolated node in the spectrum, any other waveform presents a succession of harmonic components (harmonics). Standard waveforms are specified by Equations 10, 11, 12 and 13, and are illustrated in Figure 2. These artificial waveforms are traditionally used in music for synthesis and oscillatory control of variables. They are also useful outside musical contexts [51].

The sawtooth presents all the harmonics with a decreasing energy of  $-6dB/octave$ <sup>6</sup>. The sequence of temporal samples can be described as:

$$D^f = \{d_i^f\} = \left\{ 2 \frac{i \% (\lambda_f + 1)}{\lambda_f} - 1 \right\} \quad (11)$$

The triangular waveform has only odd harmonics falling with  $-12dB/octave$ :

$$T^f = \{t_i^f\} = \left\{ 1 - \left| 2 - 4 \frac{i \% \lambda_f}{\lambda_f} \right| \right\} \quad (12)$$

The square wave also has only odd harmonics but falling at  $-6dB/octave$ :

$$Q^f = \{q_i^f\} = \begin{cases} 1 & \text{for } (i \% \lambda_f) < \lambda_f/2 \\ -1 & \text{otherwise} \end{cases} \quad (13)$$

<sup>6</sup>In musical jargon, an “octave” means a frequency and twice such frequency ( $f$  and  $2f$ ), or the bandwidth  $[f, 2f]$ .



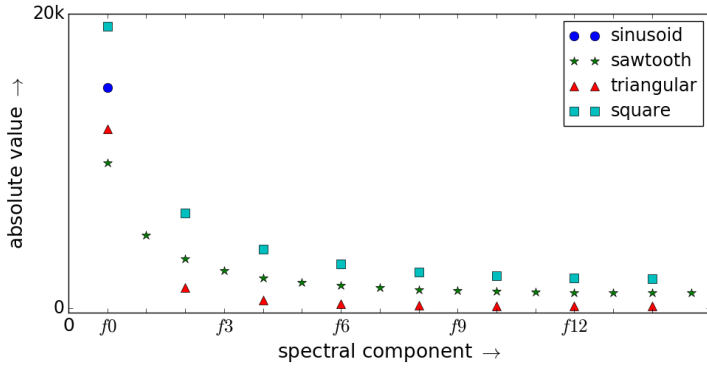


Fig. 4. Spectra of basic artificial waveforms. The isolated and exactly harmonic components of the spectra is a consequence of the fixed period. The figure exhibits the spectra described in Section 2.4: the sawtooth is the only waveform with a complete harmonic series (odd and even components); triangular and square waves have the same components (odd harmonics), decaying at  $-12\text{dB/octave}$  and  $-6\text{dB/octave}$ , respectively; the sinusoid consists of a unique node in the spectrum.

The square wave can be used in a subtractive synthesis with the purpose of mimicking a clarinet. This instrument has only the odd harmonics and the square wave is convenient with its abundant energy at high frequencies. The sawtooth is a common starting point for subtractive synthesis, because it has both odd and even harmonics with high energy. In general, these waveforms are appreciated as excessively rich in sharp harmonics, and attenuation by filtering on treble and middle parts of the spectrum is especially useful for achieving a more natural and pleasant sound. The relatively attenuated harmonics of the triangle wave makes it the more functional - among the listed possibilities - to be used in the synthesis of musical notes without any further processing. The sinusoid is often a nice choice, but a problematic one. While pleasant if not loud in a very high pitch (above  $\approx 500\text{Hz}$  it requires careful dosage), the pitch of a pure sinusoid is not accurately detected by the human auditory system, particularly at low frequencies. Also, it requires a great amplitude gain for an increase in loudness of a sinusoid if compared to other waveforms. Both particularities are understood in the scientific literature as a consequence of the nonexistence of pure sinusoidal sounds in nature [57]. The spectra of each basic waveform is illustrated in Figure 4.

## 2.5 Spectra of sampled sounds

The sinusoidal components in the discretized sound have some particularities. Considering a signal  $T$  and its corresponding Fourier decomposition  $\mathcal{F}\langle T \rangle = C = \{c_k\}_{k=0}^{\Lambda-1} = \left\{ \sum_{i=0}^{\Lambda-1} t_i e^{-ji(k\frac{2\pi}{\Lambda})} \right\}_{k=0}^{\Lambda-1}$ , the recomposition is the sum of the frequency components to yield the temporal samples<sup>7</sup>:

$$\begin{aligned} t_i &= \frac{1}{\Lambda} \sum_{k=0}^{\Lambda-1} c_k e^{j\frac{2\pi k}{\Lambda} i} \\ &= \frac{1}{\Lambda} \sum_{k=0}^{\Lambda-1} (a_k + j.b_k) [\cos(w_k i) + j.\sin(w_k i)] \end{aligned} \quad (14)$$

<sup>7</sup>The factor  $\frac{1}{\Lambda}$  can be distributed among the Fourier transform and its reconstruction, as preferred. Note that  $j$  here is the imaginary unit  $j^2 = -1$ .

where  $c_k = a_k + j.b_k$  defines the amplitude and phase of each frequency:  $w_k = \frac{2\pi}{\Lambda}k$  in radians or  $f_k = w_k \frac{f_s}{2\pi} = \frac{f_s}{\Lambda}k$  in Hertz, and are limited by  $w_k \leq \pi$  and  $f_k \leq \frac{f_s}{2}$  as given by the Nyquist Theorem.

For a sonic signal, samples  $t_i$  are real and are given by the real part of Equation 14:

$$\begin{aligned} t_i &= \frac{1}{\Lambda} \sum_{k=0}^{\Lambda-1} [a_k \cos(w_k i) - b_k \sin(w_k i)] \\ &= \frac{1}{\Lambda} \sum_{k=0}^{\Lambda-1} \sqrt{a_k^2 + b_k^2} \cos[w_k i - \arctan(b_k, a_k)] \end{aligned} \quad (15)$$

where  $\arctan(x, y) \in [0, 2\pi]$  is the inverse tangent with the right choice of the quadrant in the imaginary plane.

$\Lambda$  real samples  $t_i$  result in  $\Lambda$  complex coefficients  $c_k = a_k + j.b_k$ . The coefficients  $c_k$  are equivalent two by two, corresponding to the same frequencies and with the same contribution to its reconstruction. They are complex conjugates:  $a_{k1} = a_{k2}$  and  $b_{k1} = -b_{k2}$  and, as a consequence, the modules are equal and phases have opposite signs. Recalling that  $f_k = k \frac{f_s}{\Lambda}$ ,  $k \in \{0, \dots, \lfloor \frac{\Lambda}{2} \rfloor\}$ . When  $k > \frac{\Lambda}{2}$ , the frequency  $f_k$  is mirrored through  $\frac{f_s}{2}$  in this way:  $f_k = \frac{f_s}{2} - (f_k - \frac{f_s}{2}) = f_s - f_k = f_s - k \frac{f_s}{\Lambda} = (\Lambda - k) \frac{f_s}{\Lambda} \Rightarrow f_k \equiv f_{\Lambda-k}$ ,  $\forall k < \Lambda$ .

The same applies to  $w_k = f_k \frac{2\pi}{f_s}$  and the periodicity  $2\pi$ : it follows that  $w_k = -w_{\Lambda-k}$ ,  $\forall k < \Lambda$ . Given the cosine (an even function) and the inverse tangent (an odd function), the components in  $w_k$  and  $w_{\Lambda-k}$  contribute with coefficients  $c_k = c_{\Lambda-k}^*$  in the reconstruction of the real samples. In summary, in a decomposition of  $\Lambda$  samples, the  $\Lambda$  frequency components  $\{c_i\}_0^{\Lambda-1}$  are equivalent in pairs, except for  $f_0$ , and, when  $\Lambda$  is even, for  $f_{\Lambda/2} = f_{\max} = \frac{f_s}{2}$ . Both these components are isolated, i.e. there is one and only one component at frequency  $f_0$  or  $f_{\Lambda/2}$  (if  $\Lambda$  is even). In fact, when  $k = 0$  or  $k = \Lambda/2$  the mirror of the frequencies are themselves:  $f_{\Lambda/2} = f_{(\Lambda-\Lambda/2)=\Lambda/2}$  and  $f_0 = f_{(\Lambda-0)=\Lambda} = f_0$ . Furthermore, these two frequencies (zero and Nyquist frequency) do not have a phase offset: their coefficients are strictly real. Therefore, the number  $\tau_\Lambda$  of equivalent coefficient pairs in a decomposition of  $\Lambda$  samples is:

$$\tau_\Lambda = \frac{\Lambda - \Lambda \% 2}{2} - 2 + \Lambda \% 2 = \left\lfloor \frac{\Lambda}{2} \right\rfloor - 2 \quad (16)$$

This discussion can be summarized in the following equivalences:

$$f_k \equiv f_{\Lambda-k} \quad , \quad w_k \equiv -w_{\Lambda-k} \quad (17)$$

$$a_k = a_{\Lambda-k} \quad , \quad b_k = -b_{\Lambda-k} \quad (18)$$

$$\sqrt{a_k^2 + b_k^2} = \sqrt{a_{\Lambda-k}^2 + b_{\Lambda-k}^2} \quad (19)$$

$$\arctan(b_k, a_k) = -\arctan(b_{\Lambda-k}, a_{\Lambda-k}) \quad (20)$$

with  $\forall 1 \leq k \leq \tau_\Lambda$ ,  $k \in \mathbb{N}$ .

To express the general case for components combination in each sample  $t_i$ , one can gather the relations for the reconstruction of a real signal (Equation 15), for the number of paired coefficients (Equation 16), and for the equivalences of modules (Equation 19) and phases (Equation 20):

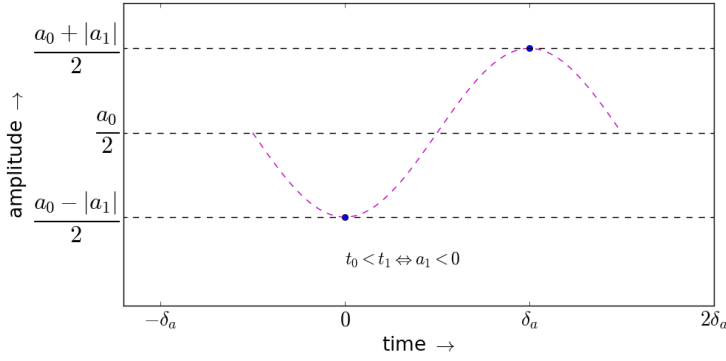


Fig. 5. Oscillation of 2 samples (maximum frequency for any  $f_s$ ). The first coefficient determines a constant detachment (called *offset*, *bias* or *DC component*) and the second coefficient specifies the oscillation amplitude.

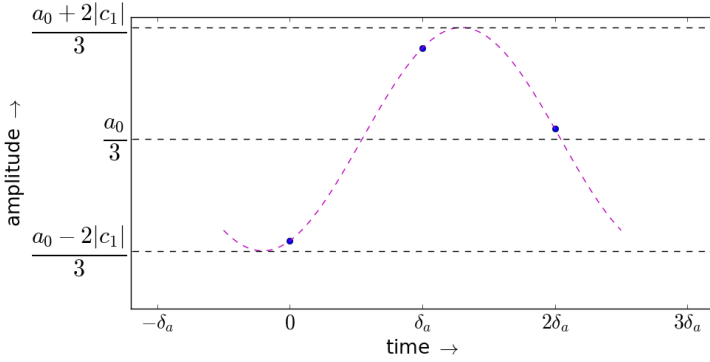


Fig. 6. Three fixed samples present only one non-null frequency.  $c_1 = c_2^*$  and  $w_1 \equiv w_2$ .

$$t_i = \frac{a_0}{\Lambda} + \frac{a_{\Lambda/2}}{\Lambda}(1 - \Lambda\%2) + \frac{2}{\Lambda} \sum_{k=1}^{\tau_{\Lambda}} \sqrt{a_k^2 + b_k^2} \cos[w_k i - \arctan(b_k, a_k)] \quad (21)$$

Figure 5 shows two samples and their spectral component. When there is only two samples, the Fourier decomposition has only one pair of coefficients  $\{c_k = a_k - j.b_k\}_0^{\Lambda-1=1}$  relative to frequencies  $\{f_k\}_0^1 = \{w_k \frac{f_s}{2\pi}\}_0^1 = \{k \frac{f_s}{\Lambda=2}\}_0^1 = \{0, \frac{f_s}{2} = f_{\max}\}$  with energies  $e_k = \frac{(c_k)^2}{\Lambda=2}$ . The role of amplitudes  $a_k$  is clearly observed with  $\frac{a_0}{2}$ , the fixed offset (also called *bias* or *DC component*), and  $\frac{a_1}{2}$  for the oscillation with frequency  $f_1 = \frac{f_s}{\Lambda=2}$ . This case has special relevance: at least 2 samples are necessary to represent an oscillation and it yields the Nyquist frequency  $f_{\max} = \frac{f_s}{2}$ , which is the maximum frequency in a sound sampled with  $f_s$  samples per second.

All fixed sequences  $T$  of only 3 samples also have just 1 frequency, since the first harmonic would have 1.5 samples and exceeds the bottom limit of 2 samples, i.e. the frequency of the harmonic would exceed the Nyquist frequency:  $\frac{2 \cdot f_s}{3} > \frac{f_s}{2}$ . The coefficients  $\{c_k\}_0^{\Lambda-1=2}$  are present in 3 frequency components. One is relative to frequency zero ( $c_0$ ), and the other two ( $c_1$  and  $c_2$ ) have the same role for reconstructing a sinusoid with  $f = f_s/3$ . This case is illustrated in Figure 6.

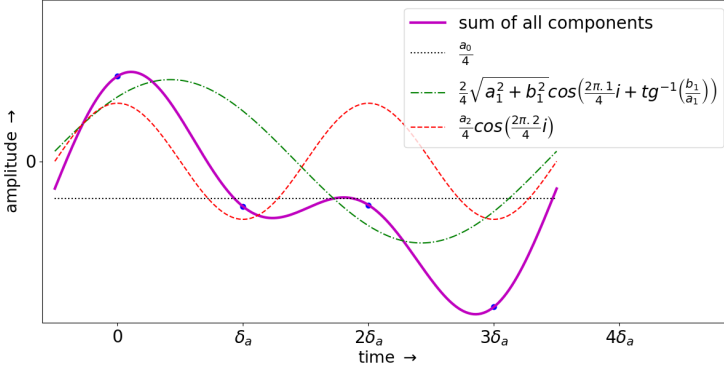


Fig. 7. Frequency components for 4 samples.

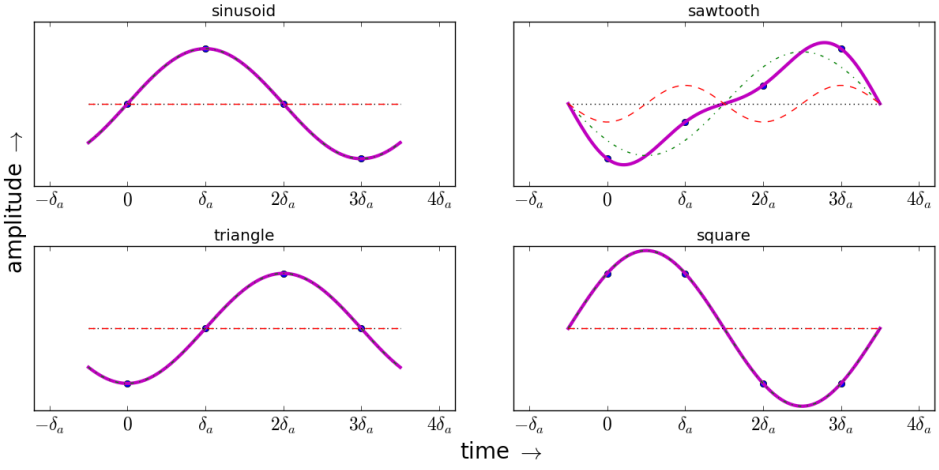
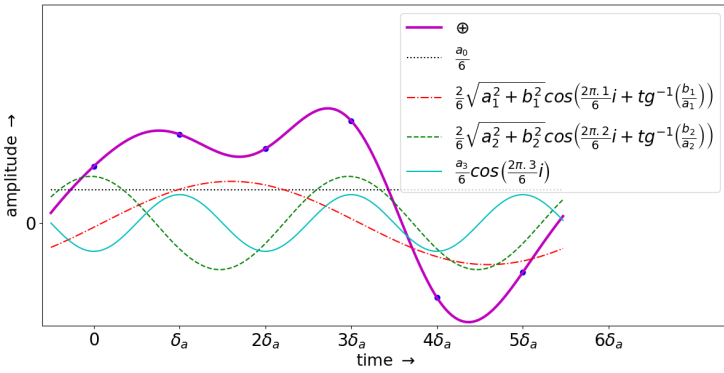


Fig. 8. Basic waveforms with 4 samples.

Fig. 9. Frequency components for 6 samples: 4 sinusoids, one of them is the *bias* with zero frequency.

With 4 samples it is possible to represent 1 or 2 frequencies with independence of magnitude and phase. Figure 7 depicts the contribution of each of the two (possible) components. The individual components sum to the original waveform and a brief inspection reveals the major curvatures resulting from the higher frequency, while the fixed offset is captured in the component with frequency  $f_0 = 0$ . Figure 8 shows the harmonics for the basic waveforms of Equations 10, 11, 12 and 13 in the case of 4 samples. There is only 1 sinusoid for each waveform, with the exception of the sawtooth, which has the even harmonics.

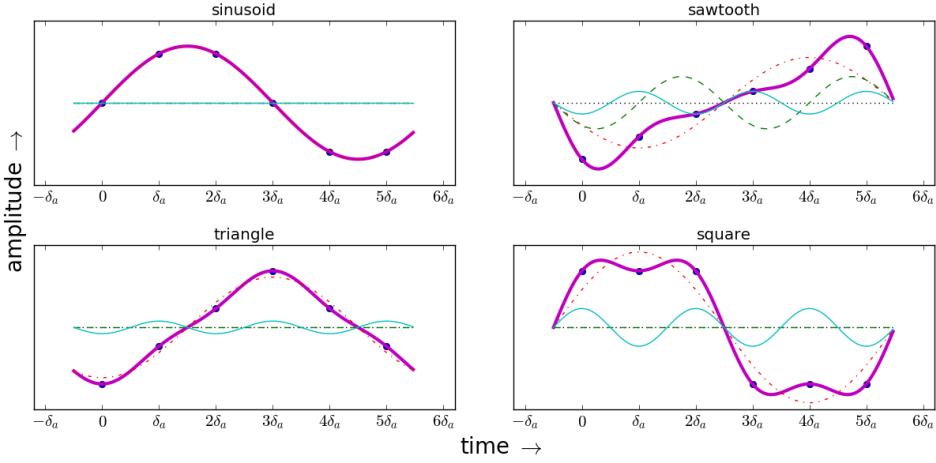


Fig. 10. Basic waveforms with 6 samples: triangular and square waveforms have odd harmonics, with different proportions and phases; the sawtooth has even harmonics.

Figure 9 exposes the sinusoidal components within 6 samples, while Figure 10 presents the decomposition of the basic waveforms: square and triangular have the same components but in different proportions, while the sawtooth has an extra component.

## 2.6 The basic note

In a nutshell, a sequence  $T$  of sonic samples separated by  $\delta_a = 1/f_s$  expresses a musical note with a frequency of  $f$  Hertz<sup>8</sup> and  $\Delta$  seconds of duration if, and only if, it has the periodicity  $\lambda_f = f_s/f$  and size  $\Lambda = \lfloor f_s \cdot \Delta \rfloor$ :

$$T^{f, \Delta} = \{t_{i \% \lambda_f}\}_{i=0}^{\Lambda-1} = \left\{ t_{i \% \left(\frac{f_s}{f}\right)}^f \right\}_{i=0}^{\Lambda-1} \quad (22)$$

Such note still does not have a timbre: it is necessary to choose a waveform for the samples  $t_i$  to have a value. Any waveform can be used to further specify the note, where  $\lambda_f = \frac{f_s}{f}$  is the number of samples in each period. Let  $L^f \in \{S^f, Q^f, T^f, D^f, R^f\}$  (as given by Equations 10, 11, 12 and 13 and let  $R_i^f$  be a sampled waveform) be the sequence that describes a period of the waveform with duration  $\delta_f = 1/f$ :

<sup>8</sup>Let  $f$  be such that it divides  $f_s$ . As mentioned before, this limitation simplifies the exposition for now and will be overcome in the next section.

$$L^f = \left\{ l_i^f \right\}_0^{\delta_f \cdot f_s^{-1}} = \left\{ l_i^f \right\}_0^{\lambda_f^{-1}} \quad (23)$$

Thereafter, the sequence  $T$  for a note of duration  $\Delta$  and frequency  $f$  is:

$$T^{f, \Delta} = \left\{ t_i^f \right\}_0^{\lfloor f_s \cdot \Delta \rfloor - 1} = \left\{ l_{i \% \left( \frac{f_s}{f} \right)}^f \right\}_0^{\Delta - 1} \quad (24)$$

## 2.7 Spatialization: localization and reverberation

A musical note is always spatialized (i.e. it is always produced within the ordinary three dimensional physical space) even though it is not one of its four basic properties in canonical music theory (duration, loudness, pitch and timbre). The consideration of this fact is the subject of the spatialization knowledge field and practice<sup>9</sup>. A note has a source which has a three dimensional position. This position is the spatial localization of the sound. It is often (modeled as) a single point but can be a surface or a volume. The reverberation in the environment in which a sound occurs is an important topic of spatialization. Both concepts, spatial localization and reverberation, are widely valued by composers, audiophiles and the music industry [48].

