

Chapter 2

Sound modeling: signal-based approaches

Giovanni De Poli and Federico Avanzini

Copyright © 2006 by Giovanni De Poli and Federico Avanzini.

All rights reserved except for paragraphs labeled as *adapted from <reference>*.

2.1 Introduction

The sound produced by acoustic musical instruments is caused by the physical vibration of a certain resonating structure. This vibration can be described by signals that correspond to the time-evolution of the acoustic pressure associated to it. The fact that the sound can be characterized by a set of signals suggests quite naturally that some computing equipment could be successfully employed for generating sounds, for either the imitation of acoustic instruments or the creation of new sounds with novel timbral properties.

A wide variety of sound synthesis algorithms is currently available either commercially or in the literature. Each one of them exhibits some peculiar characteristics that could make it preferable to others, depending on goals and needs. Technological progress has made enormous steps forward in the past few years as far as the computational power that can be made available at low cost is concerned. At the same time, sound synthesis methods have become more and more computationally efficient and the user interface has become friendlier and friendlier. As a consequence, musicians can nowadays access a wide collection of synthesis techniques (all available at low cost in their full functionality), and concentrate on their timbral properties.

Each sound synthesis algorithm can be thought of as a digital model for the sound itself. Though this observation may seem quite obvious, its meaning for sound synthesis is not so straightforward. As a matter of fact, modeling sounds is much more than just generating them, as a digital model can be used for representing and generating a whole class of sounds, depending on the choice of control parameters. The idea of associating a class of sounds to a digital sound model is in complete accordance with the way we tend to classify natural musical instruments according to their sound generation mechanism. For example, strings and woodwinds are normally seen as timbral classes of acoustic instruments characterized by their sound generation mechanism. It should be quite clear that the degree of compactness of a class of sounds is determined, on one hand, by the sensitivity of the

digital model to parameter variations and, on the other hand, on the amount of control that is necessary for obtaining a certain desired sound. As an extreme example we may think of a situation in which a musician is required to generate sounds sample by sample, while the task of the computing equipment is just that of playing the samples. In this case the control signal is represented by the sound itself, therefore the class of sounds that can be produced is unlimited but the instrument is impossible for a musician to control and play. An opposite extremal situation is that in which the synthesis technique is actually the model of an acoustic musical instrument. In this case the class of sounds that can be produced is much more limited (it is characteristic of the mechanism that is being modeled by the algorithm), but the degree of difficulty involved in generating the control parameters is quite modest, as it corresponds to physical parameters that have an intuitive counterpart in the experience of the musician.

An interested conclusion that could be already drawn in the light of what stated above is that the compactness of the class of sounds associated to a sound synthesis algorithm is somehow in contrast with the “playability” of the algorithm itself. One should remember that the “playability” is of crucial importance for the success of a specific sound synthesis algorithm as, in order for a sound synthesis algorithm to be suitable for musical purposes, the musician needs an intuitive and easy access to its control parameters during both the sound design process and the performance. Such requirements often represents the reason why a certain synthesis technique is preferred to others.

Some considerations on control parameters are now in order. Varying the control parameters of a sound synthesis algorithm can serve several purposes, the first one of which is certainly that of exploring a sound space, i.e. producing all the different sounds that belong to the class characterized by the algorithm itself. This very traditional way of using control parameters would nowadays be largely insufficient by itself. As a matter of fact, with the progress in the computational devices that are currently being employed for musical purposes, the musician’s needs have turned more and more toward problems of timbral dynamics. For example, timbral differences between soft (dark) and loud (brilliant) tones are usually obtained through appropriate parameter control. Timbral expression parameters tend to operate at a note-level time-scale. As such, they can be suitably treated as signals characterized by a rather slow rate.

Another reason for the importance of time-variations in the algorithm parameters is that the musician needs to control the musical expression while playing. For example, staccato, legato, vibrato etc. need to be obtained through parameter control. Such parameter variations operate at a phrase-level time-scale. Because of that, they can be suitably treated as sequences of symbols events characterized by a very slow rate.

In conclusion, control parameters are signals characterized by their own time-scales. Controls signals for timbral dynamics are best described as discrete-time signals with a slow sampling rate, while controls for musical expression are best described by streams of asynchronous symbols events. As a consequence, the generation of control signals can once again be seen as a problem of signal synthesis.

2.2 Signal generators

In this category we find methods that directly generate the signal. We will see periodic waveform generators and noise generators.



2.2.1 Waveform generators

2.2.1.1 Digital oscillators

In many musical sounds, pitch is a characteristic to which we are quite sensitive. In examining the temporal waveform of pitched sounds, we see a periodic repetition of the waveform without great variations. The simplest synthesis method attempts to reproduce this characteristic, generating a periodic signal through a continuous repetition of a waveform. An algorithm that implements this method is called oscillator. A first method consists in computing the appropriate value of the function for every sample. Often function, as sinusoids are approximated by polynomial or rational truncated series. For example a sinusoid of frequency f can be computed by

$$s[n] = \sin(\omega n) = p(\omega(n \bmod F_s))$$

where $\omega = 2\pi f/F_s$. More efficient algorithms will be presented in the next sections.

2.2.1.2 Table lookup oscillator

A very efficient approach is to precompute the samples of the waveform, store them in a table which is usually implemented as a circular buffer, and access them from the table whenever needed. If we store in the table a copy of one period of the desired waveform, when we cycle over the wavetable with the aid of a circular pointer, we generate a periodic waveform. When the pointer reach the end of the table, it wraps around and points again at the beginning of the table. Given a table of length L , the period of the generated waveform is given by $T_L = LT$, and its fundamental frequency by $f_0 = F_s/L$. If we want to change the frequency, we would need the same waveform stored in tables of different lengths. A better solution is to store many equidistant points of the (continuous) waveform in the table, and then read the value in correspondence of the desired abscissa. Obviously, the more numerous are the points in the table, the better the approximation will be. The oscillator cyclically searches the table to get the point nearest to the required one. In this way, the oscillator resample the table to generate a waveform with different frequency. The distance in the table between two samples at subsequent instants is called *SI* (*sampling increment*) and is proportional to the frequency f of the generated sound:

$$f = \frac{SI \cdot F_s}{L} \quad (2.1)$$

If the sampling increment SI is greater than 1, it can happen that the highest frequencies of the waveform overcome the Nyquist frequency $F_N = F_s/2$ giving rise to foldover.

M-2.1

Implement in Matlab a circular look-up from a table of length L and with sampling increment SI .

M-2.1 Solution

```
phi=mod(phi +SI,L);
s=tab[phi];
```

where `phi` is a state variable indicating the reading point in the table, `A` is a scaling parameter, `s` is the output signal sample. The function `mod(x,y)` computes the remainder of the division x/y and is used here to implement circular reading of the table. Notice that `phi` can be a non integer value. In order to use it as array index, it can be truncated, or rounded to the nearest integer. A more accurate output can be obtained by linear interpolation between adjacent table values.

2.2.1.3 Recurrent sinusoidal signal generators

Sinusoidal signal can be generated also by recurrent methods. A first method is based on the second order resonator filter (that we will introduce in Sec. 2.5.1) with the poles lying on the unit circle. The equation is

$$y[n + 1] = 2 \cos(\omega)y[n] - y[n - 1] \quad (2.2)$$

where $\omega = 2\pi f/F_s$. With initial values $y[0] = 1$ and $y[-1] = \cos \omega$ the generator produces

$$y[n] = \cos n\omega = \cos(2\pi fTn)$$

With $y[0] = 0$ and $y[-1] = -\sin \omega$ the generator produces $y[n] = \sin n\omega$. In general if $y[0] = \cos \phi$ and $y[-1] = \cos(\phi - \omega)$ the generator produces $y[n] = \cos(n\omega + \phi)$. This property can be justified remembering the trigonometric relation $\cos \omega \cdot \cos \phi = 0.5[\cos(\phi + \omega) + \cos(\phi - \omega)]$.

Another method that combines both the sinusoidal and cosinusoidal generators is the so called coupled form described by the equations

$$x[n + 1] = \cos \omega \cdot x[n] - \sin \omega \cdot y[n] \quad (2.3)$$

$$y[n + 1] = \sin \omega \cdot x[n] + \cos \omega \cdot y[n] \quad (2.4)$$

With $x[0] = 1$ and $y[0] = 0$ we have $x[n] = \cos(n\omega)$ and $y[n] = \sin(n\omega)$. This property can be verified considering that if we define a complex variable $\alpha[n] = x[n] + jy[n] = \exp(jn\omega)$, it results $\alpha[n + 1] = \exp(j\omega) \cdot \alpha[n]$. The real and imaginary parts of this relation give equations 2.4 and 2.4 respectively.

Both methods have the drawback that coefficient quantization can give rise to numerical instability, i.e. poles outside the unitary circle. The waveform will then tend to grow exponentially or to decay rapidly into silence. To avoid this problem, a periodic re-initialization is advisable. It is possible to use a slightly different set of coefficients to produce absolutely stable sinusoidal waveforms

$$x[n + 1] = x[n] - c \cdot y[n] \quad (2.5)$$

$$y[n + 1] = c \cdot x[n + 1] + y[n] \quad (2.6)$$

where $c = 2 \sin(\omega/2)$. With $x[0] = 1$ and $y[0] = c/2$ we have $x[n] = \cos(n\omega)$.

2.2.1.4 Amplitude/frequency controlled oscillators

The amplitude and frequency of a sound are usually required to be time-varying parameters. Amplitude control is needed in order to define suitable sound envelopes, or to create *tremolo* effects (quasi-periodic amplitude variations around an average value). Frequency control is needed to simulate *portamento* between two tones, or subtle pitch variations in the sound attack/release, or *vibrato* effects (quasi-periodic pitch variations around an average value), and so on.

We then want to have at our disposal a digital oscillator of the form

$$s[n] = A[n] \cdot \text{tab}[\phi[n]], \quad (2.7)$$

where $A[n]$ scales the amplitude of the signal, and the phase $\phi[n]$ does not in general increase linearly in time and is computed as a function of the instantaneous frequency. Figure 2.1 shows the symbol usually adopted to depict an oscillator with fixed waveform and varying amplitude and frequency

Many sound synthesis languages (e.g., the well known *Csound*) define control signals at *frame rate*: a frame is a time window with pre-defined length (typically 5 to 50 ms), in which the control signals can be reasonably assumed to be approximately constant. This approximation clearly helps to reduce computational loads significantly.

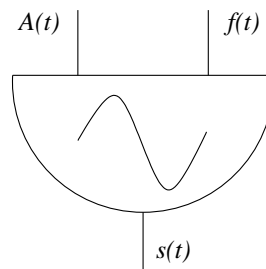


Figure 2.1: Symbol of the fixed waveform oscillator, with varying amplitude and frequency.

M-2.2

Assume that a function `sinosc(t0, a, f, ph0)` realizes a sinusoidal oscillator (t_0 is the initial time, a , f are the frame-rate amplitude and frequency vectors, and ph_0 is the initial phase). Then generate a sinusoid of length 2 s, with constant amplitude and frequency.

M-2.2 Solution

```

%%% headers %%%
global Fs;           %sample rate
global SpF;         %samples per Frame
Fs=22050;
ControlW=0.01;      %control window (in sec): 10 ms
SpF=round(Fs*ControlW);
Fc=Fs/SpF;          %control rate

%%% define controls %%%
slength=2;          %soundlength in seconds
nframes=slength*Fc; %total number of frames
a=ones(1,nframes); %constant amplitude
f=50*ones(1,nframes); %constant frequency

%%% compute sound %%%
s=sinosc(0,a,f,0); %sound signal

```

Note the structure of this simple example: in the “headers” section some global parameters are defined, that need to be known also to auxiliary functions; a second section defines the control parameters, and finally the audio signal is computed.

M-2.3

When the oscillator frequency is constant the phase is a linear function of time, $\phi(t) = 2\pi ft$, therefore in the digital domain ϕ can be computed as $\phi[n+1] = \phi[n] + 2\pi f/F_s$. In the more general case in which the frequency varies at frame rate, we have to understand how to compute the phase of the oscillator. The starting point is the equation

$$f(t) = \frac{1}{2\pi} \frac{d\phi}{dt}(t), \quad (2.8)$$

which simply says that the radian frequency $\omega(t) = 2\pi f(t)$ is the instantaneous angular velocity of the time-varying phase $\phi(t)$. If $f(t)$ is varying slowly enough (i.e. it is varying at frame rate), we can say that in the K -th frame the first-order approximation

$$\frac{1}{2\pi} \frac{d\phi}{dt}(t) = f(t) \sim f(T_K) + F_c [f(T_{K+1}) - f(T_K)] \cdot (t - T_K) \quad (2.9)$$

holds, where T_K, T_{K+1} are the initial instants of frames K and $K+1$, respectively. The term $F_c [f(T_{K+1}) - f(T_K)]$ approximates the derivative df/dt inside the K th frame. We can then find the phase by integrating equation (2.9):

$$\begin{aligned} \phi(t) &= \phi(T_k) + 2\pi f(T_k)(t - T_k) + 2\pi F_c [f(T_{K+1}) - f(T_K)] \frac{(t - T_K)^2}{2}, \\ \phi((K-1) \cdot \text{SpF} + n) &= \phi(K) + 2\pi \frac{f(K)n}{F_s} + \pi \frac{f(K+1) - f(K)}{\text{SpF} \cdot F_s} n^2, \end{aligned} \quad (2.10)$$

where $n = 0 \dots (\text{SpF} - 1)$ spans the frame. In summary, equation (2.10) computes ϕ at sample rate, given the frame rate frequencies. The key ingredient of this derivation is the linear interpolation (2.9).

Realize the `sinosc(t0, a, f, ph0)` function that we have used in M-2.2. Use equation (2.10) to compute the phase given the frequency vector `f`.

M-2.3 Solution

```
function s = sinosc(t0,a,f,ph0);

global SpF;           %samples per frame
global Fs;           %sampling rate

nframes=length(a);   %total number of frames
if (length(f)==1) f=f*ones(1,nframes); end
if (length(f)~=nframes) %check
    error('f and a must have the same length!');
end

s=zeros(1,nframes*SpF); %signal vector (initialized to 0)
lastampl=a(1);
lastfreq=f(1);
lastphase=ph0;
for i=1:nframes %cycle on the frames
    naux=1:SpF; %counts samples within frame
    ampl=lastampl +... %compute amplitudes within frame
        (a(i)-lastampl)/SpF.*naux;
    phase=lastphase+pi/Fs.*naux.* ... %compute phases within frame
        (2*lastfreq + (1/SpF)*(f(i)-lastfreq)).*naux;
    s((i-1)*SpF+1:i*SpF)=ampl.*cos(phase); %read from table
    lastampl=a(i); %save last values
    lastfreq=f(i); %of amplitude,
    lastphase=phase(SpF); %frequency, phase
end

s=[zeros(1,round(t0*Fs+1)) s]; %add initial silence of t0 sec.
```



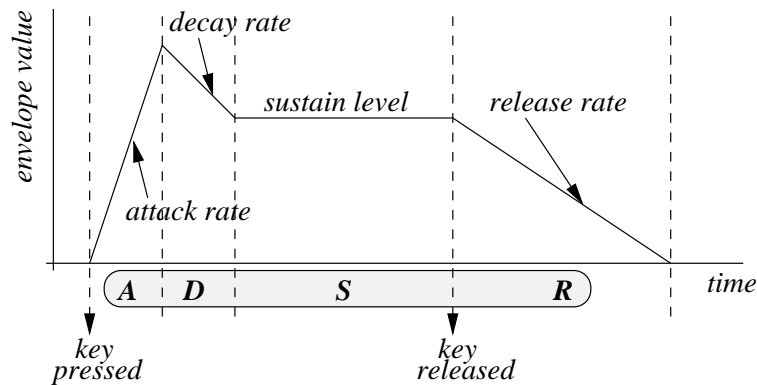


Figure 2.2: The four phases of an ADSR envelope over time.

Both the amplitude a and frequency f envelopes are defined at frame rate and are interpolated at sample rate inside the function body. Note in particular the computation of the `phase` vector within each frame.

2.2.1.5 Envelope generators

It is possible to use the same algorithm of table look-up oscillator to produce time envelopes. In this case, to generate a time envelope of d sec and scan only once the table, a sampling increment $SI = L/(F_s d)$ should be used. Often in sound synthesis, the amplitude envelope is described by a linearly varying function. A typical schema is the ADSR envelope (Fig. 2.2). The time envelope is described by four phases of *Attack*, *Decay*, *Sustain* e *Release*. When we want to change the tone duration, it is advisable to only slightly modify the attack and release, that marks the identity of the sound, while the sustain can be lengthened more freely. So the oscillator table will be read once with different sampling increments for the different parts of the generated envelope.

The use of waveform and envelope generators allows to generate quasi periodic sounds with very limited hardware and constitutes the building block of many more sophisticated algorithms.

M-2.4

Write a function that realizes a line-segment envelope generator. The input to the function are a vector of time instants and a corresponding vector of envelope values.

M-2.4 Solution

```
function env = envgen(t,a,method);           %t vector of time instants
                                           %a vector of envelope values

global SpF;                                %samples per frame
global Fs;                                  %sampling rate

if (nargin<3)
    method='linear';
end

frt=floor(t*Fs/SpF+1);                     %times instants as frame numbers
nframes=frt(length(frt));                  %total number of frames
env=interp1(frt,a,[1:nframes],method);    %linear interpolation
```

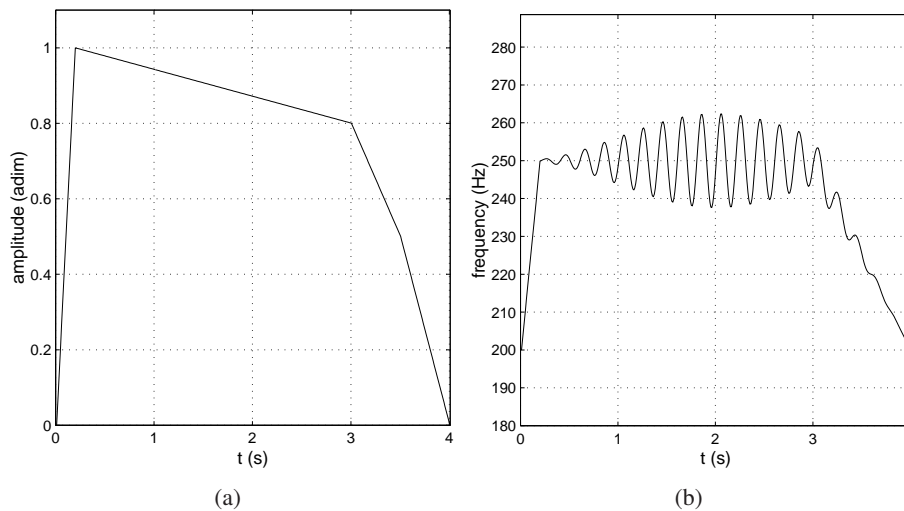


Figure 2.3: Amplitude (a) and frequency (b) control signals

The envelope shape is specified by break-points, described as couples (time instant (sec) and amplitude). The function generates the envelope at frame rate. Notice that the interpolation function `interp1` allows to easily use cubic or *spline* interpolations.

M-2.5

Synthesize a modulated sinusoid using the functions `sinosc` and `envgen`.

M-2.5 Solution

```

%%% headers %%%
%[...]

%%% define controls %%%
a=envgen([0, .2, 1, 1.5, 2], [0, 1, .8, .5, 0], 'linear'); %ADSR amp. envelope
f=envgen([0, .2, 1, 2], [200, 250, 250, 200], 'linear'); %pitch envelope
f=f+max(f)*0.05*... %pitch envelope with vibrato added
    sin(2*pi*5*(SpF/Fs)*[0:length(f)-1]).*hanning(length(f))';

%%% compute sound %%%
s=sinosc(0, a, f, 0);

```

In fig. 2.3 amplitude `a` and frequency `f` control signals are shown.

2.2.2 Noise generators

Up to now, we have considered signals whose behavior at any instant is supposed to be perfectly knowable. These signals are called deterministic signals. Besides these signals, *random signals* of unknown or only partly known behavior may be considered. For random signals, only some general characteristics, called statistical properties, are known or are of interest. The statistical properties are characteristic of an entire signal class rather than of a single signal. A set of random signals

is represented by a random process. Particular numerical procedures simulate random processes, producing sequences of random (or more precisely, pseudorandom) numbers.

Random sequences can be used both as signals (i.e., to produce white or colored noise used as input to a filter) and as control functions to provide a variety in the synthesis parameters most perceptible by the listener. In the analysis of natural sounds, some characteristics vary in an unpredictable way; their mean statistical properties are perceptibly more significant than their exact behavior. Hence, the addition of a random component to the deterministic functions controlling the synthesis parameters is often desirable. In general, a combination of random processes is used because the temporal organization of the musical parameters often has a hierarchical aspect. It cannot be well described by a single random process, but rather by a combination of random processes evolving at different rates. For example this technique is employed to generate $1/f$ noise.

2.2.2.1 Random noise models

White noise generators The spread part of the spectrum is perceived as random noise. In order to generate a random sequence, we need a random number generator. There are many algorithms that generate random numbers, typically uniformly distributed over the standardized interval $[0, 1)$. However it is hard to find good random number generators, i.e. that pass all or most criteria of randomness. The most common is the so called *linear congruential* generator. It can produce fairly long sequences of independent random numbers, typically of the order of two billion numbers before repeating periodically. Given an initial number (seed) $I[0]$ in the interval $0 \leq I[0] < M$, the algorithm is described by the recursive equations

$$\begin{aligned} I[n] &= (aI[n-1] + c) \bmod M \\ u[n] &= I[n]/M \end{aligned} \quad (2.11)$$

where a and c are two constants that should be chosen very carefully in order to have a maximal length sequence, i.e. long M samples before repetition. The actual generated sequence depends on the initial value $I[0]$; that is why the sequence is called pseudorandom. The numbers are uniformly distributed over the interval $0 \leq u[n] < 1$. The mean is $E[u] = 1/2$ and the variance is $\sigma_u^2 = 1/12$. The transformation $s[n] = 2u[n] - 1$ generates a zero-mean uniformly distributed random sequence over the interval $[-1, 1)$. This sequence corresponds to a white noise signal because the generated numbers are mutually independent. The power spectral density is given by $S(f) = \sigma_u^2$. Thus the sequence contains all the frequencies in equal proportion and exhibits equally slow and rapid variation in time.

