

# Sigma-Delta ( $\Sigma\Delta$ ) Quantization and Finite Frames

John J. Benedetto, Alexander M. Powell, and Özgür Yılmaz

**Abstract**—The  $K$ -level Sigma-Delta ( $\Sigma\Delta$ ) scheme with step size  $\delta$  is introduced as a technique for quantizing finite frame expansions for  $\mathbb{R}^d$ . Error estimates for various quantized frame expansions are derived, and, in particular, it is shown that  $\Sigma\Delta$  quantization of a unit-norm finite frame expansion in  $\mathbb{R}^d$  achieves approximation error

$$\|x - \tilde{x}\| \leq \frac{\delta d}{2N} (\sigma(F, p) + 1)$$

where  $N$  is the frame size, and the frame variation  $\sigma(F, p)$  is a quantity which reflects the dependence of the  $\Sigma\Delta$  scheme on the frame. Here  $\|\cdot\|$  is the  $d$ -dimensional Euclidean 2-norm. Lower bounds and refined upper bounds are derived for certain specific cases. As a direct consequence of these error bounds one is able to bound the mean squared error (MSE) by an order of  $1/N^2$ . When dealing with sufficiently redundant frame expansions, this represents a significant improvement over classical pulse-code modulation (PCM) quantization, which only has MSE of order  $1/N$  under certain nonrigorous statistical assumptions.  $\Sigma\Delta$  also achieves the optimal MSE order for PCM with consistent reconstruction.

**Index Terms**—Finite frames, Sigma-Delta quantization.

## I. INTRODUCTION

**I**N signal processing, one of the primary goals is to obtain a digital representation of the signal of interest that is suitable for storage, transmission, and recovery. In general, the first step toward this objective is finding an atomic decomposition of the signal. More precisely, one expands a given signal  $x$  over an at most countable dictionary  $\{e_n\}_{n \in \Lambda}$  such that

$$x = \sum_{n \in \Lambda} c_n e_n \quad (1)$$

where  $c_n$  are real or complex numbers. Such an expansion is said to be *redundant* if the choice of  $c_n$  in (1) is not unique.

Although (1) is a discrete representation, it is certainly not “digital” since the coefficient sequence  $\{c_n\}_{n \in \Lambda}$  is real or complex valued. Therefore, a second step is needed to reduce the continuous range of this sequence to a discrete, and preferably

Manuscript received August 13, 2004; revised December 16, 2005. This work was supported by the National Science Foundation under DMS Grant 0139759, NSF Grant 0219233, ONR Grant N000140210398, and by the Natural Science and Engineering Research Council of Canada. The material in this paper was presented in part at the IEEE International Conference on Acoustics, Speech and Signal Processing, Montreal, QC, Canada, May 2004.

J. J. Benedetto is with the Department of Mathematics, University of Maryland, College Park, MD 20742 USA (e-mail: jjb@math.umd.edu).

A. M. Powell is with the Department of Mathematics, Vanderbilt University, Nashville, TN 37240 USA (e-mail: alexander.m.powell@vanderbilt.edu).

Ö. Yılmaz is with the Department of Mathematics, University of British Columbia, Vancouver, B.C. V6T 1Z2, Canada (e-mail: oyilmaz@math.ubc.ca).

Communicated by S. A. Savari, Associate Editor for Source Coding.

Digital Object Identifier 10.1109/TIT.2006.872849

finite, set. This second step is called *quantization*. A *quantizer* maps each expansion (1) to an element of

$$\Gamma_{\mathcal{A}} = \left\{ \sum_{n \in \Lambda} q_n e_n : q_n \in \mathcal{A} \right\}$$

where the *quantization alphabet*  $\mathcal{A}$  is a given discrete, and preferably finite, set. The performance of a quantizer is reflected in the approximation error  $\|x - \tilde{x}\|$ , where  $\|\cdot\|$  is a suitable norm, and

$$\tilde{x} = \sum_{n \in \Lambda} q_n e_n \quad (2)$$

is the quantized expansion.

The process of reconstructing  $\tilde{x}$  in (2) from the quantized coefficients,  $q_n, n \in \Lambda$ , is called *linear reconstruction*. More general approaches to quantization, such as consistent reconstruction, e.g., [1]–[3], use nonlinear reconstruction, but unless otherwise mentioned, we shall focus on quantization using linear reconstruction, as in (2).

A simple example of quantization, for a given expansion (1), is to choose  $q_n$  to be the closest point in the alphabet  $\mathcal{A}$  to  $c_n$ . Quantizers defined this way are usually called pulse-code modulation (PCM) algorithms. If  $\{e_n\}_{n \in \Lambda}$  is an orthonormal basis for a Hilbert space  $H$ , then PCM algorithms provide the optimal quantizers in that they minimize  $\|x - \tilde{x}\|$  for every  $x$  in  $H$ , where  $\|\cdot\|$  is the Hilbert space norm. On the other hand, PCM can perform poorly if the set  $\{e_n\}_{n \in \Lambda}$  is redundant. We shall discuss this in detail in Section II-B.

In this paper, we shall examine the quantization of redundant real finite atomic decompositions (1) for  $\mathbb{R}^d$ . The signal,  $x$ , and dictionary elements  $e_n, n \in \Lambda$  are elements of  $\mathbb{R}^d$ , the index set  $\Lambda$  is finite, and the coefficients  $c_n, n \in \Lambda$  are real numbers. For some existing approaches to this problem see [2]–[7].

### A. Frames, Redundancy, and Robustness

In various applications it is convenient to assume that the signals of interest are elements of a Hilbert space  $H$ , e.g.,  $H = L^2(\mathbb{R}^d)$ , or  $H = \mathbb{R}^d$ , or  $H$  is a space of band-limited functions. In this case, one can consider more structured dictionaries, such as frames. Frames are a special type of dictionary which can be used to give stable redundant decompositions (1). Redundant frames are used in signal processing because they yield representations that are robust under

- additive noise [8] (in the setting of Gabor and wavelet frames for  $L^2(\mathbb{R})$ ), [9] (in the setting of oversampled band-limited functions), and [10] (in the setting of tight Gabor frames);
- quantization [11]–[13] (in the setting of oversampled band-limited functions), [14] (in the setting of tight

Gabor frames), and [2] (in the setting of finite frames for  $\mathbb{R}^d$ ); and

- partial data loss [4], [15] (in the setting of finite frames for  $\mathbb{R}^d$ ).

Although redundant frame expansions use a larger than necessary bit budget to represent a signal (and hence are not preferred for storage purposes where data compression is the main goal), the robustness properties listed above make them ideal for applications where data is to be transferred over noisy channels, or to be quantized very coarsely. In particular, in the case of Sigma-Delta ( $\Sigma\Delta$ ) modulation of oversampled band-limited functions  $x$ , one has very good reconstruction using only 1-bit quantized values of the frame coefficients [11], [16], [17]. Moreover, the resulting approximation  $\tilde{x}$  is robust under quantizer imperfections as well as bit-flips [11]–[13].

Another example, where redundant frames are used, this time to ensure robust transmission, can be found in the works of Goyal, Kovačević, Kelner, and Vetterli [4], [18], cf., [19]. They propose using finite tight frames for  $\mathbb{R}^d$  to transmit data over erasure channels; these are channels over which transmission errors can be modeled in terms of the loss (erasure) of certain packets of data. They show that the redundancy of these frames can be used to “mitigate the effect of the losses in packet-based communication systems,” [20], cf., [21]. Further, the use of finite frames has been proposed for generalized multiple description coding [22], [15], for multiple-antenna code design [23], and for solving modified quantum detection problems [24]. Thus, finite frames for  $\mathbb{R}^d$  or  $\mathbb{C}^d$  are emerging as a natural mathematical model and tool for many applications.

## B. Redundancy and Quantization

A key property of frames is that greater frame redundancy translates into more robust frame expansions. For example, given a unit-norm tight frame for  $\mathbb{R}^d$  with frame bound  $A$ , any transmission error that is caused by the erasure of  $e$  coefficients can be corrected as long as  $e < A$  [4]. In other words, increasing the frame bound, i.e., the redundancy of the frame, makes the representation more robust with respect to erasures. However, increasing redundancy also increases the number of coefficients to be transmitted. If one has a fixed bit budget, a consequence is that one has fewer bits to spend for each coefficient and hence needs to be more resourceful in *how* one allocates the available bits.

- When dealing with PCM, using linear reconstruction, for finite frame expansions in  $\mathbb{R}^d$ , a long-standing analysis with certain assumptions on quantization “noise” bounds the resulting mean-square approximation error by  $C_1\delta^2/A$  where  $C_1$  is a constant,  $A$  is the frame bound, and  $\delta$  is the quantizer step size [25], see Section II-B.
- On the other hand, for 1-bit first-order  $\Sigma\Delta$  quantization of oversampled *band-limited* functions, the approximation error is bounded by  $C_2/A$  pointwise [11], and the mean-square approximation error is bounded by  $C_3/A^3$  [16], [17].

Thus, if we momentarily “compare apples with oranges,” we see that  $\Sigma\Delta$  quantization algorithms for band-limited functions

utilize the redundancy of the expansion more efficiently than PCM algorithms for  $\mathbb{R}^d$ .

## C. Overview of the Paper and Main Results

Section II discusses necessary background material. In particular, Section II-A gives basic definitions and theorems from frame theory, and Section II-B presents basic error estimates for PCM quantization of finite frame expansions for  $\mathbb{R}^d$ .

In Section III, we introduce the  $K$ -level  $\Sigma\Delta$  scheme with step size  $\delta$  as a new technique for quantizing unit-norm finite frame expansions. A main theme of this paper is to show that the  $\Sigma\Delta$  scheme outperforms linearly reconstructed PCM quantization of finite frame expansions. In Section III-A, we introduce the notion of *frame variation*,  $\sigma(F, p)$ , as a quantity which reflects the dependence of the  $\Sigma\Delta$  scheme’s performance on properties of the frame. Section III-B uses the frame variation,  $\sigma(F, p)$ , to derive basic approximation error estimates for the  $\Sigma\Delta$  scheme. For example, we prove that if  $F$  is a unit-norm tight frame for  $\mathbb{R}^d$  of cardinality  $N \geq d$ , then the  $K$ -level  $\Sigma\Delta$  scheme with quantization step size  $\delta$  gives approximation error

$$\|x - \tilde{x}\| \leq \frac{\delta d}{2N}(\sigma(F, p) + 1)$$

where  $\|\cdot\|$  is the  $d$ -dimensional Euclidean 2-norm.

Section IV is devoted primarily to examples. We give examples of infinite families of frames with uniformly bounded frame variation. We compare the error bounds of Section III with the numerically observed error for these families of frames. Since  $\Sigma\Delta$  schemes are iterative, they require one to choose a *quantization order*,  $p$ , in which frame coefficients are given as input to the scheme. We present a numerical example which shows the importance of carefully choosing the quantization order to ensure good approximations.

In Section V, we derive lower bounds and refined upper bounds for the  $\Sigma\Delta$  scheme. This partially explains properties of the approximation error which are experimentally observed in Section IV. In particular, we show that in certain situations, if the frame size  $N$  is even, then one has the improved approximation error bound  $\|x - \tilde{x}\| \leq C_1(\log N)/N^{5/4}$  for an  $x$ -dependent constant  $C_1$ . On the other hand, if  $N$  is odd, we prove the lower bound  $C_2/N \leq \|x - \tilde{x}\|$  for an  $x$ -dependent constant  $C_2$ . In both cases,  $\|\cdot\|$  is the Euclidean norm.

In Section VI, we compare the mean square (approximation) error (MSE) for the  $\Sigma\Delta$  scheme with PCM using linear reconstruction. If we have a harmonic frame for  $\mathbb{R}^d$  of cardinality  $N \geq d$ , then we show that the MSE for the  $\Sigma\Delta$  scheme is bounded by an order of  $1/N^2$ , whereas the standard MSE estimates for PCM are only of order  $1/N$ . Thus, if the frame redundancy is large enough then  $\Sigma\Delta$  outperforms PCM. We present numerical examples to illustrate this. This also shows that  $\Sigma\Delta$  quantization achieves the optimal approximation order for PCM with consistent reconstruction.

## II. BACKGROUND

### A. Frame Theory

The theory of frames in harmonic analysis is due to Duffin and Schaeffer [26]. Modern expositions on frame theory can be

found in [8], [27], [28]. In the following definitions,  $\Lambda$  is an at most countable index set.

*Definition II.1:* A collection  $F = \{e_n\}_{n \in \Lambda}$  in a Hilbert space  $H$  is a frame for  $H$  if there exists  $0 < A \leq B < \infty$  such that

$$\forall x \in H, \quad A\|x\|^2 \leq \sum_{n \in \Lambda} |\langle x, e_n \rangle|^2 \leq B\|x\|^2.$$

The constants  $A$  and  $B$  are called the frame bounds.

A frame is *tight* if  $A = B$ . An important remark is that the size of the frame bound of a *unit-norm* tight frame, i.e., a tight frame with  $\|e_n\| = 1$  for all  $n$ , “measures” the redundancy of the system. For example, if  $A = 1$ , then a unit-norm tight frame must be an orthonormal basis and there is no redundancy, see [8, Proposition 3.2.1]. The larger the frame bound  $A \geq 1$  is, the more redundant a unit-norm tight frame is.