**2.7.1 Spatial localization.** It is understood that the perception of sound localization occurs in our nervous system mainly by three cues: the delay of the incoming sound (and its reflections in the surfaces) between both ears, the difference of sound intensity at each ear and the filtering performed by the human body, specially in the chest, head and ears [11, 37, 57].

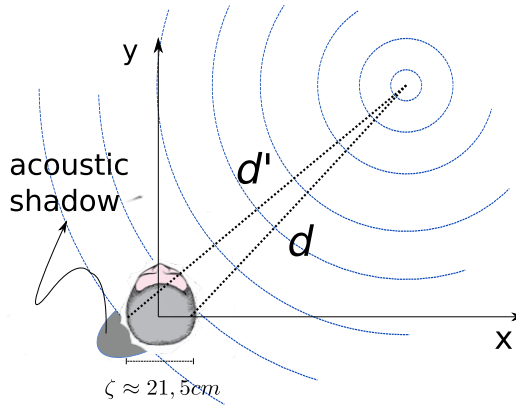


Fig. 11. Detection of sound source localization: schema used to calculate the Interaural Time Difference (ITD) and the Interaural Intensity Difference (IID).

An object placed at  $(x, y)$ , as in Figure 11, is distant of each ear by:

<sup>9</sup>By spatialization one might find both: 1) the consideration of cues in sound that derive from the environment, including the localization of the listener and the sound source; 2) techniques to produce sound through various sources, such as loudspeakers, singers and traditional musical instruments, for musical purposes. We focus in the first issue although issues of the second are also tackled and they are obviously intermingled.

$$\begin{aligned}
 d &= \sqrt{\left(x - \frac{\zeta}{2}\right)^2 + y^2} \\
 d' &= \sqrt{\left(x + \frac{\zeta}{2}\right)^2 + y^2}
 \end{aligned} \tag{25}$$

where  $\zeta$  is the distance between the ears, known to be  $\zeta \approx 21.5\text{cm}$  in an adult human. The cues for the sonic localization are not easy to calculate, but, in a very simplified model, useful for musical purposes, straightforward calculations result in the Interaural Time Difference:

$$ITD = \frac{d' - d}{v_{\text{sound at air}} \approx 343.2} \quad \text{seconds} \tag{26}$$

and in the Interaural Intensity Difference:

$$IID = 20 \log_{10} \left( \frac{d}{d'} \right) \quad \text{decibels} \tag{27}$$

$IID_a = \frac{d}{d'}$  can be used as a multiplicative constant to the right channel of a stereo sound signal together with ITD [37]:

$$\begin{aligned}
 \Lambda_{ITD} &= \left\lfloor \frac{d' - d}{343.2} f_s \right\rfloor \\
 IID_a &= \frac{d}{d'} \\
 \left\{ t'_{(i+\Lambda_{ITD})} \right\}_{\Lambda_{ITD}}^{\Lambda_{ITD}-1} &= \{ IID_a \cdot t_i \}_0^{\Lambda_{ITD}-1} \\
 \{ t'_i \}_0^{\Lambda_{ITD}-1} &= 0
 \end{aligned} \tag{28}$$

where, where  $\{t'_i\}$  are samples of the wave incident in the left ear,  $\{t_i\}$  are samples for the right ear, and  $\Lambda_{ITD} = \lfloor ITD \cdot f_s \rfloor$ . If  $\Lambda_{ITD} < 0$ , it is necessary to change  $t_i$  by  $t'_i$  and use  $\Lambda'_{ITD} = |\Lambda_{ITD}|$  and  $IID'_a = 1/IID_a$ .

Spatial localization depends considerably on other cues. By using only ITD and IID it is possible to specify solely the horizontal angle (azimuthal)  $\theta$  given by:

$$\theta = \arctan(y, x) \tag{29}$$

with  $x, y$  as presented in Figure 11. Notice that the same pair of ITD and IID (as defined in Equations 26 and 27) is related to all the points in a vertical circle parallel to the head, i.e. the source can have any horizontal component inside the circle. Such a circle is called the "cone of confusion". In general, one can assume that the source is in the same horizontal plane as the listener and at its front (because humans are prone to hearing frontal and horizontal sources). Even in such cases, there are other important cues for sound localization. Consider the acoustic shadow depicted in Figure 11: for lateral sources the inference of the azimuthal angle is especially dependent on the filtering of frequencies by the head, pinna (outer ear) and torso. Also, low frequencies diffract and arrive to the opposite ear with a greater ITD. The complete localization, including height and distance of a sound source, is given by the Head Related Transfer Function (HRTF). There are well known open databases of HRTFs, such as CIPIC, and it is possible to apply such transfer functions in a sonic signal by convolution (see Equation 43). Each human body has its own filtering and there are techniques to generate HRTFs to be universally used. [3, 9, 11, 37, 48]



**2.7.2 Reverberation.** The reverberation results from sound reflections and absorption by the environment (e.g. a room) surface where a sound occurs. The sound propagates through the air with a speed of  $\approx 343.2m/s$  and can be emitted from a source with any directionality pattern. When a sound front encounters a surface there are: 1) inversion of the propagation speed component normal to the surface; 2) energy absorption, especially in high frequencies. The sonic waves propagate until they reach inaudible levels (and further but then can often be neglected). As a sonic front reaches the human ear, it can be described as the original sound, with the last reflection point as the source, and the absorption filters of each surface it has reached. It is possible to simulate reverberations that are impossible in real systems. For example, it is possible to use asymmetric reflections with relation to the axis perpendicular to the surface, to model propagation in a space with more than three dimensions, or consider a listener located in various positions.

There are reverberation models less related to each independent reflection and that explores valuable cues to the auditory system. In fact, reverberation can be modeled with a set of two temporal and two spectral regions [69]:

- First period: 'first reflections' are more intense and scattered.
- Second period: 'late reverberation' is practically a dense succession of indistinct delays with exponential decay and statistical occurrences.
- First band: the bass has some resonance bandwidths relatively spaced.
- Second band: mid and treble have a progressive decay and smooth statistical fluctuations.

Smith III states that usual concert rooms have a total reverberation time of  $\approx 1.9$  seconds, and that the period of first reflections is around 0.1s. With these values, there are perceived wave fronts which propagate for 652.08m before reaching the ear. In addition, sound reflections made after propagation for 34.32m have incidences less distinct by hearing. These first reflections are particularly important to spatial sensation. The first incidence is the direct sound, described by ITD and IID e.g. as in Equations 26 and 27. Assuming that each one of the first reflections, before reaching the ear, will propagate at least 3 – 30m, depending on the room dimensions, the separation between the first reflections is 8 – 90ms. Also, it is experimentally verifiable that the number of reflections increases with the square of time. A discussion about the use of convolutions and filtering to favor the implementation of these phenomena is provided in Section 3.6, particularly in the paragraphs about reverberation. [69]

## 2.8 Musical usages

Once the basic note is defined, it is didactically convenient to build musical structures with sequences based on these particles. The sum of the amplitudes of  $N$  sequences with same size  $\Lambda$  results in the overlapped spectral contents of each sequence, in a process called mixing:

$$\{t_i\}_0^{\Lambda-1} = \left\{ \sum_{k=0}^{N-1} t_{k,i} \right\}_0^{\Lambda-1} \quad (30)$$

Figure 12 illustrates this overlapping process of discretized sound waves, each with 100 samples. If  $f_s = 44.1kHz$ , the frequencies of the sawtooth, square and sine wave are, respectively:  $\frac{f_s}{100/2} = 882Hz$ ,  $\frac{f_s}{100/4} = 1764Hz$  and  $\frac{f_s}{100/5} = 2205Hz$ . The duration of each sequence is very short  $\frac{f_s=44.1kHz}{100} \approx 2ms$ . One can complete the sequence with zeroes to sum (mix) sequences with different sizes.

The mixed notes are generally separated by the ear according to the physical laws of resonance and by the nervous system [57]. This process of mixing musical notes results in musical harmony,

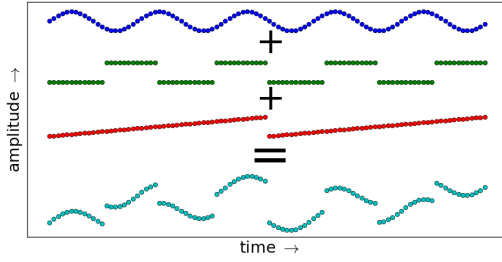


Fig. 12. Mixing of three sonic sequences. The amplitudes are directly summed sample-by-sample.

where intervals between frequencies and chords of simultaneous notes guide subjective and abstract aspects of music appreciation [62] and are addressed in Section 4.

Sequences can be concatenated in time. If the sequences  $\{t_{k,i}\}_0^{\Lambda_k-1}$  represent musical notes, their concatenation in a unique sequence  $T$  is a simple melodic sequence (or melody):

$$T = \{t_i\}_0^{\sum \Lambda_k-1} = \{t_{l,i}\}_0^{\sum \Lambda_k-1},$$

$$l \text{ smallest integer} : \Lambda_l > i - \sum_{j=0}^{l-1} \Lambda_j \quad (31)$$

This mechanism is illustrated in Figure 13 with the same sequences of Figure 12. Although the sequences are short for the usual sample rates, it is easy to visually observe the concatenation of sonic sequences. In addition, each note has a duration larger than 100ms if  $f_s < 1kHz$  (but need to oscillate faster to yield audible frequencies).

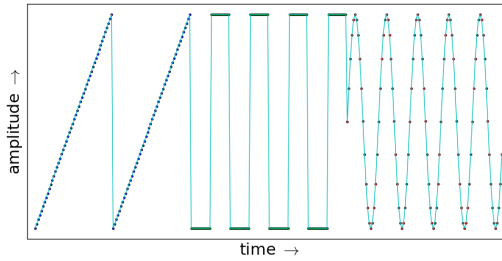


Fig. 13. Concatenation of three sounds.

The musical piece *reduced-fi* explores the temporal juxtaposition of notes, resulting in a homophonic piece. The vertical principle (mixing) is demonstrated at the *sonic portraits*, static sounds with peculiar spectrum. [27]

With the basic musical note in discrete-time audio carefully described, the next section develops the temporal evolution of its contents as in *glissandi* and intensity envelopes. Filtering of spectral components and noise generation complements the musical note as a self-contained unit. Section 4 is dedicated to the organization of these notes e.g. by using metrics and trajectories, with regards to traditional music theory.

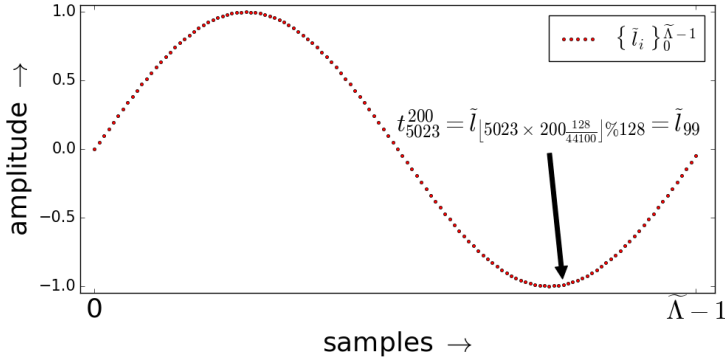


Fig. 14. Search (lookup) in a reference table (*lookup table* or LUT) to synthesize sounds at different frequencies using a unique waveform with high resolution. Each  $i$ -th sample  $t_i^f$  of a sound with frequency  $f$  is related to the samples in the table  $\tilde{L} = \{\tilde{l}_i\}_{\tilde{\Lambda}-1}^0$  by  $t_i^f = \tilde{l}_{[if \frac{\tilde{\Lambda}}{f_s}] \% \tilde{\Lambda}}$  where  $f_s$  is the sampling rate.

### 3 VARIATION IN THE BASIC NOTE

The basic digital music note defined in Section 2 has the following parameters: duration, pitch, intensity (loudness) and timbre. This is a useful and paradigmatic model, but it does not exhaust all the aspects of a musical note. First of all, characteristics of the note change along the note itself [12]. For example, a 3s piano note has intensity with an abrupt rise at the beginning and a progressive decay, has spectral variations with harmonics decaying and some others emerging along time. These variations are not mandatory, but they are used in sound synthesis for music because they reflect how sounds appear in nature. This is considered so important that there is a rule of thumb: to make a sound that incites interest by itself, arrange internal variations on it [57]. To explore all the ways by which variations occur within a note is out of the scope of any work, given the sensibility of the human ear and the complexity of human sound cognition. In this section, primary resources are presented to produce variations in the basic note. It is worthwhile to recall that all the relations in this and other sections are implemented in Python and published in public domain. The musical pieces *ParaMeter transitions*; *Shakes and wiggles*; *Tremolos, vibratos and frequency*; *Little train of impulsive hillbillies*; *Noisy band*; *Bela rugosi*; *Children choir*; and *ADa and SaRa* were made to validate and illustrate concepts of this section. The code that synthesizes these pieces is also part of the toolbox. [27]

#### 3.1 Lookup table

The *Lookup Table* (LUT) is an array for indexed operations which substitutes continuous and repetitive calculations. It is used to reduce computational complexity and for employing functions without calculating them directly, e.g. from sampled data or hand picked values. In music its usage simplifies many operations and enables the use a single wave period to synthesize sounds in the whole audible spectrum, with any waveform.

Let  $\tilde{\Lambda}$  be the wave period in samples and  $\tilde{L} = \{\tilde{l}_i\}_{\tilde{\Lambda}-1}^0$  the sample sequence with the waveform. A sequence  $T^{f, \Delta}$  with samples of a sound with frequency  $f$  and duration  $\Delta$  can be obtained by means of  $\tilde{L}$ :

$$T^{f, \Delta} = \left\{ t_i^f \right\}_0^{\lfloor f_s \cdot \Delta \rfloor - 1} = \left\{ \tilde{t}_{\gamma_i \% \tilde{\Lambda}} \right\}_0^{\Lambda - 1}, \quad \text{where } \gamma_i = \left\lfloor i f \frac{\tilde{\Lambda}}{f_s} \right\rfloor \quad (32)$$

In other words, with the right LUT indexes ( $\gamma_i \% \tilde{\Lambda}$ ) it is possible to synthesize sounds at any frequency. Figure 14 illustrates the calculation of a sample  $t_i$  from  $\left\{ \tilde{t}_i \right\}$  for  $f = 200\text{Hz}$ ,  $\tilde{\Lambda} = 128$  and adopting the sample rate of  $f_s = 44.1\text{kHz}$ . Though this is not a practical configuration (as discussed below), it allows for a graphical visualization of the procedure.

The calculation of the integer  $\gamma_i$  introduces noise which decreases as  $\tilde{\Lambda}$  increases. In order to use this calculation in sound synthesis, with  $f_s = 44.1\text{kHz}$ , the standard guidelines suggest the use of  $\tilde{\Lambda} = 1024$  samples, yielding a noise level of  $\approx -60\text{dB}$ . Larger tables might be used to achieve sounds with a greater quality. Also, a rounding or interpolation method can be used, but we advocate the use of a larger table since it does not introduce relevant computation overhead. [31]

The expression defining the variable  $\gamma_i$  can be understood as  $f_s$  being added to  $i$  at each second. If  $i$  is divided by the sample frequency,  $\frac{i}{f_s}$  is incremented by 1 at each second. Multiplied by the period, it results in  $i \frac{\tilde{\Lambda}}{f_s}$ , which covers the period in one second. Finally, with frequency  $f$  it results in  $i f \frac{\tilde{\Lambda}}{f_s}$  which completes  $f$  periods  $\tilde{\Lambda}$  in 1 second, i.e. the resulting sequence presents the fundamental frequency  $f$ .

There are important considerations here: it is possible to use practically any frequency  $f$ . Limits exist only at low frequencies when the size of table  $\tilde{\Lambda}$  is not sufficient for the sample rate  $f_s$ . The lookup procedure is virtually costless and replaces calculations by simple indexed searches (what is generally understood as an optimization process). Unless otherwise stated, this procedure will be used along all the following discussions for every applicable case. LUTs are broadly used in computational implementations for music, and are known also as wavetables. A classical usage of LUTs is known as *Wavetable Synthesis*, which generally consists of many LUTs used together to generate a quasi-periodic musical note [6, 15].

### 3.2 Incremental variations of frequency and intensity

As stated by the (Weber and) Fechner law [17], human perception holds a logarithmic relation to stimulus. That is to say, the exponential progression of a stimulus is perceived as linear. For didactic reasons, and given its use in AM and FM synthesis (Section 3.5), linear variation is discussed first.

Consider a note with duration  $\Delta = \frac{\Lambda}{f_s}$ , in which the frequency  $f = f_i$  varies linearly from  $f_0$  to  $f_{\Lambda-1}$ . Thus:

$$F = \{f_i\}_0^{\Lambda-1} = \left\{ f_0 + (f_{\Lambda-1} - f_0) \frac{i}{\Lambda - 1} \right\}_0^{\Lambda-1} \quad (33)$$

$$\begin{aligned} \Delta_{\gamma_i} = \frac{\tilde{\Lambda}}{f_s} f_i &\Rightarrow \gamma_i = \left\lfloor \sum_{j=0}^i \frac{\tilde{\Lambda}}{f_s} f_j \right\rfloor \\ \gamma_i &= \left\lfloor \sum_{j=0}^i \frac{\tilde{\Lambda}}{f_s} \left[ f_0 + (f_{\Lambda-1} - f_0) \frac{j}{\Lambda - 1} \right] \right\rfloor \end{aligned} \quad (34)$$

$$\left\{ \overline{t_i^{f_0, f_{\Lambda-1}}} \right\}_0^{\Lambda-1} = \left\{ \tilde{t}_{\gamma_i \% \tilde{\Lambda}} \right\}_0^{\Lambda-1} \quad (35)$$

where  $\Delta_{y_i} = f_i \frac{\tilde{\Lambda}}{f_s}$  is the LUT increment between two samples given the sound frequency of the first sample. There is a general rule to be noticed here: when a sound has variations in the fundamental frequency, one should account for them in the LUT indexing. The resulting indexes can be found by a cumulative sum of each indexing displacement. The equations for linear pitch transition are:

$$F = \{f_i\}_0^{\Lambda-1} = \left\{ f_0 \left( \frac{f_{\Lambda-1}}{f_0} \right)^{\frac{i}{\Lambda-1}} \right\}_0^{\Lambda-1} \quad (36)$$

$$\Delta_{y_i} = \frac{\tilde{\Lambda}}{f_s} f_i \Rightarrow y_i = \left\lfloor \sum_{j=0}^i \frac{\tilde{\Lambda}}{f_s} f_j \right\rfloor \quad (37)$$

$$y_i = \left\lfloor \sum_{j=0}^i f_0 \frac{\tilde{\Lambda}}{f_s} \left( \frac{f_{\Lambda-1}}{f_0} \right)^{\frac{j}{\Lambda-1}} \right\rfloor$$

$$\left\{ t_i^{\overline{f_0, f_{\Lambda-1}}} \right\}_0^{\Lambda-1} = \left\{ \tilde{t}_{y_i \% \tilde{\Lambda}} \right\}_0^{\Lambda-1} \quad (38)$$

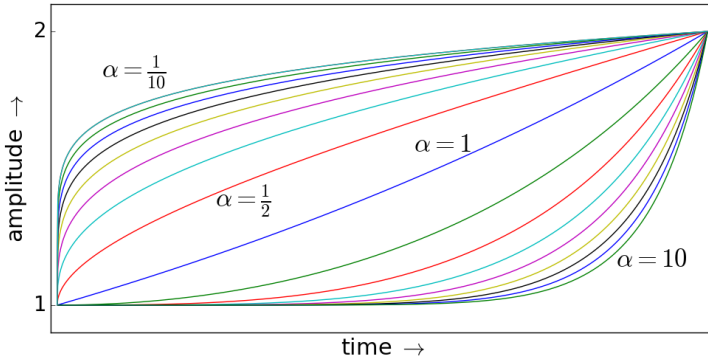


Fig. 15. Intensity transitions for different values of  $\alpha$  (see Equations 39 and 40).

The term  $\frac{i}{\Lambda-1}$  covers the interval  $[0, 1]$  and it is possible to raise it to a power  $\alpha \geq 0$  in such a way that the beginning of the transition will be smoother or steeper. This procedure is especially useful for energy variations with the purpose of changing the loudness<sup>10</sup>. Thus, for amplitude variations:

$$\{a_i\}_0^{\Lambda-1} = \left\{ a_0 \left( \frac{a_{\Lambda-1}}{a_0} \right)^{\left( \frac{i}{\Lambda-1} \right)^\alpha} \right\}_0^{\Lambda-1} = \left\{ (a_{\Lambda-1})^{\left( \frac{i}{\Lambda-1} \right)^\alpha} \right\}_0^{\Lambda-1} \quad (39)$$

where  $a_0$  is the initial amplitude factor and  $a_{\Lambda-1}$  is an amplitude factor to be reached at the end of the transition. Applying the loudness transition to a sonic sequence  $T$ :

$$T' = T \odot A = \{t_i \cdot a_i\}_0^{\Lambda-1} = \left\{ t_i \cdot (a_{\Lambda-1})^{\left( \frac{i}{\Lambda-1} \right)^\alpha} \right\}_0^{\Lambda-1} \quad (40)$$

It is often convenient to have  $a_0 = 1$  to start a new sequence with the original amplitude and then progressively change it. If  $\alpha = 1$ , the amplitude variation follows the exponential progression

<sup>10</sup>See Section 2.2 for considerations about loudness, amplitudes and decibels.

that is related to the linear variation of loudness. Figure 15 depicts transitions between values 1 and 2 and for different values of  $\alpha$ , a gain of  $\approx 6dB$  as given by Equation 4.

