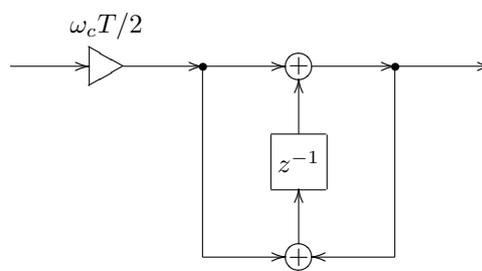


THE ART OF VA FILTER DESIGN



Vadim Zavalishin

rev. 1.1.1 (July 22, 2015)

About this book: the book covers the theoretical and practical aspects of the virtual analog filter design in the music DSP context. Only a basic amount of DSP knowledge is assumed as a prerequisite. For digital musical instrument and effect developers.

Front picture: BLT integrator.

DISCLAIMER: THIS BOOK IS PROVIDED “AS IS”, SOLELY AS AN EXPRESSION OF THE AUTHOR’S BELIEFS AND OPINIONS AT THE TIME OF THE WRITING, AND IS INTENDED FOR THE INFORMATIONAL PURPOSES ONLY.

*To the memory of Elena Golushko,
may her soul travel the happiest path. . .*

Contents

Preface	vii
1 Fourier theory	1
1.1 Complex sinusoids	1
1.2 Fourier series	2
1.3 Fourier integral	3
1.4 Dirac delta function	4
1.5 Laplace transform	5
2 Analog 1-pole filters	7
2.1 RC filter	7
2.2 Block diagrams	8
2.3 Transfer function	9
2.4 Complex impedances	12
2.5 Amplitude and phase responses	13
2.6 Lowpass filtering	14
2.7 Cutoff parametrization	15
2.8 Highpass filter	17
2.9 Poles, zeros and stability	18
2.10 LP to HP substitution	19
2.11 Multimode filter	20
2.12 Shelving filters	21
2.13 Allpass filter	24
2.14 Transposed multimode filter	26
3 Time-discretization	29
3.1 Discrete-time signals	29
3.2 Naive integration	31
3.3 Naive lowpass filter	31
3.4 Block diagrams	32
3.5 Transfer function	34
3.6 Poles	35
3.7 Trapezoidal integration	37
3.8 Bilinear transform	40
3.9 Cutoff prewarping	43
3.10 Zero-delay feedback	44
3.11 Direct forms	49
3.12 Other replacement techniques	52

3.13	Instantaneously unstable feedback	55
4	Ladder filter	61
4.1	Linear analog model	61
4.2	Linear digital model	63
4.3	Feedback shaping	64
4.4	Multimode ladder filter	64
4.5	HP and BP ladders	68
4.6	Simple nonlinear model	70
4.7	Advanced nonlinear model	72
4.8	Diode ladder	73
5	2-pole filters	81
5.1	Linear analog model	81
5.2	Linear digital model	85
5.3	Further filter types	86
5.4	LP to BP/BS substitutions	91
5.5	Nonlinear model	93
5.6	Serial decomposition	95
5.7	Transposed Sallen–Key filters	98
6	Allpass-based effects	107
6.1	Phasers	107
6.2	Flangers	110
7	Frequency shifters	113
7.1	General ideas	113
7.2	Analytic signals	115
7.3	Phase splitter	115
7.4	Implementation structure	117
7.5	Remez algorithm	119
7.6	Cutoff optimization	126
7.7	Analytical construction of phase response	130
7.8	“LP to analytic” substitution	140
7.9	Cutoff prewarping	143
	History	145
	Index	147

Preface

The classical way of presentation of the DSP theory is not very well suitable for the purposes of virtual analog filter design. The linearity and time-invariance of structures are not assumed merely to simplify certain analysis and design aspects, but are handled more or less as an “ultimate truth”. The connection to the continuous-time (analog) world is lost most of the time. The key focus points, particularly the discussed filter types, are of little interest to a digital music instrument developer. This makes it difficult to apply the obtained knowledge in the music DSP context, especially in the virtual analog filter design.

This book attempts to amend this deficiency. The concepts are introduced with the musical VA filter design in mind. The depth of theoretical explanation is restricted to an intuitive and practically applicable amount. The focus of the book is the design of digital models of classical musical analog filter structures using the *topology-preserving transform* approach, which can be considered as a generalization of bilinear transform, zero-delay feedback and trapezoidal integration methods. This results in digital filters having nice amplitude and phase responses, nice time-varying behavior and plenty of options for nonlinearities. In a way, this book can be seen as a detailed explanation of the materials provided in the author’s article “*Preserving the LTI system topology in s- to z-plane transforms.*”

The main purpose of this book is not to explain how to build high-quality emulations of analog hardware (although the techniques explained in the book can be an important and valuable tool for building VA emulations). Rather it is about how to build high-quality time-varying digital filters. The author hopes that these techniques will be used to construct new digital filters, rather than only to build emulations of existing analog structures.

The prerequisites for the reader include familiarity with the basic DSP concepts, complex algebra and the basic ideas of mathematical analysis. Some basic knowledge of electronics may be helpful at one or two places, but is not critical for the understanding of the presented materials.

The author apologizes for possible mistakes and messy explanations, as the book didn’t go through any serious proofreading.

Acknowledgements

The author would like to express his gratitude to a number of people who work (or worked at a certain time) at NI and helped him with the matters related to the creation of this book in one or another way: Daniel Haver, Mate Galic, Tom Kurth, Nicolas Gross, Maïke Weber, Martijn Zwartjes, and Mike Daliot. Special thanks to Stephan Schmitt, Egbert Jürgens, Eike Jonas, Maximilian Zagler, and Marin Vrbica.

The author is also grateful to a number of people on the KVR Audio DSP forum and the music DSP mailing list for productive discussions regarding the matters discussed in the book. Particularly to Martin Eisenberg for the detailed and extensive discussion of the delayless feedback, to Dominique Wurtz for the idea of the full equivalence of different BLT integrators, to the forum member “neotec” for the introduction of the transposed direct form II BLT integrator in the TPT context, to Teemu Voipio for his active involvement into the related discussions and research and to Urs Heckmann for being an active proponent of the ZDF techniques and actually (as far as the author knows) starting the whole avalanche of their usage. Thanks to Robin Schmidt and Richard Hoffmann for reporting a number of mistakes in the book text.

One shouldn’t underestimate the small but invaluable contribution by Helene Kolpakova, whose questions and interest in the VA filter design matters have triggered the initial idea of writing this book. Thanks to Julian Parker for productive discussions, which stimulated the creation of the book’s next revision.

Last, but most importantly, big thanks to Bob Moog for inventing the voltage-controlled transistor ladder filter.

Prior work credits

Various flavors and applications of delayless feedback techniques were in prior use for quite a while. Particularly there are works by A.Härmä, F.Avancini, G.Borin, G.De Poli, F.Fontana, D.Rocchesso, T.Serafini and P.Zamboni, although reportedly this subject has been appearing as far ago as in the 70s of the 20th century.

Chapter 1

Fourier theory

When we are talking about filters we say that filters modify the frequency content of the signal. E.g. a lowpass filter lets the low frequencies through, while suppressing the high frequencies, a highpass filter does vice versa etc. In this chapter we are going to develop a formal definition¹ of the concept of frequencies “contained” in a signal. We will later use this concept to analyse the behavior of the filters.

1.1 Complex sinusoids

In order to talk about the filter theory we need to introduce complex sinusoidal signals. Consider the complex identity:

$$e^{jt} = \cos t + j \sin t \quad (t \in \mathbb{R})$$

(notice that, if t is the time, then the point e^{jt} is simply moving along a unit circle in the complex plane). Then

$$\cos t = \frac{e^{jt} + e^{-jt}}{2}$$

and

$$\sin t = \frac{e^{jt} - e^{-jt}}{2j}$$

Then a real sinusoidal signal $a \cos(\omega t + \varphi)$ where a is the real amplitude and φ is the initial phase can be represented as a sum of two complex conjugate sinusoidal signals:

$$a \cos(\omega t + \varphi) = \frac{a}{2} \left(e^{j(\omega t + \varphi)} + e^{-j(\omega t + \varphi)} \right) = \left(\frac{a}{2} e^{j\varphi} \right) e^{j\omega t} + \left(\frac{a}{2} e^{-j\varphi} \right) e^{-j\omega t}$$

Notice that we have a sum of two complex conjugate sinusoids $e^{\pm j\omega t}$ with respective complex conjugate amplitudes $(a/2)e^{\pm j\varphi}$. So, the complex amplitude simultaneously encodes both the amplitude information (in its absolute magnitude) and the phase information (in its argument). For the positive-frequency component $(a/2)e^{j\varphi} \cdot e^{j\omega t}$, the complex “amplitude” $a/2$ is a half of the real amplitude and the complex “phase” φ is equal to the real phase.

¹More precisely we will develop a number of definitions.

1.2 Fourier series

Let $x(t)$ be a real periodic signal of a period T :

$$x(t) = x(t + T)$$

Let $\omega = 2\pi/T$ be the fundamental frequency of that signal. Then $x(t)$ can be represented² as a sum of a finite or infinite number of sinusoidal signals of harmonically related frequencies $jn\omega$ plus the *DC offset* term³ $a_0/2$:

$$x(t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n \cos(jn\omega t + \varphi_n) \quad (1.1)$$

The representation (1.1) is referred to as *real-form Fourier series*. The respective sinusoidal terms are referred to as the *harmonics* or the harmonic *partials* of the signal.

Using the complex sinusoid notation the same can be rewritten as

$$x(t) = \sum_{n=-\infty}^{\infty} X_n e^{jn\omega t} \quad (1.2)$$

where each harmonic term $a_n \cos(jn\omega t + \varphi_n)$ will be represented by a sum of $X_n e^{jn\omega t}$ and $X_{-n} e^{-jn\omega t}$, where X_n and X_{-n} are mutually conjugate: $X_n = X_{-n}^*$. The representation (1.2) is referred to as *complex-form Fourier series*. Note that we don't have an explicit DC offset partial in this case, it is implicitly contained in the series as the term for $n = 0$.

It can be easily shown that the real- and complex-form coefficients are related as

$$\begin{aligned} X_n &= \frac{a_n}{2} e^{j\varphi_n} & (n > 0) \\ X_0 &= \frac{a_0}{2} \end{aligned}$$

This means that intuitively we can use the absolute magnitude and the argument of X_n (for positive-frequency terms) as the amplitudes and phases of the real Fourier series partials.

Complex-form Fourier series can also be used to represent complex (rather than real) periodic signals in exactly the same way, except that the equality $X_n = X_{-n}^*$ doesn't hold anymore.

Thus, any real periodic signal can be represented as a sum of harmonically related real sinusoidal partials plus the DC offset. Alternatively, any periodic signal can be represented as a sum of harmonically related complex sinusoidal partials.

²Formally speaking, there are some restrictions on $x(t)$. It would be sufficient to require that $x(t)$ is bounded and continuous, except for a finite number of discontinuous jumps per period.

³The reason the DC offset term is notated as $a_0/2$ and not as a_0 has to do with simplifying the math notation in other related formulas.

1.3 Fourier integral

While periodic signals are representable as a sum of a countable number of sinusoidal partials, a nonperiodic real signal can be represented⁴ as a sum of an uncountable number of sinusoidal partials:

$$x(t) = \int_0^{\infty} a(\omega) \cos(\omega t + \varphi(\omega)) \frac{d\omega}{2\pi} \quad (1.3)$$

The representation (1.3) is referred to as *Fourier integral*.⁵ The DC offset term doesn't explicitly appear in this case.

The complex-form version of Fourier integral⁶ is

$$x(t) = \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} \frac{d\omega}{2\pi} \quad (1.4)$$

For real $x(t)$ we have a Hermitian $X(\omega)$: $X(\omega) = X^*(-\omega)$, for complex $x(t)$ there is no such restriction. The function $X(\omega)$ is referred to as *Fourier transform* of $x(t)$.⁷

It can be easily shown that the relationship between the parameters of the real and complex forms of Fourier transform is

$$X(\omega) = \frac{a(\omega)}{2} e^{j\varphi(\omega)} \quad (\omega > 0)$$

This means that intuitively we can use the absolute magnitude and the argument of $X(\omega)$ (for positive frequencies) as the amplitudes and phases of the real Fourier integral partials.

Thus, any timelimited signal can be represented as a sum of an uncountable number of sinusoidal partials of infinitely small amplitudes.

⁴As with Fourier series, there are some restrictions on $x(t)$. It is sufficient to require $x(t)$ to be absolutely integrable, bounded and continuous (except for a finite number of discontinuous jumps per any finite range of the argument value). The most critical requirement here is probably the absolute integrability, which is particularly fulfilled for the timelimited signals.

⁵The $1/2\pi$ factor is typically used to simplify the notation in the theoretical analysis involving the computation. Intuitively, the integration is done with respect to the ordinary, rather than circular frequency:

$$x(t) = \int_0^{\infty} a(f) \cos(2\pi ft + \varphi(f)) df$$

Some texts do not use the $1/2\pi$ factor in this position, in which case it appears in other places instead.

⁶A more common term for (1.4) is *inverse Fourier transform*. However the term *inverse Fourier transform* stresses the fact that $x(t)$ is obtained by computing the inverse of some transform, whereas in this book we are more interested in the fact that $x(t)$ is representable as a combination of sinusoidal signals. The term *Fourier integral* better reflects this aspect. It also suggests a similarity to the Fourier series representation.

⁷The notation $X(\omega)$ for Fourier transform shouldn't be confused with the notation $X(s)$ for Laplace transform. Typically one can be told from the other by the semantics and the notation of the argument. Fourier transform has a real argument, most commonly denoted as ω . Laplace transform has a complex argument, most commonly denoted as s .

1.4 Dirac delta function

The *Dirac delta function* $\delta(t)$ is intuitively defined as a very high and a very short symmetric impulse with a unit area (Fig. 1.1):

$$\delta(t) = \begin{cases} +\infty & \text{if } t = 0 \\ 0 & \text{if } t \neq 0 \end{cases}$$

$$\delta(-t) = \delta(t)$$

$$\int_{-\infty}^{\infty} \delta(t) dt = 1$$

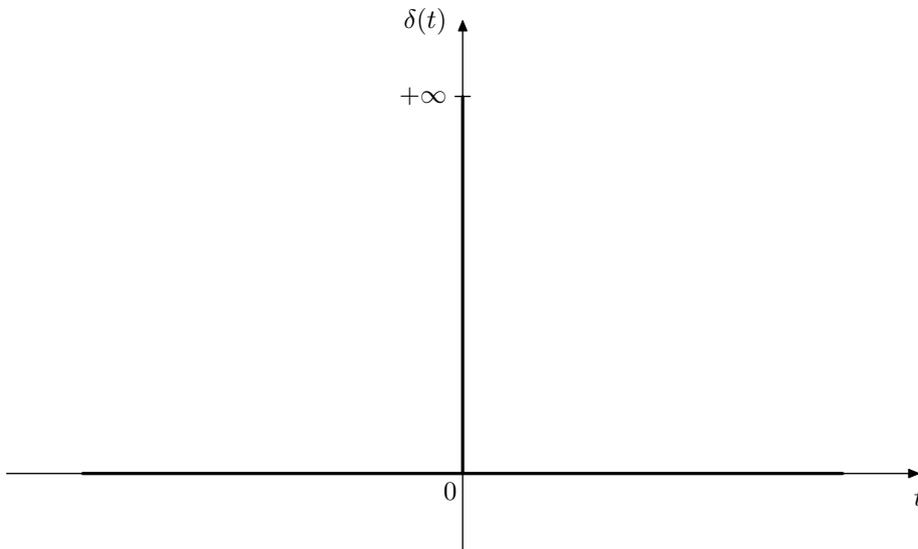


Figure 1.1: Dirac delta function.

Since the impulse is infinitely narrow and since it has a unit area,

$$\int_{-\infty}^{\infty} f(\tau)\delta(\tau) d\tau = f(0) \quad \forall f$$

from where it follows that a convolution of any function $f(t)$ with $\delta(t)$ doesn't change $f(t)$:

$$(f * \delta)(t) = \int_{-\infty}^{\infty} f(\tau)\delta(t - \tau) d\tau = f(t)$$

Dirac delta can be used to represent Fourier series by a Fourier integral. If we let

$$X(\omega) = \sum_{n=-\infty}^{\infty} 2\pi\delta(\omega - n\omega_f)X_n$$

then

$$\sum_{n=-\infty}^{\infty} X_n e^{jn\omega_f t} = \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} \frac{d\omega}{2\pi}$$

From now on, we'll not separately mention Fourier series, assuming that Fourier integral can represent any necessary signal.

Thus, most signals can be represented as a sum of (a possibly infinite number of) sinusoidal partials.

1.5 Laplace transform

Let $s = j\omega$. Then, a complex-form Fourier integral can be rewritten as

$$x(t) = \int_{-j\infty}^{+j\infty} X(s)e^{st} \frac{ds}{2\pi j}$$

where the integration is done in the complex plane along the straight line from $-j\infty$ to $+j\infty$ (apparently $X(s)$ is a different function than $X(\omega)$).⁸ For time-limited signals the function $X(s)$ can be defined on the entire complex plane in such a way that the integration can be done along any line which is parallel to the imaginary axis:

$$x(t) = \int_{\sigma-j\infty}^{\sigma+j\infty} X(s)e^{st} \frac{ds}{2\pi j} \quad (\sigma \in \mathbb{R}) \quad (1.5)$$

In many other cases such $X(s)$ can be defined within some strip $\sigma_1 < \operatorname{Re} s < \sigma_2$. Such function $X(s)$ is referred to as bilateral *Laplace transform* of $x(t)$, whereas the representation (1.5) can be referred to as *Laplace integral*.^{9 10}

Notice that the *complex exponential* e^{st} is representable as

$$e^{st} = e^{\operatorname{Re} s \cdot t} e^{\operatorname{Im} s \cdot t}$$

Considering $e^{\operatorname{Re} s \cdot t}$ as the amplitude of the complex sinusoid $e^{\operatorname{Im} s \cdot t}$ we notice that e^{st} is:

- an exponentially decaying complex sinusoid if $\operatorname{Re} s < 0$,
- an exponentially growing complex sinusoid if $\operatorname{Re} s > 0$,
- a complex sinusoid of constant amplitude if $\operatorname{Re} s = 0$.

Thus, most signals can be represented as a sum of (a possibly infinite number of) complex exponential partials, where the amplitude growth or decay speed of these partials can be relatively arbitrarily chosen.

⁸As already mentioned, the notation $X(\omega)$ for Fourier transform shouldn't be confused with the notation $X(s)$ for Laplace transform. Typically one can be told from the other by the semantics and the notation of the argument. Fourier transform has a real argument, most commonly denoted as ω . Laplace transform has a complex argument, most commonly denoted as s .

⁹A more common term for (1.5) is *inverse Laplace transform*. However the term *inverse Laplace transform* stresses the fact that $x(t)$ is obtained by computing the inverse of some transform, whereas in this book we are more interested in the fact that $x(t)$ is representable as a combination of exponential signals. The term *Laplace integral* better reflects this aspect.

¹⁰The representation of periodic signals by Laplace integral (using Dirac delta function) is problematic for $\sigma \neq 0$. Nevertheless, we can represent them by a Laplace integral if we restrict σ to $\sigma = 0$ (that is $\operatorname{Re} s = 0$ for $X(s)$).

SUMMARY

The most important conclusion of this chapter is: any signal occurring in practice can be represented as a sum of sinusoidal (real or complex) components. The frequencies of these sinusoids can be referred to as the “frequencies contained in the signal”. For complex representation, the real amplitude and phase information is encoded in the absolute magnitude and the argument of the complex amplitudes of the positive-frequency partials (where the absolute magnitude of the complex amplitude is a half of the real amplitude).

It is also possible to use complex exponentials instead of sinusoids.

Chapter 2

Analog 1-pole filters

In this chapter we are going to introduce the basic analog RC-filter and use it as an example to develop the key concepts of the analog filter analysis.

2.1 RC filter

Consider the circuit in Fig. 2.1, where the voltage $x(t)$ is the input signal and the capacitor voltage $y(t)$ is the output signal. This circuit represents the simplest 1-pole *lowpass filter*, which we are now going to analyse.

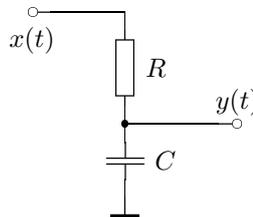


Figure 2.1: A simple RC lowpass filter.

Writing the equations for that circuit we have:

$$\begin{aligned}x &= U_R + U_C \\y &= U_C \\U_R &= RI \\I &= \dot{q}_C \\q_C &= CU_C\end{aligned}$$

where U_R is the resistor voltage, U_C is the capacitor voltage, I is the current through the circuit and q_C is the capacitor charge. Reducing the number of variables, we can simplify the equation system to:

$$x = RC\dot{y} + y$$

or

$$\dot{y} = \frac{1}{RC}(x - y)$$

or, integrating with respect to time:

$$y = y(t_0) + \int_{t_0}^t \frac{1}{RC} (x(\tau) - y(\tau)) d\tau$$

where t_0 is the *initial time moment*. Introducing the notation $\omega_c = 1/RC$ we have

$$y = y(t_0) + \int_{t_0}^t \omega_c (x(\tau) - y(\tau)) d\tau \quad (2.1)$$

We will reintroduce ω_c later as the *cutoff* of the filter.

Notice that we didn't factor $1/RC$ (or ω_c) out of the integral for the case when the value of R is varying with time. The varying R corresponds to the varying cutoff of the filter, and this situation is highly typical in the music DSP context.¹

2.2 Block diagrams

The integral equation (2.1) can be expressed in the block diagram form (Fig. 2.2).

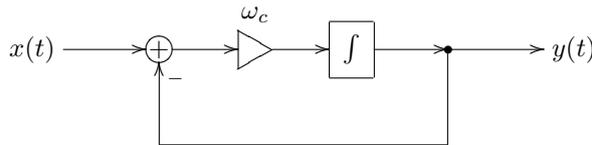


Figure 2.2: A 1-pole RC lowpass filter in the block diagram form.

The meaning of the elements of the diagram should be intuitively clear. The *gain element* (represented by a triangle) multiplies the input signal by ω_c . Notice the inverting input of the summator, denoted by “-”. The integrator simply integrates the input signal:

$$output(t) = output(t_0) + \int_{t_0}^t input(\tau) d\tau$$

The representation of the system by the integral (rather than differential) equation and the respective usage of the integrator element in the block diagram has an important intuitive meaning. Intuitively, the capacitor integrates the current flowing through it, accumulating it as its own charge:

$$q_C(t) = q_C(t_0) + \int_{t_0}^t I(\tau) d\tau$$

or, equivalently

$$U_C(t) = U_C(t_0) + \frac{1}{C} \int_{t_0}^t I(\tau) d\tau$$

One can observe from Fig. 2.2 that the output signal is always trying to “reach” the input signal. Indeed, the difference $x - y$ is always “directed” from

¹We didn't assume the varying C because then our simplification of the equation system doesn't hold anymore, since $\dot{q}_C \neq C\dot{U}_C$ in this case.

y to x . Since $\omega_c > 0$, the integrator will respectively increase or decrease its output value in the respective direction. This corresponds to the fact that the capacitor voltage in Fig. 2.1 is always trying to reach the input voltage. Thus, the circuit works as a kind of smoother of the input signal.

2.3 Transfer function

Consider the integrator:

$$x(t) \longrightarrow \boxed{\int} \longrightarrow y(t)$$

Suppose $x(t) = e^{st}$ (where $s = j\omega$ or, possibly, another complex value). Then

$$y(t) = y(t_0) + \int_{t_0}^t e^{s\tau} d\tau = y(t_0) + \frac{1}{s} e^{s\tau} \Big|_{\tau=t_0}^t = \frac{1}{s} e^{st} + \left(y(t_0) - \frac{1}{s} e^{st_0} \right)$$

Thus, a complex sinusoid (or exponential) e^{st} sent through an integrator comes out as the same signal e^{st} just with a different amplitude $1/s$ plus some DC term $y(t_0) - e^{st_0}/s$. Similarly, a signal $X(s)e^{st}$ (where $X(s)$ is the complex amplitude of the signal) comes out as $(X(s)/s)e^{st}$ plus some DC term. That is, if we forget about the extra DC term, *the integrator simply multiplies the amplitudes of complex exponential signals e^{st} by $1/s$.*

Now, the good news is: for our purposes of filter analysis we can simply *forget* about the extra DC term. The reason for this is the following. Suppose the initial time moment t_0 was quite long ago ($t_0 \ll 0$). Suppose further that the integrator is contained in a *stable* filter (we will discuss the filter stability later, for now we'll simply mention that we're mostly interested in the stable filters for the purposes of the current discussion). It can be shown that in this case the effect of the extra DC term on the output signal is negligible. Since the initial state $y(t_0)$ is incorporated into the same DC term, it also means that the effect of the initial state is negligible!²

Thus, we simply write (for an integrator):

$$\int e^{s\tau} d\tau = \frac{1}{s} e^{st}$$

This means that e^{st} is an *eigenfunction* of the integrator with the respective eigenvalue $1/s$.

Since the integrator is linear,³ not only are we able to factor $X(s)$ out of the integration:

$$\int X(s)e^{s\tau} d\tau = X(s) \int e^{s\tau} d\tau = \frac{1}{s} X(s)e^{st}$$

²In practice, typically, a zero initial state is assumed. Then, particularly, in the case of absence of the input signal, the output signal of the filter is zero from the very beginning (rather than for $t \gg t_0$).

³The linearity here is understood in the sense of the operator linearity. An operator \hat{H} is linear, if

$$\hat{H}(\lambda_1 f_1(t) + \lambda_2 f_2(t)) = \lambda_1 \hat{H}f_1(t) + \lambda_2 \hat{H}f_2(t)$$

but we can also apply the integration independently to all Fourier (or Laplace) partials of an arbitrary signal $x(t)$:

$$\begin{aligned} \int \left(\int_{\sigma-j\infty}^{\sigma+j\infty} X(s)e^{s\tau} \frac{ds}{2\pi j} \right) d\tau &= \int_{\sigma-j\infty}^{\sigma+j\infty} \left(\int X(s)e^{s\tau} d\tau \right) \frac{ds}{2\pi j} = \\ &= \int_{\sigma-j\infty}^{\sigma+j\infty} \frac{X(s)}{s} e^{s\tau} \frac{ds}{2\pi j} \end{aligned}$$

That is, the integrator changes the complex amplitude of each partial by a $1/s$ factor.

Consider again the structure in Fig. 2.2. Assuming the input signal $x(t)$ has the form e^{st} we can replace the integrator by a gain element with a $1/s$ factor. We symbolically reflect this by replacing the integrator symbol in the diagram with the $1/s$ fraction (Fig. 2.3).⁴

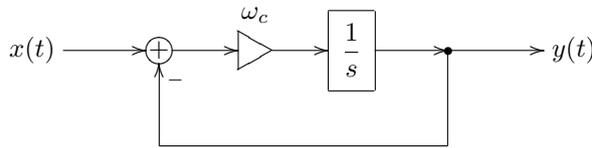


Figure 2.3: A 1-pole RC lowpass filter in the block diagram form with a $1/s$ notation for the integrator.

So, suppose $x(t) = X(s)e^{st}$ and suppose we know $y(t)$. Then the input signal for the integrator is $\omega_c(x - y)$. We now will further take for granted the knowledge that $y(t)$ will be the same signal e^{st} with some different complex amplitude $Y(s)$, that is $y(t) = Y(s)e^{st}$ (notably, this holds only if ω_c is constant, that is, if the system is *time-invariant!!!*)⁵ Then the input signal of the integrator is $\omega_c(X(s) - Y(s))e^{st}$ and the integrator simply multiplies its amplitude by $1/s$. Thus the output signal of the integrator is $\omega_c(x - y)/s$. But, on the other hand $y(t)$ is the output signal of the integrator, thus

$$y(t) = \omega_c \frac{x(t) - y(t)}{s}$$

or

$$Y(s)e^{st} = \omega_c \frac{X(s) - Y(s)}{s} e^{st}$$

or

$$Y(s) = \omega_c \frac{X(s) - Y(s)}{s}$$

from where

$$sY(s) = \omega_c X(s) - \omega_c Y(s)$$

⁴Often in such cases the input and output signal notation for the block diagram is replaced with $X(s)$ and $Y(s)$. Such diagram then “works” in terms of Laplace transform, the input of the diagram is the Laplace transform $X(s)$ of the input signal $x(t)$, the output is respectively the Laplace transform $Y(s)$ of the output signal $y(t)$. The integrators can then be seen as s -dependent gain elements, where the gain coefficient is $1/s$.

⁵In other words, we take for granted the fact that e^{st} is an eigenfunction of the entire circuit.

and

$$Y(s) = \frac{\omega_c}{s + \omega_c} X(s)$$

Thus, the circuit in Fig. 2.3 (or in Fig. 2.2) simply scales the amplitude of the input sinusoidal (or exponential) signal $X(s)e^{st}$ by the $\omega_c/(s + \omega_c)$ factor.

Let's introduce the notation

$$H(s) = \frac{\omega_c}{s + \omega_c} \quad (2.2)$$

Then

$$Y(s) = H(s)X(s)$$

$H(s)$ is referred to as the *transfer function* of the structure in Fig. 2.3 (or Fig. 2.2). Notice that $H(s)$ is a complex function of a complex argument.

For an arbitrary input signal $x(t)$ we can use the Laplace transform representation

$$x(t) = \int_{\sigma-j\infty}^{\sigma+j\infty} X(s)e^{st} \frac{ds}{2\pi j}$$

From the *linearity*⁶ of the circuit in Fig. 2.3, it follows that the result of the application of the circuit to a linear combination of some signals is equal to the linear combination of the results of the application of the circuit to the individual signals. That is, for each input signal of the form $X(s)e^{st}$ we obtain the output signal $H(s)X(s)e^{st}$. Then for an input signal which is an integral sum of $X(s)e^{st}$, we obtain the output signal which is an integral sum of $H(s)X(s)e^{st}$. That is

$$y(t) = \int_{\sigma-j\infty}^{\sigma+j\infty} H(s)X(s)e^{st} \frac{ds}{2\pi j}$$

So, the circuit in Fig. 2.3 independently modifies the complex amplitudes of the sinusoidal (or exponential) partials e^{st} by the $H(s)$ factor!

Notably, the transfer function can be introduced for any system which is linear and time-invariant. For the systems, whose block diagrams consist of integrators, summators and fixed gains, the transfer function is always a *non-strictly proper*⁷ rational function of s . Particularly, this holds for the electronic circuits, where the differential elements are capacitors and inductors, since these types of elements logically perform integration (capacitors integrate the current to obtain the voltage, while inductors integrate the voltage to obtain the current).

It is important to realize that in the derivation of the transfer function concept we used the linearity and time-invariance (the absence of parameter modulation) of the structure. If these properties do not hold, the transfer function can't be introduced! This means that all transfer function-based analysis holds only in the case of fixed parameter values. In practice, if the parameters are not changing too quickly, one can assume that they are approximately constant

⁶Here we again understand the linearity in the operator sense:

$$\hat{H}(\lambda_1 f_1(t) + \lambda_2 f_2(t)) = \lambda_1 \hat{H}f_1(t) + \lambda_2 \hat{H}f_2(t)$$

The operator here corresponds to the circuit in question: $y(t) = \hat{H}x(t)$ where $x(t)$ and $y(t)$ are the input and output signals of the circuit.

⁷A rational function is nonstrictly proper, if the order of its numerator doesn't exceed the order of its denominator.

during certain time range. That is we can “approximately” apply the transfer function concept (and the discussed later derived concepts, such as amplitude and phase responses, poles and zeros, stability criterion etc.) if the modulation of the parameter values is “not too fast”.

2.4 Complex impedances

Actually, we could have obtained the transfer function of the circuit in Fig. 2.1 using the concept of *complex impedances*.

Consider the capacitor equation:

$$I = C\dot{U}$$

If

$$\begin{aligned} I(t) &= I(s)e^{st} \\ U(t) &= U(s)e^{st} \end{aligned}$$

(where $I(t)$ and $I(s)$ are obviously two different functions, the same for $U(t)$ and $U(s)$), then

$$\dot{U} = sU(s)e^{st} = sU(t)$$

and thus

$$I(t) = I(s)e^{st} = C\dot{U} = CsU(s)e^{st} = sCU(t)$$

that is

$$I = sCU$$

or

$$U = \frac{1}{sC}I$$

Now the latter equation looks almost like Ohm’s law for a resistor: $U = RI$. The complex value $1/sC$ is called the *complex impedance* of the capacitor. The same equation can be written in the Laplace transform form: $U(s) = (1/sC)I(s)$.

For an inductor we have $U = L\dot{I}$ and respectively, for $I(t) = I(s)e^{st}$ and $U(t) = U(s)e^{st}$ we obtain $U(t) = sLI(t)$ or $U(s) = sLI(s)$. Thus, the complex impedance of the inductor is sL .

Using the complex impedances as if they were resistances (which we can do, assuming the input signal has the form $X(s)e^{st}$), we simply write the voltage division formula for the circuit in in Fig. 2.1:

$$y(t) = \frac{U_C}{U_R + U_C}x(t)$$

or, cancelling the common current factor $I(t)$ from the numerator and the denominator, we obtain the impedances instead of voltages:

$$y(t) = \frac{1/sC}{R + 1/sC}x(t)$$

from where

$$H(s) = \frac{y(t)}{x(t)} = \frac{1/sC}{R + 1/sC} = \frac{1}{1 + sRC} = \frac{1/RC}{s + 1/RC} = \frac{\omega_c}{s + \omega_c}$$

which coincides with (2.2).

2.5 Amplitude and phase responses

Consider again the structure in Fig. 2.3. Let $x(t)$ be a real signal and let

$$x(t) = \int_{\sigma-j\infty}^{\sigma+j\infty} X(s)e^{st} \frac{ds}{2\pi j}$$

be its Laplace integral representation. Let $y(t)$ be the output signal (which is obviously also real) and let

$$y(t) = \int_{\sigma-j\infty}^{\sigma+j\infty} Y(s)e^{st} \frac{ds}{2\pi j}$$

be its Laplace integral representation. As we have shown, $Y(s) = H(s)X(s)$ where $H(s)$ is the transfer function of the circuit.

The respective Fourier integral representation of $x(t)$ is apparently

$$x(t) = \int_{-\infty}^{+\infty} X(j\omega)e^{j\omega t} \frac{d\omega}{2\pi}$$

where $X(j\omega)$ is the Laplace transform $X(s)$ evaluated at $s = j\omega$. The real Fourier integral representation is then obtained as

$$\begin{aligned} a_x(\omega) &= 2 \cdot |X(j\omega)| \\ \varphi_x(\omega) &= \arg X(j\omega) \end{aligned}$$

For $y(t)$ we respectively have^{8 9}

$$\begin{aligned} a_y(\omega) &= 2 \cdot |Y(j\omega)| = 2 \cdot |H(j\omega)X(j\omega)| = |H(j\omega)| \cdot a_x(\omega) \\ \varphi_y(\omega) &= \arg Y(j\omega) = \arg (H(j\omega)X(j\omega)) = \varphi_x(\omega) + \arg H(j\omega) \end{aligned} \quad (\omega \geq 0)$$

Thus, the amplitudes of the real sinusoidal partials are magnified by the $|H(j\omega)|$ factor and their phases are shifted by $\arg H(j\omega)$ ($\omega \geq 0$). The function $|H(j\omega)|$ is referred to as the *amplitude response* of the circuit and the function $\arg H(j\omega)$ is referred to as the *phase response* of the circuit. Note that both the amplitude and the phase response are real functions of a real argument ω .

The complex-valued function $H(j\omega)$ of the real argument ω is referred to as the *frequency response* of the circuit. Simply put, the frequency response is equal to the transfer function evaluated on the imaginary axis.

Since the transfer function concept works only in the linear time-invariant case, so do the concepts of the amplitude, phase and frequency responses!

⁸This relationship holds only if $H(j\omega)$ is Hermitian: $H(j\omega) = H^*(-j\omega)$. If it weren't the case, the Hermitian property wouldn't hold for $Y(j\omega)$ and $y(t)$ couldn't have been a real signal (for a real input $x(t)$). Fortunately, for real systems $H(j\omega)$ is always Hermitian. Particularly, rational transfer functions $H(s)$ with real coefficients obviously result in Hermitian $H(j\omega)$.

⁹Formally, $\omega = 0$ requires special treatment in case of a Dirac delta component at $\omega = 0$ (arising particularly if the Fourier series is represented by a Fourier integral and there is a nonzero DC offset). Nevertheless, the resulting relationship between $a_y(0)$ and $a_x(0)$ is exactly the same as for $\omega > 0$, that is $a_y(0) = H(0)a_x(0)$. A more complicated but same argument holds for the phase.

2.6 Lowpass filtering

Consider again the transfer function of the structure in Fig. 2.2:

$$H(s) = \frac{\omega_c}{s + \omega_c}$$

The respective amplitude response is

$$|H(j\omega)| = \left| \frac{\omega_c}{\omega_c + j\omega} \right|$$

Apparently at $\omega = 0$ we have $H(0) = 1$. On the other hand, as ω grows, the magnitude of the denominator grows as well and the function decays to zero: $H(+j\infty) = 0$. This suggests the lowpass filtering behavior of the circuit: it lets the partials with frequencies $\omega \ll \omega_c$ through and stops the partials with frequencies $\omega \gg \omega_c$. The circuit is therefore referred to as a *lowpass filter*, while the value ω_c is defined as the *cutoff* frequency of the circuit.

It is convenient to plot the amplitude response of the filter in a fully logarithmic scale. The amplitude gain will then be plotted in decibels, while the frequency axis will have a uniform spacing of octaves. For $H(s) = \omega_c/(s + \omega_c)$ the plot looks like the one in Fig. 2.4.

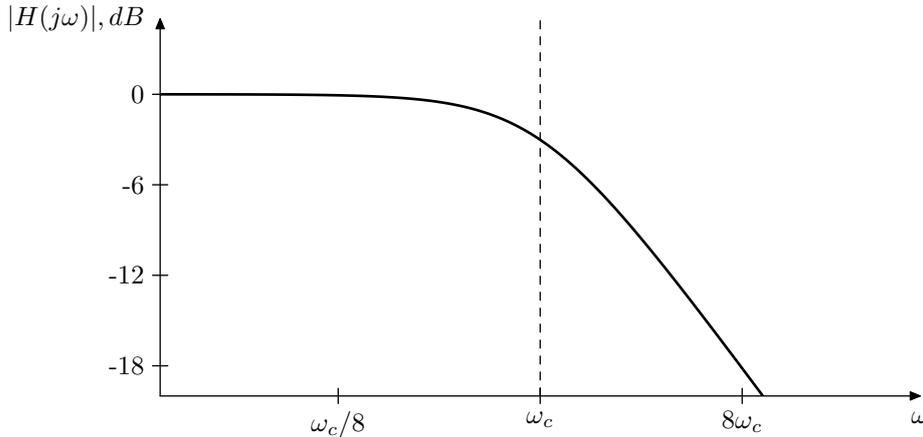


Figure 2.4: Amplitude response of a 1-pole lowpass filter.

Notice that the plot falls off in an almost straight line as $\omega \rightarrow \infty$. Apparently, at $\omega \gg \omega_c$ and respectively $|s| \gg \omega_c$ we have $H(s) \approx \omega_c/s$ and $|H(s)| \approx \omega_c/\omega$. This is a hyperbola in the linear scale and a straight line in a fully logarithmic scale. If ω doubles (corresponding to a step up by one octave), the amplitude gain is approximately halved (that is, drops by approximately 6 decibel). We say that this lowpass filter has a *rolloff* of 6dB/oct.

Another property of this filter is that the amplitude drop at the cutoff is -3 dB. Indeed

$$|H(j\omega_c)| = \left| \frac{\omega_c}{\omega_c + j\omega_c} \right| = \left| \frac{1}{1 + j} \right| = \frac{1}{\sqrt{2}} \approx -3\text{dB}$$

2.7 Cutoff parametrization

Suppose $\omega_c = 1$. Then the lowpass transfer function (2.2) turns into

$$H(s) = \frac{1}{s+1}$$

Now perform the substitution $s \leftarrow s/\omega_c$. We obtain

$$H(s) = \frac{1}{s/\omega_c + 1} = \frac{\omega_c}{s + \omega_c}$$

which is again our familiar transfer function of the lowpass filter.

Consider the amplitude response graph of $1/(s+1)$ in a logarithmic scale. The substitution $s \leftarrow s/\omega_c$ simply shifts this graph to the left or to the right (depending on whether $\omega_c < 1$ or $\omega_c > 1$) without changing its shape. Thus, the variation of the cutoff parameter doesn't change the shape of the amplitude response graph (Fig. 2.5), or of the phase response graph, for that matter (Fig. 2.6).

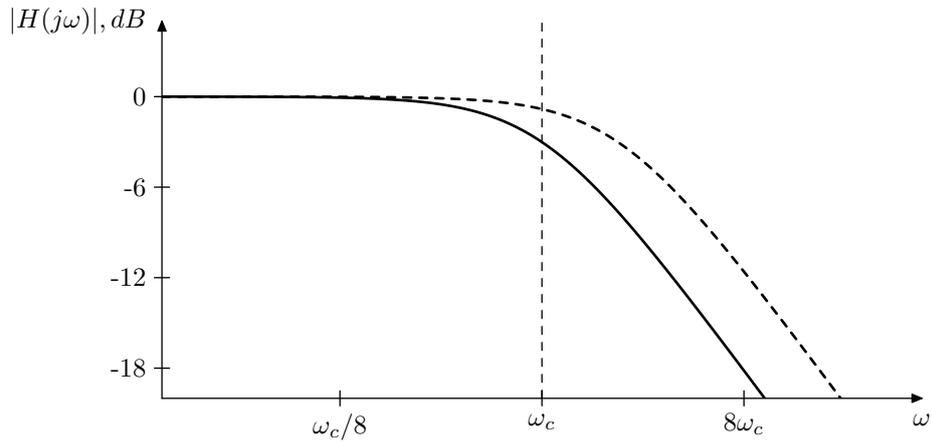


Figure 2.5: 1-pole lowpass filter's amplitude response shift by a cutoff change.

The substitution $s \leftarrow s/\omega_c$ is a generic way to handle cutoff parametrization for analog filters, because it doesn't change the response shapes. This has a nice counterpart on the block diagram level. For all types of filters we simply visually combine an ω_c gain and an integrator into a single block:¹⁰



¹⁰Notice that including the cutoff gain into the integrator makes the integrator block invariant to the choice of the time units:

$$y(t) = y(t_0) + \int_{t_0}^t \omega_c x(\tau) d\tau$$

because the product $\omega_c d\tau$ is invariant to the choice of the time units. This will become important once we start building discrete-time models of filters, where we would often assume unit sampling period.

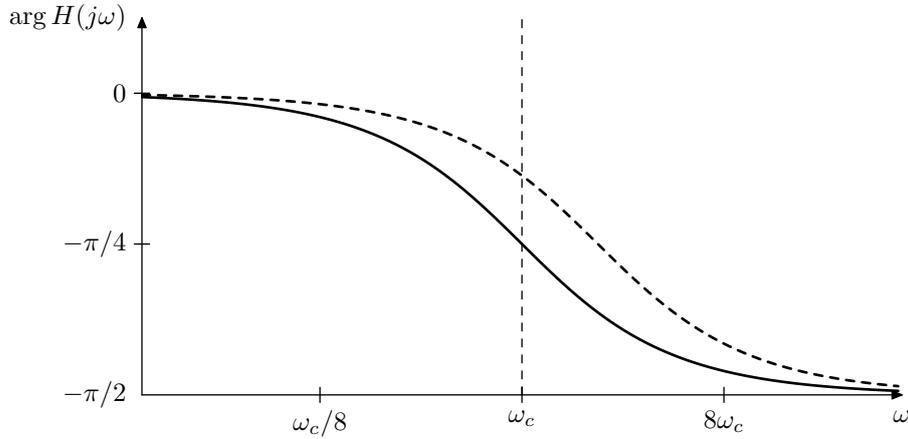
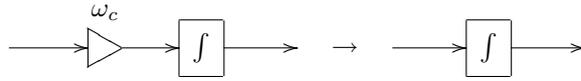


Figure 2.6: 1-pole lowpass filter's phase response shift by a cutoff change.

Apparently, the reason for the ω_c/s notation is that this is the transfer function of the serial connection of an ω_c gain and an integrator. Alternatively, we simply assume that the cutoff gain is contained inside the integrator:



The internal representation of such integrator block is of course still a cutoff gain followed by an integrator. Whether the gain should precede the integrator or follow it may depend on the details of the analog prototype circuit. In the absence of the analog prototype it's better to put the integrator *after* the cutoff gain, because then the integrator will smooth the jumps and further artifacts arising out of the cutoff modulation.

With the cutoff gain implied inside the integrator block, the structure from Fig. 2.2 is further simplified to the one in Fig. 2.7:

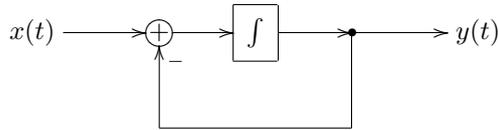


Figure 2.7: A 1-pole RC lowpass filter with an implied cutoff.

As a further shortcut arising out of the just discussed facts, it is common to assume $\omega_c = 1$ during the filter analysis. Particularly, the transfer function of a 1-pole lowpass filter is often written as

$$H(s) = \frac{1}{s + 1}$$

It is assumed that the reader will perform the $s \leftarrow s/\omega_c$ substitution as necessary.

2.8 Highpass filter

If instead of the capacitor voltage in Fig. 2.1 we pick up the resistor voltage as the output signal, we obtain the block diagram representation as in Fig. 2.8.

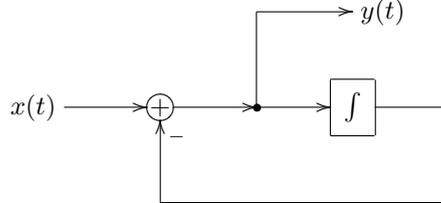


Figure 2.8: A 1-pole highpass filter.

Obtaining the transfer function of this filter we get

$$H(s) = \frac{s}{s + \omega_c}$$

or, in the unit-cutoff form,

$$H(s) = \frac{s}{s + 1}$$

It's easy to see that $H(0) = 0$ and $H(+j\infty) = 1$, whereas the biggest change in the amplitude response occurs again around $\omega = \omega_c$. Thus, we have a *highpass filter* here. The amplitude response of this filter is shown in Fig. 2.9 (in the logarithmic scale).

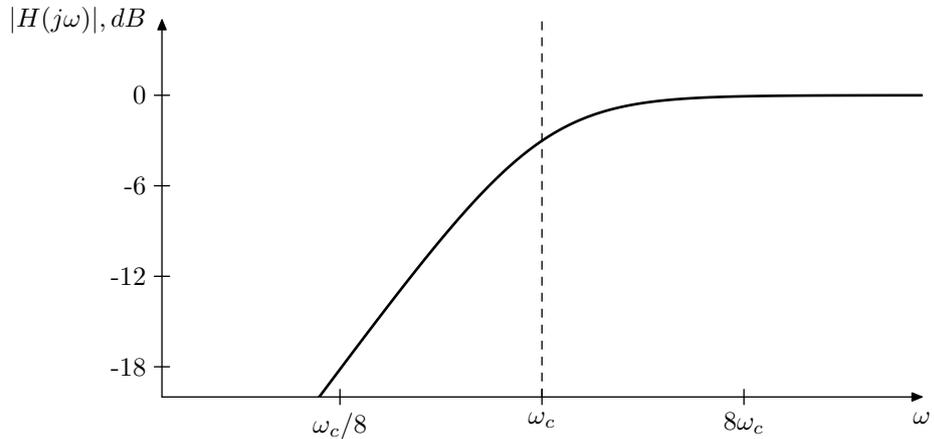


Figure 2.9: Amplitude response of a 1-pole highpass filter.