With a suitable choice of the coefficients a and b , it produces pseudorandom sequences with flat spectral density magnitude (white noise). Different spectral shapes can be obtained using white noise as input to a filter.

M-2.6

A method of generating a Gaussian distributed random sequence is based on the central limit theorem, which states that the sum of a large number of independent random variables is Gaussian. As exercise, implement a very good approximation of a Gaussian noise, by summing 12 independent uniform noise generators.

Pink noise generators If we desire that the numbers vary at a slower rate, we can generate a new random number every d sampling instants and hold the previous value in the interval (*holder*) or

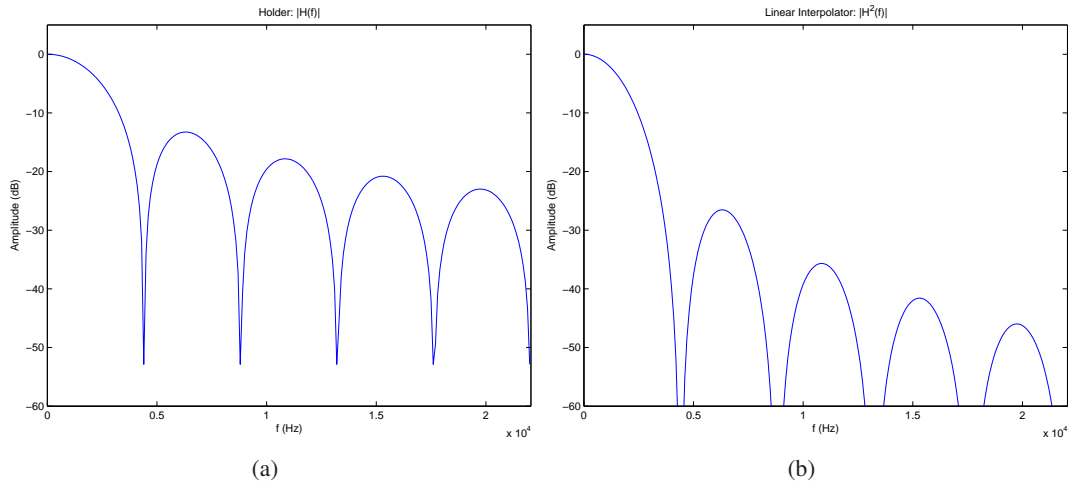


Figure 2.4: Spectral envelope $|H(f)|$ of low frequency noise generators where a new random number is generated every $d = 10$ samples: (a) hold generator; (b) linear interpolator.

interpolate between two successive random numbers (interpolator). In this case the power spectrum is given by

$$S(f) = |H(f)|^2 \frac{\sigma_u^2}{d}$$

with

$$|H(f)| = \left| \frac{\sin(\pi f d / F_s)}{\sin(\pi f / F_s)} \right|$$

for the holder (fig. 2.4(a)) and

$$|H(f)| = \frac{1}{d} \left[\frac{\sin(\pi f d / F_s)}{\sin(\pi f / F_s)} \right]^2$$

for linear interpolation (fig. 2.4(b)).

1/f noise generators A 1/f noise, also called pink noise, is characterized by a power spectrum that fall in frequency like 1/f

$$S(f) = \frac{A}{f} \quad (2.12)$$

To avoid the infinity at $f = 0$, this behaviour is assumed valid for $f \geq f_{min}$, where f_{min} is a desired minimum frequency. The spectrum is characterized by a 3 db per octave drop, i.e. $S(2f) = S(f)/2$. The amount of power contained within a frequency interval $[f_1, f_2]$ is

$$\int_{f_1}^{f_2} S(f) df = A \ln \left(\frac{f_1}{f_2} \right)$$

This implies that the amount of power in any octave is the same. 1/f noise is ubiquitous in nature and is related to fractal phenomena. In audio domain it is known as pink noise. It represents the psychoacoustic equivalent of the white noise because he approximately excites uniformly the critical bands. The physical interpretation is a phenomenon that depends on many processes that evolve on

different time scales. So a $1/f$ signal can be generated by the sum of several white noise generators that are filtered through first-order filters having the time constants that are successively larger and larger, forming a geometric progression.

M-2.7

In the Voss $1/f$ noise generation algorithm, the role of the low pass filters is played by the hold filter seen in the previous paragraph. The $1/f$ noise is generated by taking the average of several periodically held generators $y_i[n]$, with periods forming a geometric progression $d_i = 2^i$, i.e.

$$y[n] = \frac{1}{M} \sum_{i=1}^M y_i[n] \quad (2.13)$$

The power spectrum does not have an exact $1/f$ shape, but it is close to it for frequencies $f \geq F_s/2^M$. As exercise, implement a $1/f$ noise generator and use it to assign the pitches to a melody.

M-2.8

The music derived from the $1/f$ noise is closed to the human music: it does not have the unpredictability and randomness of white noise nor the predictability of brown noise. $1/f$ processes correlate logarithmically with the past. Thus the averaged activity of the last ten events has as much influence on the current value as the last hundred events, and the last thousand. Thus they have a relatively long-term memory.

$1/f$ noise is a fractal one; it exhibits self-similarity, one property of the fractal objects. In a self-similar sequence, the pattern of the small details matches the pattern of the larger forms, but on a different scale. In this case, is used to say that $1/f$ fractional noise exhibits statistical self-similarity. The pink noise algorithm for generating pitches has become a standard in algorithmic music. Use the $1/f$ generator developed in M-2.7 to produce a fractal melody.

2.3 Time-segment based models

2.3.1 Wavetable synthesis

2.3.1.1 Definitions and applications

Finding a mathematical model that faithfully imitates a real sound is an extremely difficult task. If an existing reference sound is available, however, it is always possible to reproduce it through recording. Such a method, though simple in its principle, is widely adopted by digital sampling instruments or samplers and is called wavetable synthesis or sampling. Samplers store a large quantity of examples of complete sounds, usually produced by other musical instruments. When we wish to synthesize a sound we just need to directly play one sound of the stored repertoire.

The possibility of modification is rather limited, as it would be for the sound recorded by a tape deck. The most common modification is that of somewhat varying the sampling rate (speed) when reproducing the sound, which results in a pitch deviation. On the other hand, what makes the method interesting the most is certainly the variety of sounds available.

From the implementation viewpoint, computational simplicity and limited amount of information to be stored are two contrasting needs for samplers. In fact, in order to reduce the data to be stored, it is possible to adopt looping techniques with almost any stationary portion of sounds. One method of improving the expressive possibilities of samplers is store multiple of the sounds at different pitches, and switching or interpolating between these upon synthesis. This method, called *multisampling*, also might include the storage of separate samples for loud and soft sounds. During synthesis a linear interpolation between these sampled is performed as function of the desired loudness. Infact most

instruments exhibit richer spectra for louder sounds. Sometimes a filter is used to control the spectral variation allowing to obtain softer sounds from a stored loud sound.

In most cases sampling techniques are presented as a method for reproducing natural sounds and is evaluated in comparison with the original instruments. This is the main reason why the most popular commercial digital keyboards, such as electronic pianos and organs, adopt this synthesis technique. Of course, sampling cannot feature all the expressive possibilities of the original instrument. Notice that sampled sounds can also be obtained synthetically or through the modification of other sounds, which is a way of widening the range of possibilities of application of samplers. From the composer's viewpoint, the use of samplers represents a practical approach to the so-called *musique concrète*. This kind of music, begun in Paris in late Forties by the main effort of Pierre Schaeffer, started to use, as basic sonic material of its musical compositions, any kind of sound, recorded by a microphone and eventually processed.

2.3.1.2 Transformations: pitch shift, looping

The most common modification is that of varying the sampling rate (speed) when reproducing the sound, which results in a pitch deviation. In digital domain this effect is obtained by resampling the stored waveform scanning the table with a sampling increment different than 1. The algorithm is similar to the digital oscillator by table look-up (see Sect.2.2.1.1). In this case if we want to change the frequency f_{stored} of the stored sound to a new frequency f_{new} , we will read the table with a sampling increment

$$SI = \frac{f_{new}}{f_{stored}}$$

However, substantial pitch variations are generally not very satisfactory as a temporal waveform compression or expansion results in unnatural timbral modifications, which is exactly what happens with an accelerated tape recorder. It is thus necessary to allow only pitch variations of few semitones to take place for the synthetic sound to be similar to the original one. This is especially true for sounds characterized by formants in the spectrum. For a good sampling synthesis of the whole pitch extension, many samples should be stored (e.g. three for each octave). Moreover special care should be paid to assure that adjacent sounds be similar.

M-2.9

Import a .wav file of a single instrument tone. Scale it (compress and expand) to different extents and listen to the new sounds. Up to what scaling ratio are the results acceptable?

Often it is desired to vary the sound also in function of other parameters, the most important being the intensity. To this purpose it is not sufficient to change the sound amplitude by a multiplication, but it is necessary to modify the timbre of the sound. In general louder sounds are characterized by a sharper attack and by a brighter spectrum. In this case a technique could be to use a unique sound prototype (e.g. a tone played fortissimo) and then obtaining the other intensity by simple spectral processing, as low pass filtering. A different and more effective solution, is to use a set of different sound prototype, recorder with different intensity (e.g. tones played fortissimo, mezzo forte, pianissimo) and then obtaining the other dynamic values by interpolations and further processing.

This technique is thus characterized by high computational efficiency and high imitation quality, but by low flexibility for sounds not initially included in the repertoire or not easily obtainable with simple transformations. There is a trade-off of memory size with sound fidelity.

In order to employ efficiently the memory, often the sustain part of the tone is not entirely stored but only a part (or few significant parts) and in the synthesis this part is repeated (*looping*). Naturally

the repeated part should not be too short, to avoid a static character of the resulting sound. For example to lengthen the duration of a note, first the attack is reproduced without modification, then the sustain part is cyclically repeated, with possible cross interpolation among the different selected parts, and finally the sound release stored part is reproduced. Notice that if we want to avoid artefacts in cycling, particular care should be devoted to choosing the points of the beginning and ending of the loop. Normally an integer number of periods is used for looping starting with a null value, to avoid amplitude or phase discontinuities. In fact these discontinuities are very annoying. To this purpose it may be necessary to process the recorded samples by slightly changing the phases of the partials.

M-2.10

Import a .wav file of a single instrument tone. Find the stationary (sustain) part, isolate a section, and perform the looping operation. Listen to the results, and listen to the artifacts when the looped section does not start/end at zero-crossings.

If we want a less static sustain, it is possible to individuate some different and significant sound segments, and during the synthesis interpolate (*cross-fade*) among subsequent segments. In this case the temporal evolution of the tone can be more faithfully reproduced.

2.3.2 Granular synthesis

Granular synthesis, together with additive synthesis, shares the idea of building complex sounds from simpler ones. Granular synthesis assumes that a sound can be considered as a sequence, possibly with overlaps, of elementary acoustic elements called grains. Granular synthesis constructs complex and dynamic acoustic events starting from a large quantity of grains. The features of the grains and their temporal location determine the sound's timbre. We can see it as being similar to the cinema, where a rapid sequence of static images gives the impression of objects in movement.

The initial idea of granular synthesis dates back to Gabor's work aimed at pinpointing the physical and mathematical ideas needed to understand what a time-frequency spectrum is. He considered sound as a sum of elementary Gaussian functions that have been shifted in time and frequency. Gabor considered these elementary functions as *acoustic quanta*, the basic constituents of a sound. In the scientific field these works have been rich in implications and have been the starting point for studying time-frequency representations. The usual Gabor expansion on a rectangular time-frequency lattice of a signal $x(t)$ can be expressed as a linear combination of properly shifted and modulated versions $g_{mk}(t)$ of a synthesis window $g(t)$

$$x(t) = \sum_m \sum_k a_{mk} g_{mk}(t)$$

with

$$g_{mk} = g(t - m\alpha T) e^{jk\beta\Omega t}$$

The time step αT and the frequency step $\beta\Omega$ satisfy the relationship $\Omega T = 2\pi$ and $\alpha\beta \leq 1$.

In music, granular synthesis arises from the experiences of taped electronic music. In the beginning musicians had tools that did not allow a great variation of timbre, for example fixed waveform oscillators and filters. They obtained dynamic sounds by cutting tapes into short sections and the putting together again. The rapid alternation of acoustic elements give a certain variety to the sound. The source of the sound could be electronic or live recorded sounds that were sometimes electronically processed. Iannis Xenakis developed this method in the field of analog electronic music. Starting from Gabor's theory, Xenakis considers the grains as being music quanta and suggested a method of

composition that is based on the organization of the grains by means of screen sequences, which specify the frequency and amplitude parameters of the grains at discrete points in time. In this way a common conceptual approach is used both for micro and macro musical structure.

2.3.2.1 Sound granulation

Two main approaches to granular synthesis can be identified: the former based on sampled sounds and the latter based on abstract synthesis. In the first case, *sound granulation*, complex waveforms, extracted from real sounds or described by spectra, occurs in succession with partial overlap with the method called Overlap And Add (OLA). In this way, it is possible both to reproduce accurately real sounds and modify them in their dynamic characteristics.

Let $x[n]$ and $y[n]$ be the input and output signals. The grains $g_k[i]$ are extracted from the input signal with the help of a window function $w_k[i]$ of length L_k by

$$g_k[i] = x[i + i_k] w_k[i]$$

with $i = 0 \dots L_k - 1$. The time instant i_k indicates the point where the segment is extracted; the length L_k determines the amount of signal extracted; the window waveform $w_k[i]$ should ensure fade-in and fade-out at the border of the grain and affects the frequency content of the grain. Long grains tend to maintain the timbre identity of the portion of the input signal, while short ones acquire a pulse-like quality. When the grain is long, the window has a flat top and it is used only to fade-in and fade-out the borders of the segment.

As in additive synthesis the organization of the choice of the frequencies is very important, so in granular synthesis the proper timing organization of the grain is essential to avoid artifacts produced by discontinuities. This problem makes often the control quite difficult. An example of use is the synthesis of a time varying stochastic component of a signal. In this case it is only necessary to control the spectral envelope. To this purpose, it is convenient to employ the inverse Fourier transform of a spectrum, whose magnitude is defined by the spectral envelope and the phase are generated randomly. Every frame is multiplied by a suitable window before the overlap-and-add, i.e. the sum of the various frames partially overlapped. It is possible also to use this approach for transforming sampled sounds (sound granulation). In this case grains are built selecting short segments of a sound, previously or directly recorded, and then are shaped by an amplitude envelope. These grains are used by the composer with a different order and time location or a different velocity. Notice that it is possible to extract grains from different sound files to create hybrid textures, e.g. evolving from one texture to another (fig. 2.5).

2.3.2.2 Synthetic grains

In abstract granular synthesis (second case), grains consist of arbitrary waveforms whose amplitude envelope is a short Gaussian function. Often frequency modulated Gaussian functions are used in order to localize the energy both in frequency and time domain. Grains are scattered on the frequency-time plane in the form of clouds. Notice that it is possible to recognize a certain similarity between this type of granular synthesis and the technique of mosaics, where the grains are single monochromatic tesseras and their juxtaposition produces a complex image.

In this case the waveform of the i -th grain (fig. 2.6) is given by

$$g_i[n] = w_i[n] \cdot \cos \left(2\pi \frac{f_i}{F_s} n + \phi_i \right)$$

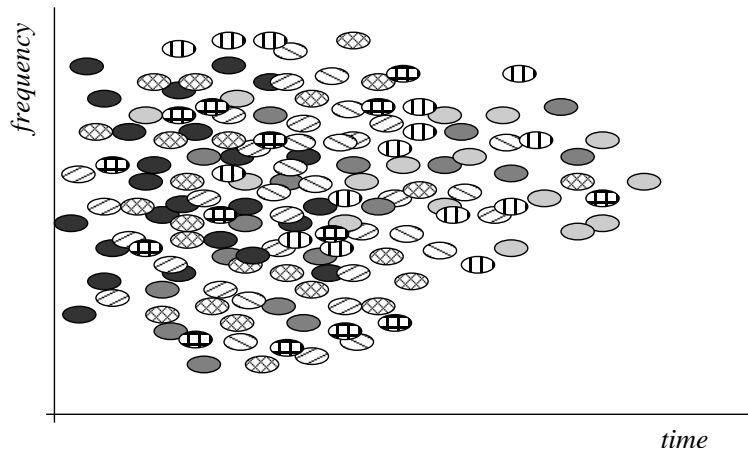


Figure 2.5: Representation of granular synthesis where grains derived from different sources are randomly mixed.

where $w_i[n]$ is a window of duration N_i samples. The synthesis expression is given by

$$s[n] = \sum_i a_i \cdot g_i[n - n_i] \quad (2.14)$$

where a_i is the amplitude coefficient of the i -th grain and n_i is its time instant. Every grain contributes to the total energy around the point (n_i, f_i) of the time frequency plane.

The most important and classic type of granular synthesis, (*asynchronous granular synthesis*), is when simple grains are irregularly distributed in the time-frequency plane, e.g. they are scattered onto a mask that delimitate a particular area of the time-frequency-amplitude space. It results a cloud of micro-sounds or a sonic texture that is time-varying. Moreover the grain density inside the mask can be controlled. In this way many sound textures and natural noisy sounds can be modelled where general statistical properties are more important than the exact sound evolution. These kind of sounds can be defined as the accumulation of more or less complex sonic grains, with their proper temporal and spectral variability. For example, the sound of rice falling onto a metal plate is composed of thousands of elementary ticks; the rain produces, in the same way, the accumulation of a large amount of water droplet micro-sounds. Scratching or cracking sounds made by the accumulation of thousands of complex micro-sounds not necessarily deterministic. In fact, in the real world, when multiple realizations of a same event, of a same phenomenon occur, we can expect these types of sounds.

Grain duration affects the sonic texture: short duration (few samples) produces a noisy, particulate disintegration effect; medium duration (tents of ms) produces fluttering, warbling, gurgling; longer durations (hundreds of ms) produce aperiodic tremolo, jittering spatial position. When the grains are distributed on a large frequency region, the texture has a massive character, while when the band is quite narrow, it result a pitched sound. Sparse densities (e.g. 5 grains per second) give rise to a pointillistic texture.

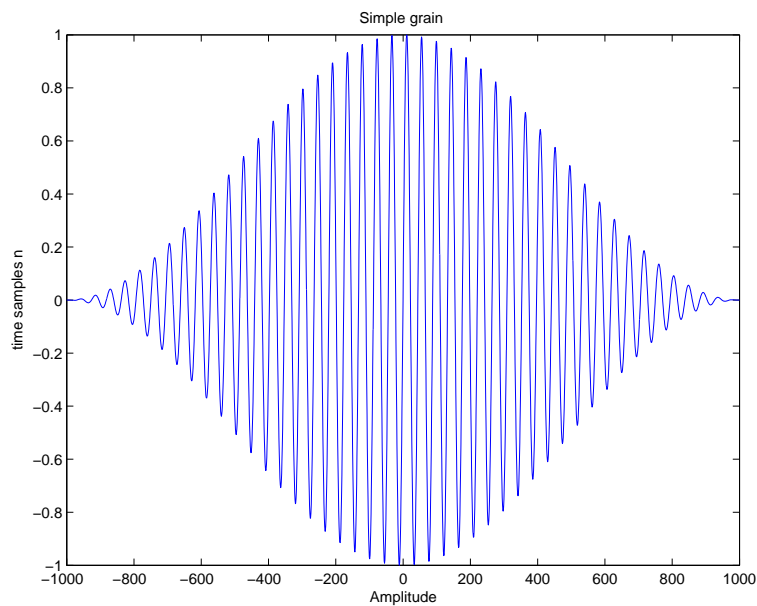


Figure 2.6: Example of a synthetic grain waveform.

2.3.3 Overlap-Add (OLA) methods

2.3.3.1 Time-domain OLA

The definition Overlap-Add (OLA) refers in general to a family of algorithms that produce a signal by properly assembling a number of signal segments. OLA methods are developed both in the time domain and in the frequency domain, here we are interested in reviewing briefly time-domain approaches.

Consider a sound signal $x[n]$. A time-frequency representation of $x[n]$ can be seen as a series of overlapping DFTs, typically obtained by windowing x in the desired frame. More precisely, we have that a frame $X_m(e^{j\omega_k})$ of the STFT is the DFT of the signal windowed segment $x_m[n] = x[n]w_a[n - mS_a]$, where $w_a[n]$ is the chosen analysis window and S_a is the *analysis hop-size*, i.e. the time-lag (in samples) between one analysis frame and the following one. If the window w_a is N samples long, then the *block size*, i.e. the length of each frame X_m , will be N . In order for the signal segments to actually overlap, the inequality $S_a \leq N$ must be verified. When $S_a = N$ the segments are exactly juxtaposed with no overlap.

Given the above signal segmentation, Overlap-Add (OLA) methods are typically used to modify and reconstruct the signal in two main steps:

1. Any desired modification is applied to the spectra (e.g. multiplying by a filter frequency response function), and modified frame spectra $Y_m(e^{j\omega_k})$ are obtained.
2. Windowed segments $y_m[n]$ of the modified signal $y[n]$ are obtained by computing the inverse DFT (IDFT) of the frames Y_m .
3. The output is reconstructed by overlapping-adding the windowed segments: $y[n] = \sum_m y_m[n]$.

Figure 2.7 illustrates the procedure.

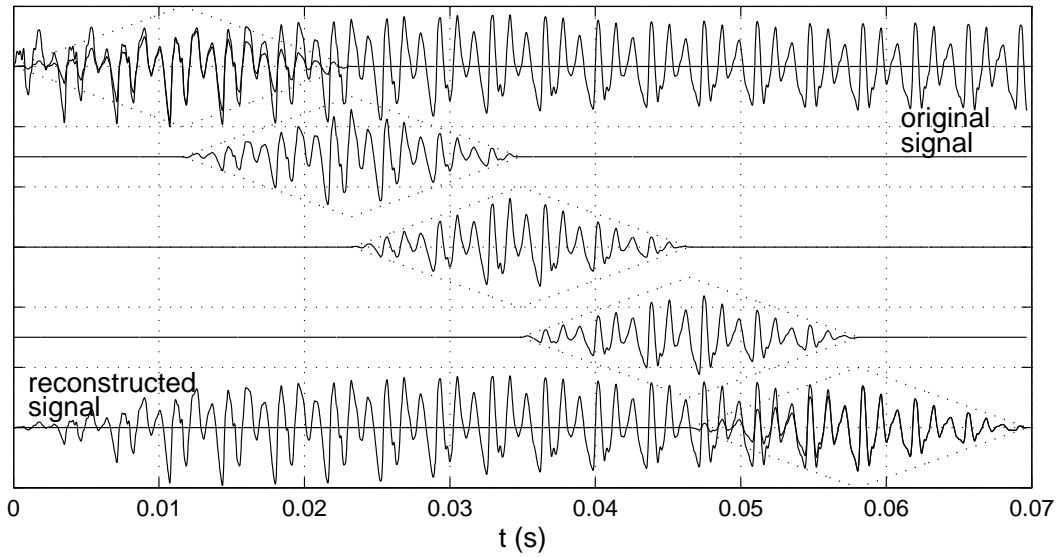


Figure 2.7: An example of Overlap-Add signal reconstruction, with triangular windowing.