Special attention should be given while considering  $a = 0$ . In Equation 39,  $a_0 = 0$  results in a division by zero and if  $a_{\Lambda-1} = 0$ , there will be a multiplication by zero. Both cases make the procedure useless, once a ratio of any number in relation to zero is not well defined for our purposes. It is possible to solve this dilemma choosing a number that is small enough like  $-80dB \Rightarrow a = 10^{\frac{-80}{20}} = 10^{-4}$  as the minimum loudness for a *fade in* ( $a_0 = 10^{-4}$ ) or for a *fade out* ( $a_{\Lambda-1} = 10^{-4}$ ). A linear fade can be used then to reach zero amplitude, if needed. Another common solution is the use of the quartic polynomial term  $x^4$ , as it reaches zero without these difficulties and gets reasonably close to the curve with  $\alpha = 1$  as it departs from zero [15].

Using Equations 7 and 40 to specify a transition of  $V_{dB}$  decibels:

$$T' = \left\{ t_i 10^{\frac{V_{dB}}{20} \left( \frac{i}{\Lambda-1} \right)^\alpha} \right\}_0^{\Lambda-1} \quad (41)$$

in the general case of amplitude variations following a geometric progression. The greater the value of  $\alpha$ , the smoother the sound introduction and more intense its end.  $\alpha > 1$  results in loudness transitions commonly called *slow fade*, while  $\alpha < 1$  results in *fast fade* [36].

For linear amplification – but not linear perception – it is sufficient to use an appropriate sequence  $\{a_i\}$ :

$$a_i = a_0 + (a_{\Lambda-1} - a_0) \frac{i}{\Lambda - 1} \quad (42)$$

The linear transitions will be used for AM and FM synthesis, while exponential transitions are proper for tremolos and vibratos, as developed in Section 3.5. A non-oscillatory exploration of these variations is in the music piece *ParaMeter transitions* [27].

### 3.3 Application of digital filters

This subsection is limited to a description of sequences processing by convolution and difference equations, and immediate applications, as a thorough discussion of filtering is beyond the scope of this study<sup>11</sup>. With this procedure it is possible to achieve reverberators, equalizers, *delays*, to name a few of a variety of other filters for sound processing used to obtain musical/artistic effects. Filter employment can be part of the synthesis process or made subsequently as part of processes commonly referred to as “acoustic/sound treatment”.

**3.3.1 Convolution and finite impulse response (FIR) filters.** Filters applied by means of convolution are known by the acronym FIR (Finite Impulse Response) and are characterized by having a finite sample representation. This sample representation is called ‘impulse response’  $\{h_i\}$ . FIR filters are applied in the time domain by means of convolution of the sound with the respective impulse response of the filter. For the purposes of this work, convolution of  $T$  with  $H$  is defined as:

<sup>11</sup>The implementation of filters encompasses an area of recognized complexity, with dedicated literature and software [51, 68].



Fig. 16. Graphical interpretation of convolution. Each resulting sample is the sum of the previous samples of a signal, with each one multiplied by the retrograde of the other sequence.

$$\begin{aligned}
 \{t'_i\}_0^{\Lambda_t + \Lambda_h - 2} &= \{(T * H)_i\}_0^{\Lambda_{t'} - 1} = \{(H * T)_i\}_0^{\Lambda_{t'} - 1} \\
 &= \left\{ \sum_{j=0}^{\min(\Lambda_h - 1, i)} h_j t_{i-j} \right\}_0^{\Lambda_{t'} - 1} \\
 &= \left\{ \sum_{j=\max(i+1-\Lambda_h, 0)}^i t_j h_{i-j} \right\}_0^{\Lambda_{t'} - 1}
 \end{aligned} \tag{43}$$

where  $t_i = 0$  for the samples not given. In other words, the sound  $\{t'_i\}$ , resulting from the convolution of  $\{t_i\}$ , with the impulse response  $\{h_i\}$ , has each  $i$ -th sample  $t_i$  overwritten by the sum of its last  $\Lambda_h$  samples  $\{t_{i-j}\}_{j=0}^{\Lambda_h-1}$  multiplied one-by-one by samples of the impulse response  $\{h_i\}_0^{\Lambda_h-1}$ . This procedure is illustrated in Figure 16, where the impulse response  $\{h_i\}$  is in its retrograde form, and  $t'_{12}$  and  $t'_{32}$  are two samples calculated using the convolution given by  $(T * H)_i = t'_i$ . The final signal always has the length of  $\Lambda_t + \Lambda_h - 1 = \Lambda_{t'}$ . It is also possible to apply the filter by multiplying the Fourier coefficients of both the sound and the impulse response, and then performing the inverse Fourier transform [51]. This application of the filter in the frequency domain is usually much faster especially when using a Fast Fourier Transform (FFT) routine.

The impulse response can be provided by physical measurement or by pure synthesis. An impulse response for a reverberation, for example, can be obtained by recording the sound of the environment when someone triggers a click which resembles an impulse, or obtained by a sinusoidal sweep whose Fourier transform approximates its frequency response. Both are impulse responses which, properly convoluted with the sonic sequence, result in the same sound with a reverberation that resembles the original environment where the measurement was made [15].



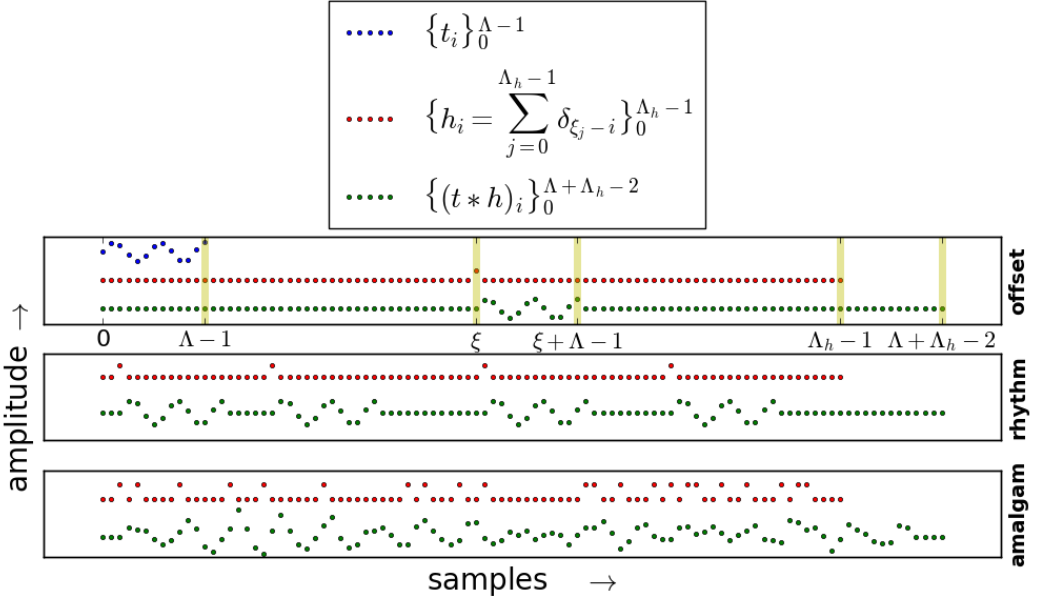


Fig. 17. Convolution with different densities of impulses: shifting (a), delay lines (b) and granular synthesis (c). The vertical axis is related to amplitude although one should keep in mind that each subplot has two or three displaced signals.

The Fourier transform of an impulse response of a FIR filter is an even and real envelope. Convolved with a sound (in the time or frequency domain), it performs the frequency filtering specified by the envelope. The greater the number of samples, the higher the envelope resolution and the computational complexity, which should often be weighted, for convolution is expensive.

An important property is the time shift caused by convolution with a shifted impulse. Despite being computationally expensive, it is possible to create *delay lines* by means of a convolution with an impulse response that has an impulse for each intended re-incidence of the sound. Figure 17 shows the shift caused by convolution with an impulse. Depending on the density of the impulses, the result is perceived as rhythm (from an impulse for each couple of seconds to about 20 impulses per second) or as pitch (from about 20 impulses per second and higher densities). In the latter case, the process is considered e.g. granular synthesis, reverberation or equalization.

**3.3.2 Infinite impulse response (IIR) filters.** This class of filters, known by the acronym IIR, is characterized by having an infinite time representation, i.e. the impulse response does not converge to zero. Its application is usually made by the following equation:

$$t'_i = \frac{1}{b_0} \left( \sum_{j=0}^J a_j t_{i-j} + \sum_{k=1}^K b_k t'_{i-k} \right) \quad (44)$$

The variables may be normalized:  $a'_j = \frac{a_j}{b_0}$  and  $b'_k = \frac{b_k}{b_0} \Rightarrow b'_0 = 1$ . Equation 44 is called ‘difference equation’ because the resulting samples  $\{t'_i\}$  are given by weighted differences between original samples  $\{t_i\}$  and previous ones in the resulting signal  $\{t'_{i-k}\}$ .

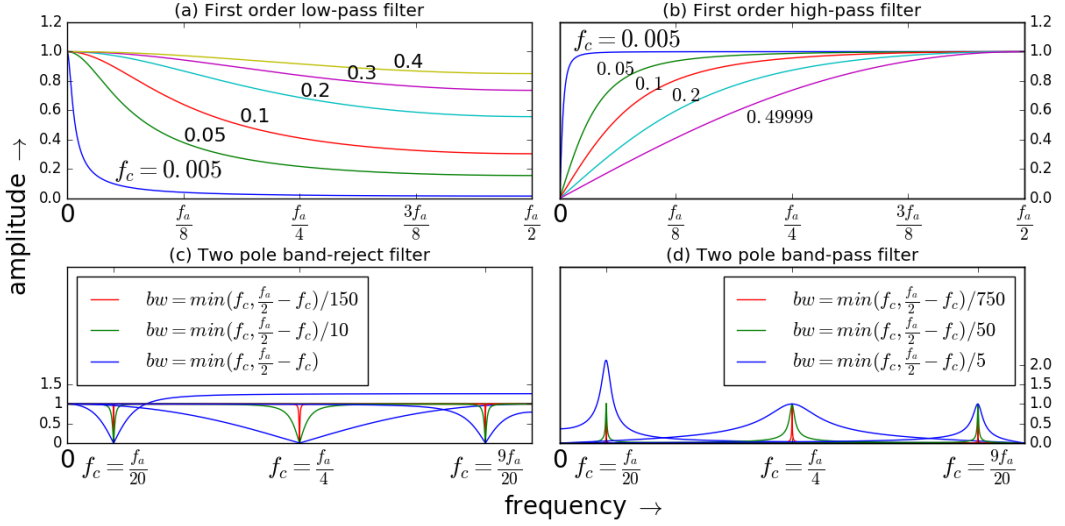


Fig. 18. Moduli for the frequency response (a), (b), (c) and (d) for IIR filters of Equations 45, 46, 48 and 49 respectively, considering different cutoff frequencies, center frequencies and bandwidth.

There are many methods and tools to obtain IIR filters. The text below lists a selection for didactic purposes and as a reference. They are well behaved filters and their main characteristics are described in Figure 18. For filters of simple order, the cutoff frequency  $f_c$  is where the filter performs an attenuation of  $-3dB \approx 0.707$  of the original amplitude. For band-pass and band-reject (or 'notch') filters, this attenuation has two specifications:  $f_c$  (in this case, the 'center frequency') and bandwidth  $bw$ . In both frequencies  $f_c \pm bw$  there is an attenuation of  $-3dB \approx 0.707$  of the original amplitude. There is sound amplification in band-pass and band-reject filters when the cutoff frequency is low and the bandwidth is large enough. In trebles, these filters present only a deviation of the expected profile, extending the envelope to the bass.

It is possible to apply filters successively in order to obtain filters with other frequency responses. Another possibility is to use a biquad 'filter recipe'<sup>12</sup> or the calculation of Chebichev filter coefficients<sup>13</sup>. Both alternatives are explored by [68, 70], and by the collection of filters maintained by the *Music-DSP* community of the Columbia University [13, 51].

- (1) Low-pass with a simple pole, module of the frequency response in the upper left corner of Figure 18. The general equation has the cutoff frequency  $f_c \in (0, \frac{1}{2})$ , fraction of the sample frequency  $f_s$  in which an attenuation of  $3dB$  occurs. The coefficients  $a_0$  and  $b_1$  of the IIR filter are given by  $x \in [e^{-\pi}, 1]$ :

$$\begin{aligned} x &= e^{-2\pi f_c} \\ a_0 &= 1 - x \\ b_1 &= x \end{aligned} \tag{45}$$

- (2) High-pass filter with a simple pole, module of its frequency responses at the upper right corner of Figure 18. The general equation with cutoff frequency  $f_c \in (0, \frac{1}{2})$  is calculated by

<sup>12</sup>Short for 'biquadratic': its transfer function has two poles and two zeros, i.e. its first direct form consists of two quadratic polynomials in the fraction:  $\mathbb{H}(z) = \frac{a_0 + a_1 \cdot z^{-1} + a_2 \cdot z^{-2}}{1 - b_1 \cdot z^{-1} - b_2 \cdot z^{-2}}$ .

<sup>13</sup>Butterworth and Elliptical filters can be considered as special cases of Chebichev filters [51, 68].

means of  $x \in [e^{-\pi}, 1]$ :

$$\begin{aligned} x &= e^{-2\pi f_c} \\ a_0 &= \frac{x+1}{2} \\ a_1 &= -\frac{x+1}{2} \\ b_1 &= x \end{aligned} \tag{46}$$

- (3) Notch filter. This filter is parametrized by a center frequency  $f_c$  and bandwidth  $bw$ , both given as fractions of  $f_s$ , therefore  $f, bw \in (0, \frac{1}{2})$ . Both frequencies  $f_c \pm bw$  have  $\approx 0.707$  of the amplitude, i.e. an attenuation of 3dB. The auxiliary variables  $K$  and  $R$  are:

$$\begin{aligned} R &= 1 - 3bw \\ K &= \frac{1 - 2R \cos(2\pi f_c) + R^2}{2 - 2 \cos(2\pi f_c)} \end{aligned} \tag{47}$$

The band-pass filter in the lower left corner of Figure 18 has the following coefficients:

$$\begin{aligned} a_0 &= 1 - K \\ a_1 &= 2(K - R) \cos(2\pi f_c) \\ a_2 &= R^2 - K \\ b_1 &= 2R \cos(2\pi f_c) \\ b_2 &= -R^2 \end{aligned} \tag{48}$$

The coefficients of band-reject filter, depicted in the lower right of Figure 18, are:

$$\begin{aligned} a_0 &= K \\ a_1 &= -2K \cos(2\pi f_c) \\ a_2 &= K \\ b_1 &= 2R \cos(2\pi f_c) \\ b_2 &= -R^2 \end{aligned} \tag{49}$$

### 3.4 Noise

Sounds without an easily recognizable pitch are generally called noise [43]. They are important musical sounds, as noise is present in real notes, e.g. emitted by a violin or a piano. Furthermore, many percussion instruments do not exhibit an unequivocal pitch and their sounds are generally regarded as noise [57]. In electronic music, including electro-acoustic and dance genres, noise has diverse uses and frequently characterizes the music style [15].

The absence of a definite pitch is due to the lack of a perceptible harmonic organization in the sinusoidal components of the sound. Hence, there are many ways to generate noise. The use of random values to generate the sound sequence  $T$  is a trivial method but not outstandingly useful because it tends to produce white noise with little or no variations [15]. Another possibility to generate noise is by using the desired spectrum, from which it is possible to perform the inverse Fourier transform. The spectral distribution should be done with care: if phases of components present prominent correlation, the synthesized sound will concentrate energy in some portions of its duration.

Some noises with static spectra are listed below. They are called *colored noise* since they are associated with colors for many reasons. Figure 19 shows the spectral profile and the corresponding

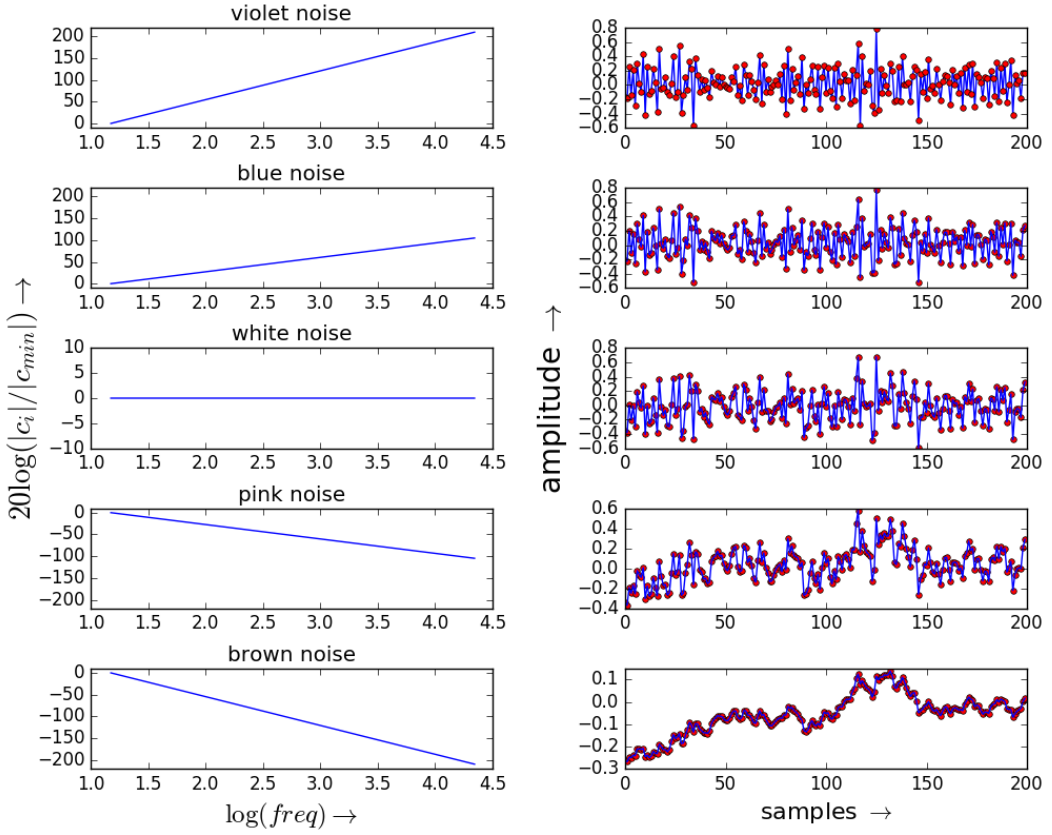


Fig. 19. Colors of noise generated by Equations 50, 51, 52, 53 and 54: spectrum and example waveforms.

sonic sequence side-by-side. All five noises were generated with the same phase for each component, making it straightforward to observe the contributions of different parts of the spectrum.

- The white noise has this name because its energy is distributed equally among all frequencies, such as the white color. It is possible to obtain white noise with the inverse transform of the following coefficients:

$$\begin{aligned}
 f_{\min} &\approx 15\text{Hz} \\
 f_i &= i \frac{f_s}{\Lambda}, \quad i \leq \frac{\Lambda}{2}, \quad i \in \mathbb{N} \\
 c_i &= 0, \quad \forall i : f_i < f_{\min} \\
 c_i &= e^{j \cdot x}, \quad x \text{ random} \in [0, 2\pi], \quad \forall i : f_{\min} \leq f_i < f_{\lceil \Lambda/2 - 1 \rceil} \\
 c_{\Lambda/2} &= 1, \quad \text{if } \Lambda \text{ even} \\
 c_i &= c_{\Lambda-i}^*, \quad \text{for } i > \frac{\Lambda}{2}
 \end{aligned} \tag{50}$$

The minimum frequency  $f_{\min}$  is chosen considering that a sound component with frequency below  $\approx 20\text{Hz}$  is usually inaudible. The exponential  $e^{j \cdot x}$  is a way to obtain unitary module

and random phase for the value of  $c_i$ . In addition,  $c_{\Lambda/2}$  is always real (as discussed in the previous section).

Other noises can be made by a similar procedure. In the following equations, the same coefficients are used and weighted using  $\alpha_i$ .

