

The Implementation of Recursive Digital Filters for High-Fidelity Audio*

JON DATTORRO

ENSONIQ Corporation, Malvern, PA 19355, USA

The problems are described which the practicing engineer encounters who unwittingly approaches the realization of IIR digital filters for the first time. It is assumed that suitable design programs are available to calculate the coefficients, and it is desired only to *implement* the filter. Elegant solutions are provided for some of the most intimidating problems typically encountered, which are 1) input scaling requirements, 2) truncation noise propagation and recirculation, and 3) accurate low critical frequency filtering. It is shown that the direct form I noncanonic topology is the best for use in digital filtering, and while 16/32-bit DSP chips such as the TMS32010 or the ADSP-2100 can be used in many applications of high-fidelity digital audio, they will *not* meet the most demanding requirements.

0 INTRODUCTION

0.1 Ground Rules

Let us assume at the outset that the theoretical design of digital filters in the z domain is a routine task. There are many books and papers written on this topic. Moorer [1] provides unity-gain, minimum-phase designs and is a superb reference for digital audio work. The low-pass bound on the transition region of 6 dB per octave per conjugate pole pair does not exist in Moorer's second-order parametric filter designs. We have tested and verified the C programs that he provides. As such, this paper will not deal with the derivation of floating-point filter coefficients. Rather, it is assumed that they are known and the only remaining task is to design the software (or hardware) that will implement the filters. It should be emphasized that the design of digital filters and the *implementation* of digital filters can be carried out as two distinct tasks. Most design procedures produce the coefficients for the direct form filter topologies. The coefficients for other topologies can usually be derived from these via suitable transformations. More specifically, a typical design program would produce one set of coefficients which are directly applicable to either of the two standard second-order sections (the direct forms, Figs. 1 and 2, or their transposes, Fig.

3) which are described in the frequency domain by the transfer

$$H(z) = \frac{a_{F0} + a_{F1}z^{-1} + a_{F2}z^{-2}}{1 - b_{F1}z^{-1} - b_{F2}z^{-2}} \quad (1)$$

For brevity in the text, we will often write equations like Eq. (1) as

$$H(z) = \frac{\sum a_{Fi}z^{-i}}{1 - \sum b_{Fi}z^{-i}} \quad (2)$$

where the subscripts F indicate that these coefficients are the floating-point values provided by some design procedure. Our implementations will generally operate with fixed-point coefficient values, and so this becomes one of the issues that bear discussion.

We will always prefer minimum-phase designs (all zeros within the unit circle) because 1) the total phase excursion is always less than 360° for a second-order section; long phase ("propagation" or "transport") delays are unacceptable; 2) all coefficients are bounded as in Pascal's triangle if the filter design is unity gain. We prefer unity-gain designs because it will be less apparent when they are kicked in.

We will purposely show our filter circuit topologies as having positive-signed coefficients at the multipliers (see Fig. 1) because this leads to positive accumulations in the resulting software. We call this a *positive to-*

* Manuscript received 1988 July 18; revised September 7, 1988.

pology. This topology is sometimes more convenient for the programmer; hence the negative denominator coefficients in Eq. (1). The implementer must be aware of the convention concerning the signs of the coefficients produced by some design procedure. If the filter blows up at the start, this is the first place to look.

This work was conceived within the dimensions of the digital signal processing (DSP) architecture of the entire TMS320 series, but the concepts and problem solutions are universal and can be readily applied on other DSP processors. The primary feature of the TMS320 architecture (besides being nearly a full featured microprocessor) is a fast-signed 16- by 16-bit multiplier having a 32-bit product which can be accumulated to 32 bits. We call this a 16/32-bit architecture. So our work will be applicable to any architecture having N -bit by N -bit multiplications with nominally $2N$ -bit products and $2N$ -bit (the same wordlength) accumulations. This is typical of the current wave of DSP chips. The TMS32010 and 20 are peculiar in that there is no provision for an unsigned multiply. This at first seems to impair double-precision coefficient performance, but a very simple method to get around this limitation will be discussed.

TMS32010 code is provided in Appendix 1, which comes from a commercial parametric filtering application. It is provided mainly as proof of principle to the concepts presented herein.

0.2 Audio

In order for a digital filter to be usable for high-quality audio, it must not degrade system performance such that the system specifications are violated. This means that the effects of the filter must be transparent to the user except for the concomitant changes in frequency and phase response. Therefore such things as the signal-to-noise ratio (SNR), the dynamic range, or the total harmonic distortion plus noise (THD + N) of the entire system should not be degraded. The digital filter itself, then, must have specifications that surpass those of the host system.

We would like there to be some advantage to the use of a digital filter in comparison to the corresponding analog filter, so we mention one advantage which is not often cited: in analog filter implementations, tran-

sient intermodulation distortion (TIMD) may be a problem. In corresponding digital filter implementations there is no physical analogy to the cause of TIMD (slew-rate limiting) [2], thus there can be no TIMD induced.

In audio work it is customary to find designs centered in the very low-frequency region. To a recording engineer there is a vast difference between a filter whose critical frequency is placed at 50 Hz and another placed at 60 Hz. Single-precision (16-bit) coefficient realizations are not satisfactory for the required control over low critical frequency. It is generally accepted now that 24-bit coefficients are adequate for high-quality digital audio work.

0.3 Truncation Noise

The first widely used digital filters conceived in hardware in the 1960s were intended for telecommunications applications such as the telephone touch-tone decoder. Leland B. Jackson, who is often called the "father of digital filters" for his pioneering work in the analyses of associated nonlinear phenomena, was employed by Bell Telephone Laboratories when he and James F. Kaiser presented a hardware scheme [3] for the implementation of such a decoder. The outstanding topological idiosyncrasy of their circuit implementations, from our present vantage point, was that all the signal paths, including the multiplier paths, were constrained to have the same wordlength. Therefore their first analyses of roundoff noise warranted by the truncation process were based on a constraint that need no longer exist. DSP chips suitable for audio routinely supply double-wordlength products and accumulations. We must still be concerned, however, with the truncation at the output *and* its effect upon the digital circuit when fed back.

In concept, if a digital filter accepts as input an integer, then an integer plus a fraction (a fixed-point number having a nonzero fractional part) is, in general, produced at the output. If the output were rounded off to the nearest integer and we consider only the effects of errors that occur *at* the output, then the output would be in error by at most 1/2 LSB (3 dB) in magnitude. This rounding (or truncation) process at the filter output is unlike the process that describes the quantization errors

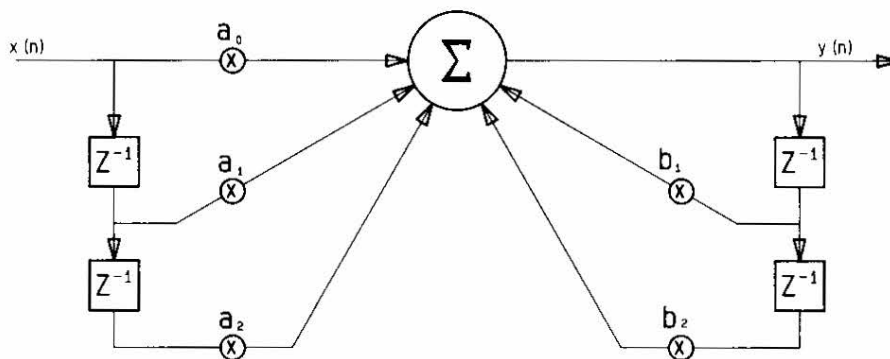


Fig. 1. Direct form I; for audio use.

inherent in analog-to-digital (A/D) conversion. Still, it is most often modeled in the literature [4, p. 413] as a random process that produces an additive white noise source at the output node. In any case, the input signals to cascaded second-order digital filters have, in addition to the signal itself, a noise source from the preceding stage which gets filtered along with the signal. We will be concerned with just how these noise sources combine through a cascade.

0.4 Filter Circuit Topology

0.4.1 Second Order

Since all of the required hardware to implement digital filters in real time has now been integrated onto a single programmable chip, there is a tendency to think of a digital filter as an amorphous software implementation as opposed to a hard circuitual implementation. (Indeed, the first digital filters were realized as models of continuous systems on general-purpose computers in high-level languages for the purpose of analyzing statistical data in nonreal time.) From the viewpoint of a programmer, a digital filter appears at first to be the solution to a difference equation, and circuit topology seems not too important. This cannot be further from the truth because the choice of circuit topology is as intrinsic to a given software realization as it is to the corresponding discrete hardware realization. One needs only to convince oneself that soft instructions solving a difference equation are, in fact, shuffling data about predetermined paths inside some architecture made up of arithmetic elements, that is, hardware. The circuit topology chosen determines the number and the order of the computations, and vice versa. Why choose a topology which might require more computation or more storage? The answer has to do with numerical inaccuracy. Certain topologies react better than others to numerical errors. So for these reasons we often choose *not* to use the most straight forward solution to the filter difference equation in our software. This is where the study of circuit topology comes in.

There are other considerations besides numerical inaccuracy which impact our choice of topology. One of the great disappointments in the past with conventional 16/32-bit digital filter implementation was the finding that filter inputs needed to be scaled (attenuated). Scal-

ing was required in order that overflow be avoided at internal nodes in a given filter topology, even though the output node itself may never have been threatened with overflow. The unfortunate side effect of scaling is that the SNR of the input audio is irrevocably worsened by an amount proportional to the scale factor. The same problem arises in analog designs. But good analog filter implementation can accommodate this reduction of input level because it can start off with as much as 120 dB SNR. For example, an analog signal with 90-dB SNR scaled by, say, 20 dB is ideally the same as an unscaled analog signal having 110-dB SNR. The analog filter can easily accommodate this SNR, and so analog scaling is generally not a problem. But the situation is different in digital. Using a 16-bit signal converter, we start out with a noise floor due to quantization in the vicinity of -90 dB, theoretically. With our 16/32-bit digital filter itself only having a 90-dB SNR, we really cannot ask it to represent a scaled signal now having a 110-dB SNR. We would end up with an output signal having an SNR of at most 70 dB.

A 24/48-bit architecture would ameliorate the scaling problem for 16-bit signals, in analogy to the 120-dB analog system mentioned. This is because using the feedback techniques presented herein, a 24/48-bit architecture digital filter can be made to operate with a 138-dB SNR. Input scaling a 16-bit signal by less than 48 dB would amount to a justification problem within the 24-bit wordlength, that is, no data loss. Even so, we are going to show how to sidestep the scaling problem altogether, regardless of the wordlength of the processor we are using, by clever choice of topology. For these reasons we venture to suggest that *any* 16/32-bit digital implementation employing input scaling is inadequate for high-fidelity digital audio.

0.4.2 High Order

Constructing higher order filters from second-order sections is done by cascading however many second-order filters are produced by a design procedure. A parallel realization would require what is called a partial fraction expansion, which transforms the given product of biquadratic transfers into sums of biquadratics. These procedures involve pole-zero pairing selection and the

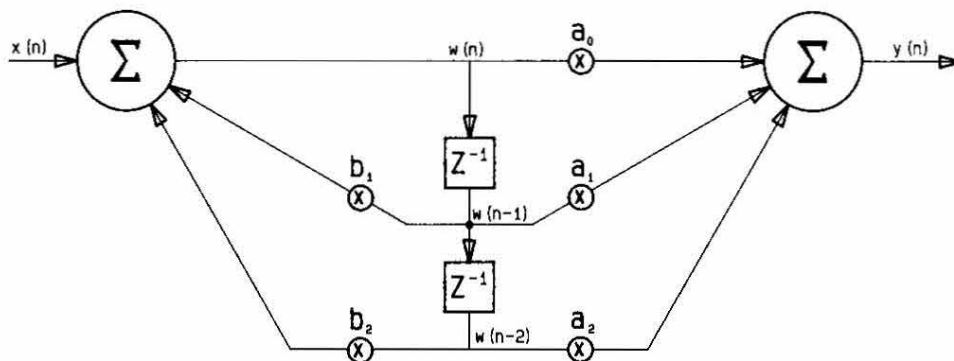


Fig. 2. Direct form II, canonic; not for audio use.

scaling of the input signal into individual sections for optimal SNR performance and the prevention of output overflow from individual sections. Since it is not possible, in general, to design high-order filters such that each second-order section is unity gain, input scaling is mandatory regardless of the section topology. If the pole-zero pairing is chosen for as many close-to-unity-gain sections as possible, and if the section topology is chosen to be direct form I for the reasons to be discussed shortly, then the scaling requirement can be minimized and limited to consideration of only the output node of each stage. (Jackson [5] gives closed-form scaling and pairing procedures.)

The construction of higher order filters from cascades or parallels of second-order sections (“biquads”) is the best approach because one soon discovers that the direct form implementation of high-order filters leads to geometric increases in the range of coefficient values, pole-zero sensitivity, and truncation noise recirculation. Based on the material presented in this paper, we will always want to choose the parallel and cascade forms I, as given by Jackson [5]. He notes that parallel realizations, while sometimes being less noisy than cascades, have very high zero sensitivity to coefficient quantization, and so cascades are often preferable. (Note that some of Jackson’s analyses still account for truncation before accumulation.) Graphic equalizers are traditionally realized as parallel second-order stages, while parametric equalizers are cascade implementations having individual second-order stages forming a complete filter.

The only direct high-order topology which is known to elude both the scaling and the truncation noise recirculation problem at once is the Gray–Markel all-pole four-multiplier ladder [6] (Fig. 4e). For simplicity’s sake, in this paper we will constrain all our filters to be biquadratic, and we want their magnitude responses to be deviations from unity gain. Therefore we will be able to show how to skirt the scaling issue totally. If we speak of a cascade of second-order sections, it will usually *not* be in reference to the construction of a high-order filter. Each second-order stage will itself be a complete filter.

Aside from this cursory discussion of higher order filters, all of the work in this paper will be limited to complete second-order digital filters, where all coefficient values are bounded in magnitude by 2 [Eq. (30)]. Later we will show examples of second-order filters having high gain and high Q factor. This comes about because of the proximity of the pole-zero pair in those designs.

We will now reintroduce a classic topology which inherently eludes the input scaling requirement, hence picking up some free SNR.

1 CIRCUIT TOPOLOGY: SCALING AND OVERFLOW

We will abruptly end our search for the best digital audio filter circuit topology with the most simple to-

poloogy, the direct form I noncanonic implementation of the second-order digital filter. This circuit is shown in Fig. 1. The second-order difference equation which describes the output $y(n)$ is given in Eq. (3) and is straightforward to derive from Fig. 1,

$$y(n) = \sum a_i x(n - i) + \sum b_i y(n - i) . \quad (3)$$

Had we implemented this second-order difference equation directly in software, Fig. 1 would represent the topological choice we had unwittingly made. The direct form II filter circuit topology shown in Fig. 2 is chosen, historically, in preference to direct form I in Fig. 1 because it has superior truncation noise performance when both are implemented with truncation *prior* to accumulation. Otherwise, the direct form I is better. The direct form II also has two less storage (delay) elements, which yields the minimum number of storage elements for the second-order topology. Hence the term *canonic topology* (for “minimal” topology) is often used synonymously. The equations that describe the second-order filter in Fig. 2 are

$$\begin{aligned} w(n) &= x(n) + \sum b_i w(n - i) \\ y(n) &= \sum a_i w(n - i) . \end{aligned} \quad (4)$$

These simultaneous equations can be simplified to Eq. (3). Implementing these equations in software represents the choice of the canonic topology. The two other standard forms, which are the circuitual transposes of the direct forms I and II, are shown in Fig. 3. There are many other topologies that embody the second-order section. The numerous topologies possess different characteristics with regard to the recirculation of truncation noise. These topologies can each be quantitatively characterized by SNRs of input signal to truncation noise power, and each topology is said to have a particular intrinsic “noise gain.” Agarwal and Burrus [7] published an excellent article on this topic, and the topic is collectively referred to as *roundoff noise* [8]. In this paper we substitute the term *truncation* since it is this and not rounding which will be performed in our filter programs (for good reason).

We choose direct form I because the truncation at the accumulator output is restricted to only the feedback paths; thus it will be easier to compensate. A casual scrutiny of this familiar topology reveals that there is only one accumulator, hence for truncation postaccumulation, one noise source. We choose the direct form I also because of its unique overflow properties. Notice in Fig. 1 that all the internal paths of this second-order section converge at this one accumulator. Further, the output of this mecca is attached directly to the output of the filter itself. Therefore by designing the filter for unity gain from input to output, we can ensure that the filter output $y(n)$ will not overflow under most real operating conditions. During the course of the calculations, that is, at intermediate steps in the accumulation of the final output, the accumulator is likely to overflow

many times. Jackson showed back in 1968 [3] that because 2's complement arithmetic is a modulo math, as long as the final result falls into the numerical range of the accumulator (the first modulo), intermediate accumulations may be allowed to overflow any number of times (traverse other moduli) without causing an erroneous accumulation. This means that the effective wordlength of the accumulator in the direct form I is much greater than its actual wordlength.

We have, therefore, infinite headroom at the internal nodes using the direct form I topology because all the internal nodes feed the one accumulator. The internal nodes which feed multipliers are not capable of overflow (by design), so there is no need to scale the input signal; hence no potential loss in the signal's SNR. Another

way to look at this is that we have an infinite headroom accumulator.

To take full advantage of modulo arithmetic we will never want the accumulator to autonomously operate in the saturation mode. If overflow is detected *at the final output*, then saturation must be performed there, but only under program control.

