

On the Robustness of Single-Loop Sigma-Delta Modulation

C. Sinan Güntürk, Jeffrey C. Lagarias, and Vinay A. Vaishampayan, *Member, IEEE*

Abstract—Sigma-delta modulation, a widely used method of analog-to-digital (A/D) signal conversion, is known to be robust to hardware imperfections, i.e., bit streams generated by slightly imprecise hardware components can be decoded comparably well. We formulate a model for robustness and give a rigorous analysis for single-loop sigma-delta modulation applied to constant signals (dc inputs) for N time cycles, with an arbitrary (small enough) initial condition u_0 , and a quantizer that may contain an offset error. The mean-square error (MSE) of any decoding scheme for this quantizer (with u_0 and the offset error known) is bounded below by $\frac{1}{96}N^{-3}$. We also determine the asymptotically best possible MSE as $N \rightarrow \infty$ for perfect decoding when $u_0 = 0$ and $u_0 = \frac{1}{2}$. The robustness result is the upper bound that a triangular linear filter decoder (with both u_0 and the offset error unknown) achieves an MSE of $\frac{40}{3}N^{-3}$. These results establish the known result that the $O(1/N^3)$ decay of the MSE with N is optimal in the single-loop case, under weaker assumptions than previous analyses, and show that a suitable linear decoder is robust against offset error. These results are obtained using methods from number theory and Fourier analysis.

Index Terms—Dynamical systems, oversampled quantization, quantization, robustness, sigma-delta modulation.

I. INTRODUCTION

MODERN techniques of high-accuracy analog-to-digital (A/D) conversion of band-limited signals is based on using single-bit quantization together with oversampling, as a practical alternative to using a multibit quantizer on a sequence sampled at the Nyquist rate. This is true for reasons related to both the quantizer and to oversampling. Single-bit quantizers are preferable to multibit quantizers because they are easier and cheaper to build. Also, single-bit sigma-delta modulation is robust against circuit imperfection, owing to the feedback which compensates for deviations in the quantization thresholds. The deviations in two-level feedback, which occurs in the case of single-bit quantization, only amounts to an offset error combined with an amplification error, while deviations from a multilevel feedback would cause irreversible harmonics. We consider this robustness viewpoint further below. Oversampling facilitates implementations in various ways, including making the job of analog filtering easier, see Candy and Temes [3] for a general discussion.

Manuscript received October 6, 1999; revised July 24, 2000.

C. S. Güntürk is with Program in Applied and Computational Mathematics, Princeton University, Princeton, NJ 08544 USA (e-mail: cgunturk@math.princeton.edu).

J. C. Lagarias and V. A. Vaishampayan are with Information Sciences Research, AT&T Labs (Shannon Laboratory), Florham Park, NJ 07932 USA (e-mail: jcl@research.att.com; vinay@research.att.com).

Communicated by P. A. Chou, Associate Editor for Source Coding.

Publisher Item Identifier S 0018-9448(01)04416-9.

Sigma-delta modulation is a widely used method for A/D conversion, see [3], [20], [22], [25]. It transforms a band-limited signal by oversampling using a single-bit quantizer with feedback, to produce a two-valued signal stream which we call the *coded signal*. This signal stream is then appropriately filtered—usually with a linear filter—to produce a (vector) quantized version of the original signal. This filtering step may be regarded as a form of decoding. The simplest version of sigma-delta modulation is the single-loop version originally introduced in [23], [1], [14], [8], but many more complicated multiloop systems have since been considered.

The sigma-delta modulator is a nonlinear system with feedback and is notoriously difficult to analyze. One of the first rigorous analyses of this system was performed by Gray for constant inputs [8], where he showed that filtering the quantization output sequence with a rectangular window of length N results in a reconstruction error bounded by $O(1/N)$, for all initial conditions and values of the constant input. Gray [9] later used a connection with ergodic theory to show that the mean-squared error (MSE) decays asymptotically as $O(1/N^3)$, as $N \rightarrow \infty$, where N is the number of taps of the filter used in decoding. The notion of MSE used here is taken over a uniform distribution of the value x to be estimated, but also requires a time average¹ of the error signal. Concerning lower bounds, Hein and Zakhor [18] and Hein, Ibrahim, and Zakhor [17] showed that for any decoding scheme for dc input, the quantization error must be at least as large as a constant times $1/N^3$, where the constant depends on the initial value u_0 of the integrator in the sigma-delta modulator at the beginning of the quantization interval. In actual practice, sigma-delta modulators are used for A/D conversion of band-limited signals and not for dc signals. Constant signals are apparently the worst case for such modulators, and engineering practice recommends adding a high-frequency dither signal to make the input vary, cf. Candy and Temes [4, p. 14]. Regarding the performance on more general classes of signals, and the use of nonlinear decoding methods, see [11], [21], [6], [27]–[29]; we discuss this further in the concluding section.

This paper studies the robustness of sigma-delta schemes against certain hardware imperfections. This seems to be one of the main reasons they are used in practice [4, p. 13]. Nonidealities for circuits using a single-bit quantizer can include offset quantizer threshold, offset quantization level, leakage in the integrator, nonunitary integrator gain, nonzero initial state, and random noise. Because of their complexity, robustness of a given

¹See Gray [9, eq. (6.1 ff)]. This time averaged quantity $M\{(\hat{x}_n - \frac{x}{2^b})^2\}$ is equivalent to an average over initial conditions that are the values of u_n occurring in an infinitely long string of samples with fixed input value x . In fact, Gray's paper is mainly concerned with the statistics of the "quantization noise" $\{q(u_n) - u_n, n \geq 1\}$.

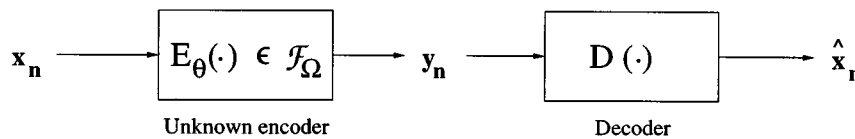


Fig. 1. Robust encoding problem.

system has generally been studied by simulation, see [21] for an example. Feely and Chua [7] present a rigorous study of the effect of integrator leakage on quantization of dc inputs. In this case there can be constant size errors which do not go to zero as the number N of output bits increases, because there is a mode-locking effect.² In this paper, we study the effects of offset quantizer thresholds and of nonzero initial state of the encoder, and rigorously establish that robust decoders do exist for constant signals.

The issue of robustness is one of information theory insofar as it concerns the design of both the encoder and decoder. From the robustness viewpoint, an A/D conversion scheme is a (necessarily nonlinear) encoding operation that produces an output stream of bits describing a signal. One has an ideal system $E(\cdot)$ used for circuit design while the actual hardware produces a system $\tilde{E}(\cdot)$ that approximates the behavior of $E(\cdot)$ but which makes certain “systematic” errors. Whenever feedback is present, the nonlinear nature of the system can potentially lead, after some time, to large discrepancies between the internal states of system $E(\cdot)$ and system $\tilde{E}(\cdot)$ given the same input, no matter how small the “systematic” errors are. Here the “errors” which relate the difference of the actual system $\tilde{E}(\cdot)$ to the ideal system are highly correlated, but can lead to extreme changes in the output bit stream (large Hamming distance). The robust encoding problem is to design an ideal encoder $E(\cdot)$ for which there exists some decoder robust against certain types of hardware imperfections $\tilde{E}(\cdot)$. Given a robust encoder, the robust decoding problem is to design a decoder $D(\cdot)$ which produces an adequate reconstruction of the signal from the digital bit streams produced by any of the (not precisely known) systems $\tilde{E}(\cdot)$. The robust coding problem is quite different from the classical model in information theory of the binary-symmetric channel, in which errors occur randomly. In the classical case, error-correcting coding is introduced in advance of transmission over the channel, but that is not available in this context. Here, the errors are not random but are systematic, caused by the uncertainty in the encoder used.