- The pink noise is characterized by a decrease of  $3dB$  per octave. This noise is useful for testing electronic devices, being prominent in nature [57].

$$\begin{aligned}\alpha_i &= \left(10^{-\frac{3}{20}}\right)^{\log_2\left(\frac{f_i}{f_{\min}}\right)} \\ c_i &= e^{j \cdot x} \alpha_i, \quad x \text{ random} \in [0, 2\pi], \quad \forall i : f_{\min} \leq f_i < f_{\lceil \Lambda/2 - 1 \rceil} \\ c_{\Lambda/2} &= \alpha_{\Lambda/2}, \text{ if } \Lambda \text{ even}\end{aligned}\tag{51}$$

- The brown noise (also Brownian noise) received this name after Robert Brown, who described the Brownian movement<sup>14</sup>. What characterizes brown noise is the decrease of  $6dB$  per octave, with  $\alpha_i$  in Equations 51 being:

$$\alpha_i = \left(10^{-\frac{6}{20}}\right)^{\log_2\left(\frac{f_i}{f_{\min}}\right)}\tag{52}$$

- In the blue noise there is a gain of  $3dB$  per octave in a band limited by the minimum frequency  $f_{\min}$  and the maximum frequency  $f_{\max}$ . Therefore (also based on the Equations 51):

$$\begin{aligned}\alpha_i &= \left(10^{\frac{3}{20}}\right)^{\log_2\left(\frac{f_i}{f_{\min}}\right)} \\ c_i &= 0, \quad \forall i : f_i < f_{\min} \text{ or } f_i > f_{\max}\end{aligned}\tag{53}$$

- The violet noise is similar to the blue noise, but its gain is  $6dB$  per octave:

$$\alpha_i = \left(10^{\frac{6}{20}}\right)^{\log_2\left(\frac{f_i}{f_{\min}}\right)}\tag{54}$$

- The black noise has higher losses than  $6dB$  for octave:

$$\alpha_i = \left(10^{-\frac{\beta}{20}}\right)^{\log_2\left(\frac{f_i}{f_{\min}}\right)}, \quad \beta > 6\tag{55}$$

- The gray noise is defined as a white noise subject to one of the ISO-audible curves. Such curves are obtained by experiments and are imperative to obtain  $\alpha_i$ . An implementation of ISO 226, which is the last established revision of these curves, is in the MASS toolbox as an auxiliary file [27].

This subsection discussed only noises with static spectra. There are also characterizations for noises with a dynamic spectrum along time, and noises which are fundamentally transient, like clicks and chirps. The former are easily modeled by an impulse relatively isolated, while a chirps is not in fact a noise, but a fast scan of some given frequency band [15].

<sup>14</sup>Although its origin is disparate with its color association, this noise became established with this specific name in musical contexts. Anyway, this association can be considered satisfactory once violet, blue, white and pink noises are more strident and associated with more vivid colors [15, 36].

### 3.5 Tremolo and vibrato, AM and FM

A vibrato is a periodic variation of pitch and a tremolo is a periodic variation of loudness<sup>15</sup>. A vibrato can be achieved by:

$$\gamma_i' = \left\lfloor i f' \frac{\tilde{\Lambda}_M}{f_s} \right\rfloor \quad (56)$$

$$t_i' = \tilde{m}_{\gamma_i' \% \tilde{\Lambda}_M} \quad (57)$$

$$f_i = f \left( \frac{f + \mu}{f} \right)^{t_i'} = f \cdot 2^{t_i' \frac{\nu}{12}} \quad (58)$$

$$\begin{aligned} \Delta_{\gamma_i} = \frac{\tilde{\Lambda}}{f_s} f_i &\Rightarrow \gamma_i = \left\lfloor \sum_{j=0}^i \frac{\tilde{\Lambda}}{f_s} f_j \right\rfloor \\ &= \left\lfloor \sum_{j=0}^i \frac{\tilde{\Lambda}}{f_s} f \left( \frac{f + \mu}{f} \right)^{t_j'} \right\rfloor \\ &= \left\lfloor \sum_{j=0}^i \frac{\tilde{\Lambda}}{f_s} f \cdot 2^{t_j' \frac{\nu}{12}} \right\rfloor \end{aligned} \quad (59)$$

$$T^{f, vibr}(f', \nu) = \left\{ t_i^{f, vibr}(f', \nu) \right\}_0^{\Lambda-1} = \left\{ \tilde{t}_{\gamma_i \% \tilde{\Lambda}} \right\}_0^{\Lambda-1} \quad (60)$$

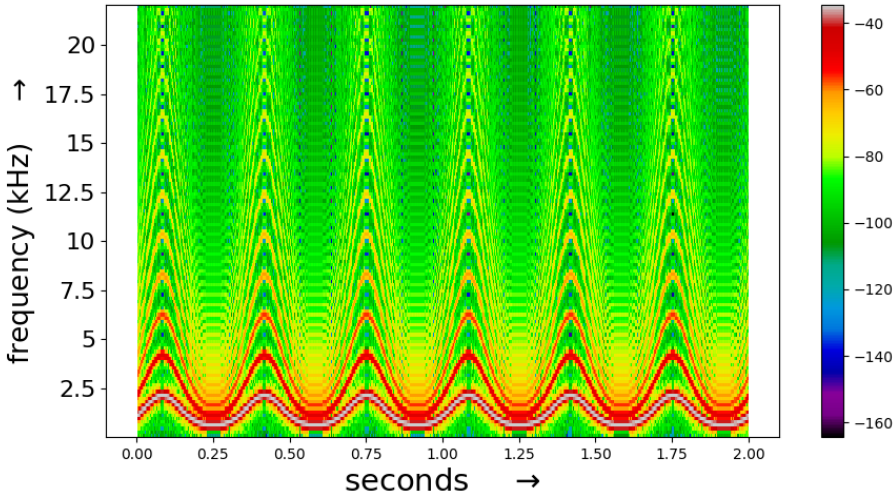


Fig. 20. Spectrogram of a sound with a sinusoidal vibrato of 3Hz and one octave of depth in a 1000Hz sawtooth wave ( $f_s = 44.1\text{kHz}$ ). The color bar is in decibels.

<sup>15</sup>The jargon may be different in other contexts. For example, in piano music, a tremolo is a vibrato in the classification used here. The definitions used in this document are usual in contexts regarding music theory and electronic music, i.e. they are based on a broader literature than the one used for a specific instrument, practice or musical tradition [43, 62].

For the proper realization of the vibrato, it is important to pay attention to both tables and sequences. Table  $\tilde{M}$  with length  $\tilde{\Lambda}_M$  and the sequence of indices  $\gamma'_i$  make the sequence  $t'_i$  which is the oscillatory pattern in the frequency while table  $\tilde{L}$  with length  $\tilde{\Lambda}$  and the sequence of indices  $\gamma_i$  make  $t_i$  which is the sound itself. Variables  $\mu$  and  $\nu$  quantify the vibrato intensity:

- $\mu$  is a direct measure of how many Hertz are involved in the upper limit of the oscillation, while
- $\nu$  is the direct measure of how many semitones (or half steps) are involved in the oscillation ( $2\nu$  is the number of semitones between the upper and lower peaks of the frequency oscillations of the sound  $\{t_i\}$ ).

It is convenient to use  $\nu = \log_2 \frac{f+\mu}{f}$  in this case because the maximum frequency increase is not equivalent to the maximum frequency decrease. The maximum semitone/pitch displacement is the invariant quantity and is called 'vibrato depth'. Most often, a vibrato depth is specified in semitones or cents (one cent =  $\frac{1}{100}$  of a semitone).

Figure 20 is the spectrogram of an artificial vibrato in a note with 1000Hz, in which the pitch deviation reaches one octave above and one below. Practically any waveform can be used to generate a sound and the vibrato oscillatory pattern, with virtually any oscillation frequency and pitch deviation. Such oscillations with precise waveforms and arbitrary amplitudes are not possible in traditional music instruments, and thus it introduces novelty in the artistic possibilities.

Tremolo is similar:  $f'$ ,  $\gamma'_i$  and  $t'_i$  remain the same. The amplitude sequence to be multiplied by the original sequence  $t_i$  is:

$$a_i = 10^{\frac{V_{dB}}{20} t'_i} = a_{\max}^{t'_i} \quad (61)$$

and, finally:

$$T^{tr(f')} = \left\{ t_i^{tr(f')} \right\}_0^{\Lambda-1} = \{ t_i \cdot a_i \}_0^{\Lambda-1} = \left\{ t_i \cdot 10^{t'_i \frac{V_{dB}}{20}} \right\}_0^{\Lambda-1} = \left\{ t_i \cdot a_{\max}^{t'_i} \right\}_0^{\Lambda-1} \quad (62)$$

where  $V_{dB}$  is the oscillation depth in decibels and  $a_{\max} = 10^{\frac{V_{dB}}{20}}$  is the maximum amplitude gain. The measurement in decibels is suitable because the maximum increase in amplitude is not equivalent to the maximum decrease, while the difference in decibels is preserved. Notice that the tremolo is applied to a preexisting sound and thus the characteristics of the tremolo do not need to be accounted for when synthesizing such sound (if it is synthesized) in contrast with making a sound with a vibrato.

Figure 21 shows the amplitude of the sequences  $\{a_i\}_0^{\Lambda-1}$  and  $\{t'_i\}_0^{\Lambda-1}$  for three oscillations of a tremolo with a sawtooth waveform. The curvature is due to the logarithmic progression of the intensity. The tremolo frequency is 1.5Hz if  $f_s = 44.1kHz$  because duration =  $\frac{t_{\max}-82000}{f_s} = 2s \Rightarrow \frac{3\text{oscillations}}{2s} = 1.5$  oscillations per second.

The musical piece *Shakes and wiggles* explores these possibilities given by tremolos and vibratos, both used in conjunction and independently (tremolos and vibratos occur many times together in a conventional music instrument), with different frequencies  $f'$ , depths ( $\nu$  and  $V_{dB}$ ), and progressive variations of parameters. Aiming at a qualitative appreciation, the piece also develops a comparison between vibratos and tremolos in logarithmic and linear scales. [27]

The proximity of  $f'$  to 20Hz generates roughness in both tremolos and vibratos. This roughness is largely appreciated both in traditional classical music and current electronic music, especially in the *Dubstep* genre. Roughness is also generated by spectral content that produces beating [52, 53].



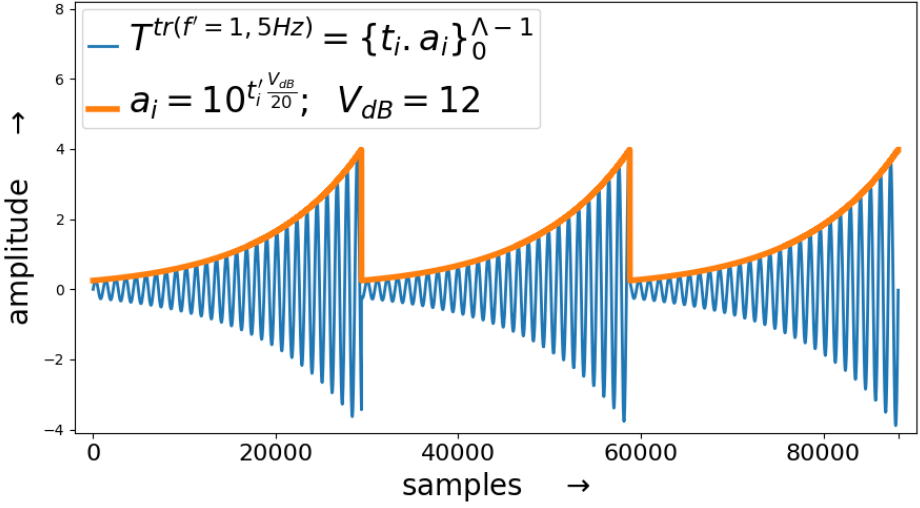


Fig. 21. Tremolo with a depth of  $V_{dB} = 12dB$ , with a sawtooth waveform as its oscillatory pattern, with  $f' = 1.5Hz$  in a sine of  $f = 40Hz$  ( $f_s = 44.1kHz$ ).

The sequence *Bela Rugosi* explores this roughness threshold with concomitant tremolos and vibratos at the same voice, with different intensities and waveforms. [27]

As the frequency increases further, these oscillations no longer remain noticeable individually. In this case, the oscillations become audible as pitch. Then,  $f'$ , the depths ( $\nu$  and  $V_{dB}$ ), and the waveform together change the audible spectrum of original sound  $T$  in different ways for tremolos and vibratos. They are called AM (*Amplitude Modulation*) and FM (*Frequency Modulation*) synthesis, respectively. These techniques are well known, with applications in synthesizers like *Yamaha DX7*, and even with applications outside music, as in telecommunications for data transfer by means of electromagnetic waves (e.g. AM and FM radios).

For musical goals, it is possible to understand FM based on the case of sines and, when other waveforms are employed, to consider the signals by their respective Fourier components (i.e. sines as well). The FM synthesis performed with a sinusoidal vibrato of frequency  $f'$  and depth  $\mu$  in a sinusoidal sound  $T$  with frequency  $f$  generates bands centered around  $f$  and far from each other by  $f'$ :

$$\begin{aligned}
 \{t'_i\} &= \left\{ \cos \left[ f \cdot 2\pi \frac{i}{f_s - 1} + \mu \cdot \sin \left( f' \cdot 2\pi \frac{i}{f_s - 1} \right) \right] \right\} = \\
 &= \left\{ \sum_{k=-\infty}^{+\infty} J_k(\mu) \cos \left[ f \cdot 2\pi \frac{i}{f_s - 1} + k \cdot f' \cdot 2\pi \frac{i}{f_s - 1} \right] \right\} = \\
 &= \left\{ \sum_{k=-\infty}^{+\infty} J_k(\mu) \cos \left[ (f + k \cdot f') \cdot 2\pi \frac{i}{f_s - 1} \right] \right\}
 \end{aligned} \tag{63}$$

where

$$J_k(\mu) = \frac{2}{\pi} \int_0^{\frac{\pi}{2}} \left[ \cos \left( \bar{k} \frac{\pi}{2} + \mu \cdot \sin w \right) \cdot \cos \left( \bar{k} \frac{\pi}{2} + k \cdot w \right) \right] dw, \quad \bar{k} = k \% 2, \quad k \in \mathbb{N} \tag{64}$$

is the Bessel function [65, 70] and specifies the amplitude of each component in an FM synthesis.

In these equations, the frequency variation introduced by  $\{t'_i\}$  does not follow the geometric progression that yields linear pitch variation, but reflects Equation 33. The result of using Equations 58 for FM is described in the Appendix D of [23], where the spectral content of the FM synthesis is calculated for oscillations in the logarithmic scale. In fact, the simple and attractive FM behavior is usually observed with linear oscillations, such as in Equation 63, which yield less strident and less noisy sounds.

For the amplitude modulation (AM):

$$\begin{aligned} \{t'_i\}_0^{\Lambda-1} &= \{(1 + a_i).t_i\}_0^{\Lambda-1} = \left\{ \left[ 1 + M. \sin \left( f'.2\pi \frac{i}{f_s - 1} \right) \right] . P. \sin \left( f.2\pi \frac{i}{f_s - 1} \right) \right\}_0^{\Lambda-1} = \\ &= \left\{ P. \sin \left( f.2\pi \frac{i}{f_s - 1} \right) + \frac{P.M}{2} \left[ \sin \left( (f - f').2\pi \frac{i}{f_s - 1} \right) + \sin \left( (f + f').2\pi \frac{i}{f_s - 1} \right) \right] \right\}_0^{\Lambda-1} \end{aligned} \quad (65)$$

The resulting sound is the original one together with the reproduction of its spectral content below and above with a displacement of  $f'$ . Again, this is achieved by variations in the linear scale (of the amplitude). The spectrum of an AM performed with oscillations in the logarithmic amplitude scale is described in Appendix D of [23]. The sequence  $T$ , with frequency  $f$ , called 'carrier', is modulated by  $f'$ , called 'modulator'. In FM and AM jargon,  $\mu$  and  $a_{max} = 10^{\frac{V_{dB}}{20}}$  are 'modulation indexes'. The following equations are defined for the oscillatory pattern of the modulator sequence  $\{t'_i\}$ :

$$\gamma'_i = \left\lfloor i f' \frac{\tilde{\Lambda}_M}{f_s} \right\rfloor \quad (66)$$

$$t'_i = \tilde{m}_{\gamma'_i \% \tilde{\Lambda}_M} \quad (67)$$

In FM, the modulator  $\{t'_i\}$  is applied to the carrier  $\{t_i\}$  by:

$$f_i = f + \mu.t'_i \quad (68)$$

$$\Delta_{\gamma_i} = f_i \frac{\tilde{\Lambda}}{f_s} \Rightarrow \gamma_i = \left\lfloor \sum_{j=0}^i f_j \frac{\tilde{\Lambda}}{f_s} \right\rfloor = \left\lfloor \sum_{j=0}^i \frac{\tilde{\Lambda}}{f_s} (f + \mu.t'_j) \right\rfloor \quad (69)$$

$$T^{f, FM(f', \mu)} = \left\{ t_i^{f, FM(f', \mu)} \right\}_0^{\Lambda-1} = \left\{ \tilde{l}_{\gamma_i \% \tilde{\Lambda}} \right\}_0^{\Lambda-1} \quad (70)$$

where  $\tilde{l}$  is the waveform period with a length of  $\tilde{\Lambda}$  samples, used for the carrier signal.

To perform AM, the signal  $\{t_i\}$  needs to be modulated with  $\{t'_i\}$  using the following equations:

$$a_i = 1 + \alpha.t'_i \quad (71)$$

$$T^{f, AM(f', \alpha)} = \left\{ t_i^{f, AM(f', \alpha)} \right\}_0^{\Lambda-1} = \{t_i.a_i\}_0^{\Lambda-1} = \{t_i.(1 + \alpha.t'_i)\}_0^{\Lambda-1} \quad (72)$$

### 3.6 Musical usages

At this point the musical possibilities are very wide. Sonic characteristics, like pitch (given by frequency), timbre (achieved by waveforms, filters and noise) and loudness (manipulated by intensity) can be considered in an absolute form or varied throughout the duration of a sound or a musical piece. The following musical usages encompass a collection of possibilities with the purpose of exemplifying types of sonic manipulations that result in musical material. Some of them are discussed more deeply in the next section.

**3.6.1 Relations between characteristics.** This is a widespread procedure used to obtain musically attractive and coherent excerpts. A possibility is to establish relations between parameters of tremolos and vibratos, and of the basic note like frequency. Let a vibrato frequency be proportional to note pitch, or a tremolo depth be inversely proportional to pitch. Therefore, with Equations 56, 58 and 61:

$$\begin{aligned} f^{vbr} &= f^{tr} = func_a(f) \\ v &= func_b(f) \\ V_{dB} &= func_c(f) \end{aligned} \quad (73)$$

with  $f^{vbr}$  and  $f^{tr}$  as  $f'$  in the referenced equations.  $v$  and  $V_{dB}$  are the respective depth values of vibrato and tremolo. Functions  $func_a$ ,  $func_b$  and  $func_c$  are arbitrary and dependent on musical intentions. The music piece *Bonds* explores such bonds and exhibits variations in the waveforms with the purpose of building a *musical language* (details in Section 4). [27]

**3.6.2 Convolution for rhythm and meter.** A musical pulse - such as specified by a BPM tempo - can be implied by an impulse at the start of each beat: the convolution with an impulse shifts the sound to impulse position, as stated in Section 3.3.1. For example, two impulses equally spaced build a binary division of the pulse. Two signals, one with 2 impulses and the other with 3 impulses, both equally spaced in the pulse duration, yield a pulse maintenance with a rhythm which eases both binary or ternary divisions. This is found in many ethnic and traditional musical styles [32]. The absolute values of the impulses entail proportions among the amplitudes of the sonic re-incidences. The use of convolution with impulses in this context is explored in the music piece *Little train of impulsive hillbillies*. These procedures also encompass the creation of 'sound amalgams' based on granular synthesis; see Figure 24. [27]

**3.6.3 Moving source and receptor, Doppler effect.** According to the discussion in Section 2.7, when an audio source (or receptor) is moving, the IID and ITD are constantly changing and are ideally updated at each sample of the digital signal (if fast computational rendering is not at stake). As given by basic theory, the audio source speed  $s_s$ , with positive values if the source moves away from receptor, and receptor speed  $s_r$ , positive when it gets closer to audio source (one might always use  $s_r = 0$  for musical purposes), relates the frequency  $f$  as perceived by the receiver and the frequency  $f_0$  emitted by:

$$f = \left( \frac{s_{sound} + s_r}{s_{sound} + s_s} \right) f_0 \quad (74)$$

Using the coordinates as in Figure 11, and Equation 25, the speed  $s_s$  can be found simply by  $s_s = f_s(d_{i+1} - d_i)$ . One should also use IID for the intensity progression of the sound, and ITD to correctly start and end the sonic sequences related to each ear. The change in pitch is antisymmetric upon the crossing of source with receptor: the same semitones (or fraction of) that are added during

the approach are decreased during the departure. Moreover, the transition is abrupt if source and receptor intersect with zero distance, otherwise, there is a smooth progression.

The musical piece *Doppeleer* explores and exemplifies the musical use of the Doppler effect. [27]

**3.6.4 Filters and noises.** With the use of filters, the possibilities are even wider. Convolve a signal to have a reverberated version of it, to remove its noise, to distort or to handle the audio aesthetically in many other ways. For example, sounds originated from an old television or telephone can be simulated with a band-pass filter, allowing only frequencies between  $1kHz$  and  $3kHz$ . By rejecting the frequency of an electric oscillation (usually  $50Hz$  or  $60Hz$ ) and the harmonics, one can remove noises caused by audio devices connected to the power supply. A more musical application is to perform filtering in specific bands and to use those bands as an additional parameter to the notes.

