

# An Exact FFT Recovery Theory: A Nonsubtractive Dither Quantization Approach with Applications

L. Cheded and S. Akhtar

*Systems Engineering Department, King Fahd University of Petroleum and Minerals,  
KFUPM Box 116, Dhahran 31261, Saudi Arabia*

Received 27 June 2004; Revised 13 September 2005; Accepted 26 September 2005

Recommended for Publication by Jar-Ferr Kevin Yang

Fourier transform is undoubtedly one of the cornerstones of digital signal processing (DSP). The introduction of the now famous FFT algorithm has not only breathed a new lease of life into an otherwise latent classical DFT algorithm, but also led to an explosion in applications that have now far transcended the confines of the DSP field. For a good accuracy, the digital implementation of the FFT requires that the input and/or the 2 basis functions be finely quantized. This paper exploits the use of coarse quantization of the FFT signals with a view to further improving the FFT computational efficiency while preserving its computational accuracy and simplifying its architecture. In order to resolve this apparent conflict between preserving an excellent computational accuracy while using a quantization scheme as coarse as can be desired, this paper advances new theoretical results which form the basis for two new and practically attractive FFT estimators that rely on the principle of 1 bit nonsubtractive dithered quantization (NSDQ). The proposed theory is very well substantiated by the extensive simulation work carried out in both noise-free and noisy environments. This makes the prospect of implementing the 2 proposed 1 bit FFT estimators on a chip both practically attractive and rewarding and certainly worthy of a further pursuit.

Copyright © 2006 Hindawi Publishing Corporation. All rights reserved.

## 1. INTRODUCTION

The vast success of the Fourier transform is amply reflected in the wide applicability it enjoys in a variety of engineering fields such as signal and image processing, control, communications, filtering, geophysics, seismics, optics, acoustics, radar, and sonar signal processing. This explosion in applications was brought about by the introduction of the now famous and ubiquitous fast Fourier transform (FFT) which has transformed the classical discrete Fourier transform (DFT) from being a mere “academic” curiosity, with limited applications, to being a powerful computational tool whose applications continue to grow unabatedly [1]. The original radix-2 structure of the FFT underwent several structural changes all aiming at further increasing the computational speed and/or adapting the original FFT algorithm to various data length characteristics (e.g., prime and composite lengths) [2]. A contemporary view as well as a review of the state of the art of the FFT can be found in [3, 4], respectively.

The numerous variations of the original radix-2 FFT algorithm were brought about through the dual use of the exploitation of symmetry properties inherent in the FFT algorithm and the principle of “divide and conquer.” How-

ever, both the original FFT algorithm and all its existing variants rely, in their conventional digital implementation, on input signals that are sufficiently highly quantized (i.e., resolution  $\geq 8$  bits). In the practical implementation of a digital signal processing (DSP) system, the user has to minimize what is commonly known as the finite wordlength effects, otherwise these will introduce noise into the designed system and lead to nonideal, if not unreliable, system responses. The processes generating these adverse effects are classified into the following 4 categories: (1) input quantization, (2) coefficient quantization, (3) overflow (or underflow) in internal arithmetic operations, and (4) rounding (or truncation) of data for storage in memory or register. In this paper, we only focus on the first process (input quantization) that is carried out by the analog-to-digital converter (ADC) and briefly discuss its effect on the accuracy of both input coding and FFT estimation. It is well known that a quantizer, which is part of an ADC, with input  $x$  and output  $x_Q$ , introduces an error known as the quantization error (or noise) and defined by:  $e_Q = x - x_Q$ . Given a  $B$ -bit quantizer with an input whose peak-to-peak (or full scale) range is  $V_{PP}$ , then the quantizer's step is given by  $q = V_{PP}/2^B$ . Provided that  $B$  is sufficiently large,  $e_Q$  will then behave like an additional white noise that

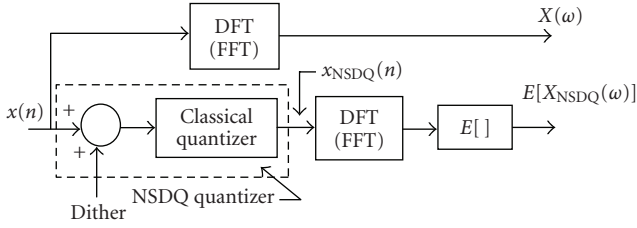


FIGURE 1: MR-FFT estimation scheme.

is uncorrelated with the quantizer's input, uniformly distributed (UD) over the range  $[-q/2, q/2]$ , is zero-mean and has a variance of  $\sigma_e^2 = q^2/12$ . One of the performance measures of a quantizer is its signal-to-quantization-noise ratio (SQNR) defined by:  $SQNR = 10 \log_{10}(P_x/\sigma_e^2)$ , where  $P_x$  is the input power. In the case of an input sinewave, it can be shown [5] that:  $SQNR = (6.02B + 1.76)$  dB. This equation which provides a good basis as a design guideline reveals the interesting fact that a 1 bit increase in the quantizer's resolution ( $B$ ) leads to a 6 dB gain in its SQNR and hence in its dynamic range. A further improvement to the quantizer's SQNR can be achieved through the use of the oversampling technique which ensures a 6 dB improvement in the SNQR for an oversampling factor of 4 (see [5] for further details). The effect of a  $B$ -bit input quantization on a decimation-in-time (DIT) radix-2 FFT algorithm of length  $N$  was studied in [5] and showed that the total noise variance is given by:  $\sigma_T^2 = (N - 1)2^{-2B}/3$ . This clearly shows that as the quantization resolution  $B$  increases, the noise variance decreases and hence the FFT estimation accuracy improves. However, this improvement is gained at a cost of an increase in system complexity, implementational cost, and computational load. If these 3 system characteristics are to be reduced to any desired level while preserving a good FFT estimation accuracy, then the conventional approach, as described above, offers no flexible solution at all since low complexity (achievable with low quantization resolutions) and good accuracy (achievable with high quantization resolutions) are clearly 2 incompatible requirements. This fact is clearly borne out by the results of Figures 3 and 4 which depict the degradation in performance of 2 FFT estimators using the lowest possible (i.e., 1 bit) quantization resolution.

Although low quantization resolutions entail an irreversible loss of accuracy which becomes more prohibitive as the resolution gets smaller, they nevertheless offer several practically attractive advantages primarily associated with the use of shorter wordlengths. Such practical advantages include structurally simple, low-cost, and fast FFT processing schemes that can only enhance the already high speed boasted by existing FFT algorithms. These advantages will in turn lead to the possibility of a fast fully parallel FFT algorithm that can be cost effectively implemented using, for example, FPGA technology. However, in order to unlock all of these important potential practical advantages, a way to reconcile two seemingly disparate requirements, namely, achieving high accuracy in FFT processing while using only coarsely quantized signals, has to be found.

The main objective of this paper is therefore to propose a new and practical solution to this problem, in the form of a new exact FFT recovery theory which forms the theoretical basis for two new and fast FFT estimators: a modified relay FFT estimator and a modified polarity coincidence FFT estimator, referred to henceforth as the MR-FFT and the MPC-FFT estimators, respectively. These 2 estimators have the unique feature of permitting signal quantization resolution as low as 1 bit while incurring only an acceptable small loss in FFT estimation accuracy. At the heart of this new solution lies the exact moment theory (EMR) which itself hinges upon a conceptually simple signal coding scheme based on the nonsubtractive dithered quantization (NSDQ) technique [6] to be described below. Other related studies discussing dithered quantization can be found in [7, 8]. However, unlike these 2 studies, our work of [6] focuses on the exact recovery of any existing finite-order moments of the dithered quantizer's input from those of its output. It is this precise feature of our work of [6] that is exploited and extended here. It is of vital importance to point out here that, in addition to being assumed stationary, all the signals used in this paper are also assumed to be ergodic so as to justify the equivalence between the ensemble averages upon which rely all of the theoretical derivations in our approach, which is essentially stochastic in nature, and the time averages used in our simulation work.

The block diagrammatic description of the MR-FFT is shown above in Figure 1. Here, only the input signal,  $x(n)$ , whose FFT spectrum is to be estimated, is fed into the NSDQ quantizer. From this figure, it is clear that the NSDQ scheme is basically equivalent to a classical uniform quantization whose input has been dithered by a dither signal with certain specific statistical characteristics to be discussed later. In order to reap the maximal benefits from this flexible architecture, we therefore need to use the crudest possible (i.e., 1 bit) NSDQ scheme. In this scheme, the 2 multiplications required are between the quantized version of the dithered input, that is,  $x_{NSDQ}(n)$ , and the 2 FFT basis functions "cos" and "sin" (not shown but included in the block-labeled DFT (FFT) in Figure 1). When 1 bit NSDQ quantization is used, as is the case in our proposed MR-FFT scheme,  $x_{NSDQ}(n)$  will simply be a random binary signal which, when multiplied with the 2 basis functions, will in effect be switching them on and off. Because this technique of implementing a multiplication as a mere switching operation is commonly found in relays, the 2 multiplications required in our proposed scheme of Figure 1 are therefore analogous to 2 relay-type multiplications. Since the switching signal  $x_{NSDQ}(n)$  is derived from a modified (here dithered) version of the input  $x(n)$ , the resulting estimator is thus called a modified relay FFT (MR-FFT).

As to the architecture of the second proposed estimator, it is shown in Figure 2 below where, as clearly shown, all of the 3 signals involved, that is, the input  $x(n)$  and the 2 real basis signals  $s(n)$  and  $c(n)$  which make up the Fourier complex kernel ( $K(n, \omega_i) = e^{-j\omega_i n}$ ), are now each NSDQ-quantized. Each of the 3 required NSDQ quantizers in Figure 2 has exactly the same internal architecture as the one shown above in Figure 1. Here too, maximal benefits are obtained when

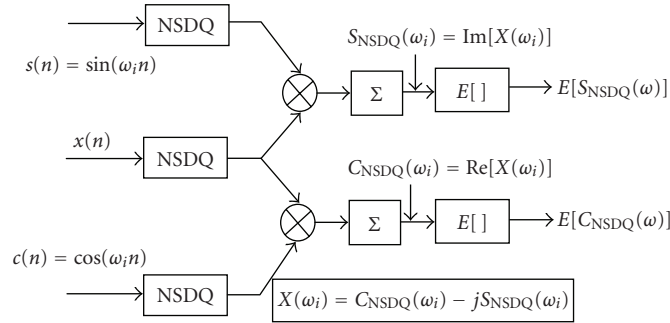


FIGURE 2: MPC-FFT estimation scheme.

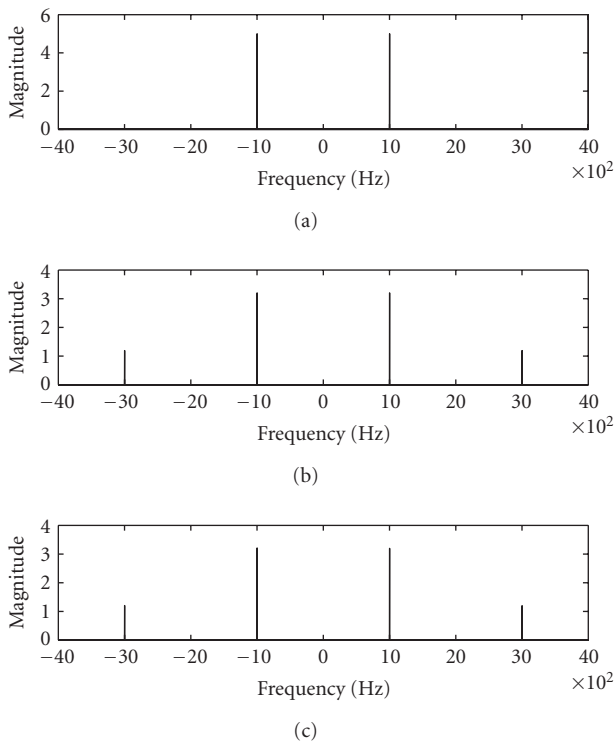


FIGURE 3: FFT magnitude spectra of a single sinusoid: original (true) spectrum (top), estimated with R-FFT estimator (middle) and with PC-FFT estimator (bottom).

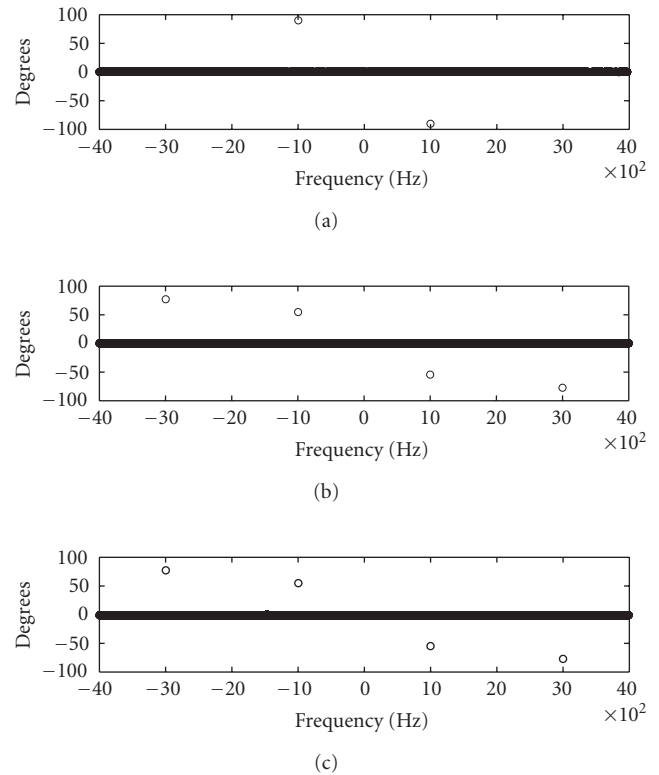


FIGURE 4: FFT phase spectra of a single sinusoid: original (true) spectrum (top), estimated with R-FFT estimator (middle) and with PC-FFT estimator (bottom).