It's not difficult to observe and not difficult to show that this response is a mirrored version of the one in Fig. 2.4. Particularly, at $\omega \ll \omega_c$ we have $H(s) \approx s/\omega_c$, so when the frequency is halved (dropped by an octave), the amplitude gain is approximately halved as well (drops by approximately 6dB). Again, we have a 6dB/oct rolloff.

2.9 Poles, zeros and stability

Consider the lowpass transfer function:

$$H(s) = \frac{\omega_c}{s + \omega_c}$$

Apparently, this function has a pole in the complex plane at $s = -\omega_c$. Similarly, the highpass transfer function

$$H(s) = \frac{s}{s + \omega_c}$$

also has a pole at $s = -\omega_c$, but it also has a zero at $s = 0$.

Recall that the transfer functions of linear time-invariant differential systems are nonstrictly proper rational functions of s . Thus they always have poles and often have zeros, the numbers of poles and zeros matching the orders of the numerator and the denominator respectively. The poles and zeros of transfer function (especially the poles) play an important role in the filter analysis. For simplicity they are referred to as the poles and zeros of the filters.

The transfer functions of real linear time-invariant differential systems have real coefficients in the numerator and denominator polynomials. Apparently, this doesn't prevent them from having complex poles and zeros, however, being roots of real polynomials, those must come in complex conjugate pairs. E.g. a transfer function with a 3rd order denominator can have either three real poles, or one real and two complex conjugate poles.

The lowpass and highpass filters discussed so far, each have one pole. For that reason they are referred to as 1-pole filters. Actually, the number of poles is always equal to the order of the filter or (which is the same) to the number of integrators in the filter.¹¹ Therefore it is common, instead of e.g. a "4th-order filter" to say a "4-pole filter".

The most important property of the poles is that a filter¹² is *stable* if and only if all its poles are located in the left complex semiplane (that is to the left of the imaginary axis). For our lowpass and highpass filters this is apparently true, as long as $\omega_c > 0$.¹³ If $\omega_c < 0$, the pole is moved to the right semiplane, the filter becomes unstable and will "explode". Also the definition of the frequency response doesn't make much sense in this case. If we put a sinusoidal signal through a stable filter we will (as we have shown) obtain an amplitude-modified and phase-shifted sinusoidal signal of the same frequency.¹⁴ If we put a sinusoidal signal through an unstable filter, the filter simply "explodes" (its

¹¹In certain singular cases, depending on the particular definition details, these numbers might be not equal to each other.

¹²More precisely a linear time-invariant system, which particularly implies fixed parameters. This remark is actually unnecessary, since, as we mentioned, the transfer function (and respectively the poles) are defined only for the linear time-invariant case.

¹³Notably, the same condition ensures the stability of the 1-pole RC lowpass and highpass filters in the time-varying case, which can be directly seen from the fact that the lowpass filter's output never exceeds the maximum level of its input.

¹⁴Strictly speaking, this will happen only after the filter has stabilized itself "to the new signal". This takes a certain amount of time. The closer the poles are to the imaginary axis (from the left), the larger is this stabilization time. The characteristic time value of the stabilization has the order of magnitude of $-1/\max\{\operatorname{Re} p_n\}$, where p_n are the poles. Actually the effects of the transition (occurring at the moment of the appearance of the sinusoidal signal) decay exponentially as $e^{t \max\{\operatorname{Re} p_n\}}$.

output grows infinitely), thus it makes no sense to talk of amplitude and phase responses.

It is also possible to obtain an intuitive understanding of the effect of the pole position on the filter stability. Consider a transfer function of the form

$$H(s) = \frac{F(s)}{\prod_{n=1}^N (s - p_n)}$$

where $F(s)$ is the numerator of the transfer function and p_n are the poles. Suppose all poles are initially in the left complex semiplane and now one of the poles (let's say p_1) starts moving towards the imaginary axis. As the pole gets closer to the axis, the amplitude response at $\omega = \text{Im } p_1$ grows. When p_1 gets onto the axis, the amplitude response at $\omega = \text{Im } p_1$ is infinitely large (since $j\omega = p_1$, we have $H(j\omega) = H(p_1) = \infty$). This corresponds to the filter getting unstable.^{15 16}

The poles and zeros also define the rolloff speed of the amplitude response. Let N_p be the number of poles and N_z be the number of zeros. Since the transfer function must be nonstrictly proper, $N_p \geq N_z$. It's not difficult to see that the amplitude response rolloff at $\omega \rightarrow +\infty$ is $6(N_p - N_z)\text{dB/oct}$. Respectively, the rolloff at $\omega \rightarrow 0$ is $6N_{z0}\text{dB/oct}$, where N_{z0} is the number of zeros at $s = 0$ (provided there are no poles at $s = 0$). Considering that $0 \leq N_{z0} \leq N_z \leq N_p$, the rolloff speed at $\omega \rightarrow +\infty$ or at $\omega \rightarrow 0$ can't exceed $6N_p\text{dB/oct}$. Also, if all zeros of a filter are at $s = 0$ (that is $N_{z0} = N_z$) then the sum of the rolloff speeds at $\omega \rightarrow 0$ and $\omega \rightarrow +\infty$ is exactly $6N_p\text{dB/oct}$.

2.10 LP to HP substitution

The symmetry between the lowpass and the highpass 1-pole amplitude responses has an algebraic explanation. The 1-pole highpass transfer function can be obtained from the 1-pole lowpass transfer function by the *LP to HP* (lowpass to highpass) *substitution*:

$$s \leftarrow 1/s$$

Applying the same substitution to a highpass 1-pole we obtain a lowpass 1-pole. The name "LP to HP substitution" originates from the fact that a number of filters are designed as lowpass filters and then are being transformed to their highpass versions.

Recalling that $s = j\omega$, the respective transformation of the imaginary axis is $j\omega \leftarrow 1/j\omega$ or, equivalently

$$\omega \leftarrow -1/\omega$$

Recalling that the amplitude responses of real systems are symmetric between positive and negative frequencies ($|H(j\omega)| = |H(-j\omega)|$) we can also write

$$\omega \leftarrow 1/\omega \quad (\text{for amplitude response only})$$

¹⁵The reason, why the stable area is the left (and not the right) complex semiplane, falls outside the scope of this book.

¹⁶The discussed 1-pole lowpass filter is actually still kind of stable at $\omega = 0$ (corresponding to the pole at $s = 0$). In fact, it has a constant output level (its state is not changing) in this case. However, strictly speaking, this case is not really stable, since all signals in a truly stable filter must decay to zero in the absence of the input signal.

Taking the logarithm of both sides gives:

$$\log \omega \leftarrow -\log \omega \quad (\text{for amplitude response only})$$

Thus, the amplitude response is flipped around $\omega = 1$ in the logarithmic scale.

The LP to HP substitutions also transforms the filter's poles and zeros by the same formula:

$$s' = 1/s$$

where we substitute pole and zero positions for s . Clearly this transformation maps the complex values in the left semiplane to the values in the left semiplane and the values in the right semiplane to the right semiplane. Thus, the LP to HP substitution exactly preserves the stability of the filters.

The LP to HP substitution can be performed not only algebraically (on a transfer function), but also directly on a block diagram, if we allow the usage of differentiators. Since the differentiator's transfer function is $H(s) = s$, replacing all integrators by differentiators will effectively perform the $1/s \leftarrow s$ substitution, which apparently is the same as the $s \leftarrow 1/s$ substitution. Shall the usage of the differentiators be forbidden, it might still be possible to convert differentiation to the integration by analytical transformations of the equations expressed by the block diagram.

2.11 Multimode filter

Actually, we can pick up the lowpass and highpass signals simultaneously from the same structure (Fig. 2.10). This is referred to as a *multimode filter*.

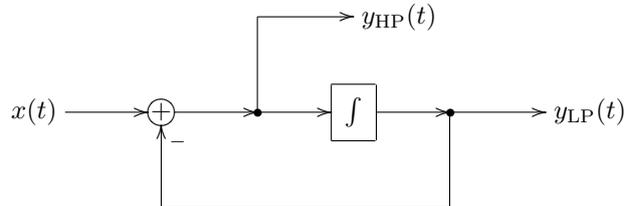


Figure 2.10: A 1-pole multimode filter.

It's easy to observe that $y_{LP}(t) + y_{HP}(t) = x(t)$, that is the input signal is split by the filter into the lowpass and highpass components. In the transfer function form this corresponds to

$$H_{LP}(s) + H_{HP}(s) = \frac{\omega_c}{s + \omega_c} + \frac{s}{s + \omega_c} = 1$$

The multimode filter can be used to implement a 1st-order differential filter for practically any given transfer function, by simply mixing its outputs. Indeed, let

$$H(s) = \frac{b_1 s + b_0}{s + a_0} \quad (a_0 \neq 0)$$

where we can eliminate the case $a_0 = 0$, because it is not defining a stable filter.¹⁷ Letting $\omega_c = a_0$ we obtain

$$H(s) = \frac{b_1 s + b_0}{s + \omega_c} = b_1 \frac{s}{s + \omega_c} + \frac{b_0}{\omega_c} \cdot \frac{\omega_c}{s + \omega_c} = b_1 H_{\text{HP}}(s) + \left(\frac{b_0}{\omega_c}\right) H_{\text{LP}}(s)$$

Thus we simply need to set the filter's cutoff to a_0 and take the sum

$$y = b_1 y_{\text{HP}}(t) + \left(\frac{b_0}{\omega_c}\right) y_{\text{LP}}(t)$$

as the output signal.

2.12 Shelving filters

By adding/subtracting the lowpass-filtered signal to/from the unmodified input signal one can build a low-shelving filter:

$$y(t) = x(t) + K \cdot y_{\text{LP}}(t)$$

The transfer function of the low-shelving filter is respectively:

$$H(s) = 1 + K \frac{1}{s + 1}$$

The amplitude response is plotted Fig. 2.11. Typically $K \geq -1$. At $K = 0$ the signal is unchanged. At $K = -1$ the filter turns into a highpass.

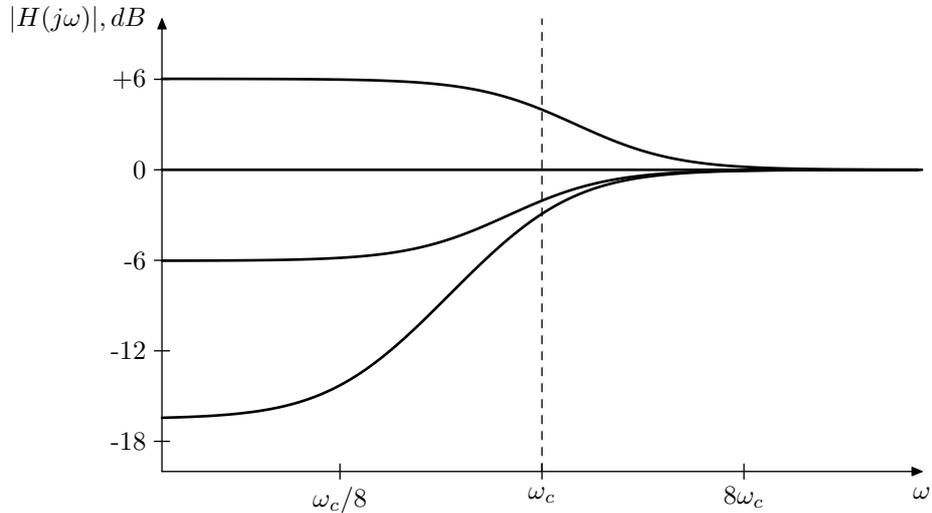


Figure 2.11: Amplitude response of a 1-pole low-shelving filter (for various K).

The high-shelving filter is built in a similar way:

$$y(t) = x(t) + K \cdot y_{\text{HP}}(t)$$

¹⁷If supporting the case $a_0 = 0$ is really desired, it can be done by introducing a gain element into the feedback path of the 1-pole filter.

and

$$H(s) = 1 + K \frac{s}{s+1}$$

The amplitude response is plotted Fig. 2.12.

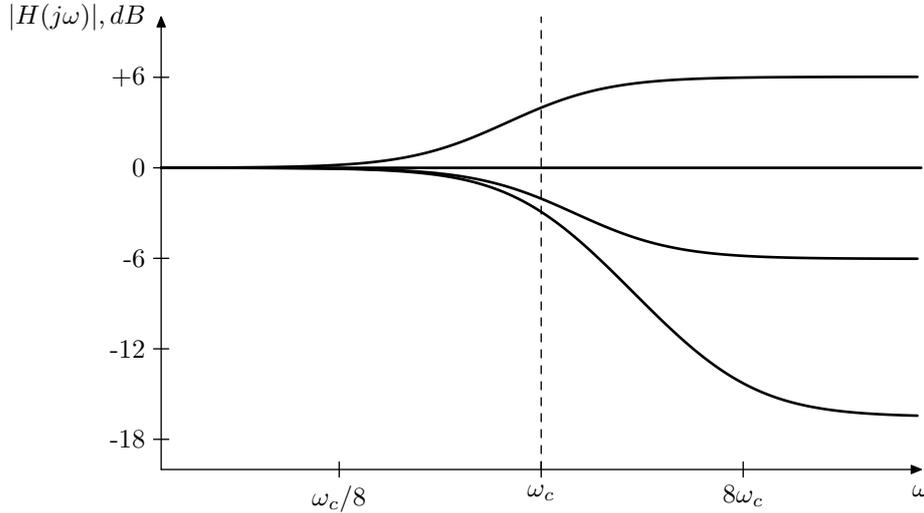


Figure 2.12: Amplitude response of a 1-pole high-shelving filter (for various K).

There are a couple of nontrivial moments here, though. The first one has to do with the fact that the amplitude boost or drop for the “shelf” is more convenient to be specified in decibels. Which requires translation of the level change specified in decibels into the K factor. It’s not difficult to realize that

$$\text{dB} = 20 \log_{10}(K + 1)$$

Indeed, e.g. for the low-shelving filter at $\omega = 0$ (that is $s = 0$) we have¹⁸

$$H(0) = 1 + K$$

We also obtain $H(+j\infty) = 1 + K$ for the high-shelving filter.

A further nontrivial moment is that the definition of the cutoff at $\omega = 1$ for such filters is not really convenient. Indeed, looking at the amplitude response graphs in Figs. 2.11 and 2.12 we would rather wish to have the cutoff point positioned exactly at the middle of the respective slopes. Let’s find where the middle is. E.g. for the lowpass (and remembering that both scales of the graph are logarithmic) we first find the mid-height, which is the geometric average of the shelf’s gain and the unit gain: $\sqrt{1+K}$. Then we need to find ω at which the amplitude response is $\sqrt{1+K}$:

$$\left| 1 + K \frac{1}{j\omega + 1} \right|^2 = \left| \frac{j\omega + 1 + K}{j\omega + 1} \right|^2 = \frac{(1+K)^2 + \omega^2}{1 + \omega^2} = 1 + K$$

¹⁸ $H(0) = 1 + K$ is not a fully trivial result here. We have it only because the lowpass filter doesn’t change the signal’s phase at $\omega = 0$. If instead it had e.g. inverted the phase, then we would have obtained $1 - K$ here.

from where

$$\omega = \sqrt{1 + K}$$

This is the frequency of the midslope point of a low-shelving filter built from a unit-cutoff lowpass. If we want the midslope point to be at $\omega = 1$ then the lowpass cutoff needs to be set to $1/\sqrt{1 + K}$. Other midslope point frequencies are obtained in a similar fashion

$$\omega_c = \frac{\omega_{\text{mid}}}{\sqrt{1 + K}} \quad (\text{low-shelving})$$

The amplitude responses for various K then begin to look like in Fig. 2.13.

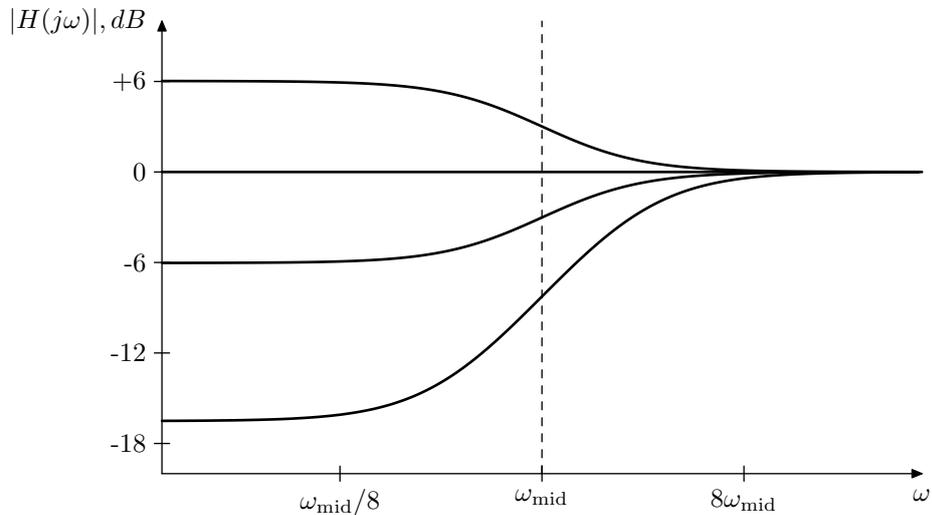


Figure 2.13: Fixed-midpoint amplitude responses of a 1-pole low-shelving filter.

Notably, the low-shelving filter's amplitude response is symmetric around the midpoint in the fully logarithmic scale plot. This can be better illustrated by starting off with a shelving filter transfer function written in a differently scaled way:

$$G(s) = \frac{s + M}{Ms + 1}$$

Clearly, the “LP to HP” substitution applied to $G(s)$ simply reciprocates it, thus $|G(j\omega)|$ is naturally symmetric in the fully logarithmic scale around the point $\omega = 1$, $|G(j)| = 1$ (so, $\omega_{\text{mid}} = 1$ for $G(s)$).

In order to establish the relationship between $G(s)$ and $H(s)$ notice that

$$G(s) = \frac{s + M}{Ms + 1} = \frac{1}{M} \cdot \frac{Ms + M^2}{Ms + 1} = \frac{1}{M} \cdot \frac{(s/\omega_c) + M^2}{(s/\omega_c) + 1}$$

where $\omega_c = 1/M$. Comparing to

$$H(s) = 1 + K \frac{1}{s + 1} = \frac{s + (1 + K)}{s + 1}$$

we have

$$H(s) = M \cdot G(s) \quad \text{where } M = \sqrt{1+K} \text{ and } \omega_c = 1/\sqrt{1+K}$$

For a high-shelving filter there is a similar symmetry. Starting with a reciprocal of the low-shelving filter's $G(s)$:

$$G(s) = \frac{Ms+1}{s+M} = \frac{1}{M} \cdot \frac{M^2(s/M)+1}{(s/M)+1} = \frac{1}{M} \cdot \frac{M^2(s/\omega_c)+1}{(s/\omega_c)+1}$$

Comparing to

$$H(s) = 1 + K \frac{s}{s+1} = \frac{(1+K)s}{s+1}$$

we have

$$H(s) = M \cdot G(s) \quad \text{where } M = \sqrt{1+K} \text{ and } \omega_c = \sqrt{1+K}$$

Respectively, for a non-unity ω_{mid} :

$$\omega_c = \omega_{\text{mid}} \sqrt{1+K} \quad (\text{high-shelving})$$

The fixed-midpoint high-shelving amplitude responses are plotted in Fig. 2.14.

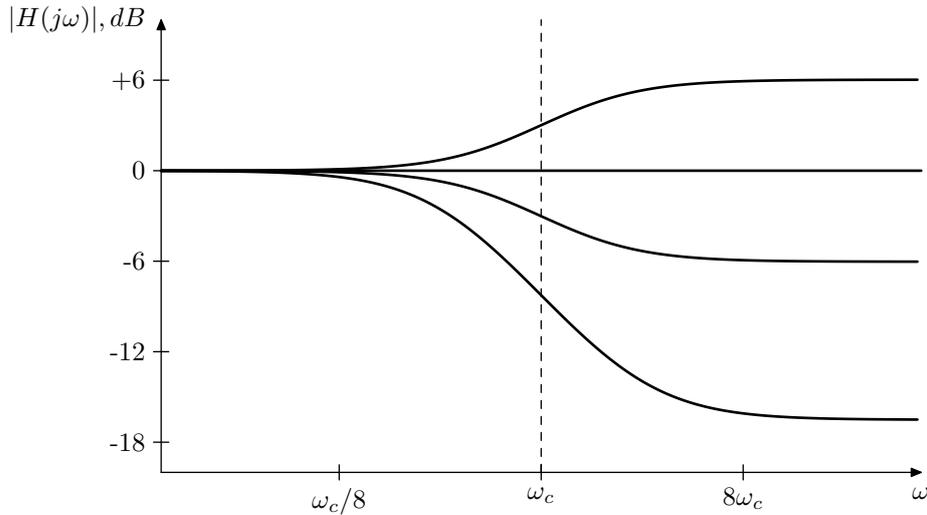


Figure 2.14: Fixed-midpoint amplitude responses of a 1-pole high-shelving filter.

2.13 Allpass filter

By subtracting the highpass output from the lowpass output of the multimode filter we obtain the *allpass filter*:

$$H(s) = H_{\text{LP}}(s) - H_{\text{HP}}(s) = \frac{1}{1+s} - \frac{s}{1+s} = \frac{1-s}{1+s}$$

The amplitude response of the allpass filter is always unity:

$$|H(j\omega)| = 1 \quad \forall \omega$$

Indeed, the numerator $1 - j\omega$ and the denominator $1 + j\omega$ of the frequency response are mutually conjugate, therefore they have equal magnitudes.

The allpass filter is used because of its phase response (Fig. 2.15). That is sometimes we wish to change the phases of the signal's partials without changing their amplitudes. The most common VA use for the allpass filters is probably in phasers.

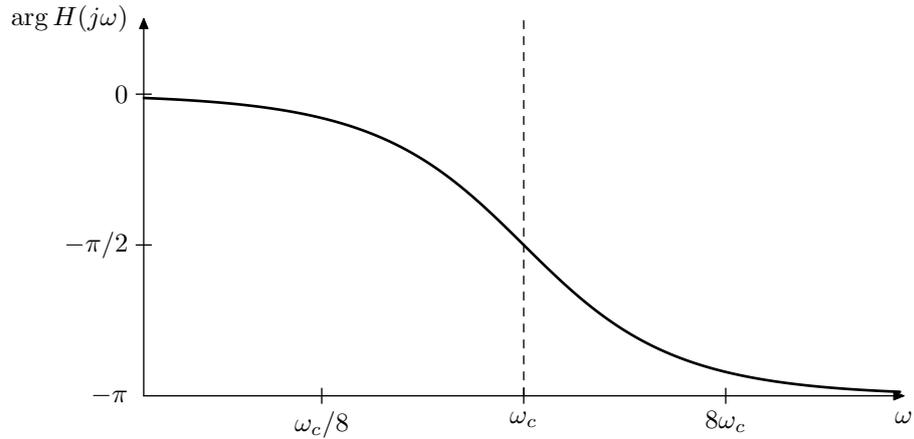


Figure 2.15: Phase response of a 1-pole allpass filter.

We could also subtract the lowpass from the highpass:

$$H(s) = \frac{s}{s+1} - \frac{1}{s+1} = \frac{s-1}{1+s}$$

Apparently the result differs from the previous one only by the inverted phase.

In regards to the unit amplitude response of the 1-pole allpass filter, we could have simply noticed that the zero and the pole of the filter are mutually symmetric relative to the imaginary axis. This is a general property of differential allpass filters: their poles and zeros always come in pairs, located symmetrically relative to the imaginary axis (since the poles of a stable filter have to be in the left complex semiplane, the zeros will be in the right complex semiplane). Expressing the transfer function's numerator and denominator in the multiplicative form, we have

$$|H(s)| = \left| \frac{\prod_{n=1}^N (s - z_n)}{\prod_{n=1}^N (s - p_n)} \right| = \frac{\prod_{n=1}^N |s - z_n|}{\prod_{n=1}^N |s - p_n|}$$

where p_n and z_n are poles and zeros. If each pair p_n and z_n is mutually symmetric relative to the imaginary axis ($p_n = -z_n^*$), then the factors $|j\omega - z_n|$

and $|j\omega - p_n|$ of the amplitude response are always equal, thus the amplitude response is always unity.

The requirement $p_n = -z_n^*$ is not only sufficient but also necessary in order for a differential system to be an allpass. Indeed, let

$$H(j\omega) = g \cdot \frac{P(\omega)}{Q(\omega)} = g \cdot \frac{\prod_{n=1}^N (j\omega - z_n)}{\prod_{n=1}^N (j\omega - p_n)} = g \cdot \frac{\prod_{n=1}^N (\omega + jz_n)}{\prod_{n=1}^N (\omega + jp_n)}$$

where $P(\omega)$ and $Q(\omega)$ are polynomials of ω and g is some unknown coefficient. We also assume that $H(s)$ doesn't have any cancellation between its poles and zeros.

From $|H(\infty)| = 1$ it follows that $|g| = 1$. The allpass property dictates that

$$|P(\omega)| = |Q(\omega)| \quad \forall \omega \in \mathbb{R}$$

or equivalently

$$P(\omega)P^*(\omega) = Q(\omega)Q^*(\omega) \quad \forall \omega \in \mathbb{R} \quad (2.3)$$

Particularly for $\omega = -jz_n$ we have $P(\omega) = 0$, therefore the left-hand side of (2.3) is zero and so the right-hand side must be zero as well, which implies either $Q(\omega) = 0$ or $Q^*(\omega) = 0$. Now, $Q(\omega) = 0$ is impossible, since this would mean that $P(\omega)$ and $Q(\omega)$ have a common root and there is pole/zero cancellation in $H(s)$. Then $\omega = -jz_n$ must be a root of $Q^*(\omega)$. Considering that for $\omega \in \mathbb{R}$

$$Q^*(\omega) = \prod_{n=1}^N (\omega - jp_n^*)$$

the value $\omega = jz_n$ being a root of $Q^*(\omega)$ implies

$$-jz_n - jp_{n'} = 0 \quad \text{for some } n'$$

That is $z_n = -p_{n'}^*$.

2.14 Transposed multimode filter

We could apply the *transposition* to the block diagram in Fig. 2.10. The transposition process is defined as reverting the direction of all signal flow, where forks turn into summatoms and vice versa (Fig. 2.16).¹⁹ The transposition keeps the transfer function relationship within each pair of an input and an output (where the input becomes the output and vice versa). Thus in Fig. 2.16 we have a lowpass and a highpass input and a single output.

The transposed multimode filter has less practical use than the nontransposed one in Fig. 2.10. However, one particular usage case is feedback shaping. Imagine we are mixing an input signal $x_{\text{in}}(t)$ with a feedback signal $x_{\text{fbk}}(t)$, and

¹⁹The inverting input of the summator in the transposed version was obtained from the respective inverting input of the summator in the non-transposed version as follows. First the inverting input is replaced by an explicit inverting gain element (gain factor -1), then the transposition is performed, then the inverting gain is merged into the new summator.

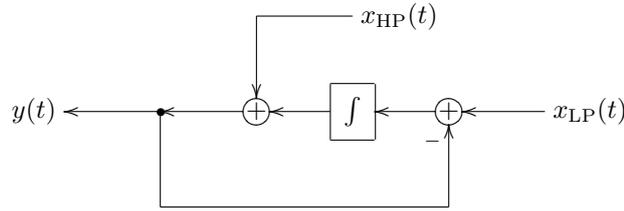


Figure 2.16: A 1-pole transposed multimode filter.

we wish to filter each one of those by a 1-pole filter, and the cutoffs of these 1-pole filters are identical. That is, the transfer functions of those filters share a common denominator. Then we could use a single transposed 1-pole multimode filter as in Fig. 2.17.

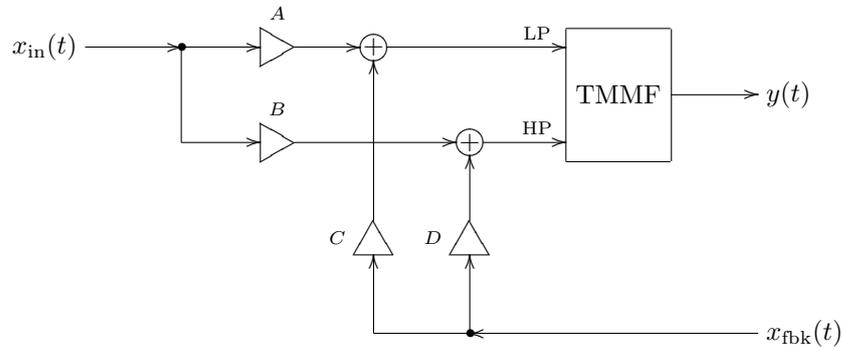


Figure 2.17: A transposed multimode filter (TMMF) used for feedback signal mixing.

The mixing coefficients A , B , C and D will define the numerators of the respective two transfer functions (in exactly the same way as we have been mixing the outputs of a nontransposed multimode filter), whereas the denominator will be $s + \omega_c$, where ω_c is the cutoff of the transposed multimode filter.

SUMMARY

The analog 1-pole filter implementations are built around the idea of the multimode 1-pole filter in Fig. 2.10. The transfer functions of the lowpass and highpass 1-pole filters are

$$H_{\text{LP}}(s) = \frac{\omega_c}{s + \omega_c}$$

and

$$H_{\text{HP}}(s) = \frac{s}{s + \omega_c}$$

respectively. Other 1-pole filter types can be built by combining the lowpass and the highpass signals.

Chapter 3

Time-discretization

Now that we have introduced the basic ideas of analog filter analysis, we will develop an approach to convert analog filter models to the discrete time.

3.1 Discrete-time signals

The discussion of the basic concepts of discrete-time signal representation and processing is outside the scope of this book. We are assuming that the reader is familiar with the basic concepts of discrete-time signal processing, such as sampling, sampling rate, sampling period, Nyquist frequency, analog-to-digital and digital-to-analog signal conversion. However we are going to make some remarks in this respect.

As many other texts do, we will use the square bracket notation to denote discrete-time signals and round parentheses notation to denote continuous-time signals: e.g. $x[n]$ and $x(t)$.

We will often assume a unit sampling rate $f_s = 1$ (and, respectively, a unit sampling period $T = 1$), which puts the Nyquist frequency at $1/2$, or, in the circular frequency terms, at π . Apparently, this can be achieved simply by a corresponding choice of time units.

Theoretical DSP texts typically state that discrete-time signals have periodic frequency spectra. This might be convenient for certain aspects of theoretical analysis such as analog-to-digital and digital-to-analog signal conversion, but it's highly unintuitive otherwise. It would be more intuitive, whenever talking of a discrete-time signal, to imagine an ideal DAC connected to this signal, and think that the discrete-time signal represents the respective continuous-time signal produced by such DAC. Especially, since by sampling this continuous-time signal we obtain the original discrete-time signal again. So the DAC and ADC conversions are exact inverses of each other (in this case). Now, the continuous-time signal produced by such DAC doesn't contain any partials above the Nyquist frequency. Thus, its Fourier integral representation (assuming $T = 1$) is

$$x[n] = \int_{-\pi}^{\pi} X(\omega) e^{j\omega n} \frac{d\omega}{2\pi}$$

and its Laplace integral representation is

$$x[n] = \int_{\sigma-j\pi}^{\sigma+j\pi} X(s)e^{sn} \frac{ds}{2\pi j}$$

Introducing notation $z = e^s$ and noticing that

$$ds = d(\log z) = \frac{dz}{z}$$

we can rewrite the Laplace integral as

$$x[n] = \oint X(z)z^n \frac{dz}{2\pi jz}$$

(where $X(z)$ is apparently a different function than $X(s)$) where the integration is done counterclockwise along a circle of radius e^σ centered at the complex plane's origin:¹

$$z = e^s = e^{\sigma+j\omega} = e^\sigma \cdot e^{j\omega} \quad (-\pi \leq \omega \leq \pi) \quad (3.1)$$

We will refer the representation (3.1) as the *z-integral*.² The function $X(z)$ is referred to as the *z-transform* of $x[n]$.

In case of non-unit sampling period $T \neq 1$ the formulas are the same, except that the frequency-related parameters get multiplied by T (or divided by f_s), or equivalently, the n index gets multiplied by T in continuous-time expressions:³

$$x[n] = \int_{-\pi f_s}^{\pi f_s} X(\omega)e^{j\omega Tn} \frac{d\omega}{2\pi}$$

$$x[n] = \int_{\sigma-j\pi f_s}^{\sigma+j\pi f_s} X(s)e^{sTn} \frac{ds}{2\pi j}$$

$$z = e^{sT}$$

$$x[n] = \oint X(z)z^n \frac{dz}{2\pi jz} \quad (z = e^{\sigma+j\omega T}, -\pi f_s \leq \omega \leq \pi f_s)$$

The notation z^n is commonly used for discrete-time complex exponential signals. A continuous-time signal $x(t) = e^{st}$ is written as $x[n] = z^n$ in discrete-time, where $z = e^{sT}$. The Laplace-integral amplitude coefficient $X(s)$ in $X(s)e^{st}$ then may be replaced by a *z-integral* amplitude coefficient $X(z)$ such as in $X(z)z^n$.

¹As with Laplace transform, sometimes there are no restrictions on the radius e^σ of the circle, sometimes there are.

²A more common term for (3.1) is the *inverse z-transform*, but we will prefer the *z-integral* term for the same reason as with Fourier and Laplace integrals.

³Formally the σ parameter of the Laplace integral (and *z-integral*) should have been multiplied by T as well, but it doesn't matter, since this parameter is chosen rather arbitrarily.

3.2 Naive integration

The most “interesting” element of analog filter block diagrams is obviously the integrator. The time-discretization for other elements is trivial, so we should concentrate on building the discrete-time models of the analog integrator.

The continuous-time integrator equation is

$$y(t) = y(t_0) + \int_{t_0}^t x(\tau) d\tau$$

In discrete time we could approximate the integration by a summation of the input samples. Assuming for simplicity $T = 1$, we could have implemented a discrete-time integrator as

$$y[n] = y[n_0 - 1] + \sum_{\nu=n_0}^n x[\nu]$$

We will refer to the above as the *naive* digital integrator.

A pseudocode routine for this integrator could simply consist of an accumulating assignment:

```
// perform one sample tick of the integrator
integrator_output := integrator_output + integrator_input;
```

It takes the current state of the integrator stored in the *integrator_output* variable and adds the current sample’s value of the *integrator_input* on top of that.

In case of a non-unit sampling period $T \neq 1$ we have to multiply the accumulated input values by T :⁴

```
// perform one sample tick of the integrator
integrator_output := integrator_output + integrator_input*T;
```

3.3 Naive lowpass filter

We could further apply this “naive” approach to construct a discrete-time model of the lowpass filter in Fig. 2.2. We will use the naive integrator as a basis for this model.

Let the x variable contain the current input sample of the filter. Considering that the output of the filter in Fig. 2.2 coincides with the output of the integrator, let the y variable contain the integrator state and simultaneously serve as the output sample. As we begin to process the next input sample, the y variable will contain the previous output value. At the end of the processing of the sample (by the filter model) the y variable will contain the new output sample. In this setup, the input value for the integrator is apparently $(x - y)\omega_c$, thus we simply have

```
// perform one sample tick of the lowpass filter
y := y + (x-y)*omega_c;
```

⁴Alternatively, we could, of course, scale the integrator’s output by T , but this is less useful in practice, because the T factor will be usually combined with the cutoff gain factor ω_c preceding the integrator.

(mind that ω_c must have been scaled to the time units corresponding to the unit sample period!)

A naive discrete-time model of the multimode filter in Fig. 2.10 could have been implemented as:

```
// perform one sample tick of the multimode filter
hp := x-lp;
lp := lp + hp*omega_c;
```

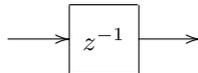
where the integrator state is stored in the lp variable.

The above naive implementations (and any other similar naive implementations, for that matter) work reasonably well as long as $\omega_c \ll 1$, that is the cutoff must be much lower than the sampling rate. At larger ω_c the behavior of the filter becomes rather strange, ultimately the filter gets unstable. We will now develop some theoretical means to analyse the behavior of the discrete-time filter models, figure out what are the problems with the naive implementations, and then introduce another discretization approach.

3.4 Block diagrams

Let's express the naive discrete-time integrator in the form of a discrete-time block diagram. The discrete-time block diagrams are constructed from the same elements as continuous-time block diagrams, except that instead of integrators they have *unit delays*. A unit delay simply delays the signal by one sample. That is the output of a unit delay comes "one sample late" compared to the input. Apparently, the implementation of a unit delay requires a variable, which will be used to store the new incoming value and keep it there until the next sample. Thus, a unit delay element has a *state*, while the other block diagram elements are obviously stateless. This makes the unit delays in a way similar to the integrators in the analog block diagrams, where the integrators are the only elements with a state.

A unit delay element in a block diagram is denoted as:



The reason for the notation z^{-1} will be explained a little bit later. Using a unit delay, we can create a block diagram for our naive integrator (Fig. 3.1). For an arbitrary sampling period we obtain the structure in Fig. 3.2. For an integrator with embedded cutoff gain we can combine the ω_c gain element with the T gain element (Fig. 3.3). Notice that the integrator thereby becomes invariant to the choice of the time units, since $\omega_c T$ is invariant to this choice.

Now let's construct the block diagram of the naive 1-pole lowpass filter. Recalling the implementation routine:

```
// perform one sample tick of the lowpass filter
y := y + (x-y)*omega_c;
```

we obtain the diagram in Fig. 3.4. The z^{-1} element in the feedback from the filter's output to the leftmost summator is occurring due to the fact that we are

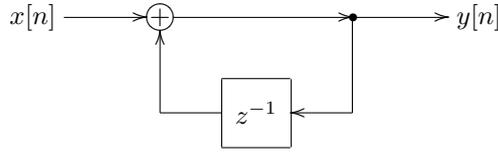


Figure 3.1: Naive integrator for $T = 1$.

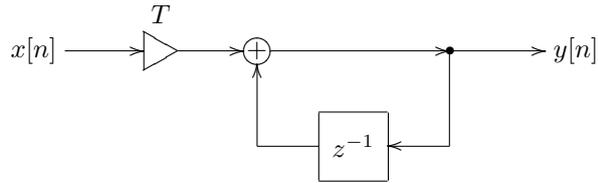


Figure 3.2: Naive integrator for arbitrary T .

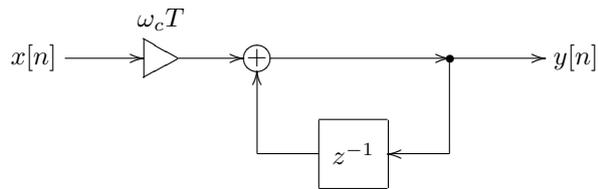


Figure 3.3: Naive integrator with embedded cutoff.

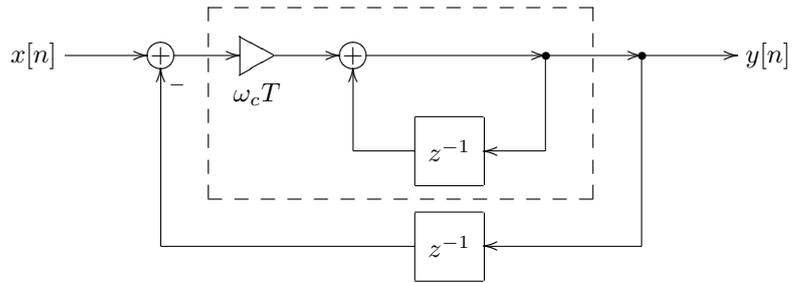


Figure 3.4: Naive 1-pole lowpass filter (the dashed line denotes the integrator).

picking up the *previous* value of y in the routine when computing the difference $x - y$.

This unit delay occurring in the discrete-time feedback is a common problem in discrete-time implementations. This problem is solvable, however it doesn't make too much sense to solve it for the naive integrator-based models, as the increased complexity doesn't justify the improvement in sound. We will address the problem of the zero-delay discrete-time feedback later, for now we'll con-

centrate on the naive model in Fig. 3.4. This model can be simplified a bit, by combining the two z^{-1} elements into one (Fig. 3.5), so that the block diagram explicitly contains a single state variable (as does its pseudocode counterpart).

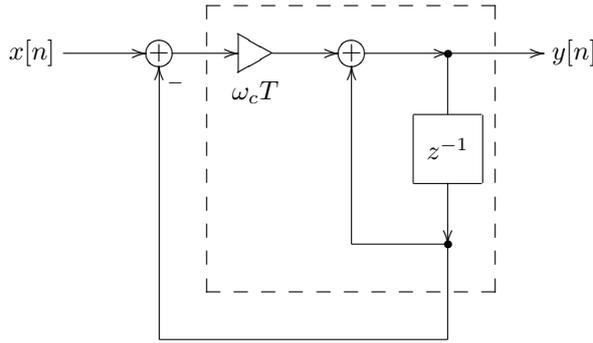
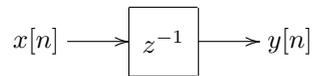


Figure 3.5: Naive 1-pole lowpass filter with just one z^{-1} element (the dashed line denotes the integrator).

3.5 Transfer function

Let $x[n]$ and $y[n]$ be respectively the input and the output signals of a unit delay:



For a complex exponential input $x[n] = e^{sn} = z^n$ we obtain

$$y[n] = e^{s(n-1)} = e^{sn} e^{-s} = z^n z^{-1} = z^{-1} x[n]$$

That is

$$y[n] = z^{-1} x[n]$$

That is, z^{-1} is the *transfer function* of the unit delay! It is common to express discrete-time transfer functions as functions of z rather than functions of s . The reason is that in this case the transfer functions are nonstrictly proper⁵ rational functions, similarly to the continuous-time case, which is pretty convenient. So, for a unit delay we could write $H(z) = z^{-1}$.

Now we can obtain the transfer function of the naive integrator in Fig. 3.1. Suppose⁶ $x[n] = X(z)z^n$ and $y[n] = Y(z)z^n$, or shortly, $x = X(z)z^n$ and $y = Y(z)z^n$. Then the output of the z^{-1} element is yz^{-1} . The output of the summator is then $x + yz^{-1}$, thus

$$y = x + yz^{-1}$$

⁵Under the assumption of causality, which holds if the system is built of unit delays.

⁶As in continuous-time case, we take for granted the fact that complex exponentials z^n are eigenfunctions of discrete-time linear time-invariant systems.

from where

$$y(1 - z^{-1}) = x$$

and

$$H(z) = \frac{y}{x} = \frac{1}{1 - z^{-1}}$$

This is the transfer function of the naive integrator (for $T = 1$).

It is relatively common to express discrete-time transfer functions as rational functions of z^{-1} (like the one above) rather than rational functions of z . However, for the purposes of the analysis it is also often convenient to have them expressed as rational functions of z (particularly, for finding their poles and zeros). We can therefore multiply the numerator and the denominator of the above $H(z)$ by z , obtaining:

$$H(z) = \frac{z}{z - 1}$$

Since $z = e^s$, the *frequency response* is obtained as $H(e^{j\omega})$. The amplitude and phase responses are $|H(e^{j\omega})|$ and $\arg H(e^{j\omega})$ respectively.⁷

For $T \neq 1$ we obtain

$$H(z) = T \frac{z}{z - 1}$$

and, since $z = e^{sT}$, the frequency response is $H(e^{j\omega T})$.

Now let's obtain the transfer function of the naive 1-pole lowpass filter in Fig. 3.5, where, for the simplicity of notation, we assume $T = 1$. Assuming complex exponentials $x = X(z)z^n$ and $y = Y(z)z^n$ we have x and yz^{-1} as the inputs of the first summator. Respectively the integrator's input is $\omega_c(x - yz^{-1})$. And the integrator output is the sum of yz^{-1} and the integrator's input. Therefore

$$y = yz^{-1} + \omega_c(x - yz^{-1})$$

From where

$$(1 - (1 - \omega_c)z^{-1})y = \omega_c x$$

and

$$H(z) = \frac{y}{x} = \frac{\omega_c}{1 - (1 - \omega_c)z^{-1}} = \frac{\omega_c z}{z - (1 - \omega_c)}$$

The transfer function for $T \neq 1$ can be obtained by simply replacing ω_c by $\omega_c T$.

The respective amplitude response is plotted in Fig. 3.6. Comparing it to the amplitude response of the analog prototype we can observe serious deviation closer to the Nyquist frequency. The phase response (Fig. 3.7) has similar deviation problems.

3.6 Poles

Discrete-time block diagrams are differing from continuous-time block diagrams only by having z^{-1} elements instead of integrators. Recalling that the transfer

⁷Another way to look at this is to notice that in order for z^n to be a complex sinusoid $e^{j\omega n}$ we need to let $z = e^{j\omega}$.

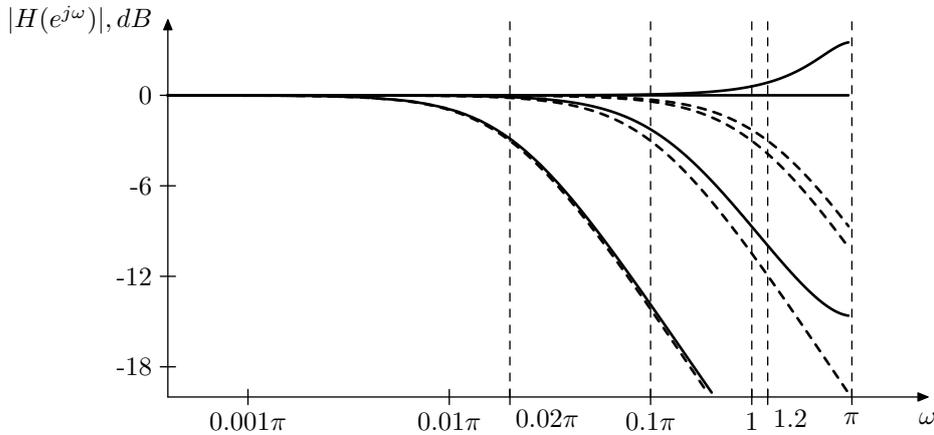


Figure 3.6: Amplitude response of a naive 1-pole lowpass filter for a number of different cutoffs. Dashed curves represent the respective analog filter responses for the same cutoffs.

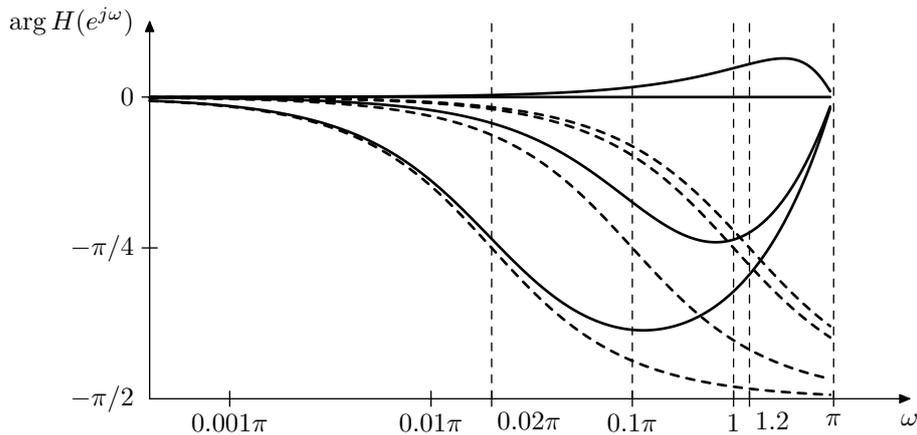


Figure 3.7: Phase response of a naive 1-pole lowpass filter for a number of different cutoffs. Dashed curves represent the respective analog filter responses for the same cutoffs.

function of an integrator is s^{-1} , we conclude that from the formal point of view the difference is purely notational.

Now, the transfer functions of continuous-time block diagrams are nonstrictly proper rational functions of s . Respectively, the transfer functions of discrete-time block diagrams are nonstrictly proper rational functions of z .

Thus, discrete-time transfer functions will have poles and zeros in a way similar to continuous-time transfer functions. Similarly to continuous-time transfer functions, the poles will define the stability of a linear time-invariant filter. Consider that $z = e^{sT}$ and recall the stability criterion $\text{Re } s < 0$ (where $s = p_n$, where p_n are the poles). Apparently, $\text{Re } s < 0 \iff |z| < 1$. We might therefore intuitively expect the discrete-time stability criterion to be $|p_n| < 1$ where p_n are the discrete-time poles. This is indeed the case, a linear time-invariant

difference system⁸ is stable if and only if all its poles are located inside the unit circle.