These advantages cannot be capitalized on using the canonic topology and, for that matter, in most other topologies because they have at least two accumulators. The overflowed output of one of them will invariably feed back from an internal node into a multiplier input, which can only operate on fixed wordlengths (Fig. 4). In order to prevent this from happening, we would somehow need to ensure that the transfer from the input

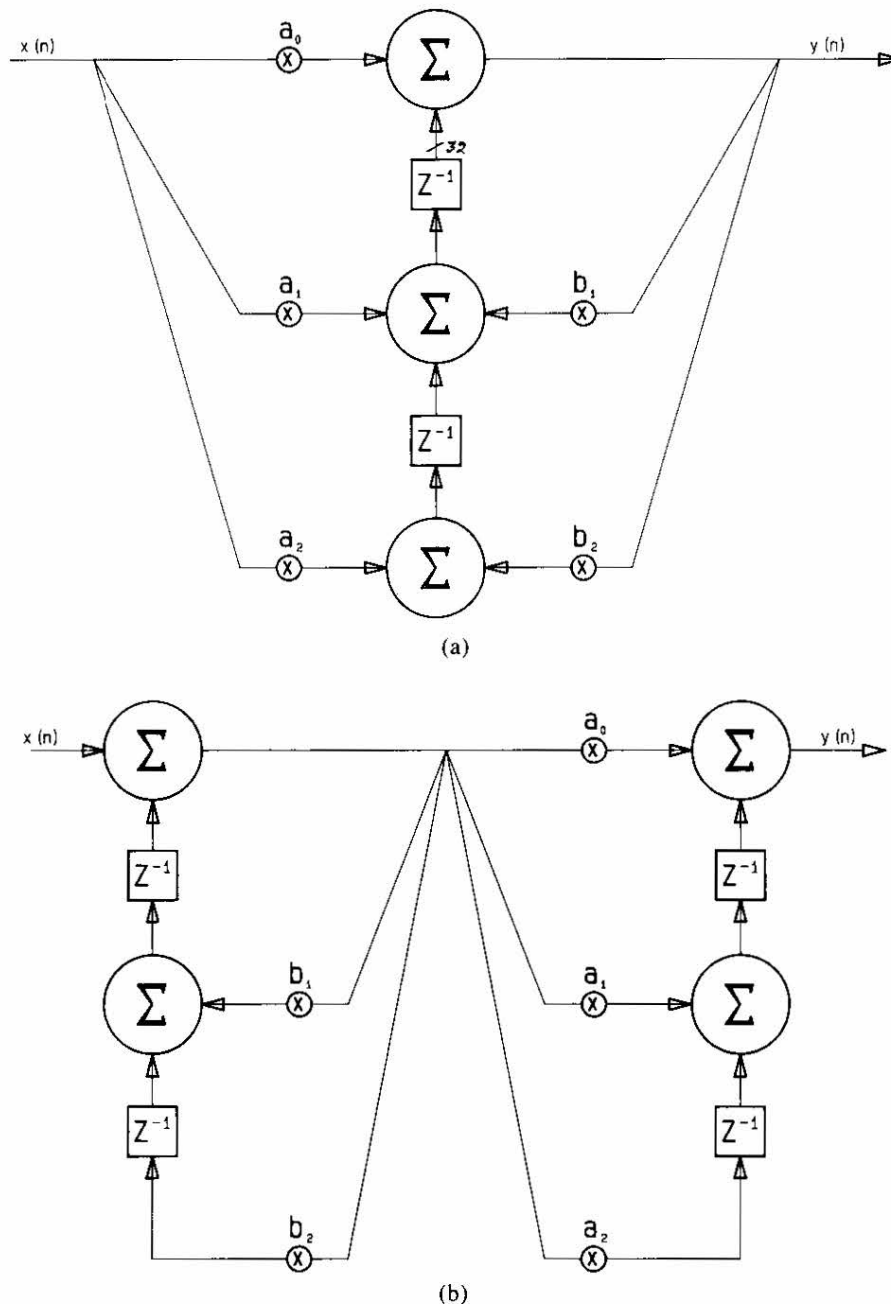


Fig. 3. (a) Direct form II transpose; for audio use. (b) Direct form I transpose; not for audio use.

to each of the sensitive internal nodes did not exceed unity, thus precluding accumulator overflow. [Notice that the *transpose* of the direct form II (canonic) topology shown in Fig. 3(a) is equally suitable for our

purposes and possesses all the desired characteristics, but we will only work with the direct form I in this paper.] Scaling the input with respect to the internal nodes is, therefore, mandatory for the canonic topology.

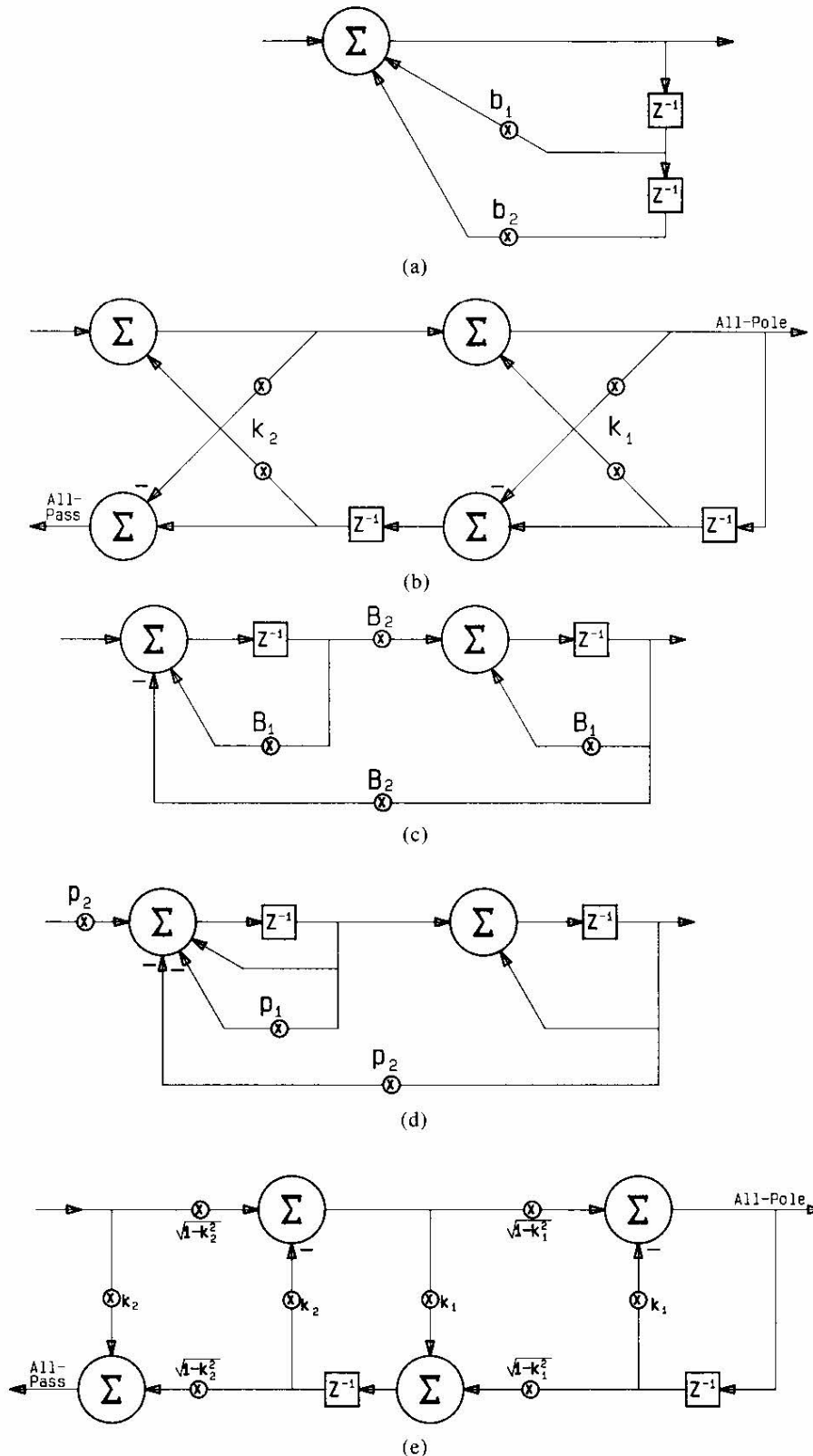


Fig. 4. Some second-order topologies. (a) Direct form I; has one accumulator. (b) Two-multiplier lattice; $k_1 - k_1k_2 = b_1$, $k_2 = b_2$. (c) Rader-Gold "coupled" or "normal" form in one guise. (d) Agarwal-Burrus. (e) Gray-Markel four-multiplier ladder.

So, we do not want to use it. Lower noise topologies such as the Agarwal–Burrus topology [7], [9] (which is a variation of the Rader–Gold topology) suffer from the same problem and thus require scaling.

1.1 Modulo Math

Two’s complement addition can be thought of in rings or moduli. The circumference of the ring is dictated by the greatest magnitude representable by a given wordlength. As an example, consider the sum of two positive numbers that exceeds the available wordlength, which is represented pictorially in Fig. 5. The two positive numbers can each be represented at a location within the first ring. For the sum, we can visualize the result as having gone into the second ring, like the minute hand of a clock after one hour has elapsed. If a third number is now added which is negative and of the proper magnitude, then the clock hand will come back somewhere into the first ring and the result will be correct, even though the intermediate result was in error.

1.2 Unity Gain

By the term *unity gain filter* we mean a filter that is designed as a deviation from unity. Several unity gain filters are shown in Fig. 6. Although we implied in the preceding discussion that designing for unity gain will ensure no output overflow, there are three qualifications to this reasoning which compel us to perform overflow detection at the output on a per-sample basis and to saturate the output on detection of overflow there.

1.2.1 Boost Filters

Unity gain design implies that we cannot design boost filters. We will dispel this idea now. To design a boost filter (having a gain that exceeds unity) is fine as long as input signals falling in the frequency region of the boost are such that they do not exceed unity *after* boosting. If the request for a boost is because of a paucity of energy in the corresponding frequency region, then the threat of overdrive is small. So just as with analog filters, it is up to the user who requests a boost to ensure that the inputs are not overdriven. If this cannot be done conveniently, then the filter could be redesigned by dropping the absolute level of the transfer function. This would be equivalent to scaling the input by some desired amount. The point here is that there need be no inherent flaw in the filter implementation itself which demands the use of scaling at all times.

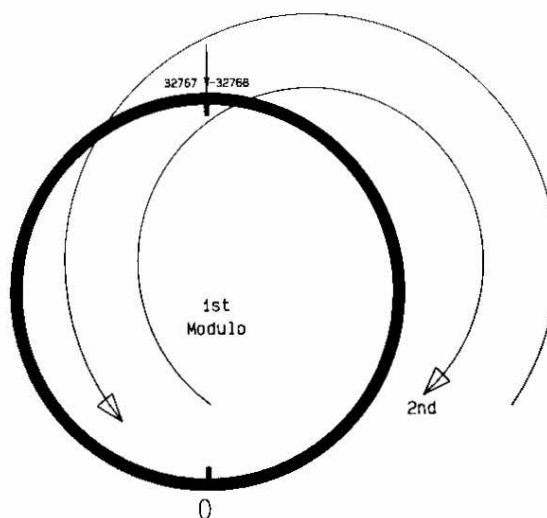
1.2.2 Obtuse Signals and the L_p Norm

It can be shown that there exist nonsinusoidal well-behaved input signals which will drive a unity gain filter output past unity. For example, the time-reversed impulse response of a filter, when used as input to that same filter, will drive the output past unity. (The interested reader is referred to the topic of “Matched filters” [10].) So designing a filter for a unity gain transfer function does not ensure the absence of overflow at the output for all input signals.

Input scaling criteria have been developed to prevent overflow at *any* node of a given topology, although the

EXAMPLE:

DESIRED RESULT	MODULO RESULT	
32512	32512	
+ 256	+ 256	
32768	-32768 (2nd Modulo)	
- 768	- 768	
32000	32000 (1st Modulo)	



Jackson’s Rule: Any number of additions and/or subtractions may occur. Intermediate results and operands may fall into any modulo. As long as the final result is made to fall into the first modulo by design, it will be representable in two’s complement at the chosen wordlength, and a valid result.

Fig. 5. Two’s complement is a modulo arithmetic. The first modulo (ring) for 16 bits is shown. The arrows help visualize the traversal from the first into the second modulo and then back again.

sensitive nodes are only at multiplier inputs, as previously discussed. One of these criteria generates scale factors referred to as L_p norms, which are derived in the frequency domain in the following way:

$$\|F\|_p = \left(\frac{1}{\omega_s} \int_0^{\omega_s} |F(\omega)|^p d\omega \right)^{1/p} \quad (5)$$

where $F(\omega)$ is the transfer from the input to the sensitive node and ω_s corresponds to the sample rate in radians per second. A particular value of p is first chosen and the expression is evaluated for every sensitive node. L_p is then assigned as the inverse of the maximum $\|F\|_p$ with respect to all the nodes, and is then used as a scalar to the input of the filter section in question. The values of p usually adopted are 1, 2, or infinity. The case $p = 1$ ensures that the mean of the magnitude over frequency at the sensitive node will not exceed that of the input, while the case $p = 2$ ensures that the average power over frequency at the sensitive node will not exceed that of the input. When $p = \infty$, this norm is the most conservative of the L_p norms because in this case $\|F\|_p$ becomes simply the maximum of $|F(\omega)|$ with respect to frequency. Intuitively, we can understand why this is conservative because the assumption is being made that all the filter input energy resides at a single frequency [5]. This choice of p (infinity) is convenient, however, for test purposes.

Not being the most severe scaling criterion, the L_p norm will not guard against the time-reversed impulse response as input. In practice, though, such obtuse signals which can potentially overdrive filter outputs are infrequent [5], [7]. Because we use the direct form I, it is necessary to think about scaling the input with respect to only the output node. Therefore, while not precluding output overflow, the unity gain design criterion is usually sufficient corresponding to a p value of infinity in the frequency regions of precisely unity gain. The responsibility to ensure that boost filters are not overdriven is, again, left up to the user [35, p. 1008]. (Recording engineers typically work at approximately 8–24 dB below saturation to give themselves adequate clipping headroom.)

1.2.3 Overflow Prevention and Practice

Overflow detection must be performed at each output stage in a cascade *on every sample output*, and upon detection, saturation of the affected output (appropriately positive or negative) must take place. Regardless of the precautions taken to prevent output overflow, detection is mandatory under absolutely any and all conditions in a real-time operating environment. If output saturation is not performed, the filter output will be a hard nonlinearity, due to 2's complement wraparound, which is never audibly pleasant. If not performed on a per-sample basis, then upon recovery from the overflow condition, the filter can enter into weakly nonlinear modes which only show up as subtle but constant deviations from nominal THD + N measurements performed on the filter output for an input

sinusoid. A theoretical explanation of these overflow oscillations is not given here, but can be found in the standard literature. We stress that saturation is implemented only for the final output of each stage, as pictorially represented in Fig. 7 (ignore the truncator Q box for now), and not during intermediate calculations.

The proper way to detect output overflow is discussed in Appendix 2.

1.3 Summary

By designing unity gain filters using direct form I structures, we can skirt the input scaling issue altogether by allowing intermediate accumulator overflow. The choice of the direct form I topology does not exclude us from considerations of the effects of recirculating output truncation noise, however, so we will have to deal with this noise problem in a different manner than by judicious choice of topology. For a good survey of standard topologies see [12].

2 TRUNCATION NOISE AND ITS PROPAGATION

The currently available DSP processors, such as the ADSP-2100 and TMS32010, routinely provide double wordlength products and accumulations. This is a tremendous advantage to the DSP algorithm designer because it obviates the need to consider truncation noise effects at any node other than the output node in the direct form I digital filter structure.

Accepting the fact that the direct form I is by no means the most quiet topology, we must now deal with truncation noise generated at its output. We must consider the effect that this noise has on filter operation when unavoidably fed back into that filter. First-time designers of digital filters tend to ignore this aspect of the implementation because it is considered a second-order effect. If the application is one that involves fidelity of reproduction or accurate frequency synthesis [23], then this aspect should not be ignored. On the other hand, if the primary interest is an implementation ending up in a low-end sound effects processor, then this topic need not be of concern.

Fig. 8 explicitly shows the direct form I topology having the nonlinear truncation operator (Q for quantizer) at its output. Fig. 7 shows how the truncator and the saturator are used in the same circuit. We omit the saturator in the remaining figures for clarity.

2.1 Low Signal Level Aberrations

The most apparent manifestation of the effect that output truncation has on the output signal of a second-order digital filter occurs when low-level signals are input. An anomaly emerges as a discrepancy between the transfer function of the filter for the full-level (unity amplitude) input as compared with the transfer function for an attenuated input signal. For ideal linear filters, the transfer function is supposed to be independent of the input signal level. Let us say that a sinusoid at an amplitude approximately 30 dB below unity is input to a 16-bit digital filter having a modest gain and low

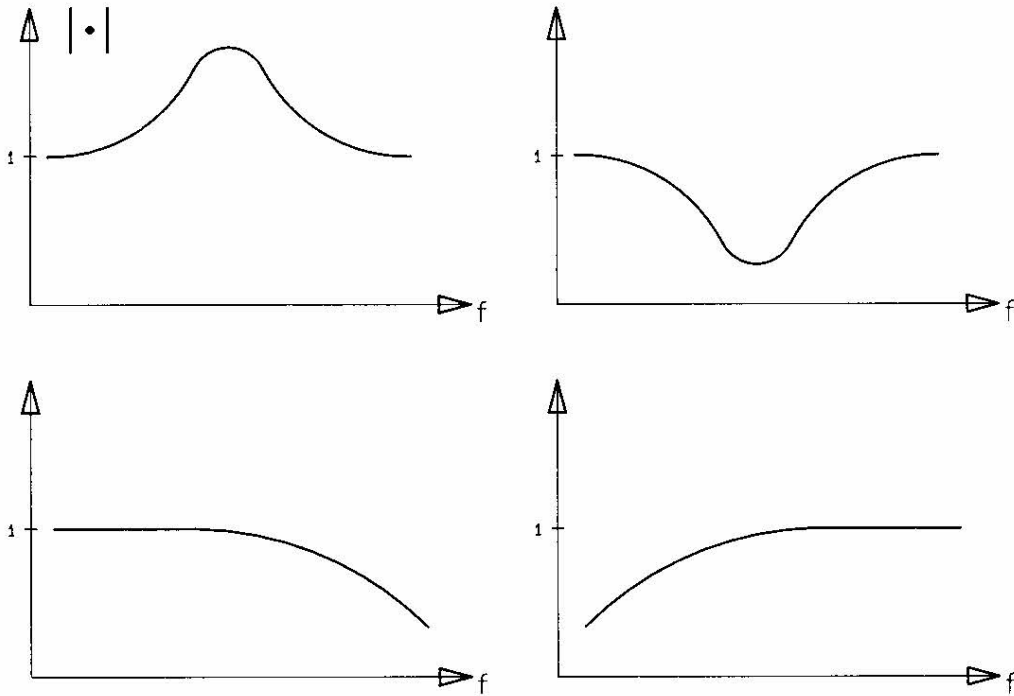


Fig. 6. Unity gain design.