1 bit resolution is used in all 3 NSDQ quantizers since in this case the 2 required multiplications are reduced to simple polarity coincidence-type of multiplications between the 2 pairs of modified (here dithered) signals, hence the name of modified polarity-coincidence FFT (MPC-FFT) given to the resulting estimator.

It is worth pointing out here that the preliminary tests of these 2 FFT estimators proved successful in both noise-free and moderately noisy environments [9–12]. Moreover, the theory underlying the 2 proposed estimators can be interpreted as a frequency-domain extension of the aforementioned EMR theory of [6] which has enjoyed other successful applications [13–15].

This paper is organized as follows: Section 2 introduces only some relevant fundamental results of the EMR theory in its 1-D setting and shows how a key theorem (Theorem 1) can be used to furnish the first proposed estimator (MR-FFT). In Section 3, the 2-D extension of the results presented in Section 1 is given, leading to another key theorem (Theorem 2) which is shown therein to lead to the second proposed estimator (MPC-FFT). Section 4 presents some simulation results which demonstrate the very good performance of the 2 proposed 1 bit FFT estimators, using both simulated signals as well as recordings of real signals. Finally, some concluding remarks are given in Section 5.

## 2. ONE-DIMENSIONAL EMR THEORY: FUNDAMENTAL RESULTS AND APPLICATION TO THE MODIFIED RELAY (MR)-FFT ESTIMATION

Some fundamental results of the EMR theory are presented here and a new theorem (Theorem 1) is derived on the moment-sense equivalence between the NSDQ-based DFT and a frequency-domain mapping to be defined below in Section 2.4.

### 2.1. Definition of the NSDQ quantization scheme

Given an input  $x$  and a (user-defined) dither signal  $D$  that is statistically independent of  $x$ , then a nonsubtractively dithered quantization (NSDQ) of  $x$  is equivalent to the classical quantization ( $Q_a$ ) of the dithered signal  $y = x + D$ , that is,

$$x \rightarrow x_{\text{NSDQ}} = \text{NSDQ}(x) = Q_a(y) = y_Q. \quad (1)$$

Here,  $Q_a$  represents the entire class of uniform classical quantizers parametrized by the step ( $q$ ) and the shift factor  $a \in [-1/2, 1/2)$ , that is,

$$y_Q = \left(a + l + \frac{1}{2}\right)q \quad \text{if } y \in [(a + l)q, (a + l + 1)q). \quad (2)$$

Note here that the 2 well-known classical quantizers, namely, the mid-riser (without any dead zone) and mid-stepper (with dead zone), correspond to  $a = 0$  and  $a = -1/2$ , respectively. The quantizers used in this study are all of the mid-riser type.

### 2.2. Definition of the $p$ th-order class of linearizing dither signals $\mathcal{D}_p$

Given an ergodic and stationary dither signal  $D$  and its characteristic function  $W_D(u)$ , then

$$D \in \mathcal{D}_p \iff W_D^{(r)}\left(\frac{2m\pi}{q}\right) = 0 \quad \forall r \in [0, p-1], m \neq 0. \quad (3)$$

A detailed discussion as to the origin of this definition can be found in [6]. However, it suffices here to use it and show that it holds the key to the solution of the problem of exactly recovering the FFT spectrum of a given signal in an NSDQ quantization setting. It is also interesting to note here that this definition requires only that the characteristic functions of the dither signal have a set of equispaced zeroes (except at the origin), with the constant spacing controlled by the uniform NSDQ quantizer's step. We cite here 3 types of member signals of  $\mathcal{D}_p$ : the basic uniformly distributed (UD) dither signal (type 1), a signal formed with any finite number of statistically independent UD dither signals (type 2), and a signal formed with a sum of at least one type-1 or type-2 member signal and any finite number of statistically independent signals that are not necessarily members of  $\mathcal{D}_p$ . The last 2 types owe their existence to the closure property discussed next.

According to the closure property of  $\mathcal{D}_p$  [6], we can say that if  $D \in \mathcal{D}_p$  and for any signal  $x$  that is statistically independent of  $D$ , then the dithered signal  $y = (x + D) \in \mathcal{D}_p$ . Note that although the proof of this property for the important case of  $p = 1$  treated in this paper can be straightforwardly carried out here, we have nevertheless provided it in Appendix D where the proof of the general  $p$ th-case version of the closure property is also carried out. The version of the closure property with more than 2 signals is discussed in [6].

### 2.3. Statistical characterization of NSDQ: the $p$ th-order moment-sense input/output function

As shown in [6], every NSDQ quantizer is statistically characterized by a special function called the  $p$ th-order moment-sense input/output function (MSIOF) and denoted by  $h_p(x)$ . The following important lemma, proved in [6], shows precisely the role played by this function in the EMR theory.

**Lemma 1.** *A uniform NSDQ quantizer of step  $q$ , dither signal  $D$  and shift factor  $a$  where  $a \in [-1/2, 1/2)$ , is equivalent, from a  $p$ th-order moment point of view, to a transformation  $h_p(x)$ , henceforth called the quantizer's MSIOF, which satisfies the following relationship:*

$$m_{\text{NSDQ}_p} \triangleq E[x_{\text{NSDQ}}^p] = E[h_p(x)] \quad \forall p \geq 1, \quad (4)$$

where

$$h_p(x) = \sum_{k=0}^p c_k x^k, \quad (5)$$

$$c_k = \sum_{t=0}^{p-k} \frac{p!}{(p-k-t+1)!k!t!} \left(\frac{q}{2}\right)^{p-k-t} E[D^t][P \oplus k \oplus t \oplus 1],$$

with  $\oplus$  denoting modulo-2 operation.

Note here that for the important case of  $p = 1$ , the resulting first-order MSIOF becomes perfectly linearized, that is,  $h_1(x) = x$ , as shown in Appendix E.

### 2.4. A key theorem on the derivation of the MR-FFT estimator

We will now state and prove a general theorem on the exact recovery of the DFT of any finite-energy signal, that is,  $x^p$ , from the DFT of its NSDQ-quantized version, that is,  $x_{\text{NSDQ}}^p$ , regardless of the quantization resolution used.

**Theorem 1.** *Given (1) a 1-D NSDQ quantizer whose  $p$ th-order MSIOF, input, output, and dither signals are given, respectively, by  $h_p(x)$ ,  $x$ ,  $x_{\text{NSDQ}}$ , and  $D$  where  $D$  is both zero-mean and statistically independent of  $x$ , and (2) the following (DFT) spectra: the quantizer's input  $p$ th-order DFT defined by:  $X_p(\omega_i) \triangleq \sum_{n=0}^{N-1} x^p(n) \cdot K(n, \omega_i)$ , and the corresponding quantizer's output  $p$ th-order DFT, also called here the 1-quantized channel  $p$ th-order DFT and defined by:  $X_{\text{NSDQ}_p}^{[1]}(\omega_i) \triangleq \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) \cdot K(n, \omega_i)$ , for all  $i \in [1, N]$ ,*



where the complex DFT kernel is defined by:  $K(n, \omega_i) \triangleq e^{-j\omega_i n}$ , then  $X_{\text{NSDQ}_p}^{[1]}(\omega_i)$  is moment-sense equivalent to a  $p$ -D frequency-domain mapping  $H_p(\omega_i)$  defined below, that is,

$$E[X_{\text{NSDQ}_p}^{[1]}(\omega_i)] = E[H_p(\omega_i)],$$

$$\text{where } H_p(\omega_i) \triangleq \text{DFT}[h_p(x(n))] = \sum_{k=0}^{p-1} c_k X_k(\omega_i), \quad (6)$$

$$\forall i \in [1, N],$$

and the coefficient  $c_k$  is as defined in (5).

*Proof (see Appendix A).* It is also important to note here that exact recovery of a spectrum of a high-order signal, say  $X_p(\omega_i)$ , for all  $p > 1$ , would require 1 NSDQ quantizer but the estimation of  $p$  different NSDQ-quantized spectra,  $X_{\text{NSDQ}_k}^{[1]}(\omega_i)$ , for all  $k \in [1, p]$ . An attractive alternative to this would instead require the use of  $p$  different NSDQ quantizers, each with its own dither signal being statistically independent of all the inputs and other dither signals, and the estimation of only 1  $p$ -D NSDQ-quantized spectrum. However this would involve the use of a  $p$ -D EMR theory which is outside the scope of this paper.  $\square$

### 2.5. Application of Theorem 1 to the MR-FFT estimation

If we now let  $p = 1$  in the two  $p$ th-order mappings,  $h_p(x)$  and its DFT  $H_p(\omega_i)$ , their resulting first-order expressions would then simplify to

$$h_1(x) = x \implies H_1(\omega_i) = \text{DFT}[h_1(x)] = X(\omega_i). \quad (7)$$

It is clear from (7) that (a) the first-order MSIOF,  $h_1(x)$ , represents nothing but the perfectly linearized average input/output (I/O) characteristics of the NSDQ quantizer. In the absence of dither,  $h_1(x)$  reduces to the well-known staircase-like I/O function of the classical (i.e., undithered) quantizer  $Q_a(x)$ . And (b), the first-order frequency-domain mapping,  $H_1(\omega_i)$ , gives directly the desired input spectrum  $X(\omega_i)$ .

Note here that, according to Theorem 1 and for  $p = 1$ , the 1-quantized channel first-order DFT is nothing but the MR-FFT spectrum of the input  $x$  since we have  $X_{\text{NSDQ}_1}^{[1]}(\omega_i) \triangleq \text{DFT}[x_{\text{NSDQ}_1}^1] = X_{\text{NSDQ}_1}^{[1]}(\omega_i)$ . Thus combining (6) and (7) gives

$$E[X_{\text{NSDQ}_1}^{[1]}(\omega_i)] = E[H_1(\omega_i)] = E[X(\omega_i)]. \quad (8)$$

Note here that when the NSDQ quantizer's input,  $x(n)$ , is deterministic, then (8) reduces to  $E[X_{\text{NSDQ}_1}^{[1]}(\omega_i)] = X(\omega_i)$ . This states that the DFT itself, rather than its average, of the NSDQ quantizer's input can be exactly recovered from the average of the DFT of the quantizer's output, irrespective of the quantization resolution used. This result has tremendous practical benefits since this exact recovery is now possible even with 1 bit resolution as tested in our simulation work. Moreover, in practice and as pointed out in Section 4, we can dispense with the use of a separate expectation operation on the (1 bit)

quantized DFT,  $X_{\text{NSDQ}_1}^{[1]}(\omega_i)$ , since the DFT operation itself involves some form of averaging. Note also that should the NSDQ quantizer's input be noisy, say  $x(n) = x_0(n) + v(n)$  where  $x_0(n)$  and  $v(n)$  are the deterministic and noisy components, respectively, and provided the noise signal  $v(n)$  is both zero-mean and statistically independent of the input and dither signals used, then exact recovery of the DFT of the deterministic component will also be possible.

## 3. TWO-DIMENSIONAL EXTENSION OF EMR THEORY: FUNDAMENTAL RESULTS AND APPLICATION TO THE MODIFIED POLARITY-COINCIDENCE (MPC)-FFT ESTIMATION

As the development of the MPC-FFT estimator would require the NSDQ quantization of 2 channels as indicated in Figure 2, there is therefore a need to extend the 1-D EMR theory of the previous section to its 2-D counterpart. This section will introduce some fundamental results that emanated from such an extension.

### 3.1. Two-dimensional definition of NSDQ

Given a 2-D input vector  $\underline{x} = (x_1, x_2)^T$  and a user-defined 2-D dither vector  $\underline{D} = (D_1, D_2)^T$  which is component-wise statistically independent of  $\underline{x}$ , then a 2-D nonsubtractively dithered quantization (NSDQ) of  $\underline{x}$  is equivalent to the classical quantization ( $Q_a$ ) of the dithered 2-D vector  $\underline{y} = \underline{x} + \underline{D}$ , that is,

$$\underline{x} \longmapsto \underline{x}_{\text{NSDQ}} = \text{NSDQ}(\underline{x}) = Q_a(\underline{x}) = \underline{x}_Q. \quad (9)$$

Here,  $Q_a$  represents the entire class of uniform classical quantizers parametrized by the 2-D uniform quantization step  $\underline{q} = (q_1, q_2)^T$  and the shift factor vector  $\underline{a} = (a_1, a_2)^T$  where  $a_i \in [-1/2, 1/2)$  for  $i = 1, 2$ , that is,

$$y_{Q_i} = \left( a_i + l_i + \frac{1}{2} \right) q_i \quad (10)$$

$$\text{if } y_i \in [(a_i + l_i)q_i, (a_i + l_i + 1)q_i) \text{ for } i = 1, 2.$$

The 2-D mid-riser and mid-stepper quantizers are defined here by  $\underline{a} = (0, 0)^T$  and  $\underline{a} = (1/2, 1/2)^T$ , respectively. Here too, all 2-D quantizers used in this paper are of the mid-riser type.