3.7 Trapezoidal integration

Instead of naive integration, we could attempt using the trapezoidal integration method ($T = 1$):

```
// perform one sample tick of the integrator
integrator_output := integrator_output +
    (integrator_input + previous_integrator_input)/2;
previous_integrator_input := integrator_input;
```

Notice that now we need two state variables per integrator: *integrator_output* and *previous_integrator_input*. The block diagram of a trapezoidal integrator is shown in Fig. 3.8. We'll refer to this integrator as a *direct form I trapezoidal integrator*. The reason for this term will be explained later.

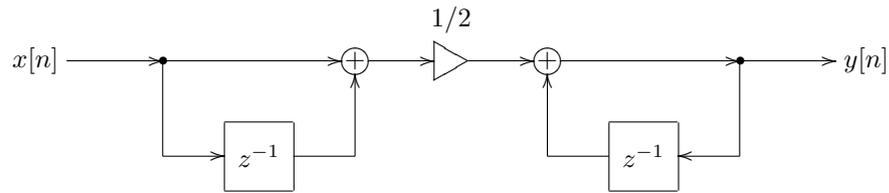


Figure 3.8: Direct form I trapezoidal integrator ($T = 1$).

We could also construct a trapezoidal integrator implementation with only a single state variable. Consider the expression for the trapezoidal integrator's output:

$$y[n] = y[n_0 - 1] + \sum_{\nu=n_0}^n \frac{x[\nu - 1] + x[\nu]}{2} \quad (3.2)$$

Suppose $y[n_0 - 1] = 0$ and $x[n_0 - 1] = 0$, corresponding to a zero initial state (recall that both $y[n_0 - 1]$ and $x[n_0 - 1]$ are technically stored in the z^{-1} elements). Then

$$\begin{aligned} y[n] &= \sum_{\nu=n_0}^n \frac{x[\nu - 1] + x[\nu]}{2} = \frac{1}{2} \left(\sum_{\nu=n_0}^n x[\nu - 1] + \sum_{\nu=n_0}^n x[\nu] \right) = \\ &= \frac{1}{2} \left(\sum_{\nu=n_0+1}^n x[\nu - 1] + \sum_{\nu=n_0}^n x[\nu] \right) = \frac{1}{2} \left(\sum_{\nu=n_0}^{n-1} x[\nu] + \sum_{\nu=n_0}^n x[\nu] \right) = \\ &= \frac{u[n - 1] + u[n]}{2} \end{aligned}$$

⁸Difference systems can be defined as those, whose block diagrams consist of gains, summaters and unit delays. More precisely those are causal difference systems. There are also difference systems with a lookahead into the future, but we don't consider them in this book.

where

$$u[n] = \sum_{\nu=n_0}^n x[\nu]$$

Now notice that $u[n]$ is the output of a naive integrator, whose input signal is $x[n]$. At the same time $y[n]$ is the average of the previous and the current output values of the naive integrator. This can be implemented by the structure in Fig. 3.9. Similar considerations apply for nonzero initial state. We'll refer to the integrator in Fig. 3.9 as a *direct form II* or *canonical* trapezoidal integrator. The reason for this term will be explained later.

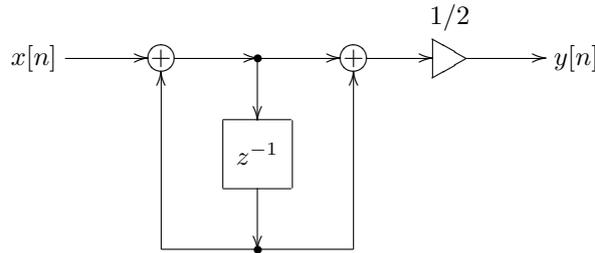


Figure 3.9: Direct form II (canonical) trapezoidal integrator ($T = 1$).

We can develop yet another form of the bilinear integrator with a single state variable. Let's rewrite (3.2) as

$$y[n] = y[n_0 - 1] + \frac{x[n_0 - 1]}{2} + \sum_{\nu=n_0}^{n-1} x[\nu] + \frac{x[n]}{2}$$

and let

$$u[n - 1] = y[n] - \frac{x[n]}{2} = y[n_0 - 1] + \frac{x[n_0 - 1]}{2} + \sum_{\nu=n_0}^{n-1} x[\nu]$$

Notice that

$$y[n] = u[n - 1] + \frac{x[n]}{2} \quad (3.3)$$

and

$$u[n] = u[n - 1] + x[n] = y[n] + \frac{x[n]}{2} \quad (3.4)$$

Expressing (3.3) and (3.4) in a graphical form, we obtain the structure in Fig. 3.10, where the output of the z^{-1} block corresponds to $u[n - 1]$. We'll refer to the integrator in Fig. 3.10 as a *transposed direct form II* or *transposed canonical* trapezoidal integrator. The reason for this term will be explained later.

The positioning of the $1/2$ gain prior to the integrator in Fig. 3.10 is quite convenient, because we can combine the $1/2$ gain with the cutoff gain into a single gain element. In case of an arbitrary sampling period we could also include the T factor into the same gain element, thus obtaining the structure in

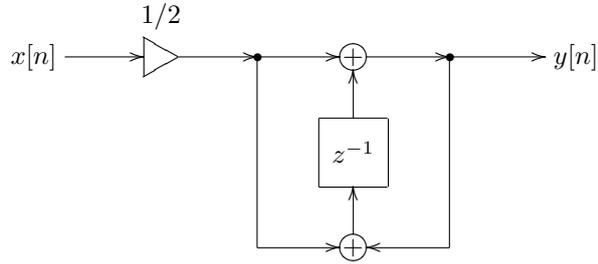


Figure 3.10: Transposed direct form II (transposed canonical) trapezoidal integrator ($T = 1$).

Fig. 3.11. A similar trick can be performed for the other two integrators, if we move the $1/2$ gain element to the input of the respective integrator. Since the integrator is a linear time-invariant system, this doesn't affect the integrator's behavior in a slightest way.

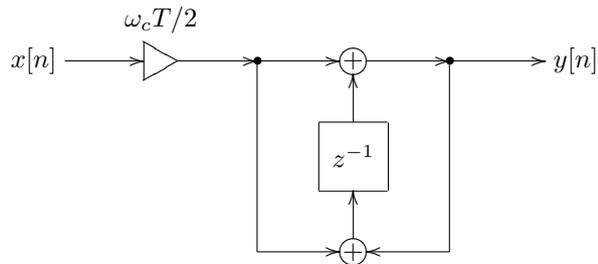


Figure 3.11: Transposed direct form II (transposed canonical) trapezoidal integrator with “embedded” cutoff gain.

Typically one would prefer the direct form II integrators to the direct form I integrator, because the former have only one state variable. In this book we will mostly use the transposed direct form II integrator, because this is resulting in slightly simpler zero-delay feedback equations and also offers a nice possibility for the internal saturation in the integrator.

The transfer functions of all three integrators are identical. Let's obtain e.g. the transfer function of the transposed canonical integrator (in Fig. 3.10). Let u be the output signal of the z^{-1} element. Assuming signals of the exponential form z^n , we have

$$u = \left(\frac{x}{2} + y\right) z^{-1}$$

$$y = \frac{x}{2} + u$$

from where

$$u = y - \frac{x}{2}$$

and

$$y - \frac{x}{2} = \left(\frac{x}{2} + y\right) z^{-1}$$

or

$$\left(y - \frac{x}{2}\right)z = \frac{x}{2} + y$$

from where

$$y(z-1) = \frac{x}{2}(z+1)$$

and the transfer function of the trapezoidal integrator is thus

$$H(z) = \frac{y}{x} = \frac{1}{2} \cdot \frac{z+1}{z-1}$$

For an arbitrary T one has to multiply the result by T , to take the respective gain element into account:

$$H(z) = \frac{T}{2} \cdot \frac{z+1}{z-1}$$

If also the cutoff gain is included, we obtain

$$H(z) = \frac{\omega_c T}{2} \cdot \frac{z+1}{z-1}$$

One can obtain the same results for the other two integrators.

What is so special about this transfer function, that makes the trapezoidal integrator so superior to the naive one, is to be discussed next.

3.8 Bilinear transform

Suppose we take an arbitrary continuous-time block diagram, like the familiar lowpass filter in Fig. 2.2 and replace all continuous-time integrators by discrete-time trapezoidal integrators. On the transfer function level, this will correspond to replacing all s^{-1} with $\frac{T}{2} \cdot \frac{z+1}{z-1}$. That is, technically we perform a substitution

$$s^{-1} = \frac{T}{2} \cdot \frac{z+1}{z-1}$$

in the transfer function expression.

It would be more convenient to write this substitution explicitly as

$$s = \frac{2}{T} \cdot \frac{z-1}{z+1} \quad (3.5)$$

The substitution (3.5) is referred to as the *bilinear transform*, or shortly BLT. For that reason we can also refer to trapezoidal integrators as *BLT integrators*. Let's figure out, how does the bilinear transform affect the frequency response of the filter, that is, what is the relationship between the original continuous-time frequency response prior to the substitution and the resulting discrete-time frequency response after the substitution.

Let $H_a(s)$ be the original continuous-time transfer function. Then the respective discrete-time transfer function is

$$H_d(z) = H_a\left(\frac{2}{T} \cdot \frac{z-1}{z+1}\right) \quad (3.6)$$

Respectively, the discrete-time frequency response is

$$\begin{aligned} H_d(e^{j\omega T}) &= H_a\left(\frac{2}{T} \cdot \frac{e^{j\omega T} - 1}{e^{j\omega T} + 1}\right) = H_a\left(\frac{2}{T} \cdot \frac{e^{j\omega T/2} - e^{-j\omega T/2}}{e^{j\omega T/2} + e^{-j\omega T/2}}\right) = \\ &= H_a\left(\frac{2}{T}j \tan \frac{\omega T}{2}\right) \end{aligned}$$

Notice that $H_a(s)$ in the last expression is evaluated on the imaginary axis!!! That is, the bilinear transform maps the imaginary axis in the s -plane to the unit circle in the z -plane! Now, $H_a\left(\frac{2}{T}j \tan \frac{\omega T}{2}\right)$ is the analog frequency response evaluated at $\frac{2}{T} \tan \frac{\omega T}{2}$. That is, the digital frequency response at ω is equal to the analog frequency response at $\frac{2}{T} \tan \frac{\omega T}{2}$. This means that the analog frequency response in the range $0 \leq \omega < +\infty$ is mapped into the digital frequency range $0 \leq \omega T < \pi$ ($0 \leq \omega < \pi f_s$), that is from zero to Nyquist!⁹ Denoting the analog frequency as ω_a and the digital frequency as ω_d we can express the argument mapping of the frequency response function as

$$\omega_a = \frac{2}{T} \tan \frac{\omega_d T}{2} \quad (3.7)$$

or, in a more symmetrical way

$$\frac{\omega_a T}{2} = \tan \frac{\omega_d T}{2} \quad (3.8)$$

Notice that for frequencies much smaller than Nyquist frequency we have $\omega T \ll 1$ and respectively $\omega_a \approx \omega_d$.

This is what is so unique about the bilinear transform. It simply warps the frequency range $[0, +\infty)$ into the zero-to-Nyquist range, but otherwise doesn't change the frequency response at all! Considering in comparison a naive integrator, we would have obtained:

$$\begin{aligned} s^{-1} &= \frac{z}{z-1} \\ s &= \frac{z-1}{z} \end{aligned} \quad (3.9)$$

$$H_d(z) = H_a\left(\frac{z-1}{z}\right)$$

$$H_d(e^{j\omega}) = H_a\left(\frac{e^{j\omega} - 1}{e^{j\omega}}\right) = H_a(1 - e^{-j\omega})$$

which means that the digital frequency response is equal to the analog transfer function evaluated on a circle of radius 1 centered at $s = 1$. This hardly defines a clear relationship between the two frequency responses.

So, by simply replacing the analog integrators with digital trapezoidal integrators, we obtain a digital filter whose frequency response is essentially the same as the one of the analog prototype, except for the frequency warping. Particularly, the relationship between the amplitude and phase responses of the filter is fully preserved, which is particularly highly important if the filter is to be used as a building block in a larger filter. Very close to perfect!

⁹A similar mapping obviously occurs for the negative frequencies.

Furthermore, the bilinear transform maps the left complex semiplane in the s -domain into the inner region of the unit circle in the z -domain. Indeed, let's obtain the inverse bilinear transform formula. From (3.5) we have

$$(z + 1) \frac{sT}{2} = z - 1$$

from where

$$1 + \frac{sT}{2} = z \left(1 - \frac{sT}{2} \right)$$

and

$$z = \frac{1 + \frac{sT}{2}}{1 - \frac{sT}{2}} \quad (3.10)$$

The equation (3.10) defines the *inverse bilinear transform*. Now, if $\operatorname{Re} s < 0$, then, obviously

$$\left| 1 + \frac{sT}{2} \right| < \left| 1 - \frac{sT}{2} \right|$$

and $|z| < 1$. Thus, the left complex semiplane in the s -plane is mapped to the inner region of the unit circle in the z -plane. In the same way one can show that the right complex semiplane is mapped to the outer region of the unit circle. And the imaginary axis is mapped to the unit circle itself. Comparing the stability criterion of analog filters (the poles must be in the left complex semiplane) to the one of digital filters (the poles must be inside the unit circle), we conclude that the bilinear transform exactly preserves the stability of the filters!

In comparison, for a naive integrator replacement we would have the following. Inverting the (3.9) substitution we obtain

$$sz = z - 1$$

$$z(1 - s) = 1$$

and

$$z = \frac{1}{1 - s}$$

Assuming $\operatorname{Re} s < 0$ and considering that in this case

$$\left| z - \frac{1}{2} \right| = \left| \frac{1}{1 - s} - \frac{1}{2} \right| = \left| \frac{1 - \frac{1}{2} + \frac{s}{2}}{1 - s} \right| = \left| \frac{1}{2} \cdot \frac{1 + s}{1 - s} \right| < \frac{1}{2}$$

we conclude that the left semiplane is mapped into a circle of radius 0.5 centered at $z = 0.5$. So the naive integrator overpreserves the stability, which is not nice, since we would rather have digital filters behaving as closely to their analog prototypes as possible. Considering that this comes in a package with a poor frequency response transformation, we should rather stick with trapezoidal integrators.

So, let's replace e.g. the integrator in the familiar lowpass filter structure in Fig. 2.2 with a trapezoidal integrator. Performing the integrator replacement, we obtain the structure in Fig. 3.12. We will refer to the trapezoidal integrator replacement method as the *topology-preserving transform* (TPT) method. This term will be explained and properly introduced later. For now, before we simply attempt to implement the structure in Fig. 3.12 in code, we should become aware of a few further issues.

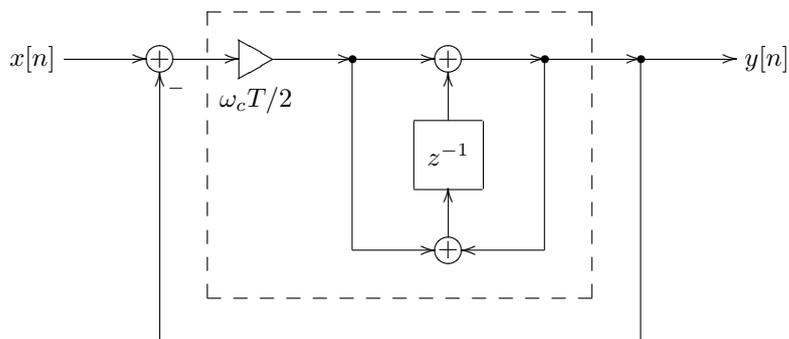


Figure 3.12: 1-pole TPT lowpass filter (the dashed line denotes the trapezoidal integrator).

3.9 Cutoff prewarping

Suppose we are using the lowpass filter structure in Fig. 3.12 and we wish to have its cutoff at ω_c . If we however simply put this ω_c parameter into the respective integrator gain element $\omega_c T/2$, our frequency response at the cutoff will be different from the expected one. Considering the transfer function of an analog 1-pole lowpass filter (2.2), at the cutoff we expect

$$H(j\omega_c) = \frac{\omega_c}{\omega_c + j\omega_c} = \frac{1}{1 + j}$$

corresponding to a -3dB drop in amplitude and a 45° phase shift. However, letting $\omega_a = \omega_c$ in (3.8) we will obtain some $\omega_d \neq \omega_c$. That is the cutoff point of the analog frequency response will be mapped to some other frequency ω_d in the digital frequency response (Fig. 3.13). This is the result of the frequency axis warping by the bilinear transform.¹⁰

However, if we desire to have the $1/(1 + j)$ frequency response exactly at $\omega_d = \omega_c$, we can simply apply (3.7) to $\omega_d = \omega_c$, thereby obtaining some ω_a . This ω_a should be used in the gain element of the integrator, that is the gain should be $\omega_a T/2$ instead of $\omega_c T/2$. This cutoff substitution is referred to as the *cutoff prewarping*.¹¹ The result of the cutoff prewarping is illustrated in Fig. 3.14.

Apparently, the importance of the cutoff prewarping grows as the cutoff values get higher. For cutoff values much lower than the Nyquist frequency the prewarping has only a minor effect.

Notice that it's possible to choose any other point for the prewarping, not necessarily the cutoff point. That is it's possible to make any single chosen point on the analog frequency response to be located at the desired digital frequency. In order to do so we first choose ω_d of interest, then use (3.7) to find the respective ω_a . Now we want a particular point on the analog frequency

¹⁰The response difference at the cutoff in Fig. 3.13 might seem negligible. However it will be even higher for cutoffs closer to Nyquist. Also for filters with strong resonance the detuning of the cutoff by frequency warping might be way more noticeable.

¹¹Since the value $\omega_a T/2$ is the one explicitly used in the integrator, it's more practical to directly use (3.8) rather than (3.7) for the prewarping.

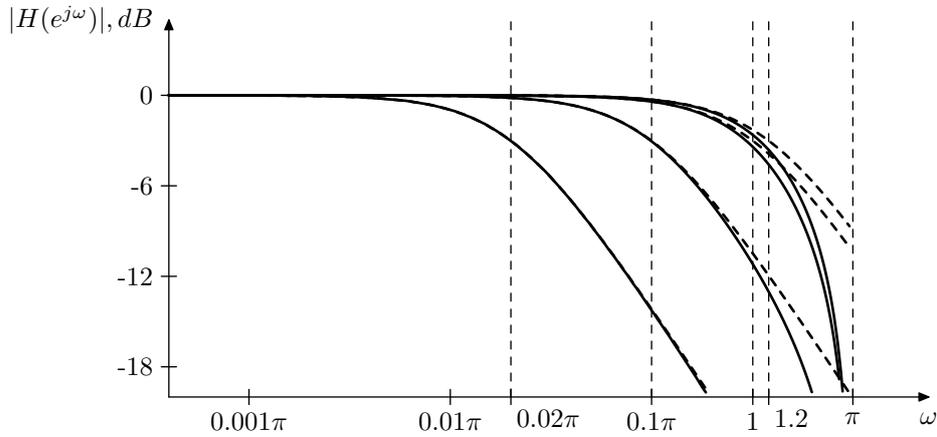


Figure 3.13: Amplitude response of an unwarped bilinear-transformed 1-pole lowpass filter for a number of different cutoffs. Dashed curves represent the respective analog filter responses for the same cutoffs. Observe the difference between the analog and digital responses at each cutoff.

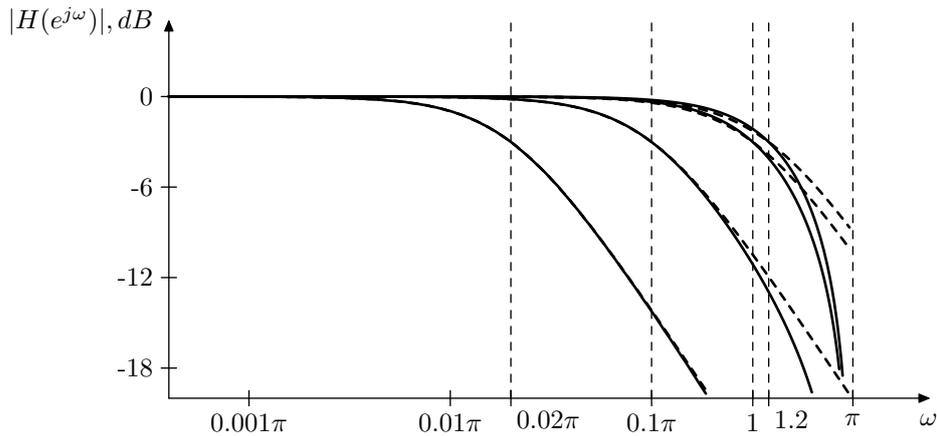


Figure 3.14: Amplitude response of a prewarped bilinear-transformed 1-pole lowpass filter for a number of different cutoffs. Dashed curves represent the respective analog filter responses for the same cutoffs. Observe the identical values of the analog and digital responses at each cutoff.

response to be located at ω_a , which can be achieved by a proper choice of the analog cutoff value. Now we put this cutoff value into the integrators and that's it!

3.10 Zero-delay feedback

There is a further problem with the trapezoidal integrator replacement in the TPT method. Replacing the integrators with trapezoidal ones introduces *delay-*

less feedback loops (that is, feedback loops not containing any delay elements) into the structure. E.g. consider the structure in Fig. 3.12. Carefully examining this structure, we find that it has a feedback loop which doesn't contain any unit delay elements. This loop goes from the leftmost summator through the gain, through the upper path of the integrator to the filter's output and back through the large feedback path to the leftmost summator.

Why is this delayless loop a problem? Let's consider for example the naive lowpass filter structure in Fig. 3.5. Suppose we don't have the respective program code representation and wish to obtain it from the block diagram. We could do it in the following way. Consider Fig. 3.15, which is the same as Fig. 3.5, except that it labels all signal points. At the beginning of the computation of a new sample the signals A and B are already known. $A = x[n]$ is the current input sample and B is taken from the internal state memory of the z^{-1} element. Therefore we can compute $C = A - B$. Then we can compute $D = (\omega_c T)C$ and finally $E = D + B$. The value of E is then stored into the internal memory of the z^{-1} element (for the next sample computation) and is also sent to the output as the new $y[n]$ value. Easy, right?

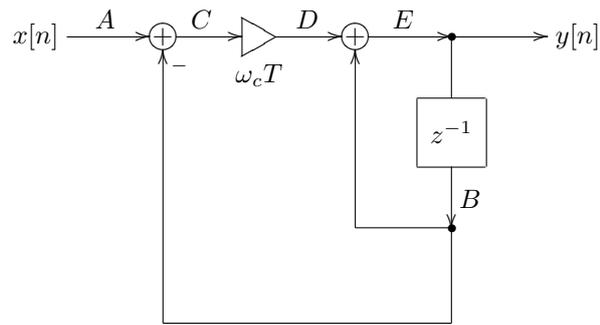


Figure 3.15: Naive 1-pole lowpass filter and the respective signal computation order.

Now the same approach doesn't work for the structure in Fig. 3.12. Because there is a delayless loop, we can't find a starting point for the computation within that loop.

The classical way of solving this problem is exactly the same as what we had in the naive approach: introduce a z^{-1} into the delayless feedback, turning it into a feedback containing a unit delay (Fig. 3.16). Now there are no delayless feedback paths and we can arrange the computation order in a way similar to Fig. 3.15. This however destroys the resulting frequency response, because the transfer function is now different. In fact the obtained result is not significantly better (if better at all) than the one from the naive approach. There are some serious artifacts in the frequency response closer to the Nyquist frequency, if the filter cutoff is sufficiently high.

Therefore we shouldn't introduce any modifications into the structure and solve the *zero-delay feedback* problem instead. The term "zero-delay feedback" originates from the fact that we avoid introducing a one-sample delay into the feedback (like in Fig. 3.16) and instead keep the feedback delay equal to zero.

So, let's solve the zero-delay feedback problem for the structure in Fig. 3.12.

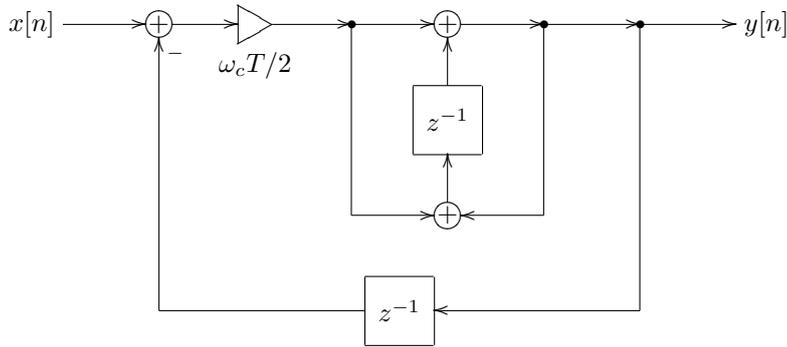


Figure 3.16: Digital 1-pole lowpass filter with a trapezoidal integrator and an extra delay in the feedback.

Notice that this structure simply consists of a negative feedback loop around a trapezoidal integrator, where the trapezoidal integrator structure is exactly the one from Fig. 3.11. We will now introduce the concept of the *instantaneous response* of this integrator structure.

So, consider the integrator structure in Fig. 3.11 and let $u[n]$ denote the input signal of the z^{-1} element, respectively its output will be $u[n - 1]$. Since there are no delayless loops in the integrator, it's not difficult to obtain the following expression for $y[n]$:

$$y[n] = \frac{\omega_c T}{2} x[n] + u[n - 1] \quad (3.11)$$

Notice that, at the time $x[n]$ arrives at the integrator's input, all values in the right-hand side of (3.11) are known (no unknown variables). Introducing notation

$$g = \frac{\omega_c T}{2}$$

$$s[n] = u[n - 1]$$

we have

$$y[n] = gx[n] + s[n]$$

or, dropping the discrete time argument notation for simplicity,

$$y = gx + s$$

That is, at any given time moment n , the output of the integrator y is a linear function of its input x , where the values of the parameters of this linear function are known. The g parameter doesn't depend on the internal state of the integrator, while the s parameter does depend on the internal state of the integrator. We will refer to the linear function $f(x) = gx + s$ as the *instantaneous response* of the integrator at the respective implied time moment n . The coefficient g can be referred to as the *instantaneous response gain* or simply *instantaneous gain*. The term s can be referred to as the *instantaneous response offset* or simply *instantaneous offset*.

Let's now redraw the filter structure in Fig. 3.12 as in Fig. 3.17. We have changed the notation from x to ξ in the $g\xi + s$ expression to avoid the confusion with the input signal $x[n]$ of the entire filter.

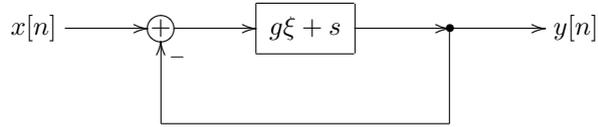


Figure 3.17: 1-pole TPT lowpass filter with the integrator in the instantaneous response form.

Now we can easily write and solve the zero-delay feedback equation. Indeed, suppose we already know the filter output $y[n]$. Then the output signal of the feedback summator is $x[n] - y[n]$ and the output of the integrator is respectively $g(x[n] - y[n]) + s$. Thus

$$y[n] = g(x[n] - y[n]) + s$$

or, dropping the time argument notation for simplicity,

$$y = g(x - y) + s \quad (3.12)$$

The equation (3.12) is the zero-delay feedback equation for the filter in Fig. 3.17 (or, for that matter, in Fig. 3.12). Solving this equation, we obtain

$$y(1 + g) = gx + s$$

and respectively

$$y = \frac{gx + s}{1 + g} \quad (3.13)$$

Having found y (that is, having predicted the output $y[n]$), we can then proceed with computing the other signals in the structure in Fig. 3.12, beginning with the output of the leftmost summator.¹²

It's worth mentioning that (3.13) can be used to obtain the instantaneous response of the entire filter from Fig. 3.12. Indeed, rewriting (3.13) as

$$y = \frac{g}{1 + g}x + \frac{s}{1 + g}$$

and introducing notations

$$G = \frac{g}{1 + g}$$

$$S = \frac{s}{1 + g}$$

we have

$$y = Gx + S \quad (3.14)$$

¹²Notice that the choice of the signal point for the prediction is rather arbitrary. We could have chosen any other point within the delayless feedback loop.

So, the instantaneous response of the entire lowpass filter in Fig. 3.12 is again a linear function of the input. We could use the expression (3.14) e.g. to solve the zero-delay feedback problem for some larger feedback loop containing a 1-pole lowpass filter.

Let's now convert the structure in Fig. 3.12 into a piece of code. We already know y from (3.14). Let's notice that the output of the $\omega_c T/2$ gain is used twice in the structure. Let v denote the output of this gain. Since $g = \omega_c T/2$, we have

$$\begin{aligned} v &= g(x - y) = g(x - Gx - S) = g\left(x - \frac{g}{1+g}x - \frac{s}{1+g}\right) = \\ &= g\left(\frac{1}{1+g}x - \frac{s}{1+g}\right) = g\frac{x-s}{1+g} \end{aligned} \quad (3.15)$$

Recall that s is the output value of the z^{-1} element and let u denote its input value. Then

$$y = v + s \quad (3.16)$$

and

$$u = y + v \quad (3.17)$$

The equations (3.15), (3.16) and (3.17) can be directly expressed in program code:

```
// perform one sample tick of the lowpass filter
v := (x-z1_state)*g/(1+g);
y := v + z1_state;
z1_state := y + v;
```

or instead expressed in a block diagram form (Fig. 3.18). Notice that the block diagram doesn't contain any delayless loops anymore.

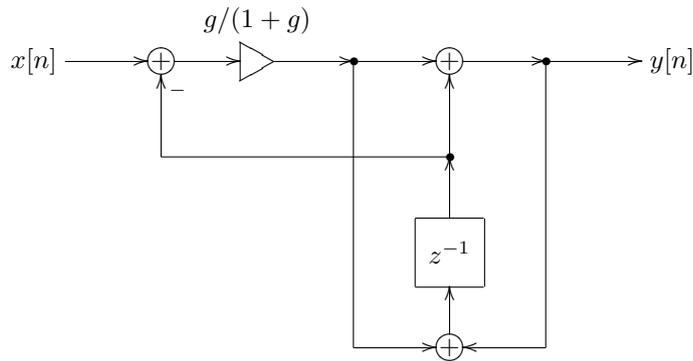


Figure 3.18: 1-pole TPT lowpass filter with resolved zero-delay feedback.

3.11 Direct forms

Consider again the equation (3.6), which describes the application of the bilinear transform to convert an analog transfer function to a digital one. There is a classical method of digital filter design which is based directly on this transformation, without using any integrator replacement techniques. In the author's experience, for music DSP needs this method typically has a largely inferior quality, compared to the TPT. Nevertheless we will describe it here for completeness and for a couple of other reasons. Firstly, it would be nice to try to analyse and understand the reasons for the problems of this method. Secondly, this method could be useful once in a while. Particularly, its deficiencies mostly disappear in the time-invariant (unmodulated or sufficiently slowly modulated) case.

Having obtained a digital transfer function from (3.6), we could observe, that, since the original analog transfer function was a rational function of s , the resulting digital transfer function will necessarily be a rational function of z . E.g. using the familiar 1-pole lowpass transfer function

$$H_a(s) = \frac{\omega_c}{s + \omega_c}$$

we obtain

$$\begin{aligned} H_d(z) &= H_a\left(\frac{2}{T} \cdot \frac{z-1}{z+1}\right) = \frac{\omega_c}{\frac{2}{T} \cdot \frac{z-1}{z+1} + \omega_c} = \\ &= \frac{\frac{\omega_c T}{2}(z+1)}{(z-1) + \frac{\omega_c T}{2}(z+1)} = \frac{\frac{\omega_c T}{2}(z+1)}{(1 + \frac{\omega_c T}{2})z - (1 - \frac{\omega_c T}{2})} \end{aligned}$$

Now, there are standard discrete-time structures allowing an implementation of any given nonstrictly proper rational transfer function. It is easier to use these structures, if the transfer function is expressed as a rational function of z^{-1} rather than the one of z . In our particular example, we can multiply the numerator and the denominator by z^{-1} , obtaining

$$H_d(z) = \frac{\frac{\omega_c T}{2}(1 + z^{-1})}{(1 + \frac{\omega_c T}{2}) - (1 - \frac{\omega_c T}{2})z^{-1}}$$

The further requirement is to have the constant term in the denominator equal to 1, which can be achieved by dividing everything by $1 + \omega_c T/2$:

$$H_d(z) = \frac{\frac{\omega_c T}{2}(1 + z^{-1})}{1 - \frac{1 - \frac{\omega_c T}{2}}{1 + \frac{\omega_c T}{2}}z^{-1}} \quad (3.18)$$

Now suppose we have an arbitrary rational nonstrictly proper transfer function of z , expressed via z^{-1} and having the constant term in the denominator equal to 1:

$$H(z) = \frac{\sum_{n=0}^N b_n z^{-n}}{1 - \sum_{n=1}^N a_n z^{-n}} \quad (3.19)$$

This transfer function can be implemented by the structure in Fig. 3.19 or by the structure in Fig. 3.20. One can verify (by computing the transfer functions of the respective structures) that they indeed implement the transfer function (3.19). There are also transposed versions of these structures, which the readers should be able to construct on their own.

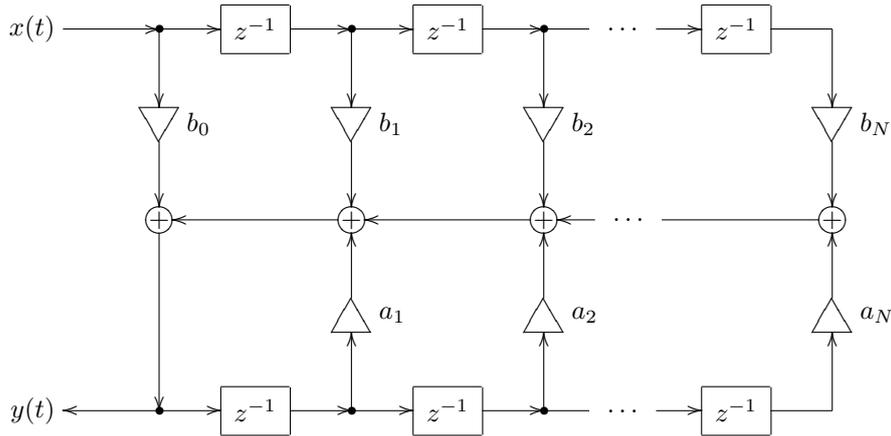


Figure 3.19: Direct form I (DF1).

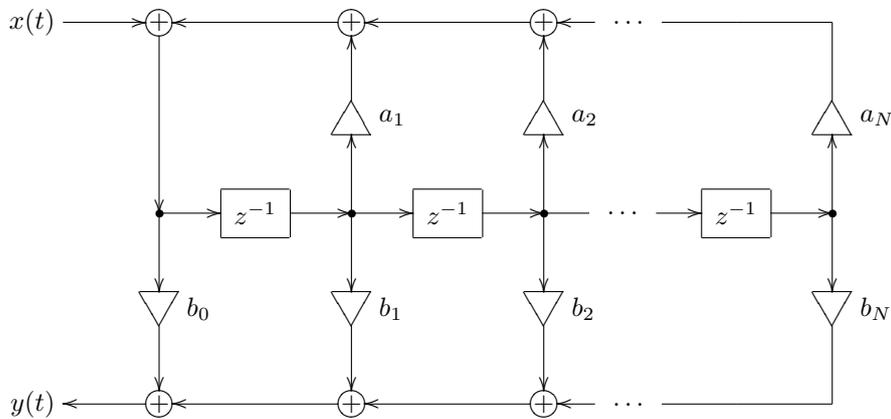


Figure 3.20: Direct form II (DF2), a.k.a. canonical form.

Let's use the direct form II to implement (3.18). Apparently, we have

$$N = 1$$

$$b_0 = b_1 = \frac{\frac{\omega_c T}{2}}{1 + \frac{\omega_c T}{2}}$$

$$a_1 = \frac{1 - \frac{\omega_c T}{2}}{1 + \frac{\omega_c T}{2}}$$

and the direct form implementation itself is the one in Fig. 3.21 (we have merged the b_0 and b_1 coefficients into a single gain element).

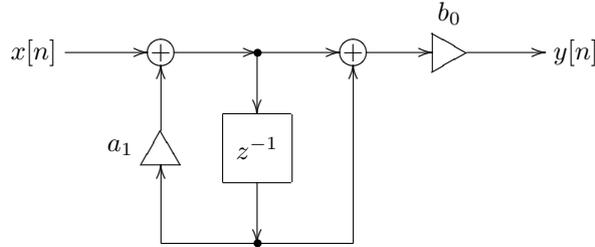
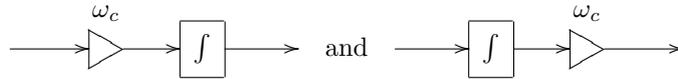


Figure 3.21: Direct form II 1-pole lowpass filter.

In the time-invariant (unmodulated) case the performance of the direct form filter in Fig. 3.21 should be identical to the TPT filter in Fig. 3.12, since both implement the same bilinear-transformed analog transfer function (2.2). When the cutoff is modulated, however, the performance will be different.

This is very easy to understand intuitively. First, consider the two following analog structures, implementing two different ways of combining a cutoff gain with an integrator:



Suppose the input signal is a sine and there is a sudden jump in the cutoff parameter. In this case there will also be a sudden jump in the input of the first integrator, however the jump will be converted into a break by the integration. In the second case the jump will not be converted, because it appears after the integrator. Ignoring the problem of a DC offset possibly introduced by such jump in the first structure (because in real stable filters this DC offset will quickly disappear with time), we should say that the first structure has a better modulation performance, since the cutoff jumps do not produce so audible clicks in the output.

Apparently the two structures behave differently in the time-varying case, even though both have the same transfer function ω_c/s . We say that the two structures have the same transfer function but different *topology* (the latter term referring to the components used in the block diagram and the way they are connected to each other). As mentioned, the transfer function is applicable only to the time-invariant case. No wonder its possible to have structures with identical transfer functions, but different time-varying behavior.

Now, returning to the comparison of implementations in Figs. 3.21 and 3.12, we notice that the structure in Fig. 3.21 contains a gain element at the output, the value of this gain being approximately proportional to the cutoff (at low cutoffs). This will particularly produce unsmoothed jumps in the output in response to jumps in the cutoff value. In the structure in Fig. 3.12, on the other hand, the cutoff jumps will be smoothed by the integrator. Thus, the difference

between the two structures is similar to the just discussed effect of the cutoff gain placement with the integrator.

We should conclude that, other things being equal, the structure in Fig. 3.21 is inferior to the one in Fig. 3.12 (or Fig. 3.18). In this respect consider that Fig. 3.12 is trying to explicitly emulate the analog integration behavior, *preserving the topology* of the original analog structure, while Fig. 3.21 is concerned solely with implementing a correct transfer function. Since Fig. 3.21 implements a classical approach to the bilinear transform application for digital filter design (which ignores the filter topology) we'll refer to the trapezoidal integration replacement technique as the *topology-preserving bilinear transform* (or, shortly, TPBLT). Or, even shorter, we can refer to this technique as simply the *topology-preserving transform* (TPT), implicitly assuming that the bilinear transform is being used.¹³

In principle, sometimes there are possibilities to “manually fix” the structures such as in Fig. 3.21. E.g. the time-varying performance of the latter is drastically improved by moving the b_0 gain to the input. The problem however is that this kind of fixing quickly gets more complicated (if being possible at all) with larger filter structures. On the other hand, the TPT method explicitly aims at emulating the time-varying behavior of the analog prototype structure, which aspect is completely ignored by the classical transform approach. Besides, if the structure contains nonlinearities, preserving the topology becomes absolutely critical, because otherwise these nonlinearities can not be placed in the digital model.¹⁴ Also, the direct forms suffer from precision loss issues, the problem growing bigger with the order of the system. For that reason in practice the direct forms of orders higher than 2 are rarely used,¹⁵ but even 2nd-order direct forms could already noticeably suffer from precision losses.

3.12 Other replacement techniques

The trapezoidal integrator replacement technique can be seen as a particular case of a more general set of replacement techniques. Suppose we have two filters, whose frequency response functions are $F_1(\omega)$ and $F_2(\omega)$ respectively. The filters do not need to have the same nature, particularly one can be an analog filter while the other can be a digital one. Suppose further, there is a frequency axis mapping function $\omega' = \mu(\omega)$ such that

$$F_2(\omega) = F_1(\mu(\omega))$$

Typically $\mu(\omega)$ should map the entire domain of $F_2(\omega)$ onto the entire domain of $F_1(\omega)$ (however the exceptions are possible).

¹³Apparently, naive filter design techniques also preserve the topology, but they do a rather poor job on the transfer functions. Classical bilinear transform approach does a good job on the transfer function, but doesn't preserve the topology. The topology-preserving transform achieves both goals simultaneously.

¹⁴This is related to the fact that transfer functions can be defined only for linear time-invariant systems. Nonlinear cases are obviously not linear, thus some critical information can be lost, if the conversion is done solely based on the transfer functions.

¹⁵A higher-order transfer function is typically decomposed into a product of transfer functions of 1st- and 2nd-order rational functions (with real coefficients!). Then it can be implemented by a serial connection of the respective 1st- and 2nd-order direct form filters.

To make the subsequent discussion more intuitive, we will assume that $\mu(\omega)$ is monotone, although this is absolutely not a must.¹⁶ In this case we could say that $F_2(\omega)$ is obtained from $F_1(\omega)$ by a frequency axis warping. Particularly, this is exactly what happens in the bilinear transform case (the mapping $\mu(\omega)$ is then defined by the equation (3.7)). One cool thing about the frequency axis warping is that it preserves the relationship between the amplitude and phase.

Suppose that we have a structure built around filters of frequency response $F_1(\omega)$, and the rest of the structure doesn't contain any memory elements (such as integrators or unit delays). Then the frequency response $F(\omega)$ of this structure will be a function of $F_1(\omega)$:

$$F(\omega) = \Phi(F_1(\omega))$$

where the specifics of the function $\Phi(w)$ will be defined by the details of the container structure. E.g. if the building-block filters are analog integrators, then $F_1(\omega) = 1/j\omega$. For the filter in Fig. 2.2 we then have

$$\Phi(w) = \frac{w}{w+1}$$

Indeed, substituting $F_1(\omega)$ into $\Phi(w)$ we obtain

$$F(\omega) = \Phi(F_1(\omega)) = \Phi(1/j\omega) = \frac{1/j\omega}{1+1/j\omega} = \frac{1}{1+j\omega}$$

which is the already familiar to us frequency response of the analog lowpass filter.

Now, we can view the trapezoidal integrator replacement as a substitution of F_2 instead of F_1 , where $\mu(\omega)$ is obtained from (3.7):

$$\omega_a = \mu(\omega_d) = \frac{2}{T} \tan \frac{\omega_d T}{2}$$

The frequency response of the resulting filter is obviously equal to $\Phi(F_2(\omega))$, where $F_2(\omega)$ is the frequency response of the trapezoidal integrators (used in place of analog ones). But since $F_2(\omega) = F_1(\mu(\omega))$.

$$\Phi(F_2(\omega)) = \Phi(F_1(\mu(\omega)))$$

which means that the frequency response $\Phi(F_2(\cdot))$ of the structure with trapezoidal integrators is obtained from the frequency response $\Phi(F_1(\cdot))$ of the structure with analog integrators simply by warping the frequency axis. If the warping is not too strong, the frequency responses will be very close to each other. This is exactly what is happening in the trapezoidal integrator replacement and generally in the bilinear transform.

Differentiator-based filters

We could have used some other two filters, with their respective frequency responses F_1 and F_2 . E.g. we could consider continuous-time systems built around

¹⁶Strictly speaking, we don't even care whether $\mu(\omega)$ is single-valued. We could have instead required that

$$F_2(\mu_2(\omega)) = F_1(\mu_1(\omega))$$

for some $\mu_1(\omega)$ and $\mu_2(\omega)$.

differentiators rather than integrators.¹⁷ The transfer function of a differentiator is apparently simply $H(s) = s$, so we could use (3.5) to build a discrete-time “trapezoidal differentiator”. Particularly, if we use the direct form II approach, it could look similarly to the integrator in Fig. 3.9. When embedding the cutoff control into a differentiator (in the form of a $1/\omega_c$ gain), it’s probably better to position it after the differentiator, to avoid the unnecessary “de-smoothing” of the control modulation by the differentiator. Replacing the analog differentiators in a structure by such digital trapezoidal differentiators we effectively perform a differentiator-based TPT.

E.g. if we replace the integrator in the highpass filter in Fig. 2.8 by a differentiator, we essentially perform a $1/s \leftarrow s$ substitution, thus we should have obtained a (differentiator-based) lowpass filter. Remarkably, if we perform a differentiator-based TPT on such filter, the obtained digital structure is fully equivalent to the previously obtained integrator-based TPT 1-pole lowpass filter.

Allpass substitution

One particularly interesting case occurs when F_1 and F_2 define two different allpass frequency responses. That is $|F_1(\omega)| \equiv 1$ and $|F_2(\omega)| \equiv 1$. In this case the mapping $\mu(\omega)$ is always possible. Especially since the allpass responses (defined by rational transfer functions of analog and digital systems) always cover the entire phase range from $-\pi$ to π .¹⁸ In intuitive terms it means: for a filter built of identical allpass elements, we can always replace those allpass elements with an arbitrary other type of allpass elements (provided all other elements are memoryless, that is there are only gains and summators). We will refer to this process as *allpass substitution*. Whereas in the trapezoidal integrator replacement we have replaced analog integrators by digital trapezoidal integrators, in the allpass substitution we replace allpass filters of one type by allpass filters of another type.

We can even replace digital allpass filters with analog ones and vice versa. E.g., noticing that z^{-1} elements *are* allpass filters, we could replace them with analog allpass filters. One particularly interesting case arises out of the inverse bilinear transform (3.10). From (3.10) we obtain

$$z^{-1} = \frac{1 - \frac{sT}{2}}{1 + \frac{sT}{2}} \quad (3.20)$$

The right-hand side of (3.20) obviously defines a stable 1-pole allpass filter, whose cutoff is $2/T$. We could take a digital filter and replace all z^{-1} elements with an analog allpass filter structure implementing (3.20). By doing this we would have performed a topology-preserving inverse bilinear transform.

We could then apply the cutoff parametrization to these underlying analog allpass elements:

$$\frac{sT}{2} \leftarrow \frac{s}{\omega_c}$$

¹⁷The real-world analog electronic circuits are “built around” integrators rather than differentiators. However, formally one still can “invert” the causality direction in the equations and pretend that $\dot{x}(t)$ is defined by $x(t)$, and not vice versa.

¹⁸Actually, for $-\infty < \omega < +\infty$, they cover this range exactly N times, where N is the order of the filter.

so that we obtain

$$z^{-1} = \frac{1 - s/\omega_c}{1 + s/\omega_c}$$

The expression s/ω_c can be also rewritten as $sT/2\alpha$, where α is the cutoff scaling factor:

$$z^{-1} = \frac{1 - sT/2\alpha}{1 + sT/2\alpha} \quad (3.21)$$

Finally, we can apply the trapezoidal integrator replacement to the cutoff-scaled analog filter, converting it back to the digital domain. By doing so, we have applied the cutoff scaling in the digital domain! On the transfer function level this is equivalent to applying the bilinear transform to (3.21), resulting in

$$\begin{aligned} z^{-1} &= \frac{1 - sT/2\alpha}{1 + sT/2\alpha} \leftarrow \frac{1 - \frac{z-1}{\alpha(z+1)}}{1 + \frac{z-1}{\alpha(z+1)}} = \\ &= \frac{\alpha(z+1) - (z-1)}{\alpha(z+1) + (z-1)} = \frac{(\alpha-1)z + (\alpha+1)}{(\alpha+1)z + (\alpha-1)} \end{aligned}$$

That is, we have obtained a discrete-time allpass substitution

$$z^{-1} \leftarrow \frac{(\alpha-1)z + (\alpha+1)}{(\alpha+1)z + (\alpha-1)}$$

which applies cutoff scaling in the digital domain.¹⁹ The allpass filter

$$H(z) = \frac{(\alpha-1)z + (\alpha+1)}{(\alpha+1)z + (\alpha-1)}$$

should have been obtained, as described, by the trapezoidal integrator replacement in an analog implementation of (3.21), alternatively we could use a direct form implementation. Notice that this filter has a pole at $z = (\alpha-1)/(\alpha+1)$. Since $|\alpha-1| < |\alpha+1| \forall \alpha > 0$, the pole is always located inside the unit circle, and the filter is always stable.