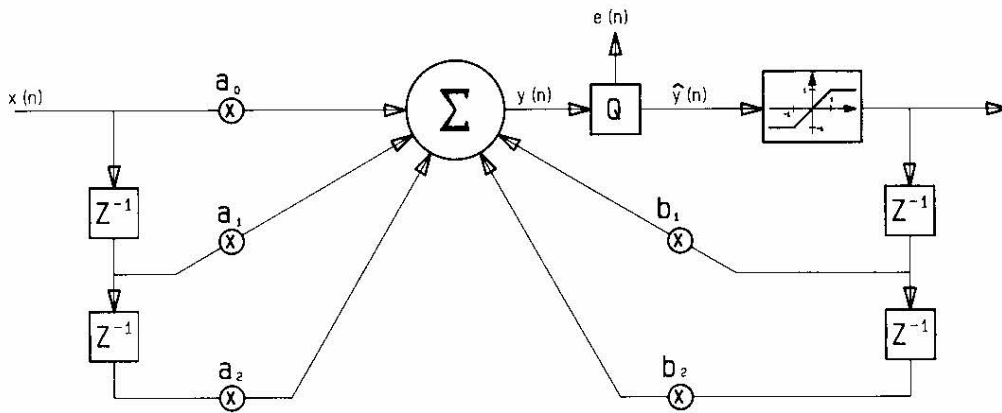


Fig. 7. Truncator Q and saturator.

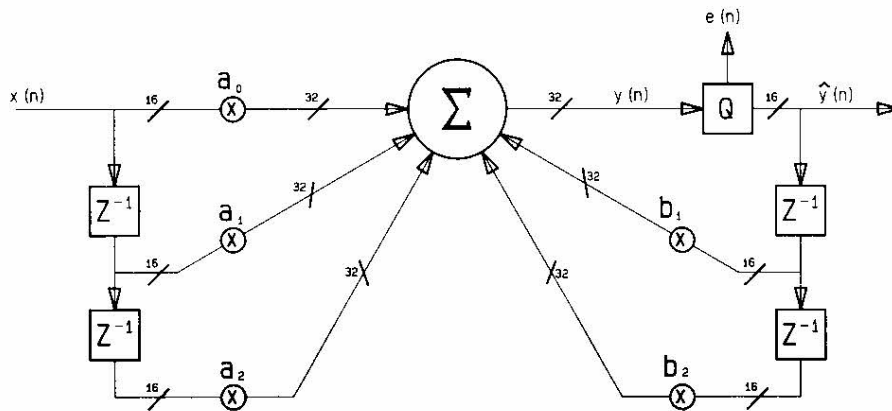


Fig. 8. Direct form I having truncator.

critical frequency. Then independent of the resolution of the coefficients, the output sinusoid will exhibit, to the eye, a variety of amplitude aberrations when viewed on an oscilloscope. These are a direct consequence of the nonlinearity introduced by truncating the accumulator output.

This effect is not hard to reproduce in the laboratory. The easiest one to see is called the *jump phenomenon*. This aberration was named by Kristiansson in 1973 during his study of it [13]. The phenomenon appears as abrupt shifts in the amplitude or dc level of a digitally filtered sinusoid.

The truncation operation is unavoidable because the inputs to the multipliers must be 16 bits while the accumulator accepts and produces 32-bit results. Also, the outside world is most likely 16 bits, so that regardless of the accumulation wordlength, the filter output must pass on 16-bit signal data. Unfortunately the amplitude of the transfer function at any one point along the curve is not a constant but is rather some function of time as caused by this nonlinearity. Needless to say, as the engineers of these digital filters, this performance falls short of our expectations.

The aberrations are even more pronounced at low input frequencies, less than 100 Hz, where one does not have to go as far down as -30 dB to start seeing these effects. They can be seen at input levels of 0 dB. If a common pen plotter is used to characterize a digital filter over frequency, then small ripples will be seen in the plot that are not the same among multiple plots of the same transfer. Further, in the case of 0-dB mid-frequency signals, even if the amplitude aberrations are not apparent on an oscilloscope or a pen plotter, then a THD + N analyzer picks up the problem by calculating a figure upward of 10 dB worse than nominal. Indeed, the gentlest filters will cause at least a 10-dB degradation in THD + N. If the amplitude aberrations are so bad as to be visible, then the degradation to THD + N will be much more than 10 dB.

There are at least two workable solutions to this problem, but let us first perform a simple analysis of the truncation error in an attempt to put a handle on it. Once this is done, the solutions will be straightforward.

2.2 Truncation Error Math

Referring to Fig. 8, we have designated the 32-bit quantity $y(n)$ as the fixed-point output of the digital filter while $\hat{y}(n)$ is the truncated output. We would like to derive a frequency domain expression for the truncated (actual) output signal $\hat{y}(n)$ in terms of the input and truncation error signals. To do this we may define the ideal output $y(n)$ as a 16-bit integer $\hat{y}(n)$ plus a signed fractional part, the truncation error signal $e(n)$,

$$y(n) \triangleq \hat{y}(n) + e(n) . \quad (6)$$

The input signal, generally a 16-bit integer, is called $x(n)$, allowing us to write

$$y(n) = \sum a_i x(n - i) + \sum b_i \hat{y}(n - i) . \quad (7)$$

This is the equation that would be used in a program to perform the filtering function. The 32-bit output $y(n)$ is then truncated by taking the integer part and storing it as $\hat{y}(n)$. We can gain more insight into Eq. (7) by first using Eq. (6) to rewrite it as

$$y(n) = \sum a_i x(n - i) + [\sum b_i y(n - i) - \sum b_i e(n - i)] . \quad (8)$$

Since we know the source of the error signal $e(n)$, it is not truly a random variable, neither is it uncorrelated to the output signal. Indeed, the error signal is deterministic because for the same input signal and coefficient set; it can be replicated exactly. Therefore we may justifiably assume that it has a frequency domain representation.

Then, using Eqs. (6) in (7) in the frequency domain,

$$\hat{Y}(z) = X(z) \frac{\sum a_i z^{-i}}{1 - \sum b_i z^{-i}} - E(z) \frac{1}{1 - \sum b_i z^{-i}} . \quad (9)$$

This important Eq. (9) states that the truncated-output spectrum of the classical direct form I digital filter topology is equal to the ideal filter transfer times the input spectrum plus some "error function," multiplying the truncation error spectrum $E(z)$. The truncation error spectrum is amplified by the poles of the filter transfer function *whether the filter boosts or cuts*, and is therefore highly correlated with the resonances of the filter. This is bad news because the gain introduced by the poles can be as large as 90 dB for some practical cases. It is interesting to note that $E(z)$ is not affected by the feed-forward paths in this circuit. This is one of the reasons we chose this topology since this is an advantage.

Eq. (9) will serve as our reference. Any improvements made to minimize the effects of truncation error will be gauged to the last term in Eq. (9) since this is what we are left with if we do nothing about the truncation error problem.

2.3 Truncation Error Feedback Math

The first solution to impact the truncation error spectrum was suggested to us by a Philips Corporation document [14], which describes the DSP involved in the design of their oversampling D/A converters which have found wide use in Compact Disc players. Fig. 9 shows the first-order *error feedback* scheme. Here $\hat{y}(n)$ is again the 16-bit truncated (integer) value which goes to the outside world, while $y(n)$ is the 32-bit nontruncated output signal. The error signal $e(n)$ is a signed quantity formulated as implicit by Eq. (6),

$$e(n) = y(n) - \hat{y}(n) . \quad (10)$$

We like to define the error in this sense because it is sometimes a facile method from the programmer's point of view, where positive-signed inputs to the accumulator

will result later. This means that the error signal is formed by subtracting the 16 bits that form the integer part of the 32-bit accumulator output, from the accumulator itself. This assigns the signed fractional value to the error signal whose magnitude can, thus, never exceed 1.

In all our error feedback schemes we will never use rounding to form $\hat{y}(n)$ simply because it uses more program steps and gains us little advantage. Recall that statistically, the difference between rounding and truncation is a 1/2 LSB dc offset in the time average of the error signal. Hence the output signal spectrum also experiences a small dc offset when truncation is used. Rounding is used when the absolute dc value of a quantity is of concern; dc is of little concern in the audio signal.

Notice in Fig. 9 that the error signal $e(n)$ is delayed and then fed back to the accumulator. We show a multiplier $K = \pm 1, \pm 2, \pm 4$ one value (to be determined) multiplying the delayed error signal. In order to avoid the use of the hardware multiplier, K is constrained to a nonzero trivial binary integer. Although we will determine the optimal value for K given this constraint, later we will relax this constraint at the expense of more intensity. The 16-bit input to the filter structure is $x(n)$. We wish once again to find a frequency domain equation that expresses the truncated output signal in terms of the input and truncation error signals. As before,

$$y(n) = \sum a_i x(n - i) + \sum b_i \hat{y}(n - i) + Ke(n - 1) \tag{11}$$

This is the equation we would use in our program. Notice the new term in Eq. (11) as compared with Eq. (7). $\hat{y}(n)$ is formed within the machine by truncating $y(n)$, while $e(n - 1)$ is formed using Eq. (10).

Using Eq. (6) in Eq. (11) we get

$$y(n) = \sum a_i x(n - i) + [\sum b_i y(n - i) - \sum b_i e(n - i)] + Ke(n - 1) \tag{12}$$

and in the frequency domain, using Eq. (6) on the left-hand side of Eq. (11),

$$\hat{Y}(z) = X(z) \frac{\sum a_i z^{-i}}{1 - \sum b_i z^{-i}} - E(z) \frac{1 - Kz^{-1}}{1 - \sum b_i z^{-i}} \tag{13}$$

Eq. (13) has the same form as Eq. (9), except that the error function, the last term in Eq. (13), has been modified slightly due to the new error feedback path. What we wish to know at this point is whether the modification to the error function has diminished the impact of the truncation error $E(z)$ at all. We have already seen how the poles of the filter amplify the truncation error, as shown by Eq. (9). Does the introduction of a zero in the numerator of this new error function make it smaller in some way? The optimum choice of a zero location would be somewhere right upon the unit circle in the z plane because this would bring the magnitude of the error function down to absolute zero at some audio frequency. The error function now has a zero at $z = K$. But since K is constrained to be a nonzero binary integer, then the only two choices lay upon the real axis at $z = \pm 1$. A zero at $z = 1$ corresponds to 0 Hz, whereas $z = -1$ would correspond to a zero at the Nyquist. The choice at first appears to be arbitrary from a theoretical point of view. The acquisition of a zero in the error function will certainly diminish its impact in either the very low or the very high frequency region. We choose a positive-valued $K (= 1)$ for empirical reasons discussed below.

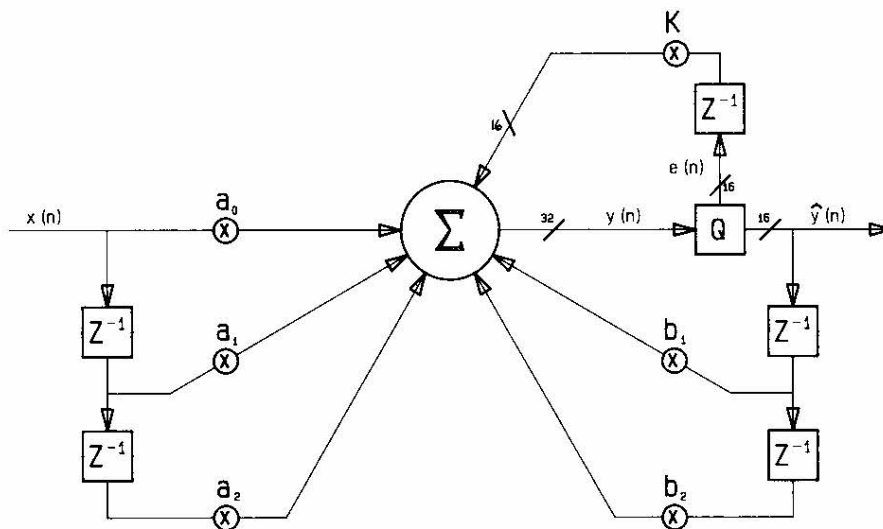


Fig. 9. First-order error feedback.

Fig. 10 shows sketches of the various error situations assuming $K = 1$. The horizontal line in Fig. 10(a) represents the magnitude of $E(z)$, here assumed to be white. Fig. 10(b) shows the error function of Eq. (9) multiplying $E(z)$, where $E(z)$ is amplified by the poles of some arbitrary second-order filter. Fig. 10(c) shows the new $1 - z^{-1}$ term by itself, from the error function of Eq. (13). Notice that there is a boost above unity in the high-frequency region in this graph. This graph shows that the new zero term will attenuate $E(z)$ at some frequencies, but will boost it at others. The attenuation at dc is complete. Fig. 10(d) shows the result of putting the zero into the numerator of the error function of Eq. (9), which results in the error function of Eq. (13) multiplying $E(z)$.

We believe that the situation shown in Fig. 10(d) is an improvement over that in Fig. 10(b), but it is not obvious why. Fig. 10 might persuade us that since the zero in the numerator acts to boost high-frequency truncation noise, we should not use this "error feedback" technique. But we have yet to contradict the classical assumption which states that the truncation noise spectrum $E(z)$ is white. A contradiction would have bearing on our decision. Our observations in the laboratory tell us that much of the energy in $E(z)$ lies in the low-frequency region. O'Donnell [15] independently observed a nonwhite truncation error spectrum on other second-order topologies and filter types years earlier. So the graph of the truncation error $E(z)$ might be more realistically portrayed in Fig. 10(a) as a low-pass signal. Knowing this, we choose to place the zero introduced by the first-order error feedback scheme at 0 Hz, which is the choice nearest to the frequency region where most of the truncation noise is found in practice.

There is a far more practical reason to choose to put

the error function zero at dc, that is, if we know ahead of time that the boost of $E(z)$ by the poles of our filter predominates in the low-frequency region. Had we designed a filter whose critical frequency were nearer Nyquist than dc, we may well have decided to try $K = -1$ for an error function zero at Nyquist, our knowledge of $E(z)$ notwithstanding.

2.3.1 First-Order Error Feedback Summary

The degree of amplification of $E(z)$ by the poles of the filter function is now mitigated by the introduction of a new term in the numerator of the error function of Eq. (9), as shown in Eq. (13). It has been found that this trivial feedback scheme really helps to linearize an inherently nonlinear digital circuit. The amplitude aberrations that were seen in low-level filtered signals not having the benefit of error feedback, actually disappeared when the feedback was introduced, to the extent that the impact on a reference THD + N measurement was negligible. (The reference THD + N measurement was made bypassing the digital filters altogether.) Further, the scheme only requires three additional one-cycle TMS32010 instructions per stage, making it computationally attractive.

2.4 Second- and Higher Order Error Feedback

Fig. 11 shows a more powerful error feedback scheme, this time using two trivial multipliers for second-order feedback. Using reasoning similar to that above, the optimal K coefficient set places zeros on the unit circle in the error function closest to the anticipated frequency region of high noise boost by the filter poles [16]. These zeros do not adversely affect the filter transfer function in any way. Table 1 summarizes the viable choices of trivial multipliers and

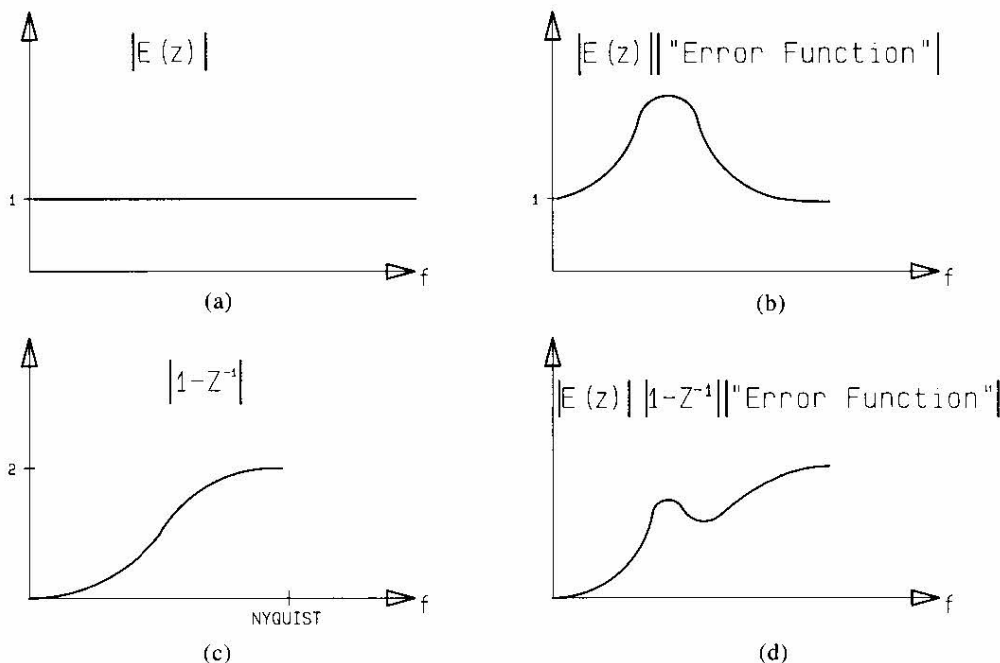


Fig. 10. Presumably white truncation noise (a) is passed through a digital filter having an error function as in (b) and an error feedback network having the transfer (c), resulting in truncation error spectrum shaping (d).

their associated frequency regions, normalized to 2π .

While not unequivocally optimal, these error feedback schemes are economical, and their power should not be underestimated. Truncation noise recirculation is worst for filters of high Q or extreme critical frequency. As it turns out in audio filter design, you get the double whammy for extreme critical frequencies because poles need to be quite close to the unit circle to counteract the proximity effects of the mirrored (conjugate) pole. This is the reason why these feedback techniques are effective. They reduce the truncation noise in the filters that need it the most.

Still more elaborate feedback circuits have been devised such that the error function spectral amplitude is pushed predominantly into the high-frequency region [17] while squelching the lower frequency noise toward zero over any desired bandwidth. Of course, these schemes proportionately boost, by many times, the power of the high-frequency noise, but it is assumed that some other mechanism will be designed to filter that noise out.

These error feedback schemes have become collectively known in the literature as error spectrum shaping (ESS) [11], [18], [19]. All these DSP techniques accomplish the same goal: they selectively shape the truncation error $E(z)$. ESS techniques were developed with an eye toward trivial error feedback coefficients, but there is no necessity for this constraint if we are willing to pay the price. We now explore another ESS scheme which is optimal over the entire band, but just a bit more intensive than the error feedback schemes discussed so far.

2.5 Truncation Error Cancellation

We constrained the error feedback coefficients in our derivations above to be integers. We relax this constraint

here. In Higgins and Munson [20] it is shown that the truncation error cancellation scheme, which follows, can be considered to be the optimal ESS filter, but with no constraints imposed on the feedback coefficients. Further, using a classical linear statistical analysis they show that the 16/32-bit truncation error cancellation scheme is equivalent in noise performance to a double-precision (32/64-bit) realization and outperforms linear state space (state variable) topologies. It is only when we go to a 24-bit processor employing truncation error cancellation that performance exceeds the double-precision realization.