### 3.2. Definition of the 2-D( $p_1, p_2$ )th-order class of linearizing dither signals $\mathcal{D}_{p_1, p_2}$

Given an ergodic and stationary dither vector  $\underline{D} = (D_1, D_2)^T$  and its characteristic function (CF)  $W_{\underline{D}}(u_1, u_2)$ , then

$$\underline{D} \in \mathcal{D}_{p_1, p_2} \iff W_{\underline{D}}^{(r_1, r_2)} \left( \frac{2m_1\pi}{q_1}, \frac{2m_2\pi}{q_2} \right) = 0 \quad (11)$$

$$\forall r_i \in [0, p_i - 1], m_i \neq 0 \text{ for } i = 1, 2.$$

Note here that if either  $n_1$  or  $n_2$  is allowed to go to infinity, then the definition of the 1-D  $p$ th-order class of linearizing dither signals  $\mathcal{D}_p$ , as given above in Section 2, is immediately obtained.

Moreover, note that if the component signals  $D_1$  and  $D_2$  are statistically independent of each other, then the 2-D class  $\mathcal{D}_{p_1, p_2}$  becomes separable, that is,  $\mathcal{D}_{p_1, p_2} = \mathcal{D}_{p_1} \times \mathcal{D}_{p_2}$ . This important case will be exploited later in our simulation work.

### 3.3. Statistical characterization of 2-D NSDQ: the 2-D $(p_1, p_2)$ th-order moment-sense input/output function

We will introduce here a new lemma (Lemma 2) which characterizes the 2-D NSDQ quantizer by a new  $(p_1, p_2)$ th-order statistical I/O function called the quantizer's  $(p_1, p_2)$ th-order moment-sense input/output function (MSIOF).

**Lemma 2.** A uniform 2-D NSDQ quantizer of step vector  $\underline{q} = (q_1, q_2)^T$  and shift factor vector  $\underline{a} = (a_1, a_2)^T$ , where  $a_i \in [-1/2, 1/2)$  for  $i = 1, 2$ , is equivalent, from a  $(p_1, p_2)$ th-order moment point of view, to a mapping  $h_{p_1, p_2}(x_1, x_2)$ , henceforth called the quantizer's  $(p_1, p_2)$ th-order MSIOF, which satisfies the following relationship:

$$\begin{aligned} m_{\text{NSDQ}_{12}} &\triangleq E[x_{\text{NSDQ}_1}^{p_1} x_{\text{NSDQ}_2}^{p_2}] \\ &= E[h_{p_1, p_2}(x_1, x_2)] \quad \forall p_1, p_2 \geq 1, \end{aligned} \quad (12)$$

where

$$\begin{aligned} h_{p_1, p_2}(x_1, x_2) &= \sum_{l_1} \sum_{l_2} \left[ \left( a_1 + l_1 + \frac{1}{2} \right) q_1 \right]^{p_1} \left[ \left( a_2 + l_2 + \frac{1}{2} \right) q_2 \right]^{p_2} \\ &\times \left\{ P_{\underline{D}}[(a_1 + l_1 + 1)q_1 - x_1, (a_2 + l_2 + 1)q_2 - x_2] \right. \\ &\quad - P_{\underline{D}}[(a_1 + l_1)q_1 - x_1, (a_2 + l_2 + 1)q_2 - x_2] \\ &\quad - P_{\underline{D}}[(a_1 + l_1 + 1)q_1 - x_1, (a_2 + l_2)q_2 - x_2] \\ &\quad \left. + P_{\underline{D}}[(a_1 + l_1)q_1 - x_1, (a_2 + l_2)q_2 - x_2] \right\}. \end{aligned} \quad (13)$$

*Proof.* (See [15, Appendix 2] and note that the summation over  $l_1$  and  $l_2$  in (13) range from  $-\infty$  to  $+\infty$ .)  $\square$

#### Important property of separability

It is important to point out at this stage that if the 2 dither signals used in the 2-D NSDQ quantizer are statistically independent of each other and of the 2 quantizer's inputs, then the 2-D  $(p_1, p_2)$ th-order MSIOF becomes separable into its two 1-D  $(p_1)$ th- and  $(p_2)$ th-order MSIOFs, whose expressions are given by (5), that is,

$$h_{p_1, p_2}(x_1, x_2) = h_{p_1}(x_1)h_{p_2}(x_2). \quad (14)$$

This important property provides the user with an easy and effective practical way of implementing any multidimensional (m-D) NSDQ quantizer with a set of m 1-D NSDQ quantizers. This property is fully exploited in our simulation work.

### 3.4. A key theorem on the derivation of the MPC-FFT estimator

We will now state and prove a new theorem which guarantees, irrespective of the quantization resolution used, the

exact recovery of the  $p$ th-order DFT of a signal from a 2-channel quantized  $p$ th-order DFT estimation scheme which involves NSDQ quantizing both the input and the DFT kernel (or equivalently the 2 basis functions). It is worth pointing out at this juncture that the MPC-FFT estimation scheme represents a quadrature estimation of the DFT as it involves 2 basis functions that have a quadrature relationship in that their phases differ by  $\pi/2$ .

**Theorem 2.** Given (1) a 2-D vector NSDQ quantizer, characterized by its 2 signal triplets  $(x_l, x_{\text{NSDQ}_l}, D_l)$ ,  $l = 1, 2$ , where the 2 dither signals  $D_1$  and  $D_2$  are both zero-mean and statistically independent of each other and of the input signals  $x_1$  and  $x_2$ , and whose 2-D  $(p_1, p_2)$ th-order MSIOF is  $h_{p_1, p_2}(x_1, x_2)$ , and (2) the NSDQ quantizer's input  $p$ th-order DFT defined by:  $X_p(\omega_i) \triangleq \sum_{n=0}^{N-1} x^p(n) \cdot K(n, \omega_i)$  and the corresponding 2-quantized channel  $p$ th order DFT, which involves quantizing both the input and the DFT kernel and which is defined by:  $X_{\text{NSDQ}_p}^{[2]}(\omega_i) \triangleq \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) \cdot K_{\text{NSDQ}}(n, \omega_i)$ , where  $K_{\text{NSDQ}}(n, \omega_i) = (e^{-j\omega_i n})_{\text{NSDQ}}$  and  $i \in [1, N]$ , then  $X_{\text{NSDQ}_p}^{[2]}(\omega_i)$  is moment-sense equivalent to a  $p$ -D frequency-domain mapping  $H_p(\omega_i)$  defined below, that is,

$$E[X_{\text{NSDQ}_p}^{[2]}(\omega_i)] = E[H_p(\omega_i)], \quad (15)$$

where  $H_p(\omega_i) \triangleq \text{DFT}[h_p(x(n))] = \sum_{k=0}^p c_k X_k(\omega_i)$ , for all  $i \in [1, N]$  and the coefficient  $c_k$  is as defined in (5).

*Proof* (see Appendix B). It is easy to see that the DFT kernel has the following Cartesian expression  $K(n, \omega_i) \triangleq e^{j\omega_i n} = c(n) - js(n)$  where  $c(n)$  and  $s(n)$  are simply the basis (cosine and sine) functions shown in Figure 2. As such, the NSDQ-quantized version of this kernel is given by  $K_{\text{NSDQ}}(n, \omega_i) \triangleq (e^{j\omega_i n})_{\text{NSDQ}} = c_{\text{NSDQ}}(n, \omega_i) - js_{\text{NSDQ}}(n, \omega_i)$ , which indicates why in practice 2 NSDQ quantizers are required to quantize this complex kernel, as clearly shown in Figure 2.

As Theorem 2 addresses the exact recovery of the DFT of a particular signal using 2 NSDQ-quantized channels, it clearly represents a 2-D generalization of Theorem 1 which addresses the same problem using only 1 NSDQ-quantized channel.  $\square$

### 3.5. Application of Theorem 2 to the MPC-FFT estimation

Proceeding along similar lines to those in Section 2.5, and since the same signal-domain and frequency-domain mappings, that is,  $h_p(x)$  and  $H_p(\omega_i)$ , respectively, are involved here as well, it then becomes clear that by letting  $p = 1$  in the general expressions of these 2 mappings, both  $h_1(x)$  and  $H_1(\omega_i)$  will assume their respective simplified expressions given in (7).

Combining (7) and (14) leads directly to the desired result:

$$E[X_{\text{NSDQ}}^{[2]}(\omega_i)] = E[H_1(\omega_i)] = E[X(\omega_i)]. \quad (16)$$

Here too, for a deterministic signal  $x(n)$ , we will have from

(15):  $E[X_{\text{NSDQ}}^{[2]}(\omega_i)] = X(\omega_i)$  which shows that in this particular case, it is the DFT itself, rather than its average, of the NSDQ quantizer's input which will be exactly recovered from the average of the DFT of the NSDQ quantizer's output, irrespective of the quantization resolution used. The same remark, made in Section 2.5, on the dispensation with the expectation operation in the estimation scheme also applies here to the MPC-FFT estimator. In the event that  $x(n)$  is noisy and provided that its noisy component is both zero-mean and statistically independent of the dither signals used, then exact recovery of the DFT of the deterministic component will also be possible. In either case, the MPC-FFT estimator offers far greater practical advantages than its MR-FFT counterpart since its practical implementation is purely digital (as opposed to the hybrid one for the MR-FFT estimator), involves the processing of 1 bit (binary) signals only and hence would require only 1 bit logic devices for its multiply-and-accumulate operation.

### 3.6. Remarks on some statistical properties of the 2 proposed estimators

#### 3.6.1. Unbiasedness and consistency

Given a random variable (RV)  $Y$ , its true mean  $\mu_Y = E[Y]$  and its sample mean estimator  $\hat{Y} \triangleq (1/K) \sum_{k=0}^{K-1} Y_k$ , it is

well known [16] that the sample mean estimator is an unbiased and consistent estimator of the true mean. In our case and for each discrete frequency  $\omega_i$ , the RVs are represented by the samples of the NSDQ-quantized spectra which are  $X_{\text{NSDQ}}^{[1]}(\omega_i)$  (for MR-FFT) and  $X_{\text{NSDQ}}^{[2]}(\omega_i)$  (for MPC-FFT). In our simulation, the true mean of these quantized RVs, that is,  $E[X_{\text{NSDQ}}^{[1]}(\omega_i)]$  and  $E[X_{\text{NSDQ}}^{[2]}(\omega_i)]$ , are respectively estimated by the following sample mean estimators,  $\hat{X}_{\text{NSDQ}}^{[1]}(\omega_i) \triangleq (1/K) \sum_{k=0}^{K-1} X_{\text{NSDQ}_k}^{[1]}(\omega_i)$  and  $\hat{X}_{\text{NSDQ}}^{[2]}(\omega_i) \triangleq (1/K) \sum_{k=0}^{K-1} X_{\text{NSDQ}_k}^{[2]}(\omega_i)$ . As pointed out above, these sample mean estimators are therefore unbiased and consistent estimators of their respective true means, namely,  $E[X_{\text{NSDQ}}^{[1]}(\omega_i)]$  and  $E[X_{\text{NSDQ}}^{[2]}(\omega_i)]$ . Moreover, since (8) and (16) show that each of these 2 true means is in fact equal to the desired true mean of the unquantized spectrum, that is,  $E[X(\omega_i)]$ , it then follows that the 2 sample mean estimators used in our simulation, that is,  $\hat{X}_{\text{NSDQ}}^{[1]}(\omega_i)$  and  $\hat{X}_{\text{NSDQ}}^{[2]}(\omega_i)$ , are unbiased and consistent estimators of the desired true mean  $E[X(\omega_i)]$ .

#### 3.6.2. Variance analysis

According to Appendix C, the variance expression for the MD-FFT and MH-FFT estimators are given by

$$\sigma_{\text{MD-FFT}}^2 = \sigma_{\text{SD-FFT}}^2 + \underbrace{\frac{1}{K} \left[ \sum_{n=0}^{N-1} E \left[ (x_{\text{NSDQ}}^p(n) | K_{\text{NSDQ}}(n, \omega_i) |)^2 - (x^p(n) | K(n, \omega_i) |)^2 \right] \right]}_{\text{MD-FFT excess variance}}, \quad (17)$$

$$\sigma_{\text{MH-FFT}}^2 = \sigma_{\text{SD-FFT}}^2 + \underbrace{\frac{1}{K} \left[ \sum_{n=0}^{N-1} E \left[ (x_{\text{NSDQ}}^p(n) | K(n, \omega_i) |)^2 - (x^p(n) | K(n, \omega_i) |)^2 \right] \right]}_{\text{MH-FFT excess variance}}. \quad (18)$$

First note that the 2 excess-variance terms, involved in both (17) and (18), account solely for the contribution of NSDQ quantization to the variance of each of the 2 quantized estimators. This fact can be easily checked from both (17) and (18) since these extra terms vanish in the absence of NSDQ quantization. These extra terms also vanish if an infinite number of spectrum estimates is used (i.e., if  $K \rightarrow \infty$ ). This last fact then reveals that both the MD-FFT and MH-FFT estimators are 2 equally asymptotically efficient estimators. However, the rate at which the variance of the 3 estimators (i.e., SD-FFT, MH-FFT, and MD-FFT) converges to zero is the smallest for the MD-FFT and highest for the SD-FFT, as expected.

In terms of the relative sizes of these variance excesses, we have obtained new results to be reported later, which show that, in the general setting of multibit, multivariable NSDQ-quantized FFT estimators, the variance excess due to the MH-FFT estimator is smaller than that due to the MD-

FFT one. This is to be expected as the MD-FFT estimator involves more quantization, and hence more distortion and quantization error, and a higher excess in variance, than the MH-FFT one.