*Definition II.2:* Let  $\{e_n\}_{n \in \Lambda}$  be a frame for a Hilbert space  $H$  with frame bounds  $A$  and  $B$ . The analysis operator

$$L : H \rightarrow l^2(\Lambda)$$

is defined by  $(Lx)_k = \langle x, e_k \rangle$ . The operator  $S = L^*L$  is called the frame operator, and it satisfies

$$AI \leq S \leq BI \quad (3)$$

where  $I$  is the identity operator on  $H$ . The inverse of  $S$ ,  $S^{-1}$ , is called the dual frame operator, and it satisfies

$$B^{-1}I \leq S^{-1} \leq A^{-1}I. \quad (4)$$

The following theorem illustrates why frames can be useful in signal processing.

*Theorem II.3:* Let  $\{e_n\}_{n \in \Lambda}$  be a frame for  $H$  with frame bounds  $A$  and  $B$ , and let  $S$  be the corresponding frame operator. Then  $\{S^{-1}e_n\}_{n \in \Lambda}$  is a frame for  $H$  with frame bounds  $B^{-1}$  and  $A^{-1}$ . Further, for all  $x \in H$

$$x = \sum_{n \in \Lambda} \langle x, e_n \rangle (S^{-1}e_n) \quad (5)$$

$$= \sum_{n \in \Lambda} \langle x, (S^{-1}e_n) \rangle e_n \quad (6)$$

with unconditional convergence of both sums.

The atomic decompositions in (5) and (6) are the first step toward a digital representation. If the frame is tight with frame bound  $A$ , then both frame expansions are equivalent and we have

$$\forall x \in H, \quad x = A^{-1} \sum_{n \in \Lambda} \langle x, e_n \rangle e_n. \quad (7)$$

When the Hilbert space  $H$  is  $\mathbb{R}^d$  or  $\mathbb{C}^d$ , and  $\Lambda$  is finite, the frame is referred to as a *finite frame* for  $H$ . In this case, it is straightforward to check if a set of vectors is a tight frame. Given a set of  $N$  vectors,  $\{v_n\}_{n=1}^N$ , in  $\mathbb{R}^d$  or  $\mathbb{C}^d$ , define the associated matrix  $L$  to be the  $N \times d$  matrix whose rows are the  $\bar{v}_n$ . The following lemma can be found in [29].

*Lemma II.4:* A set of vectors  $\{v_n\}_{n=1}^N$  in  $H = \mathbb{R}^d$  or  $H = \mathbb{C}^d$  is a tight frame with frame bound  $A$  if and only if its associated matrix  $L$  satisfies  $S = L^*L = AI_d$ , where  $L^*$  is the conjugate transpose of  $L$ , and  $I_d$  is the  $d \times d$  identity matrix.

For the important case of finite unit-norm tight frames for  $\mathbb{R}^d$  and  $\mathbb{C}^d$ , the frame constant  $A$  is  $N/d$ , where  $N$  is the cardinality of the frame [30], [4], [2], [29].

### B. PCM Algorithms and Bennett's White Noise Assumption

Let  $K \in \mathbb{N}$  and  $\delta > 0$ . Given the *midrise* quantization alphabet

$$\mathcal{A}_K^\delta = \{(-K + 1/2)\delta, (-K + 3/2)\delta, \dots, (-1/2)\delta, (1/2)\delta, \dots, (K - 1/2)\delta\}$$

consisting of  $2K$  elements, we define the  $2K$ -level midrise uniform scalar quantizer with stepsize  $\delta$  by

$$Q(u) = \arg \min_{q \in \mathcal{A}_K^\delta} |u - q|. \quad (8)$$

Thus,  $Q(u)$  is the element of the alphabet which is closest to  $u$ . If two elements of  $\mathcal{A}_K^\delta$  are equally close to  $u$  then let  $Q(u)$  be the larger of these two elements, i.e., the one larger than  $u$ . For simplicity, we only consider midrise quantizers, although many of our results are valid more generally.

Let  $\{e_n\}_{n=1}^N$  be a unit-norm tight frame for  $\mathbb{R}^d$ , so that each  $x \in \mathbb{R}^d$  has the frame expansion

$$x = \frac{d}{N} \sum_{n=1}^N x_n e_n, \quad x_n = \langle x, e_n \rangle.$$

The  $2K$ -level PCM quantizer with step size  $\delta$  replaces each  $x_n \in \mathbb{R}$  with  $q_n = Q(x_n)$ . Thus, PCM quantizes  $x$  by

$$\tilde{x} = \frac{d}{N} \sum_{n=1}^N q_n e_n.$$

It is easy to see that

$$\forall n, \quad |x_n| < (K - 1/2)\delta \implies |x_n - q_n| \leq \delta/2. \quad (9)$$

PCM quantization as defined above assumes linear reconstruction from the PCM quantized coefficients  $q_n$ . We very briefly address the nonlinear technique of consistent reconstruction in Section VI.

Fix  $\delta > 0$  and  $K \in \mathbb{N}$ , and let  $\|\cdot\|$  be the  $d$ -dimensional Euclidean 2-norm. Let  $x \in \mathbb{R}^d$  and let  $\tilde{x}$  be the quantized expansion which is obtained by using  $2K$ -level PCM quantization with step size  $\delta$ . If  $\|x\| < (K - 1/2)\delta$  then by (9), the approximation error  $\|x - \tilde{x}\|$  satisfies

$$\begin{aligned} \|x - \tilde{x}\| &= \frac{d}{N} \left\| \sum_{n=1}^N (x_n - q_n) e_n \right\| \\ &\leq \left(\frac{\delta}{2}\right) \left(\frac{d}{N}\right) \sum_{n=1}^N \|e_n\| = \left(\frac{d}{2}\right) \delta. \end{aligned} \quad (10)$$

This error estimate does not utilize the redundancy of the frame. A common way to improve the estimate (10) is to make statistical assumptions on the differences  $x_n - q_n$ , e.g., [25], [2].

*Example II.5 (Bennett's White Noise Assumption):* Let  $\{e_n\}_{n=1}^N$  be a unit-norm tight frame for  $\mathbb{R}^d$  with frame bound

$A = N/d$ , let  $x \in \mathbb{R}^d$ , and let  $x_n$ ,  $q_n$ , and  $\tilde{x}$  be defined as above. Since the “pointwise” estimate (10) is unsatisfactory, a different idea is to derive better error estimates which hold “on average” under certain statistical assumptions.

Let  $\nu$  be a probability measure on  $\mathbb{R}^d$ , and consider the random variables  $\eta_n = x_n - q_n$  with the probability distribution  $\mu_n$  induced by  $\nu$  as follows. For  $\mathcal{B} \subseteq \mathbb{R}$  measurable

$$\mu_n(\mathcal{B}) = \nu(\{x \in \mathbb{R}^d : \langle x, e_n \rangle - q_n(x) \in \mathcal{B}\}).$$

The classical approach dating back to Bennett, [25], is to assume that the quantization noise sequence  $\{\eta_n\}_{n=1}^N$  is a sequence of independent and identically distributed random variables with mean 0 and variance  $\sigma^2$ . In other words,  $\mu_n = \mu$  for  $n = 1, \dots, N$ , and the joint probability distribution  $\mu_{1, \dots, N}$  of  $\{\eta_n\}_{n=1}^N$  is given by  $\mu_{1, \dots, N} = \mu^N$ . We shall refer to this statistical assumption on  $\{\eta_n\}_{n=1}^N$  as Bennett’s white noise assumption.

It was shown in [2], that under Bennett’s white noise assumption, the mean square (approximation) error (MSE) satisfies

$$\text{MSE}_{\text{PCM}} = E(\|x - \tilde{x}\|^2) = \frac{d\sigma^2}{A} = \frac{d^2\sigma^2}{N} \quad (11)$$

where the expectation  $E(\|x - \tilde{x}\|^2)$  is defined by

$$E(\|x - \tilde{x}\|^2) = \int_{\mathbb{R}^d} \|x - \tilde{x}\|^2 d\nu(x)$$

which can be rewritten using Bennett’s white noise assumption as

$$E(\|x - \tilde{x}\|^2) = \int_{\mathbb{R}^N} \frac{d}{N} \left\| \sum_{n=1}^N \eta_n e_n \right\|^2 d\mu^N(\eta_1, \dots, \eta_N).$$

Since we are considering PCM quantization with step size  $\delta$ , and in view of (9), it is quite natural to assume that each  $\eta_n$  is a uniform random variable on  $[-\frac{\delta}{2}, \frac{\delta}{2}]$ , and hence has mean 0, and variance  $\sigma^2 = \delta^2/12$ , [31]. In this case one has

$$\text{MSE}_{\text{PCM}} = \frac{d\delta^2}{12A} = \frac{d^2\delta^2}{12N}. \quad (12)$$

Although (12) in Example II.5 represents an improvement over (10) it is still unsatisfactory for the following reasons.

- a) The MSE bound (12) only gives information about the average quantizer performance.
- b) As one increases the redundancy of the expansion, i.e., as the frame bound  $A$  increases, the MSE given in (12) decreases only as  $1/A$ , i.e., the redundancy of the expansion is not utilized very efficiently.
- c) Equation (12) is computed under assumptions that are not rigorous and, at least in certain cases, not true. See [32] for an extensive discussion and a partial deterministic analysis of the quantizer error sequence  $\{\eta_n\}$ . In Example II.6, we show an elementary setting where Bennett’s white noise assumption does not hold for PCM quantization of finite frame expansions.

Since a redundant frame has more elements than are necessary to span the signal space, there will be interdependencies between the frame elements, and hence between the frame coefficients. It is intuitively reasonable to expect that this redundancy and interdependency may violate the independence part

of Bennett’s white noise assumption. The following example makes this intuition precise.

*Example II.6 (Shortcomings of the Noise Assumption):* Consider the unit-norm tight frame for  $\mathbb{R}^2$ , with frame bound  $A = N/2$ , given by

$$\{e_n\}_{n=1}^N, \quad e_n = (\cos(2\pi n/N), \sin(2\pi n/N))$$

where  $N > 2$  is assumed to be even. Given  $x \in \mathbb{R}^2$ , and let  $x_n = \langle x, e_n \rangle$  be the corresponding  $n$ th-frame coefficient. It is easy to see that since  $N$  is even

$$\forall n, \quad e_n = -e_{n+N/2}$$

and hence,

$$\forall n, \quad x_n = -x_{n+N/2}.$$

Next, note that for almost every  $x \in \mathbb{R}^2$  (with respect to Lebesgue measure) one has

$$\forall n, \quad x_n \notin \delta\mathbb{Z}.$$

By the definition of the PCM scheme, this implies that for almost every  $x \in \mathbb{R}^2$  with  $\|x\| < 1$  one has  $q_n = -q_{n+N/2}$ , and hence,  $\eta_n = -\eta_{n+N/2}$ . This means that the quantization noise sequence  $\{\eta_n\}$  is not independent and that Bennett’s white noise assumption is violated. Thus, the MSE predicted by (12) will not be attained in this case. One can rectify the situation by applying the white noise assumption to the frame that is generated by deleting half of the points to ensure that only one of  $e_n$  and  $e_{n+N/2}$  is left in the resulting set.

In addition to the limitations of PCM mentioned above, it is also well known that PCM has poor robustness properties in the band-limited setting, [11]. In view of these shortcomings of PCM quantization, we seek an alternate quantization scheme which is well suited to utilizing frame redundancy. We shall show that the class of Sigma-Delta ( $\Sigma\Delta$ ) schemes perform exceedingly well when used to quantize redundant finite frame expansions.

### III. $\Sigma\Delta$ ALGORITHMS FOR FRAMES FOR $\mathbb{R}^d$

Sigma-Delta ( $\Sigma\Delta$ ) quantizers are widely implemented to quantize oversampled band-limited functions [33], [11]. Here, we define the fundamental  $\Sigma\Delta$  algorithm with the aim of using it to quantize finite frame expansions, see [34].

*Definition III.1:* Given  $K \in \mathbb{N}$ ,  $\delta > 0$ , and the corresponding midrise quantization alphabet and  $2K$ -level midrise uniform scalar quantizer  $Q$  with stepsize  $\delta$ . Let  $\{x_n\}_{n=1}^N \subseteq \mathbb{R}^d$ , and let  $p$  be a permutation of  $\{1, 2, \dots, N\}$ . The associated first-order  $\Sigma\Delta$  quantizer is defined by the iteration

$$\begin{aligned} u_n &= u_{n-1} + x_{p(n)} - q_n \\ q_n &= Q(u_{n-1} + x_{p(n)}) \end{aligned} \quad (13)$$

for  $n = 1, \dots, N$ , where  $u_0 = 0$ . The first-order  $\Sigma\Delta$  quantizer produces the quantized sequence  $\{q_n\}_{n=1}^N$ , and an auxiliary sequence  $\{u_n\}_{n=0}^N$  of state variables.

Thus, a first-order  $\Sigma\Delta$  quantizer is a  $2K$ -level first-order  $\Sigma\Delta$  quantizer with step size  $\delta$  if it is defined by means of (13), where  $Q$  is defined by (8). We shall refer to the permutation  $p$  as the *quantization order*. For simplicity, we have defined the initial state variable to be  $u_0 = 0$ , but it is straightforward to also consider nonzero initial conditions  $|u_0| < \delta/2$ .

The following proposition, cf., [11], shows that the first-order  $\Sigma\Delta$  quantizer is *stable*, i.e., the auxiliary sequence  $\{u_n\}$  defined by (13) is uniformly bounded if the input sequence  $\{x_n\}$  is appropriately uniformly bounded.