Besides formulating a framework for the robust encoding problem in the context of sigma-delta modulators, the object of this paper is to provide an analysis of robustness in the “simplest” situation. We give a rigorous performance analysis of single-loop sigma-delta modulation in the case of a constant signal, which includes the effect of nonzero initial state u_0 and of possible offset in the quantizer. As indicated above, the decay of the MSE for such signals is well known to be of order $O(N^{-3})$, under various hypotheses. Here we rigorously demonstrate robustness by showing that a simple linear decoder achieves MSE of order $O(N^{-3})$ with the initial condition and dc offset in the quantizer unknown. To describe the precise results, for constant signals, the sigma-delta modulator can be viewed as a (scalar) quantizer, in which the quantization

assigned to a constant signal x is the sequence of quantized bits produced over the N time periods, which will depend on x , the initial state u_0 of the sigma-delta modulator, and the quantizer used, allowing offset error δ . Our analysis for single-loop sigma-delta modulation with constant input signal x is valid for any fixed small enough initial state u_0 , for N time periods, allowing offset error in the quantizer. More precisely, for offset error δ with $|\delta| < \frac{1}{8}$ we may allow $-\frac{3}{8} < u_0 < \frac{11}{8}$. We give upper and lower bounds for the MSE of this quantizer, assuming only that the dc signal x is uniformly distributed in $[0, 1]$. In particular, no assumption is imposed on the quantizer noise statistics within the N time periods. The lower bound of $\frac{1}{96}N^{-3}$ is valid for the optimal quantizer, which assumes that both u_0 and the offset δ are known to the decoder. The proof uses the same idea as [18] and [17] but sharpens it slightly in obtaining a uniform bound independent of u_0 . We also obtain asymptotically exact bounds for the MSE of optimal quantization for the special cases $u_0 = 0$ and $u_0 = \frac{1}{2}$ as $N \rightarrow \infty$, using detailed facts about Farey fractions. The result for $u_0 = 0$ has the constant $\frac{2}{21} \frac{\zeta(2)}{\zeta(3)} - \frac{1}{16} \cong 0.06782$, which sets a limit on how much the lower bound can be improved. The robustness result is the upper bound, which is $\frac{40}{3}N^{-3}$, for the MSE using the triangular linear filter decoder, which treats both the initial integrator value u_0 and the offset δ as unknown to the decoder. The proof uses Fourier series and an estimate from elementary number theory. These MSE bounds improve on the analysis of Gray [9] in that they do not do any averaging of the signal over input values u_0 , and are valid for each fixed initial value u_0 separately. The specific constants $\frac{1}{96}$ and $\frac{40}{3}$ obtained in the analysis can be further improved, with more detailed estimates, which we do not attempt to do. A small improvement related to the upper bound is indicated at the end of Appendix C.

Compared to a multibit quantizer which can achieve an exponentially small MSE of order $O(2^{-N})$, the sigma-delta quantization output sequence is not efficient. However, this does not mean that the quantizer is very nonuniform. It is well known, e.g., [17], that the number of distinct codewords of length N produced by the first-order sigma-delta modulator on constant signals, for fixed u_0 and δ , is bounded by $O(N^2)$. Hence, an exponential rate-distortion function is still achievable with further coding of the output bit stream. The particular set of $O(N^2)$ admissible output codewords depends on the parameters u_0 and δ . In our model with no random errors, a received codeword contains some information about u_0 and δ . It is this extra information that makes robust decoding possible in this case.

II. PROBLEM FORMULATION

We first formulate the robust encoding problem as the simple block diagram given in Fig. 1. We are given a family of encoders $\mathcal{F}_\Omega = \{E_\theta: \theta \in \Omega\}$ where θ is a (vector) parameter. We may think of these as representing an ideal encoder with the parameter θ measuring the deviation from ideality of a particular hard-

²This provides a reason that constant signals are to be avoided in practice using sigma-delta modulators.

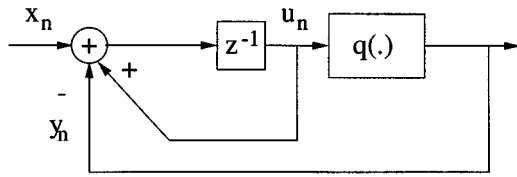


Fig. 2. Encoder of a single-stage sigma-delta modulator.

ware implementation. The parameter θ is not under any control, except to lie in a fixed compact set Ω representing the manufacturing tolerances. The decoder's performance is to be measured in the worst case against all the allowable encoders. We use MSE as a performance measure. This model scheme does not include any source of random errors, just systematic encoding errors embodied in the parameter θ . For the sigma-delta modulator we can take $\theta = (u_0, \delta)$, but as we show below, the behavior of the system really depends only on the single parameter $\tilde{\theta} = u_0 + \delta$. We study the asymptotic behavior as the number N of output bits becomes large. The possible existence of a robust decoder depends on the family of encoders \mathcal{F}_Ω considered. For example, Feely and Chua [7] consider encoders \mathcal{F}_Ω that include leaky integrators, and their results imply that for constant inputs and optimal decoding, the MSE does not go to zero with increasing N , so that a robust decoder does not exist in the asymptotic sense considered here.

We consider systems that use a single-bit quantizer. An *ideal quantizer* has a threshold at 0.5 and reconstruction levels 1 and 0, whose *quantizing map* $q(\cdot)$ is given by

$$q(u) = \begin{cases} 1, & u \geq 1/2, \\ 0, & u < 1/2. \end{cases} \quad (2.1)$$

An *offset quantizer* has a threshold at $0.5 - \delta$, where we assume $-\delta_{\max} \leq \delta \leq \delta_{\max}$ for some δ_{\max} (say, $\delta_{\max} = 0.125$), and reconstruction levels at 1 and 0, hence is the *quantizing map* $q_\delta(\cdot)$ given by

$$q_\delta(u) = \begin{cases} 1, & u \geq 1/2 - \delta \\ 0, & u < 1/2 - \delta. \end{cases} \quad (2.2)$$

A single-loop (or first-order) sigma-delta modulator is illustrated in Fig. 2. The sigma-delta modulator consists of a quantizer in a feedback loop. The behavior of the system with an ideal quantizer is described by

$$u_{n+1} = x_n + u_n - q(u_n), \quad n = 0, 1, \dots \quad (2.3)$$

while with the nonideal quantizer it is

$$\tilde{u}_{n+1} = x_n + \tilde{u}_n - q_\delta(\tilde{u}_n), \quad n = 0, 1, \dots \quad (2.4)$$

The output vector at time n is denoted $Y_n := q_\delta(\tilde{u}_n)$. The following simple fact, observed in [19], simplifies the robust quantizing problem.

Lemma 2.1: Let x_n be a fixed input sequence. The output bit sequence Y_n for the nonideal first-order sigma-delta modulator with initial value \tilde{u}_0 and offset δ is identical to the output bit sequence for the ideal first-order sigma-delta modulator with the modified input value $u_0 := \tilde{u}_0 + \delta$.

Proof: Since $q_\delta(\cdot) = q(\cdot + \delta)$, on setting $u_n := \tilde{u}_n + \delta$, the system (2.4) becomes equivalent to the system (2.3) with the initial condition $u_0 = \tilde{u}_0 + \delta$. \square

Lemma 2.1 shows that studying robustness of a first-order sigma-delta modulator against arbitrary initial value \tilde{u}_0 and offset error reduces to the special case of studying the ideal system (2.3) with arbitrary (unknown) initial condition u_0 . This reduction is special to first-order sigma-delta modulation. In higher order schemes, the initial value u_0 and offset parameter δ are independent sources of error.

Let us suppose the offset error δ satisfies $|\delta| \leq \delta_{\max} < \frac{1}{2}$. The system (2.3) with the ideal quantizer maps the interval $[-1/2, 3/2)$ into itself, so if the original initial condition \tilde{u}_0 for $q_\delta(\cdot)$ satisfies

$$-\frac{1}{2} - \delta \leq \tilde{u}_0 < \frac{3}{2} - \delta \quad (2.5)$$

then this condition is preserved under iteration (2.4). As long as

$$-\frac{1}{2} + \delta_{\max} \leq \tilde{u}_0 < \frac{3}{2} - \delta_{\max}$$

(2.5) is satisfied for all allowable values of $|\delta| \leq \delta_{\max}$, and the subsequent analysis applies.