We point out two advantages of the following scheme over a standard double-precision implementation: 1) only the signal feedback paths require double-precision bit widths, and 2) the multiplier inputs are always of the signed variety.

2.5.1 Wordlength Considerations

This new error "cancellation" scheme was presented in a thesis by Rothaar [21], where wordlength of coefficients and signal paths were of great concern. The criterion in that thesis was that no degradation past the 16-bit noise floor of the input signal would be tolerated over the entire audio band for all filter settings. It was

Table 1. Error feedback zeros.

K_1	K_2	Region θ
+2	-1	0 Twice
-2	-1	π Twice
0	+1	0 and π
+1	-1	$\pi/3$ Once
-1	-1	$2\pi/3$ Once
+1	0	0 Once
-1	0	π Once
0	-1	$\pi/2$ Once

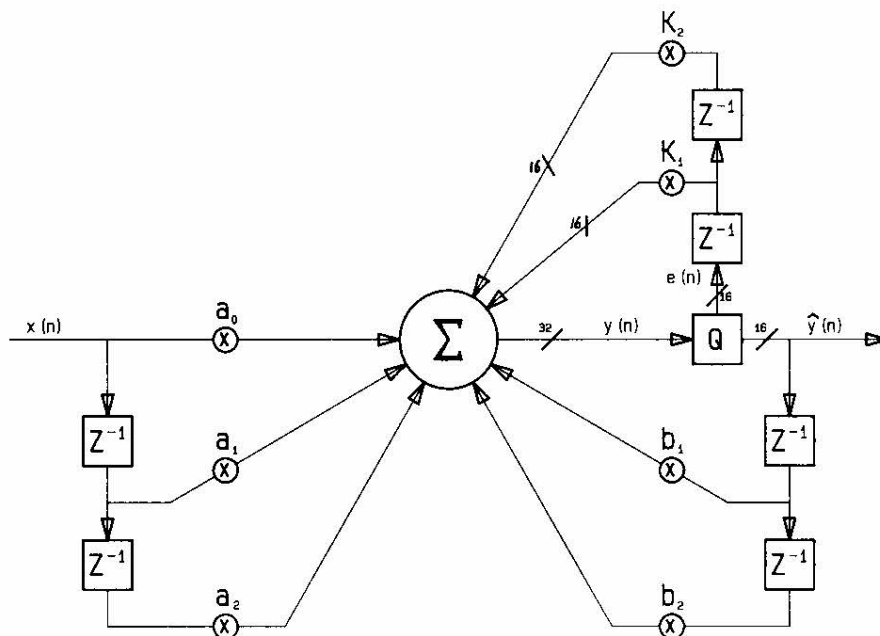


Fig. 11. Second-order error feedback.

shown that to achieve that performance from a standard second-order section *without* error cancellation required at least a 34/54-bit (multiplier/accumulator) architecture. This means 34-bit signal paths. (54 bits was the limitation on the simulator.) It was illustrated how filters having low critical frequencies place the most demands on the second-order topology. A wordlength of 24 bits was ultimately settled on for the coefficients; 16 bits was not enough. The signal data paths followed suit, resulting in a 24/48-bit system with error cancellation.

Rothaar's 24/48-bit system using error cancellation is ideal for the most demanding high-fidelity audio digital filtering applications since the 24-bit internal signal paths guarantee that the filters themselves operate with about a 138 dB SNR. It is advantageous to have the interstage signal paths exceed the 16-bit input wordlength because in this manner each stage's output-truncation noise propagating through a cascade can be kept well below the signal noise floor, since only the *last* filter in the chain need truncate its output at the 16-bit level. Assuming that truncation noise would accumulate through a cascade of N error-canceling second-order stages according to the absolute worst-case classical assumption of correlated noise sources [$10 \log(N^2\sigma^2)$], this still means that hundreds of filters could be cascaded before the DSP designer would even need to *consider* its effects on a 16-bit signal. On the other hand, and in the best case of the classical assumption concerning *uncorrelated* truncation noise propagation [$10 \log(N\sigma^2)$] the passing of 16-bit signal data between stages results in an immediate buildup above the -90-dB signal noise floor, approximately 3 dB after one stage, 6 dB after two, 9 dB after four, 12 dB after eight, and so on. (The best case, that is, uncorrelated noise sources, is usually the classical assumption taken, and although the least conservative of the classical analyses, it is the one most proven in practice. The worst case would double the decibels in the preceding list. For more on classical noise analysis, see [22, chaps. 3, 7], [5, chap. 11] [4, chaps. 9, 8].)

2.5.2 Signal Justification

If we agree that the coefficients need be 24 bits wide for sufficient control over low critical frequency filters, then given a 24/48-bit system, the question arises as to where one places the 16-bit input signal within the 24-bit register. We have infinite accumulator headroom as discussed earlier, so we need not worry about overflow. Since it is those bits to the right of the binary point that determine noise performance, we want to position the input signal with full left justification because the binary point follows the LSB of the 16-bit (Q0) input signal and we want as many fractional places as possible. We might think of this left justification as a scaling upward, by 48 dB, of the 16-bit input signal with no loss to the original quantization SNR of 90 dB. Now if the input signal absolutely must be scaled downward, as in a high-order filter design, there is some room. Regardless of the justification, we would

always want to pass the 16-bit integer plus 8-bit-fractional output signal to the next stage in a cascade for the reasons discussed above.

There is a trend in DSP chip design that provides for headroom bits in the accumulator, that is, the accumulator bit width exceeds the product register width. The section of this paper on scaling and overflow has shown this feature to be superfluous for IIR digital filtering applications.

We now show that the technique of truncation error cancellation will keep the noise level of each individual 16/32-bit second-order stage at the 16-bit boundary. The only significant noise introduced by each second-order stage itself will come when the one truncation is performed at its output.

2.6 Truncation Error Cancellation Math

Referring again to Fig. 11, we wish to write a frequency domain equation for the truncated output in terms of the input and truncation error signals. This time we will not place any constraints upon the error feedback coefficients K_i .

$$y(n) = \sum a_i x(n-i) + \sum b_i \hat{y}(n-i) + \sum K_i e(n-i) . \quad (14)$$

This is the equation that we program. The multiplications involving the fractional quantity $e(n-i)$ are performed using no scaling because it is well represented numerically, having no leading zeros. The decimal point in the error accumulation is higher up in the accumulator, though, and it must be shifted to the right before combining it with the signal accumulation.

Using Eq. (6),

$$y(n) = \sum a_i x(n-i) + [\sum b_i y(n-i) - \sum b_i e(n-i)] + \sum K_i e(n-i) . \quad (15)$$

If we set the K_i equal to the b_i , then they cancel each other and Eq. (15) becomes

$$y(n) = \sum a_i x(n-i) + \sum b_i y(n-i), \quad K_i = b_i . \quad (16)$$

The 16-bit error feedback paths in Fig. 11 now go to 32-bit widths after the multiplications because they are no longer trivial. In the frequency domain, using Eq. (6) in Eq. (16),

$$\hat{Y}(z) = X(z) \left(\frac{\sum a_i z^{-i}}{1 - \sum b_i z^{-i}} \right) - E(z) . \quad (17)$$

Comparing Eq. (17) to Eqs. (13) and (9), this is the best situation for the truncation noise that we have yet encountered since the error function is essentially gone as there is no amplification of $E(z)$ due to the poles of the filter. Eq. (17) says that the truncated output spectrum is equal to the ideal filter function of the input

spectrum plus the (negative) truncation error spectrum. The error function is decoupled from the filter transfer function.

This is indeed the best we can ever do with respect to truncation noise performance. It just does not get any better than this, no matter what circuit topology is used. Since the ideal output $y(n)$ is a fixed-point number having a nonzero fractional part, when we convert the double-precision output of the accumulator $y(n)$ to its 16-bit single-precision truncated value $\hat{y}(n)$, the truncation error spectrum unavoidably remains. But this is all that remains. The truncation error does not feed back. If we can accommodate the added computation time due to the two new feedback paths in Fig. 11, then this result is fabulous because we have a way to implement extremely quiet filters. We mention that this technique would be indispensable for high-purity digital oscillator applications [23].

Notice that the requirements for perfect error cancellation in Eq. (15) are 1) the signal feedback coefficients must equal the error feedback coefficients; nothing is said concerning the precision of those coefficients; and 2) no rounding may be performed in the formulation of the error accumulation.

2.7 Truncation Noise Spectrum $E(z)$

If $E(z)$ would add only 3 dB to the filter noise floor when $\hat{y}(n)$ is formed, then why should we be so concerned with it? The reason is because these second-order filters will most likely end up living in a cascade. (Higher order filters are often formed this way.) One interesting application is that of cascading a number of unity gain second-order filters to configure what is called a *parametric equalizer* in which each second-order filter exerts influence over a disparate part of the audio band. Typical configurations will see four cascaded second-order stages that could add a total of 9 dB [$10 \log(N\sigma^2)$] of truncation noise to the -90 -dB filter noise floor. We need in that case to consider the use of a larger wordlength processor (>16 bits) at a premium.

The assumption that the truncation error spectrum $E(z)$ is white is, more often than not, false. We reiterate that the preponderant error spectrum energy was empirically found to be concentrated in the low-frequency region. It has been said that the output truncation noise could be much like that of analog-to-digital quantization noise. This might be erroneously rationalized following the garden path: The output $y(n)$ of our second-order filter could be made ideal by using error cancellation, as shown by Eq. (16). If the input to an ideal filter is a (sampled) sinusoid, then because it is a linear system, the output must also be a pure sinusoid. Although the input sinusoid $x(n)$ was described using all integer values, the ideal $y(n)$ must be described by the filter using fixed-point (not necessarily integer) values. A truncation to an integer is required to form $\hat{y}(n)$ from $y(n)$, which should be very much like the quantization performed by an A/D converter in a successive approximation when it gets down to the LSB estimate. . . .

2.7.1 The Barnes Criterion

Barnes et al. [24] make the argument that quantization noise and multiplier roundoff noise are very different phenomena. Although he limits his scope to the output of an isolated multiplier followed by a quantizer, we should be able to extrapolate his results to the present case which sees two multipliers input to an accumulator where the quantizing operation follows the summation; see Fig. 12. (In the standard direct form I structure we have seen that the truncation error $E(z)$ is impacted only by the two feedback paths, the b_i . This is one reason why we chose that structure.) We will make no attempt here to perform a statistical analysis that would support this intuitive leap.

Barnes points out that while the conditions under which quantization noise can be considered spectrally white are quite mild, the corresponding conditions for multiplier roundoff noise are comparatively much more restrictive and highly dependent on the multiplier coefficient value. Specifically, he states that roundoff (or truncation) error will be white, uniformly distributed, and uncorrelated with the signal, in general, only if the standard deviation of the input signal is greater than approximately one half the full amplitude (unity) signal level.

Narrow-band and low-level signals each violate the criterion. In contrast, if the dynamic range of a wide-band signal exceeds only a few quanta, then its quantization noise can be accurately modeled as white in keeping with the traditional analysis. But it was demonstrated in [25], which dealt only with the topic of quantization noise, that even for full-level (unity gain) sinusoids, the ratio of sample rate to sine frequency can easily be adjusted to produce a range of quantization error characterization from purely harmonic to white spectra, the point being that the spectrum of quantization noise is highly signal dependent. Therefore we have an argument that quantization noise is not *always* white.

We hope to have supported the hypothesis, by these arguments, that the character of the truncation noise $E(z)$ is not generally white, but is highly dependent on both the input signal and the values of the feedback coefficients. It stands to reason, then, that in the particular case of a cascade of N second-order stages, where the feedback coefficients are made to be disparate, we might reasonably expect that the truncation noise would amass according to the usual classical assumption [that is, the least conservative assumption, $10 \log(N\sigma^2)$] of uncorrelated noise sources, because the contribution to the total error spectrum energy from each stage is likewise disparately splattered about the audio band.

So, let us see how our 16/32-bit truncation error feedback system fares in a little experiment.

2.8 Physical Measurement Example

In the implementation that follows we use 32-bit double-precision coefficients, as opposed to adequate 24-bit, only because it is convenient to do so. As it stands with the TMS32010, we are restricted to a 16/

32-bit system with error cancellation. So we are forced to perform the coefficient multiplications using double-precision arithmetic, which might be tricky due to the lack of an unsigned multiplication in the instruction

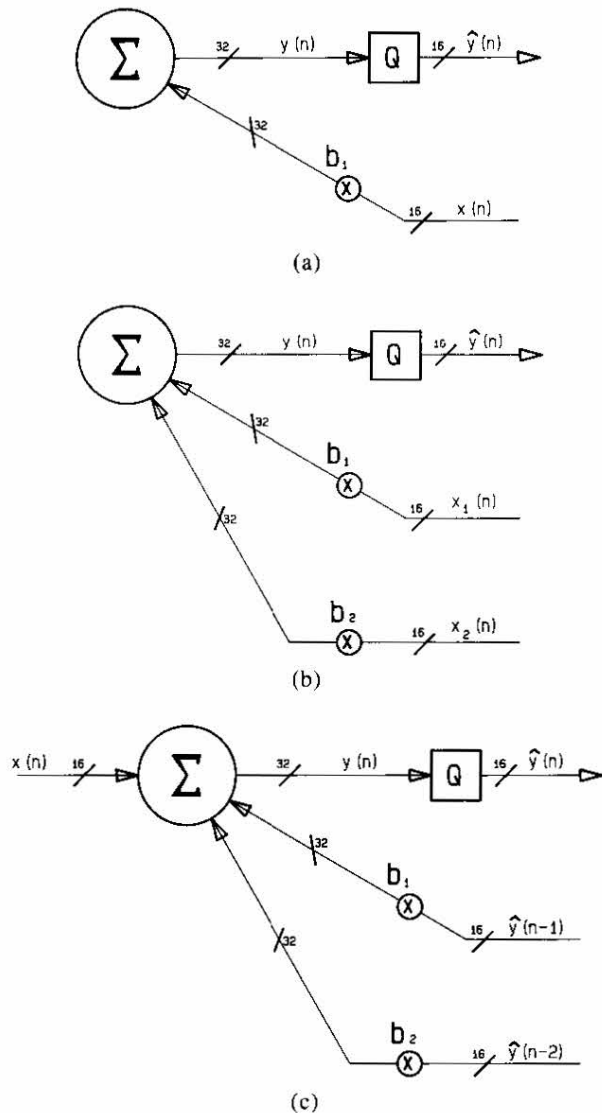


Fig. 12. (a) Barnes's equivalent formulation. (b) Logical extension. (c) Intuitive leap.

set. But an elegant method of double-precision multiplication will be illustrated when residual coding of the filter coefficients is discussed.

We would now like to substantiate the preceding work with actual measurements. Noise plus distortion was measured using a Tektronix 5000 series audio analyzer and oscillator. The A/D/A conversion process was provided by a Sony PCM 701. Table 2 shows the THD + N measurements for three digital filter test setups. Each setup transfer is sketched in Fig. 13. The first and third setups are cascades of four second-order parametrics, while the second setup consists of a single second-order section implementing a 60-Hz notch filter having an incredibly high Q of 1200 and a gain of -60 dB. Both cascade setups arbitrarily sequenced the second-order stages from lowest to highest critical frequency. All these filters are minimum-phase unity gain designs [1], use truncation error cancellation and double precision (32-bit) coefficients via residual coding, as presented later in this paper. Only the error cancellation coefficients are single precision, which could be a liability. The TMS320 code actually used is given in Appendix 1.

The reference THD + N levels were measured with the digital filters completely out of the measurement path, that is, bypassed. Since there is no signal processing going on during these reference measurements, we can attribute all the noise above the theoretical -90 -dB noise floor of a 16-bit system to quantization noise introduced by the A/D/A conversion process. Any noise later introduced above the reference level is most likely then due to either the truncation process or amplification of the quantization noise floor by boost stages. Since each THD + N measurement of a sinusoidal input encompasses the full audio bandwidth, the measurement is repeated with various input frequencies only to provide diverse situations.

2.8.1 Cut Filter Measurement

For test setups 1 and 2, the cut filter cases, the THD + N level recorded in Table 2 was arrived at by subtracting the actual level of attenuation of the unity amplitude input sinusoid from the (negative) instrument

Table 2. Digital filter noise plus distortion measurements.

Sinusoid input frequency (Hz)	Noise reference level (dB) (Filters bypassed)	THD + N (dB)		
		$Q = 20$ $A = -18$ dB $f_c = 50, 500, 5000, 15\ 000$ Hz (Test setup 1)*	$Q = 1200$ $A = -60$ dB $f_c = 60$ Hz (Test setup 2)*	$Q = 20$ $A = +18$ dB $f_c = 50, 500, 5000, 15\ 000$ Hz (Test setup 3)*
50	-80.0	-82.2	-81.8	-79.5
100	-81.2	-80.8	-81.6	-73.0
500	-82.7	-82.8	-82.0	-78.5
1000	-83.0	-82.7	-82.3	-78.0
5000	-82.8	-82.0	-82.2	-79.0
10 000	-82.5	-82.3	-82.0	-77.0
15 000	-82.6	-83.0	-82.1	-80.0

$F_{SActual} = 44\ 470.45$ Hz. "Unity" $\hat{=} 0$ dB. All setups use error cancellation.

* Test setup 1—four cascaded presence filters having negative gain; Test setup 2—notch filter; Test setup 3—four cascaded presence filters having positive gain.

reading. This legitimate compensation serves to make the measurements independent of the particular filter transfer under test. We do this because we are interested in the $S/(THD + N)$ at the filter output from *unity* to the noise floor. We presume that the noise floor is the variable here.

Looking at Table 2, we see that the attenuating cascade and the notch are doing quite well. They are introducing little noise above the reference levels in this particular system. This is shown by Table 2, which compares the noise measurements for test setups 1 and

2 with the reference noise levels. If the reference levels were closer to -90 dB, we might run into a problem here with this 16/32-bit filtering scheme. Remember that regardless of whether the filters under test boost or cut, they still *add* truncation noise. Some of the truncation noise may be getting filtered out by succeeding stages in the cascade.