The above generalized variance expressions of (17) and (18) can now be applied to the 2 proposed FFT estimators, that is, the MR-FFT and the MPC-FFT, which are merely 1 bit versions of the MH-FFT and MD-FFT estimators, respectively. If the gains of all of the 1 bit NSDQ quantizers used in the proposed estimators are set to  $(\pm q/2)$ , then the corresponding variance expressions of these 1 bit estimators are obtained as explained in the following. As the MD-FFT estimator consists of 2 channels (cosine and sine) whose estimates are uncorrelated with each other, the total impact of NSDQ quantization on the variance of this estimator, represented by the first summation term on the RHS of (17), will therefore be made of the sum of similar impacts emanating from both channels, namely,  $\sum_{n=0}^{N-1} (x_{\text{NSDQ}}^p(n) c_{\text{NSDQ}}(n, \omega_i))^2$

for the cosine channel and  $\sum_{n=0}^{N-1} (x_{\text{NSDQ}}^p(n) s_{\text{NSDQ}}(n, \omega_i))^2$  for the sine channel. Since, for the MPC-FFT estimator, all the dithered signals (the input and the 2 real basis functions) are clipped at  $\pm q/2$ , it can be easily shown that the quantization impacts from both channels are each equal to  $(q^4 N/16)$  and that the combined impact of both channels is twice that amount. In view of this, (17) now becomes

$$\sigma_{\text{MD-FFT}}^2 = \sigma_{\text{SD-FFT}}^2 + \underbrace{\frac{q^4 N}{8K} - \frac{1}{K} \sum_{n=0}^{N-1} E[(x^p(n) |K(n, \omega_i)|)^2]}_{\text{MD-FFT excess variance}}. \quad (19)$$

It is clear from (19) that the excess variance increases with the size of the quantization step  $q$  and the FFT length ( $N$ ) and decreases with the number ( $K$ ) of spectrum estimates

being averaged. The reason why this excess variance increases with  $N$  is that the amount of quantization-related distortion (and hence quantization error) introduced in the estimation process increases with the number of samples being quantized. Also, since the amplitude variation of the 2 dither signals used is fixed at  $(\pm q)$  (so as to render them optimal in the sense of minimizing this excess variance), then an increase in  $q$  will increase the power of these dither signals and hence will also increase the amount of excess variance introduced. However, in practice, the choice of  $N$  and  $q$  is dictated by the desired frequency resolution and the amplitude range of the signal, respectively. This then leaves us with only 1 free experimental parameter ( $K$ ) to use as a way of controlling the amount of excess variance introduced.

Using the fact that, in the case of the MH-FFT estimator, only the dithered input signal is clipped at  $\pm q/2$ , the variance of the MR-FFT estimator is then readily obtained from (18):

$$\sigma_{\text{MH-FFT}}^2 = \sigma_{\text{SD-FFT}}^2 + \underbrace{\frac{1}{K} \left[ \sum_{n=0}^{N-1} \left\{ \frac{q^2}{4} E[|K(n, \omega_i)|^2] - E[(x^p(n) |K(n, \omega_i)|)^2] \right\} \right]}_{\text{MH-FFT excess variance}}. \quad (20)$$

Using the fact that the Fourier kernel is a deterministic quantity and carrying out the first summation on the RHS of (20) yields

$$\sigma_{\text{MH-FFT}}^2 = \sigma_{\text{SD-FFT}}^2 + \underbrace{\frac{q^2 N}{4K} - \frac{1}{K} \sum_{n=0}^{N-1} E[(x^p(n) |K(n, \omega_i)|)^2]}_{\text{MH-FFT excess variance}}. \quad (21)$$

Here too, the excess variance is affected by the 3 parameters  $q$ ,  $N$ , and  $K$ . As pointed out above, of all 3 parameters, only  $K$  is used in practice to control the amount of excess variance introduced by NSDQ quantization.

It is to be pointed out here that the dither signals used in all of our simulation are all called “optimal” in the sense that they minimize the excess variance introduced by the NSDQ quantization. This “optimality” result is not yet published and requires that these dither signals satisfy the following criteria: (a) each dither signal is uniformly distributed over the peak-to-peak range of the input it is added to and (b) the quantizer’s gain, in each channel, is set to twice the peak value of the input to this channel. Both of these criteria have been adhered to in our simulation work.

#### 4. SIMULATION

In order to test the new theoretical developments presented in this paper and to assess the performance of the 2 proposed 1 bit FFT estimators, namely, MR-FFT and MPC-FFT, we carried out a substantial simulation work on a variety of signals, both simulated and real ones. Here we will discuss a

representative set of these results which were partly reported earlier in [9–12] along with other new results obtained in both noise-free and noisy environments.

It is important to point out at this juncture that from an implementation (or simulation) point of view and with reference to Figures 1 and 2, it can be easily shown that the discrete averaging block “ $E[\cdot]$ ,” of gain  $K^{-1}$  (say), can be subsumed in the  $N$ -point “DFT” operation, by simply oversampling the NSDQ quantizer’s input at a rate equal to  $K$  and then processing all of the resulting  $(KN)$  samples. It is also worth pointing out here that each of the dither signals used in our simulation is zero-mean, uniformly distributed over the peak-to-peak amplitude range of the signal it is added to and statistically independent of both the input and all other (if any) dither signals used.

The four simulation examples which are used here as a testbed and which are made of 2 simulated signals and 2 real ones derived from the recordings of 2 sound signals are now briefly described. In each example, both the magnitude and phase spectra of the original (i.e., unquantized and undithered) signal are used as a reference against which the performance in estimation accuracy of the 2 proposed 1 bit nonsubtractively dither-quantized (NSDQ) estimators, that is, MR-FFT and MPC-FFT, is measured. The first example involves a single sinusoid and is used primarily to demonstrate, in detail, the excellent estimation accuracy of the 2 proposed 1 bit MR-FFT and MPC-FFT schemes when compared to their 1 bit undithered counterparts, referred to here simply as relay-FFT (R-FFT) and polarity coincidence-FFT (PC-FFT), respectively. The second example builds on the success of the dithering technique employed in the first



example, by testing the FFT spectrum estimation accuracy of the 2 proposed estimators on a more general signal, namely, a multisine signal. In the third example, the proposed MR-FFT and MPC-FFT estimators are used to estimate the FFT spectrum of a real musical signal. As a final test, the 2 proposed FFT estimators are tested on the record of a sound signal obtained from the utterance of the word “*Matlab*.” The simulation work carried out here is based on the diagrammatic descriptions of the 2 proposed estimators given in Figures 1 and 2.

A detailed description of each simulation example now follows.

A sinusoidal signal of amplitude  $A = 10$  and frequency  $f = 1000$  Hz is sampled at  $f_s = 8000$  Hz and used as the input signal  $x(n)$ . A total of 80 000 points are used for the estimation of the FFT magnitude spectrum. This simulation consists of 2 parts: the first part demonstrates the deleterious effects, on the FFT spectrum estimation, of undithered 1 bit quantization, be it applied to one or both of the estimator’s channels, as shown in Figures 3 and 4. These figures show, respectively, the amplitude and phase spectra of the original (i.e., unquantized) signal and those of the undithered 1 bit quantized estimators, that is, R-FFT and PC-FFT. Figure 3 shows that: (a) at the test frequency, both of the R-FFT and PC-FFT estimators suffer from a large relative estimation error of about 60% in the FFT magnitude spectrum at the test frequency and (b) there is a noticeable presence of non-negligible spurious signal peaks located at the third harmonic (and at other not-shown odd harmonics) of the test frequency in the magnitude spectra obtained with both the R-FFT and PC-FFT estimators, thus resulting in an unwanted and well-structured error pattern which only increases the total estimation error. Note here that the relative estimation error is defined here as the estimation error normalized by the peak magnitude spectrum value at the test frequency. As to Figure 4, it shows that, with both of the R-FFT and PC-FFT estimators and in addition to the correct phase value at the test frequency, there is another non-negligible spurious phase value at the third harmonic (and at other not-shown odd harmonics) of the test frequency.

Thus it is clear from the above that both the R-FFT and PC-FFT estimators greatly suffer from the adverse effect of 1 bit quantization on the FFT spectrum estimation, thus prohibiting them from exploiting all of the practical advantages that the simple and attractive 1 bit signal coding scheme brings to them.

The second part of this simulation sets out to demonstrate the excellent performance improvement brought to both the R-FFT and PC-FFT estimators by the nonsubtractive dithering technique which, when applied, modifies both of them to the 2 proposed MR-FFT and MPC-FFT estimators. To test this fact, the input signal, a sinewave of amplitude  $A = 10$  and frequency  $f = 1000$  Hz, is first sampled at  $f_s = 8000$  Hz, then added to a zero-mean random uniformly distributed dither signal which has the input’s peak-to-peak amplitude range, and finally their sum is 1 bit quantized. This combined process of nonsubtractively dithering a signal and then 1 bit quantizing the dithered signal (i.e., the

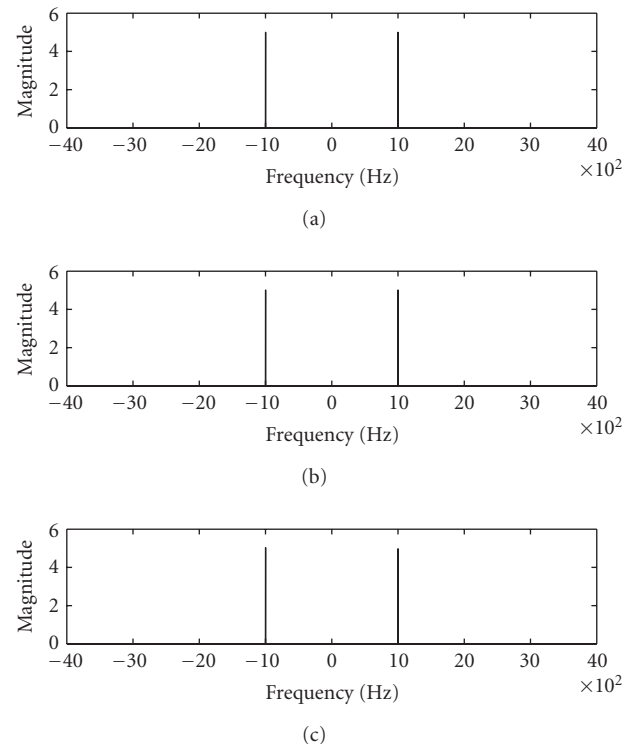


FIGURE 5: FFT magnitude spectra of a single sinusoid: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

sum signal) is what is referred to here as a 1 bit NSDQ quantization. For the MPC-FFT estimator, both the sine and cosine basis functions are also 1 bit NSDQ-quantized by using, as pointed out above, a second dither signal that is statistically independent of both the dither used for the input signal, and of the input itself. Next the FFT spectra of this 1 bit quantized signal are estimated using the proposed schemes and a total of 80 000 samples. The results, shown in Figures 5 and 6, clearly demonstrate the superior performance of the proposed MR-FFT and MPC-FFT estimators. These estimators have not only fully recovered the correct FFT magnitude and phase spectra, with a maximum relative magnitude error of at most 4%–5% for the worst-affected estimator (MPC-FFT), but have also virtually eliminated the structured harmonics-related error in the magnitude spectrum of Figure 3. It is important to note here that, in order for the MPC-FFT estimator’s performance to match that of the MR-FFT, the former estimator has to process more samples than the latter one. This fact is to be expected as the MPC-FFT estimator involves more quantization, and hence more signal distortion, since both of its channels are quantized, than does the MR-FFT one which has only one of its channels quantized. It is also worth pointing out here that, if needed, then increasing the number of samples will lead to an enhanced performance for both estimators because of the earlier-mentioned consistency of the sample mean estimators used.

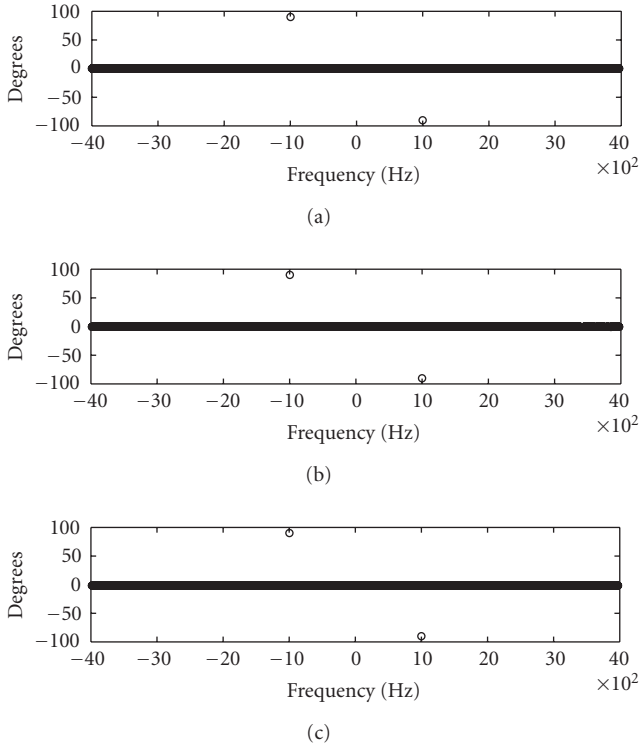


FIGURE 6: FFT phase spectra of a single sinusoid: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

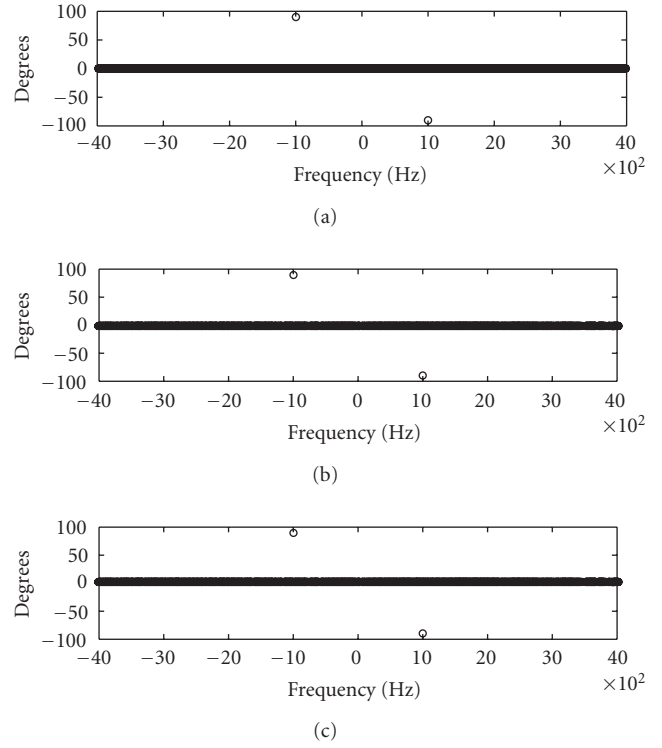


FIGURE 8: FFT phase spectra of a noisy single sinusoid: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom). SNR = 15 dB.