In view of the Lemma 2.1, in analyzing robustness against offset error, it suffices to treat the case of a first-order sigma-delta modulator with an ideal quantizer, and consider robustness against the choice of initial value u_0 , and this we do in the remainder of the paper. We treat u_0 as given, and average over the input value x , assumed uniformly distributed on $[0, 1]$ and independent of the value u_0 .

The ideal quantizer output at time n is denoted Y_n and is related to input u_n by $Y_n := q(u_n)$. From (2.3), it follows that the quantizer output is given by

$$Y_n := q(u_n) = q\left(u_0 + \sum_{i=0}^{n-1} x_i - \sum_{i=0}^{n-1} q(u_i)\right), \quad n = 0, 1, \dots \quad (2.6)$$

Equation (2.6) defines a map $f_{u_0}^N(\cdot): \mathbf{R}^N \rightarrow \{0, 1\}^N$, with input $(x_0, x_1, \dots, x_{N-1})$ and output (Y_1, Y_2, \dots, Y_N) . As the x_i vary, this map changes value at points where

$$u_0 + \sum_{i=0}^{n-1} x_i - \sum_{i=0}^{n-1} q(u_i) = 1/2, \quad n = 1, \dots, N. \quad (2.7)$$

These points constitute the boundaries of the level sets of $f_{u_0}^N$, in other words, the n -bit quantization bins. Thus, a partition of $[0, 1]^N$ is created. Note that implicit in (2.7), u_n is a nonlinear function of the input x_0, \dots, x_{n-1} . The resulting bins are very irregularly shaped.

Now suppose that we have constant input signal $x_i = x$, for $1 \leq i \leq N$. This corresponds to looking at the intersections of the sets in the partition with the principal diagonal of $[0, 1]^N$ given by $x_0 = x_1 = \dots = x_{N-1}$. This naturally induces a partition of $[0, 1]$, which we refer to as the *effective quantizer*, in order to distinguish it from the binary quantizer in the loop. The

thresholds of the effective quantizer are obtained by determining values of x in $[0, 1]$ that solve

$$u_0 + nx - \left\{ \sum_{i=0}^{n-1} q(u_i) \right\} = 1/2, \quad n = 1, 2, \dots, N. \quad (2.8)$$

Lemma 2.2: For a general u_0 in the range $-1/2 \leq u_0 < 3/2$, the quantization thresholds are given by

$$\mathcal{S}_N(u_0) := \{(j - \beta_0)/n : j = 1, \dots, n; n = 1, 2, \dots, N\} \quad (2.9)$$

where $\beta_0 = \langle u_0 + 1/2 \rangle$, and $\langle \alpha \rangle := \alpha - \lfloor \alpha \rfloor$ denotes the fractional part of α .

We give the derivation in Appendix A. For the special cases $u_0 = 0$ and $u_0 = 1$, the set of thresholds is the set

$$\mathcal{S}_N := \{(2j - 1)/2n, j = 1, \dots, n; n = 1, 2, \dots, N\}.$$

The set \mathcal{S}_N is related to the *Farey series* \mathcal{F}_N of order N , which is the set of fractions k/n in reduced form with $0 < k < n$, $1 \leq n \leq N$, with $\gcd(k, n) = 1$, on the interval $(0, 1)$, arranged in ascending order [13, Ch. 3]. To be more specific, let \mathcal{G}_{2N} be the subset of the Farey series of order $2N$ that has even denominators. Then it follows that $\mathcal{S}_N = \mathcal{G}_{2N}$. On the other hand, the special cases $u_0 = -1/2$ and $u_0 = 1/2$ have $\mathcal{S}_N(-1/2) = \mathcal{S}_N(1/2) = \mathcal{F}_N \setminus \{0\}$, the Farey series itself. The connection of breakpoints for $u_0 = 0$ to the Farey series was observed in Hein and Zakhor [20, p. 24].

Our problem is to give lower and upper bounds for the mean-squared quantization error for fixed u_0 with the constant dc input x assumed drawn from the uniform distribution on $[0, 1]$. For lower bounds, we assume optimal decoding, where u_0 is known to the decoder. The optimal MSE quantizer is described using the map $Q_{\text{opt}}(x; u_0)$, which maps x to the midpoint of the interval J that x lies in, the endpoints of J being successive elements of the thresholds of the effective quantizer with initial value u_0 . The map $Q_{\text{opt}}(x; u_0)$ is the optimal quantizer under our assumption that the quantity x being quantized is uniformly distributed in $[0, 1]$ and independent of u_0 . Our objective is to lower-bound the MSE, given by the integral

$$\text{MSE}_{u_0}(Q_{\text{opt}}) := \int_0^1 (x - Q_{\text{opt}}(x; u_0))^2 dx. \quad (2.10)$$

In the upper bound case, we suppose u_0 is fixed but unknown to the decoder, and consider a decoding algorithm which uses a particular linear filter of length N , the triangular filter. Let $Q_h(x; u_0)$ denote the triangular filtered estimate for x , which depends on x and u_0 . Our objective is to upper-bound

$$\text{MSE}_{u_0}(Q_h) := \int_0^1 (x - Q_h(x; u_0))^2 dx \quad (2.11)$$

for all u_0 . The choice of triangular filter for analysis is explained at the end of Section IV.

III. LOWER BOUND FOR MSE

In this section, we suppose that the initial value u_0 is fixed and known, with $1/2 \leq u_0 < 3/2$. The unit interval $[0, 1]$ is partitioned into subintervals $J = J(\mathbf{Y})$, where $\mathbf{Y} := (Y_1, Y_2, \dots, Y_N)$ is the data available to the decoder. The *optimal decoding algorithm*³ $Q_{\text{opt}}(x; u_0)$ assigns to the quantization data \mathbf{Y} associated to x the midpoint of the interval $J(\mathbf{Y})$. There are at most $\frac{(N+1)(N+2)}{2}$ quantization intervals determined by the values given in (2.9). For $u_0 = 0$, some of the values are repeated,⁴ and the number of distinct values is asymptotic to $\frac{3}{\pi^2}N^2$ as $N \rightarrow \infty$, using [13, Theorem 330], since the points of \mathcal{G}_{2N} can be put in one-to-one correspondence with the Farey sequence \mathcal{F}_N . It is well known that the intervals produced by the Farey sequence \mathcal{F}_{2N} range in size from $\frac{1}{2N}$ down to size $\frac{1}{4N^2}$, see [13, Theorem 35]. The interval $[0, \frac{1}{2N}]$ contributes $\frac{1}{96}N^{-3}$ all by itself to the MSE of the optimal decoding algorithm. We now show that the same bound holds for an arbitrary u_0 .

Theorem 3.1: Suppose that u_0 is fixed with $-1/2 \leq u_0 < 3/2$ and let x be drawn uniformly from $[0, 1]$. Then, single-loop oversampled sigma-delta modulation

$$u_{n+1} = x + u_n - q(u_n), \quad 0 \leq n \leq N - 1$$

with oversampling rate N , using the optimal quantizer Q_{opt} with u_0 known to the quantizer, has MSE

$$\text{MSE}_{u_0}(Q_{\text{opt}}) = \int_0^1 (Q_{\text{opt}}(x; u_0) - x)^2 dx \geq \frac{1}{96}N^{-3}. \quad (3.1)$$

Proof: The value u_0 completely determines the quantization bins. The quantization bin endpoints consist of the points

$$x = \frac{j - \beta_0}{n}, \quad 1 \leq j \leq n, 1 \leq n \leq N \quad (3.2)$$

where $\beta_0 = \langle u_0 + 1/2 \rangle$. We will show that at least one of the open intervals $(0, \frac{1}{2N})$ or $(1 - \frac{1}{2N}, 1)$ contains no quantization threshold. This interval is of length $\frac{1}{2N}$, and since an interval of length $|I|$, contributes $\frac{1}{12}|I|^3$ to the MSE, the contribution of this interval is $\frac{1}{96}N^{-3}$.