In the absence of spectral modifications, this procedure reduces to an identity ($y[n] \equiv x[n]$) if the overlapped and added analysis windows w_a sum to unity:

$$\sum_m x_m[n] = \sum_m x[n]w_a[n - mS_a] = x[n] \quad \Leftrightarrow \quad A_{w_a}[n] \triangleq \sum_m w_a[n - mS_a] \equiv 1. \quad (2.15)$$

If this condition does not hold, then the function A_{w_a} acts on the reconstructed signal as a periodic *amplitude modulation envelope*, with period S_a . This kind of frame rate distortion can be seen in the frequency domain as a series of sidebands with spacing F_s/S_a in a spectrogram of the output signal. In fact, one may prove that the condition $A_{w_a} \equiv 1$ is equivalent to the condition $W(e^{j\omega_k}) = 0$ at all harmonics of the frame rate F_s/S_a .

In their most general formulation OLA methods utilize a *synthesis* window w_s that can in general be different from the analysis window w_a . In this case the second step of the procedure outlined above is modified as follows:

2. Windowed segments $y_m[n]$ of the modified signal $y[n]$ are obtained by **(a)** computing the inverse DFT (IDFT) of the frames Y_m , **(b)** dividing by the analysis window w_a (assuming that it is non-zero for all samples), and **(c)** multiplying by the synthesis window.

This approach provides greater flexibility than the previous one: the analysis window w_a can be chosen only on the basis of its time-frequency resolution properties, but needs not to satisfy the “sum-to-unity” condition $A_{w_a} \equiv 1$. On the other hand, the synthesis window w_s is only used to cross-fade between signal segments, therefore one should only ensure that $A_{w_s} \equiv 1$. We will see in section 2.4.2 an application of this technique to frequency-domain implementation of additive synthesis.

Many digital sound effects can be obtained by employing OLA techniques. As an example, a *robotization* effect can be obtained by putting zero phase values on every FFT before reconstruction: the effect applies a fixed pitch onto a sound and moreover, as it forces the sound to be periodic, many erratic and random variations are converted into “robotic” sounds. Other effects are obtained by imposing a random phase on a time-frequency representation, with different behaviors depending

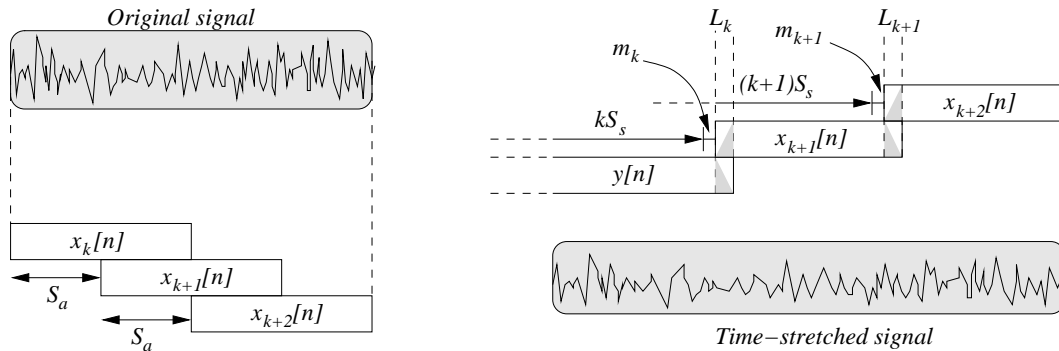


Figure 2.8: The generic SOLA algorithmic step

on the block length N : if N is large (e.g. $N = 2048$ with $F_s = 44.1$ kHz), the magnitude will represent the behavior of the partials quite well and changes in phase will produce an uncertainty over the frequency; if N is small (e.g. $N = 64$ with $F_s = 44.1$ kHz), the spectral envelope will be enhanced and this will lead to a whispering effect.

A widely studied effect is *time-stretching*, i.e. contraction or expansion of the duration of an audio signal. Time-stretching algorithms can be used in a number of applications: think about wavetable synthesis, post-synchronization of audio and video, speech technology at large, and so on. A time-stretching algorithm should ideally shorten or lengthen a sound file composed of N_{tot} samples to a new desired length $N'_{tot} = \alpha N_{tot}$, where α is the *stretching factor*. Note that a mere resampling of the sound signal does not provide the desired result, since it has the side-effect of transposing the sound: in this context resampling is the digital equivalent of playing the tape at a different speed.

What we really want is a scaling of the perceived timing attributes without affecting the perceived frequency attributes. More precisely, we want the time-scaled version of the audio signal to be perceived as the same sequence of *acoustic events* as the original signal, only distributed on a compressed/expanded time pattern. As an example, a time-stretching algorithm applied to a speech signal should change the speaking rate without altering the pitch.

Time-segment processing techniques described in this section are one possible approach to time-stretching effects. The simple OLA algorithm described above can be adapted to this problem by defining an *analysis hop size* S_a and a *synthesis hop size* $S_s = \alpha S_a$, where α is the *stretching factor* that will be applied to the output. An input signal $x[n]$ is then segmented into frames $x_k[n]$, each taken every S_a samples. The output signal $y[n]$ is produced by reassembling the same frames $x_k[n]$, each added to the preceding one every S_s samples. However this repositioning of the input segments with respect to each other destroys the original phase relationships, and constructs the output signal by interpolating between these misaligned segments. This cause pitch period discontinuities and distortions that can produce heavily audible artifacts in the output signal.

2.3.3.2 Synchronous and pitch-synchronous OLA

In order to avoid phase discontinuities at the boundaries between frames, a proper time alignment of the blocks has to be chosen. The *SOLA* (Synchronous OverLap and Add) algorithm realizes such a proper alignment, and provides a good sound quality (at least for values of α not too far from 1) while remaining computationally simple, which makes it suitable even for real-time applications. Most of commercial time-stretching software packets adopt the SOLA algorithm.

Let N be the analysis block length. In the initialization phase the SOLA algorithm copies the first N samples from $x_1[n]$ to the output $y[n]$, to obtain a minimum set of samples to work on:

$$y[j] = x_1[j] = x[j] \quad \text{for } j = 0 \dots N - 1. \quad (2.16)$$

Then, during the generic k th step the algorithm tries to find the optimal overlap between the last portion of the output signal $y[n]$ and the incoming analysis frame $x_{k+1}[n]$. More precisely, $x_{k+1}[n]$ is pasted to the output $y[n]$ starting from sample $kS_s + m_k$, where m_k is a small discrete time-lag that optimizes the alignment between y and x_k (see Fig. 2.8). Note that m_k can in general be positive or negative, although for clarity we have used a positive m_k in Fig. 2.8.

When the optimal time-lag m_k is found, a *linear crossfade* is used within the overlap window, in order to obtain a gradual transition from the last portion of y to the first portion of x_k . Then the last samples of x_k are pasted into y . If we assume that the overlap window at the k th SOLA step is L_k samples long, then the algorithmic step computes the new frame of the input y as

$$y[kS_s + j] = \begin{cases} (1 - v[j])y[kS_s + j] + v[j]x_k[j] & \text{for } m_k \leq j \leq L_k \\ x_k[j] & \text{for } L_k + 1 \leq j \leq N \end{cases} \quad (2.17)$$

where $v[j]$ is a linear *smoothing function* that realizes the crossfade between the two segments. The effect of Eq. (2.17) is a local replication or suppression of waveform periods (depending on the value of α), that eventually results in an output signal $y[n]$ with approximately the same spectral properties of the input $x[n]$, and an altered temporal evolution.

At least three techniques are commonly used in order to find the optimal value for the discrete time lag m_k at each algorithmic step k :

1. Computation of the minimum vectorial inter-frame distance in an L_1 sense (cross-AMDF)
2. Computation of the maximum cross-correlation $r_k(m)$ in a neighborhood of the sample kS_s . Let M be the width of such neighborhood, and let $y_{M_k}[i] = y[kS_s + i]$ for $i = 1 \dots M - 1$, and $x_{M_k}[i] = x_{k+1}[i]$ for $i = 1 \dots M - 1$. Then the cross-correlation $r_k(m)$ is computed as

$$r_k[m] \triangleq \sum_{i=0}^{M-m-1} y_{M_k}[i] \cdot x_{M_k}[i + m], \quad m = -M + 1, \dots, M - 1. \quad (2.18)$$

Then m_k is chosen to be the index of maximal cross-correlation: $r_k[m_k] = \max_m r_k[m]$.

3. Computation of the maximum *normalized* cross-correlation, where every value taken from the cross-correlation signal is normalized by dividing it by the product of the frame energies.

The latter technique is conceptually preferable, but the second one is often used for efficiency reasons.

M-2.11

Write a function `sola_timestretch(x, N, Sa, alpha, L)` that realizes the time-stretching SOLA algorithm through segment cross-correlation.

M-2.11 Solution

```
function y = sola_timestretch(x, N, Sa, alpha, L)

% Params: block length N; analysis hop-size Sa;
```

```

%           stretching factor alpha; overlap interval L

Ss = round(Sa*alpha); %synthesis hop-size
if (Sa > N)           disp('Sa must be less than N !!!'); end
if (Ss >= N)         disp('alpha is not correct, Ss is >= N');
elseif (Ss > N-L)    disp('alpha is not correct, Ss is > N-L'); end
if (rem(L,2)~= 0)    L = L+1; end

M = ceil(length(x)/Sa); %number of frames
x(M*Sa+N)=0;          %now x is exactly M*Sa samples
y(1:N) = x(1:N); %first frame of x is written into y;

for m=1:M-1           %%loop over frames
    frame=x(m*Sa+1:N+m*Sa); %current analysis frame
    framecorr=xcorr(frame(1:L),y(1,m*Ss:m*Ss+(L-1)));
    [corrmax,imax]=max(framecorr); %find point of max xcorrelation

    imax = L;
    xfadewin = (m*Ss-(L-1)+imax-1):length(y);
    fadein = (0:length(xfadewin)-1)/length(xfadewin); %from 0 to 1
    fadeout = 1 - fadein; %from 1 to 0

    y=[y(1,1:(xfadewin(1)-1)) ...
        (y(1,xfadewin).*fadeout +frame(1:length(xfadewin)).*fadein) ...
        frame(length(xfadewin)+1:N)];
end

```

A variation of the SOLA algorithm for time stretching is the *Pitch Synchronous Overlap-Add* (PSOLA) algorithm, which is especially used for voice processing. PSOLA assumes that the input sound is pitched, and exploits the pitch information to correctly align the segments and avoid pitch discontinuities. Here we only provide a qualitative description of the algorithm, which is composed of two phases: analysis/segmentation of the input sound, and resynthesis of the time-stretched output signal. The analysis phase works in two main steps:

1. Determination of the pitch period of the input signal $x[n]$ and determination of time instants (“pitch marks”) between which the pitch can be considered constant. This step is clearly the most critical and computationally expensive one of the whole algorithm.
2. Segmentation of the input signal, where the segments $x_m[n]$ have a block length N of two pitch periods and are centered at every “pitch mark”. Proper windowing has also to be used in order to ensure fade-in and fade-out. This procedure implies that the analysis hop-size is $S_a = N/2$.

The synthesis also uses a hop-size $S_s = N/2$. However, the segments x_m are re-assembled using a different approach with respect to SOLA. With a stretching factor $\alpha > 1$ (time expansion) some segments will *repeated*, while in the case of $\alpha < 1$ (time compression) some segments will be discarded in the resynthesis.

2.4 Spectrum based models

Since the human ear acts as a particular spectrum analyser, a first class of synthesis models aims at modeling and generating sound spectra. The Short Time Fourier Transform and other time-frequency



representations provide powerful sound analysis tools for computing the time-varying spectrum of a given sound. In this section models that can be interpreted in the frequency domain will be presented.

2.4.1 Sinusoidal model

When we analyze a pitched sound, we find that its spectral energy is mainly concentrated at a few discrete (slowly time-varying) frequencies f_k . These frequency lines correspond to different sinusoidal components called partials. If the sound is almost periodic, the frequencies of partials are approximately multiple of the fundamental frequency f_0 , ie. $f_k(t) \simeq k f_0(t)$. The amplitude a_k of each partial is not constant and its time-variation is critical for timbre characterization. If there is a good degree of correlation among the frequency and amplitude variations of different partials, these are perceived as fused to give a unique sound with its timbre identity.

The sinusoidal model assumes that the sound can be modeled as a sum of sinusoidal oscillators whose amplitude $a_k(t)$ and frequency $f_k(t)$ are slowly time-varying

$$s_s(t) = \sum_k a_k(t) \cos(\phi_k(t)), \quad (2.19)$$

$$\phi_k(t) = 2\pi \int_0^t f_k(\tau) d\tau + \phi_k(0), \quad (2.20)$$

or, digitally,

$$s_s[n] = \sum_k a_k[n] \cos(\phi_k[n]), \quad (2.21)$$

$$\phi_k[n] = 2\pi f_k[n]T_s + \phi_k[n-1], \quad (2.22)$$

where T_s is the sampling period. Notice that eq. 2.22 can also be written as

$$\phi_k(n) = \phi_{0k} + 2\pi T_s \sum_{j=1}^n f_k(j)$$

. Equations (2.19) and (2.20) are a generalization of the Fourier theorem, that states that a periodic sound of frequency f_0 can be decomposed as a sum of harmonically related sinusoids of frequency $f_i = i f_0$

$$s_s(t) = \sum_k a_k \cos(2\pi k f_0 t + \phi_k).$$

This model is also capable of representing aperiodic and inharmonic sounds, as long as their spectral energy is concentrated near discrete frequencies (spectral lines).

In computer music this model is called *additive synthesis* and is widely used in music composition. Notice that the idea behind this method is not new. As a matter of fact, additive synthesis has been used for centuries in some traditional instruments such as organs. Organ pipes, in fact, produce relatively simple sounds that, combined together, contribute to the richer spectrum of some registers. Particularly rich registers are created by using many pipes of different pitch at the same time. Moreover this method, developed for simulating natural sounds, has become the “metaphorical” foundation of a compositional methodology based on the expansion of the time scale and the reinterpretation of the spectrum in harmonic structures.

2.4.2 Spectral modeling

Spectral analysis of the sounds produced by musical instruments, or by any physical system, shows that the spectral energy of the sound signals can be interpreted as the sum of two main components: a *deterministic* component that is concentrated on a discrete set of frequencies, and a *stochastic* component that has a broadband characteristics. The deterministic –or sinusoidal– component normally corresponds to the main modes of vibration of the system. The stochastic residual accounts for the energy produced by the excitation mechanism which is not turned into stationary vibrations by the system, and for any other energy component that is not sinusoidal.

As an example, consider the sound of a wind instrument: the deterministic signal results from self-sustained oscillations inside the bore, while the residual noisy signal is generated by the turbulent flow components due to air passing through narrow apertures inside the instrument. Similar considerations apply to other classes of instruments, as well as to voice sounds, and even to non-musical sounds.

In the remainder of this section we discuss the modeling of the deterministic sound signal and introduce the main concepts of additive synthesis. Later on, in section 2.4.4 we will address the problem of including the stochastic component into the additive model.

2.4.2.1 Deterministic signal component

The term *deterministic* signal means in general any signal that is not noise. The class of deterministic signals that we consider here is restricted to sums of sinusoidal components with varying amplitude and frequency. Amplitude and frequency variations can be noticed e.g. in sound attacks: some partials that are relevant in the attack can disappear in the stationary part. In general, the frequencies can have arbitrary distributions: for quasi-periodic sounds the frequencies are approximately harmonic components (integer multiples of a common fundamental frequency), while for non-harmonic sounds (such as that of a bell) they have non-integer ratios.

The deterministic part of a discrete-time sound signal can be represented by the sinusoidal model of sect. 2.4.1. The equation is

$$s_s[n] = \sum_k a_k[n] \cos(\phi_k[n]) \quad \phi_k[n] = 2\pi f_k[n]T_s + \phi_k[n-1] . \quad (2.23)$$

This equation has a great generality and can be used to faithfully reproduce many types of sound, especially in a “synthesis-by-analysis” framework (that we discuss in section 2.4.3 below). However, as already noted, it discards completely the noisy components that are always present in real signals. Another drawback of equation (2.23) is that it needs an extremely large number of control parameters: for each note that we want to reproduce, we need to provide the amplitude and frequency envelopes for all the partials. Moreover, the envelopes for a single note are not fixed, but depend in general on the intensity.

On the other hand, additive synthesis provides a very intuitive sound representation, and this is one of the reasons why it has been one of the earliest popular synthesis techniques in computer music.¹ Moreover, sound transformations performed on the parameters of the additive representation (e.g., time-scale modifications) are perceptually very robust.

¹Some composers have even used additive synthesis as a compositional metaphor, in which sound spectra are reinterpreted as harmonic structures.

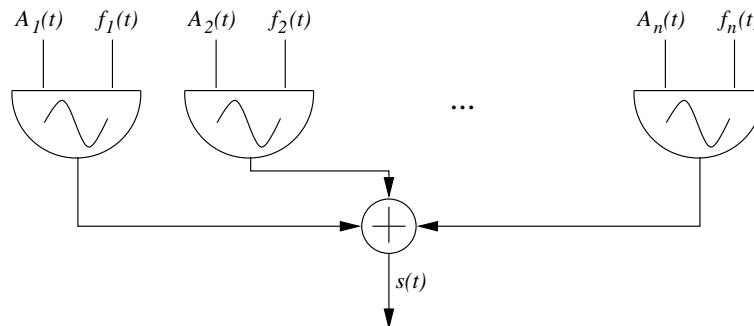


Figure 2.9: Sum of sinusoidal oscillators with time-varying amplitudes and frequencies.

2.4.2.2 Time- and frequency-domain implementations

Additive synthesis with equation (2.23) can be implemented either in the time domain or in the frequency domain. The more traditional time-domain implementation uses the digital sinusoidal oscillator in wavetable or recursive form, as discussed in section 2.2.1.1. The instantaneous amplitude and the instantaneous radian frequency of a particular partial are obtained by linear interpolation, as discussed previously. Figure 2.9 provides a block diagram of such a time-domain implementation.

M-2.12

Use the sinusoidal oscillator realized in M-2.3 to synthesize a sum of two sinusoids.

M-2.12 Solution

```

%%% headers %%%
%[...]

%%% define controls %%%
a=envgen([0, .5, 5, 10, 15, 19.5, 20], [0, 1, 1, 1, 1, 1, 0]); %fade in/out
f1=envgen([0, 20], [200, 200]); %constant freq. envelope
f2=envgen([0, 1, 5, 10, 15, 20], ... %increasing freq. envelope
          [200, 200, 205, 220, 270, 300]);

%%% compute sound %%%
s=sinosc(0, a, f1, 0)+sinosc(0, a, f2, 0);

```

The sinusoidal oscillator controlled in frequency and amplitude is the fundamental building block for time-domain implementations of additive synthesis. Here we employ it to look at the beating phenomenon. We use two oscillators, of which one has constant frequency while the second is given a slowly increasing frequency envelope. Figure 2.10 shows the f_1 , f_2 control signals and the amplitude envelope of the resulting sound signal: note the beating effect.

In alternative to the time-domain approach, a very efficient implementation of additive synthesis can be developed in the frequency domain, using the inverse FFT. Consider a sinusoid in the time-domain: its STFT is obtained by first multiplying it for a time window $w[n]$ and then performing the Fourier transform. Therefore the transform of the windowed sinusoid is the transform of the window, centered on the frequency of the sinusoid, and multiplied by a complex number whose magnitude and phase are the magnitude and phase of the sine wave:

$$s[n] = A \cos(2\pi f_0 n / F_s + \phi) \quad \Rightarrow \quad \mathcal{F}[w \cdot s](f) = A e^{i\phi} W(f - f_0). \quad (2.24)$$

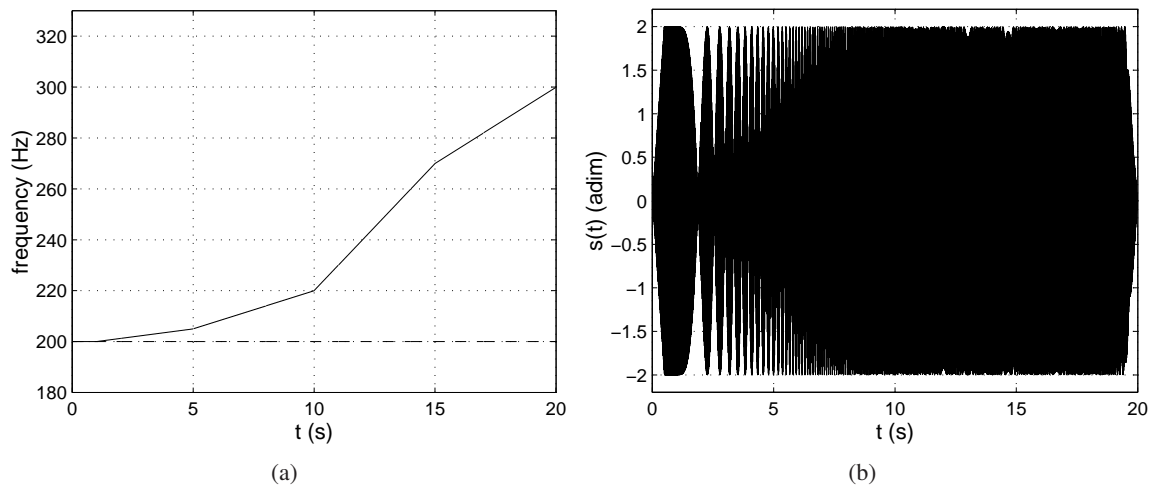


Figure 2.10: Beating effect: (a) frequency envelopes (f_1 dashed line, f_2 solid line) and (b) envelope of the resulting signal.