*Proposition III.2:* Let  $K$  be a positive integer, let  $\delta > 0$ , and consider the  $\Sigma\Delta$  system defined by (13) and (8). If

$$\forall n = 1, \dots, N, \quad |x_n| \leq (K - 1/2)\delta$$

then

$$\forall n = 1, \dots, N, \quad |u_n| \leq \delta/2.$$

*Proof:* Without loss of generality assume that  $p$  is the identity permutation. The proof proceeds by induction. The base case,  $|u_0| \leq \delta/2$ , holds by assumption. Next, suppose that  $|u_{j-1}| \leq \delta/2$ . This implies that  $|u_{j-1} + x_j| \leq K\delta$ , and hence, by (13) and the definition of  $Q$

$$|u_j| = |u_{j-1} + x_j - Q(u_{j-1} + x_j)| \leq \delta/2.$$

#### A. Frame Variation

Let  $F = \{e_n\}_{n=1}^N$  be a finite frame for  $\mathbb{R}^d$  and let

$$x = \sum_{n=1}^N x_n S^{-1} e_n, \quad x_n = \langle x, e_n \rangle \quad (14)$$

be the corresponding frame expansion for some  $x \in \mathbb{R}^d$ . Since this frame expansion is a finite sum, the representation is independent of the order of summation. In fact, recall that by Theorem II.3, any frame expansion in a Hilbert space converges unconditionally.

Although frame expansions do not depend on the ordering of the frame, the  $\Sigma\Delta$  scheme in Definition III.1 is iterative in nature, and *does* depend strongly on the order in which the frame coefficients are quantized. In particular, we shall show that changing the order in which frame coefficients are quantized can have a drastic effect on the performance of the  $\Sigma\Delta$  scheme. This, of course, stands in stark contrast to PCM schemes which are order independent. The  $\Sigma\Delta$  scheme (13) takes advantage of the fact that there are “interdependencies” between the frame elements in a redundant frame expansion. This is a main underlying reason why  $\Sigma\Delta$  schemes outperform PCM schemes, which quantize frame coefficients without considering any “interdependencies.”

We now introduce the notion of *frame variation*. This will play an important role in our error estimates and it directly reflects the importance of carefully choosing the order in which frame coefficients are quantized.

*Definition III.3:* Let  $F = \{e_n\}_{n=1}^N$  be a finite frame for  $\mathbb{R}^d$ , and let  $p$  be a permutation of  $\{1, 2, \dots, N\}$ . We define the variation of the frame  $F$  with respect to  $p$  as

$$\sigma(F, p) := \sum_{n=1}^{N-1} \|e_{p(n)} - e_{p(n+1)}\|. \quad (15)$$

Roughly speaking, if a frame  $F$  has low variation with respect to  $p$ , then the frame elements will not oscillate too much in that ordering and there is more “interdependence” between successive frame elements.

#### B. Basic Error Estimates

We now derive error estimates for the  $\Sigma\Delta$  scheme in Definition III.1 for  $K \in \mathbb{N}$  and  $\delta > 0$ . Given a frame  $F = \{e_n\}_{n=1}^N$  for  $\mathbb{R}^d$ , a permutation  $p$  of  $\{1, 2, \dots, N\}$ , and  $x \in \mathbb{R}^d$ , we shall calculate how well the quantized expansion

$$\tilde{x} = \sum_{n=1}^N q_n S^{-1} e_{p(n)}$$

approximates the frame expansion

$$x = \sum_{n=1}^N x_{p(n)} S^{-1} e_{p(n)}, \quad x_{p(n)} = \langle x, e_{p(n)} \rangle.$$

Here,  $\{q_n\}_{n=1}^N$  is the quantized sequence which is calculated using Definition III.1 and the sequence of frame coefficients,  $\{x_{p(n)}\}_{n=1}^N$ . We now state our first result on the *approximation error*  $\|x - \tilde{x}\|$ . We shall use  $\|\cdot\|_{\text{op}}$  to denote the operator norm induced by the Euclidean norm  $\|\cdot\|$  for  $\mathbb{R}^d$ .

*Theorem III.4:* Given the  $\Sigma\Delta$  scheme of Definition III.1, let  $F = \{e_n\}_{n=1}^N$  be a finite unit-norm frame for  $\mathbb{R}^d$ , let  $p$  be a permutation of  $\{1, 2, \dots, N\}$ , and let  $x \in \mathbb{R}^d$  satisfy  $\|x\| \leq (K - 1/2)\delta$ . The approximation error  $\|x - \tilde{x}\|$  satisfies

$$\|x - \tilde{x}\| \leq \|S^{-1}\|_{\text{op}} \left( \sigma(F, p) \frac{\delta}{2} + |u_N| \right) \quad (16)$$

where  $S^{-1}$  is the inverse frame operator for  $F$ .

*Proof:*

$$\begin{aligned} x - \tilde{x} &= \sum_{n=1}^N (x_{p(n)} - q_n) S^{-1} e_{p(n)} \\ &= \sum_{n=1}^N (u_n - u_{n-1}) S^{-1} e_{p(n)} \\ &= \sum_{n=1}^{N-1} u_n S^{-1} (e_{p(n)} - e_{p(n+1)}) \\ &\quad + u_N S^{-1} e_{p(N)} - u_0 S^{-1} e_{p(1)}. \end{aligned} \quad (17)$$

Since  $\|x\| \leq (K - 1/2)\delta$ , it follows that

$$\forall 1 \leq n \leq N, \quad |x_n| = |\langle x, e_n \rangle| \leq (K - 1/2)\delta.$$

Thus, by Proposition III.2

$$\begin{aligned} \|x - \tilde{x}\| &\leq \sum_{n=1}^N \frac{\delta}{2} \|S^{-1}\|_{\text{op}} \|e_{p(n)} - e_{p(n+1)}\| \\ &\quad + |u_N| \|S^{-1}\|_{\text{op}} + |u_0| \|S^{-1}\|_{\text{op}} \\ &= \|S^{-1}\|_{\text{op}} \left( \sigma(F, p) \frac{\delta}{2} + |u_N| \right). \end{aligned}$$

It is important to observe that the estimate (16) consists of two fundamentally different error terms. The  $N - 1$  term summation in (17) contributes the main error term and the remaining items are *boundary terms* resulting from the summation by parts. An analogous computation in the setting of  $\Sigma\Delta$  quantization of band-limited signals, e.g., [11], gives a

similar main error term. However, the boundary terms are a special consequence of the finite length encoding here, and are not present in the band-limited setting. For an intuitive explanation of the boundary terms note that the  $\Sigma\Delta$  scheme is a type of error diffusion algorithm. The finite nature of our problem means that one can only diffuse and compensate for errors finitely many times, leading to the possibility of a final uncompensated residual error, i.e., boundary terms.

Theorem III.4 is stated for general unit-norm frames, but since finite tight frames are especially desirable in applications, we shall restrict the remainder of our discussion to tight frames. The utility of finite unit-norm tight frames is apparent in the simple reconstruction formula (7). Note that general finite unit-norm frames for  $\mathbb{R}^d$  are elementary to construct. In fact, any finite subset of  $\mathbb{R}^d$  is a frame for its span. However, the construction and characterization of finite unit-norm tight frames is much more interesting due to the additional algebraic constraints involved [30].

*Corollary III.5:* Given the  $\Sigma\Delta$  scheme of Definition III.1, let  $F = \{e_n\}_{n=1}^N$  be a unit-norm tight frame for  $\mathbb{R}^d$  with frame bound  $A = N/d$ , let  $p$  be a permutation of  $\{1, 2, \dots, N\}$ , and let  $x \in \mathbb{R}^d$  satisfy  $\|x\| \leq (K - 1/2)\delta$ . The approximation error  $\|x - \tilde{x}\|$  satisfies

$$\|x - \tilde{x}\| \leq \frac{d}{N} \left( \sigma(F, p) \frac{\delta}{2} + |u_N| \right).$$

*Proof:* As discussed in Section II-A, a tight frame  $F = \{e_n\}_{n=1}^N$  for  $\mathbb{R}^d$  has frame bound  $A = N/d$ , and, by (4) and Lemma II.4

$$\|S^{-1}\|_{\text{op}} = \left\| \frac{d}{N} I \right\|_{\text{op}} = d/N.$$

The result now follows from Theorem III.4.

*Corollary III.6:* Given the  $\Sigma\Delta$  scheme of Definition III.1, let  $F = \{e_n\}_{n=1}^N$  be a unit-norm tight frame for  $\mathbb{R}^d$  with frame bound  $A = N/d$ , let  $p$  be a permutation of  $\{1, 2, \dots, N\}$ , and let  $x \in \mathbb{R}^d$  satisfy  $\|x\| \leq (K - 1/2)\delta$ . The approximation error  $\|x - \tilde{x}\|$  satisfies

$$\|x - \tilde{x}\| \leq \frac{\delta d}{2N} (\sigma(F, p) + 1).$$

*Proof:* Apply Corollary III.5 and Proposition III.2.

The approximation error estimate in Theorem III.4 can be made more precise if one has more information about the final state variable,  $|u_N|$ . It is somewhat surprising that for *zero sum* frames the value of  $|u_N|$  is completely determined by whether the frame has an even or odd number of elements.

*Theorem III.7:* Given the  $\Sigma\Delta$  scheme of Definition III.1. Let  $F = \{e_n\}_{n=1}^N$  be a unit-norm tight frame for  $\mathbb{R}^d$  with frame bound  $A = N/d$ , and assume that  $F$  satisfies the zero sum condition

$$\sum_{n=1}^N e_n = 0. \tag{18}$$

Then

$$|u_N| = \begin{cases} 0, & \text{if } N \text{ even} \\ \delta/2, & \text{if } N \text{ odd.} \end{cases} \tag{19}$$

*Proof:* Note that (13) implies

$$u_N = u_0 + \sum_{n=1}^N x_n - \sum_{n=1}^N q_n = \sum_{n=1}^N x_n - \sum_{n=1}^N q_n. \tag{20}$$

Next, (18) implies

$$\sum_{n=1}^N x_n = \sum_{n=1}^N \langle x, e_n \rangle = \left\langle x, \sum_{n=1}^N e_n \right\rangle = 0. \tag{21}$$

By the definition of the midrise quantization alphabet  $\mathcal{A}_K^\delta$ , each  $q_n$  is an odd integer multiple of  $\delta/2$ .

If  $N$  is even, it follows that  $\sum_{n=1}^N q_n$  is an integer multiple of  $\delta$ . Thus, by (20) and (21),  $u_N$  is an integer multiple of  $\delta$ . However,  $|u_N| \leq \delta/2$  by Proposition III.2, so that we have  $u_N = 0$ .

If  $N$  is odd, it follows that  $\sum_{n=1}^N q_n$  is an odd integer multiple of  $\delta/2$ . Thus, by (20) and (21),  $u_N$  is an odd integer multiple of  $\delta/2$ . However,  $|u_N| \leq \delta/2$  by Proposition III.2, so that we have  $|u_N| = \delta/2$ .

*Corollary III.8:* Given the  $\Sigma\Delta$  scheme of Definition III.1, let  $F = \{e_n\}_{n=1}^N$  be a unit-norm tight frame for  $\mathbb{R}^d$  with frame bound  $A = N/d$ , and assume that  $F$  satisfies the zero sum condition (18). Let  $p$  be a permutation of  $\{1, \dots, N\}$  and let  $x \in \mathbb{R}^d$  satisfy  $\|x\| \leq (K - 1/2)\delta$ . Then the approximation error  $\|x - \tilde{x}\|$  satisfies

$$\|x - \tilde{x}\| \leq \begin{cases} \frac{\delta d}{2N} \sigma(F, p), & \text{if } N \text{ even} \\ \frac{\delta d}{2N} (\sigma(F, p) + 1), & \text{if } N \text{ odd.} \end{cases} \tag{22}$$

*Proof:* Apply Corollary III.5, Theorem III.7, and Proposition III.2.

Corollary III.8 shows that as a consequence of Theorem III.7, one has smaller constants in the error estimate for  $\|x - \tilde{x}\|$  when the frame size  $N$  is even. Theorem III.7 makes an even bigger difference when deriving refined estimates as in Section V, or when dealing with higher order  $\Sigma\Delta$  schemes [35].

#### IV. FAMILIES OF FRAMES WITH BOUNDED VARIATION

One way to obtain arbitrarily small approximation error  $\|x - \tilde{x}\|$  using the estimates of the previous section is simply to fix a frame and decrease the quantizer step size  $\delta$  toward zero, while letting  $K = \lceil 1/\delta \rceil$ , where  $\lceil \cdot \rceil$  is the *ceiling function*. By Corollary III.6, as  $\delta$  goes to 0, the approximation error goes to zero. However, this approach is not always be desirable. For example, in analog-to-digital (A/D) conversion of band-limited signals, it can be quite costly to build quantizers with very high resolution, i.e., small  $\delta$  and large  $K$ , e.g., [11]. Instead, many practical applications involving A/D and digital-to-analog (D/A) converters make use of oversampling, i.e., redundant frames, and use low-resolution quantizers, e.g., [36]. To be able to adopt this type of approach for the quantization of finite frame expansions, it is important to be able to construct families of frames with uniformly bounded frame variation.