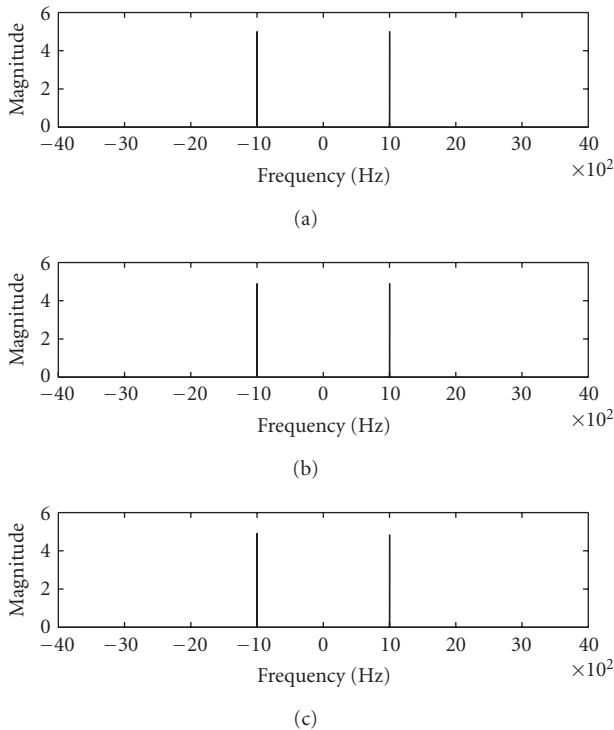


FIGURE 7: FFT magnitude spectra of a noisy single sinusoid: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom). SNR = 15 dB.

Although the performance of the 2 proposed estimators was also successfully tested in a noisy (Gaussian) environment with a SNR of 15 dB, only the results on the recovery of the FFT magnitude spectrum were reported [12]. We will now report on new results that corroborate the fact that this noise robustness is also enjoyed by the proposed estimators in the recovery of FFT phase spectra. However and as is expected with noisy environments, if the additional estimation error due to the effect of the added noise is to be reduced to a negligible level, more samples are to be processed than in noise-free environments. The test signal is a single sinewave, of amplitude  $A = 10$  and frequency  $f = 1000$  Hz, that is sampled at  $f_s = 8000$  Hz and then buried in a noisy environment characterized by a SNR of 15 dB. The total number of samples used here is 104 000 representing an excess of 24 000 samples as compared to the noise-free case discussed above. Both Figures 7 and 8 show an excellent performance by the proposed estimators in recovering both the magnitude and phase spectra at this moderate noise contamination level. Although not shown here, when the SNR is lowered to 5 dB representing a more severe noise contamination of the input signal and when the number of samples is kept unchanged, the performance of both estimators remains acceptable on the whole except for the MPC-FFT's performance in recovering the phase spectra which has been the worst affected. Nevertheless, this loss in performance can, if desired, be reduced through processing more samples.

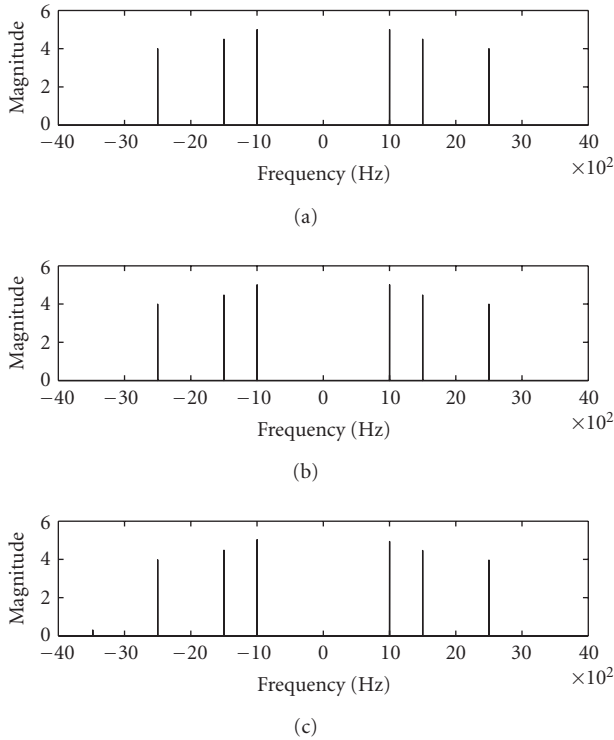


FIGURE 9: FFT magnitude spectra of a multisine signal: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

In our second example, a multisine input signal, comprised of 3 sine signals of amplitudes 10, 9, and 8 at frequencies 1000, 1500, and 2500 Hz, respectively, and sampled at  $f_s = 8000$  Hz, is considered. The results, displayed in Figures 9 and 10 and obtained with a total number of 80 000 samples, clearly show the excellent performance of the 2 proposed estimators in estimating the FFT magnitude and phase spectra, respectively. The worst-affected estimator (MPC-FFT) suffers from only a very negligible error of about 3% in estimating the magnitude spectrum and an additional spurious phase value. As explained above in the single sine example, the superior performance of the MRC-FFT estimator over the MPC-FFT one is also to be expected here.

In this case, the noise robustness of the 2 proposed estimators was also successfully tested in [12] where the magnitude spectra of a multisine signal were estimated by both proposed estimators with a very good accuracy and in a noisy environment characterized by a SNR = 15 dB. The noise robustness test in [12] showed that, as in the single sine case therein, the maximum relative estimation error of the MPC-FFT estimator at the 3 test frequencies had increased but not exceeded the value of 10. However, the noise floor peak value had increased to about 20%. Such increases in the maximum relative errors can be prevented or controlled by increasing the number of samples to be processed. This prediction is well supported by the consistency property of the estimators used and the extra computational cost involved in this case is

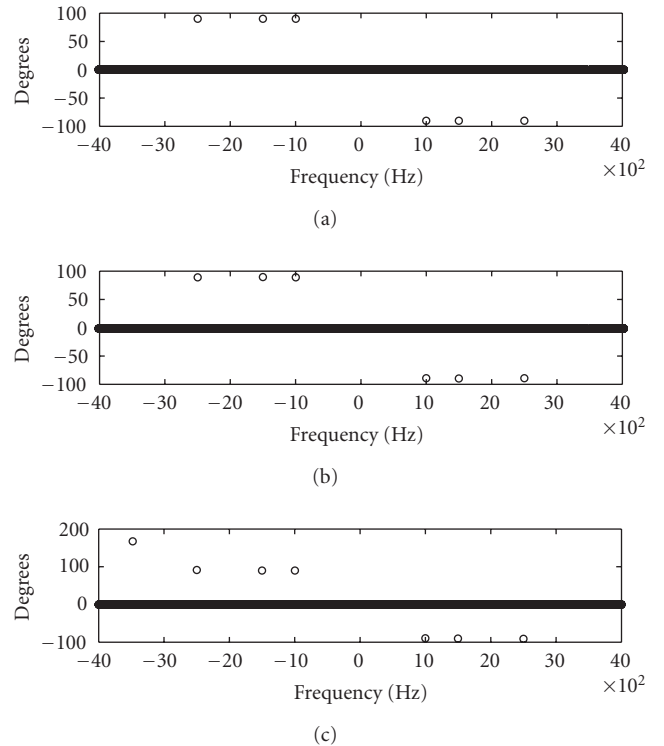


FIGURE 10: FFT phase spectra of a multisine signal: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

only small as all the samples to be processed are binary and hence can be processed very fast.

In our third example, a clarinet tune is selected as a practical example to test the performance of the two proposed estimators on a real signal. The tune was recorded at a quantization resolution of 16 bits per sample. This tune was sampled at a frequency of 16000 samples per second and its duration is 1.8 seconds. As such, a single record of this tune contains only  $N = 28800$  samples. This number of samples was found insufficient for an acceptable estimation accuracy. Since the estimation accuracy increases with the total number of samples processed (consistency property), the limited number of samples emanating from a single record of the clarinet tune data had to be replicated a number of times (i.e.,  $K = 10$ ) to increase the total number of samples to be processed in order to achieve a good estimation accuracy. This augmented data set, totalling 288 000 samples, was then used to test the 2 proposed 1 bit FFT estimators. Despite the short length of the available single record, the results of this performance test, reported in Figures 11 and 12, show that, as expected, a much better recovery of the original magnitude and phase spectra of the recorded tune was achieved by the MR-FFT estimator than by the MPC-FFT one. Moreover, due to the limited total number of samples used, the maximum level of the noisy pattern at the baseline of Figure 11 has reached about 20% of the peak magnitude value of the spectrum for the MPC-FFT estimator. With regard to the phase

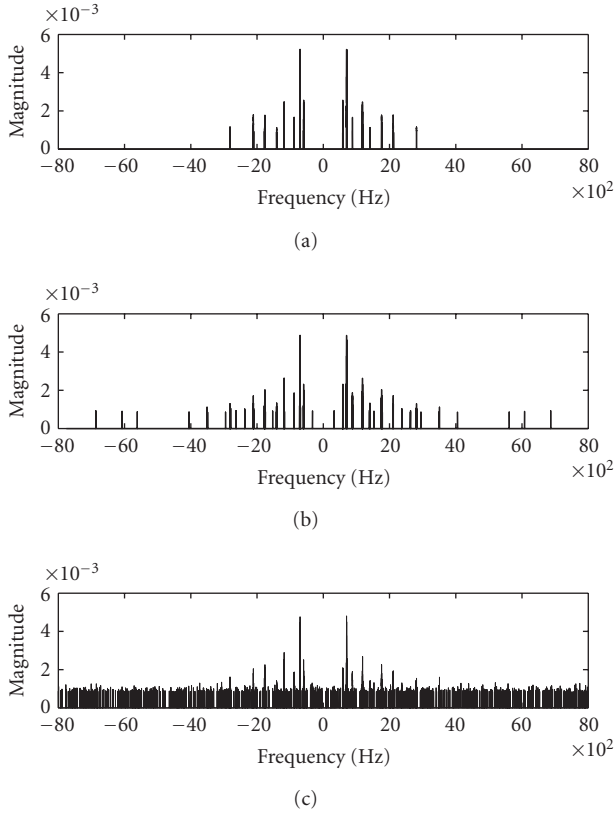


FIGURE 11: FFT magnitude spectra of a clarinet tune: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

recovery by the MPC-FFT, Figure 12 shows that there is a lot of spurious nonzero phase values, especially in the end regions where the phase should be zero. Although this may be unacceptable in certain applications, we regard this possibly unacceptable performance of the MPC-FFT estimator as being due solely to the limitation of the memory capacity of the PC used rather than to the theory underpinning the operation of the estimator itself.

In our last example, a sound recording of the utterance “*Matlab*,” of duration 0.5 second, is used as a test signal for the 2 proposed estimation schemes. The sound recording was saved at a resolution of 16 bits per sample using a sampling frequency of 8 KHz. As such, this single record accounts for  $N = 4000$  samples. As with the clarinet tune, the processing of a single record (regardless of the number of frames used) was not found sufficient for a good spectrum estimation accuracy. Hence, the available record had to be duplicated  $K = 25$  times yielding a total of 100 000 samples that was used in our simulation. Figures 13 and 14 show below the satisfactory performance of the 2 proposed 1 bit estimators in recovering both the magnitude and phase spectra. Here too, the simulation results demonstrate the superiority of the MR-FFT estimator over the MPC-FFT one. As in the previous experiment (clarinet tune), the noise floor in Figure 13 and the spurious nonzero phase values in Figure 14 that were generated by both the MR-FFT and MPC-FFT estimators can

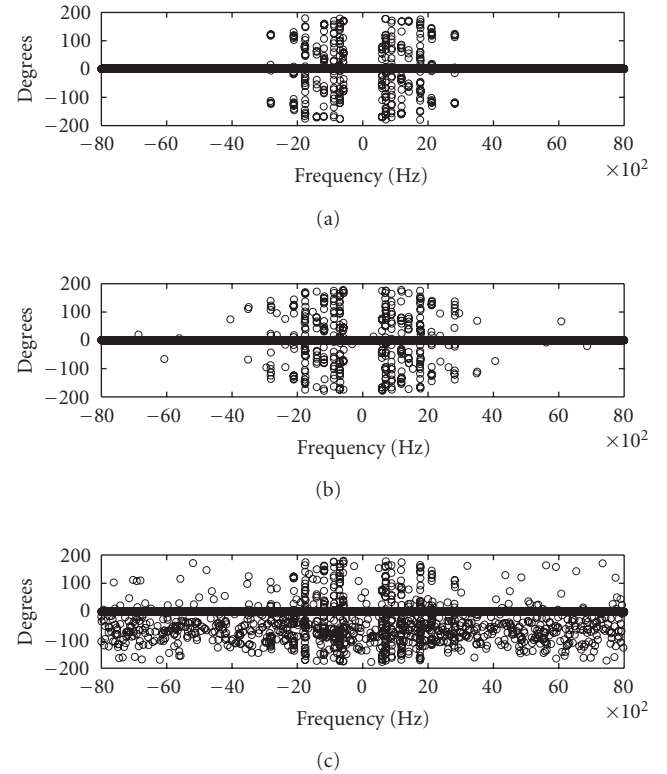


FIGURE 12: FFT phase spectra of a clarinet tune: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

be both further reduced to any acceptable level by processing further samples.