3.13 Instantaneously unstable feedback

Writing the solution (3.13) for the zero-delay feedback equation (3.12) we in fact have slightly jumped the gun. Why? Let's consider once again the structure in Fig. 3.17 and suppose g gets negative and starts growing in magnitude further in the negative direction.²⁰ When g becomes equal to -1 , the denominator of (3.13) turns into zero. Something bad must be happening at this moment.

In order to understand the meaning of this situation, let's consider the delayless feedback path as if it was an analog feedback. An analog signal value

¹⁹Differently from the analog domain, the digital cutoff scaling doesn't exactly shift the response along the frequency axis in a logarithmic scale, as some frequency axis warping is involved. The resulting frequency response change however is pretty well approximated as shifting in the lower frequency range.

²⁰Of course, such lowpass filter formally has a negative cutoff value. It is also unstable. However unstable circuits are very important as the linear basis for the analysis and implementation of e.g. nonlinear self-oscillating filters. Therefore we wish to be able to handle unstable circuits as well.

can't change instantly. It can change very quickly, but not instantly, it's always a continuous function of time. We could imagine there is a smoother unit somewhere in the feedback path (Fig. 3.22). This smoother unit has a very very fast response time. We introduce the notation \bar{y} for the output of the smoother.

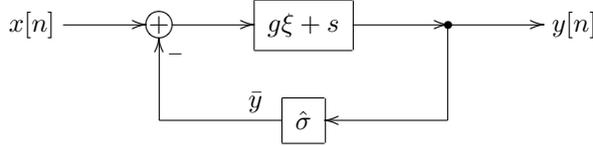


Figure 3.22: Digital 1-pole lowpass filter with a trapezoidal integrator in the instantaneous response form and a smoother unit $\hat{\sigma}$ in the delayless feedback path.

So, suppose we wish to compute a new output sample $y[n]$ for the new input sample $x[n]$. At the time $x[n]$ “arrives” at the filter’s input, the smoother still holds the old output value $y[n - 1]$. Let’s freeze the discrete time at this point (which formally means we simply are not going to update the internal state of the z^{-1} element). At the same time we will let the continuous time t run, formally starting at $t = 0$ at the discrete time moment n .

In this time-frozen setup we can choose arbitrary units for the continuous time t . The smoother equation can be written as

$$\text{sgn} \dot{\bar{y}}(t) = \text{sgn}(y(t) - \bar{y}(t))$$

That is, we don’t specify the details of the smoothing behavior, however the smoother output always changes in the direction from \bar{y} towards y at some (not necessarily constant) speed.²¹ Particularly, we can simply define a constant speed smoother:

$$\dot{\bar{y}} = \text{sgn}(y - \bar{y})$$

or we could use a 1-pole lowpass filter as a smoother:

$$\dot{\bar{y}} = y - \bar{y}$$

The initial value of the smoother is apparently $\bar{y}(0) = y[n - 1]$.

Now consider that

$$\begin{aligned} \text{sgn} \dot{\bar{y}}(t) &= \text{sgn}(y(t) - \bar{y}(t)) = \text{sgn}(g(x[n] - \bar{y}(t)) + s - \bar{y}(t)) = \\ &= \text{sgn}((gx[n] + s) - (1 + g)\bar{y}(t)) = \text{sgn}(a - (1 + g)\bar{y}(t)) \end{aligned}$$

where $a = gx[n] + s$ is constant in respect to t . First, assume $1 + g > 0$. Further, suppose $a - (1 + g)\bar{y}(0) > 0$. Then $\dot{\bar{y}}(0) > 0$ and then the value of the expression $a - (1 + g)\bar{y}(t)$ will start decreasing until it turns to zero at some t , at which point the smoothing process converges. On the other hand, if $a - (1 + g)\bar{y}(0) < 0$, then $\dot{\bar{y}}(0) < 0$ and the value of the expression $a - (1 + g)\bar{y}(t)$ will start increasing until it turns to zero at some t , at which point the smoothing process converges. If $a - (1 + g)\bar{y}(0) = 0$ then the smoothing is already in a stable equilibrium state.

²¹We also assume that the smoothing speed is sufficiently large to ensure that the smoothing process will converge at all cases where it potentially can converge (this statement should become clearer as we discuss more details).

So, in case $1 + g > 0$ the instantaneous feedback smoothing process always converges. Now assume $1 + g \leq 0$. Further, suppose $a - (1 + g)\bar{y}(0) > 0$. Then $\dot{\bar{y}}(0) > 0$ and then the value of the expression $a - (1 + g)\bar{y}(t)$ will start further increasing (or stay constant if $1 + g = 0$). Thus, $\bar{y}(t)$ will grow indefinitely. Respectively, if $a - (1 + g)\bar{y}(0) < 0$, then $\bar{y}(t)$ will decrease indefinitely. This indefinite growth/decrease will occur within the frozen discrete time. Therefore we can say that \bar{y} grows infinitely in an instant. We can refer to this as to an *instantaneously unstable* zero-delay feedback loop.

The analysis of the instantaneous stability can also be done using the analog filter stability analysis means. Let the smoother be an analog 1-pole lowpass filter with a unit cutoff (whose transfer function is $\frac{1}{s+1}$)²² and notice that in that case the structure in Fig. 3.22 can be redrawn as in Fig. 3.23. This filter has two formal inputs $x[n]$ and s and one output $y[n]$.

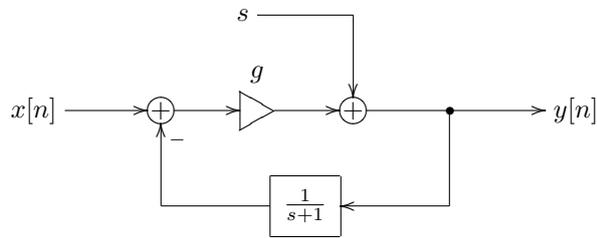


Figure 3.23: An instantaneous representation of a digital 1-pole lowpass filter with a trapezoidal integrator and an analog lowpass smoother.

We can now e.g. obtain a transfer function from the $x[n]$ input to the $y[n]$ output. Ignoring the s input signal (assuming it to be zero), for a continuous-time complex exponential input signal arriving at the $x[n]$ input, which we denote as $x[n](t)$, we have a respective continuous-time complex exponential signal at the $y[n]$ output, which we denote as $y[n](t)$:

$$y[n](t) = g \left(x[n](t) - \frac{1}{s+1} y[n](t) \right)$$

from where

$$y[n](t) = \frac{g}{1 + g \frac{1}{s+1}} x[n](t)$$

that is

$$H(s) = \frac{g}{1 + g \frac{1}{s+1}} = g \frac{s+1}{s+(1+g)}$$

This transfer function has a pole at $s = -(1 + g)$. Therefore, the structure is stable if $1 + g > 0$ and not stable otherwise.

²²Apparently, the variable s used in the transfer function $\frac{1}{s+1}$ is a different s than the one used in the instantaneous response expression for the integrator. The author apologizes for the slight confusion.

The same transfer function analysis could have been done between the s input and the $y[n]$ output, in which case we would have obtained

$$H(s) = \frac{s + 1}{s + (1 + g)}$$

The poles of this transfer function however, are exactly the same, so it doesn't matter.²³

Alright, so we have found out that the filter in Fig. 3.12 is instantaneously unstable if $g \leq -1$, but what can we do about it? Firstly, the problem typically doesn't occur, as normally $g > 0$ (not only in the 1-pole lowpass filter case, but also in other cases). Even if it can occur in principle, one can consider, whether these extreme parameter settings are so necessary to support, and possibly simply clip the filter parameters in such a way that the instantaneous instability doesn't occur.

Secondly, let's notice that $g = \omega_c T/2$. Therefore another solution could be to increase the sampling rate (and respectively reduce the sampling period T).²⁴

Unstable bilinear transform

There is yet another idea, which hasn't been tried out in practice yet.²⁵ We are going to discuss it anyway.

The instantaneous instability is occurring at the moment when one of the analog filter's poles hits the pole of the inverse bilinear transform (3.10), which is located at $s = 2/T$. On the other hand, recall that the bilinear transform is mapping the imaginary axis to the unit circle, thus kind-of preserving the frequency response. If the system is not stable, then the frequency response doesn't make sense. Formally, the reason for this is that the inverse Laplace transform of transfer functions only converges for $\sigma > \max \{\operatorname{Re} p_n\}$ where p_n are the poles of the transfer function, and respectively, if $\max \{\operatorname{Re} p_n\} \geq 0$, it doesn't converge on the imaginary axis ($\sigma = 0$). However, instead of the imaginary axis $\operatorname{Re} s = \sigma = 0$, let's choose some other axis $\operatorname{Re} s = \sigma > \max \{\operatorname{Re} p_n\}$ and use it instead of the imaginary axis to compute the "frequency response".

We also need to find a discrete-time counterpart for $\operatorname{Re} s = \sigma$. Considering that $\operatorname{Re} s$ defines the magnitude growth speed of the exponentials e^{st} we could choose a z -plane circle, on which the magnitude growth speed of z^n is the same

²³This is a common rule: the poles of a system with multiple inputs and/or multiple outputs are always the same regardless of the particular input-output pair for which the transfer function is being considered (exceptions in singular cases, arising out of pole/zero cancellation are possible, though).

²⁴Actually, the instantaneous instability has to do with the fact that the trapezoidal integration is not capable of producing reasonable approximation of the continuous-time integration, due to too extreme parameter values. Increasing the sampling rate obviously increases the precision of the trapezoidal integration as well.

The same idea can also be used to easily and reliably find out, whether the positive value of the denominator of the feedback equation's solution corresponds to the instantaneously stable case or vice versa. The sign which the denominator has for $T \rightarrow 0$ corresponds to the instantaneously stable case.

²⁵The author just got the idea while writing the book and didn't find the time yet to properly experiment with it. Sufficient theoretical analysis is not possible here due to the fact that practical applications of instantaneously unstable (or any unstable, for that matter) filters occur typically for nonlinear filters, and there's not many theoretical analysis means for the latter. Hopefully there are no mistakes in the theoretical transformations, but even if there are mistakes, at least the idea itself could maybe work.

as for $e^{\sigma t}$. Apparently, this circle is $|z| = e^{\sigma T}$. So, we need to map $\text{Re } s = \sigma$ to $|z| = e^{\sigma T}$. Considering the bilinear transform equation (3.5), we divide z by $e^{\sigma T}$ to make sure $ze^{-\sigma T}$ has a unit magnitude and shift the s -plane result by σ :

$$s = \sigma + \frac{2}{T} \cdot \frac{ze^{-\sigma T} - 1}{ze^{-\sigma T} + 1} \quad (3.22)$$

We can refer to (3.22) as the *unstable bilinear transform*, where the word “unstable” refers not to the instability of the transform itself, but rather to the fact that it is designed to be applied to unstable filters.²⁶ Notice that at $\sigma = 0$ the unstable bilinear transform turns into an ordinary bilinear transform. The inverse transform is obtained by

$$\frac{(s - \sigma)T}{2}(ze^{-\sigma T} + 1) = ze^{-\sigma T} - 1$$

from where

$$ze^{-\sigma T} \left(1 - \frac{(s - \sigma)T}{2}\right) = 1 + \frac{(s - \sigma)T}{2}$$

and

$$z = e^{\sigma T} \frac{1 + \frac{(s - \sigma)T}{2}}{1 - \frac{(s - \sigma)T}{2}} \quad (3.23)$$

Apparently the inverse unstable bilinear transform (3.23) has a pole at $s = \sigma + \frac{2}{T}$. In order to avoid hitting that pole by the poles of the filter’s transfer function (or maybe even generally avoid the real parts of the poles to go past that value) we could e.g. simply let

$$\sigma = \max\{0, \text{Re } p_n\}$$

or we could position σ midway:

$$\sigma = \max\left\{0, \text{Re } p_n - \frac{1}{T}\right\}$$

In order to construct an integrator defined by (3.22) we first need to obtain the expression for $1/s$ from (3.22):

$$\begin{aligned} \frac{1}{s} &= \frac{1}{\sigma + \frac{2}{T} \cdot \frac{ze^{-\sigma T} - 1}{ze^{-\sigma T} + 1}} = T \frac{ze^{-\sigma T} + 1}{\sigma T(ze^{-\sigma T} + 1) + 2(ze^{-\sigma T} - 1)} = \\ &= T \frac{ze^{-\sigma T} + 1}{(\sigma T + 2)e^{-\sigma T}z + (\sigma T - 2)} = T \frac{1 + e^{\sigma T}z^{-1}}{(\sigma T + 2) - (2 - \sigma T)e^{\sigma T}z^{-1}} = \\ &= \frac{T}{2 + \sigma T} \cdot \frac{1 + e^{\sigma T}z^{-1}}{1 - \frac{2 - \sigma T}{2 + \sigma T}e^{\sigma T}z^{-1}} \end{aligned}$$

That is

$$\frac{1}{s} = \frac{T}{2 + \sigma T} \cdot \frac{1 + e^{\sigma T}z^{-1}}{1 - \frac{2 - \sigma T}{2 + \sigma T}e^{\sigma T}z^{-1}} \quad (3.24)$$

²⁶Apparently, the unstable bilinear transform defines the same relationship between $\text{Im } s$ and $\text{arg } z$ as the ordinary bilinear transform. Therefore the standard prewarping formula applies.

A discrete-time structure implementing (3.24) could be e.g. the one in Fig. 3.24. Yet another approach could be to convert the right-hand side of (3.24) to the analog domain by the inverse bilinear transform, construct an analog implementation of the resulting transfer function and apply the trapezoidal integrator replacement to convert back to the digital domain. It is questionable, whether this produces better (or even different) results than Fig. 3.24.

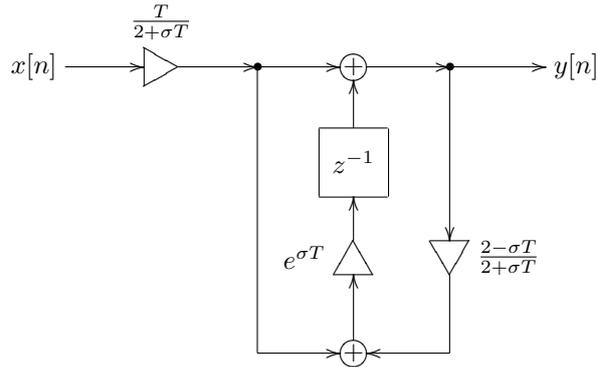


Figure 3.24: Transposed direct form II-style “unstable” trapezoidal integrator.

SUMMARY

We have considered three essentially different approaches to applying time-discretization to analog filter models: naive, TPT (by trapezoidal integrator replacement), and the classical bilinear transform (using direct forms). The TPT approach combines the best features of the naive implementation and the classical bilinear transform.

Chapter 4

Ladder filter

In this chapter we are going to discuss the most classical analog filter model: the transistor ladder filter. We will also discuss to an extent the diode ladder version.

4.1 Linear analog model

The analog transistor ladder filter is an essentially nonlinear structure. However, as the first approximation (and actually a quite good one) we will use its linearized model (Fig. 4.1).¹ The LP_1 blocks denote four identical (same cutoff) 1-pole lowpass filters (Fig. 2.2). The k coefficient controls the amount of negative feedback, which affects the filter resonance. Typically $k \geq 0$, although $k < 0$ is also sometimes used.²

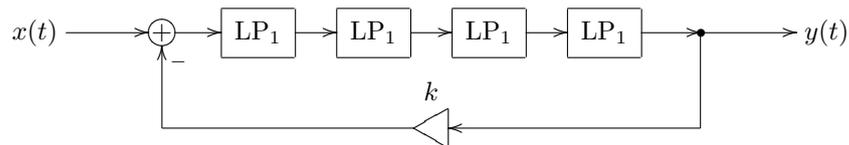


Figure 4.1: Transistor ladder filter.

Let

$$H_1(s) = \frac{1}{1+s}$$

¹A widely known piece of work describing this linear model is *Analyzing the Moog VCF with considerations for digital implementation* by T.Stilson and J.Smith.

²The reason for the negative (rather than positive) feedback is actually quite intuitively simple. Considering the phase response of four serially connected 1-pole lowpass filters at the cutoff:

$$\left(\frac{1}{1+s} \right)^4 \Big|_{s=j} = \frac{1}{(1+j)^4} = -\frac{1}{4}$$

we notice that the signal phase at the cutoff is inverted. Therefore we have to invert it once again in the feedback to achieve the resonance effect.

be the 1-pole lowpass transfer function. Assuming complex exponential x and y we write

$$y = H_1^4(s) \cdot (x - ky)$$

from where

$$y(1 + kH_1^4(s)) = H_1^4(s) \cdot x$$

and the transfer function of the ladder filter is

$$H(s) = \frac{y}{x} = \frac{H_1^4(s)}{1 + kH_1^4(s)} = \frac{\frac{1}{(1+s)^4}}{1 + k\frac{1}{(1+s)^4}} = \frac{1}{k + (1+s)^4} \quad (4.1)$$

At $k = 0$ the filter behaves as 4 serially connected 1-pole lowpass filters.

The poles of the filter are respectively

$$s = -1 + (-k)^{1/4}$$

where the raising to the 1/4th power is understood in the complex sense, therefore giving 4 different values:

$$s = -1 + \frac{\pm 1 \pm j}{\sqrt{2}} k^{1/4} \quad (k \geq 0) \quad (4.2)$$

(this time $k^{1/4}$ is understood in the real sense). Therefore, at $k = 0$ all poles are located at $s = -1$, as k grows they move apart in 4 straight lines (all going at “45° angles”). As k grows from 0 to 4 the two of the poles (at $s = -1 + \frac{1 \pm j}{\sqrt{2}} k^{1/4}$) are moving towards the imaginary axis, producing a resonance peak in the amplitude response (Fig. 4.2). At $k = 4$ they hit the imaginary axis:

$$\operatorname{Re} \left(-1 + \frac{1 \pm j}{\sqrt{2}} 4^{1/4} \right) = 0$$

and the filter becomes unstable.

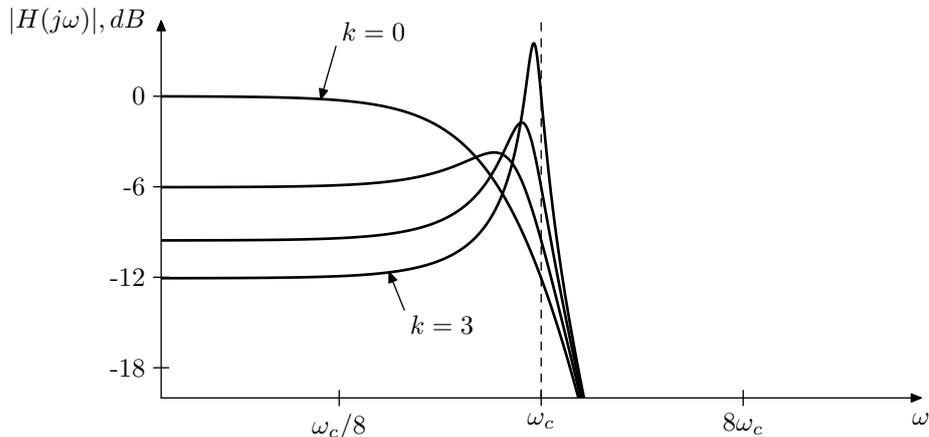


Figure 4.2: Amplitude response of the ladder filter for various k .

In Fig. 4.2 one could notice that, as the resonance increases, the filter gain at low frequencies begins to drop. Indeed, substituting $s = 0$ into (4.1) we obtain

$$H(0) = \frac{1}{1+k}$$

This is a general issue with ladder filter designs.

4.2 Linear digital model

A naive digital implementation of the ladder filter shouldn't pose any problems. We will therefore immediately skip to the TPT approach.

Recalling the instantaneous response of a single 1-pole lowpass filter (3.14), we can construct the instantaneous response of a serial connection of four of such filters. Indeed, let's denote the instantaneous responses of the respective 1-poles as $f_n(\xi) = g\xi + s_n$ (obviously, the coefficient g is identical for all four, whereas s_n depends on the filter state and therefore cannot be assumed identical). Combining two such filters in series we have

$$f_2(f_1(\xi)) = g(g\xi + s_1) + s_2 = g^2\xi + gs_1 + s_2$$

Adding the third one:

$$f_3(f_2(f_1(\xi))) = g(g^2\xi + gs_1 + s_2) + s_3 = g^3\xi + g^2s_1 + gs_2 + s_3$$

and the fourth one:

$$\begin{aligned} f_4(f_3(f_2(f_1(\xi)))) &= g(g^3\xi + g^2s_1 + gs_2 + s_3) = \\ &= g^4\xi + g^3s_1 + g^2s_2 + gs_3 + s_4 = G\xi + S \end{aligned}$$

where

$$\begin{aligned} G &= g^4 \\ S &= g^3s_1 + g^2s_2 + gs_3 + s_4 \end{aligned}$$

Using the obtained instantaneous response $G\xi + S$ of the series of 4 1-poles, we can redraw the ladder filter structure as in Fig. 4.3.

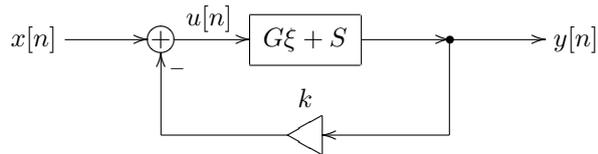


Figure 4.3: TPT ladder filter in the instantaneous response form.

Rather than solving for y , let's solve for the signal u at the feedback point. From Fig. 4.3 we obtain

$$u = x - ky = x - k(Gu + S)$$

from where

$$u = \frac{x - kS}{1 + kG} \quad (4.3)$$

We can then use the obtained value of u to process the 1-pole lowpasses one after the other, updating their state, and computing $y[n]$ as the output of the fourth lowpass.

Notice that for positive cutoff values of the underlying 1-pole lowpasses we have $g > 0$. Respectively $G = g^4 > 0$. This means that for $k \geq 0$ the denominator of (4.3) is always positive and never turns to zero, so we should be safe regarding the instantaneously unstable feedback.³

For $k < 0$ the situation is however different. Since $0 < g < 1$ (for $\omega_c > 0$), it follows that $0 < G < 1$. Thus $1 + kG > 0 \forall k \geq -1$, however at $k < -1$ we can get into an instantaneously unstable feedback case.

4.3 Feedback shaping

We have observed that at high resonance settings the amplitude gain of the filter at low frequencies drops (Fig. 4.2). An obvious way to fix this problem would be e.g. to boost the input signal by the $(1 + k)$ factor.⁴ However there's another way to address the same issue. We could “kill” the resonance at the low frequencies by introducing a highpass filter in the feedback (Fig. 4.4). In the simplest case this could be a 1-pole highpass.

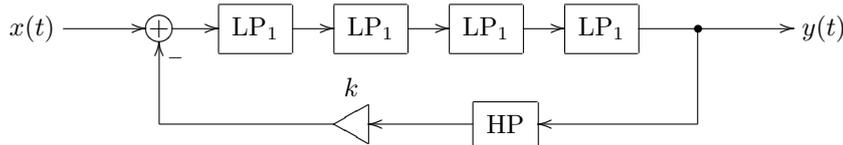


Figure 4.4: Transistor ladder filter with a highpass in the feedback.

The cutoff of the highpass filter can be static or vary along with the cutoff of the lowpasses. The static version has a nice feature that it kills the resonance effect at low frequencies regardless of the master cutoff setting, which may be desirable if the resonance at low frequencies is considered rather unpleasant (Fig. 4.5).

In principle one can also use other filter types in the feedback shaping. One has to be careful though, since this changes the positions of the filter poles. Particularly, inserting a lowpass into the feedback can easily destabilize an otherwise stable filter.

4.4 Multimode ladder filter

Warning! *The multimode functionality of the ladder filter is a rather exotic feature. If you're looking for the bread-and-butter bandpass, highpass, notch etc.*

³Strictly speaking, we should have checked the instantaneous stability using the feedback smoother approach. However typically a positive denominator of the form $1 + g$ or $1 + kG$ implies that everything is fine.

A quicker way to check for the instantaneous feedback would be to let the sampling rate infinitely grow ($T \rightarrow 0$) and then check, whether the denominator changes its sign along the way. In our case $G = g^4 = (\omega_c T/2)^4$, which means the denominator is always larger than 1 (under the assumption $k \geq 0$), regardless of T .

⁴We boost the input rather than the output signal for the same reason as when preferring to place the cutoff gains in front of the integrators.

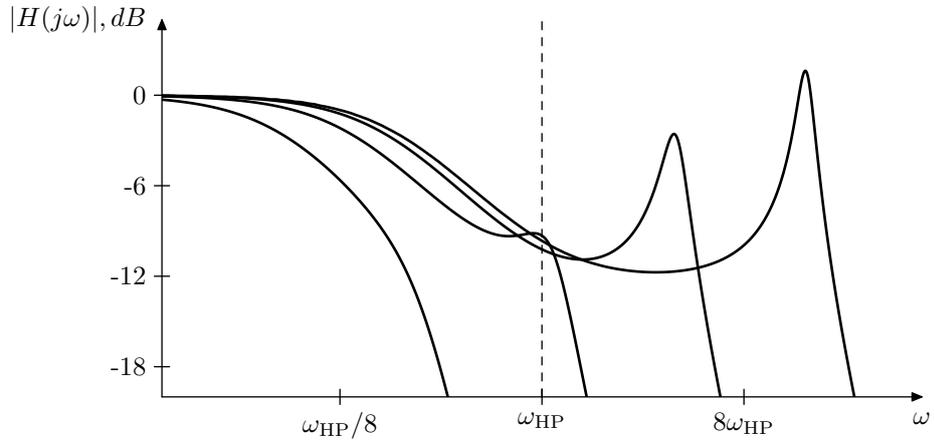


Figure 4.5: Amplitude response of the ladder filter with a static-cutoff highpass in the feedback for various lowpass cutoffs.

filters, you should first take a look at the multimode 2-pole state-variable filter discussed later in the book.

By picking up intermediate signals of the ladder filter as in Fig. 4.6 we obtain the multimode version of this filter. We then can use linear combinations of signals y_n to produce various kinds of filtered signal.⁵

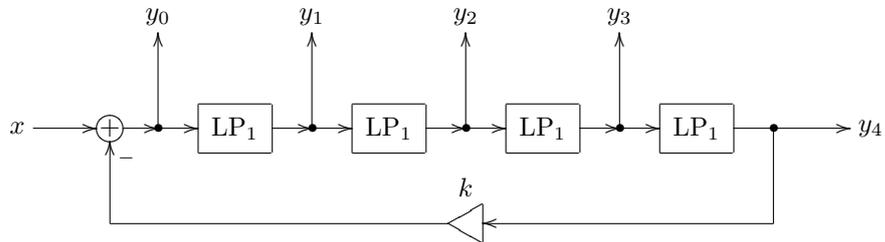


Figure 4.6: Multimode ladder filter.

Suppose $k = 0$. Apparently, in this case, the respective transfer functions associated with each of the y_n outputs are

$$H_n(s) = \frac{1}{(1 + s)^n} \quad (n = 0, \dots, 4) \tag{4.4}$$

If $k \neq 0$ then from

$$H_4(s) = \frac{1}{k + (1 + s)^4}$$

⁵ Actually, instead of y_0 we could have used the input signal x for these linear combinations. However, it doesn't matter. Since $y_0 = x - ky_4$, we can express x via y_0 or vice versa. It's just that some useful linear combinations have simpler (independent of k) coefficients if y_0 rather than x is being used.

using the obvious relationship $H_{n+1}(s) = H_n(s)/(s + 1)$ we obtain

$$H_n(s) = \frac{(1+s)^{4-n}}{k + (1+s)^4} \quad (4.5)$$

4-pole highpass mode

Considering that the 4th order lowpass transfer function (under the assumption $k = 0$) is built as a product of four 1st order lowpass transfer functions $1/(1+s)$

$$H_{LP}(s) = \frac{1}{(1+s)^4}$$

we might decide to build the 4th order highpass transfer function as a product of four 1st order highpass transfer functions $s/(1+s)$:

$$H_{HP}(s) = \frac{s^4}{(1+s)^4}$$

Let's attempt to build $H_{HP}(s)$ as a linear combination of $H_n(s)$. Apparently, a linear combination of $H_n(s)$ must have the denominator $k + (1+s)^4$, so let's instead construct

$$H_{HP}(s) = \frac{s^4}{k + (1+s)^4} \quad (4.6)$$

which at $k = 0$ will turn into $s^4/(1+s)^4$:

$$\frac{s^4}{k + (1+s)^4} = \frac{a_0(1+s)^4 + a_1(1+s)^3 + a_2(1+s)^2 + a_3(1+s) + a_4}{k + (1+s)^4}$$

that is

$$s^4 = a_0(1+s)^4 + a_1(1+s)^3 + a_2(1+s)^2 + a_3(1+s) + a_4$$

Formally replacing $s + 1$ by s (and respectively s by $s - 1$):

$$(s - 1)^4 = a_0s^4 + a_1s^3 + a_2s^2 + a_3s + a_4$$

From where immediately

$$a_0 = 1, a_1 = -4, a_2 = 6, a_3 = -4, a_4 = 1$$

The amplitude response corresponding to (4.6) is plotted in Fig. 4.7.

4-pole bandpass mode

A bandpass filter can be built as

$$H_{BP}(s) = \frac{s^2}{k + (1+s)^4} \quad (4.7)$$

The two zeros at $s = 0$ will provide for a -12dB/oct rolloff at low frequencies and will reduce the -24dB/oct rolloff at high frequencies to the same -12dB/oct . Notice that the phase response at the cutoff is zero:

$$H_{BP}(j) = \frac{-1}{k + (1+j)^4} = \frac{1}{4-k}$$

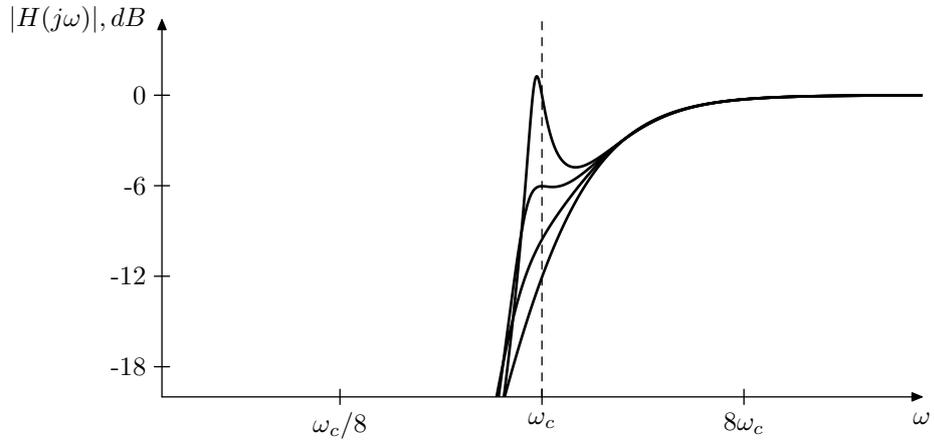


Figure 4.7: Amplitude response of the highpass mode of the ladder filter for various k .

The coefficients are found from

$$\begin{aligned} s^2 &= a_0(1+s)^4 + a_1(1+s)^3 + a_2(1+s)^2 + a_3(1+s) + a_4 \\ (s-1)^2 &= a_0s^4 + a_1s^3 + a_2s^2 + a_3s + a_4 \end{aligned}$$

The amplitude response corresponding to (4.7) is plotted in Fig. 4.8.

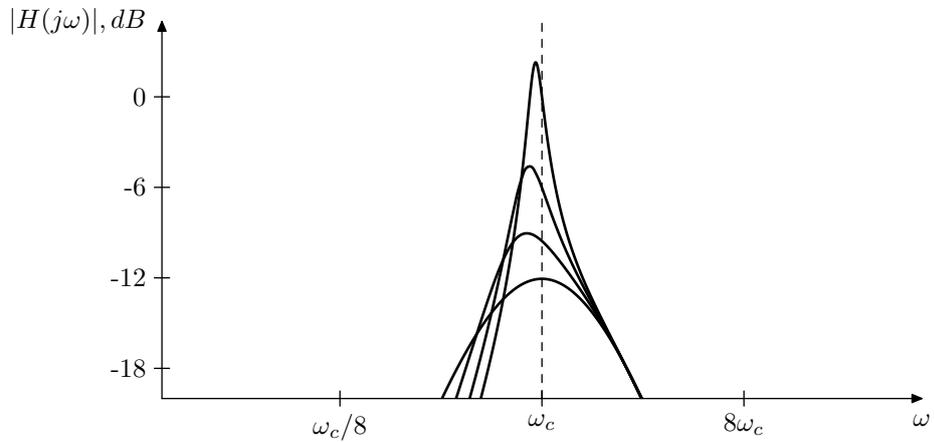


Figure 4.8: Amplitude response of the bandpass mode of the ladder filter for various k .

Lower-order modes

Recalling the transfer functions of the modal outputs y_n in the absence of the resonance (4.4), we can consider the modal signals y_n and their respective transfer functions (4.5) as a kind of “ n -pole lowpass filters with 4-pole resonance”.

“Lower-order” highpasses can be build by considering the zero-resonance

transfer functions

$$H_{\text{HP}}(s) = \frac{s^N}{(s+1)^N} = \frac{(s+1)^{4-N} s^N}{(s+1)^4}$$

which for $k \neq 0$ turn into

$$H_{\text{HP}}(s) = \frac{(s+1)^{4-N} s^N}{k + (s+1)^4}$$

In a similar way we can build a “2-pole” bandpass

$$H_{\text{BP}}(s) = \frac{s}{(s+1)^2} = \frac{(s+1)^2 s}{(s+1)^4} \quad (k=0)$$

$$H_{\text{BP}}(s) = \frac{(s+1)^2 s}{k + (s+1)^4} \quad (k \neq 0)$$

Other modes

Continuing in the same fashion we can build further modes (the transfer functions are given for $k=0$):

$\frac{s}{(s+1)^3}$	3-pole bandpass, 6/12 dB/oct
$\frac{s^2}{(s+1)^3}$	3-pole bandpass, 12/6 dB/oct
$\frac{(s+1)^4 + Ks^2}{(s+1)^4}$	band-shelving
$\frac{(s'^2 + \omega_0)^2}{(s+1)^4}$	notch ($s' = s/\sqrt{\omega_0}$, where ω_0 is the notch frequency)
$\frac{(s'^2 + 2Rs + \omega_0)^2 + (s'^2 - 2Rs + \omega_0)^2}{2(s+1)^4}$	2 notches (R affects the notch spreading, neutral setting $R=1$)
$\frac{s'^2 + 1}{(s+1)^4}$	2-pole lowpass + notch ($s' = s/\omega_0$)
$\frac{(1 + 1/s'^2)s^4}{(s+1)^4}$	2-pole highpass + notch
$\frac{(s' + 1/s')s^2}{(s+1)^4}$	2-pole bandpass + notch

etc. Notably, these modes are not necessarily perfectly matching their descriptions, the parameters may have some weird side effects.

4.5 HP and BP ladders

Performing an LP to HP transformation on the lowpass ladder filter we effectively perform it on each of the underlying 1-pole lowpasses, thus turning

them into 1-pole highpasses. Thereby we obtain a “true” highpass ladder filter (Fig. 4.9). Obviously, the amplitude response of the ladder highpass is symmetric to the amplitude response of the ladder lowpass.

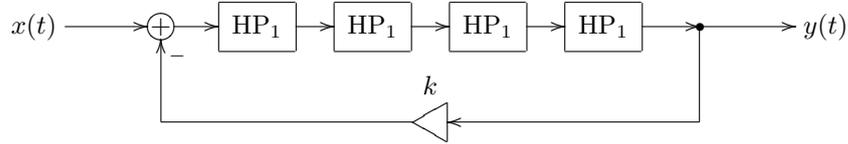


Figure 4.9: A “true” highpass ladder filter.

In order to build a “true” 4-pole bandpass ladder, we replace only half of the lowpasses with highpasses (it doesn’t matter which two of the four 1-pole lowpasses are replaced). The total transfer function of the feedforward path is thereby

$$\frac{s^2}{(1+s)^4} = \frac{s}{(1+s)^2} \cdot \frac{s}{(1+s)^2}$$

where each of the $s/(1+s)^2$ factors is built from a serial combination of a 1-pole lowpass and a 1-pole highpass:

$$\frac{s}{(1+s)^2} = \frac{s}{1+s} \cdot \frac{1}{1+s}$$

Thus $s/(1+s)^2$ is a simple 2-pole bandpass⁶ and a serial combination of two of them makes a simple 4-pole bandpass. The frequency response of $s/(1+s)^2$ at $\omega = 1$ can be easily found and is equal to $1/2$, that is there is no phase-shift. Respectively the frequency response of $s^2/(1+s)^4$ at $\omega = 1$ is $1/4$, also without a phase shift. Therefore we need to use positive rather than negative feedback (Fig. 4.10). As the lowpass and the highpass ladders, the bandpass ladder becomes unstable at $k = 4$.

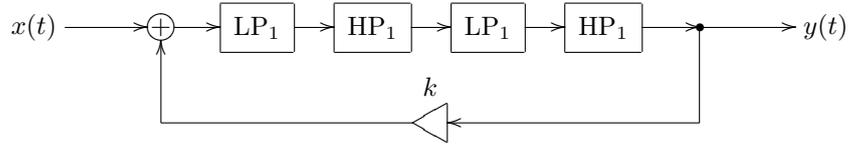


Figure 4.10: A “true” bandpass ladder filter.

Noticing that the filter structure is invariant relative to the LP to HP transformation, we conclude that its amplitude response must be symmetric (around $\omega = 1$) in the logarithmic frequency scale.

⁶A more appropriate and generic way to build a 2-pole bandpass is the multimode 2-pole state-variable filter discussed later in the book.

4.6 Simple nonlinear model

At $k \geq 4$ the ladder filter becomes unstable, as the resonance becomes too strong. We could however prevent the signal level from growing infinitely by putting a saturator into the feedback path. This will allow the filter to go into *selfoscillation* at $k > 4$. The best place for such saturator is probably at the feedback point, since then it will process both the input and the feedback signals simultaneously, applying a nice overdrive-like saturation to the input signal. A hyperbolic tangent function should provide a nice saturator (Fig. 4.11). It is transparent at low signal levels, therefore at low signal levels the filter behaves as a linear one.

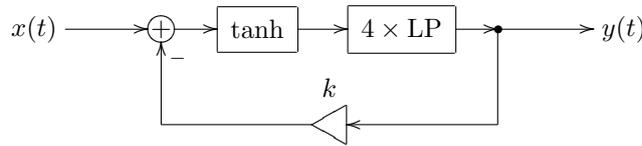


Figure 4.11: Transistor ladder filter with a saturator at the feedback point.

Another possibility is to place the saturator before the feedback point, which makes it somewhat more “transparent”, since the saturation doesn’t affect the feedforward path (Fig. 4.12). By swapping the positions of the saturator and the feedback gain k in Fig. 4.12 one can make saturation independent of the feedback amount setting k .

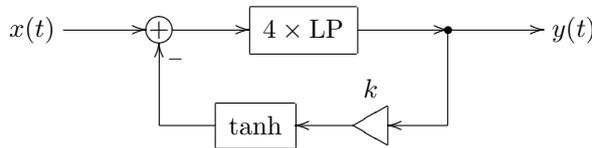


Figure 4.12: Transistor ladder filter with a saturator before the feedback point.

Other saturation curves are of course possible, where a noticeably different (smoother) saturation curve is provided by e.g. an inverse hyperbolic sine function, since this one doesn’t have horizontal asymptotes.

The introduction of the nonlinearity in the feedback path poses no problems for a naive digital model. In the TPT case however this complicates the things quite a bit. Consider Fig. 4.3 redrawn to contain the feedback nonlinearity (Fig. 4.13).

Writing the zero-delay feedback equation we obtain

$$u = x - k(G \tanh u + s) \quad (4.8)$$

Apparently, the equation (4.8) is a transcendental one. It can be solved only using numerical methods. Also, a linear zero-delay feedback equation had only one solution, but how many solutions does (4.8) have? In order to answer the

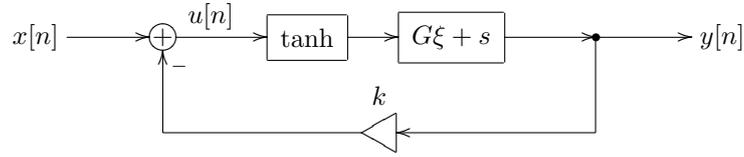


Figure 4.13: Nonlinear TPT ladder filter in the instantaneous response form.

latter question, let's rewrite (4.8) as

$$(x - ks) - u = kG \tanh u \quad (4.9)$$

If $k \geq 0$ and $G > 0$, then the right-hand side of (4.9) is a nonstrictly increasing function of u , while the left-hand side of (4.9) is a strictly decreasing function of u . Thus, (4.9) and respectively (4.8) have a single solution in this case. If $k < 0$, (4.8) can have multiple solutions (up to three). One could use the smoother paradigm introduced in the instantaneously unstable feedback discussion to find out the applicable one.

It is possible to avoid the need of solving the transcendental equation by using a saturator function which still allows analytic solution. This is particularly the case with second-order curves, such as hyperbolas. E.g. $y = \tanh x$ can be replaced by $y = x/(1 + |x|)$, while the inverse of $x = \sinh y$ can be replaced by the inverse of $x = x(1 + |x|)$. Each of these two functions consists of two second-order segments, which turns (4.9) into a quadratic equation. E.g. for $x/(1 + |x|)$ we have

$$(x - ks) - u = kG \frac{u}{1 + |u|} \quad (4.10)$$

In order to solve (4.10) one first has to check whether the solution occurs at $u > 0$ or $u < 0$, which (for $kG \geq 0$) can be done by simply evaluating the left-hand side of (4.10) at $u = 0$. Then one can replace $|u|$ by u or $-u$ respectively and solve the resulting quadratic equation.⁷

Yet another approach (which also works for multiple nonlinearities!) is to first solve the feedback equation for the linear case, and then apply the nonlinearities "on top". E.g. we use the structure in Fig. 4.3 to obtain the value of u . However then we pretend that we have found the value of u for Fig. 4.13 (or Fig. 4.11, for that matter) and proceed accordingly, putting u through the hyperbolic tangent shaper and then further through the 1-pole lowpasses. We

⁷The same kind of quadratic equation appears for any other second-order curve (hyperbola, ellipse, parabola, including their rotated versions). In solving the quadratic equation one has not only to choose the right one of the two roots of the equation. One also has to choose the right one of the two solution formulas for the quadratic equation $Ax^2 - 2Bx + C = 0$:

$$x = \frac{B \pm \sqrt{B^2 - AC}}{A} \quad \text{or} \quad x = \frac{C}{B \mp \sqrt{B^2 - AC}}$$

The reason to choose between these two different formulas is that each of them can become ill-conditioned depending on the values of A , B and C . The choice of the formula is determined by the sign of B and the relative magnitudes of the equation coefficients. More details on this approach are discussed in the author's article "Computational optimization of nonlinear zero-delay feedback by second-order piecewise approximation".

refer to this approach as the “cheap method” of applying the nonlinearities to the TPT structures. It is intuitively clear, that the cheap method is more likely to produce “wrong” results at high cutoff values.

No matter, which approach we chose to compute nonlinearities, one shouldn’t forget that nonlinear shaping introduces overtones (usually an infinite amount of those) into the signal, which in turn introduces aliasing. Meaning: the stronger are the nonlinearities in your structure, the more you might need to oversample. If the oversampling is extreme anyway, there might be little difference in quality between the naive and the TPT approach.⁸

4.7 Advanced nonlinear model

The nonlinearity introduced in Fig. 4.11 does a good job and sounds reasonably close to a hardware analog transistor ladder filter, however this is not how the nonlinearities in the hardware ladder filter “really” work. In order to describe a closer to reality nonlinear model of the ladder filter, we need to start by introducing nonlinearities into the underlying 1-pole lowpass filters (Fig. 4.14).⁹

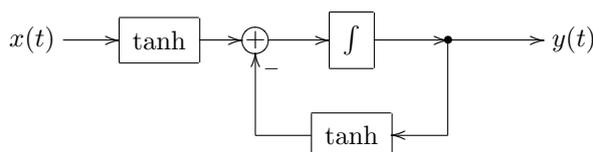


Figure 4.14: A nonlinear 1-pole lowpass element of the ladder filter.

So the equation (2.1) is transformed into

$$y = y(t_0) + \int_{t_0}^t \omega_c (\tanh x(\tau) - \tanh y(\tau)) dt$$

Which effect does this have? Apparently, the presence of the tanh function reduces the absolute value of the difference $\tanh x - \tanh y$, if the level of one or both of the signals is sufficiently high. If both x and y have large values of the same sign, it’s possible that the difference $\tanh x - \tanh y$ is close to zero, even though the difference $x - y$ is very large. This means that the filter will update its state slower than in (2.1). Intuitively this feels a little bit like “cutoff reduction” at large signal levels.

We can then connect the 1-pole models from Fig. 4.14 into a series of four 1-poles and put a feedback around them, exactly the same way as in Fig. 4.1. Notice that when connecting Fig. 4.14 filters in series, one could use a common tanh module between each of them, thereby optimizing the computation (Fig. 4.15).¹⁰

⁸Before making a final decision, it might be worth asking a few musicians with an ear for analog sound to perform a listening test, whether the differences between the naive and TPT models of a particular filter are inaudible and uncritical.

⁹A famous piece of work describing this specific nonlinear model is the DAFx’04 paper *Non-linear digital implementation of the Moog ladder filter* by Antti Huovilainen. Therefore this model is sometimes referred to as the “Antti’s model”.

¹⁰There is an issue which may appear when using simple tanh approximations having a fully

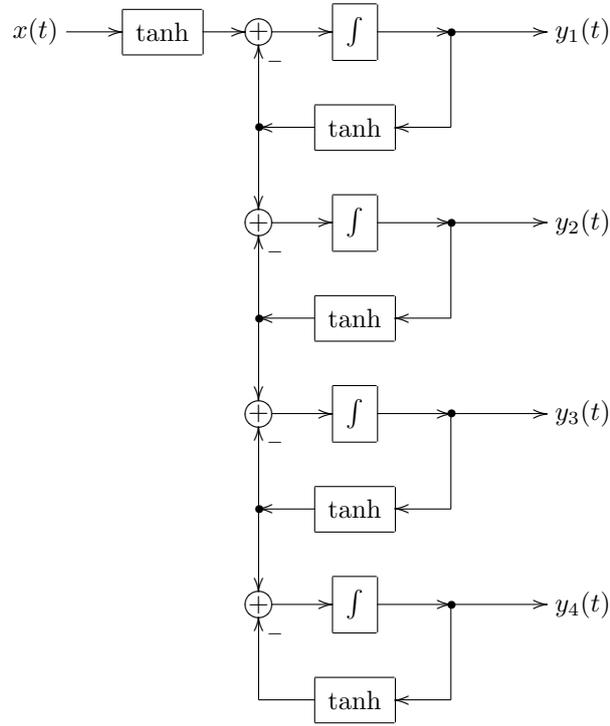


Figure 4.15: Advanced nonlinear transistor ladder (the main feedback path of the ladder filter not shown).

One could further enhance the nonlinear behavior of the ladder filter model by putting another saturator (possibly of a different type, or simply differently scaled) into the feedback path.