Let us begin by making the observation that if  $F = \{e_n\}_{n=1}^N$  is a finite unit-norm frame and  $p$  is any permutation of  $\{1, 2, \dots, N\}$  then  $\sigma(F, p) \leq 2(N - 1)$ . However, this bound

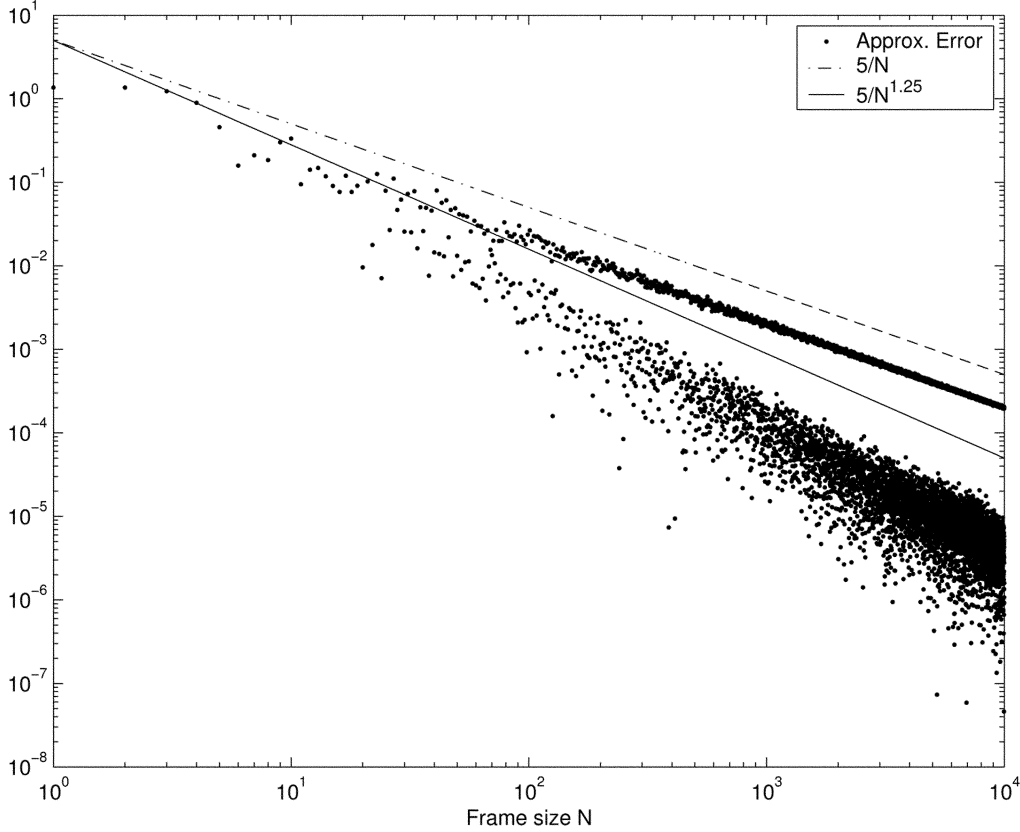


Fig. 1. The frame coefficients of  $x = (1/\pi, \sqrt{3/17})$  with respect to the  $N$ th roots of unity are quantized using the first-order  $\Sigma\Delta$  scheme. This log-log plot shows the approximation error  $\|x - \tilde{x}\|$  as a function of  $N$  compared with  $5/N$  and  $5/N^{1.25}$ .

is too weak to be of much use since substituting it into an error bound such as the even case of (22) only gives

$$\|x - \tilde{x}\| \leq \frac{\delta d(N-1)}{N}.$$

In particular, this bound does not go to zero as  $N$  gets large, i.e., as one chooses more redundant frames. On the other hand, if one finds a family of frames and a sequence of permutations, such that the resulting frame variations are uniformly bounded, then one is able to obtain an approximation error of order  $1/N$ .

*Example IV.1 (Roots of Unity):* For  $N > 2$ , let  $R_N = \{e_n^N\}_{n=1}^N$  be the  $N$ th roots of unity viewed as vectors in  $\mathbb{R}^2$ , namely

$$\forall n = 1, \dots, N, \quad e_n^N = (\cos(2\pi n/N), \sin(2\pi n/N)).$$

It is well known that  $R_N$  is a tight frame for  $\mathbb{R}^2$  with frame bound  $N/2$ , e.g., see [30]. In this example, we shall always consider  $R_N$  in its natural ordering  $\{e_n^N\}_{n=1}^N$ . Note that  $\sum_{n=1}^N e_n^N = 0$ .

Since  $\|e_n - e_{n+1}\| \leq 2\pi/N$ , it follows that

$$\forall N, \quad \sigma(R_N, p) \leq 2\pi \quad (23)$$

where  $p$  is the identity permutation of  $\{1, 2, \dots, N\}$ .

Thus, the error estimate of Corollary III.8 gives

$$\|x - \tilde{x}\| \leq \begin{cases} \frac{\delta}{N}\pi, & \text{if } N > 2 \text{ even} \\ \frac{\delta}{2N}(2\pi + 1), & \text{if } N > 2 \text{ odd.} \end{cases} \quad (24)$$

Fig. 1 shows a log-log plot of the approximation error  $\|x - \tilde{x}_N\|$  as a function of  $N$ , when the  $N$ th roots of unity are used to quantize the input  $x = (1/\pi, \sqrt{3/17})$ . The figure also shows a log-log plot of  $5/N$  and  $5/N^{1.25}$  for comparison. Note that the approximation error exhibits two very different types of behavior. In particular, for odd  $N$ , the approximation error appears to behave like  $1/N$  asymptotically, whereas for even  $N$ , the approximation error is much smaller. We shall explain this phenomenon in Section V.

The most natural examples of unit-norm tight frames in  $\mathbb{R}^d$ ,  $d > 2$  are the harmonic frames, e.g., see [4], [29], [2]. These frames are constructed using rows of the Fourier matrix.

*Example IV.2 (Harmonic Frames):* We shall show that harmonic frames in their natural ordering have uniformly bounded frame variation. We follow the notation of [29], although the terminology ‘‘harmonic frame’’ is not specifically used there. The definition of the harmonic frame  $H_N^d = \{e_j\}_{j=0}^{N-1}$ ,  $N \geq d$ , depends on whether the dimension  $d$  is even or odd.

If  $d \geq 2$  is even, let

$$e_j = \sqrt{\frac{2}{d}} \begin{bmatrix} \cos \frac{2\pi j}{N}, \sin \frac{2\pi j}{N}, \cos \frac{2\pi 2j}{N}, \sin \frac{2\pi 2j}{N}, \cos \frac{2\pi 3j}{N}, \\ \sin \frac{2\pi 3j}{N}, \dots, \cos \frac{2\pi \frac{d}{2}j}{N}, \sin \frac{2\pi \frac{d}{2}j}{N} \end{bmatrix}$$

for  $j = 0, 1, \dots, N-1$ .

If  $d > 1$  is odd let

$$e_j = \sqrt{\frac{d}{2}} \left[ \frac{1}{\sqrt{2}}, \cos \frac{2\pi j}{N}, \sin \frac{2\pi j}{N}, \cos \frac{2\pi 2j}{N}, \sin \frac{2\pi 2j}{N}, \dots, \cos \frac{2\pi \frac{d-1}{2} j}{N}, \sin \frac{2\pi \frac{d-1}{2} j}{N} \right]$$

for  $j = 0, 1, \dots, N - 1$ .

It is shown in [29] that  $H_N^d$ , as defined above, is a unit-norm tight frame for  $\mathbb{R}^d$ . If  $d$  is even then  $H_N^d$  satisfies the zero sum condition (18). If  $d$  is odd the frame is not zero sum, and, in fact

$$\sum_{j=0}^{N-1} e_j = \left( \frac{N}{\sqrt{d}}, 0, 0, \dots, 0 \right) \in \mathbb{R}^d.$$

The verification of the zero sum condition for  $d$  even follows by noting that, for each  $k \in \mathbb{Z}$  and not of the form  $k = mN$ , we have

$$\sum_{j=0}^{N-1} \cos \frac{2\pi k j}{N} = \operatorname{Re} \left[ \sum_{j=0}^{N-1} (e^{2\pi i k/N})^j \right] = 0$$

and

$$\sum_{j=0}^{N-1} \sin \frac{2\pi k j}{N} = \operatorname{Im} \left[ \sum_{j=0}^{N-1} (e^{2\pi i k/N})^j \right] = 0.$$

Note that we are simply considering one particular class of harmonic frames in this example, and that one could instead consider other families of frames which are zero sum in all dimensions.

Let us now estimate the frame variation for harmonic frames. First, suppose  $d$  even, and let  $p$  be the identity permutation. Calculating directly and using the mean value theorem in the first inequality, we have

$$\begin{aligned} & \sqrt{\frac{d}{2}} \sigma(H_N^d, p) \\ &= \sqrt{\frac{d}{2}} \sum_{j=0}^{N-2} \|e_j - e_{j+1}\| \\ &= \sum_{j=0}^{N-2} \left[ \sum_{k=1}^{d/2} \left( \cos \frac{2\pi k j}{N} - \cos \frac{2\pi k (j+1)}{N} \right)^2 \right. \\ & \quad \left. + \sum_{k=1}^{d/2} \left( \sin \frac{2\pi k j}{N} - \sin \frac{2\pi k (j+1)}{N} \right)^2 \right]^{\frac{1}{2}} \\ &\leq \sum_{j=0}^{N-2} \left[ 2 \sum_{k=1}^{d/2} \left( \frac{2\pi k}{N} \right)^2 \right]^{\frac{1}{2}} \leq 2\pi\sqrt{2} \left[ \sum_{k=1}^{d/2} k^2 \right]^{\frac{1}{2}} \\ &= 2\pi\sqrt{2} \left[ \frac{d(d/2+1)(d+1)}{12} \right]^{\frac{1}{2}} \leq 2\pi\sqrt{\frac{d}{6}}(d+1). \end{aligned}$$

If  $d$  is odd then, proceeding as above, we have

$$\sqrt{\frac{d}{2}} \sigma(H_N^d, p) \leq 2\pi\sqrt{2} \left[ \sum_{k=1}^{(d-1)/2} k^2 \right]^{\frac{1}{2}} \leq 2\pi\sqrt{\frac{d}{6}}(d+1).$$

Thus,

$$\sigma(H_N^d, p) \leq \frac{2\pi(d+1)}{\sqrt{3}} \tag{25}$$

where  $p$  is the identity permutation, i.e., we consider the natural ordering as in the definition of  $H_N^d$ .

We can now derive error estimates for  $\Sigma\Delta$  quantization of harmonic frames in their natural order. If we set  $u_0 = 0$  and assume that  $x \in \mathbb{R}^d$  satisfies  $\|x\| \leq (K - 1/2)\delta$ , then combining (25), Corollaries III.2, III.5, and III.8, and the fact that  $H_N^d$  satisfies (18) if  $N$  is even gives

$$\|x - \tilde{x}\| \leq \begin{cases} \frac{\delta d}{2N} \frac{2\pi(d+1)}{\sqrt{3}}, & \text{if } d \text{ is even and } N \text{ is even} \\ \frac{\delta d}{2N} \left[ \frac{2\pi(d+1)}{\sqrt{3}} + 1 \right], & \text{otherwise.} \end{cases}$$

Fig. 2 shows a log-log plot of the approximation error  $\|x - \tilde{x}_N\|$  as a function of  $N$ , when the harmonic frame  $H_N^4$  is used to quantize the input  $x = (1/\pi, 1/50, \sqrt{3/17}, e^{-1})$ . The figure also shows a log-log plot of  $10/N$  and  $10/N^{1.25}$  for comparison.

It is worth pointing out the different behavior of the implicit main error term and boundary terms in the  $\Sigma\Delta$  error estimate for harmonic frames in Example IV.2. First, note that the boundary term vanishes when  $N$  is even but not when  $N$  is odd. Second, note that the main error term, i.e., the frame variation term, dominates the boundary term in large dimensions, meaning that the boundary term becomes less significant in higher dimensions. This should be compared with the infinite-dimensional situation in  $\Sigma\Delta$  quantization of band-limited signals, where there is no boundary term.

As discussed earlier, the  $\Sigma\Delta$  algorithm is quite sensitive to the ordering in which the frame coefficients are quantized. In Examples IV.1 and IV.2, the natural frame order gave uniformly bounded frame variation. Let us next consider an example where a bad choice of frame ordering leads to poor approximation error.

*Example IV.3 (Order Matters):* Consider the unit-norm tight frame for  $\mathbb{R}^2$  which is given by the seventh roots of unity, viz.,  $R_7 = \{e_n\}_{n=1}^7$ , where  $e_n = (\cos(2\pi n/7), \sin(2\pi n/7))$ . We randomly choose 10 000 points in the unit ball of  $\mathbb{R}^2$ . For each of these 10 000 points, we first quantize the corresponding frame coefficients in their natural order using (13) with the alphabet

$$\mathcal{A}_4^{1/4} = \{-7/8, -5/8, -3/8, -1/8, 1/8, 3/8, 5/8, 7/8\}$$

and setting  $x_n = \langle x, e_n \rangle$ . Fig. 3 shows the histogram of the corresponding approximation errors. Next, we quantize the frame coefficients of the same 10 000 points, only this time after re-ordering the frame coefficients as  $x_1, x_4, x_7, x_3, x_6, x_2, x_5$ . Fig. 4 shows the histogram of the corresponding approximation errors in this case.