We note in passing that this -60 -dB notch filter is stable and that it would be difficult to implement in analog due to the required component tolerances. The poles of this filter lie precariously close to the unit

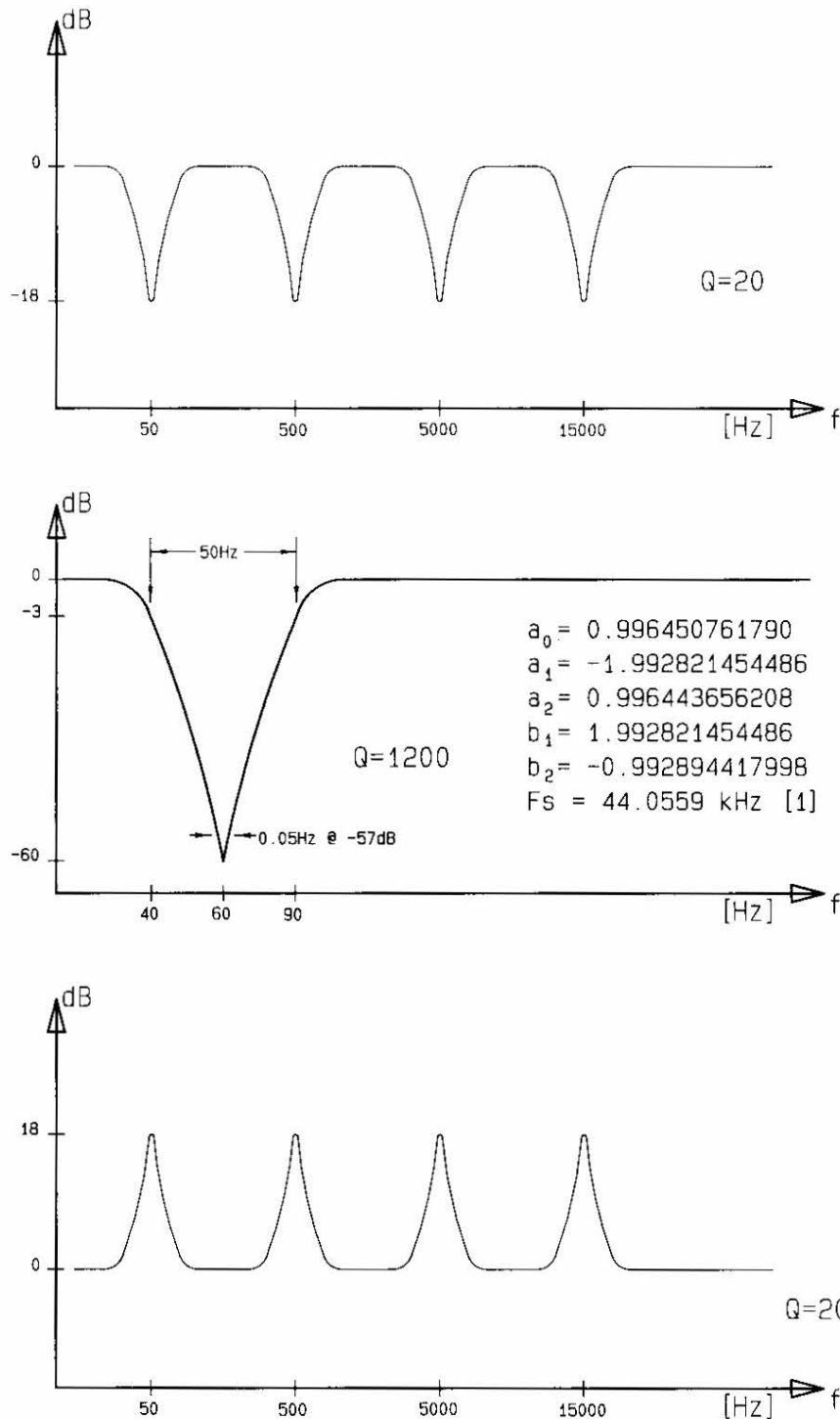


Fig. 13. Filter test setups 1, 2, 3.

circle, as this is an extreme-case design [1], [26]. The noise gain of this topology, not using the truncation error cancellation network, would be 85 dB, the gain of the poles, making this filter unusable ($\text{THD} + \text{N} \approx -5$ dB) without error cancellation. (See Eqs. (30) for $b_1 = 1.992821454486$ and $b_2 = -0.992894417998$ at a sample rate of 44.0559 kHz [1].) As it stands, we are seeing a $\text{THD} + \text{N}$ that is less than 1 dB beyond the system noise reference. This filter has been described as “surgical,” seeing that its bandwidth (at -57 dB) is only 0.05 Hz at the null, while the (-3 -dB) bandwidth at the skirts is 50 Hz. All the critical parameters of this filter were verified to within the resolution of the measurement apparatus. Resolution in frequency = 0.01 Hz, resolution in level = 0.1 dB.

2.8.2 Boost Filter Measurement

In the case of test setup 3, the amplifying cascade, the input signal level was adjusted downward below unity until the output signal level reached unity. The data in Table 2 indicate that test setup 3 is not doing too well. We’re seeing increases of up to 8.2 dB above reference. But it must be remembered that the boosting is above unity and *any* noise falling beneath each respective parametric will be amplified. Even if there were no truncation noise, the system noise floor might be boosted, thus interfering with our desired measurement of the increase in truncation noise. If we calculate the worst-case total boost of white noise across the audio band by all four parametrics, it is about 7.6 dB. Since all four boost filters have the same Q factor (already quite high at 20), the majority of the noise boost comes from the higher critical frequency filters because they span more bandwidth. We might say that it must be the quantization noise floor that is getting boosted, but the $\text{THD} + \text{N}$ is varying so much with input frequency that we are forced to conclude that it must be truncation noise that is getting boosted. This is unfortunate because it means that we need to squelch the truncation noise further. The only way this can be accomplished with the topology we have chosen is to use a 24-bit processor, but we still need the truncation error DSP.

2.9 Going Further

Had we not used truncation error cancellation in the experiment, all the $\text{THD} + \text{N}$ readings would have been horrific. The experimenter can easily verify that the filter transfers we chose are actually quite difficult designs to implement. We were more concerned about whether we needed to use a 24-bit processor than whether we needed truncation error cancellation. A bit of experience tells us that we just cannot get along without some form of error feedback, and so that was taken for granted.

Probably the best way to measure the noise produced by a particular digital filter is to implement it using a DSP chip, and then monitor the output with a noise measurement device. Computer simulations are sometimes cheaper and there are programs available. One

such program is DINAP II: A Digital Filter Analysis Program [27] from Purdue University. This program performs noise analyses of user-specified topologies in the time or frequency domain.

2.9.1 Midsummary

Thus far we have been able to alleviate any input-level scaling requirements through a discerning choice of filter topology, and the truncation noise problem has been grasped. We examined digital signal processing techniques known as error spectrum shaping (ESS), which encompass both truncation error feedback and cancellation. The primary purpose of ESS is to lower the operational noise floor within each filter that we implement. We decided that it would be best to use a 24-bit processor and some DSP in the form of ESS to improve truncation noise performance for the most demanding applications. If we have such a processor at our disposal, then 24-bit coefficients are necessary and adequate for sufficient control over low critical frequency filters. In this case we can skip the next section on residual coding.

If we are using a 16-bit processor, then we will need to find a way to gain precise control over the placement of poles and zeros in the low-frequency region. Once found, then we will have tackled the three major obstacles standing in the way of high-quality digital filtering. In the section that follows we will show how to encode double-precision coefficients on any DSP chip that does not have the ability to perform unsigned multiplications.

3 RESIDUAL CODING OF COEFFICIENTS FOR EXTREME CRITICAL FREQUENCY DESIGNS

This method of coefficient encoding is intended primarily for the TMS32010, TMS32020, or any DSP chip where unsigned multiplications are problematic. The reason for its contrivance is to elegantly perform double-precision multiplications. Specifically, we would like to use filter coefficients having 24-bit wordlengths, but our multiplier only accepts signed data having wordlengths of 16 bits. It will be convenient to use 32-bit wordlength coefficients, so we will do so.

3.1 Residual Coding in a Nutshell

Consider the purely numerical problem stated as follows:

$$y = ax$$

where a is a double-precision and x a single-precision binary integer. If we express a as having a high-order part a_H and a low-order part a_L at the bit level, then we can say

$$y = a_H x (\text{with suitable shift}) + a_L x$$

$$[\text{signed}] = [\text{signed}][\text{signed}] + [\text{unsigned}][\text{signed}] .$$

The second multiplication requires an unsigned input,

so we cannot use a multiplier that only performs signed arithmetic. We could force the MSB of a_L to be zero by shifting it 1 bit to the right before input to the signed multiplier, but we lose precision, albeit 1 bit. Instead, we redefine a_L to make it signed:

$$a_L \triangleq a - a_H .$$

Now a_L is no longer set to the low-order bits of a double-precision binary integer; it is a signed quantity. In the detailed explanation that follows, keep in mind that the complexity of the implementation is no better or worse than that of a conventional double-precision implementation using unsigned arithmetic.

3.2 The Details

Let us name the generic single-precision coefficient c_i . This is the fixed-point decimal representation of the ideal floating-point coefficient, which we have seen in our filter topology. The ideal floating-point filter coefficient c_{Fi} is calculated by a filter design routine. We define the signed difference between the ideal and the single-precision coefficient as ec_i and dub it the *residual coefficient*.

We formulate c_i in fixed-point decimal using the following:

$$c_i = \frac{\text{round}(c_{Fi}2^{qc})}{2^{qc}} \quad (18)$$

where the round function means round to the nearest integer. We use the round function here because it will later bound the residual coefficient to $\pm 0.5/2^{qc}$ which will be advantageous when we encode it. The quantity qc is the number of binary places following the binary point in the coding of the single-precision coefficient c_i .

To "code" a coefficient means to find a 2's complement representation of the fixed-point decimal number. To code c_i ,

$$\text{binary code}(c_i) = c_i2^{qc} = \text{round}(c_{Fi}2^{qc}) . \quad (19)$$

The resulting encoded integer, often expressed in hexadecimal, is used directly in assembly code, representing a 16-bit fixed-point binary number having qc binary places after the binary point. Therefore the sign and integer part must be representable in $[16 - qc]$ binary places. If qc were 12, then we would say that c_i is a Q12 coded number ($= [4.12]$); a Q12 number, for short.

The residual coefficient is formulated in fixed-point decimal as

$$ec_i = c_{Fi} - c_i . \quad (20)$$

The residual coefficient is encoded using an additional scale factor that is equal to 2^{qe} , because it is too small to be represented adequately in a 16-bit fixed-point format. Of course, the filter program must account for

the scaled residual coefficient in the calculations. Due to the rounding procedure, after scaling, the residual coefficient will have a decimal magnitude no greater than 0.5,

$$\text{binary code}(ec_i) = \text{round}(2^{qe}(ec_i2^{qc})) . \quad (21)$$

[Numerically, it has been found that it is better *not* to simplify Eq. (21) by substituting previous equations and then canceling like terms.]

The residual coefficient has qe binary places following the binary point, and qe is not necessarily equal to qc . A good choice for qe might be 15. In that case we would say that ec_i is coded as a Q15 number, meaning that there are 15 binary places following the binary point, one sign bit, but no integer bits [1.15]. For stable second-order filter sections having unity gain and minimum phase, a good choice for qc would be 14 (that is, Q14 = [2.14], having one integer bit) since the magnitude of the coefficients is bounded by 2. Due to idiosyncrasies of the TMS32010, we will set qc equal to 12 in order to save four program steps per filter stage, and qe to 14 for maximum headroom in the separate residual coefficient accumulation.

This notational method of working with fixed-point binary numbers in terms of Q follows those conventions set forth in the TMS32010 user guide [28]. The Q notation follows the same rules as exponents in a by-hand multiplication. For example, if the multiplier were a 16-bit Q12 number and the multiplicand were a 16-bit Q14 number, then the product would be a 32-bit Q26 number, that is, the sum of the two Q factors. (Note the redundant sign bit in the product of any fractional multiplication.)

3.3 Filter Calculations Using the Residual Coefficients

Now that we have encoded the double-precision coefficients, we are ready to perform the digital filtering itself. We tally five double-precision coefficients for the standard second-order section plus two single-precision coefficients for the truncation error cancellation network. This is equivalent to 12 single-precision coefficient multiplications in terms of computation time. The price we have paid for these hi-fi filters is an amount of computation that exceeds that for the basic second-order section by a factor of roughly 12/5. (In comparison, for a 24/48-bit architecture the expense factor would only be 7/5 because all the coefficients would be in single precision which, of itself, entails even less housekeeping.)

To justify the use of single-precision coefficients in the error cancellation network, refer to Fig. 14. There,

$$\begin{aligned} y(n) = & [\sum a_i x(n-i) + \sum ea_i x(n-i)] \\ & + [\sum b_i \hat{y}(n-i) + \sum eb_i \hat{y}(n-i)] \\ & + \sum b_i e(n-i) \end{aligned} \quad (22)$$

where for the sake of less nomenclature, we call the double-precision coefficient the same as the floating-point coefficient since

$$c_i + \frac{ec_i}{2^{2^c} 2^{2^c}} \approx c_{Fi}$$

(double precision on the left, floating point on the right).

Eq. (22) is the one that is actually implemented in the code in Appendix 1.

Using Eq. (6), we can rewrite Eq. (22) as

$$y(n) = \sum a_i x(n-i) + \sum e a_i x(n-i) + [\sum b_i y(n-i) - \sum b_i e(n-i)] + [\sum e b_i y(n-i) - \sum e b_i e(n-i)] + \sum b_i e(n-i) \tag{23}$$

Using Eq. (20) and simplifying, we get

$$y(n) = \sum a_{Fi} x(n-i) + \sum b_{Fi} y(n-i) - \sum e b_i e(n-i) \tag{24}$$

If we compare Eq. (24), which is the single-precision case of truncation error cancellation, with Eq. (16), aside from the use of double-precision coefficients in Eq. (24), the only difference is the (last) error term in Eq. (24). This excess error term arises because we chose *not* to use double-precision (that is, the same) coefficients in the error cancellation network. Ideally, we would like to use double-precision error cancellation coefficients. Keep in mind that perfect truncation error

cancellation, as in Eq. (15), is not dependent on the precision of the error feedback coefficients. The signal feedback coefficients must be *equal* to the error feedback coefficients. In that case, the error feedback becomes error "cancellation."

If we look at Eq. (24) in the frequency domain, using Eq. (6) to derive an expression for $\hat{Y}(z)$, then we get an excess error term at the truncated output that looks like this:

$$E(z) \frac{1 - \sum (b_{Fi} - eb_i) z^{-i}}{1 - \sum b_{Fi} z^{-i}} \tag{25}$$

If the eb_i are small, then the numerator of this error function should approximately cancel the denominator, which comprises the filter poles. The use of single-precision coefficients in the error cancellation network is probably sufficient in most cases except those demanding high precision, which encompass filters of high Q factor or (extreme) critical frequency close to dc or Nyquist. Recalling that the eb_i are bounded in magnitude by $(0.5)2^{-q_c}$, then the eb_i will be smallest when the b_i have the greatest possible precision (that is, the highest resolution). So we now have a theoretical reason to prefer the greater precision Q14 representation of the single-precision part b_i of the double-precision coefficient b_{Fi} rather than the Q12 representation that is used in our TMS code at the present time.

3.4 Residual Coefficient Implementation

When we get down to writing the actual program code for Eq. (22), since the residual coefficients are encoded using a scale factor, they must be accumulated separately. The error cancellation term must also be accumulated separately, but for a different reason. It

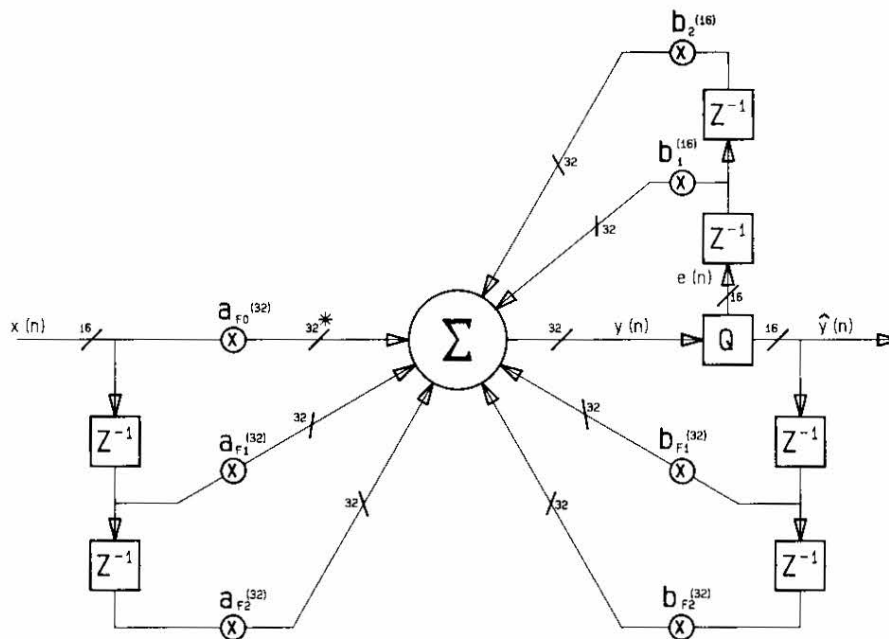


Fig. 14. Truncation error cancellation and residual coefficient coding. See TMS code in Appendix 1 for implementation of path alignment.

is because the error signal $e(n)$ is not Q0 as are the input and truncated output signals; it is Q12, the same Q value as the single-precision coefficients. The most efficient order of calculations is to first group all the residual coefficients, then perform the truncation error cancellation, and finally the single-precision coefficients.

Both the error cancellation and the residual accumulation must provide for worst-case headroom in the accumulator. Overflow is not acceptable at these stages in the calculations because they are nonlinear. In the case of the residual accumulation, since the input and output signals are Q0 and the magnitude of each of the five residual coefficients never exceeds 0.5, a potential overflow factor of 2.5 requires that there be two headroom bits. Thus the residual coefficients need to be in Q14 format ($q_e = 14$). After the residual coefficient scale factor 2^{q_c} is accounted for, the accumulator ends in Q26 format.

In the case of the error cancellation, the Q12 error signal has a maximum magnitude of 1.0, whereas the Q12 feedback coefficients are bounded in magnitude by 2.0 and 1.0, respectively. Therefore a potential overflow factor of 3.0 again requires two headroom bits. But the accumulator will be in Q24 at the end of the error cancellation accumulation. This means we can toss up to five of the MSBs; in the code we discard four for the sake of convenience.

When these two separate accumulations are combined with the accumulation of the single-precision coefficients, we may then take advantage of the infinite accumulator headroom and disregard intermediate overflow, as discussed earlier.

Refer to the TMS320 code in Appendix 1 for an explicit presentation of all the principles discussed in this paper.