Inspired by traditional music instruments, it is possible to apply a time-dependent filter [57]. Chaining such filters can be useful for performing complex and more accurate filtering routines. The musical piece *Noisy band* explores filters and many kinds and noise synthesis. [27]

A sound can be altered through different filtering processes and then mixed to create an effect known as *chorus*. Based on what happens in a choir of singers, the sound is synthesized using small and potentially arbitrary modifications of parameters like center frequency, presence (or absence) of vibrato or tremolo and its characteristics, equalization, loudness, etc. As a final result, those versions of the original sound are mixed together (see Equation 30). The musical piece *Children choir* implements a very simple chorus and applies it to structures described in the next section. [27]

**3.6.5 Reverberation.** Using the same terms of Section 2.7, the late reverberation can be achieved by a convolution with a section of pink, brown or black noise, with an exponential decay of amplitude along time. Delay lines can be added as a prefix to the noise with the decay, and this accounts for both time parts of the reverberation: the early reflections and the late reverberation. Quality can be improved by varying the geometric trajectory and filtering by each surface where the wavefront reflected before reaching the ear in the first  $100 - 200ms$  (mainly with a LP). The colored noise can be gradually introduced with a *fade-in*: the initial moment given by direct incidence of sound (i.e. without any reflection and given by ITD and IID), reaching its maximum at the beginning of the 'late reverberation', when the geometric incidences loose their relevance to the statistical properties of the decaying noise. As an example, consider  $\Delta_1$  as the duration of the first reverberation section and  $\Delta_R$  as the complete duration of the reverberation ( $\Lambda_1 = \Delta_1 f_s$ ,  $\Lambda_R = \Delta_R f_s$ ). Let  $p_i$  be the probability of a sound to be repeated in the  $i$ -th sample. Following Section 2.7, the sequence  $R^1$  with the amplitudes of the impulse response of the first period can be described as:

$$R^1 = \{r_i^1\}_{0}^{\Lambda_1-1}, \text{ where } r_i^1 = \begin{cases} 10^{\frac{V_{dB}}{20} \frac{i}{\Lambda_R-1}} & \text{with probability } p_i = \left(\frac{i}{\Lambda_1}\right)^2 \\ 0 & \text{with probability } 1 - p_i \end{cases} \quad (75)$$

where  $V_{dB}$  is the total decay in decibels, typically  $-80dB$  or  $-120dB$ . The sequence  $R^2$  with the samples of the impulse response of the second period can be obtained from a brown noise  $N^b$  (or by a pink noise  $N^p$ ) with an exponential amplitude decay of the waveform:

$$R^2 = \{r_i^2\}_{\Lambda_1}^{\Lambda_R-1} = \left\{ 10^{\frac{V_{dB}}{20} \frac{i}{\Lambda_R-1}} \cdot r_i^b \right\}_{\Lambda_1}^{\Lambda_R-1} \quad (76)$$

Finally:

$$R = \{r_i\}_0^{\Lambda_R-1}, \text{ where } r_i = \begin{cases} r_i^1 & \text{if } 0 \leq i < \Lambda_1 - 1 \\ r_i^2 & \text{if } \Lambda_1 \leq i < \Lambda_R - 1 \end{cases} \quad (77)$$

A sound with an artificial reverberation can be achieved by a simple convolution of  $R$  (called reverberation impulse response) with the sound sequence  $T$ , as described in Section 3.3. Reverberation is well known for causing great interest in listeners and to provide sonorities that are more enjoyable. Furthermore, modifications in the reverberation consist in a common technique (almost a *cliché*) to surprise and attract the listener. The musical piece *Re-verb* explores reverberations in various settings. [27]

**3.6.6 ADSR envelopes.** The variation of loudness along the duration of a sound is crucial to our timbre perception. The intensity envelope known as ADSR (*Attack-Decay-Sustain-Release*) has many implementations in both hardware and software synthesizers. A pioneering implementation can be found in the Hammond Novachord synthesizer of 1938 and some variants are mentioned below [55]. The canonical ADSR envelope is characterized by 4 parameters: attack duration (time at which the sound reaches its maximum amplitude), decay duration (follows the attack immediately), level of sustained intensity (in which the intensity remains stable after the decay) and release duration (after sustained section, this is the duration needed for amplitude to reach zero or final value). Note that the sustain duration is not specified because it is the difference between the total duration and the sum of the attack, decay and release durations.

The ADSR envelope with durations  $\Lambda_A$ ,  $\Lambda_D$  and  $\Lambda_R$ , with total duration  $\Lambda$  and sustain level  $a_S$ , given as the fraction of the maximum amplitude, to be applied to any sound sequence  $T = \{t_i\}$  (ideally also with duration  $\Lambda$ ), can be expressed as:

$$\begin{aligned} \{a_i\}_0^{\Lambda_A-1} &= \left\{ \xi \left( \frac{1}{\xi} \right)^{\frac{i}{\Lambda_A-1}} \right\}_0^{\Lambda_A-1} \quad \text{or} \\ &= \left\{ \frac{i}{\Lambda_A - 1} \right\}_0^{\Lambda_A} \\ \{a_i\}_{\Lambda_A}^{\Lambda_A+\Lambda_D-1} &= \left\{ a_S^{\frac{i-\Lambda_A}{\Lambda_D-1}} \right\}_{\Lambda_A}^{\Lambda_A+\Lambda_D-1} \quad \text{or} \\ &= \left\{ 1 - (1 - a_S) \frac{i - \Lambda_A}{\Lambda_D - 1} \right\}_{\Lambda_A}^{\Lambda_A+\Lambda_D-1} \\ \{a_i\}_{\Lambda_A+\Lambda_D}^{\Lambda-\Lambda_R-1} &= \{a_S\}_{\Lambda_A+\Lambda_D}^{\Lambda-\Lambda_R-1} \\ \{a_i\}_{\Lambda-\Lambda_R}^{\Lambda-1} &= \left\{ a_S \left( \frac{\xi}{a_S} \right)^{\frac{i-(\Lambda-\Lambda_R)}{\Lambda_R-1}} \right\}_{\Lambda-\Lambda_R}^{\Lambda-1} \quad \text{or} \\ &= \left\{ a_S - a_S \frac{i + \Lambda_R - \Lambda}{\Lambda_R - 1} \right\}_{\Lambda-\Lambda_R}^{\Lambda-1} \end{aligned} \quad (78)$$

with  $\Lambda_X = \lfloor \Lambda_X \cdot f_s \rfloor \quad \forall \quad X \in (A, D, R)$  and  $\xi$  being a small value that provides a satisfactory *fade in* and *fade out*, e.g.  $\xi = 10^{\frac{-80}{20}} = 10^{-4}$ . The lower the  $\xi$ , the slower the *fade*, similar to the  $\alpha$  illustrated in Figure 15. One might also use a linear or quartic ( $x^4$ ) fade at the beginning of the attack and the end of the release sections to reach zero amplitude (exponential fades never reach zero). Schematically, Figure 22 shows the ADSR envelope in a classical implementation that

supports many variations. For example, between attack and decay it is possible to add an extra section where the maximum amplitude remains for more than a peak. Another common example is the use of more elaborated envelopes for attack or decay. The music piece *ADa and SaRa* explores many configurations of the ADSR envelope. [27]

$$\{t_i^{ADSR}\}_0^{\Lambda-1} = \{t_i \cdot a_i\}_0^{\Lambda-1} \quad (79)$$

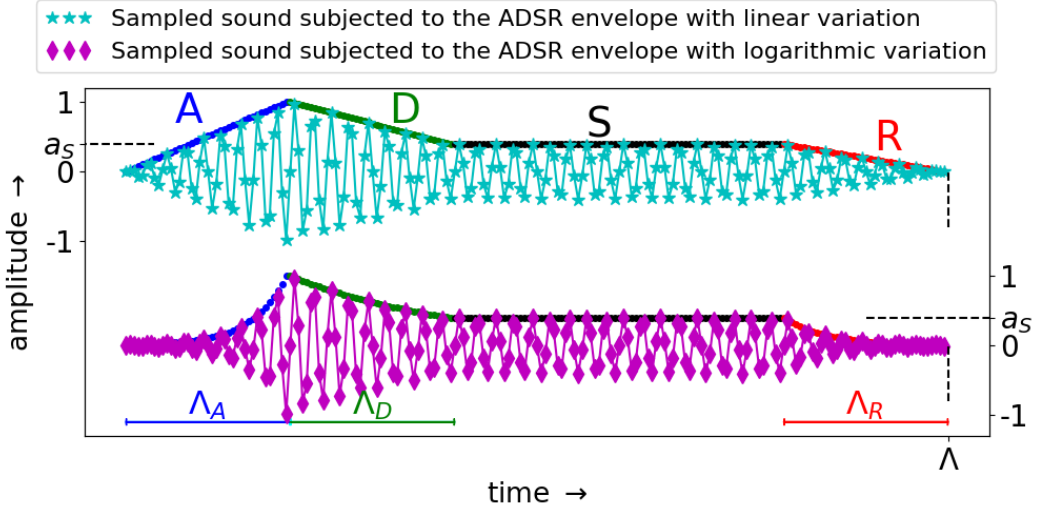


Fig. 22. An ADSR envelope (Attack, Decay, Sustain, Release) applied to an arbitrary sound sequence. The linear variation of the amplitude is above, in blue. Below the amplitude variation is exponential.

#### 4 ORGANIZATION OF NOTES IN MUSIC

Let  $S = \{s_j = T_i^j = \{t_i^j\}_{i=0}^{\Lambda_j-1}\}_{j=0}^{H-1}$  be the sequence of  $H$  musical events  $s_j$ . Consider  $S$  as a ‘musical structure’. This section is dedicated to techniques that make  $S$  interesting and enjoyable for hearing. More specifically, what follows is a summary of academic music composition theory and praxis. This section does not benefit from equations that dictate the amplitude of each sample as deeply as the previous sections. Even so, we understand that this content is very useful for synthesizing music and is not trivially integrated to Sections 2 and 3. The concepts are given algorithmic implementations in MASS [28] and can be further formalized [47], although at a cost of prompt intelligibility which we chose to avoid in this exposition.

The elements of  $S$  can be overlapped by mixing them together, as in Equation 30 and Figure 12, for building intervals and chords. This reflects the ‘vertical thought’ in music. On the other hand, the concatenation of events in  $S$ , as in Equation 31 and in Figure 13, yields melodic sequences and rhythms, which are associated with the ‘horizontal thought’. The fundamental frequency  $f$  and the starting moment (attack) are generally considered the most important characteristics of the elements  $s_j$ . These observations are convenient to describe and create music constituted by pitches and by temporal metrics and rhythms. We will start by considering such aspects of musical organization as they are more traditional in music theory and are usually easier to understand.

#### 4.1 Tuning, intervals, scales and chords

**4.1.1 Tuning.** Doubling the frequency is equivalent to ascending one octave ( $f = 2f_0$ ). The octave division in twelve pitches is the canon for classical western music. Its usage has also been observed outside western tradition, e.g. in ceremonial/religious and ethnic contexts [76]. The intervals between the pitches need not to be equivalent, as will become clear in the next paragraphs, but, roughly, the factor given by  $\varepsilon = 2^{\frac{1}{12}}$  defines a semitone, i.e. if  $f = 2^{\frac{1}{12}} f_0$ , there is a semitone between  $f_0$  and  $f$ . This entails a note grid along the spectrum in which, given a frequency  $f$ , any other fundamental frequency  $f'$  is related to  $f$  by  $f' = \varepsilon^i f$  where  $i$  is an integer. Twelve successive semitones yield an octave. Notice that equivalences of semitones and octaves are not absolute: 2 pitches related by an octave ( $f_2 = 2f_1$ ) are different at least because one is higher than the other, but are equivalent in the sense that they have similar uses and might be added or substituted in a sound without introducing much novelty or change; semitones might not be perceived as equivalent "distances" (this is dependent on context and listener), but are equivalent e.g. for transposing melodies, harmonies and other pitch-related structures.

The absolute accuracy of  $\varepsilon = 2^{\frac{1}{12}}$  is usual in computational implementations. Performances with real musical instruments, however, often present semitones that are not exactly  $2^{\frac{1}{12}}$  because the pitches yield by such grid do not match the harmonics. The fixed interval  $\varepsilon = 2^{\frac{1}{12}}$  characterizes an equally tempered tuning but there are other tunings. The first formalizations of tunings (that the scientific tradition has reported) date from around two thousand years before the advent of the equal temperament [57]. Two emblematic tunings are:

- The **just intonation**, defined by association of intervals with ratios of low-order integers, as found in the harmonic series. E.g. the white piano keys from C to C are achieved by the ratios of frequency: 1, 9/8, 5/4, 4/3, 3/2, 5/3, 15/8, 2/1. The semitone 16/15 is also often considered. There are many ways to perform the division of the 12 notes in the just intonation.
- The **Pythagorean tuning**, based on the interval 3/2 (perfect fifth). The 'white piano keys' become: 1, 9/8, 81/64, 4/3, 3/2, 27/16, 243/128, 2/1. Also often used are the 'minor second' 256/243, the 'minor third' 32/27, the 'augmented fourth' 729/512, the 'diminished fifth' 1024/729, the 'minor sixth' 128/81 and the 'minor seventh' 16/9.

In order to account for micro-tonality<sup>16</sup>, non-integer values can be used as factors of  $\varepsilon = 2^{\frac{1}{12}}$  between frequencies, or one can maintain the usage of integer values and change  $\varepsilon$ . For example, a tuning that approximates the harmonic series is proposed with the equal division of the octave in 53 notes:  $\varepsilon = 2^{\frac{1}{53}}$  [75]. Note that if  $S = \{s_i\}$  is a pitch sequence related by means of  $\varepsilon = 2^{1/\eta}$ , the sequence  $S'$  with the same notes, but related by  $\varepsilon' = 2^{1/\eta'}$ , is  $S' = \{s'_i\} = \left\{s_i \frac{\eta'}{\eta}\right\}$  because:

$$\begin{aligned}
 F &= \{f_i\} \\
 S &= \{s_i\} \Rightarrow f_i = f 2^{s_i/\eta} \\
 S' &= \{s'_i\} \Rightarrow f_i = f 2^{s'_i/\eta'} \\
 f_i &= f 2^{s_i/\eta} = f 2^{s'_i/\eta'} \Rightarrow s'_i = s_i \frac{\eta'}{\eta}
 \end{aligned} \tag{80}$$

<sup>16</sup>Micro-tonality is the use of intervals smaller than one semitone and has ornamental and structuring functionalities in music. The division of the octave in 12 notes has physical grounds but is still a *convention* adopted by western classical music. Other tunings are incident, e.g. a traditional Thai music style uses an octave division in seven notes equally spaced ( $\varepsilon = 2^{\frac{1}{7}}$ ), which allows intervals quite different than those found when  $\varepsilon = 2^{\frac{1}{12}}$  [76].



The music piece *Micro tone* exemplifies the use of microtonal features.

**4.1.2 Intervals.** Using the ratio  $\varepsilon = 2^{\frac{1}{12}}$  between note frequencies (i.e. one semitone) the intervals in the twelve tone system can be represented by integers. Table 1 summarizes the intervals: traditional notation, qualifications of consonance and dissonances, and number of semitones.

Table 1. Musical intervals: traditional notation, basic classification for dissonances and consonances, and number of semitones. Unison, fifth and octave are the perfect (P) consonances. Major (M) and minor (m) thirds and sixths are the imperfect consonances. Minor seconds and major sevenths are the harsh (also strong or sharp) dissonances. Major seconds and minor sevenths are the mild (also weak) dissonances. Perfect fourth is a special case, as it is a perfect consonance when considered as an inversion of the perfect fifth and a dissonance or an imperfect consonance otherwise. Another special case is the tritone (A4 or aug4, d5 or dim5, tri or TT). This interval is consonant in some cultures. For tonal music, the tritone indicates a dominant (chord, function or harmonic field, see Section 4.2) and seeks urgent resolution into a third or sixth, and due to this instability it is considered a dissonant interval.

<b>consonances</b>		
	traditional notation	number of semitones
perfect:	P1, P5, P8	0, 7, 12
imperfect:	m3, M3, m6, M6	3, 4, 8, 9
<b>dissonances</b>		
	traditional notation	number of semitones
strong:	m2, M7	1, 11
weak:	M2, m7	2, 10
<b>special cases</b>		
	traditional notation	number of semitones
consonance or dissonance:	P4	5
dissonance in Western tradition:	tritone, aug4, dim5	6

The nomenclature, based on conveniences for tonal music and practical aspects of manipulating notes, can be specified as follows [57, 76]:

- Intervals are inspected first by the number of steps between notes. The simple intervals (which are at most an octave wide) are: first (unison), second, third, fourth, fifth, sixth, seventh and eighth (octave). Each of these intervals are related to one step less the their names suggest: a third is an interval with two steps. As can be noticed in Table 1, one step is not one semitone. A step, in this sense, is yield by two consecutive notes in a musical scale. A scale for now can be regarded as any arbitrary monotonic sequence of pitches and will be discussed in the next section.
- The intervals are represented by numeric digits, e.g. 1, 3, 5 are a unison, a third and a fifth, respectively<sup>17</sup>.
- An interval wider than an octave (e.g. ninth, tenth, eleventh) is called a 'compound interval' and is classified in terms of the simple interval between the same notes but in the same octave. Their notation can be achieved by adding a multiple of 7 to the simple interval: P11 is an octave plus a forth ( $7 + P4 = P11$ ), M9 is an octave plus a major second ( $7 + M2 = M9$ ), m16 is two octaves and a minor second ( $2 \times 7 + m2 = m16$ ).

<sup>17</sup>Integers might also be used to express the number of semitones in an interval.

- Quality of each interval: perfect consonances – i.e. unison, fourth, fifth and octave – are ‘perfect’. The imperfect consonances – i.e. thirds and sixths – and dissonances – i.e. seconds and sevenths – can be major and minor. The tritone is an exception to this rule because it is a dissonant interval and cannot be major or minor.
- The perfect fourth can be a perfect consonance or a dissonance according to the context and theoretical background. As a general rule, it can be considered a consonance except when it is followed by a third or a fifth by the movement of the notes.
- The tritone is a dissonance in Western music because it is typical of the “dominant” chord (see Section 4.2) and represents (or yields) instability. Some cultures consider the interval a consonance and use it as a stable interval.
- A major interval decreased by one semitone results in a minor interval. A minor interval increased by one semitone results in a major interval.
- A perfect interval (P1, P4, P5, or P8), or a major interval (M2, M3, M6 or M7), increased by one semitone results in an augmented interval (e.g. aug3 has five semitones). The augmented fourth is also called tritone (aug4, tri, or TT).
- A perfect interval or a minor interval (m2, m3, m6 or m7), decreased by one semitone results in a diminished interval. The diminished fifth is also called tritone (dim5, tri, or TT).
- An augmented interval increased by one semitone results in a ‘doubly-augmented’ interval; a diminished interval decreased by one semitone results in a ‘doubly-diminished’ interval.
- Notes played simultaneously yield a harmonic interval.
- Notes played as a sequence in time yield a melodic interval. When the lowest note comes first there is an ascending interval, while a descending interval is observed when the highest note comes first.
- A simple interval is inverted if the lowest pitch is raised one octave, or if the highest pitch is lowered one octave. The sum of an interval and its inversion is 9 (e.g. m7 is inverted to M2:  $m7 + M2 = 9$ ). An inverted major interval results in a minor interval and vice-versa. An inverted augmented interval results in a diminished interval and vice-versa (inverting a doubly-augmented results in a doubly-diminished and vice-versa, etc). An inverted perfect interval is a perfect interval as well.

The augmented/diminished intervals and the doubly-augmented/doubly-diminished intervals have the same number of semitones of other intervals (e.g. minor, major or perfect) and are consequences of the tonal system. Scale notes are in fact different pitches, with specific uses and functions. Henceforth, in a *C flat* major scale, the tonic – first degree – is *C flat*, not *B*, and the leading tone – seventh degree – is *B flat*, not *A sharp* or *C double flat*. To grasp what this entails for intervals, let the second degree (second note) of a scale be one semitone from the first degree. Consider also the leading tone (i.e. the seventh degree at one ascending semitone from the first degree). There is a diminished third between the seventh and second scale degrees [43]. Notice that the dim3 is only two semitones wide, as is the major second (or e.g. an doubly-augmented unisson!).

This description summarizes the traditional nomenclature (or theory) of musical intervals [43]. The music piece *Intervals* explores these intervals in both independent and interrelated ways [27].

**4.1.3 Scales.** A scale is an ordered set of pitches. Strictly speaking, any (ordered) set of pitches can be considered a scale. The complexity of musical scales lean mostly on tradition, i.e. on the scales and their uses which result from practice throughout history. Usually, scales repeat at each octave. The ascending sequence with all notes from the octave division in 12 equal intervals ( $\varepsilon = 2^{\frac{1}{12}}$ ) is known as the chromatic scale within the equal temperament. There are 5 perfectly

symmetric divisions of the octave using the chromatic scale. These divisions are often regarded as scales themselves owing to the easy and peculiar uses they entail.