If the window $W(f)$ has a sufficiently high sidelobe attenuation, the sinusoid can be generated in the spectral domain by calculating the samples in the main lobe of the window transform, with the appropriate magnitude, frequency and phase values. One can then synthesize as many sinusoids as desired, by adding a corresponding number of main lobes in the Fourier domain and performing an IFFT to obtain the resulting time-domain signal in a frame.

By an overlap-and-add process one then obtains the time-varying characteristics of the sound. Note however that, in order for the signal reconstruction to be free of artifacts, the overlap-and-add procedure must be carried out using a window with the property that its shifted copies overlap and add to give a constant. A particularly simple and effective window that satisfies this property is the triangular window.

The FFT-based approach can be convenient with respect to time-domain techniques when a very high number of sinusoidal components must be reproduced: the reason is that the computational costs of this implementation are largely dominated by the cost of the FFT, which does not depend on the number of components. On the other hand, this approach is less flexible than the traditional oscillator bank implementation, especially for the instantaneous control of frequency and magnitude. Note also that the instantaneous phases are not preserved using this method. A final remark concerns the FFT size: in general one wants to have a high frame rate, so that frequencies and magnitudes need not to be interpolated inside a frame. At the same time, large FFT sizes are desirable in order to achieve good frequency resolution and separation of the sinusoidal components. As in every short-time based processes, one has to find a trade-off between time and frequency resolution.

2.4.3 Synthesis by analysis

As already remarked, additive synthesis allows high quality sound reproduction if the amplitude and frequency control envelopes are extracted from Fourier analysis of real sounds. Figure 2.11 shows the result of this kind of analysis, in the case of a saxophone tone. Using these data, additive resynthesis is straightforward.

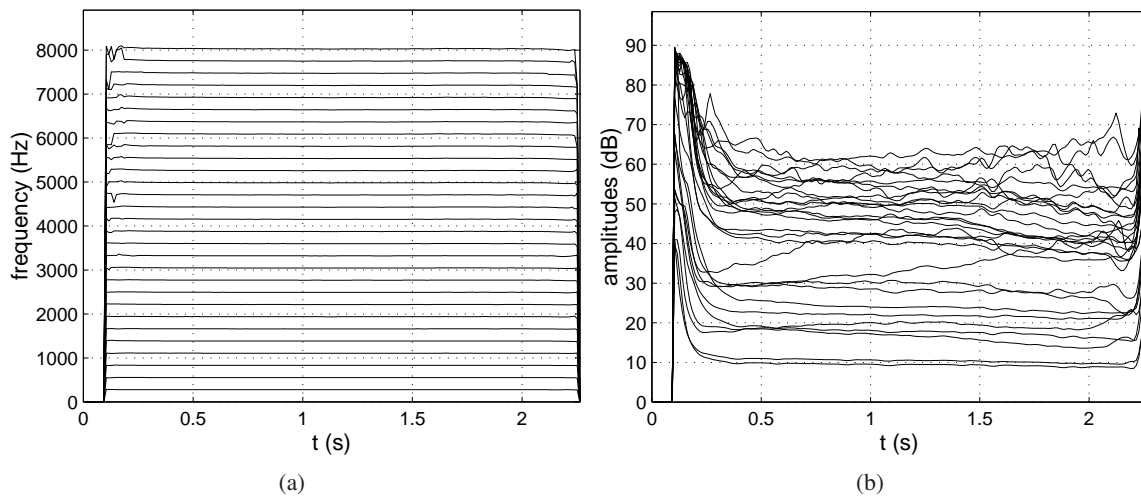


Figure 2.11: Fourier analysis of a saxophone tone: (a) frequency envelopes and (b) amplitude envelopes of the sinusoidal partials, as functions of time.

M-2.13

Assume that the script `sinan` imports two matrices `sinan_freqs` and `sinan_amps` with the partial frequency and amplitude envelopes of an analyzed sound. Resynthesize the sound.

M-2.13 Solution

```

%%% headers %%%
% [...]

%%% define controls %%%
readsan;          %import analysis matrices sinan_freqs and sinan_amps
npart=size(sinan_amps,1); %number of analyzed partials

%%% compute sound %%%
s=sinosc(...      %generate first partial
    0.5,sinan_amps(1,:),sinan_freqs(1,:),0);
for (i=2:npart) %generate higher partials and sum
    s=s+sinosc(0.5,sinan_amps(i,:),sinan_freqs(i,:),0);
end

```

2.4.3.1 Magnitude and Phase Spectra Computation

The first step of any analysis procedure that tracks frequencies and amplitudes of the sinusoidal components is the frame-by-frame computation of the sound magnitude and phase spectra. This is carried out through short-time Fourier transform. The subsequent tracking procedure will be performed in this spectral domain. The control parameters for the STFT are the window-type and size, the FFT-size, and the frame-rate. These must be set depending on the sound to be processed.

Note that the analysis step is completely independent from the synthesis, therefore the observations made in section 2.4.2.2 about FFT-based implementations (the window must overlap and add to

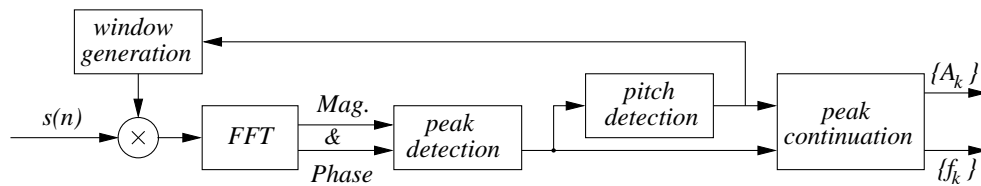


Figure 2.12: Block diagram of the sinusoid tracking process, where $s[n]$ is the analyzed sound signal and A_k , f_k are the estimated amplitude and frequency of the k th partial in the current analysis frame.

a constant) do not apply here. Good resolution of the spectrum is needed in order to correctly resolve, identify, and track the peaks which correspond to the deterministic component.

If the analyzed sound is almost stationary, long windows (i.e. windows that cover several periods) that have good side-lobe rejection can be used, with a consequent good frequency resolution. Unfortunately most interesting sounds are not stationary and a compromise is required. For harmonic sounds one can scale the actual window size as a function of pitch, thus achieving a constant time-frequency trade-off. For inharmonic sounds the size should be set according to the minimum frequency difference that exists between partials.

The question is now how to perform automatic detection and tracking of the spectral peaks that correspond to sinusoidal components. In section 2.4.3.2 below we present the main guidelines of a general analysis framework, which is summarized in figure 2.12. First, the FFT of a sound frame is computed according to the above discussion. Next, the prominent spectral peaks are detected and incorporated into partial trajectories. If the sound is pseudo-harmonic, a pitch detection step can improve the analysis by providing information about the fundamental frequency information, and can also be used to choose the size of the analysis window.

Such a scheme is only one of the possible approaches that can be used to attack the problem. Hidden Markov Models (HMMs) are another one: a HMM can optimize groups of peaks trajectories according to given criteria, such as frequency continuity. This type of approach might be very valuable for tracking partials in polyphonic sounds and complex inharmonic tones.

2.4.3.2 A sinusoid tracking procedure

We now discuss in more detail the analysis steps depicted in figure 2.12.

Peak detection. The first one is detection of the most prominent frequency peaks (i.e., local maxima in the magnitude spectrum) in the current analysis frame. Real sounds are not periodic, do not have clearly spaced and defined spectral peaks, exhibit interactions between components. Therefore, the best one can do at this point is to detect as many peaks as possible and postpone to later analysis steps the decision of which ones actually correspond to sinusoidal components. The peaks are then searched by only imposing two minimal constraints: they have to lie within a given frequency range, and above a given magnitude threshold. The detection of very soft peaks is hard: they have little resolution, and measurements are very sensitive to transformations because as soon as modifications are applied to the analysis data, parts of the sound that could not be heard in the original can become audible. Having a very clean sound with the maximum dynamic range, the magnitude threshold can be set to the amplitude of the background noise floor. In order to gain better resolution in the high frequency range, the sound may be pre-processed to introduce preemphasis, which has then to be compensated later on before the

resynthesis.

Pitch detection. After peak detection, many procedures can be used to decide whether a peak belongs to a sinusoidal partial or not. One possible strategy is to measure how close the peak shape is to the ideal sinusoidal peak (recall what we said about the transform of a windowed sinusoid and in particular equation (2.24)). A second valuable source of additional information is pitch. If a fundamental frequency is actually present, it can be exploited in two ways. First, it helps the tracking of partials. Second, the size of the analysis window can be set according to the estimated pitch in order to keep the number of periods-per-frame constant, therefore achieving the best possible time-frequency trade-off (this is an example of a *pitch-synchronous* analysis). There are many possible pitch detection strategies, which will be presented in the next chapter.

Peak continuation. A third and fundamental strategy for peak selection is to implement some sort of peak *continuation* algorithm. The basic idea is that a set of “guides” advance in time and follow appropriate frequency peaks (according to specified constraints that we discuss in the next paragraph) forming trajectories out of them. A guide is therefore an abstract entity which is used by the algorithm to create the trajectories, and the trajectories are the actual result of the peak continuation process. The guides are turned on, advanced, and finally turned off during the continuation algorithm, and their instantaneous state (frequency and magnitude) is continuously updated during the process. If the analyzed sound is harmonic and a fundamental has been estimated, then the guides are created at the beginning of the analysis, with frequencies set according to the estimated harmonic series. When no harmonic structure can be estimated, each guide is created when the first available peak is found. In the successive analysis frames, the guides modify their status depending on the last peak values. This past information is particularly relevant when the sound is not harmonic, or when the harmonics are not locked to each other and we cannot rely on the fundamental as a strong reference for all the harmonics.

The main constraints used to assign guides to spectral peaks are as follows. A peak is assigned to the guide that is closest to it and that is within an assigned frequency deviation. If a guide does not find a match, the corresponding trajectory can be turned off, and if a continuation peak is not found for a given amount of time the guide is killed. New guides and trajectories can be created starting from peaks of the current frame that have high magnitude and are not “claimed” by any of the existing trajectories. After a certain number of analysis frames, the algorithm can look at the trajectories created so far and adopt corrections: in particular, short trajectories can be deleted, and small gaps in longer trajectories can be filled by interpolating between the values of the gap edges.

One final refinement to this process can be added by noting that the sound attack is usually highly non-stationary and noisy, and the peak search is consequently difficult in this part. Therefore it is customary to perform the whole procedure backwards in time, starting from the end of the sound (which is usually a more stable part). When the attack is reached, a lot of relevant information has already been gained and non-relevant peaks can be evaluated and/or rejected.

2.4.4 “Sines-plus-noise” models

At the beginning of our discussion on additive modeling, we remarked that the spectral energy of the sound signals has a *deterministic* component that is concentrated on a discrete set of frequencies, and a *stochastic* component that has a broadband characteristics. So far we have discussed the problem

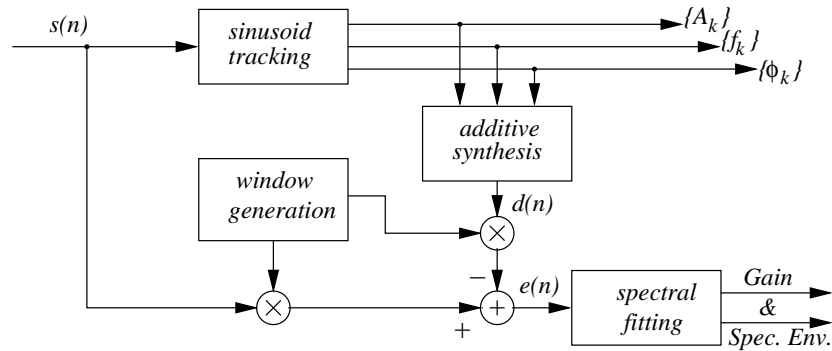


Figure 2.13: Block diagram of the stochastic analysis and modeling process, where $s[n]$ is the analyzed sound signal and A_k , f_k , ϕ_k are the estimated amplitude, frequency, and phase of the k th partial in the current analysis frame.

of modeling the deterministic –or sinusoidal– component. Now we have to include the stochastic component into the model.

A sinusoidal representation may in principle be used also to simulate noise, since noise consists of sinusoids at every frequency within the band limits. It is clear however that such a representation would be computationally very demanding. Moreover it would not be a *flexible* sound representation. Therefore the most convenient sound model is of the form

$$s[n] = s_s[n] + e[n] \quad (2.25)$$

where $s_s[n]$ represent the deterministic part (eq. 2.23 and $e[n]$ represents the stochastic component and is modeled separately from the deterministic part.

2.4.4.1 Stochastic analysis

The most straightforward approach to estimation of the stochastic component is through *subtraction* of the deterministic component from the original signal. Subtraction can be performed either in the time domain or in the frequency domain. Time domain subtraction must be done while preserving the phases of the original sound, and instantaneous phase preservation can be computationally very expensive. On the other hand, frequency-domain subtraction does not require phase preservation. However, time-domain subtraction provides much better results, and is usually favored despite the higher computational costs. For this reason we choose to examine time-domain subtraction in the remainder of this section. Figure 2.13 provides a block diagram.

Suppose that the deterministic component has been estimated in a given analysis frame, using for instance the general scheme described in section 2.4.3 (note however that in this case the analysis should be improved in order to provide estimates of the instantaneous phases as well). Then the first step in the subtraction process is the time-domain resynthesis of the deterministic component with the estimated parameters. This should be done by properly interpolating amplitude, frequency, and phase values in order to avoid artifacts in the resynthesized signal. The actual subtraction can be performed as

$$e[n] = w[n] \cdot [s[n] - d[n]], \quad n = 0, 1, \dots, N - 1, \quad (2.26)$$

where $s[n]$ is the original sound signal and $d[n]$ is the re-synthesized deterministic part. The difference $(s - d)$ is multiplied by an analysis window w of size N , which deserves some discussion.

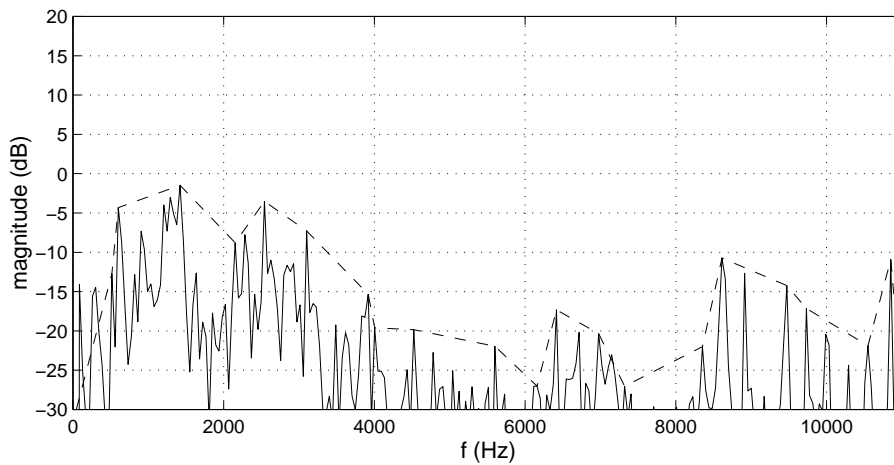


Figure 2.14: Example of residual magnitude spectrum (solid line) and its line-segment approximation (dashed line), in an analysis frame. The analyzed sound signal is the same saxophone tone used in figure 2.11.

We have seen in 2.4.3 that high frequency resolution is needed for the deterministic part, and for this reason long analysis windows are used for its estimation. On the other hand, good time resolution is more important for the stochastic part of the signal, especially in sound attacks, while frequency resolution is not a major issue for noise analysis. A way to obtain good resolutions for both the components is to use two different analysis windows. Therefore w in equation (2.26) is not in general the same window used to estimate $d[n]$, and the size N is in general small.

Once the subtraction has been performed, there is one more step than can be used to improve the analysis, namely, test can be performed on the estimated residual in order to assess how good the analysis was. If the spectrum of the residual still contains some partials, then the analysis of the deterministic component has not been performed accurately and the sound should be re-analyzed until the residual is free of deterministic components. Ideally the residual should be as close as possible to a stochastic signal, therefore one possible test is a measure of correlation of the residual samples.²

2.4.4.2 Stochastic modeling

The assumption that the residual is a stochastic signal implies that it is fully described by its amplitude and its spectral envelope characteristics. Information on the instantaneous phase is not necessary. Based on these considerations, a frame of the stochastic residual can be completely characterized by a filter that models the amplitude and general frequency characteristics of the residual. The representation of the residual for the overall sound will then be a time-varying filter.

Within a given frame we therefore assume that $e(t)$ can be modeled as

$$E(f) = H(t) \cdot U(f), \quad (2.27)$$

where U is white noise and H is the frequency response of filter whose coefficients vary on a frame-by-frame basis. The stochastic modeling step is summarized in the last block of figure 2.13.

² Note that if the analyzed sound has not been recorded in silent and anechoic settings the residual will contain not only the stochastic part of the sound, but also reverberation and/or background noise.

The filter design problem can be solved using different strategies. One approach that is often adopted uses some sort of curve fitting (line-segment approximation, spline interpolation, least squares approximation, and so on) of the magnitude spectrum of e in an analysis frame. As an example, line-segment approximation can be obtained by stepping through the magnitude spectrum, finding local maxima at each step, and connecting the maxima with straight lines. This procedure can approximate the spectral envelope with reasonable accuracy, depending on the number of points, which in turn can be set depending on the sound complexity. See figure 2.14 for an example.

Another possible approach to the filter design problem is Linear Prediction (LP) analysis, which is a popular technique in speech processing. In this context, however, curve fitting on the noise spectrum (e.g., line-segment approximation) is usually considered to be a more flexible approach and is preferred to LP analysis. We will return on Linear Prediction techniques in section 2.5.3.

The next question is how to implement the estimated time-varying filter in the resynthesis step.

2.4.4.3 Resynthesis and modifications

Figure 2.15 shows the block diagram of the synthesis process. The deterministic signal, i.e., the sinusoidal component, results from the magnitude and frequency trajectories, or their transformation, by generating a sine wave for each trajectory (additive synthesis). As we have seen, this can either be implemented in the time domain with the traditional oscillator bank method or in the frequency domain using the inverse-FFT approach.

Concerning the stochastic component, a frequency-domain implementation is usually preferred to a direct implementation of the time-domain convolution (2.27), due to its computational efficiency³ and flexibility. In each frame, the stochastic signal is generated by an inverse-FFT of the spectral envelopes. Similarly to what we have seen for the deterministic synthesis in section 2.4.2.2, the time-varying characteristics of the stochastic signal is then obtained using an overlap-and-add process.

In order to perform the IFFT, a magnitude and a phase responses have to be generated starting from the estimated frequency envelope. Generation of the magnitude spectrum is straightforwardly obtained by first linearly interpolating the spectral envelope to a curve with half the length of the FFT-size, and then multiplying it by a gain that corresponds to the average magnitude extracted in the analysis. The estimated spectral envelope gives no information on the phase response. However, since the phase response of noise is noise, a phase response can be created from scratch using a random signal generator. In order to avoid periodicities at the frame rate, new random values should be generated at every frame.

The sines-plus-noise representation is well suited for modification purposes.

- By only working on the deterministic representation and modifying the amplitude-frequency pairs or the original sound partials, many kinds of frequency and magnitude transformations can be obtained. As an example, partials can be transposed in frequency. It is also possible to decouple the sinusoidal frequencies from their amplitude, obtaining pitch-shift effects that preserve the formant structure.
- Time-stretching transformations can be obtained by resampling the analysis points in time, thus slowing down or speeding up the sound while maintaining pitch and formant structure. Given the stochastic model that we are using, the noise remains noise and faithful signal resynthesis is possible even with extreme stretching parameters.

³ In fact, by using a frequency-domain implementation for both the deterministic and the stochastic synthesis one can add the two spectra and resynthesize both the components at the cost of a single IFFT per frame.

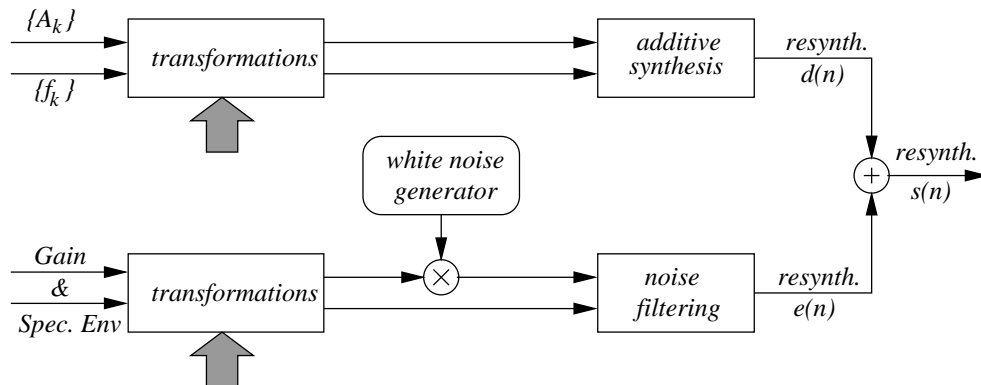


Figure 2.15: Block diagram of the sines-plus-noise synthesis process.

- By acting on the relative amplitude of the two components, interesting effects can be obtained in which either the deterministic or the stochastic parts are emphasized. As an example, the amount of “breathiness” of a voiced sound or a wind instrument tone can be adjusted in this way. One must keep in mind however that, when different transformations are applied to the two representations, the deterministic and stochastic components in the resulting signal may not be perceived as a single sound event anymore.
- Sound morphing (or *cross-synthesis* transformations can be obtained by interpolating data from two or more analysis files. This transformations are particularly effective in the case of quasi-harmonic sounds with smooth parameter curves.

2.4.5 Sinusoidal description of transients

So far we have seen how to extend the sinusoidal model by using a “sines-plus-noise” approach that explicitly describes the residual as slowly varying filtered white noise. Although this technique is very powerful, transients do not fit well into a filtered noise description, because they lose sharpness and are smeared. This consideration motivates us to handle transients separately.