3.5 Going Further Still

The purpose of residual coefficient coding is to make the actual frequency response more closely resemble the shape of the theoretical frequency response, while one purpose of truncation error feedback is to make the actual frequency response less dependent on absolute input signal level. The most straightforward way to predict the effects of coefficient quantization on transfer function is to use the desired precision coefficients in the theoretical computation of the transfer. This technique is a valid first-order approximation of the truth and will show the trends of shift in frequency response due to the drift of the pole-zero locations. While coefficient resolution primarily affects the shape of the filter transfer and stability, its impact on noise performance and other signal quantization effects is second order. Coefficient quantization alone does not make the digital filter a nonlinear system. The best method to predict all quantization effects on frequency response is to take the Fourier transform of the actual impulse response of the filter under study. Programs to do this are available; one such program is included in a filter design package from Signix Corp. [15], which can calculate

the FFT of up to 8192 points (version 1.5). Very low critical frequency filters may need more points to capture the complete impulse response.

4 FORCED OVERFLOW OSCILLATIONS

In the first three sections we have covered the most important topics concerning the implementation of digital filters. The current topic concerns the nonlinear behavior of second-order digital filters once output overflow with saturation has occurred and the input remains nonzero. While it is important that the practicing engineer be made aware of this potential problem, this topic is not essential to the understanding of the other material. This phenomenon, it must be stressed, is characteristic of many digital filter topologies and is not peculiar to our particular implementation. All the error feedback techniques we have shown have been reported to minimize this problem which we are about to discuss [29]. The simplest means of circumventing forced overflow oscillations is to ensure somehow that the filter output will not exceed unity.

4.1 Statement of the Problem

Recall that our direct form I structure is immune to overflow problems at *internal* nodes as long as the output is constrained to be less than unity. When the saturating output tries to exceed unity, all bets are off—linear analyses no longer apply and the filter output can be observed, under the proper circumstances, to “lock up” in a mode of oscillation which follows the input frequency but not its amplitude or phase. At the instant this lock-up response occurs, the output phase can be observed to jump suddenly to a new value, which is frequency dependent but difficult to predict. The input signal is present during this form of instability, which is the reason for the term *forced* overflow oscillation. The oscillation will cease when the input signal is taken away (brought to zero). Unfortunately our venerable error feedback techniques cannot help us because the output is no longer representable at the allotted bit width, causing the truncation errors to be themselves in numerical error.

If the input signal is merely reduced in amplitude, then the filter will return to its linear behavior after the input amplitude has dropped to a level that can be considerably lower than the level that elicited the response. (This return to linear behavior will happen, assuming that output overflow detection is performed on a per-sample basis, as previously discussed.) Conceptually, there is an input level hysteresis-loop function that describes this lock-up mode of the filter. The amount of hysteresis is primarily dependent on the pole positions.

A rule of thumb for sinusoidal input signals is this: if the input frequency is below the pole frequency, then there will be *no* forced overflow oscillation *hysteresis*, that is, the oscillations will stop as soon as the input stops overdriving the filter. But when the reverse is true, when the pole frequency is lower than the input

frequency, there *will* be hysteresis. We will, shortly, be more precise in our definition of exactly when the problem occurs, but this rule of thumb is good for filters having poles at less than half Nyquist.

The topic of forced overflow oscillations should not be confused with the more widely researched topics of overflow oscillations and limit cycles. These are the only two types of low-level autonomous oscillations possible in digital filters. We briefly present them here.

4.1.1 Overflow Oscillations

Autonomous overflow oscillations will occur if output saturation is not implemented upon the detection of overflow there. (Remember we do not saturate intermediate results.) This type of instability is self-sustaining once activated, and is due to the wraparound (modulo) characteristics of 2's complement arithmetic. Output saturation eliminates the problem for second-order filters, and so we can dispense with it entirely [30]. It would be nice if saturation solved the forced overflow oscillation problem, but it does not. The fundamental difference between this and forced overflow oscillation is that these oscillations will sustain under zero *and* nonzero input conditions, while forced overflow oscillations only occur under nonzero (large-signal) input conditions.

4.1.2 Limit Cycles

Limit cycles, on the other hand, are a problem due to the inherent nonlinear nature of digital filters provoked by numerical inaccuracies in the signal paths. Limit-cycle oscillations can be sustained under the condition of zero input and are dependent on initial conditions. Activation of limit cycles does not require an overflow mechanism; this is the primary conceptual difference. Nonzero inputs usually break them up, once activated. The cause of these oscillations can be understood intuitively if we realize that even when digital filters are operating in their "linear mode," they are still nonlinear devices to a second degree because of the quantization of the internal signals and the truncation errors. "Digital" filters, by definition, operate on signals that are discrete in both time and amplitude. Therefore no matter what number of bits are being used in the signal path, there will always be the possibility of limit-cycle oscillations in the standard direct forms. The amplitude of these oscillations can be quite large, although typically small and annoying, and the problem ~~worsens when rounding as opposed to truncation is used. Indeed, Chang [31] has shown that limit-cycle oscillations are impossible in first-order digital filters when truncation is used. This limit-cycle problem~~ has usually been solved in the past by increasing the number of bits in the signal paths, which reduces the amplitude of the oscillations to "acceptable levels."

In the error cancellation and feedback schemes we have presented, we diminished the impact of the truncation errors by constraining them to appear only at the output node. We accomplished this by canceling

or minimizing the truncation errors which fed back into the circuit, and by retaining full precision multiplications. Hence the internal signal paths appear much wider than they actually are because we succeeded in preventing a buildup of truncation noise along those paths. The net effect is the same as it would be had we actually increased the signal bit width. The amplitude of any potential limit cycle is greatly minimized compared to that using no form of error cancellation. Thus using cancellation/feedback (ESS) techniques, we succeed in linearizing the digital filter to a much higher degree than do conventional implementations. Indeed, Higgins and Munson [16] showed that these ESS schemes quickly approach statistical equivalence to double-precision implementations, which in itself guarantees significantly improved limit-cycle suppression. Going one step further, Chang [31] showed that for specific values of the error feedback coefficients, limit cycles could be eliminated completely. Those values include the case of truncation error cancellation discussed.

In summary, using ESS techniques, one is able to suppress limit cycles to well below the 16-bit level. We will not examine limit cycles further here because thorough treatments can be found in most standard DSP texts.

4.2 Region of Forced Overflow Oscillation

Fig. 15 shows the second-order digital filter section under study here. Only the feedback paths are shown because the feedforward paths do not contribute to this problem. The saturation arithmetic shown at the output of the accumulator is in keeping with the assumptions we have made thus far. Fig. 7 shows our implementation of the truncator and saturator operators as found in the code in Appendix 1. We believe that their ordering is not critical; this figure is only included here for reference. The truncator is not critical in this analysis and so it is omitted in Fig. 15.

Fig. 16, taken from [32], shows a typical forced overflow oscillation at the output when a sinusoid overdrives the circuit. We talked about the observation that the forced overflow oscillation seems to be worse when the input frequency is above the pole frequency, that is, it is harder to recover from it. Let us hypothesize a possible explanation for this by guessing that recovery time is somehow related to the impulse response of the filter. Let us loosely define "recovery time" as the time it takes for the circuit to stop oscillating once the input signal has been removed.

If we talk about an input sinusoid whose amplitude is just slightly above unity, then it is really just the peaks of that sinusoid which overdrive the unity gain filter. Then let us guess that the input amplitude hysteresis phenomenon has to do with the fact that a high-frequency input sinusoid has peaks which occur at a rate that is too fast for the overflowed filter to recover from. Because the filter does not have adequate time to recover from this nonlinear mode, it is necessary to bring the input amplitude far enough below unity so

that the corresponding amplitude of the filter's impulse response becomes attenuated. This in turn reduces the time required before the internal computations (convolutions) are once again representable at the allocated bit widths. Conversely, when the input sinusoid has a frequency that is below the pole frequency and its amplitude is just above unity, then the overflowed filter has plenty of time to recover and it is not necessary to bring the input amplitude back down much below unity.

Although this explanation may be oversimplified, having used the linear impulse response in the explanation of a nonlinear phenomenon, Claasen and Kristiansson [33], [34, p. 515] proved in 1975 that a necessary and sufficient condition for the existence of forced overflow oscillations is

$$|b_2 h(n)| > 1, \quad \text{for some } n = 0 \rightarrow \infty \quad (26)$$

where b_2 is the second-degree feedback coefficient and $h(n)$ is the impulse response of the ideal linear second-

order digital filter.

Using the residue theorem, the impulse response of the filter shown in Fig. 15 can be derived from its transfer

$$H(z) = \frac{1}{1 - b_1 z^{-1} - b_2 z^{-2}} = \frac{1}{(1 - p z^{-1})(1 - p^* z^{-1})} \quad (27)$$

as

$$h(n) = \frac{p^{n+1} - (p^*)^{n+1}}{p - p^*}, \quad n = 0 \rightarrow \infty \quad (28)$$

where $p = re^{j\theta}$ is the pole and p^* is the conjugate pole of this second-order system. (Note that consideration of zeros in the numerator of Eq. (27) would only change

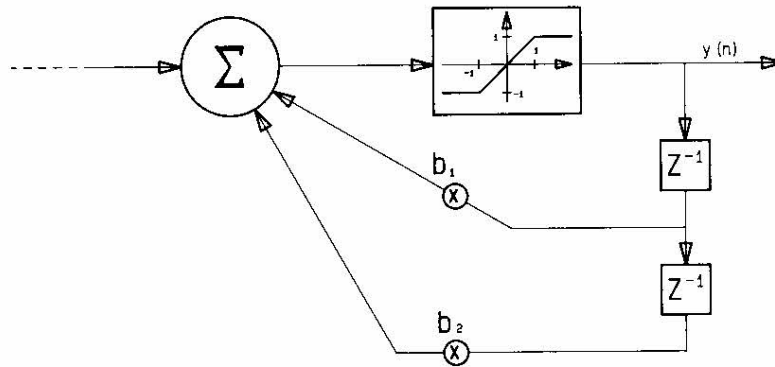


Fig. 15. Saturator, forced overflow analysis.

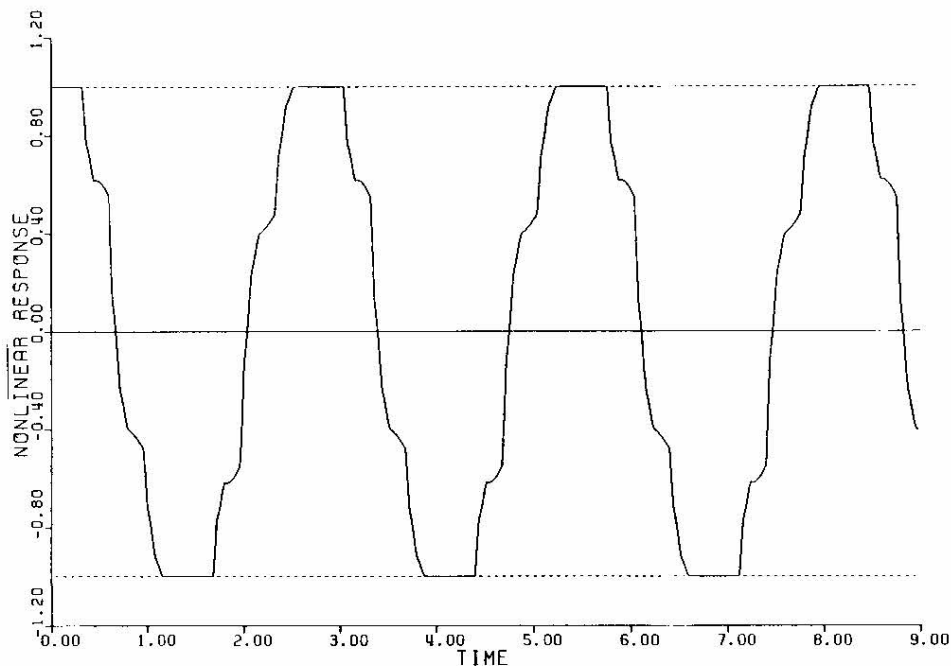


Fig. 16. Typical forced overflow response for sinusoidal input. (© 1983 IEEE [32].)

the absolute magnitude of the consequent impulse response, but not its shape.)

We can rewrite Eq. (28) in its polar form as

$$h(n) = r^n \frac{\sin(n+1)\theta}{\sin\theta} \quad (29)$$

where r is the pole radius in the z plane and θ is the pole angle in the upper half-plane. The pole radius r must be within the unit circle for stability, as is well known.

For θ very close to 0 or π ,

$$|h(n)| \approx (n+1)r^n, \quad \theta \text{ near dc or Nyquist.}$$

Relating the poles within Eq. (27) to the actual filter coefficients, it is straightforward to show that for the positive topology

$$\begin{aligned} b_1 &= 2r \cos \theta \\ b_2 &= -r^2. \end{aligned} \quad (30)$$

The filter stability expressed in terms of the filter coefficients b_1 and b_2 is illustrated by the stability triangle in Fig. 17. All coefficients within the triangle produce stable filters.

Now it is easy to see that there may easily be forced overflow oscillations for poles near dc or the Nyquist since there, by substituting Eq. (30) and the approximation to $h(n)$ into the criterion, Eq. (26),

$$|(n+1)r^{n+2}| > 1,$$

for some $n = 0 \rightarrow \infty$, θ near 0 or π

it is not hard to find an n for which it is true.

Fig. 17, taken from [32], [34], shows the coefficient regions shaded within the stability triangle for which the second-order digital filter is likely to have a forced overflow problem. The trend of the problem area is for values of b_2 close to -1 , which implies that the problem is more severe for filters whose poles are close to the unit circle. The trend as it pertains to b_1 shows problems for filters whose poles lie near dc or Nyquist. Notice that the problem areas are symmetric with respect to the b_2 axis, and that there appear nulls for particular values of b_1 for which no forced overflow oscillations are possible.

Fig. 17 suggests possible recovery schemes that involve the modification of the coefficients to bring them into safe regions. This solution implies that the filter transfers could be changed momentarily to solve the problem when it occurs. This may not be attractive in some circumstances. It may be more desirable to scale input amplitude. A nonlinear approach such as that described by Claasen and Kristiansson in [29] may be more effective. The technique they describe employs a feedback circuit that only comes into play when overflow occurs. Assuming there is some headroom in the accumulator, they feed back the difference between the saturated output and the actual value of the overflowed signal. Under normal circumstances of linear operation, this difference is zero, so nothing gets fed back.

4.3 Forced Overflow Summary

Our goal is to make the DSP engineer aware of this problem. Interestingly enough, there are only a handful of research papers on this subject, and after all is said and done, the solution typically offered is to simply reduce the input amplitude; this is now the best solution

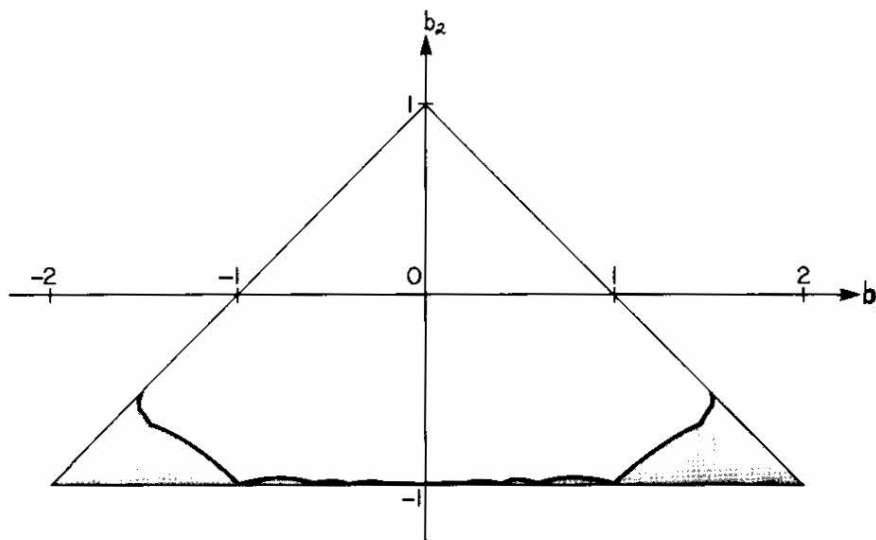


Fig. 17. Stability triangle for positive topology. Shaded region shows where forced overflow oscillations can occur. (© 1982, IEEE [32], [34].)

[35]. But whose responsibility is it to do so, the filter architect's or the end-user's?

5 CONCLUSIONS

One thing that has become eminently clear during the course of this research is that, in terms of noise performance, we could have extremely high fidelity (transparent) digital filtering if we used 24-bit signal paths. For the most demanding applications, the 16/32-bit processors just will not do, unfortunately. 24-bit signal paths will not, however, obviate the use of truncation error cancellation. 32/64-bit processing would be needed to perform brute-force digital filtering using no truncation error cancellation. But note that 24/48-bit processing with truncation error cancellation will outperform the 32/64-bit system and can be up to 8 bits (48 dB) better because the 32/64-bit system can only hold the noise down to the 16-bit boundary. We have also learned that 24-bit coefficients have sufficient accuracy for most audio frequency filtering requirements. From this perspective, a 24/48-bit version of the TMS320 series of DSP chips would be ideal.

Too often the lesson has been learned that the number and orders of the filters designed are determined not by the desired transfer characteristics, but by the sheer processing power of the implementation. For this reason, commercial filter design packages should be amended so that the principal constraint is the filter order rather than band tolerances. The program would then yield the best possible design for the given filter order.

The primary contribution of the topics we have covered is to dispel some popular myths, particularly those concerning input scaling, internal overflow, and the assumptions of "white" noise ("random") error sources. We have found that our truncation noise reduction techniques also serve to eliminate autonomous limit-cycle oscillations, while output saturation ensures freedom from autonomous overflow oscillations. Having hurdled these barriers, we have nearly reached the ideal in digital filtering. There is always room for more work, and we suggest that the most elegant solution to forced overflow oscillation is still unclear. The solution may surface in the area of chaotic system theory.

6 TENETS OF DSP FILTER IMPLEMENTATION

- 1) Overflow is not always a bad thing.
- 2) Truncation noise can be controlled at the expense of computational intensity. The degree of intensity determines the amount of control.
- 3) Signed multipliers *are* useful for double-precision calculations without loss of accuracy or efficiency.
- 4) Floating-point mathematics alone in a DSP chip does not solve the truncation noise problem [36], [37]. Numerical resolution determines noise performance. Resolution is determined by the mantissa bit width. Truncation of the mantissa to less than the required $2N$

bits prior to accumulation increases truncation noise recirculation.