## 5. CONCLUSION

In this paper, we studied the problem of improving the computational accuracy and efficiency of the FFT through the use of 1 bit NSDQ quantization. We showed that the solution to this problem resides in extending the EMR theory to the frequency domain and in the process, derived new results which provided the theoretical underpinnings for a large class of computationally efficient FFT estimators. These dithered estimators were shown to be capable, in theory at least, of exactly recovering the true FFT spectrum (or its average if the input signal is stochastic), irrespective of the quantization resolution used. This flexibility in the choice of the quantization resolution to be used was thoroughly exploited in our simulation work by considering only the most practically attractive signal coding scheme based on 1 bit NSDQ quantization. This led to the 2 proposed 1 bit MR-FFT and MPC-FFT estimators. The estimation accuracy of these 2 estimators was thoroughly tested using a variety of simulated and real signals and in both noise-free and noisy environments (as reported elsewhere). The simulation results show that the maximum relative estimation error (incurred by the worst-affected MPC-FFT estimator), although not zero as



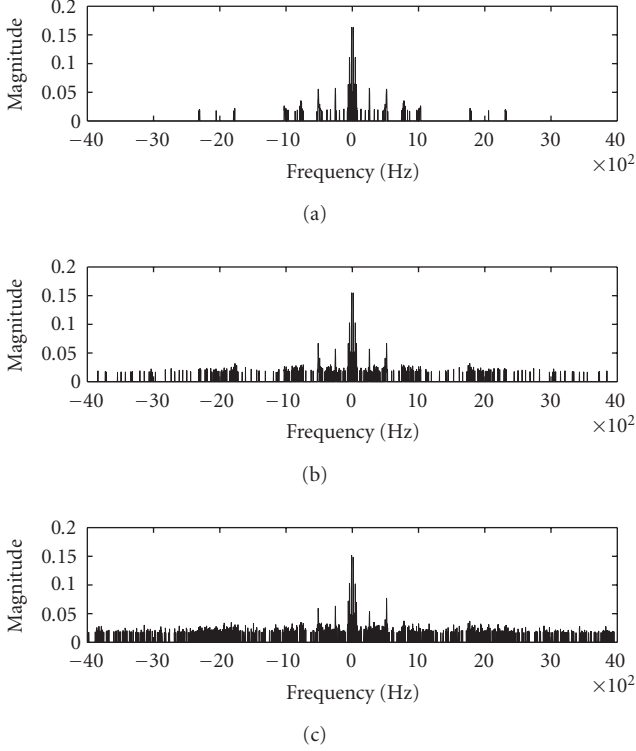


FIGURE 13: FFT magnitude spectra of the utterance “Matlab”: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

prescribed by the proposed theory since only a finite number of samples were used in estimating the different FFT spectra, remains nevertheless acceptable even in severely noisy environments and can always be reduced to any desired level through processing more samples. As such, these excellent results therefore strongly substantiate the proposed exact FFT recovery theory. The attractive practical advantages that accrue from the use of these 1 bit FFT estimators, such as simple architecture, low-cost implementation, very good accuracy, and fast and efficient computational capability, certainly provide ample encouragement not only to pursue their hardware implementation on a chip using either VLSI or FPGA technology but also to extend the 1 bit NSDQ quantization-based exact recovery theory advanced in this paper to other important transforms and to study the feasibility of parallelizing the proposed 1 bit low-cost estimation scheme for further possible computational gains. Finally, although the hardware implementation of the 2 proposed estimators is beyond the scope of this paper, it must be noted here that since the FFT used here is the one available in popular packages such as Matlab, the 1 bit nature of the input to the “FFT” block will be lost at the output of the first stage (or butterfly) of the FFT algorithm. In order to preserve this 1 bit nature throughout the FFT algorithm so as to have a purely 1 bit FFT, it is necessary to re-NSDQ-quantize the input to each subsequent stage of the FFT algorithm using statistically independent dither signals that are also members of  $D_1$ .

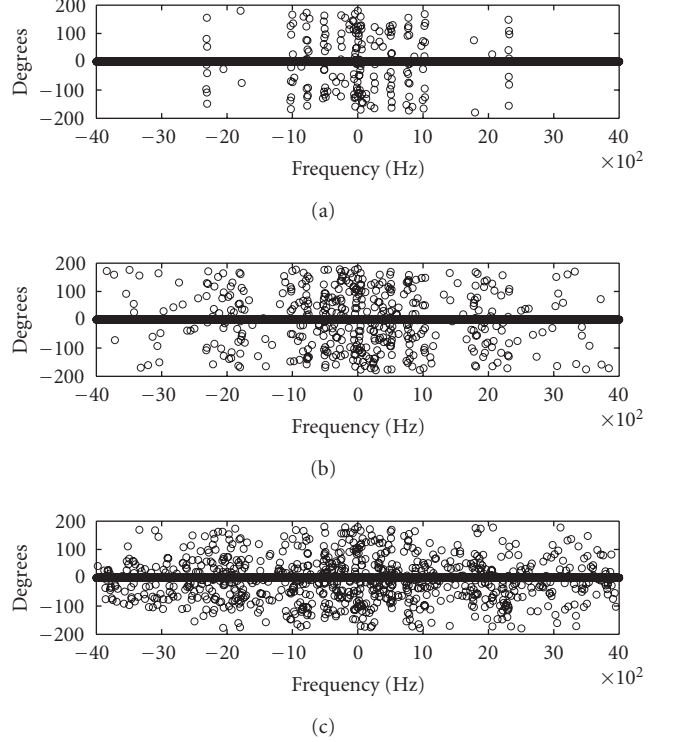


FIGURE 14: FFT phase spectra of the utterance “Matlab”: original (true) spectrum (top), estimated with MR-FFT estimator (middle) and with MPC-FFT estimator (bottom).

## APPENDICES

### A. PROOF OF THEOREM 1

Recall that, by definition, the 1 channel-quantized  $p$ th-order DFT of the NSDQ quantizer’s output is given by

$$X_{\text{NSDQ}_p}^{[1]}(\omega_i) \triangleq \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) \cdot K(n, \omega_i). \quad (\text{A.1})$$

Taking the expectation of both sides of (A.1) and using the fact that the expectation operation is a linear operation, that is, the expectation of a sum of random variables (RVs) is the sum of the expectations of these RVs, the following is obtained:

$$\begin{aligned} E[X_{\text{NSDQ}_p}^{[1]}(\omega_i)] &\triangleq E\left\{ \sum_{n=0}^{N-1} [x_{\text{NSDQ}}^p(n)K(n, \omega_i)] \right\} \\ &= \sum_{n=0}^{N-1} E[x_{\text{NSDQ}}^p(n)K(n, \omega_i)]. \end{aligned} \quad (\text{A.2})$$

Now, from (4) of Lemma 1, we know that  $x_{\text{NSDQ}}^p(n)$  is moment-sense equivalent to the  $p$ th-order polynomial mapping  $h_p(x(n))$  defined in (5). Hence, by virtue of this, we can

write that

$$E\left[x_{\text{NSDQ}}^p(n)K(n, \omega_i)\right] = E[h_p(x(n))K(n, \omega_i)]. \quad (\text{A.3})$$

Using (A.3) with (A.2) and taking the expectation operator  $E[\cdot]$  out of the summation gives

$$E\left[X_{\text{NSDQ}_p}^{[1]}(\omega_i)\right] = E\left[\underbrace{\sum_{n=0}^{N-1} h_p(x(n))K(n, \omega_i)}_{=\text{DFT}[h_p(x(n))]\triangleq H_p(\omega_i)}\right] = E[H_p(\omega_i)]. \quad (\text{A.4})$$

Now using the expression of  $h_p(x(n))$  given in (5), it is easy to show that

$$\begin{aligned} H_p(\omega_i) &= \sum_{n=0}^{N-1} \left\{ \sum_{k=0}^p c_k x^k \right\} K(n, \omega_i) \\ &= \sum_{k=0}^p c_k \underbrace{\sum_{n=0}^{N-1} x^k K(n, \omega_i)}_{\triangleq X_k(\omega_i)} = \sum_{k=0}^p c_k X_k(\omega_i). \end{aligned} \quad (\text{A.5})$$

Combining (A.4) and (A.5) leads to the final desired result, that is,

$$E\left[X_{\text{NSDQ}_p}^{[1]}(\omega_i)\right] = E[H_p(\omega_i)] = \sum_{k=0}^p c_k X_k(\omega_i). \quad (\text{A.6})$$

## B. PROOF OF THEOREM 2

Using the definition of the 2-quantized channel  $p$ th-order DFT and the Cartesian expression of the NSDQ-quantized DFT kernel, the following is obtained:

$$\begin{aligned} X_{\text{NSDQ}_p}^{[2]}(\omega_i) &\triangleq \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) \cdot K_{\text{NSDQ}}(n, \omega_i) \\ &= \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) \cdot c_{\text{NSDQ}}(n, \omega_i) \\ &\quad - j \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) \cdot s_{\text{NSDQ}}(n, \omega_i). \end{aligned} \quad (\text{B.1})$$

Taking the expectation of both sides of (B.1) and using the linearity property of the expectation operator  $E[\cdot]$  yield

$$\begin{aligned} E\left[X_{\text{NSDQ}_p}^{[2]}(\omega_i)\right] &\triangleq E\left\{ \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) c_{\text{NSDQ}}(n, \omega_i) \right. \\ &\quad \left. - j \sum_{n=0}^{N-1} x_{\text{NSDQ}}^p(n) s_{\text{NSDQ}}(n, \omega_i) \right\} \\ &= \sum_{n=0}^{N-1} E\left[x_{\text{NSDQ}}^p(n) c_{\text{NSDQ}}(n, \omega_i)\right] \\ &\quad - j \sum_{n=0}^{N-1} E\left[x_{\text{NSDQ}}^p(n) s_{\text{NSDQ}}(n, \omega_i)\right]. \end{aligned} \quad (\text{B.2})$$

Now, using (12) in Lemma 2, we know that both products  $x_{\text{NSDQ}}^p(n)c_{\text{NSDQ}}(n, \omega_i)$  and  $x_{\text{NSDQ}}^p(n)s_{\text{NSDQ}}(n, \omega_i)$  are moment-sense equivalent to the 2-D  $(p, 1)$ th-order polynomial mappings  $h_{p,1}[x(n), c(n, \omega_i)]$  and  $h_{p,1}[x(n), s(n, \omega_i)]$  which characterize the input-cosine and input-sine 2-D NSDQ quantizers, respectively (see Figure 2). These moment-sense equivalences are then given by

$$\begin{aligned} E\left[x_{\text{NSDQ}}^p(n)c_{\text{NSDQ}}(n, \omega_i)\right] &= E[h_{p,1}[x(n), c(n, \omega_i)]] \\ E\left[x_{\text{NSDQ}}^p(n)s_{\text{NSDQ}}(n, \omega_i)\right] &= E[h_{p,1}[x(n), s(n, \omega_i)]]. \end{aligned} \quad (\text{B.3})$$

Moreover, assuming that the dither signal used for the input  $x(n)$  and the common dither signal used for the 2 basis functions  $c(n)$  and  $s(n)$  are statistically independent of each other and of  $x(n)$ , this then warrants the use of the separability property of the two 2-D MSIOFs which, together with the result of (7), that is,  $h_1(x) = x$ , finally lead to the following new form of (B.3):

$$\begin{aligned} E\left[x_{\text{NSDQ}}^p(n)c_{\text{NSDQ}}(n, \omega_i)\right] &= E[h_p[x(n)]h_1[c(n, \omega_i)]] \\ &= E[h_p[x(n)]c(n, \omega_i)], \\ E\left[x_{\text{NSDQ}}^p(n)s_{\text{NSDQ}}(n, \omega_i)\right] &= E[h_p[x(n)]h_1[s(n, \omega_i)]] \\ &= E[h_p[x(n)]s(n, \omega_i)]. \end{aligned} \quad (\text{B.4})$$

Combining (B.4) with (B.2), and using the linearity property of the expectation operator  $E[\cdot]$  by taking this operator out of the summation, the following is obtained:

$$\begin{aligned} E\left[X_{\text{NSDQ}_p}^{[2]}(\omega_i)\right] &= \left\{ E\left[\sum_{n=0}^{N-1} h_p[x(n)]\{c(n, \omega_i) - js(n, \omega_i)\}\right], \right. \\ &= \left. E\left[\underbrace{\sum_{n=0}^{N-1} h_p[x(n)]\{K(n, \omega_i)\}}_{=\text{DFT}[h_p(x(n))]} \right] = E[H_p(\omega_i)]. \right. \end{aligned} \quad (\text{B.5})$$