One straightforward approach, that is sometimes used, is removing transient regions from the residual, performing the sines-plus-noise analysis, and adding the transients back into the signal. This approach obviously requires memory where the sampled transients must be stored, but since the transient residuals remain largely invariant throughout most of the range of an instrument, only a few residuals are needed in order to cover all the sounds of a single instrument. Although this approach works well, it is not flexible because there is no model for the transients. In addition, identifying transients as everything that is neither sinusoidal nor transient is not entirely correct. Therefore we look for a suitable transient model, that can be embedded in the additive description to obtain a “sines-plus-transients-plus-noise” representation.

2.4.5.1 The DCT domain

In the following we adopt a further modified version of the additive sound representation (2.23), in which the sound transients are explicitly modeled by an additional signal:

$$s[n] = s_s[n] + e_t[n] + e_r[n], \quad (2.28)$$

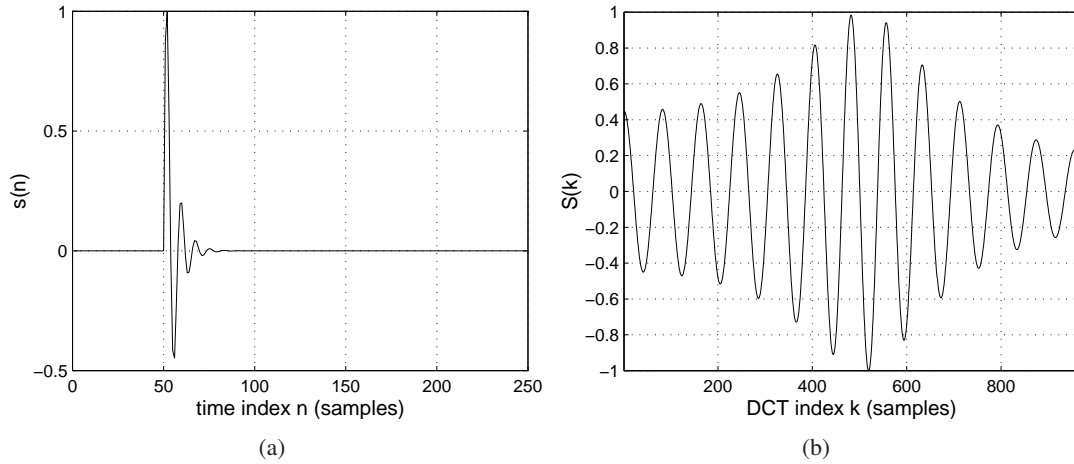


Figure 2.16: Example of DCT mapping: (a) an impulsive transient (an exponentially decaying sinusoid) and (b) its DCT as a slowly varying sinusoid.

where $s_s[n]$ is the sinusoidal component, $e_t[n]$ is the signal associated to transients and $e_r[n]$ is the noisy residual. The transient model is based on a main underlying idea: we have seen that a slowly varying sinusoidal signal is impulsive in the frequency domain, and sinusoidal models perform short-time Fourier analysis in order to track slowly varying spectral peaks (the tips of the impulsive signals) over time. Transients are very much dual to sinusoidal components: they are impulsive in the time domain, and consequently they must be oscillatory in the frequency domain. Therefore, although transient cannot be tracked by a short-time analysis (because their STFT will not contain meaningful peaks), we can track them by performing sinusoidal modeling in a properly chosen frequency domain. The mapping that we choose to use is the one provided by the discrete cosine transform (DCT):

$$S[k] = \beta[k] \sum_{n=0}^{N-1} s[n] \cos \left[\frac{(2n-1)k\pi}{2N} \right], \quad \text{for } n, k = 0, 1, \dots, N-1, \quad (2.29)$$

where $\beta[1] = \sqrt{1/N}$ and $\beta[k] = \sqrt{2/N}$ otherwise. From equation (2.29) one can see that an ideal impulse $\delta[n - n_0]$ (i.e., a Kronecker delta function centered in n_0) is transformed into a cosine whose frequency increases with n_0 . Figure 2.16(a) shows a more realistic transient signal, a one-sided exponentially decaying sine wave. Figure 2.16(b) shows the DCT of the transient signal: a slowly varying sinusoid. This considerations suggest that the time-frequency duality can be exploited to develop a transient model: the same kind of parameters that characterize the sinusoidal components of a signal can also characterize the transient components of a signal, although in a different domain.

2.4.5.2 Transient analysis and modeling

Having transformed the transient into the DCT domain, the most natural way to proceed is performing sinusoidal modeling in this domain: STFT analysis of the DCT-domain signal can be used to find meaningful peaks, and then the signal can be resynthesized in the DCT domain and back-transformed to the time domain with an inverse DCT transform (IDCT). This process is shown in figure 2.17. We now discuss the main steps involved in this block diagram.

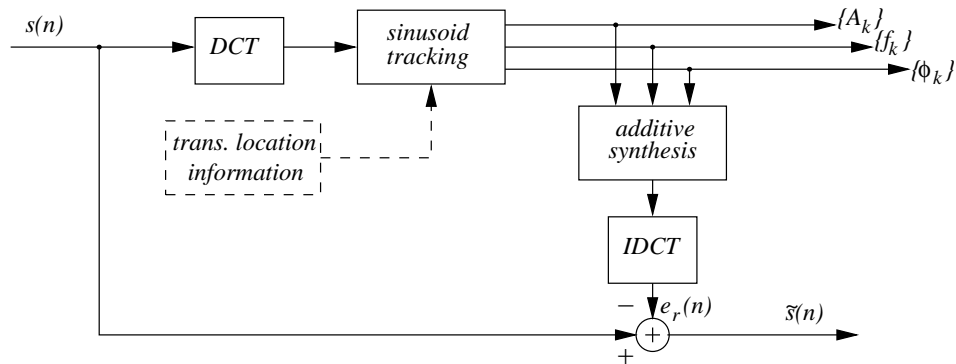


Figure 2.17: Block diagram of the transient analysis and modeling process, where $s[n]$ is the analyzed sound signal and A_k , f_k , ϕ_k are the estimated amplitude, frequency, and phase of the k th DCT-transformed transient in the current analysis frame.

First the input signal $s[n]$ is divided into non-overlapping blocks in which DCT analysis will be performed. The block length should be chosen so that a transient appears as “short”, therefore large block sizes (e.g., 1 s) are usually chosen. The block DCT is followed by a sinusoidal analysis/modeling process which is identical to what we have seen in section 2.4.3. The analysis can optionally embed some information about transient location within the block: there are many possible transient detection strategies, which we do not want to discuss here. Also, the analysis can perform better if the sinusoid tracking procedure starts from the end of the DCT-domain signal and moves backwards toward the beginning, because the beginning of a DCT frame is usually spectrally rich and this can deteriorate the performance of the analysis (similar considerations were done in section 2.4.3 when discussing sinusoid tracking in the time domain).

The analysis yields parameters that correspond to slowly varying sinusoids in the DCT domain: each transient is associated to a triplet $\{A_k, f_k, \phi_k\}$, amplitude, frequency, and phase of the k th “partial” in each STFT analysis frame within a DCT block. By recalling the properties of the DCT one can see that f_k correspond to onset locations, A_k is the amplitude of the time-domain signal also, and ϕ_k is related to the time direction (positive or negative) in which the transient evolves. Resynthesis of the transients is then performed using these parameters to reconstruct the sinusoids in the DCT domain. Finally an inverse discrete cosine transform (IDCT) on each of the reconstructed signals is used to obtain the transients in each time-domain block, and the blocks are concatenated to obtain the transients for the entire signal.

It is relatively straightforward to implement a “fast transient reconstruction” algorithm. Without entering the details, we just note that the whole procedure can be reformulated using FFT transformations only: in fact one could verify that the DCT can be implemented using an FFT block plus some post-processing (multiplication of the real and imaginary parts of the FFT by appropriate cosinusoidal and sinusoidal signals followed by a sum of the two parts). Furthermore, this kind of approach naturally leads to a FFT-based implementation of the additive synthesis step (see section 2.4.2.2).

One nice property of this transient modeling approach is that it fits well within the sines-plus-noise analysis that we have examined in the previous sections. The processing block depicted in figure 2.17 returns an output signal $\tilde{s}[n]$ in which the transient components e_t have been removed by subtraction: this signal can be used as the input to the sines-plus-noise analysis, in which the remaining components (deterministic and stochastic) will be analyzed and modeled. From the implementation

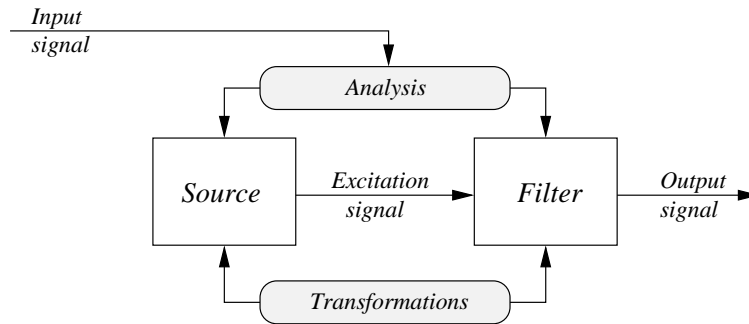


Figure 2.18: Source-filter model.

viewpoint, one advantage is that the core components of the transient-modeling algorithm (sinusoid tracking and additive resynthesis) are identical to those used for the deterministic model. Therefore the same processing blocks can be used in the two stages, although working on different domains.

2.5 Source-filter models

2.5.1 Source-filter decompositions

Some sound signals can be effectively modeled through a feed-forward source-filter structure, in which the source is in general a spectrally rich excitation signal, and the filter is a linear system that acts as a resonator and shapes the spectrum of the excitation. A typical example is the voice, where the periodic pulses or random fluctuations produced by the vocal folds are filtered by the vocal tract, that shapes the spectral envelope. The vowel quality and the voice color greatly depends on the resonance regions of the filter, usually called *formants*.

Source-filter models are typically used in an analysis-synthesis framework, in which both the source signal and the filter parameters are estimated from a target sound signal, that can be subsequently resynthesized through the identified model. Moreover, transformations can be applied to the filter and/or the excitation before the reconstruction (see Fig. 2.18). One of the most common analysis techniques is Linear Prediction, that we will address in Sec. 2.5.3.

In computer music, source-filter models are traditionally grouped under the label *subtractive synthesis*. A number of analog voltage controlled synthesizers in the '70s, e.g. Moog, made use of subtractive synthesis techniques in which audio filters are applied to spectrally rich waveforms (square waves, sawtooth waves, triangle waves, and so on).

We will assume the filter block to be linear and time-invariant (at least on a short-time scale), so that the excitation signal $x[n]$ and the output signal $s[n]$ are related through the difference equation

$$s[n] = \sum_i b_i u[n - i] - \sum_k a_k s[n - k], \quad (2.30)$$

or, in the Z-domain,

$$S(z) = H(z)X(z), \quad \text{with } H(z) = \frac{B(z)}{A(z)} = \frac{\sum_i b_i z^{-i}}{\sum_k a_k z^{-k}}. \quad (2.31)$$

Equation (2.31) shows how the features of the source and of the filter are combined: the spectral fine structure of the excitation signal (spectral lines, broadband or narrowband noise, etc.) is multi-

plied by the spectral envelope of the filter, which therefore has a *shaping* effect on the source spectrum. Therefore, it is possible to control and modify separately different features of the signal: as an example, the pitch of a speech sound depends on the excitation and can be controlled separately from the formant structure, which instead depends on the filter. When the filter coefficients are (slowly) varied over time, the frequency response changes. As a consequence, the output will be a combination of temporal variations of the input and of the filter (*cross-synthesis*).

2.5.1.1 Sources and filters

One of the simplest waveforms is the impulse train generator, which produces a sequence of unit impulses spaced by the desired fundamental period. Another simple generator for stochastic sources is the random noise generator, which produces a flat spectrum noise. We will be making use of these generators in Sec. 2.5.2.

M-2.14

Write a function `noisegen(t0, amp)` that realizes the noise generator, and a function `buzz(t0, amp, f)` that realizes the impulse generator. The parameters `(t0, amp, f)` are initial time, amplitude envelope, and frequency envelope, respectively.

M-2.14 Solution

Hint for the noise generator: simply use the function `rand()`.

Hint for the impulse generator: use additive synthesis, and sum up all the harmonic components $\cos(2k\pi f_0 t)$ ($k \in \mathbb{N}$) of the fundamental frequency f_0 , up to the Nyquist frequency $F_s/2$.

An important filter is the second order IIR filter that is very often used for modeling resonances. It is described by:

$$H(z) = \frac{b_0}{1 + a_1 z^{-1} + a_2 z^{-2}} \quad (2.32)$$

The filter coefficient can be derived by these (approximate) equations:

$$a_1 = -2r \cos(\omega_c), \quad a_2 = r^2, \quad b_0 = (1 - r) \sqrt{1 - 2r \cos(2\omega_c) + r^2}, \quad (2.33)$$

where the auxiliary parameters r and ω_c are defined as $r = e^{-\frac{\pi B}{F_s}}$ e $\omega_c = 2\pi f_c / F_s$ and correspond to the magnitude and phase of the complex conjugate poles of the transfer function. The parameter f_c is the center frequency of the resonant filter, and B is the bandwidth. The gain normalization factor b_0 is computed by normalizing the magnitude at the resonance frequency, i.e. $|H(\omega_c)| = 1$. The impulse response is a sinusoidal function with exponentially decreasing amplitude

$$h[n] = \frac{b_0}{\sin \omega_c} r^n \sin(\omega_c(n + 1)). \quad (2.34)$$

Then the 40-dB time constant τ_{40dB} (i.e. the time required for the oscillation to decay by 40 dB) is given by $\tau_{40dB} = \log(10^{-2}) / F_s \log(r) = -2 \log(r) / F_s$.

M-2.15

Write a function `baIIR2(fc, B)` that computes the coefficients of N 2nd order resonant filters, given the vectors of length N `fc` (center frequencies) and `B` (bandwidths).

M-2.15 Solution

```

function [b,a]=baIIR2(fc,B); %computes coeff. a,b of II order cells

global Fs; global Fc;
nfilters=length(fc); %no. of cells to be computed

r=exp(-(pi.*B)/Fs);
a0=ones(nfilters,1);
a1=-(2*r.*cos(2*pi*fc/Fs))';
a2=r'.^2;

b0=(1-r).*sqrt(1-2.*r.*cos(2*2*pi.*fc/Fs)+r.*r);
a=[a0 a1 a2];
b=[b0' zeros(nfilters,1) zeros(nfilters,1)];

```

Note that we have followed the Octave/Matlab convention in defining the coefficients b, a . See the help for the function `filter(b,a,in)`

2.5.1.2 An example: modal synthesis of percussive sounds

A set of N second order resonant filters R_i ($i = 1 \dots N$) of the form (2.32) can be grouped into a filterbank, where the same excitation signal x is fed to all the R_i 's, as depicted in Fig. 2.19. This specific source-filter structure is well suited to simulate percussive sounds.

In this case the excitation signal has an impulsive characteristics and represents a “striker” that a hammer or a mallet impart to a resonating object. Suitable “striker” excitation signals are e.g. a square impulse or a noise burst. The filter block represents the resonating object hit by the hammer: the center frequencies f_{ci} ($i = 1 \dots N$) of the resonant filters can be chosen to match a desired spectrum. As an example, a string will have a harmonic spectrum in which partials are regularly spaced on the frequency axis, therefore $f_{ci} = i f_{c1}$ ($i = 2 \dots N$), where f_{c1} acts as the fundamental frequency. On the other hand, the partial distribution of a bar, or a bell, or the circular membrane of a drum, will be inharmonic. As an example, it is known that the partials of an ideal bar clamped at one end are approximately $f_{c2} \sim (2.998/1.994)^2 f_{c1}$, $f_{ci} \sim [(2i + 1)/1.994]^2 f_{c1}$ ($i = 3 \dots N$).

The bandwidths B_i of the R_i 's determine the decay characteristics of each partial, according to Eq. (2.34). A first possible choice is setting the same B (i.e. the same parameter r) for every filter. An alternative choice, that better describes the behavior of such resonant objects as strings, bars, and so on, amounts to setting the same *quality factor* $Q = B_i/f_{ci}$ for all the filters.

The structure depicted in Fig. 2.19 is also an example of *modal synthesis*. We will return on modal synthesis in Chapter *Sound modeling: source based approaches*, and will provide more physically sound foundations to this sound synthesis technique.

M-2.16

Write a function `modalosc(t0,tau,amp,fc,B)` that realizes the structure of Fig. 2.19. The parameter t_0 is the total duration of the sound signal, τ, amp define duration and max. amplitude of the striker signal (e.g. a noise burst), and the vectors \mathbf{f}_c, \mathbf{B} are defined as in Example M-2.15.

2.5.2 Speech modeling

2.5.2.1 Speech production mechanism and models

Speech is an acoustic sound pressure wave created when air is expelled from the lungs through the trachea and vocal tract (fig. fig. 2.20). This tract is composed of vocal fold opening (glottis), throat,



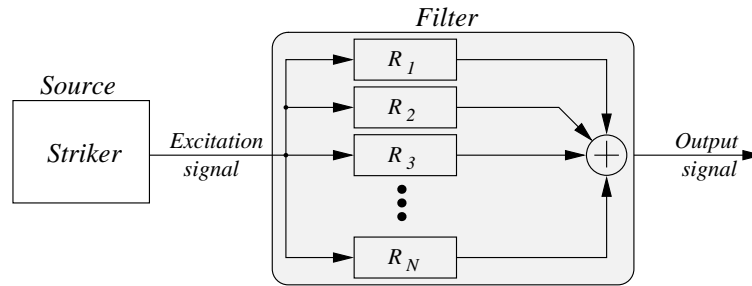


Figure 2.19: Parallel structure of digital resonators for the simulation of struck objects – the R_i 's have transfer functions of the form (2.32).

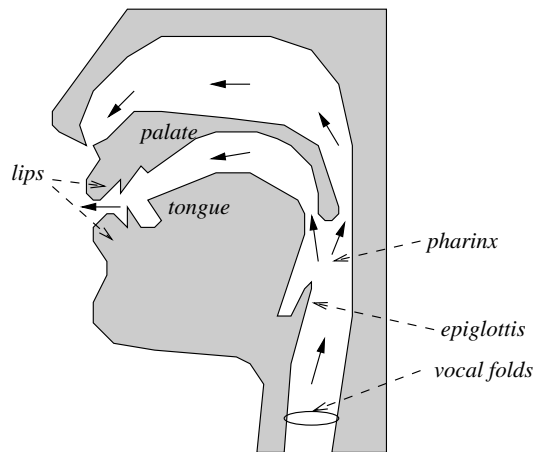


Figure 2.20: A schematic view of the phonatory system. Solid arrows indicates the direction of the airflow generated by lung pressure.

nose, mouth and lips. As the acoustic wave passes through the vocal tract, its frequency content (spectrum) is altered by the resonances of the vocal tract; vocal tract resonances are called formants. Two types of sounds characterize speech, namely voiced and unvoiced). Voiced sounds result from a periodic excitation of the vocal tract causing oscillation of the vocal chords in a quasi-periodic manner. Technically, it is useful to think voiced speech as the response of an acoustic filter, modeling the vocal tract, to a passing pressure wave. The vocal chords oscillate in a very non-sinusoidal manner at a rate that is approximately 100Hz for adult males and 120Hz for adult females (the frequency varies from speaker to speaker as well). The periodic nature of the oscillations gives rise to harmonics and the period associated with the fundamental frequency is commonly termed its pitch¹. The range of potential pitch frequencies varies from 50Hz to 250Hz for men, and from 120 to 500Hz for women. Everyone has a "habitual pitch level", which is a sort of "preferred" pitch that will be used naturally on the average. Pitch is shifted up and down in speaking in response to factors relating to stress, intonation, and emotion. Intonation is associated with the pitch contour over time and performs several functions in a language, the most important being to signal grammatical structure.

Synthesized speech can be produced by several different methods. All of these have some benefits and deficiencies. The methods are usually classified into three groups.

- *Concatenative synthesis* uses different length prerecorded samples derived from natural speech. This technique is based on the methods seen in sect. 2.3. Connecting prerecorded natural utterances is probably the easiest way to produce intelligible and natural sounding synthetic speech. However, concatenative synthesizers are usually limited to one speaker and one voice and usually require more memory capacity than other methods. The concatenative approach is becoming more and more popular in speech synthesis systems.
- *Formant synthesis* models the pole frequencies of speech signal or transfer function of vocal tract based on source-filter-model. This was the most widely used synthesis method during last decades, before the development of concatenative methods. It is based on the source-filter-model of speech described in sect. 2.5. Formant synthesis also provides infinite number of sounds which makes it more flexible than for example concatenation methods. We discuss formant synthesis in Sec. 2.5.2.2 below
- *Articulatory synthesis* attempts to model the human speech production system directly. Articulatory synthesis typically involves models of the human articulators and vocal cords. The articulators are usually modeled with a set of area functions between glottis and mouth. The articulatory control parameters may be for example lip aperture, lip protrusion, tongue tip height, tongue tip position, tongue height, tongue position and velic aperture. Phonatory or excitation parameters may be glottal aperture, cord tension, and lung pressure. The models are developed with the methods which will be seen in Chapter *Sound modeling: source based approaches*. Advantages of articulatory synthesis are that the vocal tract models allow accurate modeling of transients due to abrupt area changes, whereas formant synthesis models only spectral behavior. The articulatory method is still too complicated for high quality implementations, but may arise as a potential method in the future.

2.5.2.2 Formant synthesis

Formant synthesis of speech realizes a source-filter model in which a broadband source signal undergoes multiple filtering transformations that are associated to the action of different elements of the phonatory system. Depending on whether voiced or unvoiced speech (see above) has to be simulated, two different models are used.

If $s[n]$ is a voiced speech signal, it can be expressed in the Z-domain by the following cascaded spectral factors:

$$S(z) = g_v X(z) \cdot [F(z) \cdot V(z) \cdot R(z)] \quad (2.35)$$

where the source signal $X(z)$ is in this case a periodic pulse train whose period coincides with the pitch of the signal, g_v is a constant gain term, $F(z)$ is a filter associated to the response of the vocal folds to pitch pulses, $V(z)$ is the vocal tract filter, $R(z)$ simulates the output radiation effect of the lips.

If $s[n]$ is a unvoiced speech signal, vocal folds do not vibrate and turbulences is produced by the passage of air through a narrow constriction (such as the teeth). The turbulence is traditionally modeled as white noise. In this case, the model is expressed in the Z-domain as

$$S(z) = g_u X(z) \cdot [V(z) \cdot R(z)] \quad (2.36)$$

where the source signal $X(z)$ is in this case a driving noise sequence, while the gain term g_u is in general different from the voiced configuration gain, g_v . Note that the vocal fold response $F(z)$ is not included in the model in this case.