Case 1: $0 \leq \beta_0 < 1/2$.

For $1 \leq j \leq n$, and $1 \leq n \leq N$

$$\frac{j - \beta_0}{n} \geq \frac{j - \beta_0}{N} \geq \frac{1 - \beta_0}{N} \geq \frac{1}{2N}$$

hence $(0, \frac{1}{2N})$ contains no quantization threshold.

Case 2: $1/2 \leq \beta_0 < 1$.

For $1 \leq j \leq n$, and $1 \leq n \leq N$

$$\frac{j - \beta_0}{n} \leq \frac{n - \beta_0}{n} \leq 1 - \frac{\beta_0}{N} \leq 1 - \frac{1}{2N}$$

hence $(1 - \frac{1}{2N}, 1)$ contains no quantization threshold. \square

The optimal lower bound in Theorem 3.1 appears to have a constant on the order of five times larger than $\frac{1}{96}$ but seems hard to determine. However, we can show the following exact result.

³The optimality of this algorithm is a consequence of the assumption that the quantity x being quantized is uniformly distributed in $[0, 1]$. When conditioned on the data \mathbf{Y} , the distribution of x is uniform on the quantization interval $J(\mathbf{Y})$.

⁴That is, the values $\frac{j}{n}$ and $\frac{jk}{nk}$ for any $k \geq 2$ determine the same quantization threshold.

Theorem 3.2: Suppose that $u_0 = 0$ or $u_0 = \frac{1}{2}$ and let x be drawn uniformly from $[0, 1]$. Then single-loop oversampled sigma-delta modulation

$$u_{n+1} = x + u_n - q(u_n), \quad 0 \leq n \leq N-1$$

with oversampling rate N , using the optimal quantizer Q_{opt} with u_0 known to the quantizer, has MSE

$$\text{MSE}_{u_0}(Q_{opt}) = \alpha_{u_0} N^{-3} + O(N^{-4} \log N), \quad \text{as } N \rightarrow \infty \quad (3.3)$$

where

$$\alpha_0 := \frac{2}{21} \frac{\zeta(2)}{\zeta(3)} - \frac{1}{16} \cong 0.06782$$

and

$$\alpha_{1/2} := \frac{1}{6} \frac{\zeta(2)}{\zeta(3)} \cong 0.22807.$$

We prove this result in Appendix B. The proof is based on the explicit relation of the set of quantization thresholds in these two cases with the Farey series. Theorem 3.2 sets a limit on how much improvement is possible in the constant $\frac{1}{96}$ appearing in Theorem 3.1, showing that the best constant can be no larger than $\frac{2}{21} \frac{\zeta(2)}{\zeta(3)} - \frac{1}{16}$. Numerical simulations suggest that this bound for $u_0 = 0$ is actually close to the minimum over all initial conditions u_0 , and conceivably it might give the best constant.

IV. UPPER BOUND FOR MSE

In this section, we suppose that u_0 is viewed as fixed with $\frac{1}{2} \leq u_0 < \frac{3}{2}$, but is unknown. The quantization values $(Y_1, Y_2, \dots, Y_{N-1})$ are known to the decoder. For simplicity, we assume that $N = 2M$ is even.

We consider a triangular filter decoder

$$Q_h(x, u_0) := \sum_{n=1}^{2M-1} h_n Y_n \quad (4.1)$$

in which $\{h_n: 1 \leq n \leq 2M-1\}$ are given by

$$h_n = \begin{cases} \frac{1}{M} - \frac{(M-n)}{M^2}, & 1 \leq n \leq M \\ \frac{1}{M} - \frac{(n-M)}{M^2}, & M \leq n \leq 2M-1. \end{cases} \quad (4.2)$$

We give a detailed analysis for the case $N = 2M$ only; for the case $N = 2M + 1$ we may discard the value Y_N and use the above filter on the remaining values.

Theorem 4.1: Suppose that u_0 is fixed with $-\frac{1}{2} \leq u_0 < \frac{3}{2}$, and let x be drawn uniformly in $[0, 1]$. Then single-loop oversampled sigma-delta modulation

$$u_{n+1} = x + u_n - q(u_n), \quad 0 \leq n \leq N-1$$

at oversampling rate $N = 2M$, using quantizer Q_h , has MSE

$$\text{MSE}_{u_0}(Q_h) = \int_0^1 (Q_h(x; u_0) - x)^2 dx \leq \frac{40}{3} N^{-3}. \quad (4.3)$$

The proof uses two number-theoretic lemmas, whose proofs are given in Appendix C. In the following, (m, n) denotes the greatest common divisor of m and n .

Lemma 4.1: For fixed constant β_0 and all positive integers n and m

$$\left| \int_0^1 \langle nx + \beta_0 \rangle \langle mx + \beta_0 \rangle dx - \frac{1}{4} \right| \leq \frac{1}{12} \frac{(n, m)^2}{nm}. \quad (4.4)$$

Lemma 4.2: For all positive integers L

$$\sum_{n=1}^L \sum_{m=1}^L \frac{(n, m)^2}{nm} \leq 5L. \quad (4.5)$$

Proof of Theorem 4.1: Suppose $N = 2M$ is even. We set

$$\epsilon_N(x) := Q_h(x; u_0) - x$$

where $Q_h(x; u_0)$ is the triangular filter decoder

$$Q_h(x; u_0) := \sum_{n=1}^{2M-1} h_n Y_n \quad (4.6)$$

of filter weights (4.2), and $Y_n = q(u_n)$. We have

$$\begin{aligned} \epsilon_N(x) &= \sum_{n=1}^{2M-1} (Y_n - x) h_n \\ &= \sum_{n=1}^{2M-1} (u_n - u_{n+1}) h_n. \end{aligned}$$

Summing this by parts and substituting (A2) from Appendix A yields

$$\begin{aligned} \epsilon_N(x) &= \frac{1}{M^2} (u_1 + \dots + u_M) - \frac{1}{M^2} (u_{M+1} + \dots + u_{2M}) \\ &= \frac{1}{M^2} (w_1 + \dots + w_M) - \frac{1}{M^2} (w_{M+1} + \dots + w_{2M}). \end{aligned} \quad (4.7)$$

Then we have

$$|\epsilon_N(x)| \leq \frac{1}{M^2} \left| \sum_{n=1}^M w_n - \frac{M}{2} \right| + \frac{1}{M^2} \left| \sum_{n=M+1}^{2M} w_n - \frac{M}{2} \right| \quad (4.8)$$

which, upon taking the square yields

$$\begin{aligned} |\epsilon_N(x)|^2 &\leq \frac{2}{M^4} \left(\sum_{n=1}^M w_n - \frac{M}{2} \right)^2 \\ &\quad + \frac{2}{M^4} \left(\sum_{n=M+1}^{2M} w_n - \frac{M}{2} \right)^2. \end{aligned} \quad (4.9)$$

We now consider the MSE

$$\text{MSE}_{u_0}(Q_h) = \int_0^1 |\epsilon_N(x)|^2 dx.$$

Substituting (A4) from Appendix A into (4.9), and integrating, we get

$$\begin{aligned} \text{MSE}_{u_0}(Q_h) &\leq \frac{2}{M^4} \int_0^1 \left\{ \left(\sum_{n=0}^{M-1} \langle nx + \beta_0 \rangle - \frac{M}{2} \right)^2 \right. \\ &\quad \left. + \left(\sum_{n=M}^{2M-1} \langle nx + \beta_0 \rangle - \frac{M}{2} \right)^2 \right\} dx. \end{aligned} \quad (4.10)$$