Hence, combining (A.5) with (B.5) gives the final desired result

$$E\left[X_{\text{NSDQ}_p}^{[2]}(\omega_i)\right] = E[H_p(\omega_i)] = \sum_{k=0}^p c_k X_k(\omega_i). \quad (\text{B.6})$$

## C. GENERALIZED VARIANCE ANALYSIS

We will first define the 3 FFT estimators to be considered here, namely, the conventional sampled-data (i.e., totally

unquantized) FFT estimator (SD-FFT), the modified hybrid FFT estimator (MH-FFT) which incorporates a single multibit NSDQ-quantized channel (in our case, the input signal channel is the quantized one), and finally the modified digital (i.e., fully quantized) FFT (MD-FFT) estimator which comprises 2 multibit NSDQ-quantized channels. It is clear that the MH-FFT and MD-FFT are nothing but generalizations of the 2 proposed estimators, that is, MR-FFT and MPC-FFT, respectively:

$$\begin{aligned} \text{SD-FFT : } \hat{X}_p(\omega_i) &= \frac{1}{K} \sum_{k=0}^{K-1} \{X_{p_k}(\omega_i)\} \\ &= \frac{1}{K} \sum_{k=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} x^p(n)K(n, \omega_i) \right\}; \end{aligned} \quad (\text{C.1})$$

$$\begin{aligned} \text{MH-FFT : } \hat{X}_{MH_p}(\omega_i) &= \frac{1}{K} \sum_{k=0}^{K-1} \{X_{MH_{p_k}}(\omega_i)\} \\ &= \frac{1}{K} \sum_{k=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} x_{\text{NSDQ}}^p(n)K(n, \omega_i) \right\}; \end{aligned} \quad (\text{C.2})$$

$$\begin{aligned} \text{MD-FFT : } \hat{X}_{MD_p}(\omega_i) &= \frac{1}{K} \sum_{k=0}^{K-1} \{X_{MD_{p_k}}(\omega_i)\} \\ &= \frac{1}{K} \sum_{k=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} x_{\text{NSDQ}}^p(n)K_{\text{NSDQ}}(n, \omega_i) \right\}. \end{aligned} \quad (\text{C.3})$$

We will first define the variance of the SD-FFT estimator which will then be used as a reference against which the variances of both the MH-FFT and the MD-FFT estimators will be compared.

As pointed out in Section 3.6.1, all sample mean estimators used here, including the one for the SD-FFT estimator, are both unbiased and consistent estimators of the true spectrum  $X_p(\omega_i)$ . It is also worth pointing out here that since both the true and average spectra are both complex- and scalar-valued quantities, the usual ‘‘Hermitian’’ operation (denoted by the superscript  $H$ ) involved in the definition of the variance of a complex- and vector-valued random signal reduces here to the complex conjugation operation (denoted by the superscript  $*$ ) only. It then follows that, in this case, the definition of the variance of the SD-FFT estimator, denoted here by  $\sigma_{\text{SD-FFT}}^2$ , is expressed by

$$\begin{aligned} \sigma_{\text{SD-FFT}}^2 &\triangleq E[\{\hat{X}_p(\omega_i) - X_p(\omega_i)\}\{\hat{X}_p(\omega_i) - X_p(\omega_i)\}^*] = E[\{\hat{X}_p(\omega_i)\}\{\hat{X}_p(\omega_i)\}^*] - |X_p(\omega_i)|^2 \\ &= \frac{1}{K^2} E \left[ \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} X_{p_k}(\omega_i) X_{p_l}^*(\omega_i) \right] - |X_p(\omega_i)|^2 \\ &= \frac{1}{K^2} E \left[ \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} x^p(n)K(n, \omega_i) \right\} \left\{ \sum_{m=lN}^{(l+1)N-1} x^p(m)K(m, \omega_i) \right\}^* \right] - |X_p(\omega_i)|^2 \\ &= \frac{1}{K^2} E \left[ \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} \sum_{m=lN}^{(l+1)N-1} x^p(n)K(n, \omega_i) x^p(m)K^*(m, \omega_i) \right\} \right] - |X_p(\omega_i)|^2. \end{aligned} \quad (\text{C.4})$$

Note that in (C.4), use was made of the identity  $X_p(\omega_i)X_p^*(\omega_i) = |X_p(\omega_i)|^2$  and the complex conjugation operation ( $*$ ) was applied only to the complex Fourier kernel  $K(m, \omega_i)$  and not to the signal  $x^p(m)$  as this one is real. It is clear from (C.4) that the variance of the SD-FFT estimator can be minimized to any desired level by simply increasing the number ( $K$ ) of estimated spectra ( $X_{p_k}(\omega_i)$ ) being averaged. Note here that increasing  $K$  clearly implies processing more samples of both the signal being analyzed and the Fourier kernel.

Since the MH-FFT estimator is a special case of the MD-FFT one, we will therefore first derive the expression of the variance of the latter (MD-FFT) and then infer from it the expression of the variance of the former (MH-FFT).

By definition, the variance of the MD-FFT estimator, denoted here by  $\sigma_{\text{MD-FFT}}^2$ , is given by

$$\begin{aligned} \sigma_{\text{MD-FFT}}^2 &\triangleq E[\{\hat{X}_{MD_p}(\omega_i) - X_p(\omega_i)\}\{\hat{X}_{MD_p}(\omega_i) - X_p(\omega_i)\}^*] \\ &= E[\{\hat{X}_{MD_p}(\omega_i)\}\{\hat{X}_{MD_p}(\omega_i)\}^*] - |X_p(\omega_i)|^2 \\ &= \frac{1}{K^2} E \left[ \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} X_{MD_{p_k}}(\omega_i) X_{MD_{p_l}}^*(\omega_i) \right] - |X_p(\omega_i)|^2. \end{aligned} \quad (\text{C.5})$$

Now combining (C.5) and (C.3) yields

$$\begin{aligned}\sigma_{\text{MD-FFT}}^2 &= \frac{1}{K^2} E \left[ \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} x_{\text{NSDQ}}^p(n) K_{\text{NSDQ}}(n, \omega_i) \right\} \left\{ \sum_{m=lN}^{(l+1)N-1} x_{\text{NSDQ}}^p(m) K_{\text{NSDQ}}(m, \omega_i) \right\}^* \right] - |X_p(\omega_i)|^2 \\ &= \frac{1}{K^2} \left[ \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} \sum_{m=lN}^{(l+1)N-1} E[x_{\text{NSDQ}}^p(n) K_{\text{NSDQ}}(n, \omega_i) x_{\text{NSDQ}}^p(m) K_{\text{NSDQ}}^*(m, \omega_i)] \right\} \right] - |X_p(\omega_i)|^2.\end{aligned}\quad (\text{C.6})$$

If we now split the elements of the double summation over  $k$  and  $l$  in (C.6) into 2 sets, one for  $k \neq l$  and the other for  $k = l$ , and noticing that the conditions  $k \neq l$  and  $k = l$  also imply that  $n \neq m$  and  $n = m$ , respectively, then the following is obtained:

$$\begin{aligned}\sigma_{\text{MD-FFT}}^2 &= \frac{1}{K^2} \left[ \sum_{k \neq l} \sum_{l \neq k} \left\{ \sum_{\substack{n \neq m \\ kN}}^{(k+1)N-1} \sum_{\substack{m \neq n \\ lN}}^{(l+1)N-1} E[x_{\text{NSDQ}}^p(n) K_{\text{NSDQ}}(n, \omega_i) \right. \right. \\ &\quad \left. \left. \times x_{\text{NSDQ}}^p(m) K_{\text{NSDQ}}^*(m, \omega_i)] \right\} \right] - |X_p(\omega_i)|^2 \\ &\quad + \frac{1}{K^2} \left[ \sum_{k=0}^{K-1} \sum_{n=kN}^{(k+1)N-1} E[(x_{\text{NSDQ}}^p(n) | K_{\text{NSDQ}}(n, \omega_i)|)^2] \right].\end{aligned}\quad (\text{C.7})$$

Note that in (C.7), use was made of the following identity:  $K_{\text{NSDQ}}(n, \omega_i) K_{\text{NSDQ}}^*(n, \omega_i) = |K_{\text{NSDQ}}(n, \omega_i)|^2$ .

Note also that in (C.7), the expectation  $E[x_{\text{NSDQ}}^p(n) K_{\text{NSDQ}}(n, \omega_i) x_{\text{NSDQ}}^p(m) K_{\text{NSDQ}}^*(m, \omega_i)]$  involves 4 different signals, namely,  $x^p(n)$ ,  $K(n, \omega_i)$ ,  $x^p(m)$ , and  $K^*(m, \omega_i)$ , that are NSDQ-quantized by 4 different and statistically independent dither signals, each of which is a member of the 1-D first-order class of linearizing dither signals  $\mathcal{D}_1$ . We now need to extend (12) of Lemma 2 to the case of the above-mentioned 4 different signals and invoke the facts that in this case, (a) sta-

tistical independence between the 4 dither signals used ensures separability of the quantizer's 4-D first-order MSIOF,  $h_{1,1,1,1}(\cdot, \cdot, \cdot, \cdot)$ , into 4 different 1-D first-order MSIOF,  $h_1(\cdot)$ , and (b) membership of each of the 4 statistically independent dither signals to  $\mathcal{D}_1$  will, according to (7), perfectly linearize the first-order MSIOF,  $h_1(\cdot)$ , of the associated NSDQ  $\neg$  quantizer. Making use of these 2 facts will then yield the following:

$$\begin{aligned}h_{1,1,1,1}(x^p(n), K(n, \omega_i), x^p(m), K^*(m, \omega_i)) \\ = \{h_1(x^p(n))\} \{h_1(K(n, \omega_i))\} \\ \times \{h_1(x^p(m))\} \{h_1(K^*(m, \omega_i))\} \\ = \{x^p(n)\} \{K(n, \omega_i)\} \{x^p(m)\} \{K^*(m, \omega_i)\}.\end{aligned}\quad (\text{C.8})$$

A direct 4-D extension of (12) in Lemma 2 (see [6] for further details), combined with (C.8), leads to

$$\begin{aligned}E[x_{\text{NSDQ}}^p(n) K_{\text{NSDQ}}(n, \omega_i) x_{\text{NSDQ}}^p(m) K_{\text{NSDQ}}^*(m, \omega_i)] \\ = E[h_{1,1,1,1}(x^p(n), K(n, \omega_i), x^p(m), K^*(m, \omega_i))] \\ = E[\{h_1(x^p(n))\} \{h_1(K(n, \omega_i))\} \\ \times \{h_1(x^p(m))\} \{h_1(K^*(m, \omega_i))\}] \\ = E[x^p(n) K(n, \omega_i) x^p(m) K^*(m, \omega_i)].\end{aligned}\quad (\text{C.9})$$

Combining (C.7) and (C.9) gives

$$\begin{aligned}\sigma_{\text{MD-FFT}}^2 &= \frac{1}{K^2} \left[ \sum_{k \neq l} \sum_{l \neq k} \left\{ \sum_{\substack{n \neq m \\ kN}}^{(k+1)N-1} \sum_{\substack{m \neq n \\ lN}}^{(l+1)N-1} E[x^p(n) K(n, \omega_i) x^p(m) K^*(m, \omega_i)] \right\} \right] \\ &\quad - |X_p(\omega_i)|^2 + \frac{1}{K^2} \left[ \sum_{k=0}^{K-1} \sum_{n=kN}^{(k+1)N-1} E[(x_{\text{NSDQ}}^p(n) | K_{\text{NSDQ}}(n, \omega_i)|)^2] \right].\end{aligned}\quad (\text{C.10})$$

Let us now add and subtract the term:  $Q = (1/K^2) [\sum_{k=0}^{K-1} \sum_{n=kN}^{(k+1)N-1} E[(x^p(n) | K(n, \omega_i)|)^2]]$  to the RHS of (C.10) in the following order: first, add  $Q$  to the first term

on the RHS of (C.10) and then group these 2 terms together and next subtract  $Q$  from the last term on the RHS of (C.10). As a result of these small manipulations, (C.10) becomes



$$\begin{aligned} \sigma_{\text{MD-FFT}}^2 = & \frac{1}{K^2} E \left[ \sum_{k=0}^{K-1} \sum_{l=0}^{K-1} \left\{ \sum_{n=kN}^{(k+1)N-1} x^p(n) K(n, \omega_i) \right\} \left\{ \sum_{m=lN}^{(l+1)N-1} x^p(m) K(m, \omega_i) \right\}^* \right] - |X_p(\omega_i)|^2 \\ & + \frac{1}{K^2} \left[ \sum_{k=0}^{K-1} \sum_{n=kN}^{(k+1)N-1} E \left[ (x_{\text{NSDQ}}^p(n) |K_{\text{NSDQ}}(n, \omega_i)|)^2 - (x^p(n) |K(n, \omega_i)|)^2 \right] \right]. \end{aligned} \quad (\text{C.11})$$

Since, according to (C.4), the first 2 terms on the RHS of (C.11) are nothing but  $\sigma_{\text{SD-FFT}}^2$ , (C.11) then takes on the

following form:

$$\sigma_{\text{MD-FFT}}^2 = \underbrace{\sigma_{\text{SD-FFT}}^2 + \frac{1}{K^2} \left[ \sum_{k=0}^{K-1} \sum_{n=kN}^{(k+1)N-1} E \left[ (x_{\text{NSDQ}}^p(n) |K_{\text{NSDQ}}(n, \omega_i)|)^2 - (x^p(n) |K(n, \omega_i)|)^2 \right] \right]}_{\text{MD-FFT excess variance}}. \quad (\text{C.12})$$

It can be readily shown that the expectation term in the second term on the RHS of (C.12) can be expressed as the autocorrelation function  $E[e_{\text{NSDQ}}^2(n, \omega_i)]$  of the following NSDQ quantization error:  $e_{\text{NSDQ}}(n, \omega_i) \triangleq (x^p(n) |K(n, \omega_i)|)_{\text{NSDQ}} - (x^p(n) |K(n, \omega_i)|)$ . According to the theory of NSDQ quantization and provided that the dither used to NSDQ-quantize the product signal,  $(x^p(n) |K(n, \omega_i)|)$ , is a member of  $\mathcal{D}_1$ , then this error is both zero-mean and uncorrelated with the

product signal  $(x^p(n) |K(n, \omega_i)|)$  itself. Since an autocorrelation function is known to be positive definite, it therefore follows that the entire second term on the RHS of (C.12) is positive and thus represents an excess in variance brought about by the NSDQ quantization.