4.8 Diode ladder

In the diode ladder filter the serial connection of four 1-pole lowpass filters (implemented by the transistor ladder) is replaced by a more complicated structure of 1-pole filters (implemented by the diode ladder). The equations of the diode ladder itself (without the feedback path) are

$$\begin{aligned}\dot{y}_1 &= \omega_c (\tanh x - \tanh(y_1 - y_2)) \\ \dot{y}_2 &= \frac{\omega_c}{2} (\tanh(y_1 - y_2) - \tanh(y_2 - y_3)) \\ \dot{y}_3 &= \frac{\omega_c}{2} (\tanh(y_2 - y_3) - \tanh(y_3 - y_4)) \\ \dot{y}_4 &= \frac{\omega_c}{2} (\tanh(y_3 - y_4) - \tanh y_4)\end{aligned}$$

horizontal saturation curve. If both the input and the output signals of a 1-pole are having large values of the same sign, the tanh approximations will return two identical values and the difference $\tanh x - \tanh y$ will be approximated by zero. This might result in the filter getting “stuck” (the cutoff effectively reduced to zero).

which is implemented by the structure in Fig. 4.16 (compare to Fig. 4.15). The diode ladder itself is then built by providing the feedback path around the diode ladder, where the fourth output of the diode ladder is fed back into the diode ladder's input (Fig. 4.17).

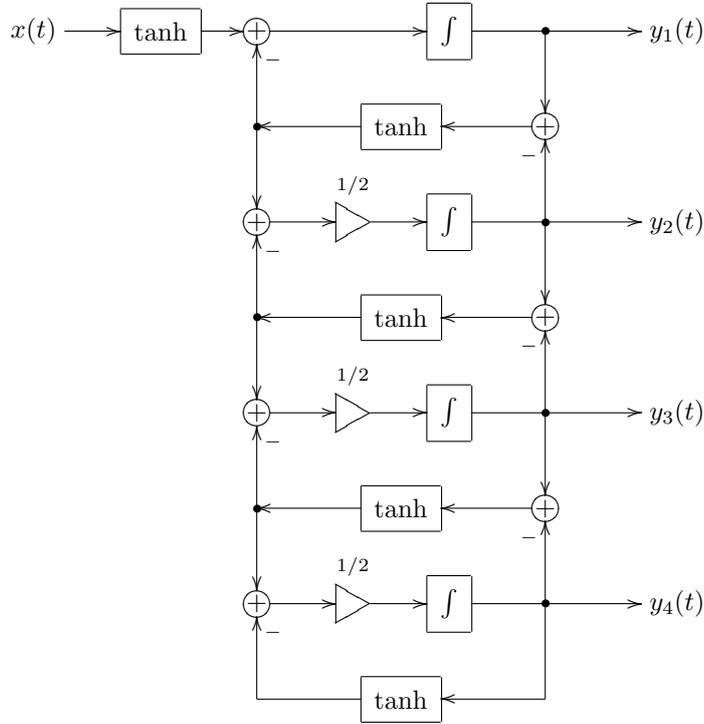


Figure 4.16: Diode ladder.

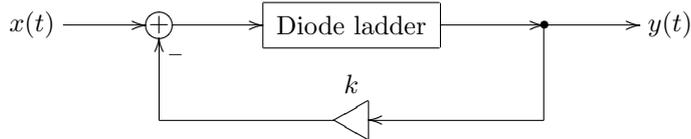


Figure 4.17: Diode ladder filter.

The linearized form of the diode ladder equations is obtained by assuming $\tanh \xi \approx \xi$, resulting in

$$\begin{aligned}\dot{y}_1 &= \omega_c((x + y_2) - y_1) \\ \dot{y}_2 &= \omega_c((y_1 + y_3)/2 - y_2) \\ \dot{y}_3 &= \omega_c((y_2 + y_4)/2 - y_3) \\ \dot{y}_4 &= \omega_c(y_3/2 - y_4)\end{aligned}$$

which apparently is representable as a serial connection of four identical 1-pole lowpass filters (all having the same cutoff ω_c) with some extra gains and feedback

from where

$$y_4 \left(1 - \frac{F^2}{1 - F^2} \right) = \frac{F^3}{1 - F^2} y_1$$

and respectively

$$y_4 = \frac{F^3}{1 - 2F^2} y_1$$

And finally,

$$y_1 = 2F \cdot (x + y_2)$$

Multiplying both sides by $F^3/(1 - 2F^2)$:

$$\begin{aligned} y_4 &= \frac{2F^4}{1 - 2F^2} \left(x + \frac{1 - F^2}{F^2} \cdot \frac{F^2}{1 - F^2} y_2 \right) = \frac{2F^4}{1 - 2F^2} \left(x + \frac{1 - F^2}{F^2} y_4 \right) = \\ &= \frac{2F^4}{1 - 2F^2} x + 2F^2 \frac{1 - F^2}{1 - 2F^2} y_4 \end{aligned}$$

from where

$$y_4 \left(1 - 2F^2 \frac{1 - F^2}{1 - 2F^2} \right) = \frac{2F^4}{1 - 2F^2} x$$

from where

$$y_4 (1 - 4F^2 + 2F^4) = 2F^4 x$$

and

$$y_4 = \frac{2F^4}{1 - 4F^2 + 2F^4} x = \frac{G^4/8}{1 - G^2 + G^4/8} x$$

That is, the transfer function $\Delta(s)$ of the diode ladder is

$$\Delta(s) = \frac{G^4/8}{G^4/8 - G^2 + 1} \quad \text{where } G(s) = \frac{1}{1 + s}$$

from where we obtain the transfer function of the entire diode ladder filter as

$$H(s) = \frac{\Delta}{1 + k\Delta}$$

The corresponding amplitude response is plotted in Fig. 4.19.

Let's find the positions of the poles of the diode ladder filter. Equating the denominator to zero:

$$1 + k\Delta = 0$$

we have

$$1 = -k\Delta = \frac{-kG^4/8}{G^4/8 - G^2 + 1}$$

that is

$$-k \frac{G^4}{8} = \frac{G^4}{8} - G^2 + 1$$

that is

$$\frac{1 + k}{8} G^4 - G^2 + 1 = 0$$

Solving for G^2 :

$$G^2 = \frac{1 \pm \sqrt{1 - \frac{1+k}{2}}}{\frac{1+k}{4}} = \frac{1 \pm \sqrt{\frac{1-k}{2}}}{\frac{1+k}{4}} \quad (4.12)$$

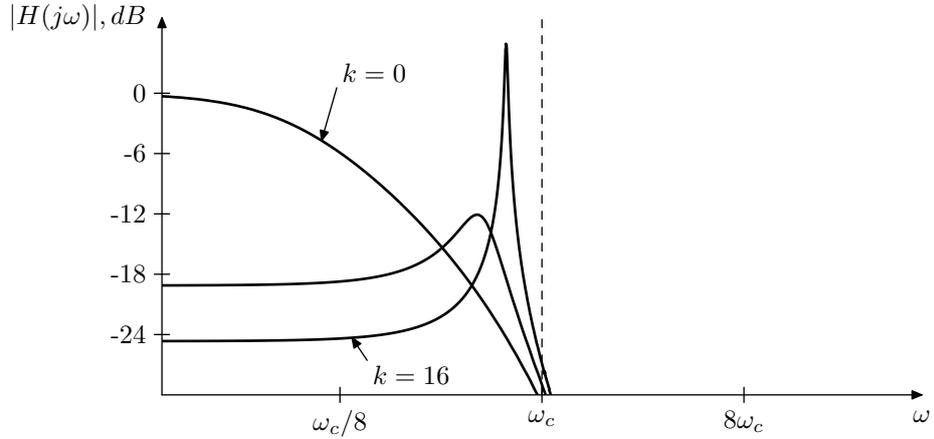


Figure 4.19: Amplitude response of the diode ladder filter for various k .

Equating (4.12) and the squared form of (4.11) we have

$$\frac{1}{(1+s)^2} = \frac{1 \pm \sqrt{\frac{1-k}{2}}}{\frac{1+k}{4}}$$

that is

$$\begin{aligned} (s+1)^2 &= \frac{\frac{1+k}{4}}{1 \pm \sqrt{\frac{1-k}{2}}} = \frac{\frac{1+k}{4} \left(1 \mp \sqrt{\frac{1-k}{2}}\right)}{1 - \frac{1-k}{2}} = \\ &= \frac{\frac{1+k}{4} \left(1 \mp \sqrt{\frac{1-k}{2}}\right)}{\frac{1+k}{2}} = \frac{1 \mp \sqrt{\frac{1-k}{2}}}{2} \end{aligned}$$

that is

$$(s+1)^2 = \frac{1}{2} \pm \frac{1}{2} \sqrt{\frac{1-k}{2}} \quad (4.13)$$

The equation (4.13) defines the positions of the poles of the diode ladder filter. Apparently, it's easily solvable. We would be interested to find out, at which k does the selfoscillation start.

If the poles are to be on the imaginary axis, then $s = j\omega$. Substituting $s = j\omega$ into (4.13) we get

$$(1 - \omega^2) + 2j\omega = \frac{1}{2} \mp \frac{j}{2} \sqrt{\frac{k-1}{2}} \quad (4.14)$$

Equating the real parts of (4.14) we obtain $1 - \omega^2 = \frac{1}{2}$ and $\omega = \pm \frac{1}{\sqrt{2}}$. Equating the imaginary parts of (4.14) and substituting $\omega = \pm \frac{1}{\sqrt{2}}$ we obtain

$$\pm \frac{2}{\sqrt{2}} = \pm \frac{1}{2} \sqrt{\frac{k-1}{2}}$$

from where

$$\pm 4 = \pm\sqrt{k-1}$$

and, since $k \in \mathbb{R}$, we have $k = 17$.

That is, given the unit cutoff of the underlying one-pole lowpass filters, the selfoscillation starts at $k = 17$, where the resonance peak is located at $\omega = 1/\sqrt{2}$.

TPT model

Converting Fig. 4.18 to the instantaneous response form we obtain the structure in Fig. 4.20. From Fig. 4.20 we wish to obtain the instantaneous response of the entire diode ladder. Then we could use this response to solve the zero-delay feedback equation for the main feedback loop of Fig. 4.17.

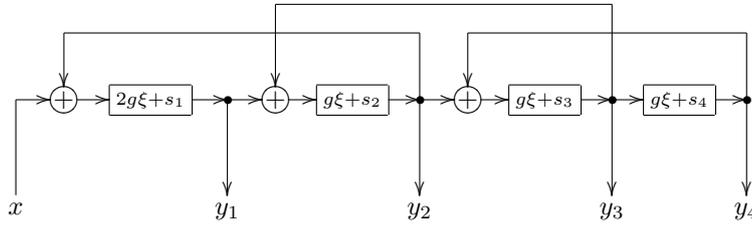


Figure 4.20: Linearized diode ladder in the instantaneous response form.

From Fig. 4.20 we have $y_4 = gy_3 + s_4$. We can rewrite it as

$$y_4 = G_4 y_3 + S_4 \quad \text{where } G_4 = g, \quad S_4 = s_4$$

Further, from Fig. 4.20 we also have

$$y_3 = g(y_2 + y_4) + s_3 = g(y_2 + G_4 y_3 + S_4) + s_3$$

from where

$$y_3 = \frac{gy_2 + gS_4 + s_3}{1 - gG_4} = G_3 y_2 + S_3 \quad \text{where } G_3 = \frac{g}{1 - gG_4}, \quad S_3 = \frac{gS_4 + s_3}{1 - gG_4}$$

Further,

$$y_2 = g(y_1 + y_3) + s_2 = g(y_1 + G_3 y_2 + S_3) + s_2$$

from where

$$y_2 = \frac{gy_1 + gS_3 + s_2}{1 - gG_3} = G_2 y_1 + S_2 \quad \text{where } G_2 = \frac{g}{1 - gG_3}, \quad S_2 = \frac{gS_3 + s_2}{1 - gG_3}$$

And ultimately

$$y_1 = 2g(x + y_2) + s_1 = 2g(x + G_2 y_1 + S_2) + s_1$$

from where

$$y_1 = \frac{2gx + 2gS_2 + s_1}{1 - 2gG_2} = G_1 x + S_1 \quad \text{where } G_1 = \frac{2g}{1 - 2gG_2}, \quad S_1 = \frac{2gS_2 + s_1}{1 - 2gG_2}$$

Thus, we have

$$y_n = G_n y_{n-1} + S_n \quad (\text{where } y_0 = x)$$

from where it's easy to obtain the instantaneous response of the entire diode ladder as

$$\begin{aligned} y_4 &= G_4 y_3 + S_4 = G_4(G_3 y_2 + S_3) + S_4 = G_4 G_3 y_2 + (G_4 S_3 + S_4) = \\ &= G_4 G_3 (G_2 y_1 + S_2) + (G_4 S_3 + S_4) = \\ &= G_4 G_3 G_2 y_1 + (G_4 G_3 S_2 + G_4 S_3 + S_4) = \\ &= G_4 G_3 G_2 (G_1 x + S_1) + (G_4 G_3 S_2 + G_4 S_3 + S_4) = \\ &= G_4 G_3 G_2 G_1 x + (G_4 G_3 G_2 S_1 + G_4 G_3 S_2 + G_4 S_3 + S_4) = Gx + S \end{aligned}$$

Notice, that we should have checked that we don't have instantaneously unstable feedback problems within the diode ladder. That is, we need to check that all denominators in the expressions for G_n and S_n don't turn to zero or to negative values. Considering that $0 < g < \frac{1}{2}$ and that $G_4 = g$, we have

$$0 < G_4 < \frac{1}{2}$$

and

$$1 - gG_4 = 1 - g^2 > \frac{3}{4}$$

Respectively

$$0 < G_3 = g/(1 - gG_4) < \frac{1/2}{3/4} = \frac{2}{3}$$

and

$$1 - gG_3 > 1 - \frac{1}{2} \cdot \frac{2}{3} = 1 - \frac{1}{3} = \frac{2}{3}$$

Then

$$0 < G_2 = g/(1 - gG_3) < \frac{1/2}{2/3} = \frac{3}{4}$$

and

$$1 - 2gG_2 > 1 - 2 \cdot \frac{1}{2} \cdot \frac{3}{4} = 1 - \frac{3}{4} = \frac{1}{4}$$

and thus all denominators are always positive.

Also

$$0 < G_1 = \frac{2g}{1 - 2gG_2} < \frac{1}{1/4} = 4$$

Thus

$$0 < G = G_4 G_3 G_2 G_1 < \frac{1}{2} \cdot \frac{2}{3} \cdot \frac{3}{4} \cdot 4 = 1$$

Using the obtained instantaneous response of the entire diode ladder we now can solve the main feedback equation for Fig. 4.17.

For a nonlinear diode ladder model we could use the structure in Fig. 4.16. However, it might be too complicated to process. Even the application of the "cheap" TPT nonlinear processing approach is not fully trivial.

One can therefore use simpler nonlinear structures instead, e.g. the one from Fig. 4.11. Also, the other ideas discussed for the transistor ladder can be applied. In regards to the multimode diode ladder filter, notice that the transfer functions corresponding to the $y_n(t)$ outputs are different from the ones of the transistor ladder, therefore the mixing coefficients which worked for the modes of the transistor ladder filter, are not going to work the same for the diode ladder.

SUMMARY

The transistor ladder filter model is constructed by placing a negative feedback around a chain of four identical 1-pole lowpass filters. The feedback amount controls the resonance. A nonlinearity in the feedback path (e.g. at the feedback point) could be used to contain the signal level, so that selfoscillation becomes possible.

Chapter 5

2-pole filters

The other classical analog filter model is the 2-pole filter design commonly referred to in the music DSP field as the *state-variable filter* (SVF). It can also serve as a generic analog model for building 2-pole filters, similarly to previously discussed 1-pole RC filter model.

5.1 Linear analog model

The block diagram of the state-variable filter is shown in Fig. 5.1. The three outputs are the highpass, bandpass and lowpass signals.¹ As usual, one can apply transposition to obtain a filter with highpass, bandpass and lowpass inputs (Fig. 5.2).

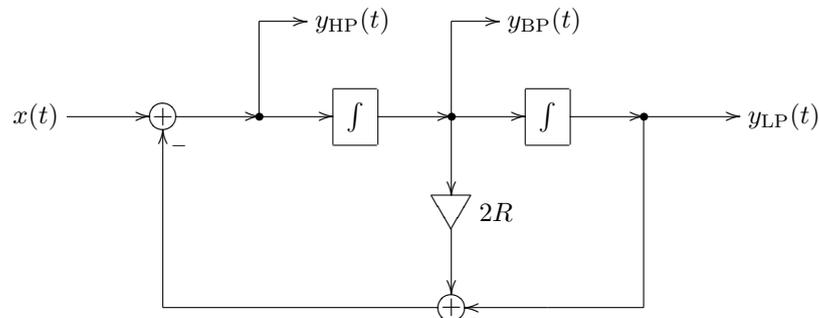


Figure 5.1: 2-pole multimode state-variable filter.

From Fig. 5.1 one can easily obtain the transfer functions for the respective signals. Assume complex exponential signals. Then, assuming unit cutoff,

$$y_{HP} = x - 2Ry_{BP} - y_{LP}$$

¹One can notice that the filter in Fig. 5.1 essentially implements an analog-domain canonical form, similar to the one in Fig. 3.20. Indeed let's substitute in Fig. 3.20 the z^{-1} elements by s^{-1} elements (integrators) and let $a_1 = -2R$, $a_2 = -1$. Then the gains b_0 , b_1 and b_2 are simply picking up the highpass, bandpass and lowpass signals respectively.

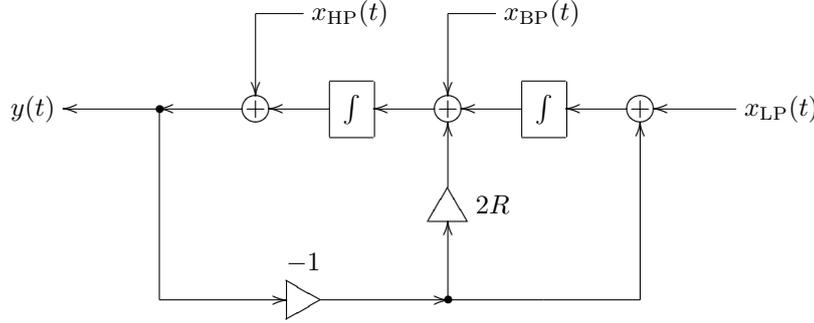


Figure 5.2: Transposed 2-pole multimode state-variable filter.

$$y_{\text{BP}} = \frac{1}{s} y_{\text{HP}}$$

$$y_{\text{LP}} = \frac{1}{s} y_{\text{BP}}$$

from where

$$y_{\text{HP}} = x - 2R \cdot \frac{1}{s} y_{\text{HP}} - \frac{1}{s^2} y_{\text{HP}}$$

from where

$$\left(1 + \frac{2R}{s} + \frac{1}{s^2}\right) y_{\text{HP}} = x$$

and

$$H_{\text{HP}}(s) = \frac{y_{\text{HP}}}{x} = \frac{1}{1 + \frac{2R}{s} + \frac{1}{s^2}} = \frac{s^2}{s^2 + 2Rs + 1}$$

Thus

$$H_{\text{HP}}(s) = \frac{s^2}{s^2 + 2Rs + 1} = \frac{s^2}{s^2 + 2R\omega_c s + \omega_c^2} \quad (\omega_c = 1)$$

$$H_{\text{BP}}(s) = \frac{s}{s^2 + 2Rs + 1} = \frac{\omega_c s}{s^2 + 2R\omega_c s + \omega_c^2} \quad (\omega_c = 1)$$

$$H_{\text{LP}}(s) = \frac{1}{s^2 + 2Rs + 1} = \frac{\omega_c^2}{s^2 + 2R\omega_c s + \omega_c^2} \quad (\omega_c = 1)$$

The respective amplitude responses are plotted in Figs. 5.3, 5.4 and 5.5. One could observe that the highpass response is a mirrored version of the lowpass response, while the bandpass response is symmetric by itself. The slope rolloff speed is apparently -12dB/oct for the low- and highpass, and -6dB/oct for the bandpass.

The relative symmetry between the lowpass and the highpass amplitude responses has a clear algebraic explanation: applying the LP to HP substitution to a 2-pole lowpass produces a 2-pole highpass and vice versa. The symmetry of the bandpass amplitude response has the same explanation: applying the LP to HP substitution to the 2-pole bandpass converts it into itself.

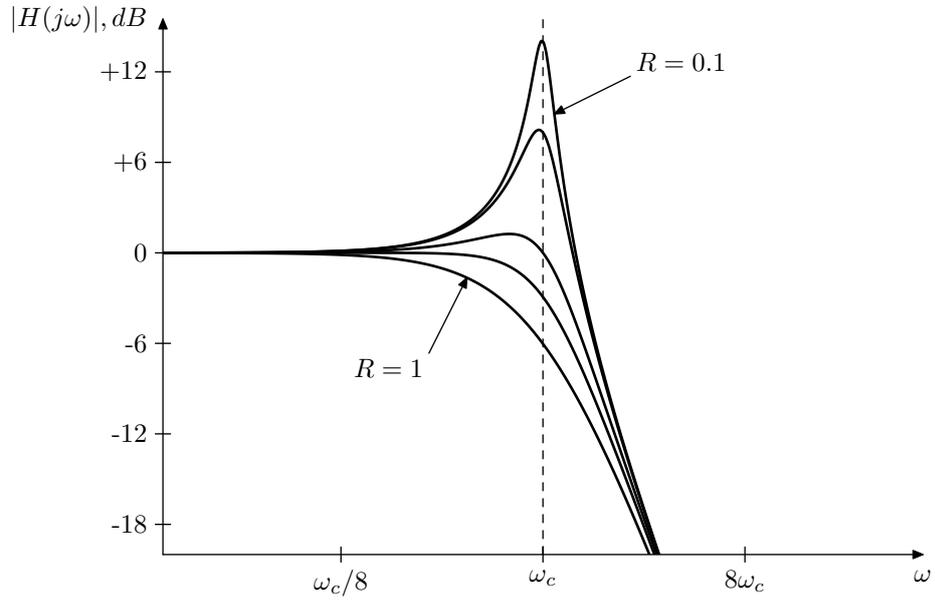


Figure 5.3: Amplitude response of a 2-pole lowpass filter.

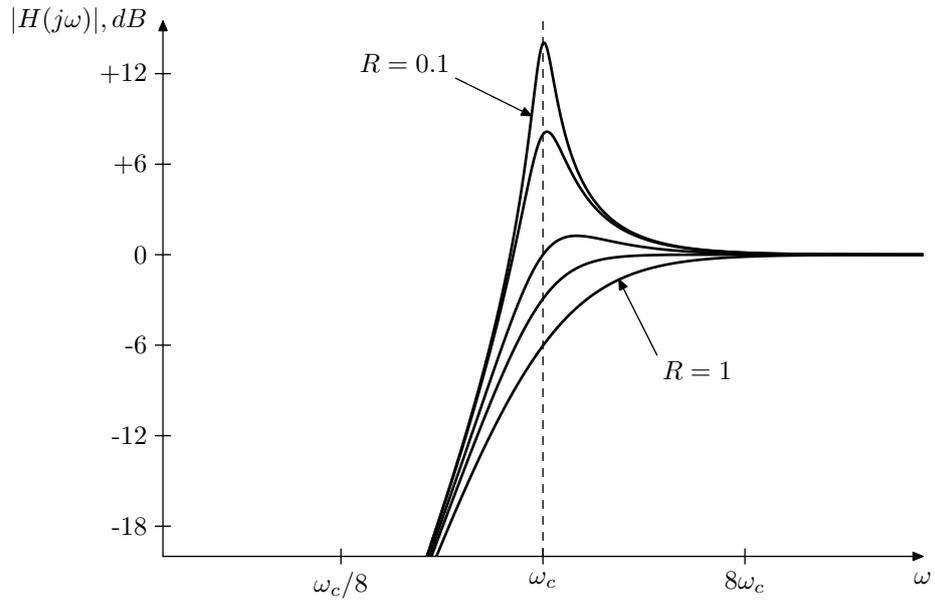


Figure 5.4: Amplitude response of a 2-pole highpass filter.

Notice that $y_{LP}(t) + 2Ry_{BP}(t) + y_{HP}(t) = x(t)$, that is, the input signal is split into lowpass, bandpass and highpass components. The same can be expressed in the transfer function form:

$$H_{LP}(s) + 2RH_{BP}(s) + H_{HP}(s) = 1$$

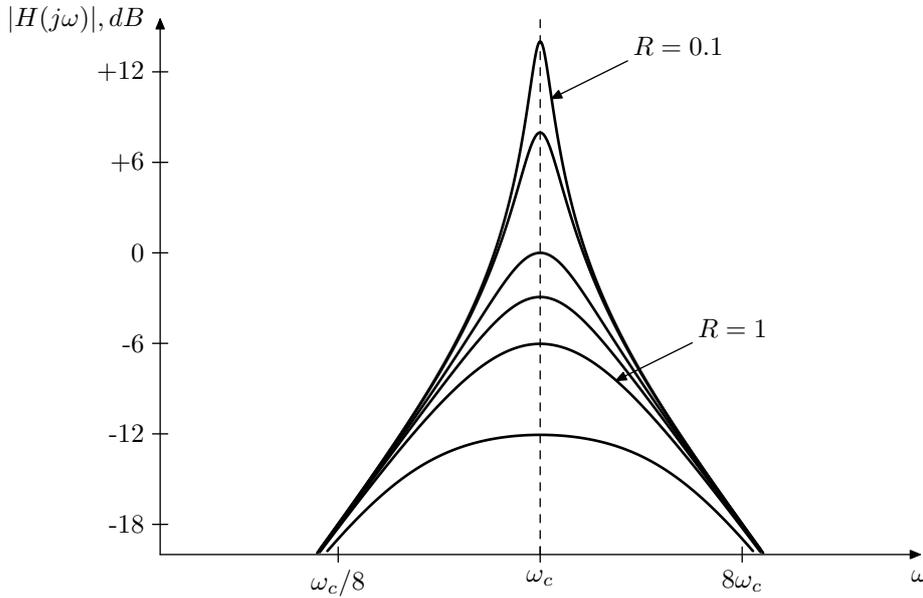


Figure 5.5: Amplitude response of a 2-pole bandpass filter.

The resonance of the filter is controlled by the R parameter. Contrarily to the ladder filter, where the resonance increases with the feedback amount, in the state-variable filter the bandpass signal feedback serves as a damping means for the resonance. In the absence of the bandpass signal feedback the filter will get unstable. The R parameter therefore may be referred to as the *damping* parameter.

Solving $s^2 + 2Rs + 1 = 0$ we obtain the poles of the filter at

$$s = -R \pm \sqrt{R^2 - 1} = \begin{cases} -R \pm \sqrt{R^2 - 1} & \text{if } R \geq 1 \\ -R \pm j\sqrt{1 - R^2} & \text{if } -1 \leq R \leq 1 \end{cases}$$

Without trying to give a precise definition of the resonance concept, we could say that at $R = 1$ there is no resonance (there are two real poles at $s = -1$). As R starts decreasing from 1 towards 0 there appear two mutually conjugate complex poles moving along the unit circle towards the imaginary axis, so the resonance slowly appears. At $R = 0$ the filter becomes unstable.²

The amplitude response at the cutoff ($\omega = 1$) is $1/2R$ for all three filter types. Except for the bandpass, the point $\omega = 1$ is not exactly the peak location but it's pretty close (the smaller the value of R , the closer is the true peak to $\omega = 1$). The phase response at the cutoff is -90° for lowpass, 0° for bandpass and $+90^\circ$ for highpass.

²The “resonance” control for the SVF filter can be introduced in a number of different ways. One common approach is to use the parameter $Q = 1/2R$, however this doesn't allow to go easily into the selfoscillation range in the nonlinear versions of this filter. Another option is using $r = 1 - R$, which differs from the resonance control parameter k of the TSK filters (discussed in section 5.7) just by a factor of 2, the selfoscillation occurring at $r = 1$. Other, more sophisticated mappings, can be used for a “more natural feel” of the resonance control.

At $|R| > 1$ the filter has two real poles and thus “falls apart” into a serial combination of two 1-pole filters:³

$$\begin{aligned} H_{\text{HP}}(s) &= \frac{s^2}{s^2 + 2Rs + 1} = \frac{s}{s - p_1} \cdot \frac{s}{s - p_2} \\ H_{\text{BP}}(s) &= \frac{s}{s^2 + 2Rs + 1} = \frac{s}{s - p_1} \cdot \frac{1}{s - p_2} \\ H_{\text{LP}}(s) &= \frac{1}{s^2 + 2Rs + 1} = \frac{1}{s - p_1} \cdot \frac{1}{s - p_2} \end{aligned}$$

where $p_1 p_2 = 1$. These 1-pole filters become visible in the amplitude responses at sufficiently large R as two different “cutoff points” (Fig. 5.6).

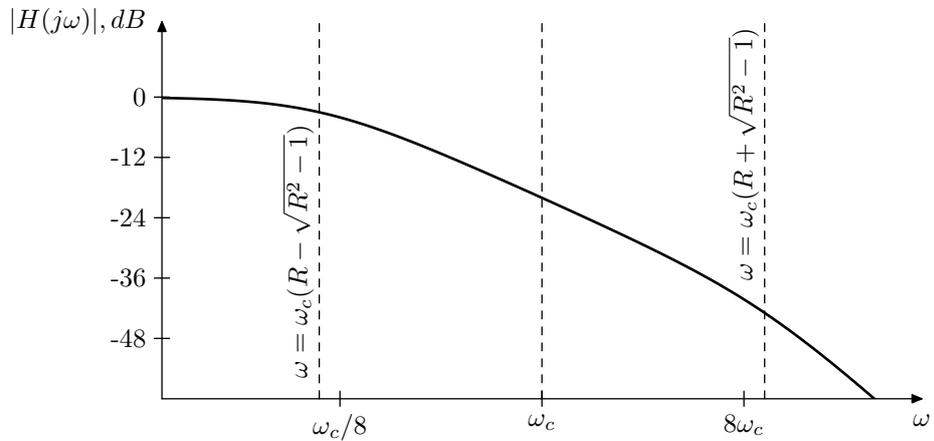


Figure 5.6: Amplitude response of a non-resonating 2-pole lowpass filter.

5.2 Linear digital model

Skipping the naive implementation, which the readers should be perfectly capable of creating and analyzing themselves by now, we proceed with the discussion of the TPT model.

Assuming $g\xi + s_n$ instantaneous responses for the two trapezoidal integrators one can redraw Fig. 5.1 to obtain the discrete-time model in Fig. 5.7.

Picking y_{HP} as the zero-delay feedback equation’s unknown⁴ we obtain from Fig. 5.7:

$$y_{\text{HP}} = x - 2R(gy_{\text{HP}} + s_1) - g(gy_{\text{HP}} + s_1) - s_2$$

from where

$$(1 + 2Rg + g^2) y_{\text{HP}} = x - 2Rs_1 - gs_1 - s_2$$

³Of course the same decomposition is formally possible for complex poles, but a 1-pole filter with a complex pole cannot be implemented as a real system.

⁴The state-variable filter has two feedback paths sharing a common path segment. In order to obtain a single feedback equation rather than an equation system we should pick a signal on this common path as the unknown variable.

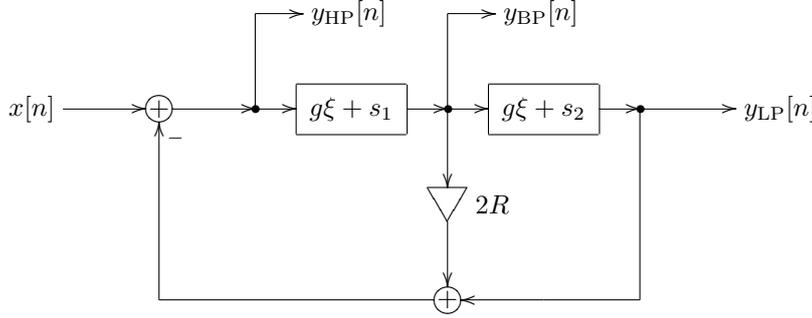


Figure 5.7: TPT 2-pole multimode state-variable filter in the instantaneous response form.

from where

$$y_{\text{HP}} = \frac{x - 2Rs_1 - gs_1 - s_2}{1 + 2Rg + g^2} \quad (5.1)$$

Using y_{HP} we can proceed defining the remaining signals in the structure.⁵

5.3 Further filter types

By mixing the lowpass, bandpass and highpass outputs one can obtain further filter types.

Unity gain (a.k.a. “normalized”) bandpass

$$H_{\text{BP1}}(s) = 2RH_{\text{BP}}(s) = \frac{2Rs}{s^2 + 2Rs + 1}$$

This version of the bandpass filter has a unity gain at the cutoff (Fig. 5.8). Notice that the unity gain bandpass signal can be directly picked up at the output of the $2R$ gain element in Fig. 5.1.

Band-shelving filter

By adding/subtracting the unity gain bandpass signal to/from the input signal one obtains the band-shelving filter (Fig. 5.9):

$$H_{\text{BS}}(s) = 1 + 2RK H_{\text{BP}}(s) = 1 + \frac{2RKs}{s^2 + 2Rs + 1}$$

Similarly to the other shelving filter types we can specify the mid-slope requirement $|H_{\text{BS}}(j\omega)| = \sqrt{1 + K}$ (for some ω), from where we obtain

$$R = \frac{|\omega - \omega^{-1}|}{2\sqrt{1 + K}} = \frac{|2^{\Delta/2} - 2^{-\Delta/2}|}{2\sqrt{1 + K}}$$

where Δ is the desired mid-slope bandwidth (in octaves) of the peak.

⁵Apparently, $1 + 2Rg + g^2 > 0 \forall g > 0$, provided $R > -1$. Thus, instantaneously unstable feedback may appear only if $R \leq -1$.

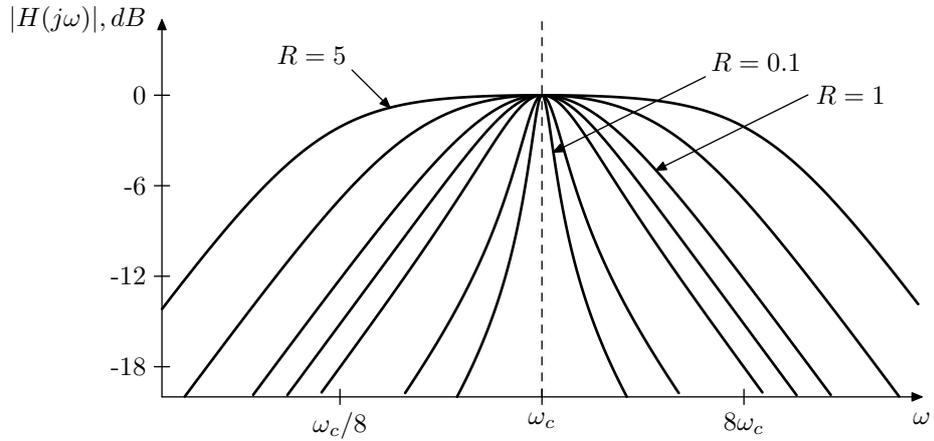


Figure 5.8: Amplitude response of a 2-pole unity gain bandpass filter.

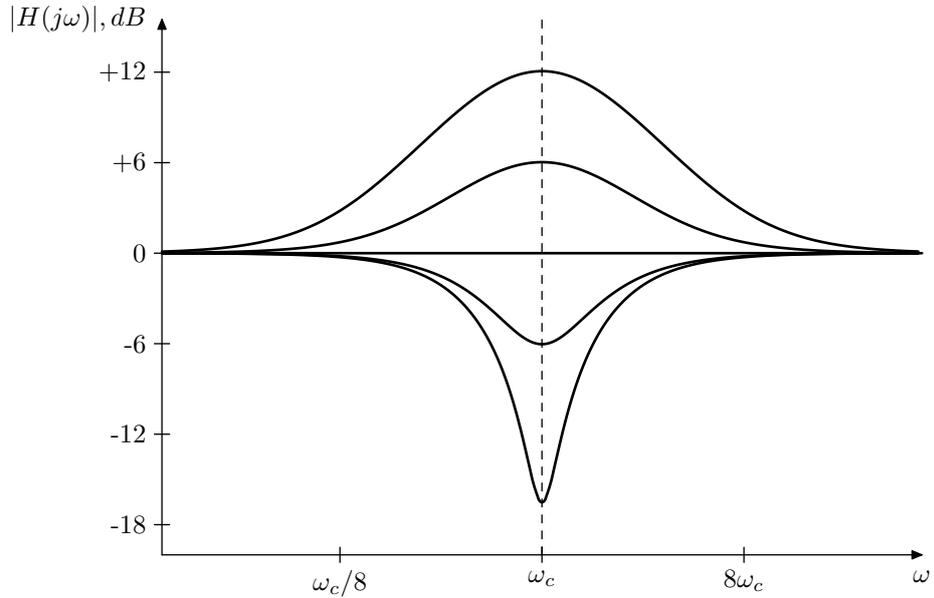


Figure 5.9: Amplitude response of a 2-pole band-shelving filter for $R = 1$ and varying K .

Low- and high-shelving filters

Attempting to obtain 2-pole low- and high-shelving filters in a straightforward fashion:

$$H_{LS}(s) = 1 + K \cdot H_{LP}(s) \quad H_{HS}(s) = 1 + K \cdot H_{HP}(s)$$

we notice that the amplitude responses of such filters have a strange dip (for $K > 0$) or peak (for $K < 0$) even at a non-resonating setting of $R = 1$ (Fig. 5.10). This peak/dip is due to a steeper phase response curve of the 2-pole lowpass

and highpass filters compared to 1-poles.

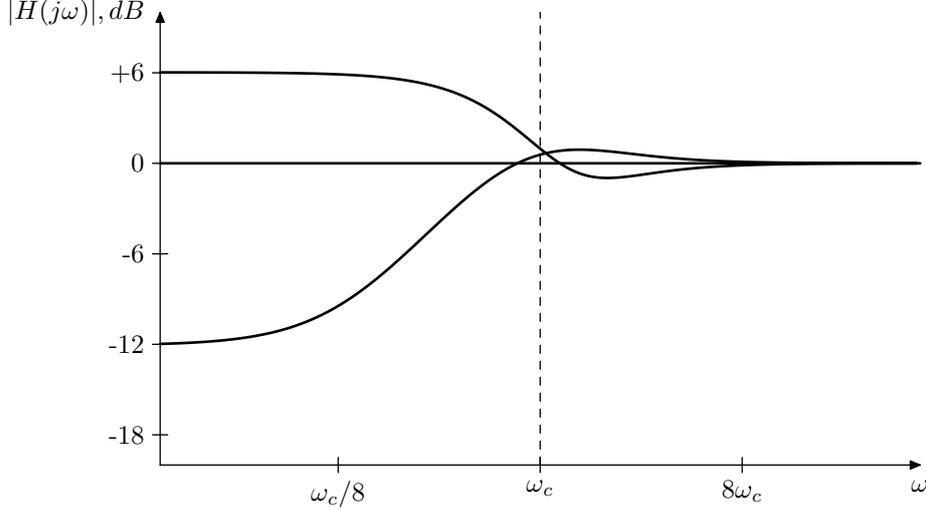


Figure 5.10: Amplitude response of a naive 2-pole low-shelving filter for $R = 1$ and varying K .

Instead, let's recall that 1-pole shelving filters have a symmetric amplitude response property. We could try to obtain 2-pole shelving filters from the same requirement of the amplitude response symmetry in respect to the "LP to HP" substitution.⁶ For a low-shelving filter we start off with:

$$\begin{aligned} G(s) &= \frac{s^2 + 2RM_s + M^2}{M^2s + 2RM_s + 1} = \frac{1}{M^2} \cdot \frac{M^2s^2 + 2RM^3s + M^4}{M^2s + 2RM_s + 1} = \\ &= \frac{1}{M^2} \cdot \frac{(s/\omega_c)^2 + 2RM^2(s/\omega_c) + M^4}{(s/\omega_c)^2 + 2R(s/\omega_c) + 1} \end{aligned}$$

where $\omega_c = 1/M$ (notice that at $R = 1$ we have a squared response of a 1-pole shelving filter). Considering the requirements $H_{\text{LS}}(0) = 1 + K$ and $H_{\text{LS}}(\infty) = 1$, we conclude that the low-shelving response is then

$$\begin{aligned} H_{\text{LS}}(s) &= M^2 \cdot G(s) = \frac{(s/\omega_c)^2 + 2RM^2(s/\omega_c) + M^4}{(s/\omega_c)^2 + 2R(s/\omega_c) + 1} = \\ &= M^4 H_{\text{LP}}(s) + M^2 \cdot 2RH_{\text{BP}}(s) + H_{\text{HP}}(s) = \\ &= 1 + KH_{\text{LP}}(s) + (\sqrt{1+K} - 1/2R) \cdot 2RH_{\text{BP}}(s) \end{aligned}$$

where $M = (1 + K)^{1/4}$ and $\omega_c = 1/M = (1 + K)^{-1/4}$. Respectively, for a non-unity midpoint frequency:

$$\omega_c = \omega_{\text{mid}}(1 + K)^{-1/4} \quad (\text{low-shelving})$$

For a high-shelving filter take a reciprocal $G(s)$:

$$G(s) = \frac{M^2s + 2RM_s + 1}{s^2 + 2RM_s + M^2} = \frac{1}{M^2} \cdot \frac{M^4(s/M)^2 + 2RM^2(s/M) + 1}{(s/M)^2 + 2RM^2(s/M) + 1} =$$

⁶The target low-shelving and high-shelving responses were taken from Robert Bristow-Johnson's *Audio EQ Cookbook*.

$$= \frac{1}{M^2} \cdot \frac{M^4(s/\omega_c)^2 + 2RM^2(s/\omega_c) + 1}{(s/\omega_c)^2 + 2R(s/\omega_c) + 1}$$

Noticing the requirements $H_{\text{HS}}(0) = 1$ and $H_{\text{HS}}(\infty) = 1 + K$:

$$\begin{aligned} H_{\text{HS}}(s) &= M^2 \cdot G(s) = \frac{M^4(s/\omega_c)^2 + 2RM^2(s/\omega_c) + 1}{(s/\omega_c)^2 + 2R(s/\omega_c) + 1} = \\ &= M^4 H_{\text{HP}}(s) + M^2 \cdot 2RH_{\text{BP}}(s) + H_{\text{LP}}(s) = \\ &= 1 + KH_{\text{HP}}(s) + (\sqrt{1+K} - 1/2R) \cdot 2RH_{\text{BP}}(s) \end{aligned}$$

where $M = (1 + K)^{1/4}$ and $\omega_c = (1 + K)^{1/4}$. Respectively, for a non-unity midpoint frequency:

$$\omega_c = \omega_{\text{mid}}(1 + K)^{1/4} \quad (\text{high-shelving})$$

The low-shelving 2-pole filter's amplitude responses are plotted in Fig. 5.11. Notice that at strong resonance there is a symmetric peak/dip pair, while at $R \gg 1$ the two different cutoff points of the two real 1-pole low-shelving filters become visible in the response.

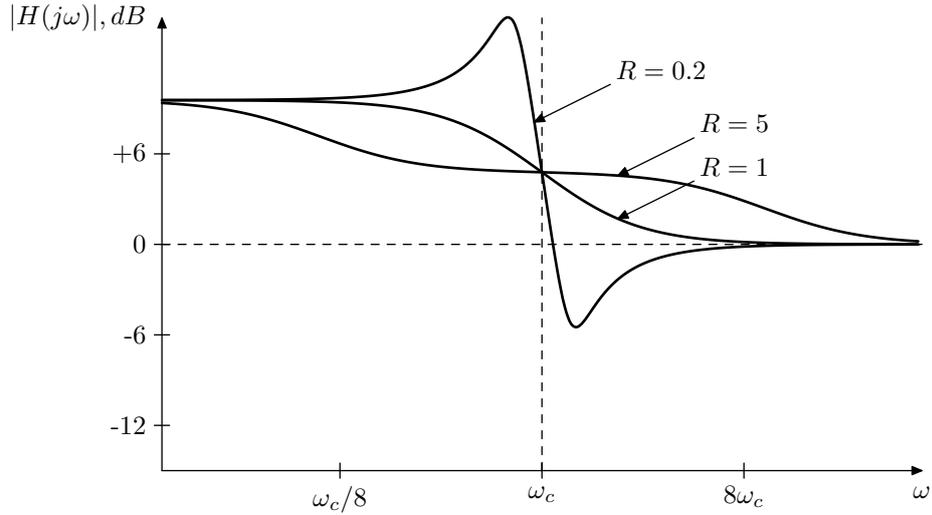


Figure 5.11: Amplitude response of a symmetric 2-pole low-shelving filter.

The high-shelving amplitude responses are similar.

Notch filter

At $K = -1$ the band-shelving filter turns into a notch (or bandstop) filter (Fig. 5.12):

$$H_{\text{N}}(s) = 1 - 2RH_{\text{BP}}(s) = \frac{s^2 + 1}{s^2 + 2Rs + 1}$$

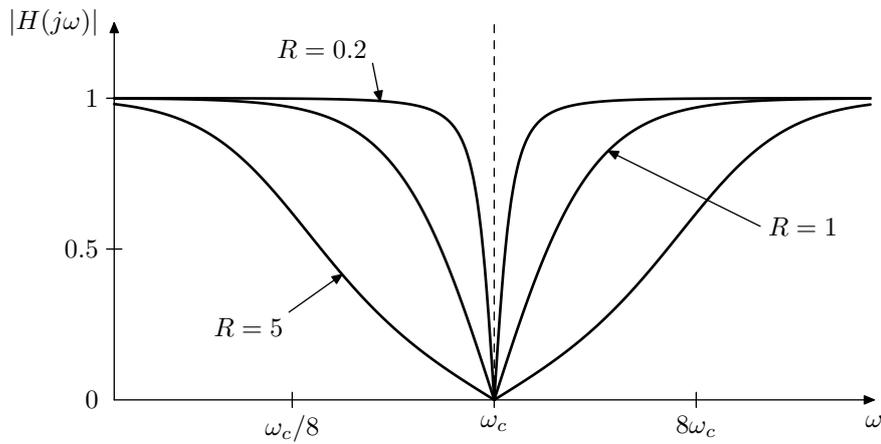


Figure 5.12: Amplitude response of a 2-pole notch filter. The amplitude scale is linear.

Allpass filter

At $K = -2$ the band-shelving filter turns into an allpass filter (Fig. 5.13):

$$H_{\text{AP}}(s) = 1 - 4RH_{\text{BP}}(s) = \frac{s^2 - 2Rs + 1}{s^2 + 2Rs + 1}$$

Notice how the damping parameter affects the phase response slope.

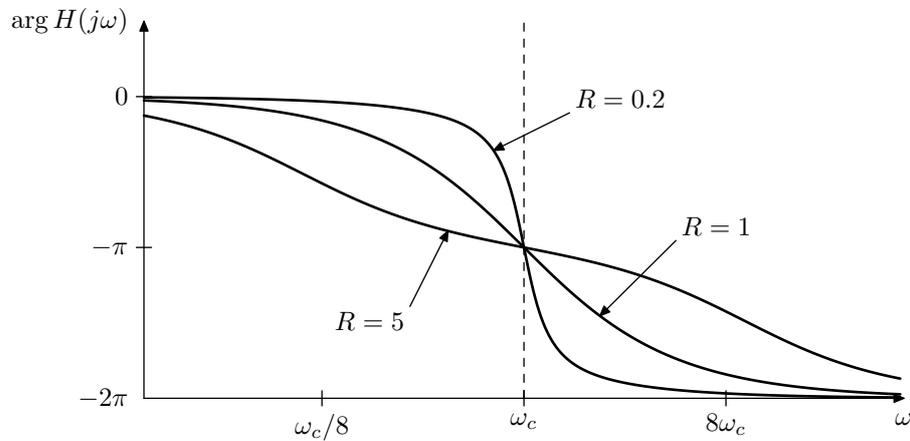


Figure 5.13: Phase response of a 2-pole allpass filter.