We expand this expression, substitute $\int_0^1 \langle nx + \beta_0 \rangle dx = 1/2$ for positive n , and rearrange to get

$$\begin{aligned} \text{MSE}_{u_0}(Q_h) &\leq \frac{2}{M^4} \left\{ \left(\beta_0 - \frac{1}{2} \right)^2 \right. \\ &+ \sum_{n=1}^{M-1} \sum_{m=1}^{M-1} \left(\int_0^1 \langle nx + \beta_0 \rangle \langle mx + \beta_0 \rangle dx - \frac{1}{4} \right) \\ &\left. + \sum_{n=M}^{2M-1} \sum_{m=M}^{2M-1} \left(\int_0^1 \langle nx + \beta_0 \rangle \langle mx + \beta_0 \rangle dx - \frac{1}{4} \right) \right\}. \quad (4.11) \end{aligned}$$

Next we apply Lemma 4.1 to (4.11) and replace the term $(\beta_0 - \frac{1}{2})^2$ by its maximum value $1/4$, to obtain

$$\begin{aligned} \text{MSE}_{u_0}(Q_h) &\leq \frac{2}{M^4} \left(\frac{1}{4} + \frac{1}{12} \sum_{n=1}^{M-1} \sum_{m=1}^{M-1} \frac{(n, m)^2}{nm} \right. \\ &\left. + \frac{1}{12} \sum_{n=M}^{2M-1} \sum_{m=M}^{2M-1} \frac{(n, m)^2}{nm} \right). \quad (4.12) \end{aligned}$$

Finally, we conclude our estimate of $\text{MSE}_{u_0}(Q_h)$ by applying Lemma 4.2 with $L = 2M - 1 = N - 1$ to (4.12) (combining the double sums) to obtain

$$\text{MSE}_{u_0}(Q_h) \leq \frac{40}{3} N^{-3} - \frac{16}{3} N^{-4} \quad (4.13)$$

which yields the desired bound. \square

Remarks:

- 1) The triangular filter was used in the analysis because of the identity (4.7) that it yields for the error expression. The “first-order” terms of size $O(1/N)$ get canceled out due to the subtraction, and this was exploited in the estimate (4.8). This is not the case for the rectangular filter; indeed, a telescoping argument gives that the error expression for this filter is equal to $(w_{N+1} - w_1)/N$, which is in general not smaller than $O(1/N)$. Other reasons based on the Fourier transform can be given, see He, Kuhlmann, and Buzo [15, Sec. IV.C].
- 2) Gray [9] determines the optimal linear filter (in the context of [9]), whose general shape is similar to the triangular filter, but differs from it slightly. Hein and Zakhov [19] later constructed an “optimal” nonlinear decoding method.
- 3) The proof of Theorem 4.1 did not determine the best constant for MSE using the triangular filter, and some improvements are possible on the constant $\frac{40}{3}$ by more careful argument. The constant in Lemma 4.2 can be improved slightly.

V. CONCLUSION

This paper gave rigorous upper and lower bounds on the MSE for single-loop sigma-delta modulation applied to constant inputs, where the quantizer may have offset error and an arbitrary

fixed initial value u_0 . It showed that a particular linear decoder is robust against such errors, and attains the optimal MSE within a multiplicative constant. In these special circumstances, a nonlinear decoder can save at most a multiplicative constant in MSE over a linear decoder, and cannot achieve further asymptotic improvement in MSE as $N \rightarrow \infty$. These results show that the redundancy built into oversampled sigma-delta modulation schemes is serving the useful purpose of permitting robust decoding by a linear decoder. It seems likely that for the first-order scheme, robust decoders should exist for a general class of non-constant signals, but that is a more difficult question which we have not addressed.

The methods of this paper exploited certain features specific to first-order sigma-delta modulation (e.g., Lemma 2.1), which do not hold for higher order sigma-delta schemes. However, the general approach of viewing higher order schemes as discrete dynamical systems is a useful one, to which Fourier-analytic methods can be successfully applied, as in Daubechies and DeVore [6] and Güntürk [10], and for these number-theoretic ideas of a more sophisticated type may also be relevant.

It would be of great interest to extend robustness analyses to higher order sigma-delta systems and to obtain bounds valid for general band-limited signals rather than constant signals. For constant signals, it is believed that a k th-order sigma-delta modulation scheme can achieve an MSE that decays like $O(1/N^{2k+1})$ for signals of length N , and that this should be best possible. An upper bound $O(1/N^{2k+1})$ is demonstrated for certain k th-order sigma-delta schemes in He, Kuhlmann, and Buzo [15, Secs. 3 and 4], [16], but their analysis treats the input x as *fixed*, and then lets $N \rightarrow \infty$; the error estimates obtained are not uniform in x (and require that x be irrational), hence their MSE bounds do not apply in the framework⁵ of this paper. Is there a similar upper bound $O(1/N^{2k+1})$ for some k th-order sigma-delta modulation scheme, using the MSE criterion of this paper, and are there such schemes that are robust against offset error in the quantizer, assuming perfect integrators are used? For signals drawn from a wider class of band-limited signals, it is believed that the achievable⁶ MSE should be $O(1/N^{2k+2})$, see Thao [27]. Demonstrating this rigorously, with or without robustness, is apparently an open problem. A rigorous lower bound of $O(1/N^{2k+2})$ for band-limited signals was obtained by Thao [27]. In the case $k = 1$, nonlinear coding schemes that experimentally achieve $O(N^{-4})$ for a class of sinusoidal signals are given in [21], [28], [29].

For general band-limited functions, it is an open problem to rigorously establish whether nonlinear decoding schemes for sigma-delta modulation schemes can offer an asymptotic improvement over linear decoding. If so, another issue would be whether there exists a nonlinear decoding achieving this improvement which is robust. It seems an important general problem to quantify the tradeoff between efficiency and robustness in such schemes, both theoretically and in practice.

⁵The framework of this paper requires integrating their bounds over x , and the resulting integral diverges.

⁶This MSE averages over a larger class of (bounded) signals, so the contribution of constant signals to the MSE is reduced.

APPENDIX A
PROOF OF LEMMA 2.1

Set $J := [-\frac{1}{2}, \frac{3}{2})$. Then, for $u_n \in J$, we have $q(u_n) = \lfloor u_n + \frac{1}{2} \rfloor$. This way, the recursion (2.3) can be rewritten as

$$u_{n+1} = x_n - \frac{1}{2} + \left\langle u_n + \frac{1}{2} \right\rangle. \quad (A1)$$

Suppose $x_n \in [0, 1]$. Then from (A1), $u_n \in J$ implies $u_{n+1} \in J$.

If $x_n = x$ is constant, then in fact, $u_n \in [x - \frac{1}{2}, x + \frac{1}{2})$, for all n . If we now define

$$w_n := u_n - x + \frac{1}{2} \quad (A2)$$

then the iteration for w_n is just rotation on the unit circle

$$w_{n+1} := \langle w_n + x \rangle \quad (A3)$$

hence

$$w_{n+1} = \langle nx + \beta_0 \rangle \quad (A4)$$

with $\beta_0 = \langle u_0 + 1/2 \rangle$, as was originally observed by Gray [8]. Thus, substituting this in (A2) one arrives at the formula

$$u_{n+1} = x - \frac{1}{2} + \langle nx + \beta_0 \rangle. \quad (A5)$$

On the other hand, one also has

$$u_{n+1} = u_0 + (n+1)x - (Y_0 + \dots + Y_n). \quad (A6)$$

Combining the two and using $Y_0 = \lfloor u_0 + 1/2 \rfloor$, it follows that

$$Z_n := Y_1 + \dots + Y_n = \lfloor nx + \beta_0 \rfloor. \quad (A7)$$

Clearly, the sequence (Y_n) is uniquely determined by the sequence (Z_n) , and *vice versa*. From (A7), it follows that Z_n increments by one at the points $\{x = (j - \beta_0)/n; j = 1, \dots, n\}$, starting with 0 at $x = 0$ and ending with n at $x = 1$. So, the N -tuple (Z_1, \dots, Z_N) and, consequently, the quantization codeword (Y_1, \dots, Y_N) attains single and distinct values on each of the subintervals defined by the threshold points $\mathcal{S}_N(u_0)$ given in (2.9), which completes the proof. \square

APPENDIX B
PROOF OF THEOREM 3.2

Lemma 2.2 shows that for $u_0 = 0$ the set of thresholds is

$$\mathcal{G}_{2N} = \left\{ \frac{2j-1}{2n} : 1 \leq j \leq n, 1 \leq n \leq N \right\}$$

while for $u_0 = \frac{1}{2}$ the set of thresholds is

$$\mathcal{F}_N = \left\{ \frac{j}{n} : 1 \leq j \leq n, 1 \leq n \leq N \right\}$$

the Farey series of order N .