Let  $e_i$  be integers indexed by  $i$  such that  $f = \varepsilon^{e_i} f_0$ , where  $f_0$  is any fixed frequency. The symmetric scales mentioned above can be expressed as:

$$\begin{aligned}
 \text{chromatic} = E^c &= \{e_i^c\}_0^{11} = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11\} = \{i\}_0^{11} \\
 \text{whole tones} = E^{wt} &= \{e_i^{wt}\}_0^5 = \{0, 2, 4, 6, 8, 10\} = \{2i\}_0^5 \\
 \text{minor thirds} = E^{mt} &= \{e_i^{mt}\}_0^3 = \{0, 3, 6, 9\} = \{3i\}_0^3 \\
 \text{major thirds} = E^{Mt} &= \{e_i^{Mt}\}_0^2 = \{0, 4, 8\} = \{4i\}_0^2 \\
 \text{tritones} = E^{tt} &= \{e_i^{tt}\}_0^1 = \{0, 6\} = \{6i\}_0^1
 \end{aligned} \tag{81}$$

For example, the third note of the whole tone scale with  $f_0 = 200\text{Hz}$  is  $f_3 = \varepsilon^{e_3^{wt}} \cdot f_0 = 2^{\frac{4}{12}} \cdot 200 \approx 251.98\text{Hz}$ . These ‘scales’, or patterns, generate stable structures by their internal symmetries and can be repeated in a sustained way which is musically effective. Section 4.7 discusses other aspects of symmetries in music. The musical piece *Crystals* uses each one of these scales, in both melodic and harmonic constructions.

The *diatonic scales* have seven notes in which the consecutive intervals include five whole tones<sup>18</sup> and two semitones (in each octave). There are seven of them:

$$\begin{aligned}
 \text{aeolian} = \text{natural minor scale} &= \\
 &= E^m = \{e_i^m\}_0^6 = \{0, 2, 3, 5, 7, 8, 10\} \\
 \text{locrian} = E^{lo} &= \{e_i^{lo}\}_0^6 = \{0, 1, 3, 5, 6, 8, 10\} \\
 \text{ionian} = \text{major scale} &= \\
 &= E^M = \{e_i^M\}_0^6 = \{0, 2, 4, 5, 7, 9, 11\} \\
 \text{dorian} = E^d &= \{e_i^d\}_0^6 = \{0, 2, 3, 5, 7, 9, 10\} \\
 \text{phrygian} = E^p &= \{e_i^p\}_0^6 = \{0, 1, 3, 5, 7, 8, 10\} \\
 \text{lydian} = E^l &= \{e_i^l\}_0^6 = \{0, 2, 4, 6, 7, 9, 11\} \\
 \text{mixolydian} = E^{mi} &= \{e_i^{mi}\}_0^6 = \{0, 2, 4, 5, 7, 9, 10\}
 \end{aligned} \tag{82}$$

They have only major, minor and perfect intervals. The unique exception is the tritone found as an augmented fourth or a diminished fifth. Diatonic scales follow a circular pattern of successive intervals *tone, tone, semitone, tone, tone, tone, semitone*. Thus, it is possible to write:

$$\begin{aligned}
 \{d_i\} &= \{2, 2, 1, 2, 2, 2, 1\} \\
 e_0 &= 0 \\
 e_i &= d_{(i+\kappa)\%7} + e_{i-1} \quad \text{for } i > 0
 \end{aligned} \tag{83}$$

with  $\kappa \in \mathbb{N}$ . For each mode (of Equations 82) there is only one  $\kappa \in [0, 6]$  for which  $\{e_i\}$  (of Equations ??) matches. For example, a brief inspection reveals that  $e_i^l = d_{(i+2)\%7} + e_{i-1}^l$ . Thus,  $\kappa = 2$  for the lydian scale.

The minor scale have two additional forms, named harmonic and melodic:

<sup>18</sup>A (whole) tone is an interval that is two semitones wide:  $f' = 2^{\frac{2}{12}} f$ .

$$\begin{aligned}
\text{natural minor} &= E^m = \{e_i^m\}_0^6 = \\
&= \{0, 2, 3, 5, 7, 8, 10\} \\
\text{harmonic minor} &= E^{mh} = \{e_i^{mh}\}_0^6 = \\
&= \{0, 2, 3, 5, 7, 8, 11\} \\
\text{melodic minor} &= E^{mm} = \{e_i^{mm}\}_0^{14} = \\
&= \{0, 2, 3, 5, 7, 9, 11, 12, 10, 8, 7, 5, 3, 2, 0\}
\end{aligned} \tag{84}$$

The different ascending and descending contours of the melodic minor scale is required in tonal music. The minor scale has one whole tone between the seventh and eighth (or first) degrees but the separation by one semitone is critical to the polarization of the first degree. This is not necessary in the descending trajectory, and therefore the scale recovers the standard form. The harmonic scale presents the modified seventh degree but does not avoid the augmented second between the sixth and seventh degrees; it does not consider the melodic trajectory and thus does not need to avoid the aug2 [62].

Although it is not a traditional scale, the harmonic series is often used as such:

$$\begin{aligned}
H = \{h_i\}_0^{19} &= \\
&= \{0, 12, 19 + 0.02, 24, 28 - 0.14, 31 + 0.2, 34 - 0.31, \\
&\quad 36, 38 + 0.04, 40 - 0.14, 42 - 0.49, 43 + 0.02, \\
&\quad 44 + 0.41, 46 - 0.31, 47 - 0.12, \\
&\quad 48, 49 + 0.05, 50 + 0.04, 51 - 0.02, 52 - 0.14\}
\end{aligned} \tag{85}$$

In this scale, the frequency of the  $i$ th note  $h_i$  is the frequency of  $i$ th harmonic  $f_i = \varepsilon^{h_i} f_0$  from the spectrum generated by  $f_0$ . Natural sounds have such frequencies (as discussed in Section 2) usually with deviations from the expected values and with noise.

Many other scales can be expressed using the framework exposed in this section, e.g. the pentatonic scales and the modes of limited transposition of Messiaen [14].

One last observation: the words *scale* and *mode* are often used as synonyms both in the literature and in colloquial discussions. The word *mode* can also be used to mean two other things:

- an unordered set of pitches (i.e. an unordered scale).
- A scale used in the context of modal harmony, in the sense presented in Section 4.2.

**4.1.4 Chords.** A musical chord is implied by the simultaneous occurrence of three or more notes. Chords are often based on triads, especially in tonal music. Triads are built by two successive thirds within 3 notes: root, third and fifth. If the lower note of a chord is the root, the chord is in the root position, otherwise it is an inverted chord. A closed position is any in which no chord note fits between two consecutive notes [43], any non-closed position is an open position. In closed and fundamental positions, and with the fundamental denoted by 0, triads can be expressed as:

$$\begin{aligned}
\text{major triad} &= A^M = \{a_i^M\}_0^2 = \{0, 4, 7\} \\
\text{minor triad} &= A^m = \{a_i^m\}_0^2 = \{0, 3, 7\} \\
\text{diminished triad} &= A^d = \{a_i^d\}_0^2 = \{0, 3, 6\} \\
\text{augmented triad} &= A^a = \{a_i^a\}_0^2 = \{0, 4, 8\}
\end{aligned} \tag{86}$$

It is commonplace to consider another successive third: it is sufficient to include 10 as the highest note to achieve a tetrad with a minor seventh, or include 11 in order to achieve a tetrad with a

major seventh. Inversions and open positions can be obtained with the addition of  $\pm 12$  to the selected note. Incomplete triadic chords, with extra notes ('dirty' chords), and non-triadic are also common. These are often interpreted as the result of further extending the succession of thirds. E.g.  $\{0, 2, 4, 7\}$  will often be understood as a major chord with a major ninth (a major ninth has 14 semitones and  $14 - 12 = 2$ ).

For general guidance:

- A fifth confirms the root (fundamental). There are theoretical discussions about why this happens, and the most usual arguments are that the fifth is the first (non-octave) harmonic of a note and that the harmonics of the fifth are in the harmonics of the fundamental. Important here is to grasp the fact that musical theory and practice assures that the fifth establishes (or helps to establish) the fundamental as the root of a chord.
- Major or minor thirds from the root entails major or minor chord qualities.
- Every tritone, especially if built between a major third and a minor seventh, tends to resolve into a third or a sixth.
- Note duplication is avoided. If duplication is needed, the preference is, in descending order: the root, fifth, third and seventh.
- Note omission is avoided in the triad. If needed, the fifth is first considered for omission, then third and then the fundamental.
- It is possible to build chords with notes different from triads, particularly if they obey a recurrent logic or sequence that justifies these different notes.
- Chords built by successive intervals different from thirds – such as fourths and seconds – are recurrent in compositions of advanced tonalism or experimental music.
- The repetition of chord successions (or of characteristics they hold) fixes a trajectory and makes it possible to introduce exotic arrangements without implying in musical incoherence.

## 4.2 Atonal and tonal harmonies, harmonic expansion and modulation

Omission of basic tonal structures is the key to achieving modal and atonal harmonies. In the absence of minimal tonal organization, harmony is (usually) considered modal if the notes match some diatonic scale (see Equations 82) or if there is only a small number of notes. If basic tonal progressions are absent and notes do not match any diatonic scale and are sufficiently diverse and dissonant (between themselves) to avoid reduction of the notes by polarization<sup>19</sup>, the harmony is atonal. In this classification, the modal harmony is not tonal or atonal and is reduced to the incidence of notes within a (most often diatonic) scale and to the absence of tonal structures. Following this conceptualization, one observes that atonal harmony is hard to be realized and, indeed, no matter how dissonant and diverse a set of notes is, tonal harmonies arise very easily if not avoided laboriously [39].

**4.2.1 Atonal harmony.** In fact, atonal music techniques avoid that a direct relation of the notes with modes and tonality be established. Manifesting such atonal structures is of such difficulty that the dodecafonism emerged. Dodecafonism consists in the use a set of notes (ideally 12 notes) and the execution of each note, one by one, in the same order. In this context, the tonic becomes difficult to be established. Nevertheless, the western listener automatically searches for tonal elements in music and obstinately finds them by unexpected and tortuous paths. The use of dissonant intervals

<sup>19</sup>By polarization we mean having some notes that are way more important than others and to which the other notes are ornaments or subordinates.

(especially tritones) without resolution reinforces the absence of tonality. In this context, while creating a musical piece, it is allowed:

- To repeat pitches. By considering immediate repetition as an extension of the previous incidence, the use of the same pitch in sequence does not add relevant information.
- To play adjacent pitches (e.g. of a dodecafonic progression) at the same time, making harmonic intervals and chords.
- Use durations and pauses with freedom, respecting pitch order.
- Vary note sequences by temporal magnification and translation; or pitch transposition and sequence inversion, retrograde and retrograde inversion. See Sections 4.5 and 4.10 or specialized literature (such as [33]) for what these terms mean.
- Make variations in orchestration, articulation, spatialization, among other possibilities in presenting the same notes.

The atonal harmony can be observed, paradigmatically, within these presented conditions (which is a simple dodecafonic model). Most of what was written by great dodecafonic composers, e.g. Alban Berg and even Schoenberg, had the purpose of mixing tonal and atonal techniques. Most frequently, atonal music is not strictly dodecafonic, but "serial", i.e. they use the same kind of techniques based in (arbitrary) sequences (called the series or row) of pitches and other sonic characteristics.

**4.2.2 Tonal harmony.** In the XX century, music with emphasis on sonorities/timbres, and rhythm, extended the concepts of tonality and harmony. Even so, tonal harmony is very often in artistic movements and commercial venues. In addition, dodecafonism itself is sometimes considered of tonal nature because it was conceived to deny tonal characteristics of polarization, i.e. it is based on tonalism. In tonal or modal music, chords – like the ones listed in Equations 86 – built with the root at each degree of a scale (as listed in Equations 82) form the pillars of harmony. Tonal (and modal) harmony deals with chord formation and progressions. Even a monophonic melody entails harmonic fields, making it possible to perceive the chord progression even in unaccompanied melodies.

In the traditional tonal music, a scale has its tonic (first degree) on any note, and can be major (with the same notes of the Ionian mode) or minor (same notes of the Eolian mode, the 'natural minor', which has both harmonic and melodic versions, as in Equations 84). The scale is the base for triads, each with its root in a degree:  $\hat{1}, \hat{2}, \hat{3}, \hat{4}, \hat{5}, \hat{6}, \hat{7}$ . To build triads, the third and the fifth notes above the root are considered together with the root (or fundamental).  $\hat{1}, \hat{3}, \hat{5}$  is the first degree chord, built on top of the scale's first degree and central for tonal music. The chords of the fifth degree  $\hat{5}, \hat{7}, \hat{2}$  ( $\hat{7}$  sharp when in a minor scale) and of the forth degree  $\hat{4}, \hat{6}, \hat{1}$  are also important. The triads built on the other degrees are less important then these and are usually understood in relation to them. The 'traditional harmony' comprises conventions and stylistic techniques to create progressions with such chords [62].

The 'functional harmony' ascribes functions to the three main chords and describes their use by means of these functions. The chord built on top of the first degree is the **tonic** chord ( $T$  or  $t$  for a major or minor tonic, respectively) and its function (role) consists on maintaining a center, usually referred to as a "ground" for the music. The chord built on the fifth degree is the **dominant** ( $D$ , the dominant is always major) and its function is to lean for the tonic (the dominant chord asks for a conclusion and this conclusion is the tonic). Thus, the dominant chord guides the music to the tonic. The triad built on the fourth degree is the **subdominant** ( $S$  or  $s$  for a major or minor subdominant, respectively) and its function is to deviate the music from the tonic. The tonal discourse aims at

confirming the tonic using the tonic-dominant-tonic progression which is expanded by using other chords in various ways.

The remaining triads are associated to these three most important chords. In the major scale, the associated relative (relative tonic  $T_r$ , relative subdominant  $S_r$  and relative dominant  $D_r$ ) is the triad built a third below, and the associated counter-relative (counter-relative tonic  $T_c$ , counter-relative subdominant  $S_c$  and the counter-relative dominant  $D_c$ ) is the triad built in a third above. In the minor scale the same happens, but the triad a third below is called counter-relative ( $tC$ ,  $sC$ ) and the triad a third above is called relative ( $tR$ ,  $sR$ ). The precise functions and musical effects of these chords are controversial but are basically the same as the chords they are associated to. Table 2 shows relations between the triads built at each degree of the major scale.

Table 2. Summary of tonal harmonic functions on the major scale. Tonic is the musical center, the dominant leans to the tonic and the subdominant moves the music away from the tonic. The three chords can, in principle, be replaced by their respective relative or counter-relative.

relative	main chord of the function	counter-relative
$\hat{6}, \hat{1}, \hat{3}$	tonic: $\hat{1}, \hat{3}, \hat{5}$	$\hat{3}, \hat{5}, \hat{7}$
$\hat{3}, \hat{5}, \hat{7}$	dominant: $\hat{5}, \hat{7}, \hat{2}$	[ $\hat{7}, \hat{2}, \hat{4}\#$ ]
$\hat{2}, \hat{4}, \hat{6}$	subdominant: $\hat{4}, \hat{6}, \hat{1}$	$\hat{6}, \hat{1}, \hat{3}$

The dominant counter-relative should form a minor chord. It explains the change in the forth degree by a semitone above  $\hat{4}\#$ . The diminished chord  $\hat{7}, \hat{2}, \hat{4}$ , is generally considered a ‘dominant seventh chord with the root omitted’ [38]. In the minor mode, there is a change in  $\hat{7}$  by an ascending semitone to achieve a separation between  $\hat{7}$  and  $\hat{1}$  of a semitone. This is important for the dominant function (which should lean to the tonic). In this way, the dominant is always major, for both major and minor scales and, therefore, in a minor scale the relative dominant remains a third below, and the counter-relative remains a third above.

**4.2.3 Tonal expansion: individual functions and chromatic mediants.** Each chord can be stressed and developed by performing their individual dominant or subdominant, which are the triads based on a fifth above or a fifth below, respectively. These individual dominants and subdominants, in the same way, have also subdominants and dominants of their own. Given a tonality, any chord can occur, no matter how distant it is from the most basic chords and from the notes of the scale. The unique (theoretical) condition is that the occurrence presents a coherent trajectory of dominants and subdominants (or their relatives and counter-relatives) to the original tonality.

There are four mediants for each chord, they are a third apart from the original chord and are simple triads, as are the relatives and counter-relatives, but retain the major/minor quality of the reference chord. The ‘chromatic mediants’ are the upper mediant, formed with the root at the third of the original chord; and the lower mediant, formed by the fifth at the third of the original chord. If two chromatic alterations exist, i.e. two notes are altered by one semitone, it is a ‘doubly-chromatic mediant’. Again, there are two forms: the upper form, with a third in the fifth of the original triad; and the lower form, with a third in the root of the original triad. This relation between chords is considered of advanced tonalism, sometimes even considered as an expansion and dissolution of tonalism, with strong and impressive effects although they are simple, consonant major/minor triads. Chromatic mediants are used since the end of Romanticism by Wagner, Lizt, Richard Strauss, among others [60, 62].



**4.2.4 Modulation.** Modulation is the change of key (tonic, or tonal center) in music, being characterized by start and end keys, and transition artifacts. Keys are always (thought of as) related by fifths and their relatives and counter-relatives. Some ways to perform modulation include:

- Transposing the discourse to a new key, without any preparation. It is a common Baroque procedure and also incident in other periods. Sometimes it is called phrasal modulation or unprepared modulation.
- Careful use of an individual dominant, and perhaps also the individual subdominant, to confirm change in key and harmonic field.
- Use of chromatic alterations to reach a chord in the new key by starting from a chord in the previous key. Called chromatic modulation.
- Featuring a unique note, possibly repeated or suspended with no accompaniment, common to start and end keys, it constitutes a peculiar way to introduce the new harmonic field.
- Changing the function, without changing the notes, of a chord. This procedure is called enharmony.
- Maintaining the tonal center and changing the key quality from major to minor (or vice-versa) is a ‘parallel modulation’. Keys with same tonic but different (major/minor) qualities are known as homonyms.

The dominant has great importance and is a natural pivot in modulations, a fact that leads to the circle of fifths [2, 38, 60, 62]. Other inventive ways to modulate are possible, to point but one common example, the minor thirds tetrad ( $E_i^t m$  in Equations 81) can be sustained to bridge between tonalities, with the facility that both its tritones can be resolved in a number of ways. The music piece *Acorde cedo* explores these chord relations [27].

### 4.3 Counterpoint

Counterpoint is a set of techniques for the conduction of simultaneous melodic lines, or “voices”. The bibliography covers systematic ways to conduct voices, leading to scholastic genres like canons, inventions and fugues [30, 63]. It is possible to summarize the rules of scholastic counterpoint, and it is known that Beethoven – among others – also outlined such a digest of counterpoint.

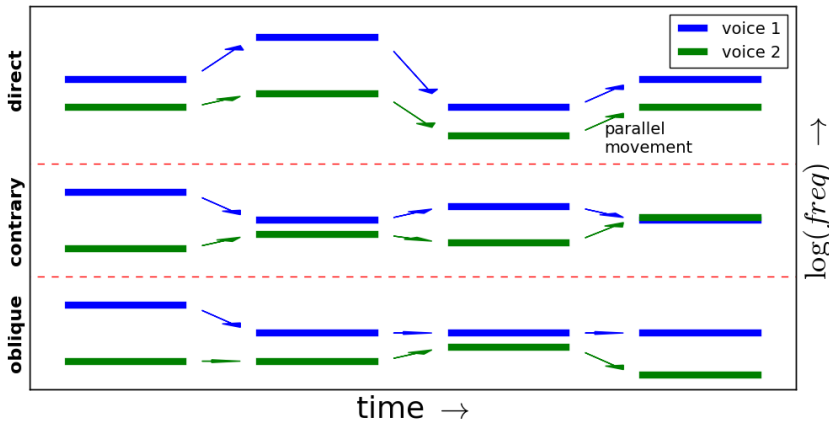


Fig. 23. Different motions of counterpoint aiming to preserve independence between voices. There are 3 types of motion: direct, contrary and oblique. The parallel motion is a type of direct motion. A voice is often called a ‘melodic line’ or a melody.



The main purpose of (scholastic) counterpoint is to conduct voices in a way that they sound independent. In order to do that, the relative motion of voices (in pairs) is crucial and categorized as: direct, oblique and contrary, as depicted in Figure 23. The parallel motion is a direct motion in which the starting and final intervals are the same. The golden rule here is to be careful with the direct motions, avoiding them when ending in a perfect consonance. The parallel motion should occur only between imperfect consonances and no more than three consecutive times. Dissonances can be forbidden or used only when followed and preceded by consonances of neighbor degrees, i.e. adjacent notes in a scale. The motions that lead to a neighbor note in the scale sound coherent and are prioritized. When having 3 or more voices, the melodic relevance lies mainly in the highest and then in the lowest of the voices [30, 63, 71].

These rules were used in the musical piece *Count point* [27].

#### 4.4 Rhythm

Table 3. Durations heard as rhythm, as pitch and the transition.

perception of durations as rhythm											
duration (s)	32,	16,	8,	4,	2,	1,	1/2,	1/4,	1/8,	...	
frequency (Hz)	1/32,	1/16,	1/8,	1/4,	1/2,	1,	2,	4,	8,	... transition	
-											
transition											
duration (s)	rhythm ... $\left\  \frac{1}{16} = 62.5ms, \frac{1}{20} = 50ms \right\ $ ... pitch										
frequency (Hz)	16, 20										
transition											
perception of durations as pitch											
duration (s)	transition ... $\left\  \frac{1}{40}, \frac{1}{80}, \frac{1}{160}, \frac{1}{320}, \frac{1}{640} \right\ $ ...										
frequency (Hz)	40 80 160 320 640										

Rhythmic notion is dependent on events separated by durations [43]. Such events can be heard individually if their onsets are spaced by at least 50 – 63ms. For the temporal separation between them to be perceived as a duration, the period should be even a bit larger, around 100ms [56]. It is possible to summarize the durations heard as rhythm or pitch as in Table 3 [16, 56].