Clearly, the average approximation error for the new ordering is significantly larger than the average approximation error associated with the original ordering. This is intuitively explained by the fact that the natural ordering has significantly smaller frame variation than the other ordering. In particular, let  $p_1$  be the identity permutation and let  $p_2$  be the permutation corresponding



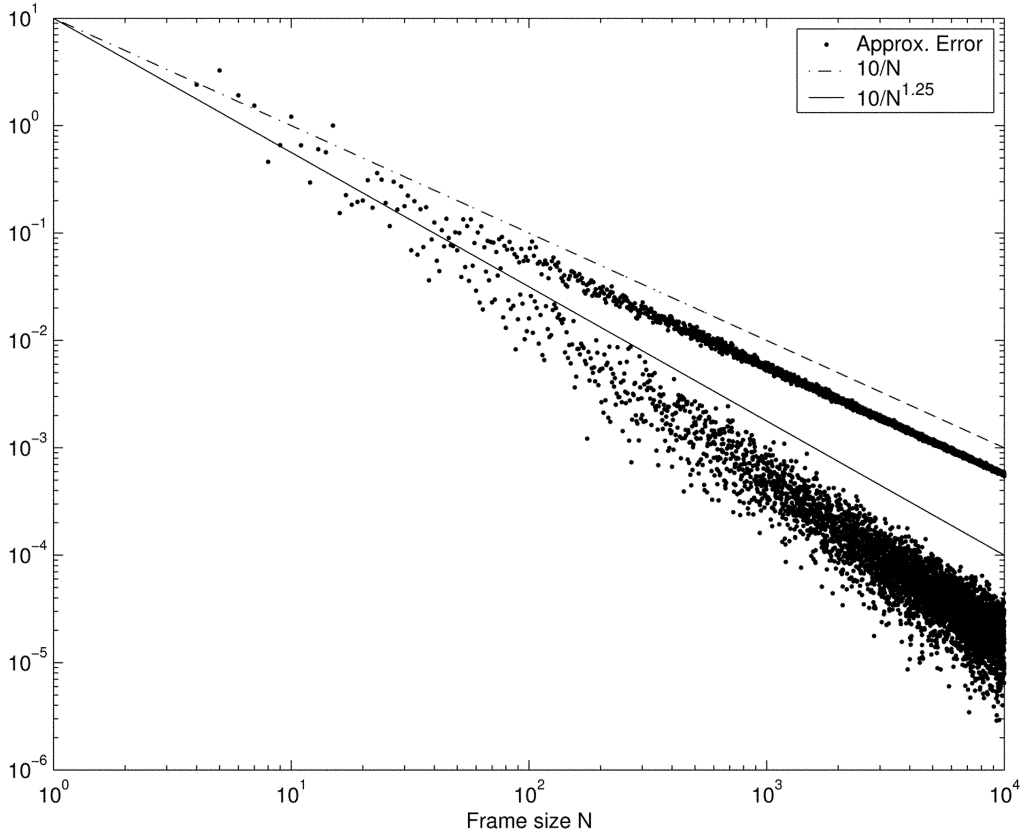


Fig. 2. The frame coefficients of  $x = (1/\pi, 1/50, \sqrt{3/17}, e^{-1})$  with respect to the harmonic frame  $H_N^A$  are quantized using the first-order  $\Sigma\Delta$  scheme. This log-log plot shows the approximation error  $\|x - \tilde{x}\|$  as a function of  $N$  compared with  $10/N$  and  $10/N^{1.25}$ .

the reordered frame coefficients used above. A direct calculation shows that

$$\sigma(F, p_1) \approx 5.2066 \quad \text{and} \quad \sigma(F, p_2) \approx 11.6991.$$

In view of this example, it is important to choose carefully the order in which frame coefficients are quantized. In  $\mathbb{R}^2$ , there is always a simple good choice.

*Theorem IV.4:* Let  $F_N = \{e_n\}_{n=1}^N$  be a unit-norm frame for  $\mathbb{R}^2$ , where  $e_n = (\cos(\alpha_n), \sin(\alpha_n))$ , and  $0 \leq \alpha_n < 2\pi$ . If  $p$  is a permutation of  $\{1, 2, \dots, N\}$  such that  $\alpha_{p(n)} \leq \alpha_{p(n+1)}$  for all  $n \in \{1, 2, \dots, N-1\}$ , then  $\sigma(F_N, p) \leq 2\pi$ .

*Proof:* Is is easy to verify that

$$\|e_{p(n)} - e_{p(n+1)}\| \leq |\alpha_{p(n)} - \alpha_{p(n+1)}|.$$

By the choice of  $p$ , and since  $0 \leq \alpha_n < 2\pi$ , it follows that

$$\sigma(F_N, p) = \sum_{n=1}^{N-1} \|e_{p(n)} - e_{p(n+1)}\| \leq 2\pi.$$

## V. REFINED ESTIMATES AND LOWER BOUNDS

In Fig. 1 of Example IV.1, we saw that the approximation error appears to exhibit very different types of behavior depending on whether  $N$  is even or odd. In the even case, the approximation error appears to decay better than the  $1/N$  estimate given by the results in Section III-B; in the odd case, it appears that the  $1/N$  actually serves as a lower bound, as well as an upper bound, for the approximation error. This dichotomy goes beyond Corollary III.8, which only predicts different con-

stants in the even/odd approximations as opposed to different orders of approximation. In this section, we shall explain this phenomenon.

Let  $\{F_N\}_{N=d}^\infty$  be a family of unit-norm tight frames for  $\mathbb{R}^d$ , with  $F_N = \{e_n^N\}_{n=1}^N$ , so that  $F_N$  has frame bound  $N/d$ . If  $x \in \mathbb{R}^d$ , then  $\{x_n^N\}_{n=1}^N$  will denote the corresponding sequence of frame coefficients with respect to  $F_N$ , i.e.,  $x_n^N = \langle x, e_n^N \rangle$ . Let  $\{q_n^N\}_{n=1}^N$  be the quantized sequence which is obtained by running the  $\Sigma\Delta$  scheme, (13), on the input sequence  $\{x_n^N\}_{n=1}^N$ , and let  $\{u_n^N\}_{n=0}^N$  be the associated state sequence. Thus, if  $x \in \mathbb{R}^d$  is expressed as a frame expansion with respect to  $F_N$ , and if this expansion is quantized by the first-order  $\Sigma\Delta$  scheme, then the resulting quantized expansion is

$$\tilde{x}_N = \frac{d}{N} \sum_{n=1}^N q_n^N e_n^N.$$

Let us begin by rewriting the approximation error in a slightly more revealing form than in Section III-B. Starting with (17), specifying  $u_0^N = 0$ , and specializing to the tight frame case where  $S^{-1} = \frac{d}{N}I$ , we have

$$\begin{aligned} x - \tilde{x}_N &= \frac{d}{N} \left( \sum_{n=1}^{N-1} u_n^N (e_n^N - e_{n+1}^N) + u_N^N e_N^N \right) \\ &= \frac{d}{N} \left( \sum_{n=1}^{N-2} v_n^N (f_n^N - f_{n+1}^N) \right. \\ &\quad \left. + v_{N-1}^N f_{N-1}^N + u_N^N e_N^N \right) \end{aligned} \quad (26)$$

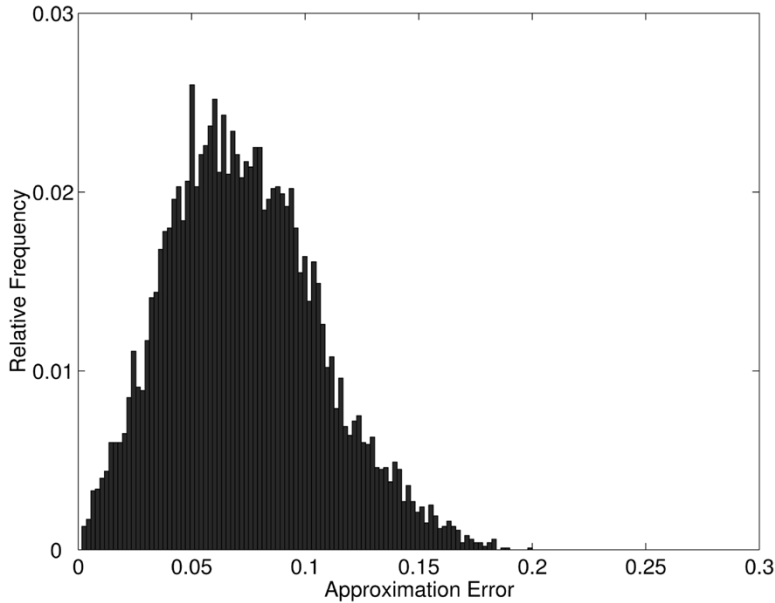


Fig. 3. Histogram of approximation error in Example IV.3 for the natural ordering.

where we have defined

$$f_n^N = e_n^N - e_{n+1}^N, \quad v_n^N = \sum_{j=1}^n u_j^N, \quad \text{and } v_0^N = 0. \quad (27)$$

When working with the approximation error written as (26), the main step toward finding improved upper error bounds, as well as lower bounds, for  $\|x - \tilde{x}\|$ , is to find a good estimate for  $|v_n|$ .

Let  $\mathcal{B}_\Omega$  be the class of  $\Omega$ -band-limited functions consisting of all functions in  $L^\infty(\mathbb{R})$  whose Fourier transforms (as distributions) are supported in  $[-\Omega, \Omega]$ . We shall work with the Fourier transform which is formally defined by  $\hat{f}(\gamma) = \int f(t)e^{-2\pi i t \gamma} dt$ . By the Paley–Wiener theorem, elements of  $\mathcal{B}_\Omega$  are restrictions of entire functions to the real line.

*Definition V.1:* Let  $f \in \mathcal{B}_\Omega$  and let  $\{z_j\}_{j=1}^{n^*}$  be the finite set of zeros of  $f'$  contained in  $[0, 1]$ . We say that  $f \in \mathcal{M}_\Omega$  if

$$\forall j = 1, \dots, n^*, \quad f''(z_j) \neq 0.$$

For simplicity and to avoid having to keep track of too many different constants, we shall use the notation  $A \lesssim B$  to mean that there exists a fixed constant  $C > 0$  such that  $A \leq CB$ . When necessary, we shall point out the dependence of  $C$  on other parameters. The following theorem relies on the uniform distribution techniques utilized by Güntürk in [16]. We briefly collect the necessary background on discrepancy and uniform distribution in Appendix I.

*Theorem V.2:* Let  $\mathcal{F} = \{F_N\}_{N=d}^\infty$  be a family of unit-norm tight frames for  $\mathbb{R}^d$ , with  $F_N = \{e_n^N\}_{n=1}^N$ . Fix  $x \in \mathbb{R}^d$  such that  $\|x\| \leq (K - 1/2)\delta$ , and let  $\{x_n^N\}_{n=1}^N$  be the sequence of frame coefficients of  $x$  with respect to  $F_N$ . If, for some  $\Omega > 0$ , there exists  $h = h_{\mathcal{F},x} \in \mathcal{M}_\Omega$  such that

$$\forall N \text{ and } 1 \leq n \leq N, \quad x_n^N = h(n/N)$$

and if  $N$  is sufficiently large, then

$$|v_n^N| \lesssim \delta \left( \frac{n}{N^{1/4}} + N^{3/4} \log N \right) \lesssim \delta N^{3/4} \log N. \quad (28)$$

The implicit constants are independent of  $N$  and  $\delta$ , but they do depend on  $h_{\mathcal{F},x}$ . The value of what constitutes a sufficiently large  $N$  depends on  $\delta$ .

*Proof:* Let  $u_n^N$  be the state variable of the  $\Sigma\Delta$  scheme and define  $\tilde{u}_n^N = u_n^N/\delta$ . By the definition of  $v_n^N$  (see (27)), and by applying Koksma’s inequality (see Appendix I), one has

$$\begin{aligned} |v_j^N| &= \delta \left| \sum_{n=1}^j \tilde{u}_n^N \right| = j\delta \left| \frac{1}{j} \sum_{n=1}^j \tilde{u}_n^N - \int_{-1/2}^{1/2} y \, dy \right| \\ &\leq j\delta \text{Disc} \left( \{\tilde{u}_n^N\}_{n=1}^j \right) \end{aligned} \quad (29)$$

where  $\text{Disc}(\cdot)$  denotes the *discrepancy* of a sequence as defined in Appendix I. Therefore, we need to estimate  $D_j^N = \text{Disc}(\{\tilde{u}_n^N\}_{n=1}^j)$ . Using the Erdős–Turán inequality (see Appendix I)

$$\forall K, \quad D_j^N \leq \frac{1}{K} + \frac{1}{j} \sum_{k=1}^K \frac{1}{k} \left| \sum_{n=1}^j e^{2\pi i k \tilde{u}_n^N} \right| \quad (30)$$

we see that it suffices for us to estimate  $|\sum_{n=1}^j e^{2\pi i k \tilde{u}_n^N}|$ .

Proceeding as in [16, Proposition 1], for each  $N$  there exists an analytic function  $X_N \in \mathcal{B}_{\Omega/N}$  such that

$$u_n^N = X_N(n) \text{ modulo } [-\delta/2, \delta/2] \quad (31)$$

and

$$|X'_N(t) - h(t/N)| \lesssim 1/N. \quad (32)$$

Bernstein’s inequality gives

$$\left| X''_N(t) - \frac{1}{N} h'(t/N) \right| \lesssim 1/N^2. \quad (33)$$

In (32) and (33), the implicit constants are independent of  $N$  and  $\delta$ , but do depend on  $h = h_{\mathcal{F},x}$ .