Peaking filter

By subtracting the highpass signal from the lowpass signal (or also vice versa) we obtain the peaking filter (Fig. 5.14):

$$H_{\text{PK}}(s) = H_{\text{LP}}(s) - H_{\text{HP}}(s) = \frac{1 - s^2}{s^2 + 2Rs + 1}$$

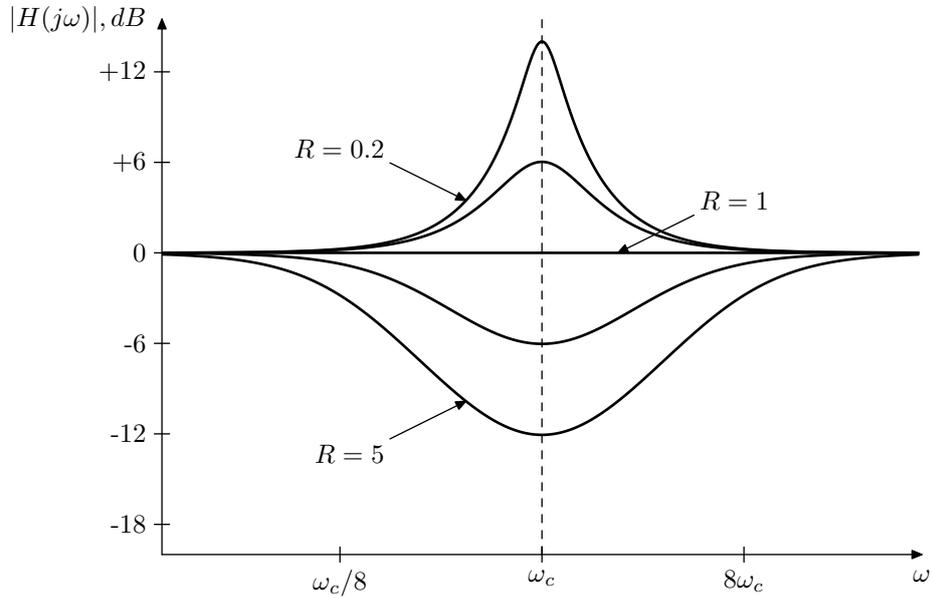


Figure 5.14: Amplitude response of a 2-pole peaking filter.

5.4 LP to BP/BS substitutions

The 2-pole unity gain bandpass response can be obtained from the lowpass response by a so-called *LP to BP* (lowpass to bandpass) *substitution*:

$$s \leftarrow \frac{1}{2R} \cdot \left(s + \frac{1}{s} \right) \quad (5.2)$$

Since s and $1/s$ are used symmetrically within the right-hand side of (5.2), it immediately follows that the result of the substitution is invariant relative to the LP to HP substitution $s \leftarrow 1/s$. Therefore the result of the LP to BP substitution has an amplitude response which is symmetric in the logarithmic frequency scale.

Using $s = j\omega$, we obtain

$$j\omega \leftarrow \frac{1}{2R} \cdot \left(j\omega + \frac{1}{j\omega} \right)$$

or

$$\omega \leftarrow \frac{1}{2R} \cdot \left(\omega - \frac{1}{\omega} \right)$$

Denoting the new ω as ω' we write

$$\omega = \frac{1}{2R} \cdot \left(\omega' - \frac{1}{\omega'} \right) \quad (5.3)$$

Instead of trying to understand the mapping of ω to ω' it is easier to understand the inverse mapping from ω' to ω , as explicitly specified by (5.3). Furthermore, it is more illustrative to express ω' in the logarithmic scale:

$$\omega = \frac{1}{2R} \cdot \left(e^{\ln \omega'} - e^{-\ln \omega'} \right) = \frac{1}{R} \sinh \ln \omega' \quad \text{if } \omega > 0$$

$$\omega = -\frac{1}{2R} \cdot \left(e^{\ln|\omega'|} - e^{-\ln|\omega'|} \right) = -\frac{1}{R} \sinh \ln|\omega'| \quad \text{if } \omega < 0$$

Thus

$$\omega = \frac{1}{R} \sinh(\operatorname{sgn} \omega' \cdot \ln|\omega'|) \quad (5.4)$$

Since $\ln|\omega'|$ takes up the entire real range of values in each of the cases $\omega > 0$ and $\omega < 0$ and respectively so does $\sinh(\operatorname{sgn} \omega' \cdot \ln|\omega'|)$,

$$\begin{aligned} \omega' \in (0, +\infty) &\iff \omega \in (-\infty, +\infty) \\ \omega' \in (-\infty, 0) &\iff \omega \in (-\infty, +\infty) \end{aligned}$$

This means that the entire range $\omega \in (-\infty, +\infty)$ is mapped once onto the positive frequencies ω' and once onto the negative frequencies ω' . Furthermore, the mapping and its inverse are strictly increasing on each of the two segments $\omega > 0$ and $\omega < 0$, since $d\omega/d\omega' > 0$. The unit frequencies $\omega' = \pm 1$ are mapped from $\omega = 0$.

Since we are often dealing with unity-cutoff transfer functions ($\omega_c = 1$), it's interesting to see to which frequencies ω'_c the unity cutoff is mapped. Recalling that the entire bipolar range of ω is mapped to the positive range of ω' , we need to include the negative cutoff point ($\omega_c = -1$) into our transformation. On the other hand, we are interested only in positive ω'_c , since the negative-frequency range of the amplitude response is symmetric to the positive-frequency range anyway. Under these reservations, from (5.4) we have:

$$\frac{1}{R} \sinh \ln \omega'_c = \pm 1$$

from where $\ln \omega'_c = \sinh^{-1} R$, or, changing the logarithm base:

$$\log_2 \omega'_c = \pm \frac{\sinh^{-1} R}{\ln 2}$$

The distance in octaves between the two ω'_c points can be defined as the bandwidth of the transformation:

$$\Delta = \frac{2}{\ln 2} \sinh^{-1} R$$

Respectively, given the bandwidth Δ , the damping is

$$R = \sinh \frac{\Delta \cdot \ln 2}{2} = \frac{2^{\Delta/2} - 2^{-\Delta/2}}{2}$$

The transformation of the poles and zeros by the LP to BP transformation can be obtained from

$$s = \frac{1}{2R} \cdot \left(s' + \frac{1}{s'} \right) \quad (5.5)$$

resulting in

$$s' = Rs \pm \sqrt{R^2 s^2 - 1}$$

Regarding the stability preservation consider that the sum $(s' + 1/s')$ in (5.5) is located in the same complex semiplane (left or right) as s' . Therefore, as long as $R > 0$, the original value s is located in the same semiplane as its images

s' , which implies that the stability is preserved. On the other hand, negative values of R “flip” the stability.

As for performing the LP to BP substitution in a block diagram, differently from the LP to HP substitution, here we don't need differentiators. The substitution can be performed by replacing all (unity-cutoff) integrators in the system with the structure in Fig. 5.15, thereby substituting $2Rs/(s^2 + 1)$ for $1/s$, which is algebraically equivalent to (5.2).⁷

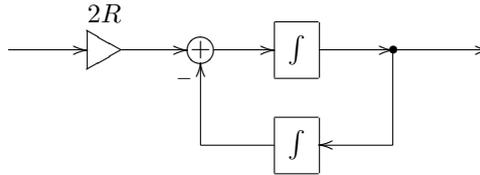


Figure 5.15: “LP to BP” integrator.

The *LP to BS* (lowpass to bandstop) *substitution*⁸ is obtained as a series of LP to HP substitution followed by an LP to BP substitution. Indeed, applying the LP to BP substitution to a 1-pole highpass, we obtain the 2-pole notch (“bandstop”) filter. Therefore, applying a series of LP to HP and LP to BP substitutions to a 1-pole lowpass we also obtain the 2-pole notch filter.

Combining the LP to HP and LP to BP substitutions expressions in the mentioned order gives an algebraic expression for the LP to BS substitution:

$$\frac{1}{s} \leftarrow \frac{1}{2R} \cdot \left(s + \frac{1}{s} \right) \quad (5.6)$$

The bandwidth considerations of the LP to BS substitution are pretty much equivalent to those of LP to BP substitution and can be obtained by considering the LP to BS substitution as an LP to BP substitution applied to a result of the LP to HP substitution.

The block-diagram form of the LP to BS substitution can be obtained by directly implementing the right-hand expression in (5.6) as a replacement for the integrators. This however requires a differentiator for the implementation of the s term of the sum.

5.5 Nonlinear model

In the ladder filter the resonance was created as the result of the feedback. Therefore by limiting the feedback level (by a saturator) we could control the resonance amount and respectively prevent the filter from becoming unstable.

The feedback in the SVF has a more complicated structure. Particularly, the bandpass path is responsible for damping the resonance. We could therefore

⁷For a differentiator, a similar substitution structure (containing an integrator and a differentiator) is trivially obtained from the right-hand side of (5.2).

⁸Notice that BS here stands for “bandstop” and not for “band-shelving”. The alternative name for the substitution could have been “LP to Notch”, but “LP to bandstop” seems to be commonly used, so we'll stick to that one.

use the direct form I-style integrator (Fig. 3.8) resulting in the integrator in Fig. 5.17. Or one could equivalently use the transposed direct form II-style integrator (Fig. 3.10) resulting in the integrator in Fig. 5.18.

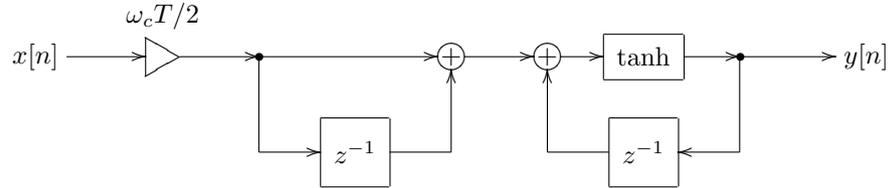


Figure 5.17: Saturating direct form I trapezoidal integrator.

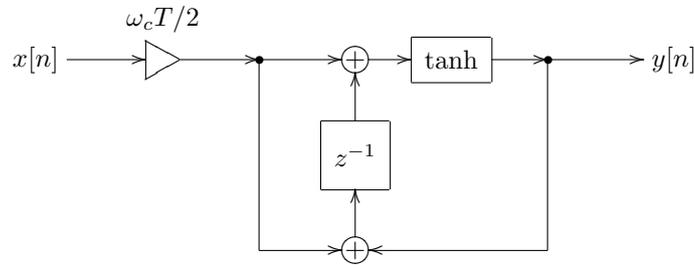


Figure 5.18: Saturating transposed direct form II trapezoidal integrator.

5.6 Serial decomposition

Recall that a 1-pole multimode filter can be used to implement any 1st-order rational transfer function. Similarly, a multimode SVF can be used to implement practically any 2nd-order rational transfer function. Indeed, consider

$$H(s) = \frac{b_2 s^2 + b_1 s + b_0}{s^2 + a_1 s + a_0}$$

where we assume $a_0 > 0$.¹¹ Then

$$\begin{aligned} H(s) &= \frac{b_2 s^2 + b_1 s + b_0}{s^2 + 2\frac{a_1}{2\sqrt{a_0}}\sqrt{a_0}s + \sqrt{a_0}^2} = \frac{b_2 s^2 + b_1 s + b_0}{s^2 + 2R\omega_c s + \omega_c^2} = \\ &= b_2 \frac{s^2}{s^2 + 2R\omega_c s + \omega_c^2} + \frac{b_1}{\omega_c} \cdot \frac{\omega_c s}{s^2 + 2R\omega_c s + \omega_c^2} + \frac{b_0}{\omega_c^2} \cdot \frac{\omega_c^2}{s^2 + 2R\omega_c s + \omega_c^2} = \end{aligned}$$

¹¹If $a_0 < 0$, this means that $H(s)$ has two real poles of different signs. If $a_0 = 0$ then at least one of the poles is at $s = 0$. In either case, this filter is already unstable, which means, if we are practically interested in its implementation, most likely there is a nonlinear analog prototype, and we simply can apply TPT to this prototype to obtain a digital structure. If we insist on using the SVF structure (why would we?), we can also extend it to the canonical form by introducing a gain element into the lowpass feedback path.

$$= b_2 H_{\text{HP}}(s) + \frac{b_1}{\omega_c} H_{\text{BP}}(s) + \frac{b_0}{\omega_c^2} H_{\text{LP}}(s)$$

This further allows to implement practically any given stable transfer function by a serial connection of a number of 2-pole (and possibly 1-pole) filters. Indeed, simply factor the numerator and the denominator into 2nd- and possibly 1st-order factors (where the 2nd-order real factors will necessarily appear for complex conjugate pairs of roots and optionally for pairs of real roots). Any pair of 2nd-order factors (one in the numerator, one in the denominator) can be implemented by a 2-pole multimode SVF. Any pair of 1st-order factors can be implemented by a 1-pole multimode. If there are not enough 2nd-order factors in the numerator or denominator, a pair of 1st order factors in the numerator or denominator can be combined into a 2nd-order factor.

The serial decomposition is not the only way to decompose a transfer function into transfer functions of lower orders. E.g. one could use partial fraction expansion to represent a transfer function as a sum of lower-order transfer functions. However partial fraction expansion becomes ill-conditioned if the poles of the transfer function are getting close together and is therefore generally less useful than serial decomposition.

Serial decomposition of lowpass ladder filter

As an example demonstrating the serial decomposition technique we will obtain a serial decomposition of a lowpass ladder filter. A lowpass ladder filter has no zeros and two pairs of conjugate pole pairs for $k > 0$. By considering two coinciding poles on a real axis also as mutually conjugate, we can assume $k \geq 0$.

Since there are no zeros, we simply need a 2-pole lowpass SVF for each conjugate pair of poles. Since

$$\frac{1}{(s-p)(s-p^*)} = \frac{1}{s^2 - 2\operatorname{Re} p \cdot s + |p|^2}$$

a 2-pole lowpass SVF expressed in terms of the complex conjugate poles is

$$\begin{aligned} H_{\text{LP2}}(s) &= \frac{|p|^2}{(s-p)(s-p^*)} = \frac{|p|^2}{s^2 + 2\frac{-\operatorname{Re} p}{|p|}|p| \cdot s + |p|^2} = \\ &= \frac{1}{\left(\frac{s}{|p|}\right)^2 + 2\frac{-\operatorname{Re} p}{|p|} \cdot \frac{s}{|p|} + 1} = \frac{1}{\left(\frac{s}{\omega_c}\right)^2 + 2R \cdot \frac{s}{\omega_c} + 1} \end{aligned}$$

Thus

$$\omega_c = |p| \quad R = -\frac{\operatorname{Re} p}{\omega_c}$$

Let p_1, p_1^*, p_2, p_2^* be the poles of the ladder filter. According to (4.2)

$$p_{1,2} = -1 + \frac{\pm 1 + j}{\sqrt{2}} k^{1/4} \quad (5.7)$$

Therefore the cutoffs of the 2-pole lowpasses are defined by

$$(\omega_{1,2})^2 = |p_{1,2}|^2 = \left(-1 \pm \frac{k^{1/4}}{\sqrt{2}}\right)^2 + \left(\frac{k^{1/4}}{\sqrt{2}}\right)^2 \quad (5.8)$$

Respectively the transfer function of the ladder filter can be represented as

$$\begin{aligned}
 H(s) &= g \frac{1}{\left(\frac{s}{|p_1|}\right)^2 + 2\frac{-\operatorname{Re} p_1}{|p_1|} \cdot \frac{s}{|p_1|} + 1} \cdot \frac{1}{\left(\frac{s}{|p_2|}\right)^2 + 2\frac{-\operatorname{Re} p_2}{|p_2|} \cdot \frac{s}{|p_2|} + 1} = \\
 &= g \frac{1}{\left(\frac{s}{\omega_1}\right)^2 + 2R_1\frac{s}{\omega_1} + 1} \cdot \frac{1}{\left(\frac{s}{\omega_2}\right)^2 + 2R_2\frac{s}{\omega_2} + 1} \quad (5.9)
 \end{aligned}$$

The unknown gain coefficient g can be found by evaluating (4.1) at $s = 0$, obtaining the condition $H(0) = 1/(1+k)$. Evaluating (5.9) at $s = 0$ yields $H(0) = g$. Therefore

$$g = \frac{1}{1+k}$$

Since the pole expressions (4.2) and respectively (5.7) apply only to the unit-cutoff ladder filter, in order to construct a ladder filter of an arbitrary cutoff ω_c we need to respectively change the cutoffs ω_1 and ω_2 of the 2-poles in the decomposition (5.9) according to

$$\begin{aligned}
 \omega'_1 &= \omega_1 \omega_c \\
 \omega'_2 &= \omega_2 \omega_c
 \end{aligned} \quad (5.10)$$

Notice that the ratio of the cutoffs is invariant relatively to the cutoff changes:

$$\left(\frac{\omega'_1}{\omega'_2}\right)^2 = \left(\frac{\omega_1}{\omega_2}\right)^2 = \frac{\left(-1 + \frac{k^{1/4}}{\sqrt{2}}\right)^2 + \left(\frac{k^{1/4}}{\sqrt{2}}\right)^2}{\left(-1 - \frac{k^{1/4}}{\sqrt{2}}\right)^2 + \left(\frac{k^{1/4}}{\sqrt{2}}\right)^2} \quad (5.11)$$

Prewarping of decomposed filters

Expressing the decomposition (5.9) in the form of a block diagram we obtain the simple structure in Fig. 5.19 where H_1 and H_2 are defined by the respective 2-pole factors of (5.9).

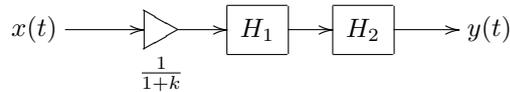


Figure 5.19: Serial decomposition of a ladder filter.

That seems to be it, but the simplicity of Fig. 5.19 can be somewhat misleading, if a discrete-time implementation of the structure is implied. What might be not immediately obvious is that the cutoffs of H_1 and H_2 defined by (5.10) are the *analog* rather than digital cutoffs. Normally one doesn't need to remember about the distinction between analog and digital cutoffs, because prewarping takes care of mapping the latter to the former (or back). However, if the analog cutoffs need to be in a special relationship (like the one in (5.11)), this relationship may be destroyed by the prewarping.

Suppose we simply take two digital lowpass filters H_1 and H_2 , stick them together and set their digital cutoffs ω_{1d} and ω_{2d} to the values defined by (5.10):

$$\begin{aligned}\omega_{1d} &= \omega_1\omega_c \\ \omega_{2d} &= \omega_2\omega_c\end{aligned}$$

where ω_1 and ω_2 are defined by (5.8). Each filter will prewarp its cutoff separately by applying (3.7), therefore the respective analog cutoffs will be

$$\begin{aligned}\omega_{1a} &= \frac{2}{T} \tan \frac{\omega_1\omega_c T}{2} \\ \omega_{2a} &= \frac{2}{T} \tan \frac{\omega_2\omega_c T}{2}\end{aligned}$$

Clearly

$$\frac{\omega_{1a}}{\omega_{2a}} \neq \frac{\omega_{1d}}{\omega_{2d}} = \frac{\omega_1}{\omega_2}$$

This means that the cutoff property (5.11) is not preserved and respectively the analog transfer function is not the one that we expect.

How critical is this deviation from the expected transfer function depends on the specifics of the filter usage case and is also, to a certain degree, subjective. However, it is important to realize that the usual thinking “bilinear transform simply warps the frequency axis according to (3.7) and doesn’t change the frequency response of the filter in any other way” doesn’t apply anymore in the just discussed situation, because each of the 2-pole filters was prewarped independently.

The solution is to use a common prewarping for both H_1 and H_2 . The natural choice for the prewarping point is ω_c . Since ω_c is naturally chosen as the prewarping point for the ladder filter, by using the same prewarping point we can make the digital structure in Fig. 5.19 to have exactly the same transfer function as the corresponding digital implementation of the ladder filter. Thus we obtain the analog prewarped cutoff ω_{ca} :

$$\omega_{ca} = \frac{2}{T} \tan \frac{\omega_c T}{2}$$

and respectively the analog prewarped cutoffs of the 2-poles:

$$\begin{aligned}\omega_{1a} &= \omega_1\omega_{ca} \\ \omega_{2a} &= \omega_2\omega_{ca}\end{aligned}$$

5.7 Transposed Sallen–Key filters

Attempting to build a 2-pole lowpass ladder filter (Fig. 5.20) we don’t end up with a useful filter.

Indeed, the transfer function of this filter is

$$H(s) = \frac{1}{k + (1 + s)^2}$$

and the poles are respectively

$$s = -1 \pm \sqrt{-k} = -1 \pm j\sqrt{k} \quad (k \geq 0)$$

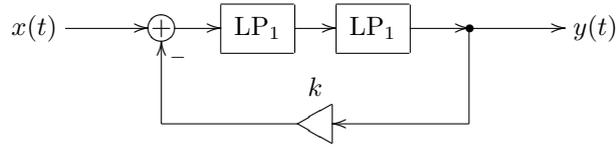


Figure 5.20: 2-pole ladder filter (not very useful).

Comparing to the transfer function and the placement of the poles of the SVF, we notice that the corresponding SVF cutoff and damping settings are

$$\begin{cases} \omega_c = |-1 \pm j\sqrt{k}| = \sqrt{1+k} \\ R = \frac{-\operatorname{Re}(-1 \pm j\sqrt{k})}{|-1 \pm j\sqrt{k}|} = \frac{1}{\sqrt{1+k}} \end{cases}$$

Thus, firstly, there is coupling between the feedback amount and the effective cutoff of the filter. Secondly, as k grows, R stays strictly positive, thus the filter never goes into selfoscillation (and, as with the 4-pole ladder filter the selfoscillation would be quite desired once we introduce a saturator into the filter). So, all in all, not a very useful structure.

Rather than giving up, let's try to introduce a second feedback path into the overall structure (Fig. 5.21) and let's try to figure some useful settings for the feedback amounts k_1 and k_2 . We also introduce the multimode pickups at the same time, to see if we can make any use of them.

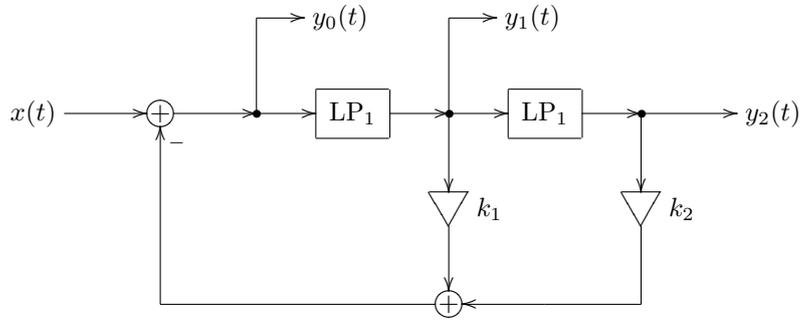


Figure 5.21: 2-pole ladder filter with two feedback paths.

Computing the transfer function we have

$$y_0 = x - k_1 \frac{y_0}{s+1} - k_2 \frac{y_0}{(s+1)^2}$$

from where

$$(s+1)^2 y_0 = (s+1)^2 x - k_1 (s+1) y_0 - k_2 y_0$$

from where

$$((s+1)^2 + k_1(s+1) + k_2) y_0 = (s+1)^2 x$$

and

$$y_0 = \frac{(s+1)^2}{(s+1)^2 + k_1(s+1) + k_2} x$$

Respectively

$$y_2 = \frac{1}{(s+1)^2} y_0 = \frac{1}{(s+1)^2 + k_1(s+1) + k_2} x$$

Rewriting the denominator we get

$$(s+1)^2 + k_1(s+1) + k_2 = s^2 + (2+k_1)s + (k_2 + k_1 + 1)$$

We wish our denominator to be of the form $s^2 + 2Rs + 1$, which is ensured if $k_2 + k_1 = 0$. Letting $k_2 = k$ and $k_1 = -k$ we have

$$(s+1)^2 + k_1(s+1) + k_2 = s^2 + 2\left(1 - \frac{k}{2}\right)s + 1$$

(so the selfoscillation occurs at $k = 2$). The corresponding structure is shown in Fig. 5.22. This structure happens to be a transpose of the Sallen–Key filter, therefore we will refer to it as the *transposed Sallen–Key* (TSK) filter.¹²

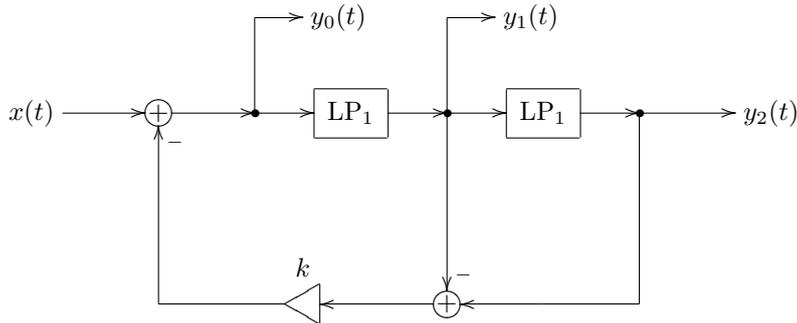


Figure 5.22: Transposed Sallen–Key filter (lowpass).

The transfer functions corresponding to y_0 , y_1 and y_2 are respectively

$$H_0(s) = \frac{(s+1)^2}{s^2 + 2(1 - k/2)s + 1}$$

$$H_1(s) = \frac{s+1}{s^2 + 2(1 - k/2)s + 1}$$

$$H_2(s) = \frac{1}{s^2 + 2(1 - k/2)s + 1}$$

So y_2 is having the familiar 2-pole lowpass filter transfer function of the SVF, where $k = 1$ corresponds to the self-oscillation point.

¹²The author has used the works of Tim Stinchcombe as the information source on the Sallen–Key filter. The idea to introduce TSK filters as a systematic concept arose from the discussions with Julian Parker.

As for y_1 , notice that

$$H_1(s) = \frac{s+1}{s^2+2s+1} = \frac{1}{s+1} \quad \text{for } k=0$$

That is, at the zero resonance setting y_1 is a 1-pole lowpass. So, $y_{\text{LP1}} = y_1$ can be considered as a kind of 1-pole lowpass “with resonance”:

$$H_{\text{LP1}}(s) = H_1(s)$$

In order to obtain further useful modes from the filter we need to introduce two additional pickups, connected to the highpass multimode components of the underlying 1-pole filters (labelled “LP₁” in the diagrams). So, let \bar{y}_1 be the highpass output of the first “LP₁” block and let \bar{y}_2 be the highpass output of the second “LP₁” block. Then

$$\begin{aligned} \bar{H}_1(s) &= sH_1(s) = \frac{(s+1)s}{s^2+2(1-k/2)s+1} \\ \bar{H}_2(s) &= sH_2(s) = \frac{s}{s^2+2(1-k/2)s+1} \end{aligned}$$

Thus $y_{\text{BP}} = \bar{y}_2$ is a 2-pole bandpass:

$$H_{\text{BP}}(s) = \bar{H}_1(s)$$

Considering that

$$\bar{H}_1(s) - \bar{H}_2(s) = \frac{s^2}{s^2+2(1-k/2)s+1}$$

we have $y_{\text{HP2}} = \bar{y}_1 - \bar{y}_2$:

$$H_{\text{HP2}}(s) = \bar{H}_1(s) - \bar{H}_2(s)$$

Noticing that

$$\bar{H}_1(s) = \frac{(s+1)s}{s^2+2s+1} = \frac{s}{s+1} \quad \text{for } k=0$$

we can define $y_{\text{HP1}} = \bar{y}_1$ as a kind of 1-pole highpass “with resonance”:

$$H_{\text{HP1}}(s) = \bar{H}_1(s)$$

Nonlinear version

The structure of the feedback in the TSK filter is very much like the one of the feedback in the ladder filter. The feedback and resonance amount grows with k . Therefore this filter can successfully accommodate a saturator immediately before or after the feedback point, e.g. like in Fig. 5.23.

As the transfer functions of both the lowpass TSK filter and the lowpass SVF are completely identical, the better (and different) accommodation of nonlinearities is probably the main reason to use a TSK filter at all, at least in the digital domain, where a linear TPT TSK filter is somewhat more computationally expensive than a TPT SVF.

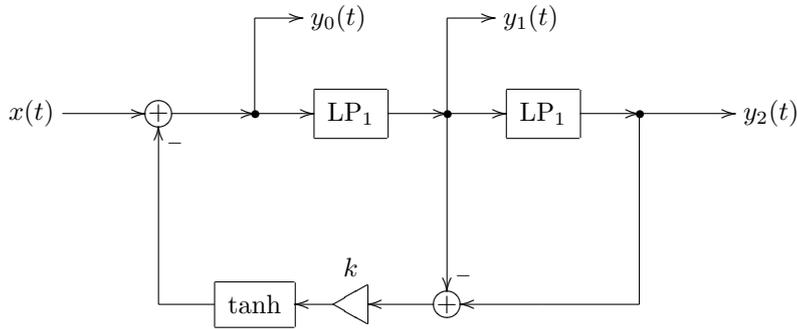


Figure 5.23: Lowpass TSK filter with saturator.

Highpass TSK filter

Instead of using 1-pole lowpass filters as the basis for the TSK filter, one could use 1-pole highpass filters. Effectively this performs an *LP to HP* substitution ($s \leftarrow 1/s$), thereby turning y_2 into a highpass output. The modal outputs get transformed respectively.¹³

Bandpass TSK filter

By using one 1-pole lowpass and one 1-pole highpass (in either order) one can produce a 2-pole bandpass signal at the y_2 output. This however requires letting $k_1 = 0$ and $k_2 = -k$ in Fig. 5.21 (this can be verified by a direct computation or obtained in the same way as we have obtained the coefficients for the lowpass TSK), resulting in the structure in Fig. 5.24.¹⁴

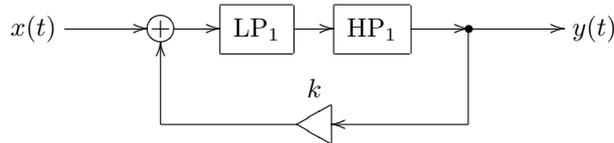


Figure 5.24: Bandpass TSK filter.

The transfer function is

$$H(s) = \frac{s}{s^2 + 2(1 - k/2) + 1}$$

Allpass TSK filter

By using two allpass filters one can produce an allpass TSK filter. This requires different values of k_1 and k_2 and the usage of the modal mixture. Considering

¹³A saturator could be less appropriate in the highpass TSK filter (especially in a digital model), since the overtones created by the nonlinearity will not be dampened back by the lowpass filtering.

¹⁴Of course, one could argue that calling the filter in Fig. 5.24 a TSK filter is a little bit far-fetched. However, considering the procedure of obtaining this filter, the name seems reasonable enough.

Fig. 5.21 with two 1-pole allpass filters with $(1-s)/(1+s)$ transfer functions, we obtain

$$y_0 = x - k_1 \frac{1-s}{1+s} y_0 - k_2 \frac{(1-s)^2}{(1+s)^2} y_0$$

from where

$$y_0 \left((1+s)^2 + k_1(1-s)(1+s) + k_2(1-s)^2 \right) = (1+s)^2 x$$

and

$$y_0 = \frac{(1+s)^2}{(1+s)^2 + k_1(1-s)(1+s) + k_2(1-s)^2} x$$

Considering the denominator alone:

$$(1+s)^2 + k_1(1-s)(1+s) + k_2(1-s)^2 = (1-k_1+k_2)s^2 + 2(1-k_2)s + (1+k_1+k_2)$$

Apparently the denominator can be turned into the form $s^2 + 2Rs + 1$ only by letting $k_1 = k_2 = 0$, but this immediately restricts it to $R = 1$. However, instead we can require that the coefficient at s^2 and the free term's coefficient are equal. This requires $k_1 = 0$. Letting $k_2 = k$ we obtain

$$y_0 = \frac{(1+s)^2}{(1+k)s^2 + 2(1-k)s + (1+k)} x$$

and

$$y_2 = \frac{(1-s)^2}{(1+k)s^2 + 2(1-k)s + (1+k)} x$$

Let $y = b_0 y_0 + b_2 y_2$ be a mixture of y_0 and y_2 with unknown coefficients b_0 and b_2 . We wish

$$y = b_0 y_0 + b_2 y_2 = \frac{(1+k)s^2 - 2(1-k)s + (1+k)}{(1+k)s^2 + 2(1-k)s + (1+k)} x = \frac{s^2 - 2\frac{1-k}{1+k}s + 1}{s^2 + 2\frac{1-k}{1+k}s + 1} x$$

from where

$$b_0(1+s)^2 + b_2(1-s)^2 = (1+k)s^2 - 2(1-k)s + (1+k)$$

or

$$(b_0 + b_2)s^2 + 2(b_0 - b_2)s + (b_0 + b_2) = (1+k)s^2 - 2(1-k)s + (1+k)$$

from where $b_0 = 1$ and $b_2 = k$, which corresponds to the structure in Fig. 5.25. It is easy to notice that this structure is very similar to the one of a phaser with some specific dry/wet mixing ratio.¹⁵

The damping parameter R is equal to $(1-k)/(1+k)$ so that for $k = -1 \dots +\infty$ the damping varies from $+\infty$ to -1 . The stable range $R = +\infty \dots 0$ corresponds to $k = -1 \dots 1$.

¹⁵The same structure can be obtained from a direct form II 1-pole allpass filter by the allpass substitution $z^{-1} \leftarrow (1-s)^2/(1+s)^2$. It is also interesting to notice that, applying the allpass substitution principle to the structure in Fig. 5.25, we can replace the series of the two 1-pole allpass filters in Fig. 5.25 by any other allpass filter, and the modified structure will still be an allpass filter.

The transposed Sallen–Key (TSK) filter is a kind of “2-pole ladder filter with two feedback signals”. This structure can nicely accommodate saturating nonlinearities and can output multiple filter modes. Similar structures can be built based on highpass filters, a mixture of a highpass and a lowpass and based on allpass filters.

Chapter 6

Allpass-based effects

Phasers are essentially LFO-modulated ladder filters built around allpass filters instead of lowpass filters. Flangers can be obtained from phasers by an allpass substitution. For these reasons both types belong to the VA filter discussion.

6.1 Phasers

The simplest *phaser* is built by mixing the unmodified (dry) input signal with an allpass-filtered (wet) signal as in Fig. 6.1, where the allpass filter's cutoff is typically modulated by an LFO.¹ The allpass filter can be rather arbitrary, except than it has to be a differential filter.²

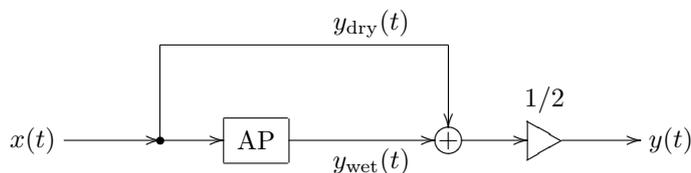


Figure 6.1: The simplest phaser.

At the points where the allpass filter's phase response is 180° , the wet and the dry signals will cancel each other, producing a notch. At the points where the allpass filter's phase response is 0° the wet and the dry signals will boost each other, producing a peak (Fig. 6.2).

The phaser structure in Fig. 6.1 contains no feedback, therefore there is no difference between naive and TPT digital implementations (except that the underlying allpass filters should be better constructed in a TPT way).

¹In the absence of LFO modulation the structure should be rather referred to as a (multi-) notch filter.

²Phasers typically use differential allpass filters or their digital counterparts. If e.g. a delay (which is not a differential filter, but *is* an allpass) is used as the allpass, the structure should be rather referred to as a *flanger*.

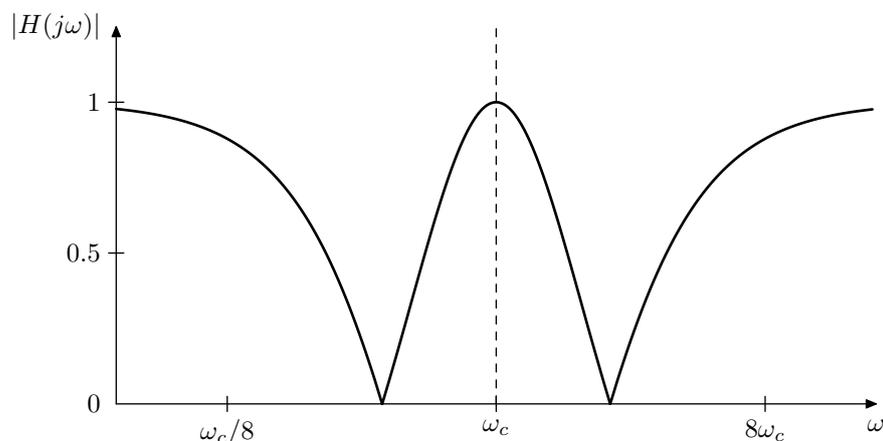


Figure 6.2: Amplitude response of the simplest phaser from Fig. 6.1 (using four identical 1-pole allpass filters with the same cutoff as the allpass core of the phaser).

Mixing at arbitrary ratios

Instead of mixing at the 50/50 ratio we can mix at any other ratio, where the sum of the dry and wet mixing gains should amount to unity. This will affect the depth of the notches and the height of the peaks. For the phaser in Fig. 6.1 the mixing ratio higher than 50/50 (where the wet signal amount is more than 50%) hardly makes sense.

Instead of mixing y_{dry} and y_{wet} at different ratios we could simply crossfade the output signal between $x(t)$ and $y(t)$, where the latter are defined as in Fig. 6.1. This second approach will also become much handier than the first one once we introduce the feedback as in Fig. 6.4.

Wet signal inversion

By inverting the wet signal, one swaps the peaks and the notches. Notice that the phase response of differential allpasses at $\omega = 0$ can be either 0° or 180° , the same holds for the phase response at $\omega = +\infty$. For that reason the possibility to swap the peaks and the notches might be handy.

Notch spacing

In the simplest case one uses a series of identical 1-pole allpasses inside a phaser. In order to control the notch spacing in an easy and nice way, one should rather use a series of identical 2-pole allpasses. As mentioned earlier, by changing the resonance amount of the 2-pole allpasses one controls the phase slope of the filters. This affects the spacing of the notches (Fig. 6.3).

Feedback

We can also introduce feedback into the phaser. Similarly to the case of the ladder filter modes, the dry signal is better picked up after the feedback point (Fig. 6.4) The feedback changes the shape of the peaks and notches (Fig. 6.5).

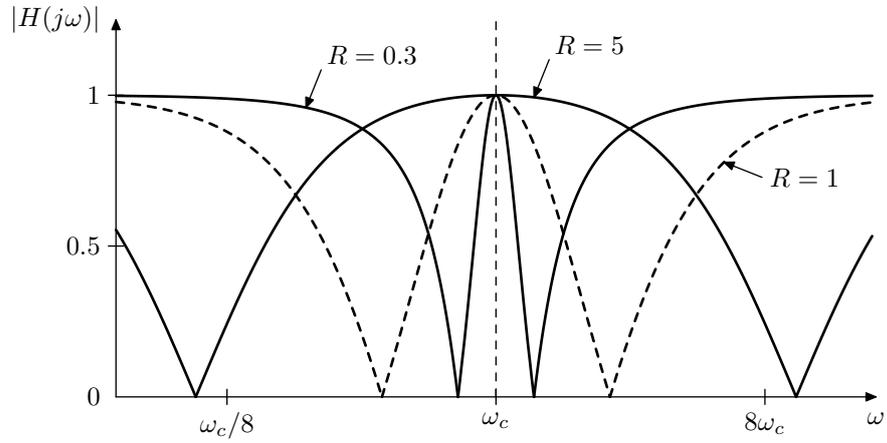


Figure 6.3: Effect of the allpass resonance on the notch spacing (using two 2-pole allpass filters as the allpass core of the phaser).

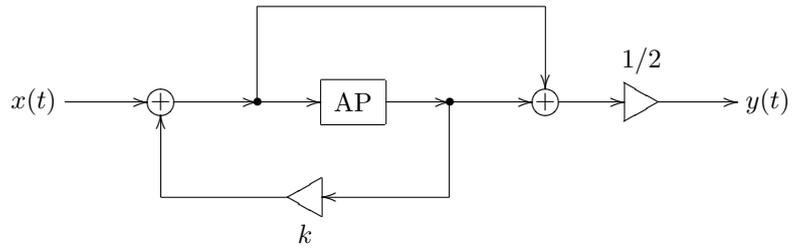


Figure 6.4: Phaser with feedback.

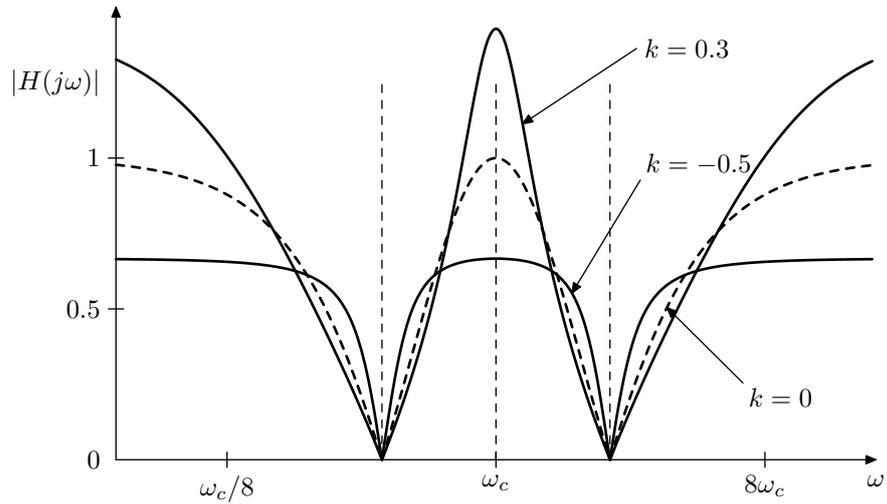


Figure 6.5: Effect of the feedback amount in Fig. 6.4 on the notch and peak shapes.

With the introduction of feedback we have a zero-delay feedback loop in the phaser structure. It can be solved using typical TPT means.³

6.2 Flangers

Flangers can be obtained from phasers by an allpass substitution.

A delay is a linear time-invariant allpass. It even has a transfer function $H(s) = e^{-sT}$ where T is the delay time. Obviously $|H(s)| = |e^{-sT}| = 1$. However it is not a differential filter, for that reason the transfer function is not a rational function of s . Digital delay models are typically built using interpolation techniques, the details of which fall outside the scope of this book.

Using the allpass substitution principle we can replace the allpass filter chain in a phaser by a delay. This produces a *flanger*.⁴ The discussion of the phasers mostly didn't assume any details about the underlying allpass, therefore most of it is applicable to flangers.

The main difference with using a delay is that the 0° and 180° phase response points are evenly spaced in the linear frequency scale (Fig. 6.6), whereas the spacing of the same points in responses of differential allpasses is not that regular. Also, a delay's phase response can easily have lots of 0° and 180° points (the larger the delay time is, the more of those points it has within the audible frequency range), while the number of those points in a differential allpass filter's phase response is limited by the filter's order.

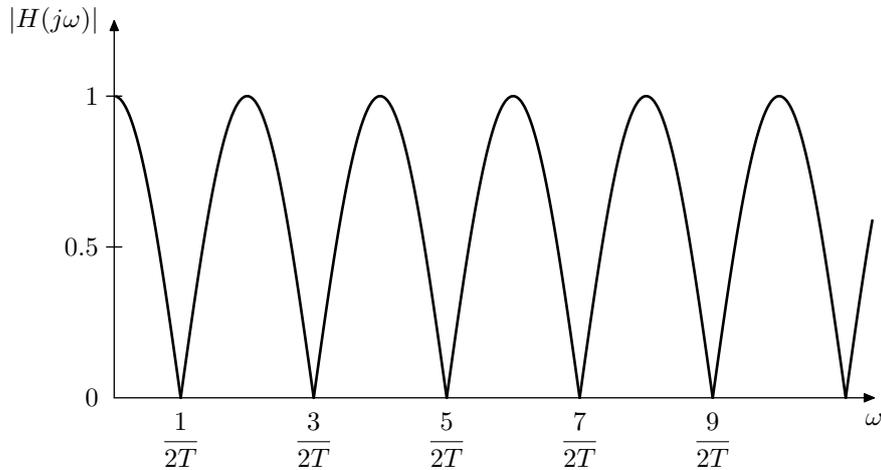


Figure 6.6: Amplitude response of the simplest flanger using the structure from Fig. 6.1.

Rather than modulating the delay time linearly by an LFO, one should consider that a filter's cutoff should be typically modulated in the logarithmic frequency scale (a.k.a. the pitch scale), therefore one in principle should do the same for the delay in a flanger. The delay's cutoff for that purpose can be simply defined as $\omega_c = 2\pi/T$, where T is the delay time.

³Inserting a unit delay in the feedback produces subtle but rather unpleasant artifacts in the phasing response, one should better use the TPT approach here.

⁴In the absence of an LFO the structure is referred to as a *comb filter*.

SUMMARY

A phaser is made of an allpass differential filter connected in parallel with the dry signal path. This creates notches at the points of 180° phase difference and peaks at 0° points. The allpass cutoff should be modulated by an LFO. Using a delay instead of a differential allpass creates a flanger. Feedback can be used to change the shape of the peaks and notches in the amplitude response.

Chapter 7

Frequency shifters

Frequency shifters are a musical application of the radio transmission technique known as *single-sideband modulation*. Despite sounding simple from the name, the construction of frequency shifters (or more specifically, the computation of the coefficients of filters used therein) is somewhat complicated. Therefore a somewhat higher math skill level is generally required by the materials discussed in this chapter compared to the other chapters in this book.

The text also refers in a few places to Butterworth and elliptic filters as well as to elliptic rational functions, whose introduction is beyond the scope of this book. A few reading suggestions regarding elliptic filters and elliptic rational functions are made at the end of the chapter.

7.1 General ideas

According to Fourier transform properties, the shifting of the signal in the time domain corresponds to modulation by a complex sinusoid in the frequency domain. The dual of that property is that the shifting of the signal in the frequency domain corresponds to modulation by a complex sinusoid in the time domain. More specifically, let

$$x(t) = \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} \frac{d\omega}{2\pi} \quad (7.1)$$

Then, a frequency-shifted version of $x(t)$ is

$$\int_{-\infty}^{\infty} X(\omega) e^{j(\omega+\Delta\omega)t} \frac{d\omega}{2\pi} = \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} e^{j\Delta\omega t} \frac{d\omega}{2\pi} = e^{j\Delta\omega t} \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} \frac{d\omega}{2\pi}$$

Thus in order to shift the frequencies of all partials of the signal by $\Delta\omega$ we need to simply multiply the signal by $e^{j\Delta\omega t}$. The problem is, however, that if $x(t)$ was originally a real signal (that is $X(\omega)$ was Hermitian), then after the multiplication by a complex sinusoid $e^{j\Delta\omega t}$ it won't be real anymore (corresponding to the fact that a shifted Hermitian spectrum is not Hermitian anymore).

So, how do we frequency-shift a real signal, so that the resulting signal is real as well? Let

$$x(t) = \int_0^{\infty} a(\omega) \cos(\omega t + \varphi(\omega)) \frac{d\omega}{2\pi}$$

We wish to obtain

$$y(t) = \int_0^{\infty} a(\omega) \cos((\omega + \Delta\omega)t + \varphi(\omega)) \frac{d\omega}{2\pi} \quad (7.2)$$

Notably, if $\Delta\omega < 0$, then some of the frequencies $\omega + \Delta\omega$ in (7.2) will be negative and will alias with the positive frequencies of the same absolute magnitude. This can be either ignored, or $x(t)$ can be prefiltered to make sure it doesn't contain frequencies below $-\Delta\omega$. So, except for the just mentioned highpass prefiltering option, the possible aliasing of the negative frequencies doesn't affect the subsequent discussion.