The transition span in Table 3 is minimized because the limits are not well defined. In fact, the duration where someone begins to perceive a fundamental frequency, or a separation between occurrences, depends on the listener and sonic characteristics [56, 57]. The rhythmic metric is commonly based on a key duration called pulse, which is typically between 0.25 and 1.5s (240 and 40BPM, respectively<sup>20</sup>). In music education and cognitive studies, it is common to associate this range of frequencies with the durations of the heart beat, movements of respiration and steps of a walking or running person [43, 57].

The pulse is subdivided into equal parts and is also repeated in sequence. These relations (division and concatenation) usually follow relations of small integers. By far, the most often musical pulse divisions (and their sequential groupings), in written and ethnic music, are: 2, 4 and 8; 3, 6 (2 groups of 3 or 3 groups of 2), 9 and 12 (3 and 4 groups of 3); and then 5 and 7, completing 1-9 and 12.

<sup>20</sup>BPM stands for Beats Per Minute and is just a frequency measure like Herz, but is the number of incidences per minute instead of second. BPM is often used as a measure of musical tempo and of heart rate.

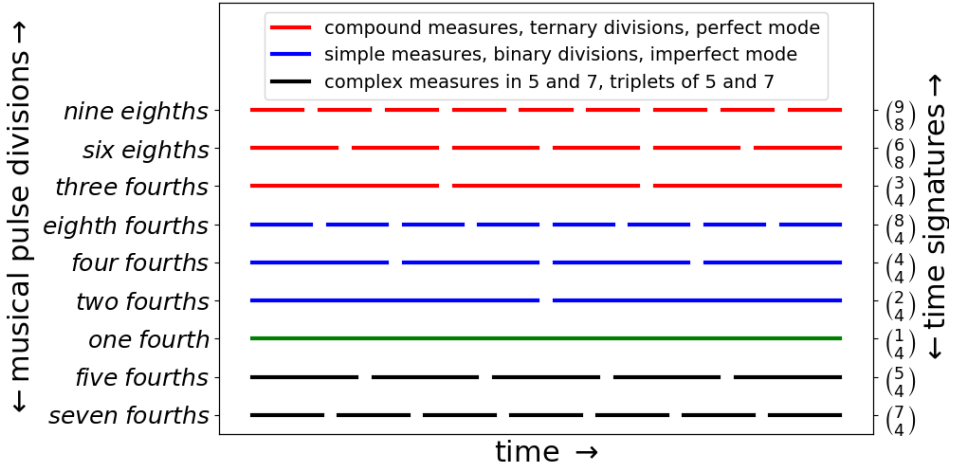


Fig. 24. Divisions and groupings of the musical pulse for establishing a metric. Divisions of the quarter note, regarded as the pulse, is presented on the left. The time signature yield by the corresponding grouping of the music pulse is presented on the right.

Other metrics are less common, like division or grouping in 13, 17, etc, and are mainly used in experimental or concert music of the XX and XXI centuries. No matter how complex they seem, metrics are almost always compositions and decompositions of 1-9 equal parts [32, 57]. This is illustrated in Figure 24.

Binary divisions are frequent in dance rhythms and celebrations, and are called “imperfect”. Ternary relations are typical of ritualistic and sacred music and are called “perfect”. Strong units (accents) fall in the ‘head’ of the units (the first subdivision) and are called downbeats. In binary divisions (2, 4 and 8), strong units alternate with weak units (e.g. division in 4 is: strong, weak, average strong, weak). In ternary divisions (3, 6 and 9) two weak units succeed the downbeat (e.g. division in 3 is: strong, weak, weak). Division in 6 is considered compound but can also occur as a binary division. Binary division units which suffer a ternary division yields two units divided into three units each: strong (subdivided in strong, weak, weak) and weak (also subdivided in strong, weak, weak). Another way to perform the division in 6 is with a ternary division whose units subdivide as binary, resulting in: a strong unit (subdivided in strong and weak) and two weak units (subdivided in strong and weak each).

An accent in the weak beat is a ‘backbeat’, whereas a note starting on a weak beat and persisting across a strong beat is a ‘syncope’. These are often found in ethnic and popular music and was used with parsimony in classical music before the XX century.

Notes can occur inside and outside of these divisions of the ‘musical metric’. In most well-behaved cases, notes occur exactly on these divisions, with greater incidence on strong beats. In extreme cases, rhythmic metric cannot be perceived [57]. Noteworthy is that (usually small or progressive) variations along the temporal grid are crucial for musical interpretation styles [15].

Let the pulse be the grouping level  $j = 0$ , the first pulse subdivision be level  $j = -1$ , the first pulse agglomeration be level  $j = 1$  and so on. Accordingly, let  $P_i^j$  be the  $i$ -th unit at grouping level  $j$ :  $P_{10}^0$  is the tenth pulse,  $P_3^1$  is the third grouped unit (possibly the third measure),  $P_2^{-1}$  is the second part of pulse subdivision. The limits of  $j$  are of special interest: pulse divisions are durations perceivable as rhythm; furthermore, the pulses sum, at its maximum, a music or a cohesive set of musical pieces.

In other words, a duration given by  $P_i^{min(j)}$ ,  $\forall i$ , should be greater than 50ms and the durations summed together  $\sum_{\forall i} P_i^{max(j)}$  should be less than a few minutes or, at most, a few hours. These limits might be extrapolated in extreme cases and with aesthetic goals.

Each level  $j$  has some parts  $i$ . When  $i$  has three different values (or multiple of three) there is a perfect (i.e. ternary or compound) relation. When  $i$  has only two, four or eight possible values, than there is an imperfect relation (i.e. binary or simple), as shown in Figure 24. Any unit can be specified as:

$$P_{\{i_k\}}^{\{j_k\}} \quad (87)$$

where  $j_k$  is the grouping level and  $i_k$  is the unit itself.

As an example, consider  $P_{3,2,2}^{-1,0,1}$  as the third subdivision  $P_3^{-1}$  of the second pulse  $P_2^0$  and of the second pulse group  $P_2^1$  (possibly second measure). Each unit  $P_i^j$  can be associated with a sequence of temporal samples  $T$  that constitutes e.g. a note. In practice, there is an underlying reference duration, usually associated with the pulse, e.g.  $d_r = 1$  second, and the durations of each segment are specified by:

- a ‘temporal notation’: where each entry is a relative duration to be multiplied by the reference duration. E.g.  $durs = \{1, 0.5, 4\}$  is mapped to  $\{d_i d_r\} = \{1d_r, 0.5d_r, 4d_r\}$ . Or:
- a ‘frequential notation’: where each entry is how many the entry that fits a same duration. E.g.  $durs = \{4, 2, 16\}$  is mapped to  $\{d_r/d_i\} = \{d_r/4, d_r/2, d_r/16\}$ . This notation might be less intuitive but it is more tightly related to traditional music theory, where e.g. the duration related to the number 4 is twice the duration related to the number 8.

See the function `rhythmToDurations` in file `src/aux/functions.py` for an implementation of both notations that allows the specification of tuplets (the use of arbitrary divisions of a reference duration). The music piece *Poli Hit Mia* uses different metrics. [27]

#### 4.5 Repetition and variation: motifs and larger units

Given both frequential (chords and scales) and rhythmic (simple, compound and complex beat divisions and agglomerations) musical strata, it is natural to present them in a coherent and meaningful manner [5]. The concept of an arc is essential in this context: by departing from a context and returning, an arc is made. One important, and maybe trivial, case is the arc from and to the absence of a unit: from the beginning (unit did not exist before) to the end (unit will not exist from thereon). The audition of melodic and harmonic lines is permeated by arcs due to the cognitive nature of the musical hearing: as the mind divides an excerpt, and groups excerpts, each of the units yields an arc. Accordingly, the note can be considered the smallest (relevant) arc, and each motif and melody as an arc as well. Each beat and subdivision, each measure and musical section, constitutes an arc. Music in which the arcs do not present consistency with one another can be understood as music with no coherence. Coherence impression comes, mostly, from the skilled handling of arcs in a music piece.

Musical arcs are abstract structures and amenable to basic operations. A spectral arc, like a chord, can be inverted, magnified and permuted, to mention just a few possibilities. Temporal arcs, like a melody, a motif, a measure or a note, are also prone to variations. Let  $S = \{s_j = T^j = \{t_i^j\}_0^{\Lambda_j-1}\}_0^{H-1}$  be a sequence of  $H$  musical events  $s_j$ , each event with its  $\Lambda_j$  samples  $t_i^j$  (refer to the beginning of this Section 4 if needed). Bellow is a list of basic techniques for variation.

- Temporal translation is a displacement  $\delta$  of a specific material to another instant  $\Gamma' = \Gamma + \delta$  of the music. It is a variation that changes temporal localization in a music:  $\{s_j'\} = \{s_j^{\Gamma'}\} =$

- $\{s_j^{\Gamma+\delta}\}$  where  $\Gamma$  is the duration between the beginning of the piece (or another reference) and the first event  $s_0$  of the original structure  $S$ , and  $\delta$  is the time offset of the displacement.
- Temporal expansion or contraction is a change in duration of each arc by a factor  $\mu$  :  $s_j'^\Delta = s_j^{\mu j \cdot \Delta}$ . Possibly,  $\mu_j = \mu$  is constant.
  - Temporal reversion consists on generating a sequence with elements in the reverse order of the original sequence  $S$ , thus:  $S' = \{s'_j\}_0^{H-1} = \{s_{(H-j-1)}\}_0^{H-1}$ .
  - Pitch translation, or transposition, is a displacement  $\tau$  of the pitches. It is a variation that changes pitch localization:  $\{s'_j\} = \{s_j^{\Xi'}\} = \{s_j^{\Xi+\tau}\}$  where  $\Xi$  is a reference value, such as the pitch of a section  $S$  or of the first event  $s_0$ . If  $\tau$  is given in semitones, the transposition displaces a frequency  $f$  to  $\tau_f = f2^{\frac{\tau}{12}}$ , and the pitch  $\Xi_i$  to  $\Xi'_i = \Xi_i + 12 \log_2 \left( \frac{f'_i}{f_i} \right)$ .<sup>21</sup>
  - Interval inversion is either: 1) the inversion of note pitch order, within the octave equivalence, such as described in Section 4.1.2; or 2) the inversion of interval orientation. In the former case, the number of semitones is preserved in the “strict inversion”:  $f'_i = 2^{-e} f_i$  where  $e$  is a positive constant; the inversion is said tonal if the distances are considered in terms of the diatonic scale  $E_k$ :  $f'_i = f \cdot 2^{\left( \frac{12 - e(7-j_e)}{12} \right)}$  where  $j_e$  is the index in  $E$  (as in Equation 83).
  - Rotation of musical elements is the translation of all elements a number of positions ahead or behind, with the care to fill empty positions with events which are out of the slots. Thus, a rotation of  $\tilde{n}$  positions is  $s'_n = s_{(n+\tilde{n})\%H}$ . If  $\tilde{n} < 0$ , it is sufficient to use  $\tilde{n}' = H - \tilde{n}$ . It is usual to associate  $\tilde{n} > 0$  (events advance) with the clockwise rotation and  $\tilde{n} < 0$  (elements delay) with the anti-clockwise rotation. Additional information concerning rotations is given in Section 4.7.
  - The insertion and removal of material in  $S$  can be ornamental or structural:  $S' = \{s'_j\} = \{s_j \text{ if condition A, otherwise } r_j\}$ , for any  $r_j$ , including silence. Elements can be inserted at the beginning, like a prefix for  $S$ ; at the end, as a suffix; or in the middle, splitting  $S$  into both a prefix and a suffix. Both materials can be mixed in a variety of ways.
  - Changes in articulation, orchestration and spatialization, or  $s'_j = s_j^{*j}$ , where  $*_j$  is the new characteristic incorporated by element  $s'_j$ .
  - Accompaniment. Musical material presented when  $S$  occurs can be modified to yield a variation.

From these processes, many others are derived, such as the inverted retrograde, the temporal contraction with an external suffix, etc. Variations are often thought about in the terms above but are also often very loose, such as an arbitrary shuffle of the notes in a melody which the composer or performer finds interesting. As a result, a whole process of mental and neurological activity is unleashed for relating the arcs, responsible for feelings, memories and imaginations, typical of a diligent musical listening. This cortical activity is critical to musical therapy, known by its utility in cases of depression and neurological injury. Also, it is known that regions of the human brain responsible for sonic processing are also used for other activities, such as for performing verbal discourse and mathematics. [57, 59]

Paradigmatic structures guide the creation of new musical material. One of the most established structures is the tension/relaxation dipole. Other traditional dipoles include tonic/dominant,

<sup>21</sup>In the MIDI protocol,  $\Xi_f = 55\text{Hz}$  when pitch  $\Xi = 33$  (an  $A1$  note). Another good MIDI reference is  $\Xi_f = 440\text{Hz}$  and  $\Xi = 69$  ( $A4$ ). The difference  $(\Xi_1 - \Xi_2)$  is in semitones.  $\Xi$  is not a measure in semitones:  $\Xi = 1$  is not a semitone, it is a note with an audible frequency as rhythm, with less than 9 periods each second (see Table 3).

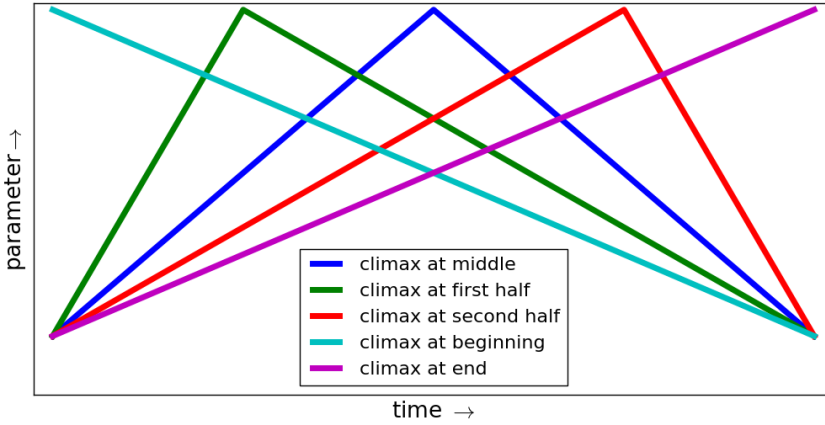


Fig. 25. Canonical distinctions of musical climax in a given melody and other arcs. The possibilities considered are: climax at the beginning, at the first half, in the middle, in the second half and at the end. The x and y-axis parameters can be non-existent and yield only a reference structure.

repetition/variation, consonance/dissonance, coherence/rupture, symmetry/asymmetry, equivalence/difference, arrival/departure, near/far, and stationary/moving. All these dipoles are often thought of as parallel or even as equivalent. Ternary constructions tend to relate to the circle and to unification. The ‘transcendental’ ternary communion, ‘modus perfectus’, opposes to the ‘passionate’ dichotomic, ‘modus imperfectus’ [4]. For a scholastic discussion on the composition of motives, phrases, melodies, themes and musical form (such as rondo, ternary, theme and variations), see [64].

#### 4.6 Directional structures

Musical arcs may be thought of as having two sections: the first reaches the apex and the second returns from apex to the initial region. This apex is called climax by traditional music theory. It is usual to distinguish between arcs whose climax is at the beginning, middle, end, or the first or second half of the duration. These structures are shown in Figure 25. The varying parameter can be non-existent, a case in which the arc consists only of a reference structure, a case which resembles a note without the fundamental frequency. [64]

Consider the sequence  $S = \{s_i\}_0^{H-1}$  with a monotonic variation of a characteristic. The sequence  $R = \{r_i\}_0^{2H-2} = \{s_{(H-1-|H-1-i|)}\}_0^{2H-2}$  presents perfect specular symmetry, i.e. the second half is the mirrored version of the first. In musical terms, the climax is in the middle of the sequence. It is possible to achieve different results by using sequences with different sizes. All the mathematics of sequences, already well established and taught routinely in calculus courses, can be used to generate these arcs [35, 64]. Theoretically, when applied to any characteristic of musical events, these sequences produce arcs, since they imply a deviation and return of an initial context (parametrization, state, etc). Henceforth, it is possible for a given sequence to have numerous distinct arcs, with different sizes and climaxes. This is an interesting and useful resort, and the correlation of arcs yields coherence [60].

In practice, and historically, there is special incidence and use of the golden ratio. The Fibonacci sequence might be generalized as follows in order for any two numbers to be used and approximate the golden ratio. Given any two numbers  $x_0$  and  $x_1$ , define the elements of the sequence  $\{x_n\}$  as:  $x_n = x_{n-1} + x_{n-2}$ . The greater  $n$  is, the more  $\frac{x_n}{x_{n-1}}$  approaches the golden ratio (1.61803398875...).

The sequence converges fast even with discrepant initial values. E.g. let  $x_0 = 1$ ,  $x_1 = 100$  and  $y_n = \frac{x_n}{x_{n+1}}$ , the error for the first values with respect to the golden ratio is, approximately,  $\{e_n\} = \{100 \frac{y_n}{1.61803398875} - 100\}_1^{10} = \{6080.33, -37.57, 23, -7.14, 2.937, -1.09, 0.42, -0.1601, 0.06125, -0.02338\}$ . The Fibonacci sequence presents the same error progression, but starts at the second step of a more discrepant initial setting ( $\frac{1}{1} \approx \frac{100+1=101}{100}$ ). One might benefit from the On-Line Encyclopedia of Integer Sequences (OEIS [66]) for exploring various sequences.

The musical piece *Dirracional* exposes the use of arcs in explicit directional structures. [27]

#### 4.7 Cyclic structures

The philosophical understanding that human thought is founded on the recognition of similarities and differences (e.g. as perceived in stimuli), places symmetries at the core of cognition [18]. Mathematically, it is commonplace to express symmetries as algebraic groups, and a finite group is always isomorphic to a permutation group (by Cayley's theorem). In a way, this states that permutations can express any symmetry in a finite system [8]. Also, any permutation set can be used as a generator of an algebraic group [7]. In music, permutations are ubiquitous in scholastic techniques, which confirms their central role. The successive application of permutations generates cyclic arcs [7, 19, 78] and e.g. these two academic documents report on the generation of musical structures using permutation groups: [20, 21]. The properties defining a group  $G$  are:

$$\begin{aligned}
 \forall p_1, p_2 \in G &\Rightarrow p_1 \bullet p_2 = p_3 \in G && \text{(closure property)} \\
 \forall p_1, p_2, p_3 \in G &\Rightarrow (p_1 \bullet p_2) \bullet p_3 = p_1 \bullet (p_2 \bullet p_3) && \text{(associativity property)} \\
 \exists e \in G : p \bullet e = e \bullet p, \quad \forall p \in G &&& \text{(existence of the identity element)} \\
 \forall p \in G, \exists p^{-1} : p \bullet p^{-1} = p^{-1} \bullet p = e &&& \text{(existence of the inverse element)}
 \end{aligned} \tag{88}$$

From the first property follows that two permutations act as one permutation. In fact, it is possible to apply a permutation  $p_1$  and another permutation  $p_2$ , and, comparing both initial and final orderings, observe another permutation  $p_3$ . Every element  $p$  operated with itself a sufficient number of times  $n$  reaches the identity element  $p^n = e$  (otherwise the group generated by  $p$  would be infinite). The order  $n$  of an element  $p$  is the lowest  $n : p^n = e$ . Thus, a finite permutation  $p$ , successively applied, reaches the initial ordering of its elements, and yields a cycle. This cycle, if used for parameters of notes or other musical structures, yields a cyclic arc.

These arcs can be established by using one or a set of permutations. As a historical example, the *change ringing* tradition conceives music through bells played one after another and then played again, but in a different order. This process is repeated until it reaches the initial ordering. The sequence of different orderings is a *peal*. Table 4 presents a traditional *peal*, named "Plain Change" [19], for 3 bells (1, 2 and 3), which explores all possible orderings. Each line indicates one bell ordering to be played. Permutations occur between each line. In this case, the musical structure consists of permutations that entail a cyclic behavior.

The use of permutations in music can be summarized in the following way: let  $S = \{s_i\}$  be a sequence of musical events  $s_i$  (e.g. notes), and  $p$  a permutation.  $S' = p(s_i)$  comprises the same elements of  $S$  but in a different order. Permutations have two notations: cyclic and natural. The natural notation basically indicates the original indexes in the order that results from the

Table 4. Change Ringing: a traditional *peal* for 3 bells. Permutations occur between each line. Each line is a bell ordering and each ordering is played at a time.

1	2	3
2	1	3
2	3	1
3	2	1
3	1	2
1	3	2
1	2	3

permutation. Thus, given the original ordering of the sequence by its indexes [0 1 2 3 4 5 ...], the permutation is noted by the sequence of indexes it produces (e.g. [1 3 7 0 ...]). In the cyclic notation, a permutation is expressed by swaps of elements and its successors. E.g. (1, 2, 5)(3, 4) in cyclic notation is equivalent to [0, 2, 5, 4, 3, 1] in natural notation.

In the auralization of a permutation, it is not necessary to permute elements of  $S$ , but only some characteristic. Thus, if  $p$  is a permutation and  $S$  is a sequence of basic notes as in the end of Section 2.6, the sequence  $S' = p^f(S) = \{s_i^{p(f)}\}$  consists of the same musical notes, following the same order and maintaining the same characteristics, but with the fundamental frequencies permuted according to  $p$ .