By hypothesis,  $h \in \mathcal{M}_\Omega$  satisfies  $x_n^N = h(n/N)$ . Let  $\{z_\ell\}_{\ell=1}^{n^*}$  be the set of zeros of  $h$  in  $[0, 1]$ , and let  $0 < \alpha < 1$  be

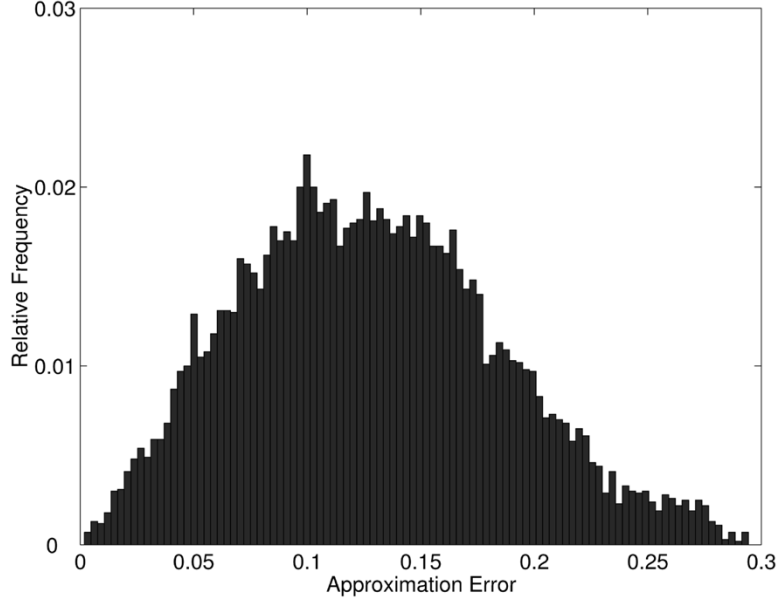


Fig. 4. Histogram of approximation error in Example IV.3 for an ordering giving higher variation.

a fixed constant to be specified later. Define the intervals  $I_\ell$  and  $J_\ell$  by

$$\begin{aligned} \forall \ell = 1, \dots, n^*, \quad I_\ell &= (Nz_\ell - N^\alpha, Nz_\ell + N^\alpha) \\ \forall \ell = 1, \dots, n^* - 1, \quad J_\ell &= [Nz_\ell + N^\alpha, Nz_{\ell+1} - N^\alpha] \end{aligned}$$

and

$$J_0 = [1, Nz_1 - N^\alpha] \text{ and } J_{n^*} = [Nz_{n^*} + N^\alpha, N].$$

In the case where either 0 or 1 is a zero of  $h'$ , one no longer needs the corresponding endpoint interval  $J_\ell$ , but needs to modify the corresponding interval  $I_\ell$  to have 0 or 1 as its appropriate endpoint. Note that if  $N$  is sufficiently large then

$$J_0 \cup I_1 \cup J_1 \cup \dots \cup I_{n^*} \cup J_{n^*} = [1, N].$$

It follows from the properties of  $h \in \mathcal{M}_\Omega$  that if  $N$  is sufficiently large then

$$\forall n \in \mathbb{N} \cap J_\ell, \quad \frac{1}{N^{1-\alpha}} = \frac{N^\alpha}{N} \lesssim |h'(n/N)|. \quad (34)$$

Thus, by (33), we have that

$$\forall n \in \mathbb{N} \cap J_\ell, \quad \frac{k}{\delta N^{2-\alpha}} \lesssim \left| \frac{k}{\delta} X_N''(n) \right|. \quad (35)$$

Also, since  $h \in \mathcal{M}_\Omega \subseteq L^\infty(\mathbb{R})$ , and by (32), we obtain

$$\forall n \in \mathbb{N} \cap J_\ell, \quad \left| \frac{k}{\delta} X_N'(n) \right| \lesssim \frac{k}{\delta}. \quad (36)$$

In (34)–(36) the implicit constants do not depend on  $N$  and  $\delta$ . Using (35), (36), [37, Theorem 2.7], and since  $0 < \delta < 1$ , we have that for  $1 \leq k$

$$\begin{aligned} \left| \sum_{n \in \mathbb{N} \cap J_\ell} e^{2\pi i k \tilde{u}_n^N} \right| &= \left| \sum_{n \in \mathbb{N} \cap J_\ell} e^{2\pi i (k/\delta) X_N(n)} \right| \\ &\lesssim (2k/\delta + 2) \left( \frac{\delta^{1/2} N^{1-\frac{\alpha}{2}}}{k^{1/2}} + 1 \right) \\ &\lesssim \frac{k^{1/2} N^{1-\frac{\alpha}{2}}}{\delta^{1/2}} + \frac{k}{\delta}. \end{aligned}$$

Also, we have the trivial estimate

$$\left| \sum_{n \in \mathbb{N} \cap I_\ell} e^{2\pi i k \tilde{u}_n^N} \right| \leq 2N^\alpha.$$

Thus,

$$\left| \sum_{n=1}^j e^{2\pi i k \tilde{u}_n^N} \right| \lesssim N^\alpha + \frac{k^{1/2} N^{1-\frac{\alpha}{2}}}{\delta^{1/2}} + \frac{k}{\delta}.$$

Here, the implicit constant is independent of  $N$  and  $\delta$ , but it does depend on  $h = h_{\mathcal{F},x}$  due to the role of  $n^*$ . Set  $\alpha = 3/4$  and  $K = N^{1/4}$ . By (30), we have that if  $N$  is sufficiently large compared to  $\delta$  then

$$\begin{aligned} D_j^N &\leq \frac{1}{K} + \frac{N^\alpha \log(K)}{j} + \frac{K^{1/2} N^{1-\frac{\alpha}{2}}}{\delta^{1/2} j} + \frac{K}{\delta j} \\ &\lesssim \frac{1}{N^{1/4}} + \frac{N^{3/4} \log(N)}{j} + \frac{N^{3/4}}{\delta^{1/2} j} + \frac{N^{1/4}}{\delta j} \\ &\lesssim \frac{1}{N^{1/4}} + \frac{N^{3/4} \log(N)}{j}. \end{aligned}$$

Thus, by (29), we have

$$|v_n^N| \leq \frac{\delta n}{N^{1/4}} + \delta N^{3/4} \log N \lesssim \delta N^{3/4} \log N,$$

and the proof is complete.

Combining Theorem V.2 and (26) gives the following improved error estimate. Although this estimate guarantees approximation on the order of  $\frac{\log N}{N^{5/4}}$  for even  $N$ , it is important to emphasize that the implicit constants depend on  $x$ . For comparison, note that Corollary III.8 only bounds the error by the order of  $\frac{1}{N}$ , but has explicit constants independent of  $x$ .

*Corollary V.3:* Let  $\mathcal{F} = \{F_N\}_{N=d}^\infty$  be a family of unit-norm tight frames for  $\mathbb{R}^d$ , for which each  $F_N = \{e_n^N\}_{n=1}^N$  satisfies the zero sum condition (18). Fix  $x \in \mathbb{R}^d$  such that  $\|x\| \leq$

$(K-1/2)\delta$ , let  $\{x_n^N\}_{n=1}^N$  be the frame coefficients of  $x$  with respect to  $F_N$ , and suppose there exists  $h = h_{\mathcal{F},x} \in \mathcal{M}_\Omega$ ,  $\Omega > 0$ , such that

$$\forall N \text{ and } 1 \leq n \leq N, \quad x_n^N = h(n/N).$$

Additionally, suppose that  $f_n^N = e_n^N - e_{n+1}^N$  satisfies

$$\forall N, n = 1, \dots, N, \quad \|f_n^N\| \lesssim \frac{1}{N} \text{ and } \|f_n^N - f_{n+1}^N\| \lesssim \frac{1}{N^2}.$$

If  $N$  is even and sufficiently large we have

$$\|x - \tilde{x}_N\| \lesssim \frac{\delta \log N}{N^{5/4}}.$$

If  $N$  is odd and sufficiently large we have

$$\frac{\delta}{N} \lesssim \|x - \tilde{x}_N\| \leq \frac{\delta d}{2N} (\sigma(F_N, p_N) + 1).$$

The implicit constants are independent of  $\delta$  and  $N$ , but do depend on  $h_{\mathcal{F},x}$ .

*Proof:* By Theorem V.2

$$\left\| \frac{2}{N} \left( \sum_{n=1}^{N-2} v_n^N (f_n^N - f_{n+1}^N) + v_{N-1}^N f_{N-1}^N \right) \right\| \lesssim \frac{\delta \log N}{N^{5/4}}. \quad (37)$$

Thus, by Theorem III.7, (37), and (26),  $N$  being even implies

$$\|x - \tilde{x}_N\| \lesssim \frac{\delta \log N}{N^{5/4}}.$$

If  $N$  is odd, then by Theorem III.7, (26), and (37) we have

$$\frac{\delta}{N} = \frac{2 \|u_N^N\| \|e_N^N\|}{N} \lesssim \|x - \tilde{x}_N\| + \frac{\delta \log N}{N^{5/4}}.$$

Combining this with (22) completes the proof.

Applying Corollary V.3 to the quantization of frame expansions given by the roots of unity explains the different error behavior for even and odd  $N$  seen in Fig. 1.

*Example V.4 (Refined Estimates for  $R_n$ ):* Let  $R_N = \{e_n^N\}_{n=1}^N$  be as in Example IV.1, i.e.,  $R_N$  is the unit-norm tight frame for  $\mathbb{R}^2$  given by the  $N$ th roots of unity. Suppose  $x \in \mathbb{R}^2$ ,  $0 < \|x\| \leq (K-1/2)\delta$ , and that  $N$  is sufficiently large with respect to  $\delta$ . The frame coefficients of  $x = (a, b) \in \mathbb{R}^2$  with respect to  $R_N$  are given by

$$\{x_n^N\}_{n=1}^N = \{h(n/N)\}_{n=1}^N$$

where  $h(t) = a \cos(2\pi t) + b \sin(2\pi t)$ .

It is straightforward to show that  $f_n^N = e_n^N - e_{n+1}^N$  satisfies

$$\|f_n^N\| \leq \frac{2\pi}{N} \quad \text{and} \quad \|f_n^N - f_{n+1}^N\| \leq \frac{(2\pi)^2}{N^2}$$

and that  $h \in \mathcal{M}_1$ . Therefore, by Corollary V.3 and (23), if  $N$  is even then

$$\|x - \tilde{x}\| \lesssim \frac{\delta \log N}{N^{5/4}},$$

and if  $N$  is odd then

$$\frac{\delta}{N} \lesssim \|x - \tilde{x}\| \leq \frac{\delta(2\pi+1)}{N}.$$

The implicit constants are independent of  $\delta$  and  $N$ , but do depend on  $x$ .

It is sometimes also possible to apply Corollary V.3 to harmonic frames.

*Example V.5 (Refined Estimates for  $H_N^d$ ):* Let the dimension  $d$  be even, and let  $H_N^d = \{e_n^N\}_{n=1}^N$  be as in Example IV.2, i.e.,  $H_N^d$  is an harmonic frame for  $\mathbb{R}^d$ . Suppose  $x \in \mathbb{R}^d$ ,  $0 < \|x\| \leq (K-1/2)\delta$ , and that  $N$  is sufficiently large with respect to  $\delta$ . The frame coefficients of  $x = (a_1, b_1, \dots, a_{d/2}, b_{d/2}) \in \mathbb{R}^d$  with respect to  $H_N^d$  are given by  $\{x_n^N\}_{n=1}^N = \{h(n/N)\}_{n=1}^N$ , where

$$h(t) = \sqrt{\frac{2}{d}} \left( \sum_{j=1}^{d/2} a_j \cos(2\pi jt) + \sum_{j=1}^{d/2} b_j \sin(2\pi jt) \right).$$

Fig. 2 in Example IV.2 shows the approximation error when the point

$$x = (1/\pi, 1/50, \sqrt{3/17}, e^{-1}) \in \mathbb{R}^4$$

is represented with the harmonic frames  $\{H_N^4\}_{N=4}^\infty$  and quantized using the  $\Sigma\Delta$  scheme. For this choice of  $x$  it is straightforward to verify that  $h \in M_{d/2}$ . A direct estimate also shows that  $f_n^N = e_n^N - e_{n+1}^N$  satisfies

$$\|f_n^N\| \lesssim \frac{1}{N} \quad \text{and} \quad \|f_n^N - f_{n+1}^N\| \lesssim \frac{1}{N^2}.$$

Therefore, by Corollary V.3 and (23), if  $N$  is even, then

$$\|x - \tilde{x}\| \lesssim \frac{\delta \log N}{N^{5/4}}$$

and if  $N$  is odd, then

$$\frac{\delta}{N} \lesssim \|x - \tilde{x}\| \leq \frac{\delta d}{2N} \left( \frac{10\pi}{\sqrt{3}} + 1 \right).$$

The implicit constants are independent of  $\delta$  and  $N$ , but do depend on  $x$ .

## VI. COMPARISON OF $\Sigma\Delta$ WITH PCM

In this section, we shall compare the MSE given by  $\Sigma\Delta$  quantization of finite frame expansions with that given by PCM schemes. We shall show that the  $\Sigma\Delta$  scheme gives better MSE estimates than PCM quantization when dealing with sufficiently redundant frames. Throughout this section, let  $F_N = \{e_n^N\}_{n=1}^N$  be a family of unit-norm tight frames for  $\mathbb{R}^d$ , and let

$$x = \frac{d}{N} \sum_{n=1}^N x_n^N e_n^N \quad \text{and} \quad \tilde{x}_N = \frac{d}{N} \sum_{n=1}^N q_n^N e_n^N$$

be corresponding frame expansions and quantized frame expansions, where  $x_n^N = \langle x, e_n^N \rangle$  are the frame coefficients of  $x \in \mathbb{R}^d$  with respect to  $F_N$ , and where  $q_n^N$  are quantized versions of  $x_n^N$ .

In Example II.5, we showed that if one uses PCM quantization to produce the quantized frame expansion  $\tilde{x}_N$ , then under Bennett's white noise assumption, the PCM scheme has MSE

$$\text{MSE}_{\text{PCM}} = \frac{d^2 \delta^2}{12N}. \quad (38)$$

However, as illustrated in Example II.6, this estimate is not rigorous since Bennett's white noise assumption is not mathematically justified and may in fact fail dramatically in certain circumstances.

If one uses  $\Sigma\Delta$  quantization to produce the quantized frame expansion  $\tilde{x}_N$ , then one has the error estimate

$$\|x - \tilde{x}_N\| \leq \frac{\delta d}{2N} (\sigma(F, p) + 1) \quad (39)$$

given by Corollary III.5. Here,  $p$  is a permutation of  $\{1, \dots, N\}$  which denotes the order in which the  $\Sigma\Delta$  scheme is run. This immediately yields the following MSE estimate for the  $\Sigma\Delta$  scheme.