We can rewrite (7.2) as

$$\begin{aligned} y(t) &= \int_0^{\infty} a(\omega) \cos((\omega + \Delta\omega)t + \varphi(\omega)) \frac{d\omega}{2\pi} = \\ &= \int_0^{\infty} a(\omega) \cos((\omega + \Delta\omega)t + \varphi(\omega)) \frac{d\omega}{2\pi} = \\ &= \int_0^{\infty} a(\omega) \cos(\Delta\omega t + \omega t + \varphi(\omega)) \frac{d\omega}{2\pi} = \\ &= \int_0^{\infty} a(\omega) \left(\cos \Delta\omega t \cos(\omega t + \varphi(\omega)) - \sin \Delta\omega t \sin(\omega t + \varphi(\omega)) \right) \frac{d\omega}{2\pi} = \\ &= \cos \Delta\omega t \cdot \int_0^{\infty} a(\omega) \cos(\omega t + \varphi(\omega)) \frac{d\omega}{2\pi} - \\ &\quad - \sin \Delta\omega t \cdot \int_0^{\infty} a(\omega) \sin(\omega t + \varphi(\omega)) \frac{d\omega}{2\pi} = \\ &= \cos \Delta\omega t \cdot \int_0^{\infty} a(\omega) \cos(\omega t + \varphi(\omega)) \frac{d\omega}{2\pi} - \\ &\quad - \sin \Delta\omega t \cdot \int_0^{\infty} a(\omega) \cos\left(\omega t + \varphi(\omega) - \frac{\pi}{2}\right) \frac{d\omega}{2\pi} = \\ &= x(t) \cos \Delta\omega t - x_{-90}(t) \sin \Delta\omega t \end{aligned} \quad (7.3)$$

where

$$x_{-90}(t) = \int_0^{\infty} a(\omega) \cos\left(\omega t + \varphi(\omega) - \frac{\pi}{2}\right) \frac{d\omega}{2\pi}$$

is a signal obtained from $x(t)$ by phase-shifting all partials by -90° .

If (7.1) is a complex spectrum of a real $x(t)$ then the complex spectrum of $x_{-90}(t)$ must be

$$x_{-90}(t) = \int_0^{\infty} j^{-1} X(\omega) e^{j\omega t} \frac{d\omega}{2\pi} + \int_{-\infty}^0 j X(\omega) e^{j\omega t} \frac{d\omega}{2\pi}$$

that is all positive frequency partials need to be shifted by -90° and all negative frequency partials need to be shifted by $+90^\circ$. That is

$$X_{-90}(\omega) = X(\omega) \cdot j \operatorname{sgn} \omega$$

or, in the Laplace transform notation

$$X_{-90}(j\omega) = X(j\omega) \cdot j \operatorname{sgn} \omega$$

We wonder whether we could construct a filter, whose frequency response is $H(j\omega) = j \operatorname{sgn} \omega$. If we succeed, it would be trivial to use this filter to build a frequency shifter.

7.2 Analytic signals

There is an alternative way to look at the same topic. Consider a complex signal

$$\begin{aligned}
 v(t) &= x(t) + jx_{-90}(t) = \\
 &= \int_{-\infty}^{\infty} X(\omega)e^{j\omega t} \frac{d\omega}{2\pi} + j \int_0^{\infty} j^{-1} X(\omega)e^{j\omega t} \frac{d\omega}{2\pi} + j \int_{-\infty}^0 j X(\omega)e^{j\omega t} \frac{d\omega}{2\pi} = \\
 &= \int_0^{\infty} 2X(\omega)e^{j\omega t} \frac{d\omega}{2\pi} + 0 \tag{7.4}
 \end{aligned}$$

So, $v(t)$ doesn't contain any negative frequency partials. A signal which contains only positive frequency partials is called *analytic*. Clearly, the spectrum of an analytic signal is not Hermitian, thus an analytic signal cannot be purely real.

The transformation that converts a real part of an analytic signal to the imaginary one is called *Hilbert transform*. Two real signals are said to form a *Hilbert transform pair* if they are a real and an imaginary part of some analytic signal. So $x(t)$ and $x_{-90}(t)$ form a Hilbert transform pair.

A signal processing algorithm which performs Hilbert transform is referred to as a *Hilbert transformer*. Common usages of Hilbert transforms are in envelope followers to estimate the signal's amplitude and in *frequency shifters*. A number of discrete-time Hilbert transform algorithms are available. However, since this book concentrates on the VA filters, it is appropriate to consider a method of designing a Hilbert transformer in continuous time.

According to our previous discussion there are two equivalent approaches to build a Hilbert transformer:

- build a complex filter defined by

$$V(j\omega) = X(j\omega) \cdot 2H_{>0}(j\omega) \quad \text{where } H_{>0}(j\omega) = \begin{cases} 1 & \text{for } \omega > 0 \\ 0 & \text{for } \omega < 0 \end{cases} \tag{7.5}$$

and take the imaginary part of the output¹

- build a $\pm 90^\circ$ phase shifter, thereby immediately obtaining the imaginary part of the analytic signal:

$$X_{-90}(j\omega) = X(j\omega)H_{-90}(j\omega) \quad \text{where } H_{-90}(j\omega) = -j \operatorname{sgn} \omega \tag{7.6}$$

Apparently the two approaches are related by the decomposition of $H_{>0}(s)$ into its real and imaginary parts:

$$2H_{>0}(s) = 1 + jH_{-90}(s) \tag{7.7}$$

7.3 Phase splitter

So far we have obtained the two expressions (7.5) and (7.6) for an ideal Hilbert transformer. Such implementations are however not possible in practice. The

¹This idea was taken by the author from *Design of multiplierless elliptic IIR halfband filters and Hilbert transformers* by M.D.Lutovac and L.D.Milic (proc. Eusipco-98), where it is further attributed to *Special Filter Design* by P.A.Regalia.

problem is that neither (7.5) nor (7.6) can be implemented by a finite-order filter, since there are no rational transfer functions satisfying (7.5) or (7.6). So, practical implementations of Hilbert transformers are always approximations thereof.

Furthermore, considering (7.5), we can notice that it cannot be implemented by a stable differential system, since stable rational transfer functions cannot have zero phase response. Indeed, a zero phase response implies that $H_{>0}(j\omega)$ is a real rational function of ω and therefore its complex poles should be mutually conjugate in the complex ω -plane. Since $s = j\omega$, each such conjugate pole pair in the ω -plane corresponds to a pair consisting of a stable and an unstable pole in the s -plane.

A counterpart to that issue is that the phase response (7.6) doesn't correspond (even approximately) to a stable allpass. Indeed, any stable allpass can be represented as a serial combination of stable 1- and 2-pole allpasses (or even just 1-poles, if we allow complex 1-pole allpasses). Recall the phase response of a stable 1-pole allpass (Fig. 7.1). In the vicinity of -90° this phase response has the steepest slope. Connecting more 1-pole allpasses in series can only make the situation worse. Using 2-poles is even further worse, since their phase responses (Fig. 7.2) are even steeper than those of 1-poles.

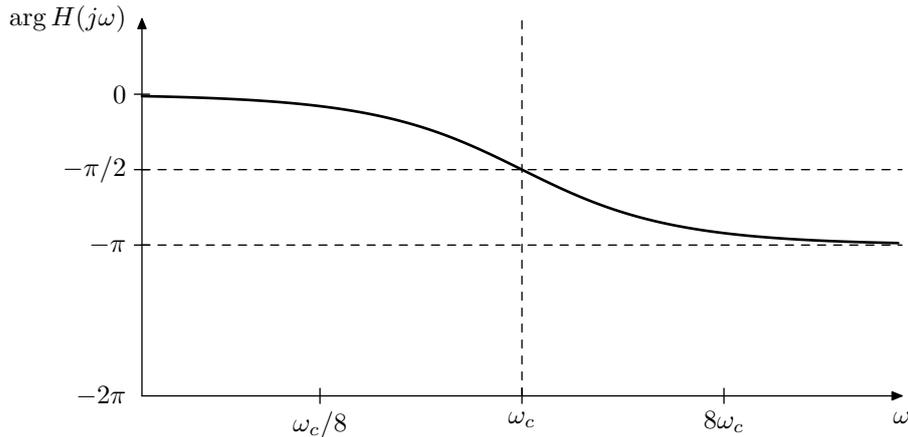


Figure 7.1: Phase response of a 1-pole allpass filter.

Thus stable approximations of (7.5) and (7.6) are not possible. A nonstable approximation of (7.5) can be however converted into a stable one by introducing an additional allpass transformation of the signal. Indeed, let p_{n+} be the unstable poles of $H_{>0}(s)$. Consider an unstable allpass $H_+(s)$ whose poles are p_{n+} :

$$H_+(s) = \prod_n \frac{-p_{n+}^* - s}{p_{n+} - s}$$

Respectively $H_+^{-1}(s)$ is a stable allpass. Considering the product $H_+^{-1}(s)H_{>0}(s)$ we notice that it defines a stable filter, since the unstable poles of $H_{>0}(s)$ are cancelled by the zeros of $H_+^{-1}(s)$. So, while $H_{>0}$ is not a stable filter, we could build the stable filter $H_+^{-1}H_{>0}$ differing from $H_{>0}$ only by shifted phases in the output signal.

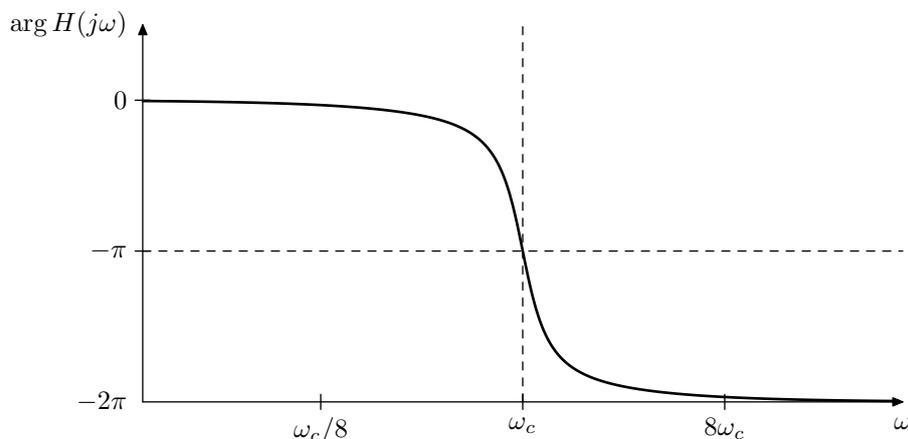


Figure 7.2: Phase response of a 2-pole allpass filter.

Multiplying both sides of (7.7) by $H_+^{-1}(s)$ we obtain

$$2H_{>0}(s)H_+^{-1}(s) = H_+^{-1}(s) + jH_+^{-1}(s)H_{-90}(s) \quad (7.8)$$

Noticing that (7.7) implies that $H_{>0}(s)$ and $H_{-90}(s)$ have identical poles, we conclude that an unstable allpass $H_{-90}(s)$ is thereby converted into a stable allpass $H_+^{-1}(s)H_{-90}(s)$. That is a pair of *stable* allpasses $H_+^{-1}(s)$ and $H_+^{-1}(s)H_{-90}(s)$ produces the real and imaginary parts of an analytical signal, where the latter differs from the analytical signal in (7.6) just by phase shifting.

Practically this means the following. Given an unstable allpass $H_{-90}(s)$ we decompose it into a product of stable and unstable parts:

$$H_{-90}(s) = H_-(s)H_+(s) = \frac{H_-(s)}{H_+^{-1}(s)} \quad (7.9)$$

The inverted unstable part H_+^{-1} then produces the real part of the analytic signal and the stable part H_- produces the imaginary part of the analytic signal (Fig. 7.3). The structure in Fig. 7.3 is referred to as the *phase splitter*.

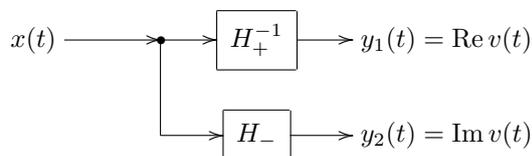


Figure 7.3: Phase splitter.

7.4 Implementation structure

Using (7.3) we can turn the phase splitter in Fig. 7.3 into a frequency shifter. By using $y_1(t)$ and $y_2(t)$ instead of $x(t)$ and $x_{-90}(t)$ the equation (7.3) is turned

into

$$y(t) = y_1(t) \cos \Delta\omega t - y_2(t) \sin \Delta\omega t$$

corresponding to the structure in Fig. 7.4

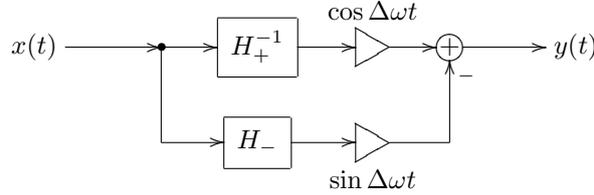


Figure 7.4: Frequency shifter.

Notably, replacing $\Delta\omega$ by $-\Delta\omega$ in (7.3) we obtain

$$\int_0^\infty a(\omega) \cos((\omega - \Delta\omega)t + \varphi(\omega)) \frac{d\omega}{2\pi} = x(t) \cos \Delta\omega t + x_{-90}(t) \sin \Delta\omega t \quad (7.10)$$

This means that we can extend the frequency shifter in Fig. 7.4 to a one that shifts simultaneously in both directions, obtaining the diagram in Fig. 7.5.

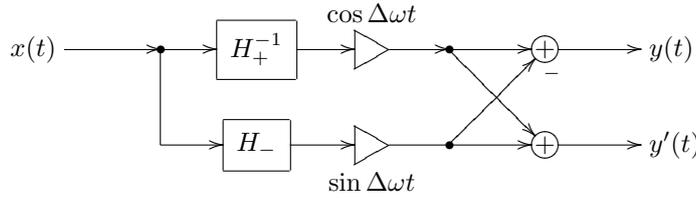


Figure 7.5: A bidirectional frequency shifter.

Adding together the frequency-shifted signals from (7.3) and (7.10) we notice that

$$\begin{aligned} & \int_0^\infty a(\omega) \cos((\omega + \Delta\omega)t + \varphi(\omega)) \frac{d\omega}{2\pi} + \int_0^\infty a(\omega) \cos((\omega - \Delta\omega)t + \varphi(\omega)) \frac{d\omega}{2\pi} = \\ & = x(t) \cos \Delta\omega t - x_{-90}(t) \sin \Delta\omega t + x(t) \cos \Delta\omega t + x_{-90}(t) \sin \Delta\omega t = \\ & = 2x(t) \cos \Delta\omega t \end{aligned}$$

That is, the sum of y and y' in Fig. 7.5 produces a ring modulation between $y_1(t)$ (which is a phase-shifted version of $x(t)$) and $\cos(\Delta\omega t)$. So frequency-shifting and ring-modulation by a sinusoid seem are very closely related. The same can be analyzed in the complex spectral domain:

$$\begin{aligned} \cos \Delta\omega t \cdot x(t) &= \frac{e^{j\Delta\omega t} + e^{-j\Delta\omega t}}{2} \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} \frac{d\omega}{2\pi} = \\ &= \frac{1}{2} \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} e^{j\Delta\omega t} \frac{d\omega}{2\pi} + \frac{1}{2} \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} e^{-j\Delta\omega t} \frac{d\omega}{2\pi} = \\ &= \frac{1}{2} \int_{-\infty}^{\infty} X(\omega) e^{j(\omega + \Delta\omega)t} \frac{d\omega}{2\pi} + \frac{1}{2} \int_{-\infty}^{\infty} X(\omega) e^{j(\omega - \Delta\omega)t} \frac{d\omega}{2\pi} \end{aligned}$$

Thus in the case of the ring modulation by a sinusoid, the partials are frequency-shifted in both directions.

So now we know how to construct a frequency shifter, except that we still do not know how to obtain $H_{>0}(s)$ or $H_{-90}(s)$, so that we can obtain from them the allpasses H_+ and H_- , and this is exactly what we shall discuss next. Analytical expressions for obtaining $H_{-90}(s)$ involve the usage of *elliptic rational functions* which are a relatively esoteric subject. This was more or less the only reasonable approach when the computational powers of modern personal computers were not available. However, nowadays a straightforward minimax optimization approach can be taken instead, as long as runtime recomputation of the coefficients is not required. We will discuss both methods.

7.5 Remez algorithm

Suppose we are given a function $f(x)$ and its approximation $\tilde{f}(x)$. There are different ways to measure the quality of the approximation. One way to measure this quality is the maximum error of the approximation on the given interval of interest $x \in [a, b]$:

$$E = \max_{[a,b]} |\tilde{f}(x) - f(x)| \quad (7.11)$$

We therefore wish to minimize the value of E . That is we want to minimize the maximum error of the approximation. Such approximations are hence called *minimax* approximations.²

Gradient search methods do not work well for minimax optimizations. Therefore a different method, called *Remez algorithm*,³ needs to be used. As of today, internet resources concerning the Remez algorithm are quite scarce, nor does this method seem to be a subject of common math textbooks. This might suggest that Remez algorithm belongs to a rather esoteric math area. The algorithm itself, however, is very simple. We will therefore cover the essentials of that algorithm in this book.⁴

Suppose $\tilde{f}(x)$ is a polynomial:

$$\tilde{f}(x) = \sum_{n=0}^N a_n x^n \quad (7.12)$$

Apparently, there are $N + 1$ degrees of freedom in the choice of $\tilde{f}(x)$, each degree corresponding to one of the coefficients a_n . Therefore we can force the function $\tilde{f}(x)$ to take arbitrarily specified values at $N + 1$ arbitrarily chosen points \bar{x}_n . Particularly, we can require

$$\tilde{f}(\bar{x}_n) = f(\bar{x}_n) \quad n = 0, \dots, N$$

²The maximum of the absolute value of a function is also the L_∞ norm of the function. Therefore minimax approximations are optimizations of the L_∞ norm.

³The Remez algorithm should not be confused with the Parks–McClellan algorithm. The latter is a specific restricted version of the former. For whatever reason, the Parks–McClellan algorithm is often referred to as the Remez algorithm in the signal processing literature.

⁴The author's primary resource for the information about the Remez algorithm was the documentation for the math toolkit of the *boost* library by J.Maddock, P.A.Bristow, H.Holin and X.Zhang.

or equivalently require the error to be zero at \bar{x}_n :

$$\tilde{f}(\bar{x}_n) - f(\bar{x}_n) = 0 \quad n = 0, \dots, N \quad (7.13)$$

(notice that the equations (7.13) are linear in respect to the unknowns a_n and therefore are easily solvable). If the points \bar{x}_n are approximately uniformly spread over the interval of interest $[a, b]$ then intuitively we can expect $\tilde{f}(x)$ to be a reasonably good approximation of $f(x)$ (Fig. 7.6).

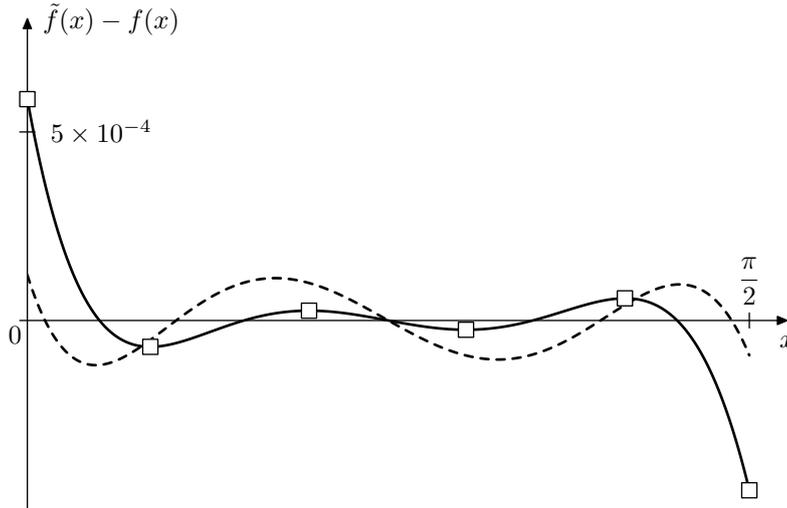


Figure 7.6: The error of the 4-th order polynomial approximations of $\sin x$ on $[0, \pi/2]$. The approximation with uniformly spaced zeros at 9° , 27° , 45° , 63° , 91° (solid line) and the one with Chebyshev zeros (dashed line). The empty square-shaped dots at the extrema of the error are the control points of the Remez algorithm.

This based on the uniform zero spacing approximation is however not the best one. Indeed, instead let \bar{x}_n equal the (properly scaled) zeros of the Chebyshev polynomial of order $N + 1$:

$$\bar{x}_n = \frac{a+b}{2} + \frac{b-a}{2}z_n \quad \bar{x}_n \in (a, b) \quad z_n \in (-1, 1)$$

$$T_{N+1}(z_n) = \cos((N+1) \arccos z_n) = 0$$

$$z_n = -\cos\left(\frac{1}{2} + n\right) \frac{\pi}{N+1} \quad n = 0, \dots, N$$

where $T_N(x) = \cos(N \arccos x)$ is the Chebyshev polynomial of order N and where the minus sign in front of the cosine ensures that z_n are in ascending order. Comparing Chebyshev zeros approximation (the dashed line in Fig. 7.6) to the uniform zeros approximation, we can see that the former is much better than the latter, at least in the minimax sense.

A noticeable property of the Chebyshev zeros approximation clearly observable in in Fig. 7.6 is that the extrema of the approximation error (counting the extrema at the boundaries of the interval $[a, b]$!) are approximately equal in absolute magnitude and have alternating signs. This is a characteristic trait of

minimax approximations: the error extrema are equal in magnitude and alternating in sign.

So, we might attempt to build a minimax approximation by trying to satisfy the *equiripple error oscillation* requirement. That is, instead of seeking to minimize the maximum error, we simply seek an error which oscillates between the two boundaries of opposite sign and equal absolute value. Somewhat surprisingly, this is a much simpler task.

Intuitive description of Remez algorithm

Consider the solid line graph in Fig. 7.6. Intuitively, imagine a “control point” at each of the extrema. Now we “take” the control point which has the largest error (the one at $x = 0$) and attempt to move it towards the x axis, reducing the error value at $x = 0$. Since there are 6 control points (4 at local extrema plus 2 at the boundaries), but only 5 degrees of freedom (corresponding to the coefficients a_n), at least one of the other control points needs to move (or several or all of them can move). Intuitively it’s clear that if we lower the error at $x = 0$, then it will grow at some other points of $[a, b]$. However, since we have the largest error at $x = 0$ anyway, we can afford the error growing elsewhere on $[a, b]$, at least for a while. Notice that during such change the x positions of control points will also change, since the extrema of the error do not have to stay at the same x coordinates.

As the error elsewhere at $[a, b]$ becomes equal in absolute magnitude to the one at $x = 0$, we have two largest-error control points which need to be moved simultaneously from now on. This can be continued until only one “free” control point remains. Simultaneously reducing the error at 5 of 6 control points we thereby increase the error at the remaining control point. At some moment both errors will become equal in absolute magnitude, which means that the error at all control points is equal in absolute magnitude. Since the control points are located at the error extrema, we have thereby an equiripple oscillating error.

Remez algorithm for polynomial approximation

Given $\tilde{f}(x)$ which is a polynomial (7.12), the process of “pushing the control points towards zero” has a simple algorithmic expression. Indeed, we seek $\tilde{f}(x)$ which satisfies

$$\tilde{f}(\hat{x}_n) + (-1)^n \varepsilon = f(\hat{x}_n) \quad n = 0, \dots, N + 1 \quad (7.14)$$

where \hat{x}_n are the (unknown) control points (including $\hat{x}_0 = a$ and $\hat{x}_{N+1} = b$) and ε is the (unknown) signed maximum error. Thus, the unknowns in (7.14) are a_n (the polynomial coefficients), \hat{x}_n (the control points at the extrema) and ε (the signed maximum error). Notice that the equations (7.14) are linear in respect to a_n and ε , which leads us to the following idea.

Suppose we already have some initial guess for $\tilde{f}(x)$, like the uniform zero polynomial in Fig. 7.6 (or the Chebyshev zero polynomial, which is even better). Identifying the extrema of $\tilde{f}(x) - f(x)$ we obtain a set of control points \hat{x}_n . Now, given these \hat{x}_n , we simply solve (7.14) for a_n and ε (where we have $N + 2$ equations and $N + 2$ unknowns in total), thereby obtaining a new set of a_n . In a way this is cheating, because \hat{x}_n are not the control points anymore, since they are not anymore the extrema of the error (and if they were, we

would already have obtained a minimax approximation by simply finding these new a_n). However, the polynomial defined by the new a_n has a much better maximum error (Fig. 7.7)!

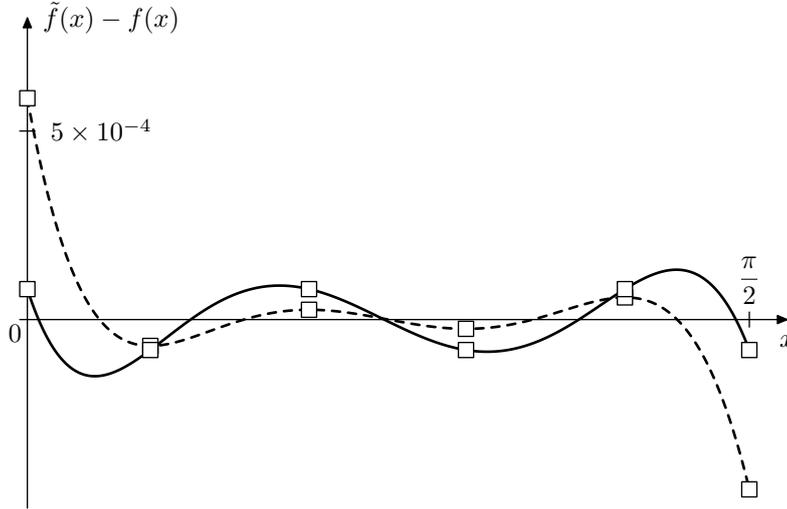


Figure 7.7: The approximation error before (dashed line) and after (solid line) a single step of the Remez polynomial approximation algorithm. The empty square-shaped dots are the control points.

So we simply update the control points \hat{x}_n to the new positions of the extrema and solve (7.14) again. Then again update the control points and solve (7.14) and so on. This is the Remez algorithm for polynomial approximation. We still need to refine some details about the algorithm though.

- The function $f(x)$ should be reasonably well-behaved (whatever that could mean) in order for Remez algorithm to work.
- As a termination condition for the iteration we can simply check the equiripple property of the error at the control points. That is, having obtained the new a_n , we find the new control points \hat{x}_n and then compute the errors $\varepsilon_n = \tilde{f}(\hat{x}_n) - f(\hat{x}_n)$. If the absolute values of ε_n are equal up to the specified precision, this means that we have an approximation which is minimax up to the specified error, and the algorithm may be stopped.
- The initial approximation $\tilde{f}(x)$ needs to have the alternating sign property. This is more or less ensured by using (7.13) to construct the initial approximation. A good choice for \tilde{x}_n (as demonstrated by Fig. 7.6) are the roots of the Chebyshev polynomial of order one higher than the order of the approximating polynomial $\tilde{f}(x)$.⁵
- The control points \hat{x}_n are the zeros of the error derivative $(\tilde{f} - f)'$ (except for $\hat{x}_0 = a$ and $\hat{x}_{N+1} = b$). There is exactly one local extremum on each interval $(\tilde{x}_n, \tilde{x}_{n+1})$ between the zeros of the error. Therefore, \hat{x}_{n+1} can

⁵This becomes kind of intuitive after considering Chebyshev polynomials as *some kind* of minimax approximations of the zero constant function $f(x) \equiv 0$ on the interval $[-1, 1]$.

be simply found as the zeros of the error derivative by bisection of the intervals $(\bar{x}_n, \bar{x}_{n+1})$.

- After having obtained new a_n , the old control points \hat{x}_n are not the extrema anymore, however the errors at \hat{x}_n are still alternating in sign. Therefore the new zeros \bar{x}_n (needed to find the new control points by bisection) can be found by bisection of the intervals $(\hat{x}_n, \hat{x}_{n+1})$.

Restrictions and variations

Often it is desired to obtain a function which is odd or even, or has some other restrictions. This can be done by simply fixing the respective a_n , thereby reducing the number of control variables a_n and reducing the number of control points \hat{x}_n and zero crossings \bar{x}_n accordingly.

Remez algorithm can also be easily modified to accommodate a weight function in the minimax norm (7.11):

$$E = \max_{[a,b]} \left(W(x) \cdot \left| \tilde{f}(x) - f(x) \right| \right) \quad W(x) > 0$$

The error function therefore turns into $W(x)(\tilde{f}(x) - f(x))$, while the minimax equations (7.14) turn into

$$\tilde{f}(\hat{x}_n) + (-1)^n W^{-1}(\hat{x}_n) \varepsilon = f(\hat{x}_n) \quad n = 0, \dots, N + 1$$

(where $W^{-1}(x)$ is the reciprocal of $W(x)$).

Remez algorithm for rational approximation

Instead of using a polynomial $\tilde{f}(x)$, better approximations can be often achieved by rational $\tilde{f}(x)$:

$$\tilde{f}(x) = \frac{\sum_{n=0}^N a_n x^n}{1 + \sum_{n=1}^M b_n x^n} \quad (7.15)$$

Besides being able to deliver better approximations in certain cases, rational functions can be often useful for obtaining approximations on infinite intervals such as $[a, +\infty)$, because by varying the degrees of the numerator and denominator the asymptotic behavior of $\tilde{f}(x)$ at $x \rightarrow \infty$ can be controlled.

For a rational $\tilde{f}(x)$ defined by (7.15) the minimax equations (7.14) become nonlinear in respect to the unknowns ε and b_n , although they are still linear in respect to the unknowns a_n :

$$\sum_{i=0}^N a_i \hat{x}_n^i + (-1)^n \left(1 + \sum_{i=1}^M b_i \hat{x}_n^i \right) \varepsilon = \left(1 + \sum_{i=1}^M b_i \hat{x}_n^i \right) f(\hat{x}_n) \quad (7.16)$$

$$n = 0, \dots, N + M + 1$$

Notice that the number of degrees of freedom is now $N + M + 1$. The equations (7.16) can be solved using different numeric methods for nonlinear equation

solution, however there is one simple trick.⁶ Rewrite (7.16) as

$$\sum_{i=0}^N a_i \hat{x}_n^i + (-1)^n \varepsilon \sum_{i=1}^M b_i \hat{x}_n^i + (-1)^n \varepsilon = \left(1 + \sum_{i=1}^M b_i \hat{x}_n^i \right) f(\hat{x}_n)$$

Now we pretend we don't know the free term ε , but we do know the value of ε before the sum of $b_i \hat{x}_n^i$:

$$\sum_{i=0}^N a_i \hat{x}_n^i + (-1)^n \varepsilon_0 \sum_{i=1}^M b_i \hat{x}_n^i + (-1)^n \varepsilon = \left(1 + \sum_{i=1}^M b_i \hat{x}_n^i \right) f(\hat{x}_n) \quad (7.17)$$

where ε_0 is this "known" value of ε . The value of ε_0 can be estimated e.g. as the average absolute error at the control points \hat{x}_n . Then (7.17) are linear equations in respect to a_n , b_n and ε and can be easily solved. Having obtained the new a_n and b_n , we can obtain a new estimation for ε_0 and solve (7.17) again. We repeat until the errors $\tilde{f}(\hat{x}_n) - f(\hat{x}_n)$ at the control points \hat{x}_n become equal in absolute value up to a necessary precision. At this point we can consider the solution of (7.16) as being obtained to a sufficient precision and proceed with the usual Remez algorithm routine (find the new \bar{x}_n , new \hat{x}_n etc.)

Here are some further notes.

- In principle the solution of (7.16) doesn't need to be obtained to a very high precision, except in the final step of the Remez algorithm. However, in order to know whether the current step is the final one or not, we need to know the true control points, so that we can estimate how well the equiripple condition is satisfied. Ultimately, this is a question of the computational expense of finding the new control points vs. computing another iteration of (7.17).
- Sometimes, if the equations are strongly nonlinear, the trick (7.17) may fail to converge. In this case one could attempt to use the discussed below more general Newton–Raphson approach (7.23), where the damping parameter may be used to mitigate the convergence problems.
- In regards to the problem of choice of the initial $\tilde{f}(x)$ for the rational Remez approximation, notice that the zero error equations (7.13) take the form

$$\sum_{n=0}^N a_n \bar{x}^n = f(\bar{x}_n) \left(1 + \sum_{n=1}^M b_n \bar{x}^n \right)$$

which is fully linear in respect to a_n and b_n , and can be easily solved.

Other kinds of approximating functions

In certain cases one could use even more complicated forms of $\tilde{f}(x)$, which are neither polynomial nor rational. In the general case such function $\tilde{f}(x)$ is controlled by a number of parameters a_n :

$$\tilde{f}(x) = \tilde{f}(x, a_1, a_2, \dots, a_N)$$

⁶This trick is adapted from the *boost* library documentation and sources.

(notice that this time the numbering of a_n is starting at one, so that there are N parameters in total, giving N degrees of freedom). The minimax equations (7.14) become

$$\tilde{f}(\hat{x}_n, a_1, a_2, \dots, a_N) + (-1)^n \varepsilon = f(\hat{x}_n) \quad n = 0, \dots, N \quad (7.18)$$

Introducing functions

$$\phi_n(a_1, a_2, \dots, a_N, \varepsilon) = \tilde{f}(\hat{x}_n, a_1, a_2, \dots, a_N) + (-1)^n \varepsilon - f(\hat{x}_n)$$

we rewrite the equations (7.18) as

$$\phi_n(a_1, a_2, \dots, a_N, \varepsilon) = 0 \quad n = 0, \dots, N \quad (7.19)$$

Introducing vector notation

$$\begin{aligned} \mathbf{\Phi} &= (\phi_0 \ \phi_1 \ \dots \ \phi_N)^T \\ \mathbf{a} &= (a_1 \ a_2 \ \dots \ a_N \ \varepsilon)^T \end{aligned}$$

we rewrite (7.19) as

$$\mathbf{\Phi}(\mathbf{a}) = 0 \quad (7.20)$$

Apparently, (7.20) is a vector form of (7.14), except that now we consider it as a generally nonlinear equation. Both the function's argument \mathbf{a} and the function's value $\mathbf{\Phi}(\mathbf{a})$ have the dimension $N + 1$, therefore the equation (7.20) is fully defined.

Different numeric methods can be applied to solving (7.20). We will be particularly interested in the application of multidimensional Newton–Raphson method. Expanding $\mathbf{\Phi}(\mathbf{a})$ into Taylor series at some fixed point \mathbf{a}_0 we transform (7.20) into:

$$\mathbf{\Phi}(\mathbf{a}_0) + \frac{\partial \mathbf{\Phi}}{\partial \mathbf{a}}(\mathbf{a}_0) \cdot \Delta \mathbf{a} + o(\Delta \mathbf{a}) = 0 \quad (7.21)$$

where $\partial \mathbf{\Phi} / \partial \mathbf{a}$ is the Jacobian matrix and $\mathbf{a} = \mathbf{a}_0 + \Delta \mathbf{a}$. By discarding the higher order terms $o(\Delta \mathbf{a})$, the equation (7.21) is turned into

$$\Delta \mathbf{a} = - \left(\frac{\partial \mathbf{\Phi}}{\partial \mathbf{a}}(\mathbf{a}_0) \right)^{-1} \cdot \mathbf{\Phi}(\mathbf{a}_0) \quad (7.22)$$

The equation (7.22) implies the Newton–Raphson iteration scheme

$$\mathbf{a}_{n+1} = \mathbf{a}_n - \alpha \cdot \left(\frac{\partial \mathbf{\Phi}}{\partial \mathbf{a}}(\mathbf{a}_0) \right)^{-1} \cdot \mathbf{\Phi}(\mathbf{a}_0) \quad (7.23)$$

where the damping factor α is either set to unity, or to a lower value, if the nonlinearity of $\mathbf{\Phi}(\mathbf{a})$ is too strong and prevents the iterations from converging. The initial value \mathbf{a}_0 is obtained from the initial settings of the parameters a_n and the estimated initial value of ε . As for the rational $\tilde{f}(x)$, the initial value of ε can be estimated e.g. as the average error at the control points.

Similarly to the rational approximation case, the solution of (7.20) doesn't need to be obtained to a very high precision during the intermediate steps of the Remez algorithm. However the same tradeoff between computing the iteration step (7.23) and finding the new control points applies.

The choice of the initial $\tilde{f}(x)$ can be done based on the same principles. The zero error equations (7.13) turn into

$$\phi_n(a_1, a_2, \dots, a_N, 0) = 0 \quad n = 1, \dots, N$$

(notice that compared to (7.19) we have set ε to zero and we have N rather than $N + 1$ equations). Letting

$$\begin{aligned} \bar{\Phi} &= (\phi_1 \ \phi_2 \ \dots \ \phi_N)^T \\ \bar{\mathbf{a}} &= (a_1 \ a_2 \ \dots \ a_N)^T \end{aligned}$$

we have an N -dimensional nonlinear equation

$$\bar{\Phi}(\bar{\mathbf{a}}) = 0$$

which can be solved by the same Newton–Raphson method:

$$\bar{\mathbf{a}}_{n+1} = \bar{\mathbf{a}}_n - \alpha \cdot \left(\frac{\partial \bar{\Phi}}{\partial \bar{\mathbf{a}}}(\bar{\mathbf{a}}_0) \right)^{-1} \cdot \bar{\Phi}(\bar{\mathbf{a}}_0) \quad (7.24)$$

7.6 Cutoff optimization

We are now going to use Remez algorithm to build an approximation of the phase shifter defined by (7.6). Apparently $H_{-90}(s)$ defined by (7.6) is an allpass. We will retain the mentioned allpass property in the approximation. Using serial decomposition the allpass $H(s)$ can be decomposed into series of 2- and 1-pole allpasses. Since we aim to have $H(s)$ with as flat (actually, constant in the range of interest) phase response as possible, 2-poles seem to be less useful than 1-poles, due to the steeper phase responses of the former (Figs. 7.1 and 7.2).

Restricting ourselves to using just 1-poles we have:

$$H(s) = \prod_{n=1}^N A_n(s) = \prod_{n=1}^N \frac{\omega_n - s}{\omega_n + s} \quad (7.25)$$

where ω_n are the cutoffs of the 1-pole allpasses $A_n(s)$. Notice that the specific form of specifying $H(s)$ in (7.25) ensures $H(0) = 1 \ \forall N$, that is we wish to have a 0° rather than -180° phase response at $\omega = 0$.

Now the idea is the following. Suppose $N = 0$ in (7.25) (that is we have no 1-pole allpasses in the serial decomposition yet). Adding the first allpass A_1 at the cutoff ω_1 we make the phase response of (7.25) equal to the one of a 1-pole allpass (Fig. 7.1). From $\omega = 0$ to $\omega = \omega_n$ the phase response is kind of what we expect it to be: it starts at $\arg H(0) = 0$ and then decreases to $\arg H(j\omega_n) = -\pi/2$. However, after $\omega = \omega_n$ it continues to decrease, which is not what we want. Therefore we insert another allpass A_2 with a *negative cutoff* $-\omega_2$:

$$H(s) = \frac{\omega_1 - s}{\omega_1 + s} \cdot \frac{-\omega_2 - s}{-\omega_2 + s} \quad 0 < \omega_1 < \omega_2$$

Clearly, A_2 is unstable. However, we already know that unstable components of $H(s)$ are not a problem, since they simply go into the H_+^{-1} part of the phase splitter.

The phase response of a negative-cutoff allpass (Fig. 7.8) is the inversion of Fig. 7.1. Therefore, given sufficient distance between ω_1 and ω_2 , the phase response of H will first drop below $-\pi/2$ (shortly after $\omega = \omega_1$) and then at some point turn around and grow back again (Fig. 7.9). Then we insert another positive-cutoff allpass A_3 , then a negative-cutoff allpass A_4 etc., obtaining if not an equiripple approximation of -90° phase response, then something of a very similar nature (Fig. 7.10).

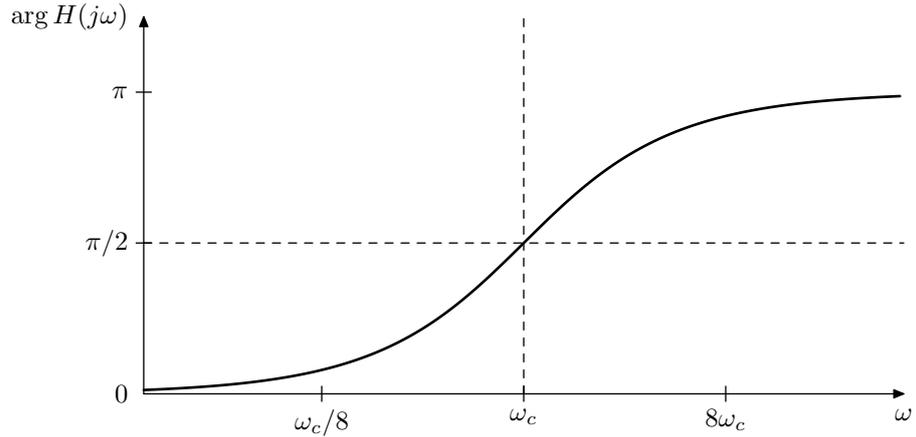


Figure 7.8: Phase response of a negative-cutoff 1-pole allpass filter.

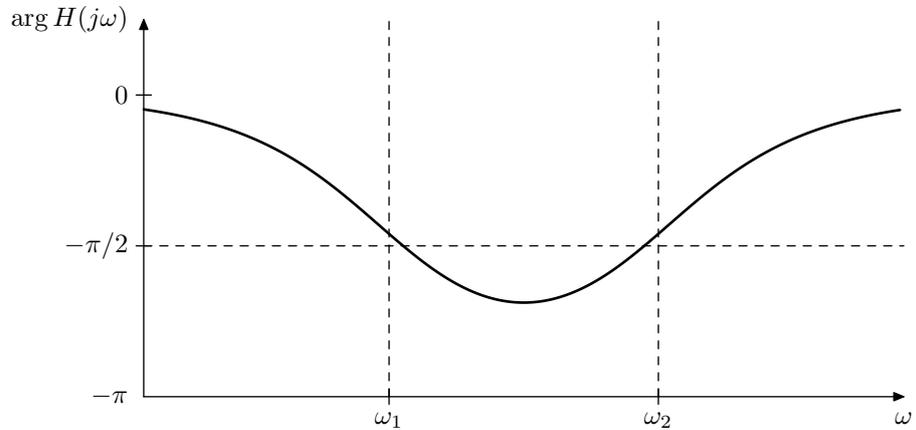


Figure 7.9: Phase response of a pair of a positive-cutoff and a negative-cutoff 1-pole allpass filters. Frequency scale is logarithmic.

The curve in Fig. 7.10 has two obvious problems. The ripple amplitude is way too large. Furthermore, in order to obtain this kind of curve, we need to position the cutoffs ω_n pretty wide apart (4 octaves between the neighboring cutoffs is a safe bet). We would like to position the cutoffs closer together, thereby reducing the ripple amplitude, however the uniform spacing of the cutoffs doesn't work very well for denser spacings of the cutoffs. We need to find a way to identify

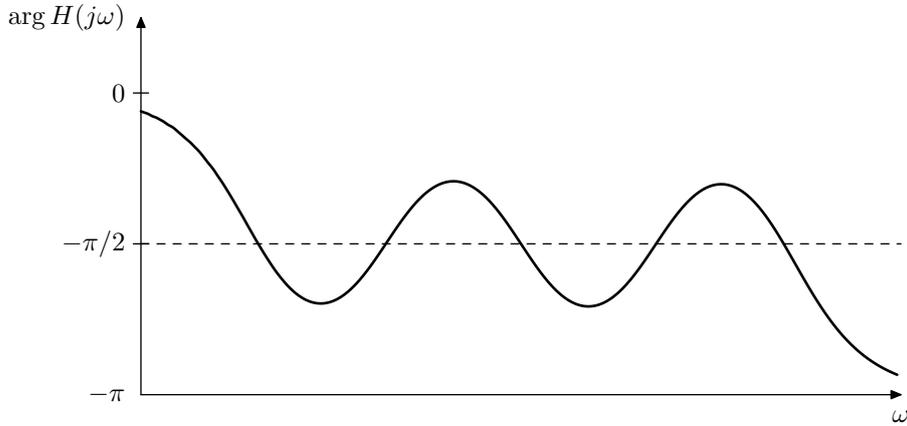


Figure 7.10: Phase response of a series of alternating positive-cutoff and negative-cutoff 1-pole allpass filters. Frequency scale is logarithmic.

the optimum cutoff positions.

Using cutoffs of alternating signs, we rewrite the transfer function expression (7.25) as

$$H(s) = \prod_{n=1}^N A_n(s) = \prod_{n=1}^N \frac{(-1)^{n+1} \omega_n - s}{(-1)^{n+1} \omega_n + s} \quad 0 < \omega_1 < \omega_2 < \dots < \omega_N \quad (7.26)$$

(the cutoff of A_1 needs to be positive in order for the phase response of H to have a negative derivative at $\omega = 0$). Considering that the phase response of a 1-pole allpass with cutoff ω_c is

$$H(j\omega) = -2 \arctan \frac{\omega}{\omega_c}$$

the phase response of the serial decomposition (7.26) is

$$\varphi(x) = \arg H(j\omega) = 2 \sum_{n=1}^N (-1)^n \arctan \frac{\omega}{\omega_n} = 2 \sum_{n=1}^N (-1)^n \arctan e^{x-a_n} \quad (7.27)$$

$$\begin{aligned} \omega &= e^x \\ \omega_n &= e^{a_n} \end{aligned}$$

where x and a_n are the logarithmic scale counterparts of ω and ω_n (essentially these are the pitch-scale values, we have just used e rather than 2 as the base to simplify the expressions of the derivatives of φ). The reason to use the logarithmic scale in (7.27) is that the phase responses of 1-pole allpasses are symmetric in the logarithmic scale, therefore the entire problem gets certain symmetry and uniformity.

Now we are in a position to specify the minimax approximation problem of construction of the phase shifter H_{-90} . We wish to find the minimax approximation of $f(x) \equiv -\pi/2$ on the specified interval $x \in [x_{\min}, x_{\max}]$, where the approximating function $\varphi(x)$ needs to be of the form (7.27).

The approximating function $\varphi(x)$ has N parameters:

$$\varphi(x) = \varphi(x, a_1, a_2, \dots, a_N)$$

which can be found by using the Remez algorithm for approximations of general form. Notably, for larger N and smaller intervals $[x_{\min}, x_{\max}]$ the problem becomes more and more nonlinear, requiring smaller damping factors α in (7.23) and (7.24). The damping factors may be chosen by restricting the lengths $|\mathbf{a}_{n+1} - \mathbf{a}_n|$ and $|\bar{\mathbf{a}}_{n+1} - \bar{\mathbf{a}}_n|$ in (7.23) and (7.24).

In order to further employ the logarithmic symmetry of the problem (although this is not a must), we may require $x_{\min} + x_{\max} = 0$ corresponding to $\omega_{\min}\omega_{\max} = 1$. Then the following applies.

- Due to the symmetry $\omega_{\min}\omega_{\max} = 1$ the obtained cutoffs ω_n will also be symmetric: $\omega_n\omega_{N+1-n} = 1$. (Actually they will be symmetric relatively to $\sqrt{\omega_{\min}\omega_{\max}}$ no matter what the ω_{\min} and ω_{\max} are, but it's convenient to have this symmetry more explicitly visible.)
- Using this symmetry the number of cutoff parameters can be halved (for odd N the middle cutoff $\omega_{(N+1)/2}$ is always at unity and therefore can be also excluded from the set of varying parameters). Essentially we simply restrict $\varphi(x)$ to be an odd (for odd N) or even (for even N) function of x .
- The obtained symmetric range $[\omega_{\min}, \omega_{\max}]$ can be scaled by an arbitrary constant A by scaling the allpass cutoffs by the same constant:

$$\begin{aligned} [\omega_{\min}, \omega_{\max}] &\leftarrow [A\omega_{\min}, A\omega_{\max}] \\ \omega_n &\leftarrow A\omega_n \end{aligned}$$

Figs. 7.11 and 7.12 contain example approximations of $H_{-90}(s)$ obtained by cutoff optimization (for the demonstration purposes, the approximation orders have been chosen relatively low, giving the phase ripple amplitude of an order of magnitude of 1°).