We first treat the case $u_0 = \frac{1}{2}$ and estimate

$$\text{MSE}_{1/2}(Q_{\text{opt}}) = \text{MSE}(\mathcal{F}_N) = \sum_{I \in \mathcal{F}_N} \frac{1}{12} |I|^3$$

where $I \in \mathcal{F}_N$ means I is an interval determined by \mathcal{F}_N , and $|I|$ denotes its length. We will show that

$$\text{MSE}(\mathcal{F}_N) = \left(\sum_{q=1}^{\infty} \frac{\phi(q)}{q^3} \right) \frac{1}{6N^3} + O\left(\frac{\log N}{N^4}\right) \quad (B1)$$

in which $\phi(q)$ is Euler's ϕ -function, which counts the number of integers k with $1 \leq k < q$ which are relatively prime to q . The intervals $[0, \frac{1}{N}]$ and $[\frac{N-1}{N}, 1]$ each contribute $\frac{1}{12} \frac{1}{N^3}$. For each $2 \leq q \leq N$, there are $\phi(q)$ fractions $s = \frac{p}{q}$ in lowest terms. Any two adjacent Farey fractions $\frac{p}{q}, \frac{p'}{q'} \in \mathcal{F}_N$ have $q+q' > N$, for otherwise their mediant $\frac{p+p'}{q+q'} \in \mathcal{F}_N$ and falls between $\frac{p}{q}$ and $\frac{p'}{q'}$, contradicting their being adjacent. Thus, if

$$\frac{p^-}{q^-} < \frac{p}{q} < \frac{p^+}{q^+}$$

are the two neighboring fractions in \mathcal{F}_N , $q \leq \frac{N}{2}$ implies that $q^- > N/2$ and $q^+ > N/2$. A well-known property of Farey fractions is that $|pq' - p'q| = 1$ for adjacent fractions, hence the interval $I^- = [\frac{p^-}{q^-}, \frac{p}{q}]$ has length $l^-(s) = \frac{1}{qq^-}$ and $I^+ = [\frac{p}{q}, \frac{p^+}{q^+}]$ has length $l^+(s) = \frac{1}{qq^+}$. All these intervals are disjoint and they contribute

$$S'_N = \sum_{s \in \mathcal{F}_N \cap q \leq N/2} \frac{1}{12} \left(\frac{1}{(qq^-)^3} + \frac{1}{(qq^+)^3} \right)$$

to $\text{MSE}(\mathcal{F}_N)$. Since $N-q < q^+$, $q^- \leq N$ we obtain the bounds

$$\begin{aligned} \frac{1}{6N^3} + \frac{1}{12} \left(\sum_{q=2}^{N/2} \frac{\phi(q)}{q^3} \frac{2}{N^3} \right) \\ \leq S'_N \leq \frac{1}{6N^3} + \frac{1}{12} \sum_{q=2}^{N/2} \left(\frac{\phi(q)}{q^3} \frac{2}{(N-q)^3} \right). \end{aligned}$$

Now

$$\frac{1}{(N-q)^3} - \frac{1}{N^3} = \frac{3N^2q - 3Nq^2 + q^3}{N^3(N-q)^3}$$

so that

$$\sum_{q=2}^{N/2} \frac{\phi(q)}{q^3} \left(\frac{1}{(N-q)^3} - \frac{1}{N^3} \right) = O\left(\frac{\log N}{N^4}\right)$$

on using $(N-q)^3 \geq \frac{1}{8}N^3$, and $\frac{\phi(q)}{q^3} \leq \frac{1}{q^2}$. These bounds imply

$$\begin{aligned} \frac{1}{6N^3} + \frac{1}{12} \left(\sum_{q=2}^{N/2} \frac{\phi(q)}{q^3} \right) \frac{2}{N^3} \\ \leq S'_N \leq \frac{1}{6N^3} + \frac{1}{12} \left(\sum_{q=2}^{N/2} \frac{\phi(q)}{q^3} \right) \frac{2}{N^3} + O\left(\frac{\log N}{N^4}\right). \end{aligned}$$

Since

$$\sum_{q=N/2}^{\infty} \frac{\phi(q)}{q^3} = O\left(\frac{1}{N}\right)$$

we obtain

$$S'_N = \left(\sum_{q=1}^{\infty} \frac{\phi(q)}{q^3} \right) \frac{1}{6N^3} + O\left(\frac{\log N}{N^4}\right). \quad (\text{B2})$$

Let S''_N denote the contribution to $\text{MSE}(\mathcal{F}_N)$ coming from all the remaining intervals. Each of them has endpoints $\frac{p}{q}, \frac{p'}{q'}$ with $q, q' \geq N/2$, hence their length $|I| = \frac{1}{qq'} \leq \frac{4}{N^2}$. There are at most N^2 such intervals, hence

$$S''_N \leq N^2 \left(\frac{4}{N^2} \right)^3 \leq \frac{64}{N^4}.$$

Combining this with (B2) establishes (B1).

We next consider the case $u_0 = 0$ and estimate

$$\text{MSE}_0(Q_{\text{opt}}) = \sum_{I \in \mathcal{G}_{2N}} \frac{1}{12} |I|^3. \quad (\text{B3})$$

We will show that

$$\begin{aligned} \text{MSE}_0(Q_{\text{opt}}) &= \left(\sum_{q=1}^{\infty} \frac{\phi(q)}{q^3} \right) \frac{1}{24N^3} \\ &+ \left(\sum_{m=1}^{\infty} \frac{\phi(2m+1)}{(2m+1)^3} \right) \frac{1}{16N^3} + O\left(\frac{\log N}{N^4}\right). \end{aligned} \quad (\text{B4})$$

To begin, for \mathcal{G}_{2N} we note that $\frac{1}{2}$ is a threshold and the set of thresholds is symmetric about $1/2$, so that

$$\text{MSE}_0(Q_{\text{opt}}) = 2 \sum_{\substack{I \in [0, 1/2] \\ I \in \mathcal{G}_{2N}}} \frac{1}{12} |I|^3. \quad (\text{B5})$$

Next, we rescale the thresholds in $\mathcal{G}_{2N} \cap [0, 1/2]$ by a factor of 2 to obtain the *modified Farey series of order N*

$$\mathcal{F}_N^* := \left\{ \frac{p}{q} : \frac{p}{q} \in \mathcal{F}_N \text{ and } p \equiv 1 \pmod{2} \right\}.$$

Thus,

$$\text{MSE}_0(Q_{\text{opt}}) = \frac{1}{4} \text{MSE}(\mathcal{F}_N^*) \quad (\text{B.6})$$

where

$$\text{MSE}(\mathcal{F}_N^*) := \sum_{I \in \mathcal{F}_N^*} \frac{1}{12} |I|^3. \quad (\text{B7})$$

The modified Farey series \mathcal{F}_N^* is obtained from the Farey series \mathcal{F}_N by removing all points $s = \frac{p}{q}$ with $2|p$, which we call “even” Farey points. Let the neighboring Farey points to the left and right of such an s be $\frac{p^-}{q^-} < \frac{p}{q} < \frac{p^+}{q^+}$, and note that neither of the neighboring points is “even.” The associated Farey intervals $[\frac{p^-}{q^-}, \frac{p}{q}]$ and $[\frac{p}{q}, \frac{p^+}{q^+}]$ have lengths $l^+(s) = \frac{1}{qq^+}$ and $l^-(s) = \frac{1}{qq^-}$. In \mathcal{F}_N^* , these intervals are combined into a single interval of length $l^+(s) + l^-(s)$. The contribution to $\text{MSE}(\mathcal{F}_N^*)$ for this interval is $\frac{1}{12}(l^+(s) + l^-(s))^3$, rather than the two contributions $\frac{1}{12}l^+(s)^3 + \frac{1}{12}l^-(s)^3$ in $\text{MSE}(\mathcal{F}_N)$. Using the identity $(x+y)^3 - (x^3 + y^3) = 3xy(x+y)$, we have

$$\text{MSE}(\mathcal{F}_N^*) = \text{MSE}(\mathcal{F}_N) + T_N, \quad (\text{B8})$$

where

$$T_N = \frac{1}{12} \sum_{\substack{s \in \mathcal{F}_N \\ s \text{ “even”}}} 3l^+(s)l^-(s)(l^+(s) + l^-(s)). \quad (\text{B9})$$

To estimate T_N , we split it into two subsums T'_N and T''_N , where T'_N sums over all “even” $s = \frac{p}{q}$ with $1 < q < \frac{N}{2}$, and T''_N sums over those “even” s with $\frac{N}{2} \leq q \leq N$.