*Theorem VI.1:* Given the  $\Sigma\Delta$  scheme of Definition III.1, let  $F = \{e_n\}_{n=1}^N$  be a unit-norm tight frame for  $\mathbb{R}^d$ , and let  $p$  be a permutation of  $\{1, 2, \dots, N\}$ . For each  $x \in \mathbb{R}^d$  satisfying  $\|x\| \leq (K - 1/2)\delta$ ,  $\tilde{x}$  shall denote the corresponding quantized output of the  $\Sigma\Delta$  scheme. Let

$$B \subseteq \{x \in \mathbb{R}^d : \|x\| \leq (K - 1/2)\delta\}$$

and define the MSE of the  $\Sigma\Delta$  scheme over  $B$  by

$$\text{MSE}_{\Sigma\Delta} = \int_B \|x - \tilde{x}\|^2 d\mu(x)$$

where  $\mu$  is any probability measure on  $B$ . Then

$$\text{MSE}_{\Sigma\Delta} \leq \frac{\delta^2 d^2}{4N^2} (\sigma(F, p) + 1)^2.$$

*Proof:* Square (39) and integrate.

One may analogously derive MSE bounds from any of the error estimates in Section III-B; we shall examine this in the subsequent example. The above estimate is completely deterministic; namely, it does not depend on statistical assumptions such as the analysis for PCM using Bennett's white noise assumption.

It is also possible to derive MSE estimates for  $\Sigma\Delta$  schemes by making empirically reasonable statistical assumptions similar to Bennett's white noise assumption for PCM. The classical approach, e.g., see [31], [38], is to assume that the state variables  $u_n$  in the  $\Sigma\Delta$  scheme (13) are independent and identically distributed uniform random variables with zero mean and variance  $\delta^2/12$ . We shall refer to this as the classical  $\Sigma\Delta$  white noise assumption.

The classical  $\Sigma\Delta$  white noise assumption yields MSE estimates in a manner similar to the PCM setting. Let us illustrate this for the case where  $\{e_n\}_{n=1}^N$  is the  $N$ th roots of unity

unit-norm tight frame for  $\mathbb{R}^2$  given in Example IV.1. Specializing the error term (17) from Theorem III.4 to this particular frame and taking  $u_0 = 0$  gives

$$x - \tilde{x} = \frac{2}{N} \left( \sum_{n=1}^{N-1} u_n (e_n - e_{n+1}) + u_N e_N \right)$$

where, by Theorem III.7,  $u_N = 0$  when  $N$  is even, and  $|u_N| = \delta/2$  when  $N$  is odd. Using the classical  $\Sigma\Delta$  white noise assumption for  $u_n, n = 1, \dots, N-1$ , a computation for  $\text{MSE}_{\Sigma\Delta WN} = E(\|x - \tilde{x}\|^2)$  similar to that in [2] yields

$$\text{MSE}_{\Sigma\Delta WN} = \frac{4}{N^2} \frac{\delta^2}{12} \sum_{n=1}^{N-1} \|e_n - e_{n+1}\|^2 + \frac{4}{N^2} |u_N|^2.$$

Since

$$\|e_n - e_{n+1}\|^2 = 2(1 - \cos(2\pi/N)) = \frac{4\pi^2}{N^2} + \mathcal{O}(N^{-4})$$

it follows that if  $N$  is even then

$$\text{MSE}_{\Sigma\Delta WN} = \frac{4\pi^2 \delta^2}{3N^3} + \mathcal{O}(N^{-4}) \quad (40)$$

whereas if  $N$  is odd then

$$\text{MSE}_{\Sigma\Delta WN} = \frac{4\pi^2 \delta^2}{3N^3} + \frac{\delta^2}{N^2} + \mathcal{O}(N^{-4}). \quad (41)$$

Analogous MSE estimates may also be derived for general classes of frames. For the  $N$ th roots of unity frame and  $N$  even, the estimate (40) shows that if one is justified in making the classical  $\Sigma\Delta$  white noise assumption then one obtains better MSE estimates than given by Theorem VI.1. On the other hand, for the  $N$ th roots of unity frame with  $N$  odd, the estimate (41) has the same order of approximation, but with a better constant, as the bound in Theorem VI.1 which was made without any statistical assumptions.

In Section IV, we saw that it is possible to choose families of frames,  $F_N = \{e_n^N\}_{n=1}^N$ , for  $\mathbb{R}^d$ , and permutations  $p = p_N$ , such that the resulting frame variation  $\sigma(F_N, p_N)$  is uniformly bounded. Whenever this is the case, Theorem VI.1 yields the MSE bound  $\text{MSE}_{\Sigma\Delta} \lesssim 1/N^2$ , which is better than the PCM bound (38) by one order of approximation. For example, if one quantizes harmonic frame expansions in their natural order, then, by (25), Theorem VI.1 gives  $\text{MSE}_{\Sigma\Delta} \lesssim 1/N^2$ . Thus, for the quantization of harmonic frame expansions one may summarize the difference between  $\Sigma\Delta$  and PCM as

$$\text{MSE}_{\Sigma\Delta} \lesssim 1/N^2 \quad \text{and} \quad 1/N \lesssim \text{MSE}_{\text{PCM}} \lesssim 1/N.$$

This says that  $\Sigma\Delta$  schemes utilize redundancy better than PCM.

Let us remark that for the class of *consistent reconstruction* schemes considered in [2], Goyal, Vetterli, and Thao bound the MSE from below by  $b/A^2$ , where  $b$  is some constant and  $A = N/d$  is the redundancy of the frame. Thus, the MSE estimate derived in Theorem VI.1 for the  $\Sigma\Delta$  scheme achieves this same optimal MSE order.

Returning to classical PCM (with linear reconstruction), it is important to note that although  $\text{MSE}_{\Sigma\Delta} \lesssim 1/N^2$  is much better than  $\text{MSE}_{\text{PCM}} \lesssim 1/N$  for large  $N$ , it is still possible to have  $\text{MSE}_{\text{PCM}} \leq \text{MSE}_{\Sigma\Delta}$  if  $N$  is small, i.e., if the frame has

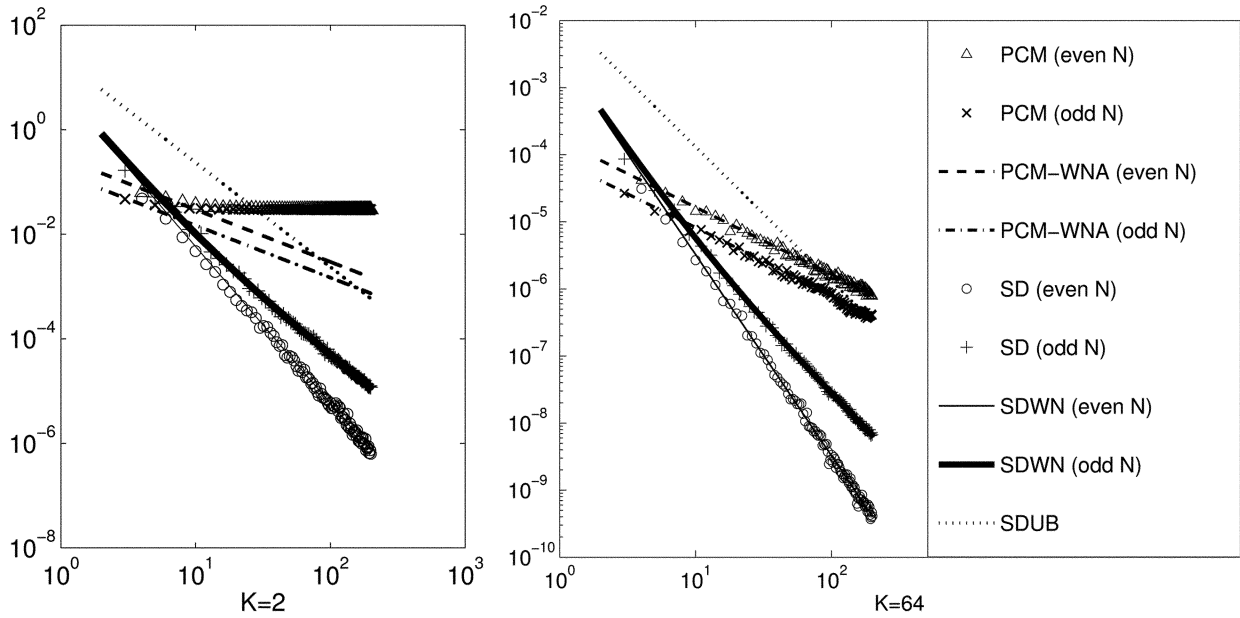


Fig. 5. Comparison of the MSE for  $2K$ -level PCM algorithms and  $2K$ -level first-order  $\Sigma\Delta$  quantizers with step size  $\delta = 1/(K - 1/2)$ . Frame expansions of 100 randomly selected points in  $\mathbb{R}^2$  for frames obtained by the  $N$ th roots of unity were quantized. In the figure legend, PCM and SD correspond to the MSE for PCM and the MSE for first-order  $\Sigma\Delta$  obtained experimentally, respectively. As the asymptotic behavior of the MSE depends on the parity of  $N$  in both cases, the MSEs for even  $N$  and odd  $N$  were plotted using different markers. In the legend, the bound on the MSE for PCM, computed with white noise assumption, is denoted by PCM-WNA, which again depends on the parity of  $N$  as discussed in Example II.6. SDWN in the legend denotes MSE bound on the performance of  $\Sigma\Delta$  given by (40) and (41) for even and odd  $N$ , respectively. Finally, SDUB in the legend stands for the MSE bound for  $\Sigma\Delta$  that follows from (24).

low redundancy. For example, if the frame being quantized is an orthonormal basis, then PCM schemes certainly offer better MSE than  $\Sigma\Delta$  since in this case there is an isometry between the frame coefficients and the signal they represent. Nonetheless, for sufficiently redundant frames,  $\Sigma\Delta$  schemes provide better MSE than PCM.

*Example VI.2 (Unit-Norm Frames for  $\mathbb{R}^2$ ):* In view of Theorem IV.4, it is easy to obtain uniform bounds for the frame variation of frames  $F$  for  $\mathbb{R}^2$ . In particular, one can always find a permutation  $p$  such that  $\sigma(F, p) \leq 2\pi$ .

A simple comparison of the MSE error bounds for PCM and  $\Sigma\Delta$  discussed above shows that the MSE corresponding to first-order  $\Sigma\Delta$  quantizers is less than the MSE corresponding to PCM algorithms for unit-norm tight frames for  $\mathbb{R}^2$  in the following cases when the redundancy  $A$  satisfies the specified inequalities:

- $A > 1.5(2\pi)^2 \approx 59$  if the unit-norm tight frame for  $\mathbb{R}^2$  has even length, is zero sum, is ordered as in Theorem IV.4, and we set  $u_0 = 0$ , see Corollary III.8;
- $A > 1.5(2\pi + 1)^2 \approx 80$  for any unit-norm tight frame for  $\mathbb{R}^2$ , as long as the frame elements are ordered as described in Theorem IV.4, and  $u_0$  is chosen to be 0, see Corollary III.5;
- $A > 1.5(2\pi + 2)^2 \approx 103$  for any unit-norm tight frame for  $\mathbb{R}^2$ , as long as the frame elements are ordered as described in Theorem IV.4, see Corollary III.6.

Fig. 5 shows the MSE achieved by  $2K$ -level PCM algorithms and  $2K$ -level first-order  $\Sigma\Delta$  quantizers with step size  $\delta = 1/K$  for several values of  $K$  for unit-norm tight frames for  $\mathbb{R}^2$  obtained by the  $N$ th roots of unity. The plots suggest that if the

frame bound is larger than approximately 10, the first-order  $\Sigma\Delta$  quantizer outperforms PCM.

*Example VI.3 (7th Roots of Unity):* Let  $x = (1/3, 1/2) \in \mathbb{R}^2$ , and let  $R_7 = \{e_n\}_{n=1}^7$  be the unit-norm tight frame for  $\mathbb{R}^2$  given by

$$e_n = (\cos(2\pi n/7), \sin(2\pi n/7)), \quad n = 1, \dots, 7.$$

The point  $x$  has the frame expansion

$$x = \frac{2}{7} \sum_{n=1}^7 x_n e_n, \quad x_n = \langle x, e_n \rangle.$$

One may compute that

$$(x_1, x_2, x_3, x_4, x_5, x_6, x_7) \approx (0.5987, 0.4133, -0.0834, -0.5173, -0.5616, -0.1831, 0.3333).$$

If we consider the 1-bit alphabet  $\mathcal{A}_1^2 = \{-1, 1\}$  then the quantization problem is to replace  $x$  by an element of

$$\Gamma = \left\{ \frac{2}{7} \sum_{n=1}^7 q_n e_n : q_n \in \mathcal{A}_1^2 \right\}.$$

Fig. 6 shows the elements of  $\Gamma$  denoted by solid dots, and shows the point  $x$  denoted by a “ $\times$ ” symbol. Note that  $x \notin \Gamma$ .