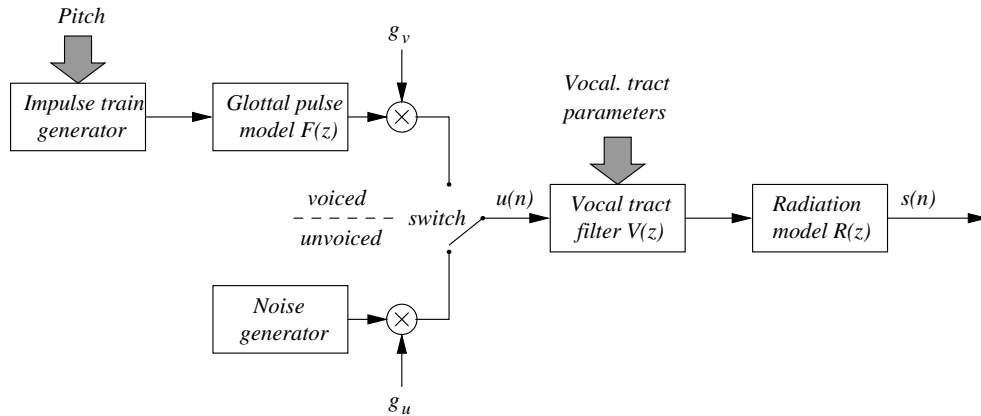


Figure 2.21: A general model for formant synthesis of speech.

Any voiced or unvoiced sound can be characterized by either Eq. (2.35) or (2.36); the complete transfer function characterizing the speech apparatus is defined as

$$H(z) = \frac{S(z)}{X(z)} \quad (2.37)$$

which may or may not include vocal fold response to pitch depending on whether the sound is voiced or unvoiced. Such a general discrete-time model for speech is shown in Fig. 2.21.

The filter $F(z)$ that shapes the glottal pulses can be modeled as

$$F(z) = \frac{1}{[1 - \exp(-c/F_s)z^{-1}]^2}.$$

The radiation filter $R(z)$ can be thought of as a load that converts the airflow signal at the lips into an outgoing pressure wave (the signal $s[n]$). Under very idealized hypothesis, $R(z)$ can be approximated by a differentiator:

$$R(z) = 1 - z^{-1}.$$

We now focus on the filter $V(z)$ that models the vocal tract formants. A single frequency formant associated to vocal tract can be modeled with a two-pole resonator, (see Sec. 2.5.1) which enables both the formant frequency (pole-pair frequency) and its bandwidth to be specified. We will denote the filter associated to the i th formant as $V_i(z)$, having center frequency f_i and bandwidth B_i (see Eq. (2.32)). At least three vocal tract formants are generally required to produce intelligible speech and up to five formants are needed to produce high quality speech.

Two basic structures, parallel and cascade, can be used in general, but for better performance some kind of combination of these is usually adopted. A cascade formant synthesizer consists of band-pass resonators connected in series, i.e. the output of each formant resonator is applied to the input of the following one. The cascade vocal tract model is given by

$$V(z) = g \prod_{i=1}^K V_i(z).$$

The cascade structure needs only formant frequencies as control information. The main advantage of this structure is that the relative formant amplitudes for vowels do not need individual controls. A

cascade model of the vocal tract is considered to provide good quality in the synthesis of vowels, but is not flexible enough for the synthesis of nasals.

A parallel formant synthesizer consists of resonators connected in parallel, i.e. the same input is applied to each formant filter. The resulting vocal tract model is then given by

$$V(z) = \sum_{i=1}^K a_i \cdot V_i(z).$$

Sometimes extra resonators for nasals are used. The excitation signal is applied to all formants simultaneously and their outputs are summed. The parallel structure enables controlling of bandwidth and gain for each formant individually and thus needs also more control information. The parallel structure has been found to be better for nasals, fricatives, and stop-consonants, but some vowels can not be modeled with parallel formant synthesizer as well as with the cascade one.

M-2.17

Using the functions `buzz` and `baIIR2`, realize a parallel formant synthesizer. Use 3 2nd order IIR cells, corresponding to the first 3 vowel formant frequencies.

M-2.17 Solution

```
global Fs; global SpF;
Fs=22050;
ControlW=0.01;           %control window (in sec): 10 ms
SpF=round(Fs*ControlW);

%%% define controls %%%
amp=envgen([0, .2, 1, 1.8, 2], [0, 1, .8, 1, 0], 'linear'); % amp. envelope
f=envgen([0, .2, 1.8, 2], [200, 250, 250, 200], 'linear'); % pitch envelope
f=f+max(f)*0.05*...      % add vibrato
    sin(2*pi*5*(SpF/Fs)*[0:length(f)-1]).*hanning(length(f));%'

[b_i,a_i]=baIIR2([300 2400 3000],[200 200 500]); %spec. envelope /i/
[b_a,a_a]=baIIR2([700 1200 2500],[200 300 500]); %spec. envelope /a/
[b_e,a_e]=baIIR2([570 1950 3000],[100 100 800]); %spec. envelope /e/

%%% compute sound %%%
s=buzz(0,amp,f); %impulse source
si=filter(b_i(1,:),a_i(1,:),s)+... %synthesize /i/
    filter(b_i(2,:),a_i(2,:),s)+filter(b_i(3,:),a_i(3,:),s);
sa=filter(b_a(1,:),a_a(1,:),s)+... %synthesize /a/
    filter(b_a(2,:),a_a(2,:),s)+filter(b_a(3,:),a_a(3,:),s);
se=filter(b_e(1,:),a_e(1,:),s)+... %synthesize /e/
    filter(b_e(2,:),a_e(2,:),s)+filter(b_e(3,:),a_e(3,:),s);
```

Note the use of the `filter` function. Figure 2.22 shows the spectrum of the original source signal, the spectral envelope defined by the filtering elements, and the spectrum of the final signal for two different pitches

2.5.3 Linear prediction

The problem of extracting a spectral envelope from a signal spectrum is generally an ill-posed problem. If the signal contains harmonic partials only, one could state that the spectral envelope is the



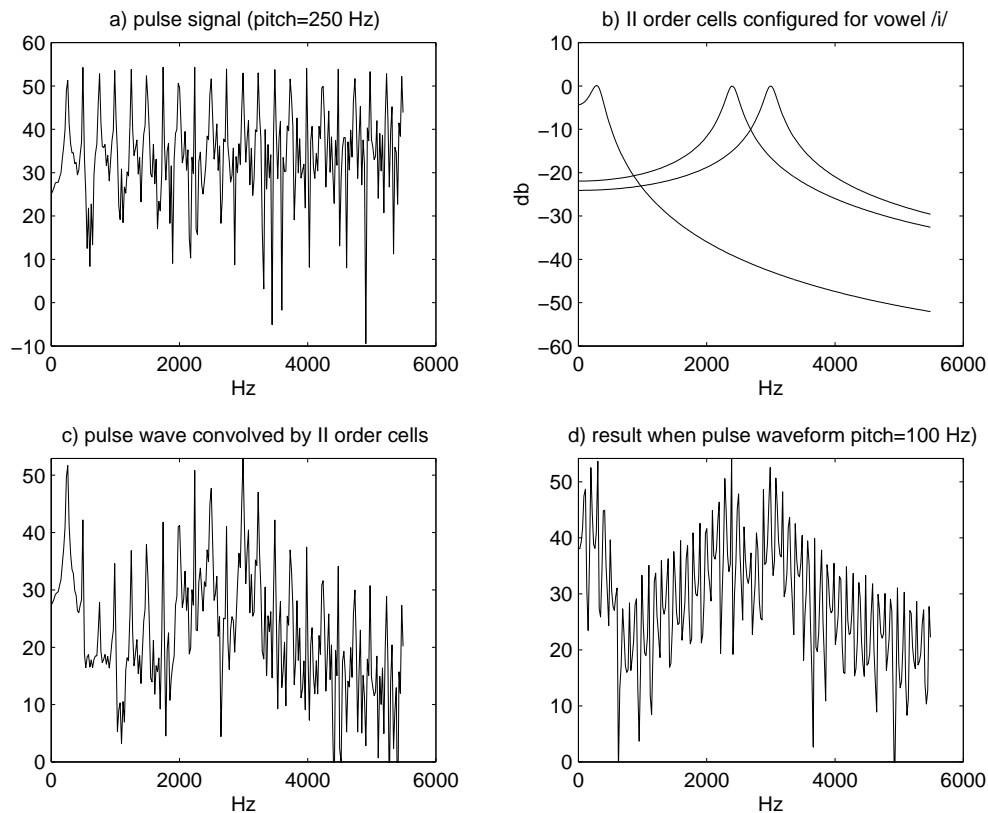


Figure 2.22: Spectrum of the original source signal, the spectral envelope defined by the filtering elements, and the spectrum of the final signal for two different pitches when a parallel synthesis structure is used.

curve that passes through the partial peaks. This implies that 1) the peak values have to be retrieved, and 2) an interpolation scheme should be (arbitrarily) chosen for the completion of the curve in between the peaks. If the sound contains inharmonic partials or a noisy part, then the notion of a spectral envelope becomes completely dependent on the definition of what belongs to the “source” and what belongs to the “filter”.

Three techniques with many variants can be used for the estimation of the spectral envelope:

- The channel vocoder uses frequency bands and performs estimations of the amplitude of the signal inside these bands and thus the spectral envelope.
- Linear prediction estimates an all-pole filter that matches the spectral content of a sound. When the order of this filter is low, only the formants are taken, hence the spectral envelope.
- Cepstrum techniques perform smoothing of the logarithm of the FFT spectrum (in decibels) in order to separate this curve into its slow varying part (the spectral envelope) and its quickly varying part (the source signal).

In this section we present the basics of Linear Prediction (LP) techniques. We will return on cepstral analysis in Chapter *From sound to content*.

2.5.3.1 Linear prediction equations

Consider a general linear system that describes a source-filter model:

$$S(z) = gH(z)X(z), \quad \text{with } H(z) = \frac{1 + \sum_{k=1}^q b_k z^{-k}}{1 - \sum_{k=1}^p a_k z^{-k}}, \quad (2.38)$$

where g is a gain scaling factor and $X(z)$ and $S(z)$ are the Z-transforms of the source signal $x[n]$ and the output signal $s[n]$, respectively. This is often termed an *ARMA*(p, q) (Auto-Regressive Moving Average) model, in which the output is expressed as a linear combination of p past samples and $q + 1$ input values. LP analysis works on an approximation of this system, namely on an all-pole model:

$$S(z) = gH_{LP}(z)X(z), \quad \text{with } H_{LP}(z) = \frac{1}{1 - \sum_{k=1}^p a_k z^{-k}}. \quad (2.39)$$

The time-domain version of this equation reads $s[n] = gx[n] + \sum_{k=1}^p a_k s[n-k]$. Therefore the output $s[n]$ can be predicted using a linear combination of its p past values, plus a weighted input term. In statistical terminology, the output *regresses* on itself, therefore system (2.39) is often termed an *AR*(p) (Auto-Regressive) model

One justification of this approximation is that the input signal $x[n]$ is generally unknown together with the filter $H(z)$. A second more substantial reason is that any filter $H(z)$ of the form (2.38) can be written as $H(z) = h_0 H_{min}(z) H_{ap}(z)$, where h_0 is a constant gain, $H_{min}(z)$ is a minimum-phase filter, and H_{ap} is an all-pass filter (i.e. $|H_{ap}(e^{2\pi jf})| = 1, \forall f$). Moreover, the minimum-phase filter $H_{min}(z)$ can be expressed as an all-pole system of the form (2.39). Therefore we can say that LP analysis ideally represents the all-pole minimum-phase portion of the general system (2.38), and therefore yields at least a correct estimate of the magnitude spectrum.

Given an output signal $s[n]$, Linear Prediction analysis provides a method for determining the “best” estimate $\{\tilde{a}_i\}$ ($k = 1, \dots, p$) for the coefficients $\{a_i\}$ of the filter (2.39). The method can be interpreted and derived in many ways, here we propose the most straightforward one. Given an estimate $\{\tilde{a}_i\}$ of the filter coefficients, we define the *linear prediction* $\tilde{s}[n]$ of the output $s[n]$ as $\tilde{s}[n] = \sum_{k=1}^p \tilde{a}_k s[n-k]$. In the Z-domain we can write

$$\tilde{S}(z) = P(z)S(z), \quad \text{with } P(z) = \sum_{k=1}^p a_k z^{-k}, \quad (2.40)$$

and we call the FIR filter $P(z)$ a *prediction filter*. We then define the *prediction error* or *residual* $e[n]$ as the difference between the output $s[n]$ and the linear prediction $\tilde{s}[n]$. In the Z-domain, the prediction error $e[n]$ is expressed as

$$E(z) = A(z)S(z), \quad \text{with } A(z) = 1 - \sum_{k=1}^p a_k z^{-k}. \quad (2.41)$$

Comparison of Eqs. (2.39) and (2.41) shows that, if the speech signal obeys the model (2.39) exactly, and if $\tilde{a}_k = a_k$, then the residual $e[n]$ coincides with the unknown input $x[n]$ times the gain factor g , and $A(z)$ is the inverse filter of $H_{LP}(z)$. Therefore LP analysis provides an estimate of the inverse system of (2.39):

$$e[n] = gx[n], \quad A(z) = [H_{LP}(z)]^{-1}. \quad (2.42)$$

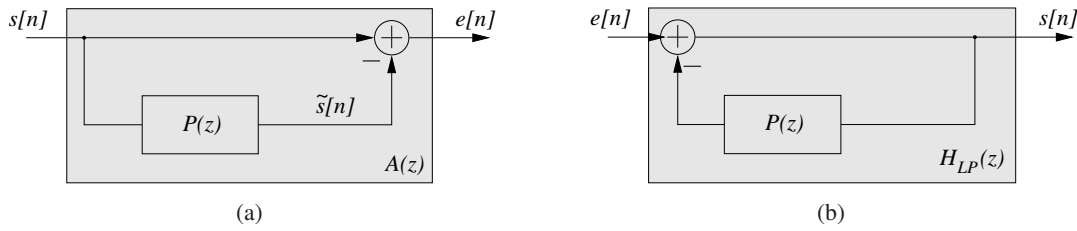


Figure 2.23: LP analysis: (a) the inverse filter $A(z)$, and (b) the prediction error $e[n]$ interpreted as the unknown input $gx[n]$.

This interpretation is illustrated in Fig. 2.23.

$A(z)$ is also called a *whitening* filter, because it produces a residual with a flat power spectrum. Two kinds of residuals, both having a flat spectrum, can be identified: the pulse train and the white noise. If LP is applied to speech signals, the pulse train represents the idealized vocal-fold excitation for voiced speech, while white noise represents the idealized excitation for unvoiced speech.

The roots of $A(z)$ (i.e., the poles of $H_{LP}(z)$) are representative of the formant peak frequencies. In other words, the angles of the roots, expressed in terms of analog frequencies, can be used as an estimate of the formant frequencies. The frequencies f_i and bandwidths B_i of the poles are extracted from complex conjugate roots r_i of the corresponding polynomial as:

$$f_i = F_s \frac{\arg(r_i)}{2\pi}, \quad B_i = F_s \frac{\log(|r_i|)}{\pi}. \quad (2.43)$$

We now describe the heart of LP analysis and derive the equations that determine the “best” estimate $\{\tilde{a}_i\}$ ($k = 1, \dots, p$). In this context “best” means best in a least-square sense: we seek the $\{\tilde{a}_i\}$ s that minimize the energy $E\{e\} = \sum_{m=-\infty}^{+\infty} e^2[m]$ of the residual, i.e. we set to zero the partial derivatives of $E\{e\}$ with respect to the a_i s:

$$0 = \frac{\partial E\{e\}}{\partial a_i} = 2 \sum_{m=-\infty}^{+\infty} e[m] \frac{\partial e[m]}{\partial a_i} = -2 \sum_{m=-\infty}^{+\infty} \left\{ \left[s[m] - \sum_{k=1}^p a_k s[m-k] \right] s[m-i] \right\}, \quad (2.44)$$

for $i = 1 \dots p$. If one defines the temporal autocorrelation of the signal $s[n]$ as the function $r_s[i] = \sum_{m=-\infty}^{+\infty} s[m]s[m-i]$, then the above equation can be written as

$$\sum_{k=1}^p a_k r_s[i-k] = r_s[i], \quad \text{for } i = 1 \dots p. \quad (2.45)$$

The system (2.45) is often referred to as the *normal equations*. Solving this system in the p unknowns a_i yields the desired estimates \tilde{a}_i .

2.5.3.2 Short-time LP analysis

In many applications of interest, and in particular analysis and resynthesis of speech signals, the coefficients a_k are not constant but slowly time-varying. Therefore the description (2.39) is only valid in a short-time sense, i.e. the a_k s can be assumed constant during an analysis frame. Therefore short-time analysis has to be used, in which the coefficients and the residual are determined from windowed sections, or frames, of the signal.

There are various efficient methods to compute the filter coefficients, the most common ones being the autocorrelation method, the covariance method, and the Burg algorithm. In this section we briefly describe the autocorrelation method, that simply amounts to substitute the autocorrelation function r_s of Eq. (2.45) with its short-time approximation:

$$r_s[i] \sim \sum_{m=1}^N u[m]u[m-i], \quad \text{where } u[m] = s[m]w[m] \quad (2.46)$$

is a windowed version of $s[m]$ in the considered frame ($w[m]$ is typically a Hamming window), and N is the length of the frame. Then the system (2.45) is solved within each frame. An efficient solution is provided by the so-called *Levinson-Durbin recursion*, an algorithm for solving the problem $\mathbf{Ax} = \mathbf{b}$, with \mathbf{A} Toeplitz, symmetric, and positive definite, and \mathbf{b} arbitrary. System (2.45) is an instance of this general problem.

M-2.18

Write a function `lp_coeffs` that computes the LP coefficients of a signal $s[n]$ given the desired prediction order p .

M-2.18 Solution

```
% Compute LP coeffs using the autocorrelation method
% s is the signal, p is the prediction order
% a are the computed LP coefficients, g is the gain (sqrt of residual variance)

function [a,g] = lpc_coeffs(s,p)

R=xcorr (s,p) ; % autocorrelation sequence R(k) with k=-p,...,p
R(1:p)=[];      % delete the first p samples
if norm(R) ~= 0
    [a,v] = levinson(R,p); % Levinson-Durbin recursion
                        % a=[1,-a1,-a2,...,-ap], v = variance of the residual
else
    a=[1, zeros(1,p)];
end
g=sqrt(sum(a' .*R')) ; % gain factor (= sqrt(v))
```

Note that we are using the native function `levinson`, that computes the filter coefficients (as well as the variance of the residual) given the autocorrelation sequence and the prediction order.

Figure 2.24 shows an example of LP analysis and resynthesis of a single frame of a speech signal. As shown in Fig. 2.24(a), the analyzed frame is a portion of voiced speech and $s[n]$ is pseudo-periodic. Correspondingly, the estimated source signal $x[n]$ is a pulse train. Figure 2.24(b) shows the magnitude responses of the target signal and the estimated transfer function $gH_{LP}(z)$. A typical feature of LP spectral modeling can also be observed from this figure: the LP spectrum matches the signal spectrum much more closely in the region of large signal energy (i.e. near the spectrum peaks) than near the regions of low energy (spectrum valley).

M-2.19

Write an example script that analyzes and resynthesizes a frame of speech using the LP model (2.39).

M-2.19 Solution



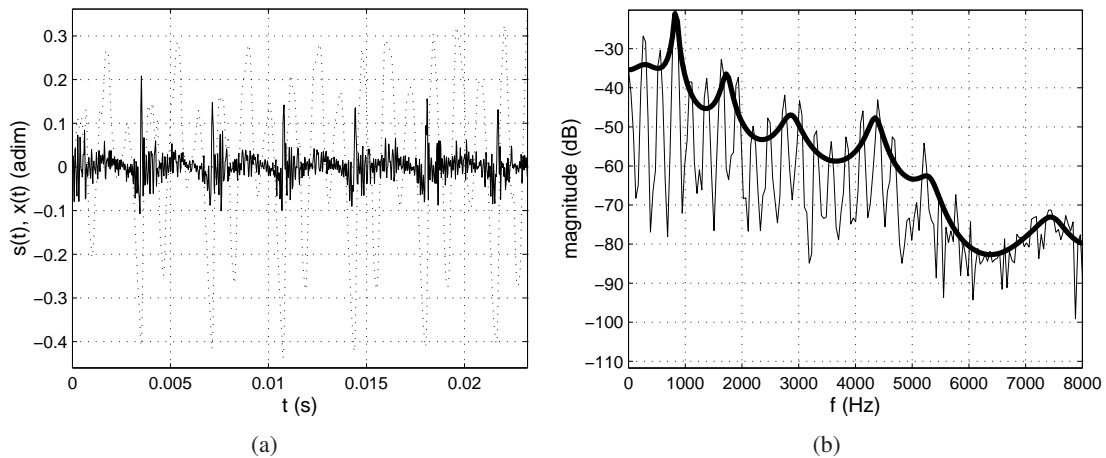


Figure 2.24: Example of LP analysis/synthesis, with prediction order $p = 50$; (a) target signal $s[n]$ (dotted line) and unit variance residual $x[n]$ (solid line); (b) magnitude spectra $|S(f)|$ (thin line) and $|gH_{LP}(f)|$ (thick line).

```

fname = 'la.wav'; %file to be processed

%% analysis parameters
N=1024; %block length
Nfft=1024; % FFT length
p=28 ; %prediction order
pre=p; %filter order= no. of samples required before n0

n0=5000; %start index
n1=n0+N-1; %end index

[s, Fs] = wavread(fname);
s=s(n0-pre:n1,1)'; %' row vector of left channel
win=hamming(N)' ; %' window for input block

[a,g]=lp_coeffs(s((1:N)+pre).*win, p); % compute LP coeffs

x=filter(a,g,s); % X(z)=A(z)/g * S(z), i.e. x[n] = e[n]/g
s_new=filter(g,a,x); %resynthesized signal

```

Note that we have used the function `lp_coeffs` written in example M-2.18. The signals plotted in Fig. 2.24 have been computed from this script.