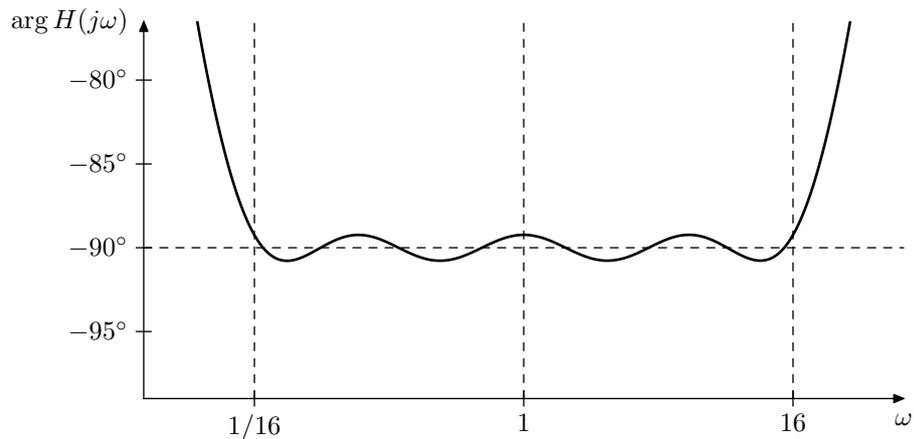


Figure 7.11: 8th-order minimax approximation of the ideal $H_{-90}(s)$.

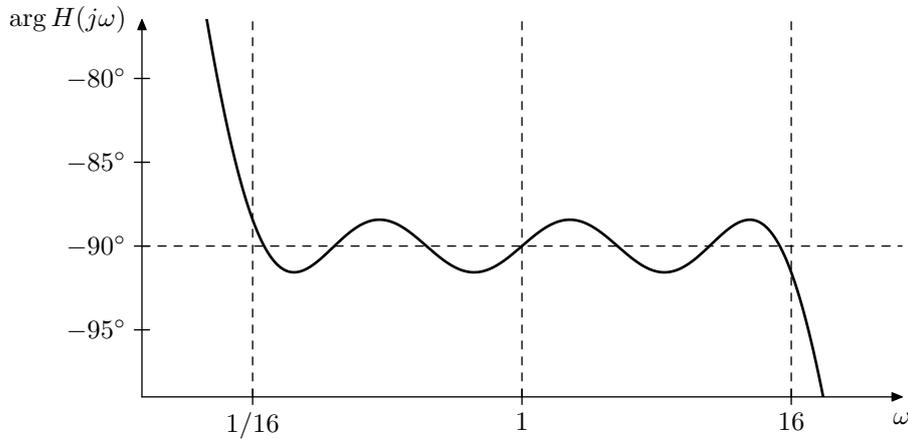


Figure 7.12: 7th-order minimax approximation of the ideal $H_{-90}(s)$.

Instead of solving the initial approximation equation (7.24) there is a different approach, which generally results in the nonlinearity of $\Phi(\mathbf{a})$ not so strongly affecting the algorithm convergence. We could take the manually constructed (7.26) with 4-octave spaced cutoffs $\omega_{n+1} = 16\omega_n$ as our initial approximation. The formal range of interest could contain two additional octaves on each side: $\omega_{\min} = \omega_1/4$, $\omega_{\max} = 4\omega_N$. Employing the logarithmic symmetry, we center the whole range around $\omega = 1$, so that $\omega_{\min}\omega_{\max} = 1$.

Using (7.23) (in the logarithmic scale x) we refine the initial approximation to the ripples of equal amplitude. Then we simply shrink the range a little bit. An efficient shrinking substitution is using the geometric averages:

$$\begin{aligned}\omega_{\min} &\leftarrow \sqrt{\omega_{\min}\omega_1} \\ \omega_{\max} &\leftarrow \sqrt{\omega_{\max}\omega_N}\end{aligned}\tag{7.28}$$

The substitution (7.28) doesn't affect the control points \hat{x}_n or the zeros \bar{x}_n of the Remez algorithm. Therefore after the substitution the Remez algorithm can be simply run again. Then the substitution is performed again, and so on, until we shrink the interval $[\omega_{\min}, \omega_{\max}]$ to the exact desired range.⁷

Notice that the approximations on the intermediate ranges $[\omega_{\min}, \omega_{\max}]$ do not need to be obtained with a very high precision, since their only purpose is to provide a starting point for the next application of the Remez algorithm on a smaller range. It is only the Remez algorithm on the exact desired range, which needs to be run to a high precision. This can noticeably improve the algorithm's running time.

7.7 Analytical construction of phase response

The cutoff optimization by Remez algorithm is a useful option if the filter coefficients are fixed and need to be obtained only once during the filter design

⁷Of course at the last step we simply set ω_{\min} and ω_{\max} to the desired values, rather than perform the substitution (7.28).

process. Should Remez algorithm fail to converge, it is possible to reduce the magnitude of the damping coefficient in the Newton–Raphson step, or make some other attempts to address the problem. However, if the coefficients need to be computed at runtime, it would be much more reliable if we had analytical expressions for the filter coefficients. This is what we are going to obtain next.

Since $H_{-90}(s)$ is an allpass, it can be written as

$$H_{-90}(j\omega) = e^{j\varphi}$$

where $\varphi = \varphi(\omega)$ is some (yet) unknown function of ω . Consider

$$e^{j\varphi} = \frac{1 + j \tan \frac{\varphi}{2}}{1 - j \tan \frac{\varphi}{2}}$$

We could attempt to approximate the desired $\tan(\varphi/2)$ by some rational function $F(\omega)$:

$$H_{-90}(j\omega) = e^{j\varphi} = \frac{1 + jF(\omega)}{1 - jF(\omega)} = \frac{j - F(\omega)}{j + F(\omega)} \quad (7.29)$$

Since phase responses of real filters must be real odd functions, so must be $F(\omega)$:

$$F(-\omega) = -F(\omega) \quad (7.30)$$

Given a real odd $F(\omega)$, the equation (7.29) will deliver a real $H_{-90}(s)$. Indeed, compare

$$H_{-90}(s) = \frac{j - F(-js)}{j + F(-js)} = \frac{j + F(js)}{j - F(js)} \quad (7.31)$$

(where the odd property of $F(\omega)$ has been used) and

$$\begin{aligned} H_{-90}(s^*) &= \frac{j + F(js^*)}{j - F(js^*)} = \frac{j + F(-j^*s^*)}{j - F(-j^*s^*)} = \frac{j - F((js)^*)}{j + F((js)^*)} = \\ &= \frac{-j^* - F^*(js)}{-j^* + F^*(js)} = \left(\frac{j + F(js)}{j - F(js)} \right)^* = H_{-90}^*(s) \end{aligned}$$

From (7.6) we obtain the requirement for an ideal $F(\omega)$:

$$F(\omega) \approx \begin{cases} -1 & \text{if } \omega > 0 \\ 1 & \text{if } \omega < 0 \end{cases} \quad (7.32)$$

Arctangent scale

When dealing with rational functions of real variable it is sometimes convenient and intuitive to plot (or at least imagine) the function graphs in the *arctangent scale*. That is, given the function $y = f(x)$, we use the variables

$$\begin{aligned} x' &= \arctan x \\ y' &= \arctan y \end{aligned}$$

to define the geometrical coordinates of the points on the graph image (in exactly the same way as we use $x' = \log x$ and $y' = \log y$ to create logarithmic scale plots).

Essentially, the arctangent scale is a visual representation of the “Riemann circle”, where the latter is the real counterpart of the complex Riemann sphere: in the same way as Riemann sphere represents the extended complex plane $\mathbb{C} \cup \{\infty\}$, the “Riemann circle” represents the extended real line $\mathbb{R} \cup \{\infty\}$. The arctangent scale is therefore “periodic”. The usage of the arctangent scale is illustrated by Figs. 7.13 and 7.14.

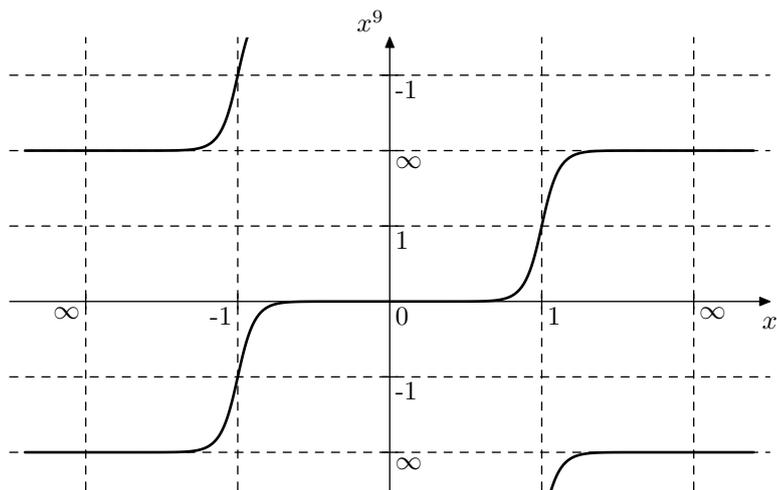


Figure 7.13: Power function x^9 in the arctangent scale.

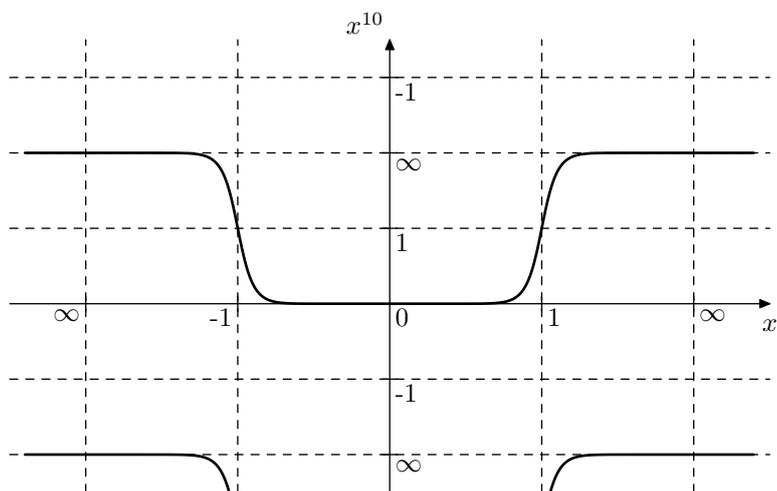


Figure 7.14: Power function x^{10} in the arctangent scale.

First-order rational transformations of the real line, when considered in the arctangent scale, look like (non-uniform) contractions/dilations, cyclic shifts and reflections of the line as well as combinations thereof. Particularly the following transformations have simple visual representations in the arctangent scale.

- Since

$$\tan\left(x' \pm \frac{\pi}{4}\right) = \frac{\tan x' \pm \tan \frac{\pi}{4}}{1 \mp \tan x' \tan \frac{\pi}{4}} = \frac{\tan x' \pm 1}{1 \mp \tan x'}$$

the transformations of the form

$$x \leftarrow \frac{x \pm 1}{1 \mp x}$$

are essentially just 45° cyclic shifts in the arctangent scale.

- Since

$$\tan\left(\frac{\pi}{2} - x\right) = 1/\tan x$$

the reciprocation of the real line

$$x \leftarrow 1/x$$

is a reflection relatively the $x' = \pi/4$ in the arctangent scale. Respectively, a pair of reciprocally symmetric values $x_1 x_2 = 1$ turns into a pair of values x'_1, x'_2 , which are symmetric relatively to $\pi/4$ (or $-\pi/4$, or any value of the form $\pi/4 + \pi n$). This explains the symmetry of the graphs of the power function in Figs. 7.13 and 7.18.

The benefit of this scale is that it treats the infinity as any other point on the line, and one can speak of points “before” or “after” the infinity. This is a convenient visual representation when dealing with rational function. E.g. one can visually observe the oscillations of elliptic rational functions not only around zero, but also around the infinity (Figs. 7.15, 7.18, 7.19 and 7.20).

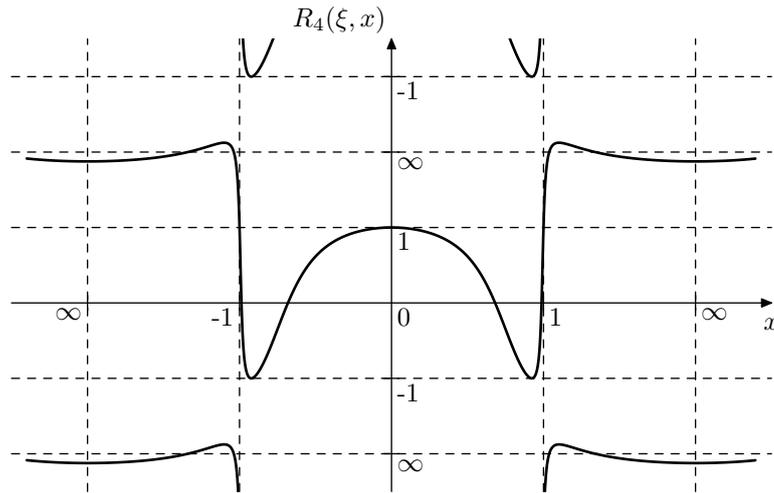


Figure 7.15: Elliptic rational function of the 4th order in the arctangent scale.

Phase responses based on power function x^N

Plotting the ideal $F(\omega)$ defined by (7.32) in the arctangent scale we obtain the graph in Fig. 7.16. The continuous graphs at the points of discontinuity ($\omega = 0$ and $\omega = \infty$) reflect the understanding of (7.32) as the limiting case of non-ideal $F(\omega)$. This is not the only possible way to draw a continuous graph of (7.32), but it will do for now. Notice that $F(0) = 0$ in Fig. 7.16, that is we wish the approximations of (7.32) to have zero phase response at $\omega = 0$.

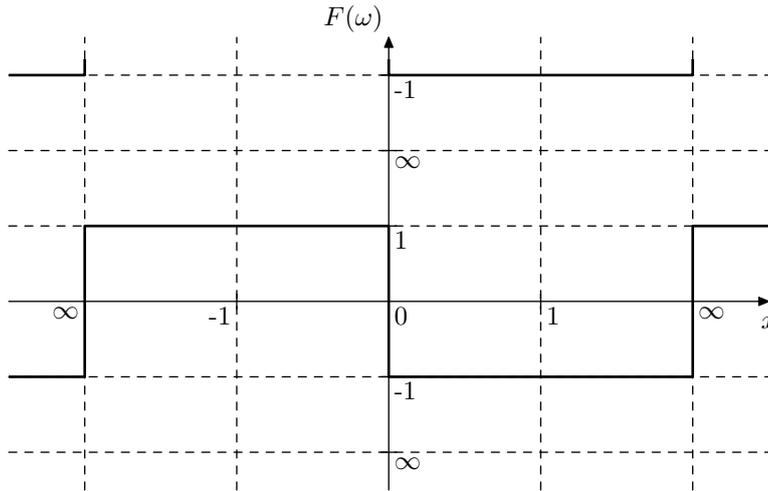


Figure 7.16: The ideal $F(\omega)$ in the arctangent scale.

Comparing the graphs in Fig. 7.16 and 7.14 we notice that they differ solely by 45° arctangent scaled shifts in both axes:

$$F = \frac{F' - 1}{F' + 1} \quad (7.33a)$$

$$\omega' = \frac{\omega - 1}{\omega + 1} \quad (7.33b)$$

where $F'(\omega') = \omega'^N$, where N is even. Notice that the requirements (7.32) therefore become

$$F'(\omega') \approx \begin{cases} 0 & \text{if } |\omega'| < 1 \\ \infty & \text{if } |\omega'| > 1 \end{cases} \quad (7.34)$$

which are fully satisfied by $F'(\omega') = \omega'^N$.

Being a 45° cyclic shift of the arctangent scale, the transformation (7.33) converts the odd symmetry requirement (7.30) into the reciprocal symmetry requirement

$$F'(1/\omega')F'(\omega') = 1 \quad (7.35)$$

since the latter in the arctangent scale is simply describing the reflection symmetry relative to $\omega' = 1$ and $F' = 1$. Apparently this requirement is also satisfied by $F'(\omega') = \omega'^N$. Thus, given $F'(\omega') = \omega'^N$ (where N is even) one could use (7.33) to directly construct $F(\omega)$, which in turn can be converted to $H_{-90}(s)$ by using (7.31).

The expression for $H_{-90}(s)$ is however not exactly what we want. What we want are the poles and zeros of $H_{-90}(s)$, so that we can represent it as a series of allpasses. According (7.29) the poles of $H_{-90}(s)$ are given by solving the equation

$$F(\omega) = -j \quad (7.36)$$

Using (7.33a), the equation (7.36) can be converted into the equivalent equation

$$F'(\omega') = -j \quad (7.37)$$

which (for $F'(\omega') = \omega'^N$) can be easily solved in terms of complex ω' . Having obtained the solutions of (7.37) in terms of ω' , we apply (7.33b) to convert ω' to ω , thereby obtaining the poles of $H_{-90}(s)$ in terms of the complex ω , which can be converted to the s -plane poles by $s = j\omega$.

The zeros of $H_{-90}(s)$ are respectively given by the equations

$$\begin{aligned} F(\omega) &= j \\ F'(\omega') &= j \end{aligned} \quad (7.38)$$

Since $F(\omega)$ is real, one doesn't need to solve (7.38), because the solutions of (7.38) and (7.36) must be mutually conjugate in the complex ω -plane. Respectively the poles and zeros of $H_{-90}(s)$ in the s -plane will be symmetric relatively to the imaginary axis, which is in agreement with the allpass property of $H_{-90}(s)$.

Half of the poles will be unstable, those will need to go into H_+^{-1} according to (7.9). A more interesting observation is that all poles will be real. Indeed, for $F'(\omega') = \omega'^N$ the equation (7.37) takes the form $\omega'^N = -j$. Its solutions ω' are lying on the unit circle in the complex ω' -plane. It's not difficult to see that (7.33b) transforms $|\omega'| = 1$ to $\text{Re } \omega = 0$ and respectively $\text{Im } s = 0$. This means that $H_{-90}(s)$ can be decomposed into a series of purely real 1-pole allpasses.

Odd-order power functions ω'^N also can be used. Instead of (7.16) consider the other way to continuously express (7.32) (while still retaining the property $F(0) = 0$), as shown in Fig. 7.17. Comparing Fig. 7.17 to Fig. 7.13 we realize that odd-order power function need to be simply negated before being subject to the 45° shifting (7.33). That is we let $F'(\omega') = -\omega'^N$ for N odd.

Phase responses based on elliptic rational functions

One could easily notice that the phase responses constructed using power functions are inferior (in the sense of larger deviation from the ideal phase response) to those obtained by the minimax optimization of allpass cutoffs.

In order to improve the phase response we could in principle perform a minimax approximation of (7.32). However, besides the mentioned earlier potential convergence issues which could be critical for the runtime application, there are two further problems associated with this approach.

- $F(\omega)$ is not the phase response $\varphi(\omega)$ itself, but its half-angle tangent $\tan(\varphi/2)$. Since the values of interest are $\varphi/2 \approx \pm\pi/4$, due to asymmetric nonlinearity of tangent in this range, the equiripple behavior of $\tan(\varphi/2)$ is not exactly the same as the equiripple behavior of φ . Although, for small ripple amplitudes this asymmetry could in principle be ignored.

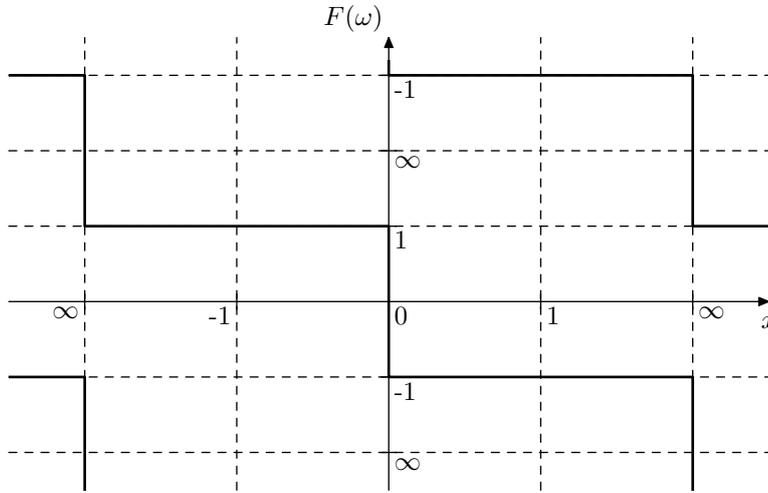


Figure 7.17: The ideal $F(\omega)$ in the arctangent scale (the other possibility).

Alternatively, one could rewrite the minimax equations (7.14) taking into account the asymmetry of the tangent around $\varphi/2 = \pm\pi/4$:

$$\tilde{f}(\hat{x}_n) \cdot (1 + \varepsilon)^{(-1)^n} = f(\hat{x}_n)$$

- After performing the minimax optimization in terms of the coefficients of the rational function $F(\omega)$, we will need to numerically obtain the solutions of the pole and zero equations (7.36) and (7.38). While this is in principle doable,⁸ wouldn't it be simpler just to run the minimax optimization of the cutoffs instead?

Fortunately, the minimax approximations of (7.32) do have analytical expressions via elliptic rational functions. Elliptic rational functions have the reciprocal symmetry

$$R_N(\xi, \xi/x)R_N(\xi, x) = L_N(\xi) \quad (7.39)$$

By scaling the elliptic rational function $R_N(\xi, x)$ by the square root of the selectivity factor ξ in the argument scale and by the square root of the discrimination factor $L_N(\xi)$ in the function's value scale:

$$F'(\omega') = L_N^{-1/2}(\xi)R_N(\xi, \xi^{1/2}\omega') \quad (7.40)$$

the reciprocal symmetry (7.39) is normalized into the symmetry (7.35). The functions (7.40) are analytical solutions of the minimax optimization of the norm

$$E = \max \left\{ \max_{|\omega'| \leq 1/\sqrt{\xi}} |F'(\omega')|, \max_{|\omega'| \geq \sqrt{\xi}} |1/F'(\omega')| \right\}$$

Respectively, the arctangent scale graphs of (7.40) look like the ones in Figs. 7.18, 7.19 and 7.20.

⁸The numerical search for the solutions of the equations (7.36) and (7.38) can be replaced by the numerical search of the poles and zeros of $H_{-90}(s)$, which in fact will be lying on the real axis. $H_{-90}(s)$ can be obtained from $F(\omega)$ by an explicit transformation of the coefficients.

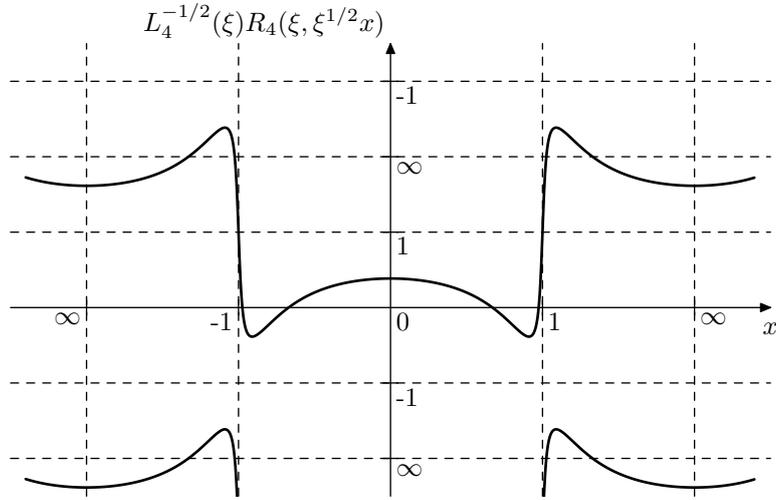


Figure 7.18: Elliptic rational function of the 4th order, scaled by $\sqrt{\xi}$ and $\sqrt{L_N(\xi)}$, in the arctangent scale.

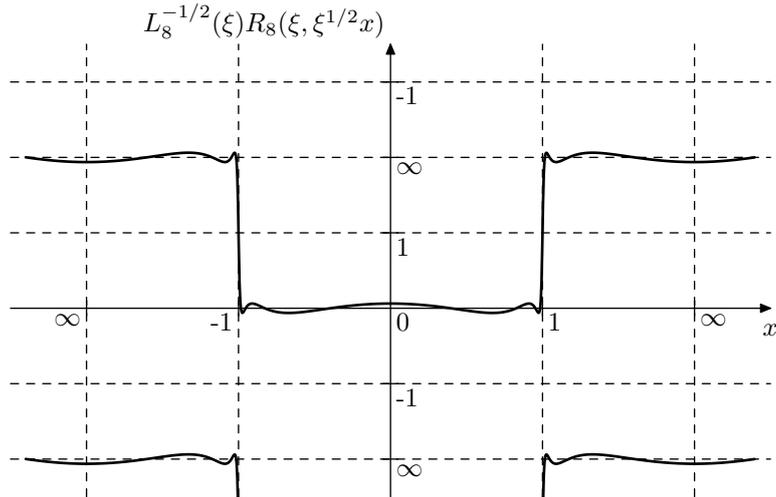


Figure 7.19: Elliptic rational function of the 8th order, scaled by $\sqrt{\xi}$ and $\sqrt{L_N(\xi)}$, in the arctangent scale.

Thus the functions (7.40) can be used with (7.33) in exactly the same way as the power functions ω'^N , where we have to remember to invert the function's value sign for odd N :

$$F'(\omega') = (-1)^N L_N^{-1/2}(\xi) R_N(\xi, \xi^{1/2} \omega') \quad (7.41)$$

The equations (7.36) and (7.38) can also be solved analytically in this case. The solutions of (7.36) are given by

$$\omega' = \frac{\text{cd } x + j \text{sn } x}{1 + j k \text{sn } x \text{cd } x} \quad (7.42)$$

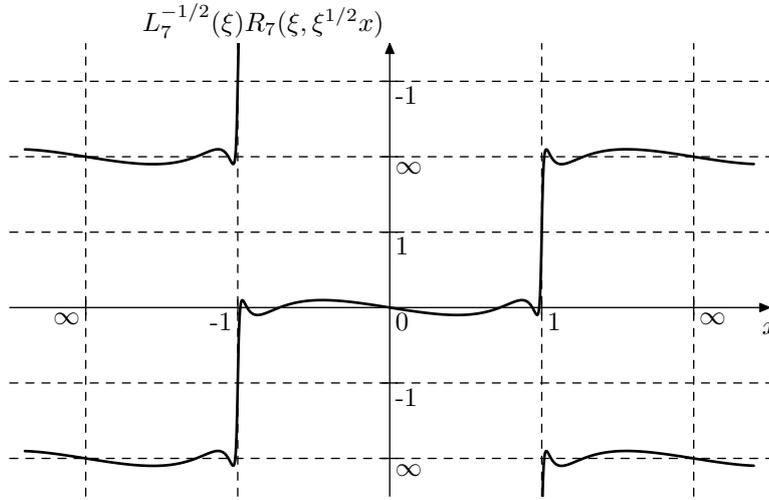


Figure 7.20: Elliptic rational function of the 7th order, scaled by $\sqrt{\xi}$ and $\sqrt{L_N(\xi)}$, in the arctangent scale.

where the elliptic modulus k is the reciprocal of the selectivity factor

$$k = 1/\xi$$

and

$$x = \frac{4n + 2 + (-1)^N K(k)}{N} \quad n = 0, 1, 2, \dots, N - 1$$

The solutions of (7.38) are the complex conjugates of (7.42). The formula (7.42) implies $|\omega'| = 1$, that is the solutions ω' of the pole equation are lying on the unit circle and respectively the poles of $H_{-90}(s)$ are real and are given by

$$s = \frac{(1 - k \operatorname{cd}^2 x) \operatorname{sn} x}{(1 + k \operatorname{sn}^2 x) \operatorname{cd} x - (1 + k^2 \operatorname{sn}^2 x \operatorname{cd}^2 x)} \quad (7.43)$$

where k and x are the same as in (7.42) and where the stable poles are given by $n < N/2$ for even N and $n < (N + 1)/2$ for odd N . The allpass transfer functions $H_{-90}(s)$ obtained from (7.43) are identical to the ones obtained by the cutoff optimization method (Figs. 7.11 and 7.12).

Bandwidth control

When using the power functions $(-\omega')^N$ as $F'(\omega')$, there are no parameters to play with, except the order N of the function. However with the elliptic rational function (7.41) there is the selectivity factor ξ , which can be freely chosen. The selectivity factor affects the width of the transition regions of the elliptic rational function. These transition regions will be mapped to the transition bands of the allpass $H_{-90}(s)$, where the phase response $\varphi(\omega)$ is changing from $+90^\circ$ to -90° (around $\omega = 0$) or from -90° to $+90^\circ$ (around $\omega = \infty$). Therefore, the selectivity factor controls the width of the transition bands of $H_{-90}(s)$. Apparently, the widths of the transition bands and the widths of the “passbands”

(the bands where the phase response is $\pm 90^\circ$) are complementary, therefore the selectivity factor also controls the passband width of $H_{-90}(s)$.

Noticing that

$$\begin{aligned}\varphi = -\frac{\pi}{2} &\iff F' = 0 \\ \varphi = \frac{\pi}{2} &\iff F' = \infty\end{aligned}$$

we identify the equiripple “regions of interest” of $F'(\omega')$:

$$\varphi \approx -\frac{\pi}{2} \iff F' \approx 0 \iff |\omega'| \leq \xi^{-1/2} \iff |F'| \leq L_N^{-1/2}(\xi) \quad (7.44a)$$

$$\varphi \approx \frac{\pi}{2} \iff F' \approx \infty \iff |\omega'| \geq \xi^{1/2} \iff |F'| \geq L_N^{1/2}(\xi) \quad (7.44b)$$

Concentrating on the range where $\varphi \approx -\pi/2$ (the other range is fully symmetric to the first one anyway) and using (7.33b), we have

$$|\omega'| \leq \xi^{-1/2} \iff \frac{1 - \xi^{-1/2}}{1 + \xi^{-1/2}} \leq \omega \leq \frac{1 + \xi^{-1/2}}{1 - \xi^{-1/2}}$$

Taking the natural logarithm, we have

$$|\ln \omega| < \ln \frac{1 + \xi^{-1/2}}{1 - \xi^{-1/2}} = 2 \tanh^{-1}(\xi^{-1/2})$$

or

$$|\log_2 \omega| \leq 2 \frac{\tanh^{-1}(\xi^{-1/2})}{\ln 2} = \frac{\Delta}{2}$$

where \tanh^{-1} is the inverse hyperbolic tangent and Δ is the octave bandwidth of the allpass $H_{-90}(s)$. Therefore, given the bandwidth Δ , the selectivity factor is defined by

$$\xi = \frac{1}{\tanh^2\left(\frac{\Delta}{4} \ln 2\right)} \quad (7.45)$$

The ripple amplitude is obtained by applying (7.33a) to (7.44a). Apparently, (7.33a) implies

$$F' = \tan\left(\frac{\varphi}{2} + \frac{\pi}{4}\right) = \tan\left(\frac{\varphi + \pi/2}{2}\right)$$

(which is in agreement with the fact that (7.33a) is a cyclic shift by 45° in the arctangent scale). Then, from (7.44a),

$$\begin{aligned}|F'| \leq L_N^{-1/2}(\xi) &\iff \left| \tan\left(\frac{\varphi + \pi/2}{2}\right) \right| \leq L_N^{-1/2}(\xi) \iff \\ &\iff \left| \varphi + \frac{\pi}{2} \right| \leq \arctan L_N^{-1/2}(\xi)\end{aligned}$$

Thus the ripple amplitude is

$$|\Delta\varphi|_{\max} = \arctan L_N^{-1/2}(\xi)$$

7.8 “LP to analytic” substitution

For the sake of a theoretical exercise (which is going to have practical implications) let's convert $H_{-90}(s)$ obtained from (7.29) into the respective analytic filter $H_{>0}(s)$. Substituting (7.33a) into (7.29) we obtain

$$H_{-90}(s) = \frac{j - F}{j + F} = \frac{j - \frac{F' - 1}{F' + 1}}{j + \frac{F' - 1}{F' + 1}} = j \frac{F' - j}{F' + j}$$

Substituting this into (7.7) yields

$$2H_{>0}(s) = 1 + jH_{-90}(s) = 1 - \frac{F' - j}{F' + j} = 2 \frac{j}{F' + j}$$

Or, using the stable version (7.8) of (7.7),

$$\begin{aligned} 2H_{>0}(s)H_+^{-1}(s) &= H_+^{-1}(s) + jH_+^{-1}(s)H_{-90}(s) = \\ &= H_+^{-1}(s) + jH_-(s) = 2H_+^{-1}(s) \frac{j}{F' + j} \end{aligned} \quad (7.46)$$

where $H_+(s)$ can be obtained as the unstable allpass component of $H_{>0}(s)$. The unstable poles arising out of the right-semiplane solutions of $F' + j = 0$ will be therefore cancelled by the zeros of $H_+^{-1}(s)$.

Explicitly writing the argument ω' of F' in (7.46) we have

$$2H_{>0}(j\omega)H_+^{-1}(j\omega) = 2H_+^{-1}(j\omega) \frac{j}{F'(\omega') + j}$$

Further rewriting the entire expression in terms of ω' we have

$$2H'_{>0}(j\omega')H_+^{-1}(j\omega') = 2H_+^{-1}(j\omega') \frac{j}{F'(\omega') + j}$$

where $H'_{>0}(s)$ and $H_+^{-1}(s)$ are obtained from $H_{>0}(s)$ and $H_+(s)$ according to (7.33b). Computing the squared amplitude response of $H'_{>0}H_+^{-1}$:

$$|H'_{>0}(j\omega')H_+^{-1}(j\omega')|^2 = |H_+^{-1}(j\omega')|^2 \cdot \left| \frac{j}{F'(\omega') + j} \right|^2 = \frac{1}{1 + F'^2(\omega')} \quad (7.47)$$

we notice the following.

- For $F'(\omega') = (-\omega')^N$ the equation (7.47) defines the amplitude response of a unit-cutoff Butterworth filter. The poles of such filter are lying on the unit circle in the s -plane.
- For the elliptic rational function (7.41) the equation (7.47) defines the amplitude response of an elliptic minimum Q-factor (EMQF) filter. The poles of such filter are also lying on the unit circle in the s -plane.⁹

⁹EMQF filter is simply an elliptic filter where the gain of the elliptic rational function has been set to the reciprocal of the square root of the discrimination factor. Depending on the definition of the EMQF filter's cutoff, its poles may be lying on a circle of some other radius. In this case the absolute magnitudes of the poles can be simply normalized, obtaining a EMQF filter with poles on the unit circle.

Actually, it’s not just that the amplitude response of $H'_{>0}H'^{-1}_+$ is the one of a Butterworth or an EMQF filter. In fact, $H'_{>0}H'^{-1}_+$ is a Butterworth or an EMQF filter. Indeed, the poles of $F'(\omega')$ (if any) are corresponding to the zeros of $H'_{>0}H'^{-1}_+$. The poles of $H'_{>0}H'^{-1}_+$ are the poles of $H'_{>0}$ and H'^{-1}_+ , where the unstable poles of $H'_{>0}$ are cancelled by the zeros of H'^{-1}_+ . The poles of $H'_{>0}$ are obtained from the equation

$$F'(\omega') + j = 0 \quad (7.48a)$$

(which is a rewritten (7.37)). Since the zeros of H'^{-1}_+ coincide with the unstable poles of $H'_{>0}$, they must be obtained from the same equation. Now, since H'^{-1}_+ is an allpass, its zeros are symmetric to its poles relatively to the imaginary axis in the s' -plane (where $s' = j\omega'$). Respectively, in the complex ω' plane they are conjugate-symmetric and thus the poles of H'^{-1}_+ are obtained from the equation

$$F'(\omega') - j = 0 \quad (7.48b)$$

Thus, the stable poles of $H'_{>0}H'^{-1}_+$ are combined from the stable poles of $H'_{>0}$, which are the stable solutions of the equation (7.48a) and the stable poles of H'^{-1}_+ , which are the stable solutions of the equation (7.48b). That is they are simply the odd and even poles of the Butterworth or the EMQF filter.¹⁰

So, the zeros and the stable poles of $H'_{>0}H'^{-1}_+$ are exactly those of a Butterworth or EMQF filter. According to (7.47) their squared amplitude responses are also equal. Any potential remaining difference between their transfer functions can be only by a constant allpass gain factor $|g| = 1$, which doesn’t matter, since our phase splitter anyway performs an allpass transformation of the input signal by H_+ .

The relationship between $H_{>0}(s)H^{-1}_+(s)$ and a Butterworth or EMQF filter creates yet another possibility to construct a phase splitter.¹¹ We could start off with a Butterworth or EMQF filter $H'_{>0}H'^{-1}_+$ (with a properly normalized cutoff, so that the poles are lying on the unit circle) and apply the relationship (7.33b) to convert it into an analytic filter $2H_{>0}H^{-1}_+$. Using the relationship

$$2H_{>0}(s)H^{-1}_+(s) = H^{-1}_+(s) + jH_-(s)$$

the analytic filter is then decomposed into the real and imaginary allpass filters H^{-1}_+ and H_- .

Notice that in principle, to do this decomposition we need to keep track of the odd and even poles of the original Butterworth or EMQF filter (so that we know which poles of the obtained analytic filter go into H^{-1}_+ and which into H_-). However there will be a simple rule to identify which is which.

Now we discuss the same steps in a little bit more detail. Applying $s = j\omega$ and $s' = j\omega'$ to (7.33b) we obtain the relationship between s and s' :

$$-js' = \frac{-js - 1}{-js + 1} \quad s = j\frac{1 - js'}{1 + js'} \quad (7.49)$$

¹⁰Notice that by multiplying the equations (7.48a) and (7.48b) we obtain the common pole equation

$$F'^2(\omega') + 1 = 0$$

where the odd and even poles are mixed together.

¹¹In fact, this is *the* classical method to construct phase splitters.

By cyclically shifting the ω axis according to (7.33b), the substitution (7.49) converts an s' -plane unit-cutoff lowpass filter into an s -plane analytic filter and therefore can be referred to as the "*LP to analytic*" substitution. This substitution particularly has the following properties.

- The substitution preserves the order of the filter.
- The substitution maps the imaginary axis onto itself according to (7.33b).
- The substitution maps each of the left and the right semiplanes back onto itself, thereby preserving the stability of the filter (this is easiest seen from (7.33b) by considering complex ω and ω' , where changing from one semiplane to the other corresponds to complex conjugation of ω and ω').
- Since the substitution (7.33b) applies a 45° cyclic arctangent scale shift to the imaginary axis, it doesn't preserve the Hermitian property of the frequency response. Thus it will turn real filters into complex ones. Conversely, some of the real filters produced by the substitution may have complex originals.
- The conjugate pairs are mapped to reciprocal conjugates, in both directions:

$$\begin{aligned} s_1 = s_2^* &\iff s'_1 \cdot s'^*_2 = 1 \\ s'_1 = s'^*_2 &\iff s_1 \cdot s_2^* = 1 \end{aligned}$$

Notice that $s_1 \cdot s_2^* = 1 \iff |s_1| \cdot |s_2| = 1 \wedge \arg s_1 = \arg s_2$.

- The points on the unit circle satisfy $s \cdot s^* = 1$, therefore their image points satisfy $s' = s'^*$, that is the unit circle is mapped to the real axis (in both directions). More specifically

$$s' = -je^{j\alpha} \iff s = \tan \frac{\alpha}{2} \quad (7.50)$$

that is the range $s' = -j \dots -1 \dots +j$ of the unit circle on the s' -plane is mapped to the range $s = 0 \dots -1 \dots -\infty$ of the negative real semiaxis of the s -plane.

- Mutually conjugate points on the unit circle are mapped to reciprocally-symmetric points on the real axis (in both directions).

Thus, given a Butterworth or an EMQF filter $H'_{>0}H'^{-1}_+$ we can apply the transformation (7.49) to its poles. Since the poles for each of these two filter types are located on the stable half of the unit circle, they will be transformed to the poles on the negative real axis according to (7.50). The interleaving of the odd and even poles will be therefore preserved due to the monotonicity of the transformation (7.50). Then the odd stable poles are used to construct $H_-(s)$ while the even stable poles are used to construct $H_+^{-1}(s)$, thereby obtaining the phase splitter.

There is a very simple rule to identify, which poles are even and which are odd. Assuming $H_{-90}(0) = 1$ (which can be achieved by ensuring $H_-(0) = H_+^{-1}(0) = 1$) the phase response of $H_{-90} = H_-/H_+^{-1} = H_-H_+$ should have a negative derivative at $\omega = 0$, which implies that the pole which is closest to

the origin $\omega = 0$ should belong to H_- . The second closest to the origin pole therefore belongs to H_+^{-1} , the third one to H_- etc.¹²

7.9 Cutoff prewarping

In constructing the discrete-time version of $H_{-90}(s)$ implemented as a series of 1-pole allpasses, the cutoff prewarping is subject to the considerations discussed in section 5.6. That is, the theoretically correct way would be to prewarp a single chosen frequency ω_c :

$$\omega_{ca} = \frac{2}{T} \tan \frac{\omega_{cd}T}{2} \quad (\omega_{cd} = \omega_c)$$

and obtain the analog 1-pole allpass cutoffs ω_n by multiplying ω_{ca} by the respective frequency ratios. This however implies that the width of the equiripple phase response band will be shrunk by the arctangent function.

Indeed, the bandwidth Δ which we specify during the design of the $\pm 90^\circ$ phase shifter is the analog bandwidth:

$$\Delta = \log_2 \frac{\omega_{\max}}{\omega_{\min}} = \log_2 \frac{\omega_{\max,a}}{\omega_{\min,a}} = \Delta_a$$

The respective digital bandwidth is

$$\Delta_d = \log_2 \frac{\omega_{\max,d}}{\omega_{\min,d}} = \log_2 \frac{\frac{2}{T} \arctan \frac{\omega_{\max,a}T}{2}}{\frac{2}{T} \arctan \frac{\omega_{\min,a}T}{2}} = \log_2 \frac{\arctan \frac{\omega_{\max,a}T}{2}}{\arctan \frac{\omega_{\min,a}T}{2}}$$

Clearly, $\Delta_d < \Delta_a$. If the upper bound of the equiripple band happens to be close to Nyquist frequency, the bandwidth reduction will be quite noticeable.

One way around this is to prewarp each of the cutoffs ω_n independently, which should almost completely eliminate the effect of bandwidth reduction. This will also destroy the equiripple minimax property of the phase response of $H_{-90}(s)$, but the resulting non-optimality of the phase response may be tolerable.

The other option is to define Δ_a based on the specified discrete-time equiripple band $[\omega_{\min,d}, \omega_{\max,d}]$. Given the sampling period T we can apply (3.7) to obtain the respective analog frequency equiripple band $[\omega_{\min,a}, \omega_{\max,a}]$. Then we can compute the analog 1-pole cutoffs by minimax optimization on that range. Alternatively we use the obtained analog equiripple band $[\omega_{\min,a}, \omega_{\max,a}]$ to define the analog bandwidth $\Delta_a = \log_2(\omega_{\max,a}/\omega_{\min,a})$, determine the necessary selectivity factor ξ from (7.45) and obtain the cutoffs analytically using (7.43). Apparently, for the same digital bandwidth, the phase response ripple amplitude will be larger at lower sampling rates, since the respective analog bandwidth will be larger.

¹²In principle, the application of the ‘‘LP to analytic’’ substitution approach to construct phase splitters is not restricted to Butterworth and EMQF filters. Other lowpass filters can be used as the prototypes. However, these lowpass filters need to satisfy a number of restrictions. Essentially, we need that $F'(\omega')$ satisfies the reciprocal symmetry property (7.35).

SUMMARY

A frequency shifter can be built by multiplying an analytic signal by a complex sinusoid. Technically this implies the usage of a 90° phase splitter, whose output signals are multiplied by phase-locked sine and cosine signals and then are subtracted or added together. A 90° phase splitter can be built as a series of 1-pole allpasses, whose cutoff coefficients can be found numerically by min-max optimization or obtained analytically from Butterworth or EMQF filters (or simply directly from power and elliptic rational functions).

Further reading

S.J.Orfanidis, *Lecture notes on elliptic filter design* (available on the author's webpage).

M.Kleehammer, *Mathematical development of the elliptic filter* (available in QSpace online repository).

Elliptic filter (Wikipedia article).

L.M.Milne-Thomson, *Jacobian elliptic functions and theta functions* (in *Handbook of mathematical functions* by M.Abramowitz and I.A.Stegun, available on the internet).

History

The revision numbering is major.minor.bugfix. Pure bugfix updates are not listed here.

1.0.2 (May 18, 2012)

first public revision

1.1.0 (June 7, 2015)

- TSK filters
- frequency shifters
- further minor changes

Index

- 1-pole filter, 7
- 2-pole filter, 81
- 4-pole filter, 61

- allpass filter, 24, 90
- allpass substitution, 54
- amplitude response, 13, 35
- analytic signal, 115
- arctangent scale, 131

- bandpass filter, 81, 86
- bilinear transform, 40
 - inverse, 42
 - topology-preserving, 52
 - unstable, 59
- BLT, 40
- BLT integrator, *see* trapezoidal integrator

- canonical form, 50
- cheap TPT method, 72
- complex exponential, 5
- complex impedances, 12
- complex sinusoid, 1
- cutoff, 8, 14
 - parametrization of, 15

- damping
 - in SVF, 84
- DC offset, 2
- delayless feedback, 44
- DF1, 50
- DF2, 50
- differentiator, 54
- diode ladder filter, 73
- Dirac delta, 4
- direct form, 50

- eigenfunction, 9
- elliptic rational function, 119

- filter
 - 1-pole, 7
 - 2-pole, 81
 - 4-pole, 61
 - allpass, 24, 90
 - bandpass, 81, 86
 - highpass, 17, 64, 81
 - ladder, 61
 - lowpass, 7, 61, 81
 - multimode, 20, 65, 81
 - notch, 89
 - peaking, 90
 - Sallen–Key, 98
 - shelving, 21, 86
 - stable, 18
 - TSK, 98

- flanger, 110
- Fourier integral, 3
- Fourier series, 2
- Fourier transform, 3
- frequency response, 13, 35
- frequency shifter, 113

- gain element, 8

- harmonics, 2
- Hermitian, 3
- highpass filter, 17, 64, 81
- Hilbert transform, 115
- Hilbert transform pair, 115

- instantaneous gain, 46
- instantaneous offset, 46
- instantaneous response, 46
- instantaneously unstable
 - feedback, 57
- integrator, 8
 - BLT, *see* integrator, trapezoidal
 - naive, 31
 - trapezoidal, 37

- ladder filter, 61
 - diode, 73

- Laplace integral, 5
- Laplace transform, 5
- linearity, 11
- lowpass filter, 7, 14, 61, 81
 - “LP to analytic” substitution, 140
- LP to BP substitution, 91
- LP to BS substitution, 93
- LP to HP substitution, 19

- minimax approximation, 119
- multimode filter, 20, 65, 81

- naive integrator, 31
- nonstrictly proper, 11
- notch filter, 89

- partials, 2
- peaking filter, 90
- phase response, 13, 35
- Phase splitter, 115
- phaser, 107
- pole, 18
- prewarping, 43

- Remez algorithm, 119
- rolloff, 14

- Sallen–Key filter, 98
- shelving filter, 21, 86
- stable filter, 18
- state-variable filter, 81
- substitution
 - LP to analytic, 140
 - LP to BP, 91
 - LP to BS, 93
 - LP to HP, 19
- sumimator, 8
- SVF, 81

- time-invariant, 10
- topology, 51
- topology-preserving transform, 42, 52
- TPBLT, 52
- TPT, 42, 52
 - cheap, 72
- transfer function, 11, 34
- transposition, 26
- trapezoidal integrator, 37
- TSK filter, 98

- unit delay, 32
- z -integral, 30
- z -transform, 30
- zero, 18
- zero-delay feedback, 45