In the first subsum, we have $N - q \leq q^+, q^- \leq N$, hence

$$\begin{aligned} &\frac{1}{4} \left(\sum_{\substack{q=3 \\ q \text{ odd}}}^{N/2} \frac{\frac{1}{2}\phi(q)}{q^3} \frac{2}{N^3} \right) \\ &\leq T'_N \leq \frac{1}{4} \left(\sum_{\substack{q=3 \\ q \text{ odd}}}^{N/2} \frac{\frac{1}{2}\phi(q)}{q^3} \frac{2}{(N-q)^3} \right). \end{aligned} \quad (\text{B10})$$

Here s is “even,” so its denominator q is odd, and there are exactly $\frac{1}{2}\phi(q)$ “even” numerators p . We obtain

$$T'_N = \frac{1}{4} \left(\sum_{m=1}^{\infty} \frac{\phi(2m+1)}{(2m+1)^3} \right) \frac{1}{N^3} + O\left(\frac{\log N}{N^4}\right)$$

in a similar fashion to the estimate (B2).

To bound the second sum T''_N , we note that the mediant $\frac{p^+ + p^-}{q^+ + q^-}$ lies inside $[\frac{p^-}{q^-}, \frac{p^+}{q^+}]$. The only point of \mathcal{F}_N inside this interval is s , so we conclude that

$$q^+ + q^- \geq q \geq \frac{N}{2}.$$

Thus, at least one of q^+ and q^- exceeds $\frac{N}{4}$, hence

$$l^+(s)l^-(s)(l^+(s) + l^-(s)) \leq \frac{1}{(N/2)^3} \frac{1}{N/4} \frac{2}{q_0} \leq \frac{16}{N^4} \frac{1}{q_0^2}$$

where $q_0 = \min(q^+, q^-)$. Each value q_0 can occur at most $2\phi(q_0)$ times, hence

$$T''_N \leq \frac{1}{6} \sum_{q_0=1}^N \frac{\phi(q_0)}{q_0^2} \frac{16}{N^4} \leq O\left(\frac{\log N}{N^4}\right).$$

Combining these estimates gives⁷

$$T_N = T'_N + T''_N = \left(\sum_{m=1}^{\infty} \frac{\phi(2m+1)}{(2m+1)^3} \right) \frac{1}{4N^3} + O\left(\frac{\log N}{N^4}\right). \quad (\text{B11})$$

Then combining (B1), (B6), (B8), and (B11) yields the desired formula (B4).

It remains to obtain explicit formulas for the coefficients in the formulas (B1) and (B4) above, i.e., to determine the constants α_{u_0} . We use the fact [30, p. 6] that

$$\frac{\zeta(s-1)}{\zeta(s)} = \sum_{q=1}^{\infty} \frac{\phi(q)}{q^s}.$$

One easily calculates that

$$\sum_{m=1}^{\infty} \frac{\phi(2m+1)}{(2m+1)^s} = \left(\frac{1-2^{1-s}}{1-2^{-s}} \right) \frac{\zeta(s-1)}{\zeta(s)} - 1$$

⁷The possible “even” point $s = \frac{N-1}{N}$ contributes only $O(\frac{1}{N^4})$ to T_N and goes in the error term.

since $\phi(2m) = \phi(m)$ if m is odd and $\phi(2m) = 2\phi(m)$ if m is even. We take $s = 3$ and obtain

$$\begin{aligned} \text{MSE}_{1/2}(Q_{\text{opt}}) &= \left(\frac{\zeta(2)}{\zeta(3)}\right) \frac{1}{6N^3} + O\left(\frac{\log N}{N^4}\right) \\ &= \frac{\alpha_{1/2}}{N^3} + O\left(\frac{\log N}{N^4}\right) \end{aligned}$$

where $\alpha_{1/2} = \frac{1}{6} \frac{\zeta(2)}{\zeta(3)}$, and

$$\begin{aligned} \text{MSE}_0(Q_{\text{opt}}) &= \left(\frac{\zeta(2)}{\zeta(3)}\right) \frac{1}{24N^3} + \left(\frac{6\zeta(2)}{7\zeta(3)} - 1\right) \frac{1}{16N^3} \\ &\quad + O\left(\frac{\log N}{N^4}\right) \\ &= \frac{\alpha_0}{N^3} + O\left(\frac{\log N}{N^4}\right) \end{aligned}$$

where $\alpha_0 = \frac{2}{21} \frac{\zeta(2)}{\zeta(3)} - \frac{1}{16}$. \square

APPENDIX C

PROOF OF LEMMAS 4.1 AND 4.2

Proof of Lemma 4.1: We shall prove the lemma by establishing the formula

$$\int_0^1 \langle nx + \beta_0 \rangle \langle mx + \beta_0 \rangle dx = \frac{1}{4} + \frac{1}{12} \frac{(n, m)^2}{nm} \phi_{n, m} \quad (\text{C1})$$

where $|\phi_{n, m}| \leq 1$. Denote the expression on the left-hand side by $c_{n, m}$. We substitute $nx + \beta_0$ and $mx + \beta_0$ for x in the Fourier series expansion

$$\langle x \rangle = \sum_{k=-\infty}^{\infty} a_k e^{2\pi i k x}$$

where $a_k = (-2\pi i k)^{-1}$ for $k \neq 0$ and $a_0 = 1/2$. This Fourier series is only conditionally convergent, and is to be interpreted as the limit as $N \rightarrow \infty$ of the sum taken from $-N$ to N . However, its partial sums are uniformly bounded

$$\left| \sum_{|k| \leq N} a_k e^{2\pi i k x} \right| \leq C, \quad \text{for all } x \text{ and all } N. \quad (\text{C2})$$

In fact, one can take $C = \frac{1}{2} + \frac{5}{4\pi}$ [24, Example 4, p. 22]. Hence, using the bounded convergence theorem, one can change the order of integration and double sum to obtain

$$c_{n, m} = \sum_{k \in \mathbf{Z}} \sum_{l \in \mathbf{Z}} a_k \bar{a}_l e^{2\pi i (k-l)\beta_0} \int_0^1 e^{2\pi i (kn-lm)x} dx. \quad (\text{C3})$$

Summing up (C3) over the nonzero indexes k, l given by $kn = lm$ and straightforward manipulations result in

$$c_{n, m} = \frac{1}{4} + \frac{1}{4\pi^2} \frac{(n, m)^2}{nm} \sum_{d \neq 0} \frac{1}{d^2} e^{2\pi i d \beta_0 (m-n)/(n, m)} \quad (\text{C4})$$

which, in turn, implies using $\sum_{d \neq 0} \frac{1}{d^2} = \frac{\pi^2}{3}$

$$c_{n, m} = \frac{1}{4} + \frac{1}{12} \frac{(n, m)^2}{nm} \phi_{n, m} \quad (\text{C5})$$

for some $|\phi_{n, m}| \leq 1$. \square

Remark: Using the formula

$$\sum_{d=1}^{\infty} \frac{1}{d^2} \cos d\theta = \frac{1}{4} (\theta - \pi)^2 - \frac{1}{12} \pi^2, \quad 0 \leq \theta \leq 2\pi \quad (\text{C.6})$$

the exact value of $\phi_{n, m}$ is easily found to be

$$\phi_{n, m} = \frac{3}{2} \left(2 \left\langle \beta_0 \frac{m-n}{(n, m)} \right\rangle - 1 \right)^2 - \frac{1}{2}.$$