When formant parameters are extracted on a frame-by-frame basis, a lot of discontinuities and local estimation observation errors are found. Therefore, proper techniques have to be used in order to determine smooth formant trajectories over analysis frames. We have already encountered a conceptually similar problem in Sec. 2.4.3, when we have discussed a “sinusoid tracking” procedure.

M-2.20

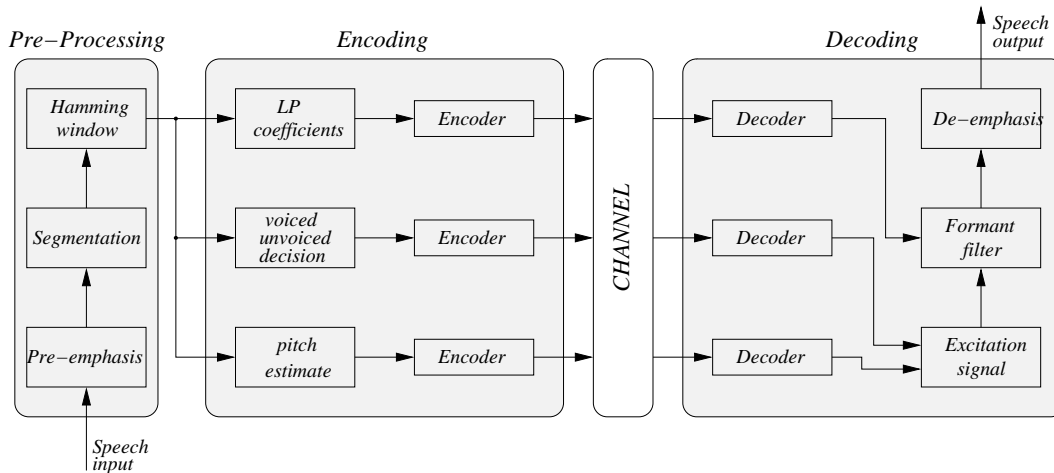


Figure 2.25: General scheme of a simple LP based codec.

Plot the formant frequencies as a function of the frame number, i.e., of time, in order to observe the time-evolution of the vocal tract filter. To this purpose, segment a speech signal $s[n]$ into M Hamming windowed frames $s_m[n]$, with a block length N and a hop-size $S_a = N/2$. Then, for each frame: a) compute the LP coefficients; b) find the formant frequencies using Eq. (2.43); c) discard roots whose magnitude is less than 0.8, as these are unlikely to be formants.

2.5.3.3 Linear Predictive Coding (LPC)

One of the most successful applications of LP analysis is in speech coding and synthesis, in particular for mobile phone communication. Figure 2.25 depicts a synthetic block diagram of a simple encoder-decoder system based on LP analysis. Speech is segmented in frames (typical frame lengths can range from 10 to 20 ms). In this phase a pre-emphasis processing can also be applied: since the lower formants contain more energy, they are preferentially modeled with respect to higher formants, and a pre-emphasis filter can compensate for this by boosting the higher frequencies (when reconstructing the signal the inverse filter should be used).

In its simplest formulation the encoder provides, for every frame, the coefficients a_k of the prediction filter, the gain g , a flag that indicates whether the frame corresponds to voiced or unvoiced speech, and the pitch (only in the case of voiced speech). The decoder uses this information to re-synthesize the speech signal. In the case of unvoiced speech, the excitation signal is simply white noise, while in the case of voiced speech the excitation signal is a pulse train whose period is determined by the encoded pitch information.

It is clear that most of the bits of the encoded signal are used for the a_k parameters. Therefore the degree of compression is strongly dependent on the order p of the LP analysis, which in turn has a strong influence on the degree of smoothness of the estimated spectral envelope, and consequently on the quality of the resynthesis (see fig 2.26). A commonly accepted operational rule for achieving reasonable intelligibility of the resynthesized speech is

$$p = \begin{cases} F_s + 4 & \text{for voiced speech} \\ F_s & \text{for unvoiced speech} \end{cases} \quad (2.47)$$

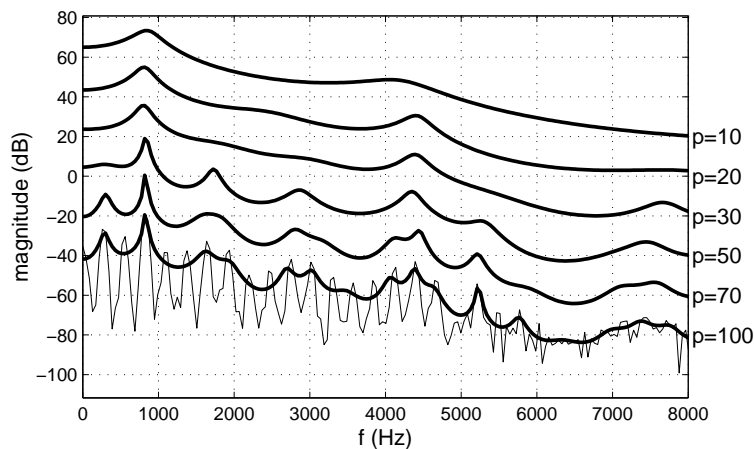


Figure 2.26: Example of LP spectra for increasing prediction orders p (the target signal is a frame of voiced speech). For the sake of clarity each spectrum is plotted with a different offset.

LPC-10 is an example of a standard that basically implements the scheme of Fig. 2.25. This standard uses an order $p = 10$ (hence the name), a sample rate $F_s = 8$ kHz (which is a common choice in telephone speech applications since most of the energy in a speech signal is the range [300, 3400 Hz]), and a frame length of about 22.5 ms. With these values, intelligible speech can be resynthesized.

However *LPC-10*, and in general similar early codecs, produced speech with very poor quality due to many artifacts: “buzzy” noise through parameter updates, jitter in the excitation signal, wide formant bandwidths, and so on. More recent and commonly used codecs are able to provide natural sounding speech at relatively low bit rates, thanks to an improved description of the excitation signal. Instead of applying a simple two-state voiced/unvoiced model, these codecs estimate the excitation signal through an analysis-by-synthesis approach: excitation waveforms are passed through the formant filter, and the excitation which gives the minimum weighted error between the original and the reconstructed speech is then chosen by the encoder and used to drive the synthesis filter at the decoder. It is this ‘closed-loop’ determination of the excitation which allows these codecs to produce good quality speech at low bit rates, at the expense of a much higher complexity of the coding stage.

Within this family of analysis-by-synthesis codecs, many different techniques have been developed for the estimation of the excitation signal. Historically, the first one is the Multi-Pulse Excited (MPE) codec. Later the Regular-Pulse Excited (RPE), and the Code-Excited Linear Predictive (CELP) codecs were introduced. The “Global System for Mobile communications” (GSM), a digital mobile radio system which is extensively used throughout Europe, and also in many other parts of the world, makes use of a RPE codec.

2.5.3.4 LP based audio effects

There are many LP based audio effects which use a synthesis filter calculated from a speech input signal. The LP filter represents the time-varying properties of the vocal tract by modelling its formant frequencies. For a proper vocal tract model the prediction order should be slightly higher than the sampling rate in kHz which is equivalent to modelling one resonance frequency each kHz since each resonance frequency requires two conjugate complex poles. Zeros in the spectral model may be

approximated by a small number of poles. Thus for a good model the prediction order should be a small amount higher than the sampling rate in kHz.

Some effects are based on the LP of one sound. Thus the standard LP analysis/synthesis is used, but modifications are applied either to the excitation or to the synthesis filter or to both of them.

The excitation for ideal signal recovery can be modelled as a signal consisting of pulses plus noise, as used in the speech production model. For voiced sounds the periodicity of the pulses determine the pitch of the sound while for unvoiced sounds the excitation is noise-like. The modelling of the excitation requires an algorithm for voiced/unvoiced separation, which can crudely be a switch between pulses and random noise. Such decisions are not so easy to make automatically and a good pitch detector should be used. But then all manipulations are possible.

The synthesis filter can be modified by taking a warped version of the initial filter which moves the formants and will give a donald duck voice without altering its pitch. The time duration of a sound can be modified by time-stretching the excitation and updating the synthesis filter at a different rate. This preserves the formant structure. To modify the pitch of the resulting signal, for voiced parts the pitch of the modelled excitation can be changed. This effect can also be combined with the frequency warping to independently change the formants induced by the synthesis filter.

An frequently used effect is cross-synthesis between two sounds. It uses two different sounds ($x_1[n]$ and $x_2[n]$) and take the residual of $x_2[n]$ as the excitation to the LP filter of $x_1[n]$. The cross-synthesis gives good results if a speech signal is used to compute the synthesis filter which results for example in the talking orchestra. For musical satisfactory results the two used sounds have to be synchronized. Thus, for the case of mixing speech and music, the played instrument must fit to the rhythm of the syllables of the speech. The performance is improved if either speech or music is coming from a prerecorded source and the other sound is produced to match to the recording.

2.6 Non linear processing models

The transformations seen in sect. 2.5, since they are linear, cannot change the frequency of the components that are present. Instead, when non linear transformations are used, frequencies can be even drastically changed and new components are created. Thus, it is possible to vary substantially the nature of the input sounds.

There are two main effects related to nonlinear transformations: spectrum enrichment and spectrum shift. The first effect is due to non linear distortion of the signal and allows for controlling the brightness of a sound, while the second is due to its multiplication by a sinusoid and moves the spectrum to the vicinity of the carrier signal, altering the harmonic relationship between the modulating signal lines. The possibility of shifting the spectrum is very intriguing in when applied to music. From simple components, harmonic and inharmonic sounds can be created, and various harmonic relations among the partials can be established. The first effect try to reproduce the nonlinearities and saturations found on real systems e.g. analog amplifiers, electronic valves. The second one instead derives from abstract mathematical properties of trigonometric functions as used in modulation theory applied to music signal. Therefore, it inherits, in part, the analogic interpretation as used in electronic music and is a new metaphor for computer musicians.

2.6.1 Memoryless non linear modelling

When a sinusoidal input sound $x[n] = \cos(2\pi f_1 T n)$ passes through a linear system (filter) with transfer function $H(f)$, the output signal $y[n]$ is still a sinusoid with the same frequency f_1 and



amplitude and phase depending on the transfer function, i.e. $y[n] = |H(f_1)| \cos(2\pi f_1 T n + \angle H(f_1))$. Instead if it passes through a non linear amplifier, the waveform is modified and various harmonics are created. Normally we want to avoid distortions in amplifiers, but sometimes, as in amplifiers for electric guitars, we may be interested in emulating the warm sound of valves. In general the output value of a non linear system, depends on present and past values of the input and can be described by the Volterra series expansion. This kind of system are used to compensate the non linear behaviour found in real system, e.g. to linearize the loudspeakers. But is quite to complicate to use and to control. Thus is not suitable for musicians. For this reason in music signal processing often non linearities without memory, i.e. that depends only on the present input value and not on the past values, are used. In this case the system is described by a non linear curve $F(x)$, called *distortion function*, and the output is given by

$$y[n] = F(x[n]) \quad (2.48)$$

In analog domain it is difficult to have an amplifier with a precise and variable distortion characteristic. In digital domain the distortion function can be previously computed and stored in a table. During the processing, all that is necessary is to look up the desired value in the table, with an eventual interpolation between adjacent points, as for the table look up oscillator. In general the distortion function produces infinite partials giving rise to annoying foldover. If $F(x)$ is a polynomial or its Taylor series expansion can be truncated, the bandwidth remains limited. However, non linear distortions in digital signals can easily surpass $F_s/2$. In this case $x[n]$ is suitably low pass filtered before the non linear processing, or the non linear computation is done on a oversampled version of the signal.

For example in order to simulate the overdrive effect of guitar amplifiers, we can use the function

$$F(x) = \begin{cases} 2x & \text{for } 0 \leq x \leq 1/3 \\ \frac{3-(2-3x)^2}{3} & \text{for } 1/3 \leq x \leq 2/3 \\ 1 & \text{for } 2/3 \leq x \leq 1 \end{cases}$$

that produces a symmetrical soft clipping of the input and realizes a smooth transition of the linear behaviour for low level signal to a high saturation for high level sounds. Overdrive has a warm and smooth sound. Asymmetric functions are used for tube simulations. More nonlinear functions are employed for distortion simulation producing tones starting beyond tube warmth to buzz saw effects.

The same technique was used for sound synthesis. In this case we have a sinusoidal input and rich harmonic sound is produced using a distortion function

$$F(x) = p(a \cdot x + b)$$

where $p(x)$ is a polynomial (or a rational) function, a and b are parameters that are used to scale and shift the distortion function.

2.6.2 Synthesis by frequency modulation

This technique does not derive from models of sound signals or sound production; instead it is based on an abstract mathematical description. The definition of “FM synthesis” denotes an entire family of techniques in which the instantaneous frequency of a periodic signal (*carrier*) is itself a signal that varies at sample rate (*modulating*).

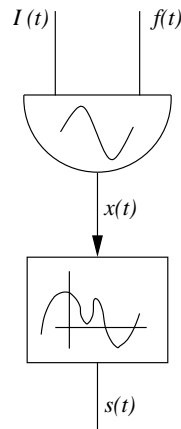


Figure 2.27: Sintesi per distorsione non lineare

M-2.21

We have already seen in section 2.2.1.4 how to compute the signal phase when the frequency is varying at frame rate. We now face the problem of computing $\phi[n]$ when the frequency varies at audio rate. A way of approximating $\phi[n]$ is through a first-order expansion:

$$\phi[n] = \phi[n-1] + \frac{d\phi}{dt}(n-1) \cdot \frac{1}{F_s}. \quad (2.49)$$

Recalling equation (2.8), that relates phase and instantaneous frequency, we can approximate it as

$$\frac{d\phi}{dt}(n-1) = 2\pi \left[\frac{f[n] + f[n-1]}{2} \right], \quad (2.50)$$

where the frequency $f(t)$ has been approximated as the average of $f[n]$ at two consecutive instants. Using the two equations above, $\phi[n]$ is finally written as

$$\phi[n] = \phi[n-1] + \frac{\pi}{F_s} (f[n] + f[n-1]). \quad (2.51)$$

Write a function `FMosc(t0, a, f, ph0)` that realizes a FM sinusoidal oscillator (the parameters `(t0, a, ph0)` are defined as in M-2.3, while `f` is now the sample-rate frequency vector).

M-2.21 Solution

```
function s=FMosc(t0,a,f,ph0)

global SpF;           %samples per frame
global Fs;           %sampling rate

nframes=length(a);   %total number of frames

s=zeros(1,nframes*SpF); %signal vector (initialized to 0)

lastfreq=f(1);
lastphase=ph0;
for (i=1:nframes)    %cycle on the frames
    phase=zeros(1,SpF); %phase vector in a frame
    for(k=1:SpF)     %cycle through samples in a frame
        phase(k)=lastphase+... %compute phase at sample rate
            pi/Fs*(f((i-1)*SpF+k)+lastfreq);
```

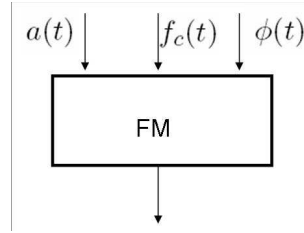


Figure 2.28: FM basic computing module

```

    lastphase=phase(k);
    lastfreq=f((i-1)*SpF+k);
end
s((i-1)*SpF+1:i*SpF)=a(i).*cos(phase);
end

s=[zeros(1,round(t0*Fs+1)) s]; %add initial silence of t0 sec.

```

Compare this function with the `sinosc` function in M-2.3. The only difference is that in this case the frequency is given at audio rate. Consequently the phase computation differs.

Although early realizations of FM synthesis were implemented in this fashion, in the next sections we will follow an equivalent “phase-modulation” formulation. According to such formulation, the FM oscillator is written as:

$$s(t) = \sin[2\pi f_c t + \phi(t)], \quad (2.52)$$

where $\phi(t)$ is the input modulating signal and f_c is the carrier frequency.

For sound synthesis, an amplitude envelope $a(t)$ should be applied. Thus we define a basic FM module that computes

$$y(t) = \mathbf{FMmodule}[a(t), f_c(t), \phi(t)] = a(t) \cdot \sin[2\pi f_c(t)t + \phi(t)] \quad (2.53)$$

and is often represented as in fig. 2.28. The algorithm is given by

$$\varphi[n] = \varphi[n-1] + \frac{2\pi}{F_s} f_c[n] + \phi[n] \quad (2.54)$$

$$y[n] = a[n] \cdot \sin(\varphi[n]) \quad (2.55)$$

where $\varphi[n]$ is a state variable representing the instantaneous phase of the oscillator. Notice that when the oscillator is implemented by a wavetable, the phases φ and ϕ vary in the interval $0 \dots (L-1)$ and the algorithm becomes

$$\varphi_L[n] = \varphi_L[n-1] + \frac{L}{F_s} f_c[n] + \phi_L[n] \quad (2.56)$$

$$y[n] = a[n] \cdot \text{tabsin}(\varphi_L[n] \bmod L) \quad (2.57)$$

2.6.2.1 Simple modulation

When the modulating signal $\phi(t)$ is a sinusoid with amplitude I (modulation index) and frequency f_m , i.e.

$$\phi(t) = I \sin(2\pi f_m t)$$

the simple modulation gives

$$s(t) = \sin [2\pi f_c t + I \sin(2\pi f_m t)] \quad (2.58)$$

$$= \sum_{k=-\infty}^{\infty} J_k(I) \sin [2\pi (f_c + k f_m) t] \quad (2.59)$$

where $J_k(I)$ is the Bessel function of first kind and k -th order computed in I . From equation 2.59 we can see that the resulting spectrum is composed of partials at frequencies $|f_c \pm k f_m|$ with amplitude given by $J_k(I)$. Notice that negative frequencies, being sine waves, are folded changing the sign. Apparently there are infinite partials, so theoretically the signal bandwidth is not limited. However it is practically limited. In the Bessel function behaviour, only few low-order functions are significant for small index values. When the index increases, the number and the order of significant function increase too. Usually in the bandwidth definition of the FM signal, The number M of lateral spectral lines (sidebands) greater than 1/100 of the nonmodulated signal is given by $M(I) = I + 2.4 \cdot I^{0.27}$, that can be approximates as $M(I) = 1.5 * I$. In this way, varying I it is possible to directly control the bandwidth around f_c . The resulting effect is similar to a low pass filter with varying cut-off frequency. Moreover the amplitude of the partials varies in a smooth way, maintaining constant the overall signal energy.

For the synthesis we can use two basic FM modules in cascade

$$y(t) = \mathbf{FMmodule} [a(t), f_c(t), \mathbf{FMmodule} [I(t), f_m(t), 0]] \quad (2.60)$$

2.6.2.2 Spectra $|f_c \pm k f_m|$

Equation 2.59 shows a spectrum with lines at frequencies $|f_c \pm k f_m|$, with $k = 0, 1, \dots$. This kind of spectra are characterized by the ratio f_c/f_m , sometimes also called c/m ratio. When this ratio is rational, it can be expressed as an irreducible fraction $f_c/f_m = N_1/N_2$ with N_1 and N_2 as integers that are prime between themselves. In this case the resulting sound is periodic, since all the partials are a multiple of a fundamental frequency f_0 according to integer factors

$$f_0 = \frac{f_c}{N_1} = \frac{f_m}{N_2}, \quad (2.61)$$

and f_c, f_m coincides with the N_1 -th and N_2 -th harmonic:

$$f_c = N_1 f_0, \quad f_m = N_2 f_0. \quad (2.62)$$

If $N_2 = 1$, all the harmonics are present and the sideband components with $k < -N_1$, i.e. with negative frequency, overlap some components with positive k . If $N_2 = 2$, only odd harmonics are present, and sidebands superimpose, after foldunder. coincide. If $N_2 = 3$, the harmonics that are multiple of 3 are missing. In general the N_1/N_2 ratio can be considered as an index of the harmonicity of the spectrum. The sound is more harmonious intuitively, when the N_1/N_2 ratio is simple and formally when the $N_1 \cdot N_2$ is smaller.

The ratios can be grouped in families. All ratios of the type $|f_c \pm k f_m|/f_m$ can produce the same components that f_c/f_m produces. Only the partial coinciding with the carrier f_c changes. Remember that $f_c = N_1 f_0$. For example the ratios $2/3, 5/3, 1/3, 4/3, 7/3$ and so on belong to the same family. Only the harmonics that are multiple of 3 are missing (see $N_2 = 3$) and the carrier is respectively the second, fifth, first, fourth, seventh harmonic. The ratio that distinguish a family is defined in normal

form when it is $\leq 1/2$. In the previous example, it is $1/3$. Each family is characterized by a ratio in normal form. Similar spectra can be produced using ratios from the same family. We can notice that the denominator N_2 characterizes the spectrum.

When the f_c/f_m quotient is irrational, the resulting sound is aperiodic and hence inharmonic. This possibility is used to easily create inharmonic sounds as bells. For example if $f_c/f_m = 1/\sqrt{2}$, the sound contains partials with frequency $f_c \pm k\sqrt{2}$. No implied fundamental pitch is audible. A similar behaviour can be obtained with complex ratios as $f_c/f_m = 5/7$.

Of particular interest is the case of an f_c/f_m ratio approximating a simple rational value, that is,

$$\frac{f_c}{f_m} = \frac{N_1}{N_2} + \epsilon. \quad (2.63)$$

Here the sound is no longer rigorously periodic. The fundamental frequency is still $f_0 = f_m/N_2$ and the harmonics are shifted from their exact value by $\pm\epsilon f_m$. Thus a small shift of the carrier frequency, does not change the pitch, even if it slightly spread the partials and makes the sound more lively. Notice that the same shift of f_m changes the pitch.

M-2.22

Synthesize a frequency modulated sinusoid, in the case of sinusoidal modulation. Plot the signal spectrum for increasing values of the modulation index.

M-2.22 Solution

```

%%% headers %%%
Fs=22050;          % sampling frequency

%%% define controls %%%
fc=700;           %carrier freq.
fm=100;           %modulating freq.
I=2;              %modulation index
t=0:(1/Fs):3;     %time vector (in s)

%%% compute sound %%%
s=sin(2*pi*fc*t+I*sin(2*pi*fm*t));

```

Figure 2.6.2.2 shows the signal spectrum for 3 values of the modulation index. Note that as the index increases the energy of the carrier frequency is progressively transferred to the lateral bands, according to the predicted behaviour.