5) Dynamic range is not the same as signal-to-noise ratio (SNR). Dynamic range is the ratio of the largest to the smallest *perceptible* signal level. Since signals well below the noise floor are perceptible by humans, the dynamic range specification often exceeds SNR. It is often made commensurate with the 6-dB-per-bit PCM criterion, but even this can be exceeded by proper dithering [38]. Dynamic range, then, is determined by the digital signal processing. SNR is the ratio of the largest signal representable to the noise floor with no signal present. In digital systems it is probably more meaningful to have a term for that ratio when signal is present: this we call THD + N. Included in it are harmonic distortion, the noise floor, and all other artifacts introduced by the digitization of the signal itself.

6) Coefficients of second-order sections may be linearly interpolated for smooth transitions to target filters without intermediate instabilities. The same is not guaranteed for direct higher order filters (except ladders). The interpolated coefficients must be encoded in double precision for smooth transitions.

7) The propagation delays of digital filters are no worse, in general, than for the corresponding analog implementations.

8) DSP chips should have the capability of shifting the entire accumulator in place and at once in either direction to facilitate the alignment of the intermediate accumulations in the realization of the techniques presented herein. The program-controlled realignment of the product register into the accumulator is also useful. (TMS32020 and TMS320C25 have these capabilities.) DSP chips must be designed to handle the access of long delay lines of any modulo and with nonequispaced taps. DSP chips need to be able to detect overflow out of several of the high-order accumulator bits below the MSB, not just from the MSB.

7 ACKNOWLEDGMENT

The author is indebted to the following people for their contributions to this paper: Richard Cabot of Audio Precision, Hayley Greenberg of Analog Devices, Tom Hegg of Lexicon, and Dana Massie of Next.

8 REFERENCES

- [1] J. A. Moorer, "The Manifold Joys of Conformal Mapping: Applications to Digital Filtering in the Studio," *Audio Eng. Soc.*, vol. 31, pp. 826-841 (1983 Nov.).
- [2] E. M. Cherry, "Transient Intermodulation Distortion—Part I: Hard Nonlinearity," *IEEE Trans. Acoust., Speech, Signal Proc.*, vol. ASSP-29, p. 137 (1981 Apr.).
- [3] L. B. Jackson, J. F. Kaiser, and H. S. McDonald, "An Approach to the Implementation of Digital Filters," *IEEE Trans. Audio Electroacoust.*, vol. AU-16, p. 413 (1968 Sept.).

- [4] A. V. Oppenheim and R. W. Schaffer, *Digital Signal Processing* (Prentice-Hall, Englewood Cliffs, NJ, 1975).
- [5] L. B. Jackson, *Digital Filters and Signal Processing*, 1st ed. (Kluwer Academic Publishers, Boston, MA, 1986).
- [6] J. D. Markel and A. H. Gray, Jr., "Fixed-Point Implementation Algorithms for a Class of Orthogonal Polynomial Filter Structures," *IEEE Trans. Acoust., Speech, Signal Proc.*, vol. ASSP-23, p. 486 (1975 Oct.).
- [7] R. C. Agarwal and C. S. Burrus, "New Recursive Digital Filter Structures Having Very Low Sensitivity and Roundoff Noise," *IEEE Trans. Circuits Sys.*, vol. CAS-22, p. 921 (1975 Dec.).
- [8] L. B. Jackson, "Roundoff Noise Bounds Derived from Coefficient Sensitivities for Digital Filters," *IEEE Trans. Circuits Sys.*, vol. CAS-23, p. 481 (1976 Aug.).
- [9] K. W. Martin and M. T. Sun, "Adaptive Filters Suitable for Real-Time Spectral Analysis," *IEEE Trans. Circuits Sys.*, vol. CAS-33, p. 218 (1986 Feb.).
- [10] A. B. Carlson, *Communication Systems, An Introduction to Signals and Noise in Electrical Communication*, 2nd ed. (McGraw-Hill, New York, 1975).
- [11] P. S. R. Diniz and A. Antoniou, "Low-Sensitivity Digital-Filter Structures Which Are Amenable to Error-Spectrum Shaping," *IEEE Trans. Circuits Sys.*, vol. CAS-32, p. 1000 (1985 Oct.).
- [12] H. T. Nagle, Jr., and V. P. Nelson, "Digital Filter Implementation on 16-Bit Microcomputers," *IEEE Micro Mag.*, p. 23 (1981 Feb.).
- [13] L. Kristiansson, "Jump Phenomenon in Digital Filters," *Electron. Lett.*, vol. 10, p. 14 (1974 Jan. 24).
- [14] Philips Corp., "News Report on Philips 16-Bit Digital to Analog Conversion System," Elcoma Marketing Communications Group, Bldg. BA, Eindhoven, The Netherlands, News Rep. 82906 (1982 Apr. 16).
- [15] J. O'Donnell, "Digital Signal Processing Software—For the Evaluation of Digital Filters," Signix Corp., 19 Pelham Island Road, Wayland, MA 01778, DISPRO version 1.5 (1986).
- [16] W. E. Higgins and D. C. Munson, Jr., "Optimal and Suboptimal Error Spectrum Shaping for Cascade-Form Digital Filters," *IEEE Trans. Circuits Sys.*, vol. CAS-31, p. 429 (1984 May).
- [17] H. A. Spang III, and P. M. Schultheiss, "Reduction of Quantizing Noise by Use of Feedback," *IRE Trans. Commun. Sys.*, vol. CS-10, p. 373 (1962 Dec.). Reprinted in *Waveform Quantization and Coding* (IEEE Press, New York, 1976).
- [18] C. T. Mullis and R. A. Roberts, "An Interpretation of Error Spectrum Shaping in Digital Filters," *IEEE Trans. Acoust., Speech, Signal Proc.*, vol. ASSP-30, p. 1013 (1982 Dec.).
- [19] D. C. Munson, Jr., and B. Liu, "Narrow-Band Recursive Filters with Error Spectrum Shaping," *IEEE Trans. Circuits Sys.*, vol. CAS-28, p. 160 (1981 Feb.).
- [20] W. E. Higgins and D. C. Munson, Jr., "Noise Reduction Strategies for Digital Filters: Error Spectrum Shaping versus the Optimal Linear State-Space Formulation," *IEEE Trans. Acoust., Speech, Signal Proc.*, vol. ASSP-30, p. 963 (1982 Dec.).
- [21] B. C. Rothaar, "A Digital Audio Tone Control," M.S. Thesis, Dept. of Computer Science, University of Utah, Provo (1982 June).
- [22] G. R. Cooper and C. D. McGillem, *Probabilistic Methods of Signal and System Analysis* (Holt, Rinehart & Winston, New York, 1971).
- [23] A. I. Abu-El-Haija and M. M. Al-Ibrahim, "Improving Performance of Digital Sinusoidal Oscillators by Means of Error Feedback Circuits," *IEEE Trans. Circuits Sys.*, vol. CAS-33, p. 373 (1986 Apr.).
- [24] C. W. Barnes, B. N. Tran, and S. H. Leung, "On the Statistics of Fixed-Point Roundoff Error," *IEEE Trans. Acoust., Speech, Signal Proc.*, vol. ASSP-33, p. 595 (1985 June).
- [25] J. M. Halbert and R. A. Belcher, "Selection of Test Signals for DSP-Based Testing of Digital Audio Systems," *J. Audio Eng. Soc. (Engineering Reports)*, vol. 34, pp. 546–555 (1986 July/Aug.).
- [26] G. T. Yan, "New Digital Notch Filter Structures with Low Coefficient Sensitivity," *IEEE Trans. Circuits Sys.*, vol. CAS-31, p. 825 (1984 Sept.).
- [27] S. C. Bass, J. W. Grundmann, and S. E. Belter, "DINAP II: A Digital Filter Analysis Program," School of Elec. Eng., Purdue University, West Lafayette, IN, TR-EE 78-14 (1978 Mar.).
- [28] Texas Instruments, "TMS32010 User's Guide," Digital Signal Processor Products Division, Houston, TX, Doc. SPRU001B (1985).
- [29] T. Claasen and L. Kristiansson, "Improvement of Overflow Behaviour of 2nd-Order Digital Filters by Means of Error Feedback," *Electron. Lett.*, vol. 10, p. 240 (1974 June 13).
- [30] P. M. Ebert, J. E. Mazo, and M. G. Taylor, "Overflow Oscillations in Digital Filters," *Bell Sys. Tech. J.*, vol. 48, p. 2999 (1969 Nov.).
- [31] T. L. Chang, "Suppression of Limit Cycles in Digital Filters Designed with One Truncation Quantizer," *IEEE Trans. Circuits Sys.*, vol. CAS-28, p. 107 (1981 Feb.).
- [32] H. Samuelli and A. N. Willson, Jr., "Nonperiodic Forced Overflow Oscillations in Digital Filters," *IEEE Trans. Circuits Sys.*, vol. CAS-30, p. 709 (1983 Oct.).
- [33] T. Claasen and L. Kristiansson, "Necessary and Sufficient Conditions for the Absence of Overflow Phenomena in a Second-Order Recursive Digital Filter," *IEEE Trans. Acoust., Speech, Signal Proc.*, vol. ASSP-23, p. 509 (1975 Dec.).
- [34] H. Samuelli and A. N. Willson, Jr., "Almost Period P Sequences and the Analysis of Forced Overflow Oscillations in Digital Filters," *IEEE Trans. Circuits Sys.*, vol. CAS-29, p. 510 (1982 Aug.).
- [35] P. K. Sim and K. K. Pang, "Effects of Input-Scaling on the Asymptotic Overflow-Stability Properties of Second-Order Recursive Digital Filters," *IEEE Trans. Circuits Sys.*, vol. CAS-32, p. 1008 (1985 Oct.).
- [36] J. A. Moorer, "The Audio Signal Processor: The Next Step in Digital Audio," *Collected Papers*

from the AES Premiere Conference on Digital Audio (Rye, NY, 1982 June 3-6), pp. 205-215.

[37] B. Liu and T. Kaneko, "Error Analysis of Digital Filters Realized with Floating-Point Arithmetic," *Proc. IEEE*, vol. 57, p. 1735 (1969 Oct.).

[38] J. Vanderkooy and S. P. Lipshitz, "Resolution Below the Least Significant Bit in Digital Systems with Dither," *J. Audio Eng. Soc.*, vol. 32, pp. 106-113 (1984 Mar.).

**APPENDIX 1
TMS320 CODE**

A tested real-time TMS32010 source code follows. This code runs at 25.60576 MHz, ±0.01%, on a TMS320C10-25 for a 50.4-kHz maximum sample rate. There are two filters maximum, executable per sample having double-precision coefficients, single-precision truncation error cancellation, and overflow detection and saturation. This leaves a headroom of 15 instruction cycles out of a maximum of 127 to perform noncritical functions in a fragmented outer executive loop.

```

***** filtering *****
FILTER
  BIOC  FILTER      * JOND - 1987
  IN    XNL,2
  IN    INDEK,3     * GET COEFFICIENT
  LAC   INDEK       * FOR OUTER LOOP
  AND   MSK4
  SACL  INDEK
  LAR   ARO,INDEK
  IN    *,2
  OUT   YNR,2       * 14 INSTR.

***** LEFT CHANNEL *****
* PERFORM RESIDUAL COEFFICIENT CALCULATIONS
ZAC
LT    XNM2L        * Q0
MPY   EA2L        * Q14
LTA   XNM1L
MPY   EA1L
LTA   XNL
MPY   EA0L
LTA   YNM2L
MPY   EB2L
LTA   YNM1L
MPY   EB1L
APAC
* 26 INSTR.
* WANT HEADROOM IN ACCUMULATOR
* DIVIDE BY 2**12 => RESID IS Q10
* END RESIDUAL COEFFICIENT CALCULATIONS

ZAC
LT    EONM2L      *** PERFORM ERROR FEEDBACK CALCULATIONS
MPY   B2L        * STORED SUCH THAT IS NEGATIVE Q12
LTD   EOFNL      * Q12
MPY   B1L        * REALLY EONM1L
APAC          * P REGISTER RESULT IS Q24
SACH  TEMP,4     * INTERMEDIATE ERROR TERM Q12

LAC   TEMP       *** NOW DO SINGLE PRECISION COEFFICIENTS
LT    XNM2L      * Q12 * 35 INSTR.
MPY   A2L        * Q0
LTD   XNM1L      * Q12
MPY   A1L
LTD   XNL
MPY   A0L
LTA   YNM2L
MPY   B2L
LTA   YNM1L
MPY   B1L
APAC
* 46 INSTR.

ADD  RESID,2     *** FORMULATE NEW TRUNCATION ERROR AND OUTPUT SAMPLE
SACH YNL,4       * Q12
SACH OVERFL     * YNL Q0, ACCUMULATOR Q12
SUB  YNL,12     * STORE MSBs OF ACCUMULATOR
SACL EOFNL      * Q12 * 51 INSTR.

* ***** CHECK FOR OUTPUT OVERFLOW *****
LAC  OVERFL,1   * SIGN EXTENSION HERE.
SACH TEMP,4     * 5 MSBs OF Q12 ACC NOW ISOLATED.
SACH SIGN       * BEST BET AS TO OUTPUT SIGN.
ZALS TEMP       * NO SIGN EXTENSION HERE.
XOR  SIGN       * ARE MSBs DIFFERENT?
BZ   IIRL       * RESULT IS 0 IF ALL SAME.
LAC  SIGN
XOR  CEILIN    * MAGIC NUMBER => -32767 (NOT 8)
SACL YNL       * 61 INSTR., SATURATE

IIRL
  
```

```

***** RIGHT CHANNEL (in cascade with LEFT) *****
* PERFORM RESIDUAL COEFFICIENT CALCULATIONS
ZAC
LT    YNM2L        * Q0
MPY   EA2R        * Q14
LTA   YNM1L
MPY   EA1R
LTA   YNL
MPY   EA0R
LTA   YNM2R
MPY   EB2R
LTA   YNR
MPY   EB1R
APAC
* 73 INSTR.
* WANT HEADROOM IN ACCUMULATOR
* DIVIDE BY 2**12 => RESID IS Q10
* END RESIDUAL COEFFICIENT CALCULATIONS

ZAC
LT    EONM2R      *** PERFORM ERROR FEEDBACK CALCULATIONS
MPY   B2R        * STORED SUCH THAT IS NEGATIVE Q12
LTD   EOFNR      * Q12
MPY   B1R        * REALLY EONM1R
APAC          * P REGISTER RESULT IS Q24
SACH  TEMP,4     * INTERMEDIATE ERROR TERM Q12

LAC   TEMP       *** NOW DO SINGLE PRECISION COEFFICIENTS
LT    YNM2L      * Q12 * 82 INSTR.
MPY   A2R        * Q0
LTD   YNM1L      * Q12
MPY   A1R
LTD   YNL
MPY   A0R
LTA   YNM2R
MPY   B2R
LTD   YNR
MPY   B1R
APAC
* 93 INSTR.

ADD  RESID,2     *** FORMULATE NEW TRUNCATION ERROR AND OUTPUT SAMPLE
SACH YNR,4       * Q12
SACH OVERFR     * YNR Q0, ACCUMULATOR Q12
SUB  YNR,12     * STORE MSBs OF ACCUMULATOR
SACL EOFNR      * Q12 * 98 INSTR.

* ***** CHECK FOR OUTPUT OVERFLOW *****
LAC  OVERFR,1   * SIGN EXTENSION HERE.
SACH TEMP,4     * 5 MSBs OF Q12 ACC NOW ISOLATED.
SACH SIGN       * BEST BET AS TO OUTPUT SIGN.
ZALS TEMP       * NO SIGN EXTENSION HERE.
XOR  SIGN       * ARE MSBs DIFFERENT?
BZ   IIR        * RESULT IS 0 IF ALL SAME.
LAC  SIGN
XOR  CEILIN    * MAGIC NUMBER => -32767 (NOT 8)
SACL YNR       * 108 INSTR., SATURATE

* ***** UNIVERSAL RETURN *****
IIR  RET        * 110 INSTR. + 2 (call) => 15 LEFT/127
                          FOR OUTER LOOP
* ***** AT 25.6MHz/4/50.4 KHz => 127 INSTR. MAX *****
  
```

**APPENDIX 2
HOW TO DETECT OUTPUT OVERFLOW**

We store the Q0 output $\hat{y}(n)$ from a 32-bit accumulator which is in Q12 format. The accumulator is Q12 because we chose to perform the single-precision (Q12) coefficient accumulation last. (See the TMS320 code in Appendix 1, right before "CHECK FOR OUTPUT OVERFLOW.") Since $\hat{y}(n)$ is a 16-bit Q0 number, we are left with 4 bits of visible headroom in the accumulator above the sign bit of $\hat{y}(n)$. If no output overflow has occurred, then these "headroom" bits should each be a copy of the sign bit of $\hat{y}(n)$. This suggests a simple way to detect output overflow, then. If any of the 5 MSBs in the accumulator differ from the others at the time that $\hat{y}(n)$ is stored, then output overflow must have occurred. Since we really do not have a clue as to the extent of the overflow, our best guess is that the proper sign of the output should be the same as the MSB of the accumulator. Using that information, we should saturate the output under program control. The 4 headroom bits will allow us to detect output overflow up to 24 dB past unity.

The TMS320 code to perform output overflow detection/saturation is shown in the program. Notice that the use of the XOR operating with the magic number CEILIN (= 32766) on the sign of the Q12 accumulator obviates the need for conditional branching to saturate

at the proper polarity. This detection/saturation routine takes 10 TMS320 instruction cycles per filter stage. This is a significant portion of the code and is almost as long as the basic IIR filter calculation itself. Unfortunately this code is inextricable for high-quality

audio work. The ADSP-2100, in contrast, has 8 hard overflow bits annexing its 32-bit accumulator. It has a single instruction monitoring the sameness of the 9 MSBs. We should expect quite a code reduction using it.

THE AUTHOR



Jon Dattorro was born in Providence, Rhode Island, in 1954. He studied piano from the age of eight with Salvatore Fransosi. From 1972 to 1977 he attended the New England Conservatory of Music where he studied piano and electronic music. While in Boston he studied composition with Pozzi Escot and was a soloist with the Boston Symphony Orchestra for a series of Children's Concerts. He holds a B.S.E.E. degree with honors from the University of Rhode Island (1981) and an M.S.E.E. with emphasis on digital signal processing from Purdue University (1984).

Mr. Dattorro was with Lexicon Inc. until 1987, where he was the principal engineer for the Model 2400, their premier time-compression device. His responsibilities also included the digital filter design for OPUS, their recently announced digital audio workstation.