Two subtleties of this procedure should be commented upon. First, a permutation  $p$  is not restricted to involve all elements of  $S$ , i.e. it can operate in a subset of  $S$ . Second, not all elements  $s_i$  need to be executed at each access to  $S$ . To exemplify, let  $S$  be a sequence of music notes  $s_i$ . If  $i$  goes from 0 to  $n$ , and  $n > 4$ , at each sequence of 4 notes it is possible to execute e.g. only the first 4 notes. The other notes of  $S$  can occur in other events where permutations allocate such notes to the first four events. The execution of disjoint sets of  $S$  is the same as modifying the permutation and executing the first  $n$  notes.

In summary, to each permutation  $p$ , we have to determine: 1) note characteristics where it operates (frequency, duration, *fades*, intensity, timbre, etc); and 2) the period of incidence (how many times  $S$  is used before a permutation is applied).

The PPEPPS/FIGGS and the 3 *Trios* present respectively a computational implementation and an instrumental musical piece that use permutations to achieve musical structures [20, 21, 24, 25, 27].

#### 4.8 Serialism and post-serial techniques

Recapitulating concepts from Sections 4.2.1 and 4.5, sequences of characteristics can be predefined and used throughout a musical piece. These sequences can be of intensities, timbre, durations, density of events, etc. Sequences can be used very strictly or loosely, such as by skipping some elements. The sequences can be of different sizes, yielding arcs until the initial condition is reached again (i.e. cycles as in Section 4.7). One paradigmatic case is the “total serialism” where all the musical characteristics are serialized [46]. Although the use of sequences is inherent to music (e.g. scales, metric pulses), their use with greater emphasis than tonal (or modal) elements in western music, as an artistic trend, took place only in the first half of the twentieth century and is called “serialism”. Post-serial techniques are numerous, and brief descriptions of important concepts to exemplify them are:

- Spectralism: consists on the use the (Fourier) spectrum of a sound or the harmonic series for musical composition, such as to obtain harmonies, sequences of pitches or a temporal evolution of the overall spectrum. For example, the most prominent frequencies can be



used as pitches, real notes can be used to mimic an original spectrum (e.g. use piano notes to mimic the spectrum of a spoken sentence) or portions of the spectrum made to vary. [34]

- Spectromorphology: can be considered a spectral music (spectralism) theoretical framework [61, 67] that examines the relation between sound spectra and their temporal evolution. The theory poses e.g. different onsets, continuations and terminations; characteristics of (sonic) “motion”; and spectral density.
- Stochastic music: the use of random variables to describe musical elements are extensively considered in stochastic music [77]. In summary, one can use probability distributions for the synthesis of basic sounds and for obtaining larger scale musical structures. Changes in these distributions or in other characteristics yield the discourse.
- Textures: sounds may be assembled in terms of a “sonic texture”. Sonic textures are often thought about very abstractly as a sonic counterpart of visual textures. Examples of parameters that are used: range between highest and lowest note, density of notes, durations of notes, motives, number of voices. If the sounds are small enough (typically < 100ms) the process can be thought of in terms of *granular synthesis* [56].

#### 4.9 Musical idiom?

In numerous studies and aesthetic endeavors, there are models, discussions and exploitation of a ‘musical language’. Some of them are linguistic theories applied to music and some discern different ‘musical idioms’ [16, 44, 60, 62]. Simply put, a musical idiom or language is the result of chosen materials together with variation techniques and relations established between elements along a music piece. In these matters, dichotomies are prominent, as explained in Section 4.5: repetition and variation, relaxation and tension, stability and instability, consonance and dissonance, etc. A thorough discussion of what can be considered a musical language is out of the scope of this article, but this brief consideration of the subject is useful as a convergence of all the previous content.

#### 4.10 Musical usages

The basic note was defined and characterized in quantitative terms in Section 2. Next, the internal note composition was addressed within both internal transitions and elementary sonic treatment (Section 3). Finally, this section aims at organizing these notes in music. The numerous resources and consequent infinitude of praxis possibilities is typical and highly relevant for artistic contexts [62, 74].

There are studies and further developments for each of the presented resources. For example, it is possible to obtain ‘dirty’ triadic harmonies (with notes out of the triad) by superposition of perfect fourths. Another interesting example is the superimposition of rhythms in different metrics, constituting what is called *polyrhythm*. The music piece *Poli-hit my* [27] explores these simultaneous metrics by impulse trains convolved with notes.

Microtonal scales are important for 20th century music [75] and yielded diverse remarkable results throughout history, e.g. fourths of a tone ( $\epsilon = 2^{\frac{1}{24}}$ ) are often used in some genres of Indian and Arabic music. The musical sequence *Micro Tone* [27] explores these possibilities, including microtonal melodies and harmonies with many pitches in a very reduced frequency bandwidth.

As in Section 3.6, relations between parameters are powerful to achieve musical pieces. E.g. the number of permuted notes can vary during the music, a relationship between permutations and the piece duration. Harmonies can be obtained from triads (Equations 86) with duplicated notes at each octave and more numerous duplication when the depth and frequency of vibratos are lower (Equations 56, 57, 58, 59, 60). Incontestably, the possibilities are very wide, which is made evident by the numerous musical pieces and styles.



The symmetries at octave divisions (Equations 81) and the symmetries presented as permutations (Table 4 and Equations 88) can be used together. In the music piece *3 trios*, this association is performed in a systematic way in order to achieve a specific style. This is an instrumental piece, not included as a source code but available online [25].

*PPEPPS* (Pure Python EP: Project Solvent) is an EP (Extended Play) synthesized using resources presented in this document. With minimal parametrization, the scripts generate complete musical pieces, allowing easy composition of sets of music. [24] A simple script of a few lines specifies music delivered as 16 bit 44.1kHz PCM files (WAVE). This facility and technological arrangement creates aesthetic possibilities for both sharing and education.

## 5 CONCLUSIONS AND FURTHER DEVELOPMENTS

In our understanding, this article is effective in relating musical elements to digital audio. We aimed at achieving a concise presentation of the subject because it involves many knowledge fields, and therefore can very easily blast into thousands of pages. Some readers might benefit from the text alone, but the *scripts* in the MASS toolbox, where all the equations and concepts are directly and simply implemented as software (in Python), are very helpful for one to achieve elaborated implementations and deeper understandings. The scripts include routines that render musical pieces to illustrate the concepts in practical contexts. This is valuable since art (music) can involve many non-trivial processes and is often deeply glamorized, which results in a nearly unmanageable terrain for a newcomer. Moreover, this didactic report and the supplied open source scripts should facilitate the use of the framework. One of the Supporting Information documents [28] holds listings of sections, equations, figures, tables, scripts and other documents. Another Supporting Information document [29] holds a PDF presentation of the code related to each section because many readers might not find it easy to browse source code files.

The possibilities provided by this exposition pour from both the organization of knowledge and the ability to achieve sounds which are extremely true to the models. For example, one can produce noises with an arbitrary resolution of the spectrum and a musical note can be synthesized with the parameters (e.g. of a vibrato) updated sample-by-sample. Furthermore, software for synthesis and processing of sounds for musical purposes by standard restricts the bit depth to 16 or 24. This is achievable in this framework but by standard Python uses more bits per floating point number. These “higher fidelity” characteristics can be crucial e.g. for psychoacoustic experiments or to generate high quality musical sounds or pieces. Simply put, it is compelling for many scientific and artistic purposes. The didactic potential of the framework is evident when noticed that:

- the integrals and derivatives, ubiquitous in continuous signal processing, are all replaced, in discrete signals, by summations, which are more intuitive and does not require fluency in calculus.
- The equations and concepts are implemented in a simple and straightforward manner as software which can be easily assembled and inspected.

In fact, this framework was used in a number of contexts, including courses, software implementations and for making music [23, 40, 73]. As far as the authors know, such detailed analytical descriptions have not been covered before in the literature, such as testified in the literature review (Appendix G of [23], where books, articles and open software are related to this framework).

The free software license, and online availability of the content, facilitate collaborations and the generation of sub-products in a co-authorship fashion, new implementations and development of musical pieces. The scripts can be divided in three groups: implementation of all the equations and topics of music theory covered here; routines for rendering musical pieces that illustrate the concepts; scripts that render the figures of this article and the article itself.

This framework favored the formation of interest groups in topics such as musical creativity and computer music. In particular, the project [labMacambira.sourceforge.net](http://labMacambira.sourceforge.net) groups Brazilian and foreign co-workers in diverse areas that range from digital direct democracy and georeferencing to art and education. This was only possible because of the usefulness of audiovisual abilities in many contexts, in particular because of the knowledge and mastery condensed in the MASS framework.<sup>22</sup>

Future work might include application of these results in artificial intelligence for the generation of attractive artistic materials. Some psychoacoustic effects were detected, which need validation and should be reported, specially with [22].<sup>23</sup> Other foreseen advances are: enhancement of the Python package written using MASS [26], a JavaScript version of the toolbox, better hypermedia deliverables of this framework, user guides for different goals (e.g. musical composition, psychophysics experiments, sound synthesis, education), creation of more musical pieces, open experiments to be studied with EEG recordings, a linked data representation of the knowledge in MASS through SKOS and OWL to tackle the issues exposed in Section 1.2, data sonification routines, and further analytical specification of musical elements in the discrete-time representation of sound as feedback is received from the community.

## ACKNOWLEDGMENTS

This work was supported by Capes, CNPq and FAPESP (project 17/05838-3).

## REFERENCES

- [1] 2017. Public GMane archive of the metareciclagem email list. (2017). <http://arquivos.metareciclagem.org/>
- [2] E. Aldwell, C. Schachter, and A. Cadwallader. 2010. *Harmony & voice leading*. Wadsworth Publishing Company.
- [3] V.R. Algazi, R.O. Duda, D.M. Thompson, and C. Avendano. 2001. The cipc hrtf database. In *Applications of Signal Processing to Audio and Acoustics, 2001 IEEE Workshop on the*. IEEE, 99–102.
- [4] Willi Apel. 1961. *The notation of polyphonic music, 900-1600*. Number 38. Medieval Academy of Amer.
- [5] Pierre Boulez. 1972. *A música hoje*.
- [6] R. Bristow-Johnson. 1996. Wavetable synthesis 101, a fundamental perspective. In *Proc. AES Convention*, Vol. 101.
- [7] F.J. Budden and F.J. Budden. 1972. *The fascination of groups*. Cambridge University Press London.
- [8] Francis James Budden. 1972. *The fascination of groups*. Cambridge Univ. Press.
- [9] B. carty and V. Iazzarini. 2009. binaural hrtf based spatialisation: new approaches and implementation. In *dafx 09 proceedings of the 12th international conference on digital audio effects, politecnico di milano, como campus, sept. 1-4, como, italy*. Dept. of Electronic Engineering, Queen Mary Univ. of London,, 1–6.
- [10] S. Chacon, J.C. Hamano, and S. Pearce. 2009. *Pro Git*. Vol. 288. Apress.
- [11] C.I. Cheng and G.H. Wakefield. 2012. Introduction to head-related transfer functions (HRTFs): Representations of HRTFs in time, frequency, and space. *Watermark* 1 (2012).
- [12] John M. Chowning. 2000. Digital sound synthesis, acoustics and perception: A rich intersection. *Proceedings of the COST G-6 Conference on Digital Audio Effects (DAFX-00)* (Dec 2000).
- [13] B.D. Class, FFT Java, A. Wah, L.F. Crusher, P. Waveshaper, S. Enhancer, and S.F.R.V.T. Matrix. 2010. musicdsp.org source code archive. (2010).
- [14] J. Clough, N. Engebretsen, and J. Kochavi. 1999. Scales, sets, and interval cycles: A taxonomy. *Music Theory Spectrum* (1999), 74–104.
- [15] Perry R. Cook. 2002. *Real sound synthesis for interactive applications*. A K Peters, Natick, Massachusetts.
- [16] Gustavo Oliveira Alfaix de Assis. *Em busca do som*. Editora UNESP.
- [17] S. Dehaene. 2003. The neural basis of the Weber–Fechner law: A logarithmic mental number line. *Trends in cognitive sciences* 7, 4 (2003), 145–147.

<sup>22</sup>There are more than 700 videos, written documents, original software applications and contributions in well-known software (such as Firefox, Scilab, LibreOffice, GEM/Puredata, to name just a few) [40–42]. Some of these efforts are available online [23]. It is evident that all these contributions are a consequence of more than just MASS, but it is also evident to the authors that MASS had a primary role in converging interests and attracting collaborators.

<sup>23</sup>The sonic portraits were sent to a public mailing list [1] and the fifth piece was reported by some individuals to induce a state in which noises from the own tongue, teeth and jaw of the individual echoed for some seconds (the GMane archives with the descriptions of the effect by listeners in the public email list is unfortunately offline at the moment).

- [18] G. Deleuze. 1968. *Difference and Repetition*. Continuum.
- [19] R. Duckworth and F. Stedman. 2007. *Tintinnalogia, or, the Art of Ringing*. Echo Library.
- [20] Renato Fabbri e Adolfo Maia Jr. 2007. Applications of Group Theory on Granular Synthesis. (2007).
- [21] Renato Fabbri e Adolfo Maia Jr. 2008. Applications of Group Theory on Sequencing and Spatialization of Granular Sounds. (2008).
- [22] Renato Fabbri. 2012. Sonic pictures. (2012). <https://soundcloud.com/le-poste-tche/sets/sonic-pictures>
- [23] Renato Fabbri. 2013. Music in digital audio: psychophysical description and software toolbox. (2013). <http://www.teses.usp.br/teses/disponiveis/76/76132/tde-19042013-095445/en.php>
- [24] Renato Fabbri. 2013. PPEPPS (Pure Python EP - Project Solvent), and FIGGUS (Finite Groups in Granular and Unit Synthesis). (2013). <https://github.com/ttm/figgus/>
- [25] Renato Fabbri. 2017. 3 Trios para oboé, flauta e fagote. (2017). <https://soundcloud.com/le-poste-tche/sets/3-trios>
- [26] R. Fabbri. 2017. Music: a Python package for rendering music (based on MASS). (2017). <https://github.com/ttm/music/>
- [27] Renato Fabbri. 2017. Public Git repository for the MASS framework. (2017). <https://github.com/ttm/maass/>
- [28] R. Fabbri et al. 2017. Equations, scripts, figures, tables and documents in the MASS framework. (2017). <https://github.com/ttm/maass/raw/master/doc/listings.pdf>
- [29] R. Fabbri et al. 2017. PDF presentation of the Python implementations in the MASS framework. (2017). <https://github.com/ttm/maass/raw/master/doc/code.pdf>
- [30] J.J. Fux and A. Mann. 1965. *The study of counterpoint from Johann Joseph Fux's Gradus ad Parnassum*. Vol. 277. WW Norton & Company.
- [31] Gnter Geiger. 2006. Table lookup oscillators using generic integrated wavetables. *Proc. of the 9th Int. Conference on Digital Audio Effects (DAFx-06)* (September 2006).
- [32] J.E. Gramani. 1996. *Rítmica viva: a consciência musical do ritmo*. UNICAMP.
- [33] Morag Josephine Grant. 2005. *Serial music, serial aesthetics: compositional theory in post-war Europe*. Vol. 16. Cambridge University Press.
- [34] Gérard Grisey and Joshua Fineberg. 2000. Did you say spectral? *Contemporary music review* 19, 3 (2000), 1–3.
- [35] H.L. Guidorizzi. 2001. *Um curso de cálculo*. Livros Técnicos e Científicos Editora.
- [36] P. Guillaume. 2006. *Music and acoustics: from instrument to computer*. Iste.
- [37] David Heeger. 2012. Perception Lecture Notes: Auditory Pathways and Sound Localization. (2012).
- [38] H.J. Koellreutter. 1986. *Harmonia funcional*. (1986).
- [39] S.M. Kostka, J.P. Clendinning, R. Ottman, and J. Phillips. 1995. *Tonal Harmony: With an Introduction to Twentieth*. McGraw-Hill.
- [40] #labmacambira @ Freenode. 2011. Canal Vimeo do Lab Macambira (mais de 700 videos). (2011). <https://vimeo.com/channels/labmacambira>
- [41] #labmacambira @ Freenode. 2013. Página principal do Lab Macambira. (2013). <http://labmacambira.sourceforge.net>
- [42] #labmacambira @ Freenode. 2017. Wiki do Lab Macambira. (2017). [http://wiki.nosdigitais.teia.org.br/Lab\\_Macambira](http://wiki.nosdigitais.teia.org.br/Lab_Macambira)
- [43] Osvaldo Lacerda. 1966. *Complndio de Teoria Elementar da Msica* (9.a ediflo ed.). Ricordi Brasileira.
- [44] Fred Lerdahl and Ray Jackendoff. 1983. *A Generative Theory of Tonal Music*. MIT Press.
- [45] L. Lessig. 2002. Free culture. *Retrieved February 5* (2002), 2006.
- [46] William LOVELOCK. 1972. *A concise history of music. Reprinted with revised record list*.
- [47] Guerino Mazzola. 2012. *The topós of music: geometric logic of concepts, theory, and performance*. Birkhäuser.
- [48] Florivaldo Menezes. 2004. *A Acústica Musical em Palavras e Sons*. Ateliê Editorial.
- [49] Jan Newmarch. 2017. Sound Codecs and File Formats. In *Linux Sound Programming*. Springer, 11–14.
- [50] T.E. Oliphant. 2006. *A Guide to NumPy*. Vol. 1. Trelgol Publishing USA.
- [51] A.V. Oppenheim and Shafer Ronald. 2009. *Discrete-time signal processing* (3 ed.). Pearson.
- [52] A.T. Porres and A.S. Pires. 2009. Um External de Aspereza para Puredata & MAX/MSP. In *Proceedings of the 12th Brazilian Symposium on Computer Music*.
- [53] TA Porres, J. Manzolli, and F. Furlanete. 2006. Análise de Dissonância Sensorial de Espectros Sonoros. In *Congresso da ANPPOM*, Vol. 16.
- [54] E.S. Raymond. 2004. *The art of Unix programming*. Addison-Wesley Professional.
- [55] C. Roads. 1996. *The computer music tutorial*. MIT press.
- [56] C. Roads. 2004. *Microsound*. MIT press.
- [57] Juan G. Roederer. 2008. *The Physics and Psychophysics of Music: An Introduction* (fourth edition ed.). Springer.
- [58] G. Van Rossum and F. L. Drake Jr. 1995. *Python tutorial*. Odense Universitet, Institut for Matematik og Datalogi.
- [59] O. Sacks. 2008. *Musophilia: Tales of Music and the Brain, Revised and Expanded Edition*. Knopf Doubleday Publishing Group.
- [60] F. Salzer. 1962. *Structural hearing: Tonal coherence in music*. Vol. 1. Dover publications.
- [61] Pierre Schaeffer. 2017. *Treatise on Musical Objects: Essays Across Disciplines*. Vol. 20. Univ of California Press.

- [62] A. Schoenberg and M. Maluf. 1999. *Harmonia*. Ed. UNESP.
- [63] A. Schoenberg and L. Stein. 1963. *Preliminary exercises in counterpoint*. Faber & Faber.
- [64] Arnold Schoenberg and Leonard Stein. 1967. *Fundamentals of musical composition*. London: Faber.
- [65] Bill Schottstaedt. 2017. An introduction to FM (Snd Manual). (2017). <https://ccrma.stanford.edu/software/snd/snd/fm.html>
- [66] Neil JA Sloane et al. 2010. The on-line encyclopedia of integer sequences. (2010). <https://oeis.org/>
- [67] Denis Smalley. 1997. Spectromorphology: explaining sound-shapes. *Organised sound* 2, 2 (1997), 107–126.
- [68] S.W. Smith. 2009. The Scientist and Engineer's Guide to Digital Signal Processing, 1999. (2009).
- [69] Julius O. Smith III. 2006. *Physical Audio Signal Processing (for virtual musical instruments and audio effects)*. <https://ccrma.stanford.edu/~jos/pasp/>
- [70] Julius O. Smith III. 2012. *Mathematics of the discrete fourier transform (dft) with audio applications* (second ed.). [https://ccrma.stanford.edu/~jos/log/FM\\_Spectra.html](https://ccrma.stanford.edu/~jos/log/FM_Spectra.html)
- [71] L. Tragtenberg. 2002. *Contraponto: uma arte de compor*. Edusp.
- [72] G. Van Rossum and F.L. Drake Jr. 1995. *Python reference manual*. Centrum voor Wiskunde en Informatica.
- [73] V. Vieira, G. Lunhany, G.M.C. Rocha Junior, C.M. Luporini, D. Penalva, R. Fabbri, and R. Fabbri. 2017. Vivace: a collaborative live coding language and platform. In *Proceedings of the 16th Brazilian Symposium on Computer Music*.
- [74] Anton Webern. 1963. *The Path To The New Music*. Theodore Presser Company.
- [75] S.R. Wilkinson, K. Laubach, R.S. Schiff, and J. Eiche. 1988. *Tuning in: Microtonality in Electronic Music: a Basic Guide to Alternate Scales, Temperaments, and Microtuning Using Synthesizers*. H. Leonard Books.
- [76] Jos Miguel Wisnik. 1999. *O som e o sentido*. Companhia das Letras.
- [77] Iannis Xenakis. 1992. *Formalized music: thought and mathematics in composition*. Number 6. Pendragon Press.
- [78] Joaquim Zamacois. 2002. *Curso de formas musicales: Con numerosos ejemplos musicales*. Barcelona : Idea Books.

Received XXXXX; revised XXXXX; accepted XXXXX