2.6.2.3 Compound modulation

There are many variation of the basic scheme. If the modulating signal is composed of N sinusoids, the following relation hold:

$$s(t) = \sin \left[2\pi f_c t + \sum_{i=1}^N I_i \sin(2\pi f_i t) \right] \quad (2.64)$$

$$= \sum_{k=-\infty}^{\infty} \prod_{i=1}^N J_k(I_i) \sin \left[2\pi \left(f_c + \sum_{k_i} k_i f_i \right) t \right] \quad (2.65)$$

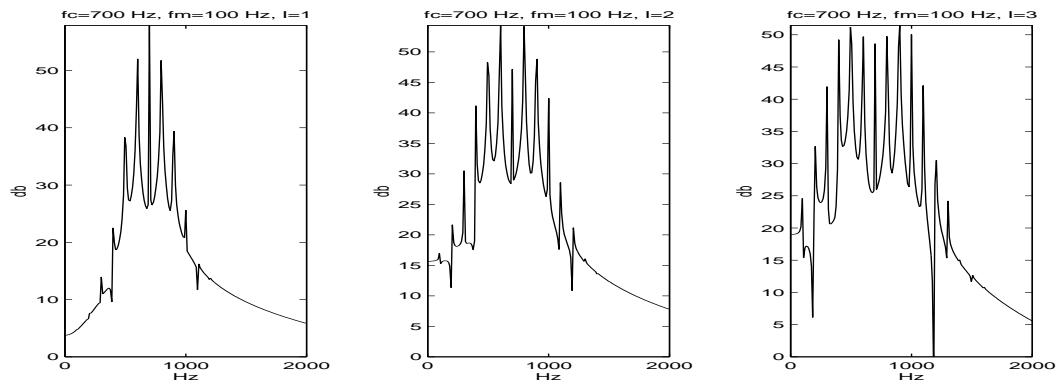


Figure 2.29: Spectrum of a simple modulation with $f_c = 700$ Hz, $f_m = 100$ Hz modulation index I varying from 1 to 3

The generated sound will have components of frequency $|f_c \pm k_1 f_1 \pm \dots \pm k_N f_N|$ with amplitudes given by the product of N Bessel functions. A very complex spectrum results. If the relations among the frequencies f_i are simple, that is, if the sum of the modulating waves is periodic with frequency f_m , then the spectrum is of the type $|f_c \pm k f_m|$. The frequency f_m can be computed as the greatest common divisor among the modulating frequencies f_i ($i = 1, \dots, N$). Otherwise the sonorities are definitely inharmonic and particularly noisy for high indexes.

For example Schottstaedt uses two modulators to simulate the piano tones, setting $f_1 \simeq f_c$ and $f_2 \simeq 4f_c$. It results that we can compute $f_m \simeq f_c$ and it results a pitch $f_0 = f_m \simeq f_c$. In this way the small inharmonicity of the piano strings is simulated. Moreover modulation indexes decrease for higher f_0 values. In this way the lower tones are richer in harmonics than higher ones.

M-2.23

Synthesize a frequency modulated sinusoid in the case of composite modulation. Plot the signal spectrum.

M-2.23 Solution

```

%%% headers %%%
%[...]

%%% define controls %%%
fc1=700;      %carrier freq.
fm=700;      %modulating freq. 1
fm=2800;     %modulating freq. 2
I1=1;        %modulation index 1
I2=1;        %modulation index 2
t=0:(1/Fs):3;%time vector (in s)

%%% compute sound %%%
s=sin(2*pi*fc*t+...      %sound signal
      I1*sin(2*pi*fm1*t)+I2*sin(2*pi*fm2*t));

```

Figure 2.30 shows the spectrum of an FM oscillator with sinusoidal carrier and a composite modulation made of two sinusoids. Note that in the first case the simple ratio between f_{m1} and f_{m2} produces a spectrum of the form $|f_c \pm k f_m|$ (in which $f_m = 100$ Hz, max. common divisor between

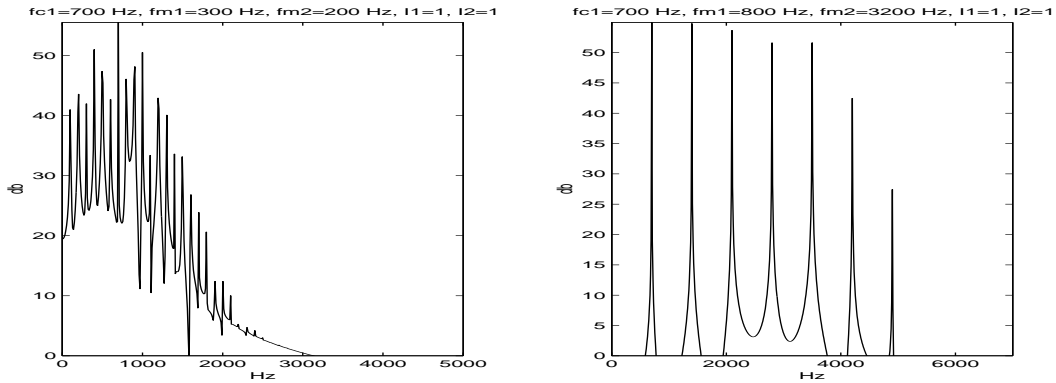


Figure 2.30: Two examples of compound modulation made of two sinusoids.

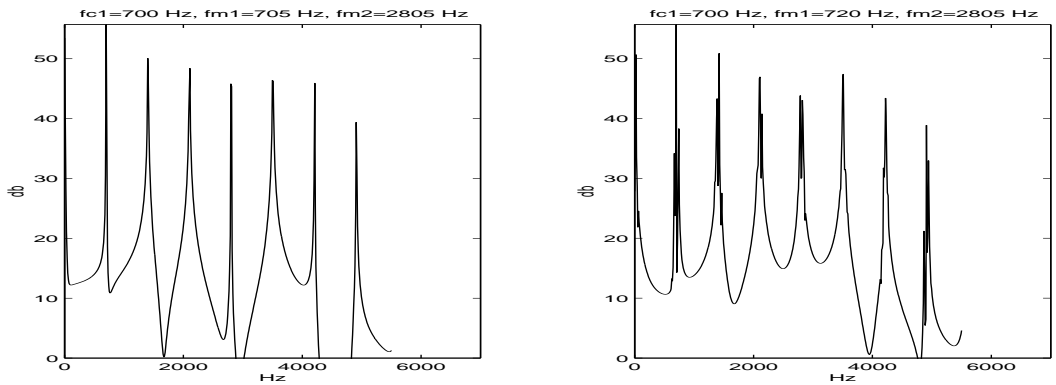


Figure 2.31: Quasi harmonic sounds produced by double modulators with $f_1 \simeq f_c$ and $f_2 \simeq 4f_c$.

f_{m1} and f_{m2} , is the fundamental). In the second case, the values $f_{m1} = f_c$ and $f_{m2} = 4f_c$ are chosen in such a way that the fundamental coincides with f_c and that upper partials are harmonic (since $f_{m1} = f_c$ coincides with the max. common divisor between f_{m1} and f_{m2}).

Figure 2.31 shows a double modulator with $f_1 \simeq f_c$ and $f_2 \simeq 4f_c$, with increasing deviations from the exact values multiple of f_c .

2.6.2.4 Nested and feedback modulation

When a sinusoidal modulator is phase modulated by another sinusoids we have

$$\phi(t) = I_1 \sin(2\pi f_1 t + I_2 \sin(2\pi f_2 t))$$

and the resulting signal is thus defined by:

$$\begin{aligned} s(t) &= \sin[2\pi f_c t + I_1 \sin(2\pi f_1 t + I_2 \sin(2\pi f_2 t))] \\ &= \sum_k J_k(I_1) \sin[2\pi(f_c + k f_1)t + k I_2 \sin(2\pi f_2 t)] \\ &= \sum_k \sum_n J_k(I_1) \cdot J_n(k I_2) \sin[2\pi(f_c + k f_1 + n f_2)t] \end{aligned}$$

The result can be interpreted as if each partial produced by the modulator f_1 were modulated in his turn by f_2 with modulation index kI_2 . The spectral structure is similar to that produced by two sinusoidal modulators, but with larger bandwidth.

As final variation of the basic technique let us consider the case that the past output value is used as modulating signal. This is the so called feedback modulation. This method is described by

$$\begin{aligned}\phi[n] &= \beta s[n-1] \\ s[n] &= \sin \left[2\pi \frac{f_c}{F_s} n + \phi[n] \right]\end{aligned}$$

where β is the feedback factor and acts as scale factor or feedback modulation index. For increasing values of β the resulting signal is periodic of frequency f_c and changes, in a continuous way, from a sinusoid to a saw-tooth waveform. The resulting spectrum has a increasing number of harmonics and it results

$$s(t) = \sum_k \frac{2}{k\beta} J_k(k\beta) \sin(2\pi k f_c t)$$

2.6.2.5 Discussion

Basically FM synthesis is a versatile method for producing many types of sounds. As of yet, however, no algorithm has been found for deriving the parameters of an FM model from the analysis of a given sound, and no intuitive interpretation can be given to the parameter choice as this synthesis technique does not evoke any previous musical experience of the performer. Its main qualities, i.e. great timbral dynamics with just a few parameters to control and to low computational costs, are progressively losing popularity when compared with other synthesis techniques which, though more expensive, can be controlled in a more natural and intuitive fashion. The FM synthesis, however, still preserves the attractiveness of its own peculiar timbral space and, though it is not particularly suitable for the simulation of natural sounds, it offers a wide range of original synthetic sounds that are of considerable interest for computer musicians.

2.6.3 Multiplicative synthesis

The simplest nonlinear transformation consists of the multiplication of two signals. In analog domain it is often called ring modulation (RM) and it is quite difficult to produce in a precise way. Let $x_1(t)$ e $x_2(t)$ be two input signals, the resulting signal is

$$s(t) = x_1(t) \cdot x_2(t) \quad (2.66)$$

and its spectrum $S(f)$ is obtained from the convolution of the two input signal spectra, i.e. $Y(f) = X_1(f) * X_2(f)$.

Usually one of the two input signals is sinusoidal with frequency f_c and is called carrier

$$c(t) = \cos(2\pi f_c t + \phi_c) \quad (2.67)$$

and the other input is the transformed signal and is called modulating signal $m(t)$ (or modulator). Equation 2.66 can be rewritten as

$$s(t) = m(t) \cdot c(t) = m(t) \cdot \cos(2\pi f_c t + \phi_c) \quad (2.68)$$



and the resulting spectrum is

$$S(f) = \frac{1}{2} \left[M(f - f_c) e^{j\phi_c} + M(f + f_c) e^{-j\phi_c} \right] \quad (2.69)$$

The spectrum of $s(t)$ is composed of two copies of the spectrum of $m(t)$: a lower sideband (LSB), reversed in frequency and an upper sideband(USB). The two sidebands are symmetric around f_c . When the bandwidth of $m(t)$ is greater than f_c , part of the LSB extends to the negative region of the frequency axis, and this part is folded around the origin (foldunder). Notice that the phase has to be taken into account while summing components of identical frequencies.

Let consider a sinusoidal carrier (2.67) and a periodic modulating signal of frequency f_m with N harmonics

$$m(t) = \sum_{k=1}^N b_k \cos(2\pi k f_m t + \phi)$$

We obtain

$$s(t) = \sum_{k=1}^N \frac{b_k}{2} [\cos [2\pi(f_c + k f_m)t + \phi_k] - \cos [2\pi(f_c - k f_m)t + \phi_k]] \quad (2.70)$$

The multiplicative synthesis causes every harmonic spectral line $k f_m$ to be replaced by two spectral lines, one in the LSB and the other in the USB, with frequency $f_c - k f_m$ e $f_c + k f_m$. Notice that the spectral lines of LSB with frequency $k f_m > f_c$, i.e. with $f_c/f_m < k \leq N$, are folded around zero. The resulting spectrum has components of frequency $|f_c \pm k f_m|$ with $k = 1, \dots, N$, where the absolute value is used to take into account the possible foldunder. The acoustic and perceptual properties of this kind of spectra will be discussed in sect. 2.6.2.2. If the carrier has many sinusoidal components, we will obtain two sidebands around each component and the resulting audio effect is less intuitive.

A variant of this method is amplitude modulation

$$s(t) = [1 + \delta m(t)] \cdot c(t) \quad (2.71)$$

where δ is the amplitude modulation index, that control the amplitude of the sidebands. The spectrum is as in eq. 2.69 plus the carrier spectral line

$$S(f) = C(f) + \frac{\delta}{2} \left[M(f - f_c) e^{j\phi_c} + M(f + f_c) e^{-j\phi_c} \right]$$

2.7 Key concepts

↪ Sound models

Definition of sound model. Sound synthesis algorithms as instrument families. Control parameters and playability of a model. Approaches to sound model classification: structure (direct generation, feed-forward models, feed-back models), or level of representation (signal models, source models).

↪ Signal generators

Table look-up oscillators, resampling and sampling increment. Recursive generators of sinusoidal signals. Definition of *control-rate*, and oscillators controlled in amplitude/frequency; in

particular, computation of the signal phase from the frame-rate frequency. Envelope generators and the ADSR amplitude envelope. Noise generators: definitions of white noise and pink ($1/f$) noise.

Sampling techniques: advantages and limitations. Sample processing: pitch shifting and looping.

↪ Granular synthesis

Definition of “grain”. Granulation of real signals, and the Overlap and Add (OLA) method. Related problems (discontinuities, control). Synthetic grains: frequency modulated gaussian grains, and visualization in the time-frequency plane. Synchronous and asynchronous granular synthesis, time-frequency masks.

↪ Additive synthesis

Basic formulation (equation (2.23)) and its relation to Fourier analysis. The problem of control signals. Additive synthesis by analysis: extraction of frequency and amplitude envelopes through *STFT*, sinusoid tracking procedure (figure 2.12).

“Sines-plus-noise” models: analysis and modeling of the stochastic component, resynthesis and transformations. “Sines-plus-transients-plus-noise” models: DCT representation, analysis and modeling of transients.

↪ Source-filter models

Basic block scheme (figure 2.18), and characteristics of the two blocks. Possible (physical) interpretations of this structure. Notable source blocks: noise generators and pulse generators. Notable filtering blocks: 2nd order resonators (equation (2.32)), center frequency and bandwidth. Time-varying filters and related problems.

↪ LP voice modeling

Functioning of the vocal system, schematization as a source-filter system. Vocalized versus non-vocalized sounds, and corresponding source signals. A simple subtractive model for the voice (equation (2.35), glottis, vocal tract, and lip radiation blocks. LP analysis, equivalence between extraction of the residual and reconstruction of the source signal.

↪ FM synthesis

Definition of carrier and modulating signals. Computation of the signal phase from the sample-rate frequency. Equivalent formulations: frequency modulation and phase modulation. The simplest case: sinusoidal carrier and sinusoidal modulation. Analysis of $f_c \pm kf_m$ spectrums (equation 2.59), classification in terms of the ratio f_c/f_m . Other FM structures: composite modulation, composite carrier, cascade and feedback modulation.

Advantages (compact structure, low number of control parameters) and drawbacks (non-intuitive control, lack of general model identification techniques) of FM.

2.8 Commented bibliography

Among the plethora of available sound synthesis languages, one of the most widely used (and one of the most important historically) is Csound, developed by Barry Vercoe at the Massachusetts Institute of Technology. Csound descends from the family of Music-N languages created by Max Mathews at Bell Laboratories. See [Vercoe, 1993].

A second influential sound synthesis programming paradigm was developed starting from the early 1980's, mainly by Miller Puckette, and is today represented in three main software implementation: Max/MSP, jmax, and Pd. The “Max paradigm” (so named in honor of Max Mathews) is described by Puckette [Puckette, 2002] as a way of combining pre-designed building blocks into sound-processing “patches”, to be used in real-time settings. This includes a scheduling protocol for both control- and audio-rate computations, modularization and component intercommunication, and a graphical environment to represent and edit patches.

A discussion on recursive generators of sinusoidal signals is found e.g. in [Orfanidis, 1996]. Models for fractal signals are also partially discussed in [Orfanidis, 1996].

About granular synthesis: the most widely treated case is (*asynchronous granular synthesis*), where simple grains are distributed irregularly. A classic introduction to the topic is [Roads, 1991]. In particular, figure 2.5 in this chapter is based on an analogous figure in [Roads, 1991]. In another classic work, Truax [Truax, 1988] describe the granulation of recorded waveforms.

Additive synthesis was one of the first sound modeling techniques adopted in computer music and has been extensively used in speech applications as well. The main ideas of the synthesis by analysis techniques that we have reviewed date back to the work by McAulay and Quatieri [McAulay and Quatieri, 1986]. In the same period, Smith and Serra started working on “sines-plus-noise” representations, usually termed *SMS* (Spectral Modeling Synthesis) by Serra. A very complete coverage of the topic is provided in [Serra, 1997]. The extension of the additive approach to a “sines-plus-transients-plus-noise” representation is more recent, and has been proposed by Verma and Meng [Verma and Meng, 2000].

About source-filter models: a tutorial about filter design techniques, including normalization approaches that use L^1 , L^2 , and L^∞ norms of the amplitude response, is [Dutilleux, 1998] Introductions to LP analysis techniques and their applications in speech technology can be found in many textbook. See e.g. [Rabiner and Schafer, 1978].

The first paper about applications of FM techniques to sound synthesis was [Chowning, 1973]. It was later reprinted in [Roads and Strawn, 1985].

References

- J. Chowning. The synthesis of complex audio spectra by means of Frequency Modulation. *J. Audio Engin. Soc.*, 21(7), 1973. Reprinted in Roads and Strawn [1985].
- P. Dutilleux. Filters, Delays, Modulations and Demodulations: A Tutorial. In *Proc. COST-G6 Conf. Digital Audio Effects (DAFx-98)*, pages 4–11, Barcelona, 1998.
- R. McAulay and T. F. Quatieri. Speech Analysis/Synthesis Based on a Sinusoidal Speech Model. *IEEE Trans. Acoustics, Speech, and Signal Process.*, 34:744–754, 1986.
- S. J. Orfanidis. *Introduction to Signal Processing*. Prentice Hall, 1996.
- M. Puckette. Max at seventeen. *Computer Music J.*, 26(4):31–43, 2002.
- L. R. Rabiner and R. W. Schafer. *Digital Processing of Speech Signals*. Prentice-Hall, Englewood Cliffs, NJ, 1978.

- C. Roads. Asynchronous granular synthesis. In G. De Poli, A. Piccialli, and C. Roads, editors, *Representations of Musical Signals*, pages 143–186. MIT Press, 1991.
- C. Roads and J. Strawn, editors. *Foundations of Computer Music*. MIT Press, 1985.
- X. Serra. Musical sound modeling with sinusoids plus noise. In C. Roads, S. Pope, A. Piccialli, and G. De Poli, editors, *Musical Signal Processing*, pages 91–122. Swets & Zeitlinger, 1997. <http://www.iaa.upf.es/~xserra/articles/msm/>.
- B. Truax. Real-time granular synthesis with a digital signal processor. *Computer Music J.*, 12(2):14–26, 1988.
- B. Vercoe. Csound: A manual for the audio processing system and supporting programs with tutorials. Technical report, Media Lab, M.I.T., Cambridge, Massachusetts, 1993. Software and Manuals available from <ftp://ftp.maths.bath.ac.uk/pub/dream/>.
- T. S. Verma and T. H. Y. Meng. Extending Spectral Modeling Synthesis with Transient Modeling Synthesis. *Computer Music J.*, 24(2):47–59, 2000.

Contents

2	Sound modeling: signal-based approaches	2.1
2.1	Introduction	2.1
2.2	Signal generators	2.2
2.2.1	Waveform generators	2.3
2.2.1.1	Digital oscillators	2.3
2.2.1.2	Table lookup oscillator	2.3
2.2.1.3	Recurrent sinusoidal signal generators	2.4
2.2.1.4	Amplitude/frequency controlled oscillators	2.4
2.2.1.5	Envelope generators	2.7
2.2.2	Noise generators	2.8
2.2.2.1	Random noise models	2.9
2.3	Time-segment based models	2.11
2.3.1	Wavetable synthesis	2.11
2.3.1.1	Definitions and applications	2.11
2.3.1.2	Tranformations: pitch shift, looping	2.12
2.3.2	Granular synthesis	2.13
2.3.2.1	Sound granulation	2.14
2.3.2.2	Synthetic grains	2.14
2.3.3	Overlap-Add (OLA) methods	2.16
2.3.3.1	Time-domain OLA	2.16
2.3.3.2	Synchronous and pitch-synchronous OLA	2.18
2.4	Spectrum based models	2.20
2.4.1	Sinusoidal model	2.21
2.4.2	Spectral modeling	2.22
2.4.2.1	Deterministic signal component	2.22
2.4.2.2	Time- and frequency-domain implementations	2.23
2.4.3	Synthesis by analysis	2.24
2.4.3.1	Magnitude and Phase Spectra Computation	2.25
2.4.3.2	A sinusoid tracking procedure	2.26
2.4.4	“Sines-plus-noise” models	2.27
2.4.4.1	Stochastic analysis	2.28
2.4.4.2	Stochastic modeling	2.29
2.4.4.3	Resynthesis and modifications	2.30
2.4.5	Sinusoidal description of transients	2.31
2.4.5.1	The DCT domain	2.31
2.4.5.2	Transient analysis and modeling	2.32

2.5	Source-filter models	2.34
2.5.1	Source-filter decompositions	2.34
2.5.1.1	Sources and filters	2.35
2.5.1.2	An example: modal synthesis of percussive sounds	2.36
2.5.2	Speech modeling	2.36
2.5.2.1	Speech production mechanism and models	2.36
2.5.2.2	Formant synthesis	2.38
2.5.3	Linear prediction	2.40
2.5.3.1	Linear prediction equations	2.42
2.5.3.2	Short-time LP analysis	2.43
2.5.3.3	Linear Predictive Coding (LPC)	2.46
2.5.3.4	LP based audio effects	2.47
2.6	Non linear processing models	2.48
2.6.1	Memoryless non linear modelling	2.48
2.6.2	Synthesis by frequency modulation	2.49
2.6.2.1	Simple modulation	2.51
2.6.2.2	Spectra $ f_c \pm k f_m $	2.52
2.6.2.3	Compound modulation	2.53
2.6.2.4	Nested and feedback modulation	2.55
2.6.2.5	Discussion	2.56
2.6.3	Multiplicative synthesis	2.56
2.7	Key concepts	2.57
2.8	Commented bibliography	2.59