Mr. Dattorro is now employed by Ensoniq Corp. in Malvern, Pennsylvania, where he has been the project engineer for the Model SQ-80 synthesizer, and where he is currently working on new product research and development. He is a member of the AES and the IEEE.

Asahi Kasei Microsystems Corporation, Atsugi, Japan.

Susumu Hara received a Bachelor's degree in physics from Kyushu University, Fukuota, Japan, in 1982, after which he joined Asahi Chemical Industry Corporation, Tokyo, Japan. From 1982 to 1983 he worked on the research and development of synthetic latex. Since 1983 he has been with Asahi Kasei Microsystems Corporation, Tokyo, Japan, where he is involved in CMOS analog integrated-circuit designs for custom and new-products integrated circuits.

Kaoru Takasuka was born in Hiroshima, Japan, in 1947. He received B.S. and M.S. degrees in instrumentation engineering from Kyushu Institute of Technology, Kyushu, Japan, in 1970 and 1972, respectively.

He joined the Asahi Chemical Industry Corporation, Tokyo, Japan in 1972, where he has been engaged in the research and design of A/D LSIs using CMOS technology. His current research interest is in digital filters and high-resolution A-to-D and D-to-A converters.

He is currently with the LSI Design Center, Asahi Kasei Microsystems Corporation, Atsugi, Japan.

LETTERS TO THE EDITOR

COMMENTS ON "THE IMPLEMENTATION OF RECURSIVE DIGITAL FILTERS OF HIGH-FIDELITY AUDIO"*

I feel comment is required on the above paper by Jon Dattorro.¹ In particular I would like to clarify the relationship between Mr. Dattorro's work as presented and the digital filter technology actually implemented in Lexicon's OPUS digital audio production system. The biographic summary states that "Mr. Dattorro was with Lexicon Inc. until 1987, where he was the principal engineer for the Model 2400, their premier time-compression device. His responsibilities also included the digital filter design for OPUS, their recently announced digital audio workstation." This is somewhat misleading. Although the bulk of the work presented in this paper was compiled by the author as part of the early stages of OPUS system development at Lexicon (and published without Lexicon's review or approval), this work is in no way related to the actual digital equalization implementation as introduced for OPUS at the conventions of the Society for Motion Picture and Television Engineers and the Audio Engineering Society this past fall.

Insofar as this is a potential source of confusion to present and prospective OPUS users we request that this clarification be published at your earliest possible convenience.

LOUIS R. EAGLE
Digital Signal Processing
Lexicon, Inc.
Waltham, MA 02154, USA

Author's Reply²

The above paper¹ does not purport to be a design review of the Lexicon OPUS digital parametric equalization filters. Although the equalization subsystem for OPUS was my responsibility until 1987, I never indicated that the paper represented that which was offered to the customers. The system, OPUS, is only

mentioned in the biography because my design experience is relevant to the paper's topic. I also designed the original OPUS mainframe backplane and I shared the console filter-strip facade and electronics design. I imagine that much has changed during the interim, between my completion of the task in 1987 and the introduction of the equalization subsystem this past November at the 85th Convention of the Audio Engineering Society. Indeed, Mr. Eagle was not a Lexicon employee or consultant prior to my departure. Mr. Eagle has declined to share any technical information since I joined ENSONIQ, so I have no knowledge of the extent to which my work was used.

I would like to note several errors which appeared in the above paper.

1) On page 858, in the last paragraph in Sec. 1.2.2, line 11 should have read: "gain or loss. The . . ."

2) Table 1, page 863, should have read as follows:

Table 1. Error feedback zeros.

K_1	K_2	Region θ
+2	-1	0 Twice
-2	-1	π Twice
0	+1	0 and π
+1	-1	$\pi/3$ Once
-1	-1	$2\pi/3$ Once
+1	0	0 Once
-1	0	π Once
0	-1	$\pi/2$ Once

3) In Sec. 2.5.1, page 864, line 26 should have started with $(N^2\sigma^2)$.

4) The first equation at the top left of page 870 should have read

$$c_i + \text{binary code}(ec_i)/(2^q 2^{q_e}) \approx c_{Fi}$$

5) An asterisk should have appeared in the caption for Fig. 14 preceding 'See TMS . . .'

6) In Sec. 4.1.2, on page 872, lines 21 through 24 in the left-hand column, "worsens when . . . limit-cycle problem," should have been deleted.

7) On page 874 the approximation to the $h(n)$ equation should take the absolute value on the left-hand side.

JON DATTORRO
ENSONIQ Corporation

* Manuscript received 1989 March 6.

¹ J. Dattorro, *J. Audio Eng. Soc.*, vol. 36, pp. 851-878 (1988 Nov.).

measurements and could lead to erroneous results. Experience with practical horns has shown that these reflections are not only common in most designs, but can be very strong. A second point, which may be of importance, is the effect of the cutoff phenomenon associated with many horn shapes. At frequencies near cutoff, gross phase dispersion effects occur, leading to group delays which, in many designs, can be more than a few milliseconds. This dispersion may significantly affect measurements taken using time windows of very short duration, as is necessary for the elimination of mouth reflections.

Some examples of our measurements along with theoretical predictions using a one-parameter, finite exponential element model based on Webster's horn equation are shown in Figs. 1 and 2. These results form part of a technical report which is near completion, an abridged form of which will be submitted shortly for possible publication.

KEITH HOLLAND
FRANK FAHY
PHILIP NEWELL

*Institute of Sound and Vibration Research
The University
Southampton SO9 5NH, UK*

COMMENTS ON "THE IMPLEMENTATION OF RECURSIVE DIGITAL FILTERS FOR HIGH-FIDELITY AUDIO"¹¹

In the above paper¹² Dattorro discussed in Secs. 2.5 and 2.6 a digital filtering scheme he calls "truncation error cancellation." I think the scheme has been explained and interpreted incorrectly, and I will offer my

¹¹ Manuscript received 1989 February 6.

¹² J. Dattorro, *J. Audio Eng. Soc.*, vol. 36, pp. 851-878 (1988 Nov.).

explanation and interpretation and show where the paper is not complete.

The author puts forth two conditions on the structure in Fig. 11 for implementing truncation error cancellation. The first is that the error feedback coefficients K_1 and K_2 are equal to b_1 and b_2 , respectively. The second condition is given in Sec. 2.6, where the author states that the error accumulation must be shifted right (in the case of Fig. 11, shifted right 16 bits) before adding to the signal accumulation. Under these two conditions, the filter can be redrawn as shown in Fig. 11'.

This filter has a 32-bit feedback delay path and a 48-bit accumulator for the feedback multiplication results. This gives the filter a wide dynamic range internally. The filter has two truncation points in it and thus has two error signals to consider. The first, e_1 , is produced by truncating the 48-bit accumulation to 32 bits before the accumulation enters the feedback path. The second, e_2 , is produced by truncating the filter output to 16 bits. Using the author's convention, the filter transfer function is

$$\hat{Y}(z) = X(z) \frac{\sum a_i z^{-i}}{1 - \sum b_i z^{-i}} - E_1(z) \frac{1}{1 - \sum b_i z^{-i}} - E_2(z) \quad (1)$$

This equation is different from Eq. (17) in the paper because the truncation from 48 bits to 32 bits was ignored. It should have been taken into consideration in Sec. 2.6. Paragraph 2, sentence 2, states that the error feedback math of Fig. 11 "must be shifted to the right before combining it with the signal accumulation." Explicitly, the 32-bit accumulator results from the error feedback accumulations must be shifted right 16 bits. But this truncates the 16 least significant bits. This must be considered as another error source.

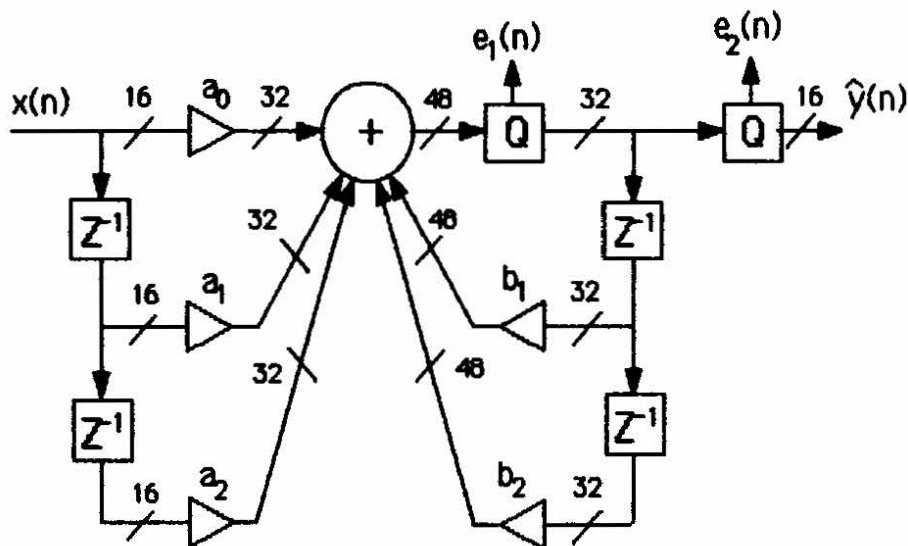


Fig. 11'.

This implementation being considered is most simply described as a filter that uses double-precision arithmetic in the feedback portion. The term "truncation error cancellation" is not proper. Truncation error is always present when using finite word length arithmetic in a recursive filter. The truncation error has simply been reduced by using a longer word length in the recursive section of the filter.

PAUL E. NEYRINCK
San Francisco, CA 94118, USA

Author's Reply¹³

Mr. Neyrinck has pointed out an additional truncation error source in the truncation error cancellation circuit, which was first brought to my attention by Tom Hegg of Lexicon, and later by Richard Cabot of Audio Precision. I consider that error source to have impact only to a second degree on that circuit's performance, and so I chose to ignore any discussion of it, hoping to simplify the mathematics for the reader. Mr. Neyrinck's astute observation, however, has made me regret that omission. It is important to note that this second-degree source of truncation error is nonexistent in the truncation error *feedback* topologies discussed in Secs. 2.3 and 2.4. This is because the error feedback coefficients were trivial multipliers there.

As for my interpretation of the truncation error cancellation circuit, which I believe is correct, it is im-

portant to understand that truncation error cancellation is a special case of error spectrum shaping (ESS), which is truly a truncation error feedback technique and *not* a means of increasing purely numeric precision. The benefits of ESS are established in the frequency domain as a means of reducing truncation noise buildup. Truncation error cancellation may well be viewed as a double-precision implementation. In Sec. 2.5.0, a *subtle* distinction between truncation error cancellation and standard double precision was expressed in terms of the double-precision realization: 1) the signal feedforward paths remain in single precision, and 2) the multiplier inputs are never unsigned. I acknowledged the fact that other engineers had found the noise performance of the truncation error cancellation circuit to be equivalent to that of a standard double-precision implementation.

The two conditions under which the truncation error cancellation will be perfect were set forth in the last paragraph of Sec. 2.6: 1) the signal feedback coefficients b_1 and b_2 must be precisely equal to the error feedback coefficients K_1 and K_2 , respectively, and 2) no rounding (truncation) may be performed in the formulation of the error accumulation.

If the error accumulation must be shifted right before combining it with the signal accumulation (Sec. 2.6, paragraph 2), and if significant digits are lost as a result, then the second criterion is clearly violated.

The material in Sec. 2.6 stands by itself as it is. The analysis that follows would, hypothetically, be added as Sec. 2.6.1, where we would go to this second level of depth into the truncation error assessment. Fig. 11"

¹³ Manuscript received 1989 September 21.

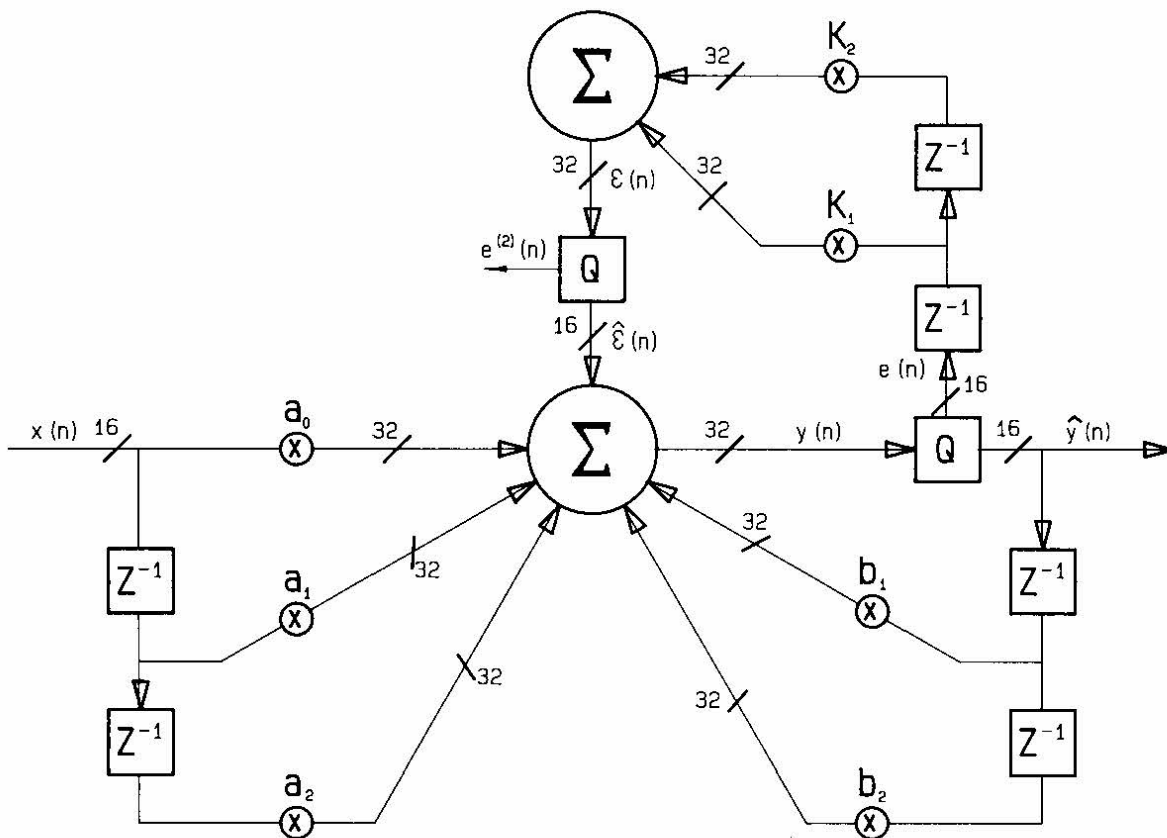


Fig. 11". Second-order truncation error cancellation showing all truncation errors.

here shows the truncation error situation in every detail for a $1\frac{1}{2}$ -bit processor having a 16-by-16-bit multiplier with 32-bit products and accumulations. Referring to the figure,

$$y(n) = \hat{y}(n) + e(n) \quad (2)$$

as before. We define a new variable, ϵ , which is simply the accumulation of the truncation errors,

$$\epsilon(n) = \sum K_i e(n - i) \quad (3)$$

In the text, Eq. (14) uses the accumulation ϵ directly. A new truncation error source $e^{(2)}(n)$ results when the truncation error accumulation ϵ is shifted right (truncated) before it is combined with the signal accumulation,

$$\epsilon(n) = \hat{\epsilon}(n) + e^{(2)}(n) \quad (4)$$

where $|e^{(2)}(n)| \ll |e(n)|$, and where the 2 raised in parentheses connotes the second-degree error. This second-degree truncation actually occurs in the TMS32010 program in Appendix 1 of the paper on the lines that read

SACH TEMP.4 * INTERMEDIATE ERROR TERM Q12

The equation which in reality becomes implemented in the program is then

$$y(n) = \sum a_i x(n - i) + \sum b_i \hat{y}(n - i) + \hat{\epsilon}(n) \quad (5)$$

Contrast this with Eq. (14) in the text. Only the last term is different because we are adding a truncated version of the truncation error accumulation. It is important to realize that both the first and second terms in Eq. (5) here and in Eq. (14) in the text retain full numerical precision because none of the operands need to be truncated prior to multiplication. (Indeed, $\hat{y}(n - i)$ is already truncated and accounted for.) What is the effect of the truncation of ϵ , the truncated truncation error accumulation, the double truncation?

In the frequency domain,

$$\hat{Y}(z) = X(z) \frac{\sum a_i z^{-i}}{1 - \sum b_i z^{-i}} - E(z) - \frac{E^{(2)}(z)}{1 - \sum b_i z^{-i}} \quad (6)$$

which agrees with Mr. Neyrinck's Eq. (1), except for the nomenclature. Contrast Eq. (6) here with Eq. (17)

in the text; Eq. (6) has an additional (third) term. Eq. (6) says that the truncated output spectrum is comprised of the ideal filter transfer minus the nonamplified output truncation minus the truncated truncation error accumulation $E^{(2)}(z)$. $E^{(2)}(z)$ is much smaller in magnitude than $E(z)$ but is amplified by the poles of the system.

Assuming a 16-bit signal and a -90 -dB signal noise floor, then in a 16-bit system $E^{(2)}(z)$ is ideally 16 bits lower in magnitude than $E(z)$. (This would be possible using a processor such as the TMS320C25, where the products can be shifted left out of the multiplier prior to accumulation.) This means that, depending on the implementation, the noise power associated with $E^{(2)}(z)$ can be 186 dB below unity. In a 24-bit system the noise power associated with $E^{(2)}(z)$ is at least 24 bits lower in magnitude than that of $E(z)$, and so 234 dB below unity. The bottom line is that this second-degree source of truncation error will only be problematic if the amplification of it by the system poles raises it above the signal noise floor. For this to happen, the amplification would need to be excessive; certainly, the use of a 24-bit processor makes this much less of a likelihood.

JON DATTORRO
ENSONIQ Corporation
Malvern, PA 19355, USA

COMMENTS ON "THE IMPLEMENTATION OF RECURSIVE DIGITAL FILTERS FOR HIGH-FIDELITY AUDIO"¹⁴

In the above paper¹² Appendix 1 needs to be amended as follows:

Addendum to Appendix 1

The first nine instructions in the TMS320 code are for I/O and must be changed or deleted to suit your purposes.

Note that the operands associated with the LTD instructions must be allocated storage in ascending order so that the program functions properly. For example,

XNL: EQU 7
XNM1L: EQU 8
XNM2L: EQU 9
etc.

JON DATTORRO
ENSONIQ Corporation
Malvern, PA 19355, USA

¹² J. Dattorro, *J. Audio Eng. Soc.*, vol. 36, pp. 851-878 (1988 Nov.).

¹⁴ Manuscript received 1989 October 16.