Proof of Lemma 4.2: We have

$$\begin{aligned} \sum_{n=1}^L \sum_{m=1}^L \frac{(n, m)^2}{nm} &= \sum_{d=1}^L \sum_{\substack{1 \leq n, m \leq L \\ (n, m) = d}} \frac{(n, m)^2}{nm} \\ &\leq \sum_{d=1}^L \sum_{j=1}^{\lfloor L/d \rfloor} \sum_{k=1}^{\lfloor L/d \rfloor} \frac{1}{jk} \\ &= \sum_{d=1}^L \left(\sum_{j=1}^{\lfloor L/d \rfloor} \frac{1}{j} \right)^2 \\ &\leq \sum_{d=1}^L \left(1 + \log \frac{L}{d} \right)^2. \quad (\text{C7}) \end{aligned}$$

However, this last expression is bounded by

$$(1 + \log L)^2 + \int_1^L (1 + \log(L/y))^2 dy = 5L - 4 - 2 \log L$$

which proves (4.5). \square

Remark: The constant 5 appearing in (4.5) can be improved to $\gamma^2 + 5\gamma/2 + 7/3 \cong 4.1$ by using the inequality

$$\sum_{j=1}^L \frac{1}{j} \leq \gamma + \log L + \frac{1}{2L} \quad (\text{C8})$$

where γ is the Euler–Mascheroni constant defined by

$$\gamma = \lim_{N \rightarrow \infty} \left(\sum_{j=1}^N \frac{1}{j} - \log N \right) = 0.5772 \dots \quad (\text{C9})$$

Numerical experiments suggest the optimal constant to be ~ 3 .

ACKNOWLEDGMENT

The authors wish to thank Z. Cvetković and especially N. T. Thao for helpful comments and references.

REFERENCES

- [1] J. C. Candy, "A use of double integration in sigma delta modulation," *IEEE Trans. Commun.*, vol. COM-33, pp. 249–258, Mar. 1985.
- [2] J. C. Candy and O. J. Benjamin, "The structure of quantization noise from sigma-delta modulation," *IEEE Trans. Commun.*, vol. COM-34, pp. 1316–1323, Sept. 1981.
- [3] J. C. Candy and G. C. Temes, *Oversampling Delta-Sigma Data Converters: Theory, Design, and Simulation*. New York: IEEE Press, 1992.
- [4] —, "Oversampling methods for A/D and D/A conversion," in *Oversampling Delta-Sigma Data Converters: Theory, Design, and Simulation*. New York: IEEE Press, 1992, pp. 1–29.
- [5] W. Chou, P. W. Wong, and R. M. Gray, "Multistage sigma-delta modulation," *IEEE Trans. Inform. Theory*, vol. 35, pp. 784–796, July 1989.
- [6] I. Daubechies and R. DeVore, "Reconstructing a bandlimited function from very coarsely quantized data. A family of stable sigma-delta modulators of arbitrary order," paper. In preparation.

- [7] O. Feely and L. O. Chua, "The effect of integrator leak in $\Sigma - \Delta$ modulation," *IEEE Trans. Circuits Syst.*, vol. 38, pp. 1293–1305, Nov. 1991.
- [8] R. M. Gray, "Oversampled sigma-delta modulation," *IEEE Trans. Commun.*, vol. COM-35, pp. 481–489, May 1987.
- [9] —, "Spectral analysis of quantization noise in a single-loop sigma-delta modulator with dc input," *IEEE Trans. Commun.*, vol. 37, pp. 588–599, June 1989.
- [10] S. Güntürk, "Reconstructing a bandlimited function from very coarsely quantized data: II. Improving the error estimate for first order sigma-delta modulators," paper. In preparation.
- [11] —, "Improved error estimates for first order sigma-delta modulation," in *Sampling Theory and Applications, SampTA'99*, Loen, Norway, Aug. 1999.
- [12] —, "Number theoretical error estimates in a quantization scheme for bandlimited signals," in *Unusual Applications of Number Theory*, ser. DIMACS Volume. Providence, RI: Amer. Math. Soc., to be published.
- [13] G. H. Hardy and E. M. Wright, *An Introduction to the Theory of Numbers*. Oxford, U.K.: Oxford Univ. Press, 1979.
- [14] T. Hayashi, Y. Inabe, K. Uchimure, and T. Kimura, "A multistage delta-sigma modulator without double integration loop," in *ISSCC Dig. Tech. Papers*, Feb. 1986, pp. 182–183.
- [15] N. He, F. Kuhlmann, and A. Buzo, "Double-loop sigma-delta modulation with dc input," *IEEE Trans. Commun.*, vol. 38, pp. 487–495, 1990.
- [16] —, "Multiloop sigma-delta quantization," *IEEE Trans. Inform. Theory*, vol. 38, pp. 1015–1027, May 1992.
- [17] S. Hein, K. Ibrahim, and A. Zakhor, "New properties of sigma-delta modulators with dc inputs," *IEEE Trans. Commun.*, vol. 40, pp. 1375–1387, Aug. 1992.
- [18] S. Hein and A. Zakhor, "Lower bounds on the MSE of the single and double loop sigma delta modulators," in *Proc. IEEE Int. Symp. Circuits and Systems*, vol. 3, 1990, pp. 1751–1755.
- [19] —, "Optimal decoding for data acquisition applications of Sigma Delta modulators," *IEEE Trans. Signal Processing*, vol. 41, pp. 602–616, Feb. 1993.
- [20] —, *Sigma Delta Modulators: Nonlinear Decoding Algorithms and Stability Analysis*. Dordrecht, The Netherlands: Kluwer, 1993.
- [21] —, "Reconstruction of oversampled band-limited signals from $\Sigma\Delta$ encoded binary sequences," *IEEE Trans. Signal Processing*, vol. 42, pp. 799–811, Apr. 1994.
- [22] D. F. Hoschele, Jr., *Analog-to-Digital and Digital-to-Analog Conversion Techniques*. New York: Wiley, 1994.
- [23] H. Inose and Y. Yasuda, "A unity bit coding method by negative feedback," in *Proc. IEEE*, vol. 51, Nov. 1963, pp. 1524–1535.
- [24] N. Y. Katznelson, *An Introduction to Harmonic Analysis*. New York: Wiley, 1968. Reprinted by Dover.
- [25] S. R. Norsworthy, R. Schrier, and G. C. Temes, Eds., *Delta-Sigma Data Converters: Theory, Design and Simulation*. Piscataway, NJ: IEEE Press, 1996.
- [26] S. K. Tewksbury and R. W. Hallock, "Oversampled, linear predictive and noise shaping coders of order $N > 1$," *IEEE Trans. Circuits Syst.*, vol. CAS-25, pp. 436–447, July 1978.
- [27] N. T. Thao, "Vector quantization analysis of sigma-delta modulation," *IEEE Trans. Signal Processing*, vol. 44, pp. 808–817, Apr. 1996.
- [28] —, "Deterministic analysis of sigma-delta modulation for linear and nonlinear signal reconstruction," *Int. J. Circuit Theory Appl., Special Issue on Delta-Sigma Modulators and Noise Shaping*, pp. 369–392, Sep.–Oct. 1997.
- [29] N. T. Thao and M. Vetterli, "Deterministic analysis of oversampled A/D conversion and decoding improvement based on consistent estimates," *IEEE Trans. Signal Processing*, vol. 42, pp. 519–531, Mar. 1994.
- [30] E. C. Titchmarsh and D. R. Heath-Brown, *The Theory of the Riemann Zeta-Function*, 2nd ed. Oxford, U.K.: Oxford Univ. Press, 1986.