The first-order  $\Sigma\Delta$  scheme with two-level alphabet  $\mathcal{A}_1^2$  and natural ordering  $p$  quantizes  $x$  by  $x_{\Sigma\Delta} \approx (0.5854, 0.5571) \in \Gamma$ . This corresponds to replacing  $(x_1, x_2, x_3, x_4, x_5, x_6, x_7)$  by

$$(q_1, q_2, q_3, q_4, q_5, q_6, q_7) = (1, 1, -1, -1, -1, 1, -1).$$

The two-level PCM scheme quantizes  $x$  by  $x_{\text{PCM}} \approx (0.8006, 1.0039) \in \Gamma$ . This corresponds to replacing  $(x_1, x_2, x_3, x_4, x_5, x_6, x_7)$  by

$$(q_1, q_2, q_3, q_4, q_5, q_6, q_7) = (1, 1, -1, -1, -1, -1, 1).$$

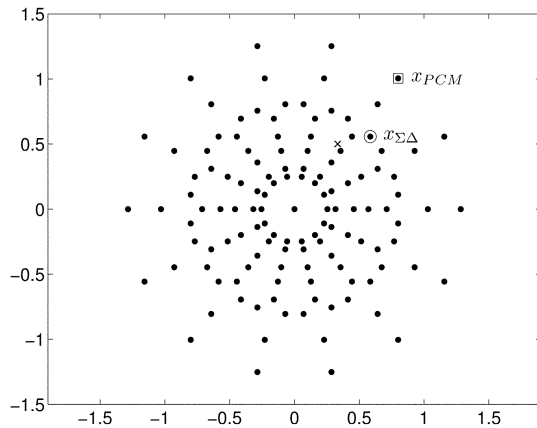


Fig. 6. The elements of  $\Gamma$  from Example VI.3 are denoted by solid dots, and the point  $x = (1/3, 1/2)$  is denoted by “x.” Note that  $x \notin \Gamma$ . “ $x_{\Sigma\Delta}$ ” is the quantized point in  $\Gamma$  obtained using first-order  $\Sigma\Delta$  quantization, and “ $x_{PCM}$ ” is the quantized point in  $\Gamma$  obtained by PCM quantization.

The points  $x_{\Sigma\Delta}$  and  $x_{PCM}$  are shown in Fig. 6 and it is visually clear that  $\|x - x_{\Sigma\Delta}\| < \|x - x_{PCM}\|$ .

The sets  $\Gamma$  corresponding to more general frames and alphabets than in Example VI.3 possess many interesting properties. This is a direction of ongoing work of the authors together with Yang Wang.

## VII. CONCLUSION

We have introduced the  $K$ -level  $\Sigma\Delta$  scheme with stepsize  $\delta$  as a technique for quantizing finite frame expansions for  $\mathbb{R}^d$ . In Section III, we have proven that if  $F$  is a unit-norm tight frame for  $\mathbb{R}^d$  of cardinality  $N$ , and  $x \in \mathbb{R}^d$ , then the  $K$ -level  $\Sigma\Delta$  scheme with stepsize  $\delta$  has approximation error

$$\|x - \tilde{x}\| \leq \frac{\delta d}{2N}(\sigma(F, p) + 1)$$

where the frame variation  $\sigma(F, p)$  depends only on the frame  $F$  and the order  $p$  in which frame coefficients are quantized. As a corollary, for harmonic frames  $H_N^d = \{e_n\}_{n=1}^N$  for  $\mathbb{R}^d$  this gives the approximation error estimate

$$\|x - \tilde{x}\| \leq \frac{\delta d}{2N}(2\pi(d+1) + 1).$$

In Section V, we showed that there are certain cases where the above error bounds can be improved to

$$\|x - \tilde{x}\| \lesssim (\log N)/N^{\frac{5}{4}}$$

where the implicit constant depends on  $x$ . Section VI compares MSE for  $\Sigma\Delta$  schemes and PCM schemes. A main consequence of our approximation error estimates is that

$$\text{MSE}_{\Sigma\Delta} \lesssim 1/N^2 \text{ whereas } 1/N \lesssim \text{MSE}_{\text{PCM}}$$

when linear reconstruction is used, see (39) and (38). This shows that first-order  $\Sigma\Delta$  schemes outperform the standard PCM scheme if the frame being quantized is sufficiently redundant. We have also shown that  $\Sigma\Delta$  quantization with linear reconstruction achieves the same order  $1/N^2$  MSE as PCM with consistent reconstruction.

Our error estimates for first-order  $\Sigma\Delta$  schemes make it reasonable to hope that second-order  $\Sigma\Delta$  schemes can perform even better. This is, in fact, the case, but the analysis of second-

order schemes becomes much more complicated and is considered separately in [35].

## APPENDIX

### DISCREPANCY AND UNIFORM DISTRIBUTION

Let  $\{u_n\}_{n=1}^N \subseteq [-1/2, 1/2)$ , where  $[-1/2, 1/2)$  is identified with the torus  $\mathbb{T}$ . The *discrepancy* of  $\{u_n\}_{n=1}^N$  is defined by

$$\text{Disc}(\{u_n\}_{n=1}^N) = \sup_{I \subset \mathbb{T}} \left| \frac{\#\left(\{u_n\}_{n=1}^N \cap I\right)}{N} - |I| \right|$$

where the sup is taken over all subarcs  $I$  of  $\mathbb{T}$ .

The *Erdős–Turan inequality* allows one to estimate discrepancy in terms of exponential sums

$$\forall K, \quad \text{Disc}(\{u_n\}_{n=1}^j) \leq \frac{1}{K} + \frac{1}{j} \sum_{k=1}^K \frac{1}{k} \left| \sum_{n=1}^j e^{2\pi i k u_n} \right|.$$

*Koksma’s inequality* states that for any function  $f : [-1/2, 1/2) \rightarrow \mathbb{R}$  of bounded variation

$$\left| \frac{1}{N} \sum_{n=1}^N f(u_n) - \int_{-1/2}^{1/2} f(t) dt \right| \leq \text{Var}(f) \text{Disc}(\{u_n\}_{n=1}^N).$$

## ACKNOWLEDGMENT

The authors would like to thank Ingrid Daubechies, Sinan Güntürk, and Nguyen Thao for valuable discussions on the material. The authors also thank Götz Pfander for sharing insightful observations on finite frames. Finally, the authors are especially grateful to the referees for their thoughtful and constructive comments.

## REFERENCES

- [1] N. Thao and M. Vetterli, “Deterministic analysis of oversampled A/D conversion and decoding improvement based on consistent estimates,” *IEEE Trans. Signal Processing*, vol. 42, no. 3, pp. 519–531, Mar. 1994.
- [2] V. Goyal, M. Vetterli, and N. Thao, “Quantized overcomplete expansions in  $\mathbb{R}^n$ : Analysis, synthesis, and algorithms,” *IEEE Trans. Info. Theory*, vol. 44, no. 1, pp. 16–31, Jan. 1998.
- [3] Z. Cvetković, “Resilience properties of redundant expansions under additive noise and quantization,” *IEEE Trans. Inf. Theory*, vol. 49, no. 3, pp. 644–656, Mar. 2003.
- [4] V. Goyal, J. Kovačević, and J. Kelner, “Quantized frame expansions with erasures,” *Appl. Comput. Harmon. Anal.*, vol. 10, pp. 203–233, 2001.
- [5] B. Beferull-Lozano and A. Ortega, “Efficient quantization for overcomplete expansions in  $\mathbb{R}^N$ ,” *IEEE Trans. Inf. Theory*, vol. 49, no. 1, pp. 129–150, Jan. 2003.
- [6] N. J. A. Sloane and B. Beferull-Lozano, “Quantizing using lattice intersections,” in *Discrete and Computational Geometry*, ser. Algorithms Combin., Berlin: Springer-Verlag, 2003, vol. 25, pp. 799–824.
- [7] H. Bölcskei and F. Hlawatsch, “Noise reduction in oversampled filter banks using predictive quantization,” *IEEE Trans. Inf. Theory*, vol. 47, no. 1, pp. 155–172, Jan. 2001.
- [8] I. Daubechies, *Ten Lectures on Wavelets*. Philadelphia, PA: SIAM, 1992.
- [9] J. Benedetto and O. Treiber, “Wavelet frames: Multiresolution analysis and extension principles,” in *Wavelet Transforms and Time-Frequency Signal Analysis*, L. Debnath, Ed. Basel, Switzerland, 2001.
- [10] J. Munch, “Noise reduction in tight Weyl–Heisenberg frames,” *IEEE Trans. Inf. Theory*, vol. 38, no. 2, pp. 608–616, Mar. 1992.
- [11] I. Daubechies and R. DeVore, “Reconstructing a band-limited function from very coarsely quantized data: A family of stable sigma-delta modulators of arbitrary order,” *Ann. Math.*, vol. 158, no. 2, pp. 679–710, 2003.
- [12] C. Güntürk, J. Lagarias, and V. Vaishampayan, “On the robustness of single-loop Sigma–Delta modulation,” *IEEE Trans. Inf. Theory*, vol. 47, no. 5, pp. 1735–1744, Jul. 2001.

- [13] Ö. Yılmaz, "Stability analysis for several sigma-delta methods of coarse quantization of band-limited functions," *Constructive Approx.*, vol. 18, pp. 599–623, 2002.
- [14] —, "Coarse quantization of highly redundant time-frequency representations of square-integrable functions," *Appl. Comput. Harmonic Anal.*, vol. 14, pp. 107–132, 2003.
- [15] T. Strohmer and R. W. Heath Jr, "Grassmannian frames with applications to coding and communications," *Appl. Comput. Harmonic Anal.*, vol. 14, no. 3, pp. 257–275, 2003.
- [16] C. Güntürk, "Approximating a band-limited function using very coarsely quantized data: Improved error estimates in Sigma-Delta modulation," *J. Amer. Math. Soc.*, vol. 17, no. 1, pp. 229–242, 2004.
- [17] W. Chen and B. Han, "Improving the Accuracy Estimate for the First Order Sigma-Delta Modulator," preprint, 2003.
- [18] V. Goyal, J. Kovačević, and M. Vetterli, "Quantized frame expansions as source-channel codes for erasure channels," in *Proc. IEEE Data Compression Conf.*, Snowbird, UT, Mar. 1999, pp. 326–335.
- [19] G. Rath and C. Guillemot, "Performance analysis and recursive syndrome decoding of DFT codes for bursty erasure recovery," *IEEE Trans. Signal Process.*, vol. 51, no. 5, pp. 1335–1350, May 2003.
- [20] P. Casazza and J. Kovačević, "Equal-norm tight frames with erasures," *Adv. Comput. Math.*, vol. 18, no. 2/4, pp. 387–430, Feb. 2003.
- [21] G. Rath and C. Guillemot, "Recent advances in DFT codes based on quantized finite frames expansions for erasure channels," *Digital Sig. Process.*, vol. 14, no. 4, pp. 332–354, 2004.
- [22] V. Goyal, J. Kovačević, and M. Vetterli, "Multiple description transform coding: Robustness to erasures using tight frame expansions," in *Proc. IEEE Int. Symp. Information Theory*, Cambridge, MA, Aug. 1998, p. 408.
- [23] B. Hochwald, T. Marzetta, T. Richardson, W. Sweldens, and R. Urbanke, "Systematic design of unitary space-time constellations," *IEEE Trans. Inf. Theory*, vol. 46, no. 6, pp. 1962–1973, Sep. 2000.
- [24] Y. Eldar and G. Forney, "Optimal tight frames and quantum measurement," *IEEE Trans. Inf. Theory*, vol. 48, no. 3, pp. 599–610, Mar. 2002.
- [25] W. Bennett, "Spectra of quantized signals," *Bell Syst. Tech. J.*, vol. 27, pp. 446–472, Jul. 1948.
- [26] R. Duffin and A. Schaeffer, "A class of nonharmonic Fourier series," *Trans. Amer. Math. Soc.*, vol. 72, pp. 341–366, 1952.
- [27] J. J. Benedetto and M. W. Frazier, Eds., *Wavelets: Mathematics and Applications*. Boca Raton, FL: CRC, 1994.
- [28] O. Christensen, *An Introduction to Frames and Riesz Bases*. Boston, MA: Birkhäuser, 2003.
- [29] G. Zimmermann, "Normalized tight frames in finite dimensions," in *Recent Progress in Multivariate Approximation*, K. Jetter, W. Haussmann, and M. Reimer, Eds. Boston, MA: Birkhäuser, 2001.
- [30] J. Benedetto and M. Fickus, "Finite normalized tight frames," *Adv. Comput. Math.*, vol. 18, no. 2/4, pp. 357–385, Feb. 2003.
- [31] J. Candy and G. Temes, Eds., *Oversampling Delta-Sigma Data Converters: Theory, Design, and Simulation*. New York: Wiley-IEEE Press, 1992.
- [32] R. Gray, "Quantization noise spectra," *IEEE Trans. Inf. Theory*, vol. 36, no. 6, pp. 1220–1244, Nov. 1990.
- [33] S. Norsworthy, R. Schreier, and G. Temes, Eds., *Delta-Sigma Data Converters*. Piscataway, NJ: IEEE Press, 1997.
- [34] J. Benedetto, A. M. Powell, and Ö. Yılmaz, "Sigma-delta quantization and finite frames," in *Proc. Int. Conf. Acoustics, Speech and Signal Processing*, vol. 3, Montreal, QC, Canada, May 2004, pp. 937–940.
- [35] —, "Second order sigma-delta ( $\Sigma\Delta$ ) quantization of finite frame expansions," *Appl. Comput. Harmonic Anal.*, vol. 20, no. 1, pp. 126–148, 2006.
- [36] E. Janssen and D. Reefman, "Super-audio CD: An introduction," *IEEE Signal Process. Mag.*, vol. 20, no. 4, pp. 83–90, Jul. 2003.
- [37] L. Kuipers and H. Niederreiter, *Uniform Distribution of Sequences*. New York: Wiley-Interscience, 1974.
- [38] R. Schreier and G. Temes, *Understanding Delta-Sigma Data Converters*. New York: Wiley-IEEE Press, 2004.