Moreover, since the signal and the Fourier kernel are both assumed to be ergodic and stationary, the excess-variance term in (C.12) can be further simplified, leading to

$$\sigma_{\text{MD-FFT}}^2 = \underbrace{\sigma_{\text{SD-FFT}}^2 + \frac{1}{K} \left[ \sum_{n=0}^{N-1} E \left[ (x_{\text{NSDQ}}^p(n) |K_{\text{NSDQ}}(n, \omega_i)|)^2 - (x^p(n) |K(n, \omega_i)|)^2 \right] \right]}_{\text{MD-FFT excess variance}}. \quad (\text{C.13})$$

It is clear from (C.13) that no matter how large the NSDQ-induced excess variance is, it can always be reduced to any desired level by choosing a sufficiently large value of  $K$ , or equivalently by processing a sufficiently large number of samples.

The expression of the variance of the MH-FFT estimator can now be readily inferred from (C.13) by simply replacing the quantized Fourier kernel,  $K_{\text{NSDQ}}(n, \omega_i)$ , with its unquantized counterpart,  $K(n, \omega_i)$ , in the “excess-variance” term on the RHS of (C.13). This yields

$$\sigma_{\text{MH-FFT}}^2 = \underbrace{\sigma_{\text{SD-FFT}}^2 + \frac{1}{K} \left[ \sum_{n=0}^{N-1} E \left[ (x_{\text{NSDQ}}^p(n) |K(n, \omega_i)|)^2 - (x^p(n) |K(n, \omega_i)|)^2 \right] \right]}_{\text{MH-FFT excess variance}}. \quad (\text{C.14})$$

Here too, the “excess-variance” term can be reduced to any desired level in the manner described above for the MD-FFT estimator.

## D. PROOF OF THE CLOSURE PROPERTY (SECTION 2.2)

Let the input signal be  $x$ , the dither signal  $D$ , and their sum signal  $y = x + D$  and let their respective characteristic functions be  $W_x(u)$ ,  $W_D(u)$ , and  $W_y(u)$ . Assume further that  $x$  and  $D$  are statistically independent of each other, which leads to

$$W_y(u) = W_x(u)W_D(u), \quad (\text{D.1})$$

where  $q$  is the already defined uniform step of the NSDQ quantizer.

We first start with the proof of the first-order ( $p = 1$ ) version of this property which is rather straightforward. It is clear that if we express (D.1) at  $u = 2m\pi/q$ ,  $m \neq 0$ , we get  $W_y(2m\pi/q) = W_x(2m\pi/q)W_D(2m\pi/q)$ ,  $m \neq 0$ . If we now use the definition of the first-order class  $D_1$ , that is,  $D \in D_1 \Leftrightarrow W_D(2m\pi/q) = 0$ ,  $m \neq 0$ , then it follows that  $W_y(2m\pi/q) = W_x(2m\pi/q)W_D(2m\pi/q) = 0$ , for  $m \neq 0$ , and for all  $x$  which leads to  $y \in D_1$ .

We will now prove the general  $p$ th case by first differentiating (D.1),  $r$  times, with  $r \in [0, p - 1]$  and  $p \geq 1$ , at the point  $u = 2m\pi/q$ ,  $m \neq 0$ . The following is then obtained:

$$W_y^{(r)}\left(\frac{2m\pi}{q}\right) = \sum_{k=0}^r C_k^r W_x^{(r-k)}\left(\frac{2m\pi}{q}\right) W_D^{(k)}\left(\frac{2m\pi}{q}\right), \quad (\text{D.2})$$

where  $C_k^r \triangleq \binom{r}{k} = (r!)/[(k!)(r - k)!]$ .

Recall that by definition,  $D$  is a member of  $\mathcal{D}_p$  if and only if its characteristic function satisfies (3) which is restated here for convenience:

$$D \in D_p \Leftrightarrow W_D^{(r)}\left(\frac{2m\pi}{q}\right) = 0 \quad \forall r \in [0, p - 1], m \neq 0. \quad (\text{D.3})$$

It is clear from (D.2) that since we have  $0 \leq k \leq r \leq p - 1$ , it then follows from (D.3) that if  $D \in D_p$ , then the relationship  $W_D^{(k)}(2m\pi/q) = 0$ ,  $m \neq 0$ , is true not only within the range for all  $k \in [0, p - 1]$  but also within any of the  $(p - 1)$  subranges defined by: for all  $k \in [0, (p - \lambda) - 1]$  with  $\lambda \in [1, p]$ , that is,

$$W_D^{(k)}\left(\frac{2m\pi}{q}\right) = 0 \quad \forall k \in [0, (p - \lambda) - 1] \quad (\text{D.4})$$

with  $\lambda \in [1, p]$ ,  $m \neq 0$ .

Since the subrange “for all  $k \in [0, r]$  with for all  $r \in [0, p - 1]$ ” is one subrange mentioned above in (D.4), it then follows that

$$W_D^{(k)}\left(\frac{2m\pi}{q}\right) = 0, \quad m \neq 0, \quad \forall k \in [0, r]. \quad (\text{D.5})$$

Using (D.5) in (D.2) yields the following desired result:

$$W_y^{(r)}\left(\frac{2m\pi}{q}\right) = 0 \quad \forall r \in [0, p - 1], m \neq 0. \quad (\text{D.6})$$

By virtue of the definition of the  $p$ th-order class  $D_p$ , (D.6) leads to  $y = (x + D) \in D_p$ .

It is worth pointing out here that (D.4) reveals an interesting property of the class  $D_p$  which states that if  $D$  belongs to a  $p$ th-order class  $D_p$ , it will then belong to any  $(p - 1)$  subclasses of lower orders, that is, if  $D \in D_p \Rightarrow D \in D_s$  for all  $s \in [1, p - 1]$ . This is the property of inclusion enjoyed by the class  $D_p$  of linearizing dither signals (see [6]).

## E. PROOF OF THE LINEARIZATION OF THE FIRST-ORDER MSIOF $h_1(x)$

Setting  $p = 1$  in (5) gives the direct expression of the first-order MSIOF of the uniform NSDQ quantizer, that is,

$$h_1(x) = c_0 + c_1 x, \quad (\text{E.1})$$

where the coefficients  $c_0$  and  $c_1$  are derived as follows.

The expression of  $c_0$  is obtained directly from the general expression  $c_k$  in (5) by setting  $k = 0$  and  $p = 1$  therein. After a direct substitution followed by a simplification, we get

$$\begin{aligned} c_0 &= \sum_{t=0}^{1-0} \frac{1!}{(1-0-t+1)0!t!} \left(\frac{q}{2}\right)^{1-0-t} E[D^t][1 \oplus 0 \oplus t \oplus 1] \\ &= \sum_{t=0}^1 \frac{1}{(2-t)!t!} \left(\frac{q}{2}\right)^{1-t} E[D^t][0 \oplus t]. \end{aligned} \quad (\text{E.2})$$

Since  $t$  is limited here to the values of 0 or 1 and since the factor  $[0 \oplus t]$  is equal to 1 for odd values of  $t$  only, it then follows that this factor is equal to 1 for  $t = 1$  and is zero, otherwise. Thus, we have

$$c_0 = \frac{1}{(2-1)!1!} \left(\frac{q}{2}\right)^{1-1} E[D^1][0 \oplus 1] = E[D] = 0 \quad (\text{E.3})$$

since the dither signal is assumed to be zero-mean.

Repeating the same process for  $c_1$  by setting  $k = 1$  and  $p = 1$  in the general expression of  $c_k$  in (5) yields the following:

$$\begin{aligned} c_1 &= \sum_{t=0}^{1-1} \frac{1!}{(1-1-t+1)1!t!} \left(\frac{q}{2}\right)^{1-1-t} E[D^t][1 \oplus 1 \oplus t \oplus 1] \\ &= \sum_{t=0}^1 \frac{1}{(1-t)!t!} \left(\frac{q}{2}\right)^{-t} E[D^t][1 \oplus t]. \end{aligned} \quad (\text{E.4})$$

Since we only have one value of  $t$  here, that is,  $t = 0$ , and since in this case  $[1 \oplus t] = 1$ , we then have

$$c_1 = \frac{1}{(1-0)!0!} \left(\frac{q}{2}\right)^{-0} E[D^0][1 \oplus 0] = 1. \quad (\text{E.5})$$

Combining (E.1), (E.3), and (E.5) leads directly to the desired result  $h_1(x) = x$ .

## ACKNOWLEDGMENT

The authors would like to acknowledge KFUPM for its support in carrying out this research work.

## REFERENCES

- [1] J. W. Cooley and J. W. Tukey, "An algorithm for the machine computation of complex Fourier series," *Mathematics of Computation*, vol. 19, no. 90, pp. 297–301, 1965.
- [2] C. S. Burrus and T. W. Parks, *DFT/FFT and Convolution Algorithms*, John Wiley & Sons, New York, NY, USA, 1985, see also C. S. Burrus: "Notes on the FFT", <http://www.dsp.rice.edu>.
- [3] A. Ganapathiraju, J. Hamaker, J. Picone, and A. Skjellum, "Contemporary view of FFT algorithms," in *Proceedings of the IASTED International Conference on Signal and Image Processing (SIP '98)*, pp. 130–133, Las Vegas, Nev, USA, October 1998.
- [4] P. Duhamel and M. Vetterli, "Fast Fourier transforms: a tutorial review and a state of the art," *Signal Processing*, vol. 19, no. 4, pp. 259–299, 1990.
- [5] S. M. Kuo and W.-S. S. Gan, *Digital Signal Processors: Architectures, Implementations, and Applications*, Prentice-Hall, Upper Saddle River, NJ, USA, 2005.
- [6] L. Cheded, "Exact recovery of higher order moments," *IEEE Transactions on Information Theory*, vol. 44, no. 2, pp. 851–858, 1998.
- [7] R. M. Gray and T. G. Stockham Jr., "Dithered quantizers," *IEEE Transactions on Information Theory*, vol. 39, no. 3, pp. 805–812, 1993.
- [8] R. A. Wannamaker, S. P. Lipshitz, J. Vanderkooy, and J. N. Wright, "A theory of nonsubtractive dither," *IEEE Transactions on Signal Processing*, vol. 48, no. 2, pp. 499–516, 2000.
- [9] L. Cheded and S. Akhtar, "On the FFT of 1-bit dither-quantized signals," in *Proceedings of 10th IEEE Technical Exchange Meeting (TEM '03)*, Dhahran, Saudi Arabia, March 2003.
- [10] L. Cheded, "On the exact recovery of the FFT of noisy signals using a non-subtractively dither-quantized input channel," in *Proceedings of 7th International Symposium on Signal Processing and Its Applications (ISSPA '03)*, vol. 2, pp. 539–542, Paris, France, July 2003.
- [11] L. Cheded and S. Akhtar, "A new, fast and low-cost FFT estimation scheme of signals using 1-bit non-subtractive dithered quantization," in *Proceedings of the 6th Nordic Signal Processing Symposium (NORSIG '04)*, pp. 236–239, Espoo, Finland, June 2004.
- [12] L. Cheded and S. Akhtar, "A novel and fast 1-bit FFT scheme with two dither-quantized channels," in *Proceedings of 12th European Signal Processing Conference (EUSIPCO '04)*, Vienna, Austria, September 2004.
- [13] L. Cheded, "On the exact recovery of cumulants," in *Proceedings of 4th International Conference on Signal Processing (ICSP '98)*, vol. 1, pp. 423–426, Beijing, China, October 1998.
- [14] L. Cheded and S. Akhtar, "A new and fast frequency response estimation technique for noisy systems," in *Proceedings of 35th Asilomar Conference on Signals, Systems and Computers (Asilomar '01)*, vol. 2, pp. 1374–1378, Pacific Grove, Calif, USA, November 2001.
- [15] L. Cheded, "Theory for fast and cost-effective frequency response estimation of systems," *IEE Proceedings - Vision, Image, & Signal Processing*, vol. 151, no. 6, pp. 467–479, 2004.

- [16] H. Stark and J. W. Woods, *Probability, Random Processes, and Estimation Theory for Engineers*, Prentice-Hall, Englewood Cliffs, NJ, USA, 2nd edition, 1994.

**L. Cheded** gained his B.S. (honors) in physics from Oran university (Algeria) in 1975, his M.S. in electronic control engineering from Salford University (UK) in 1979, and his Ph.D. in signal processing from UMIST, Manchester (UK), in 1988. While at UMIST, he taught physics to first-year students in the textile technology department (1980–1981) and was a Research Assistant in the DIAS department (1981–1984). He has been with the Systems Engineering Department of the King Fahd University of Petroleum & Minerals University in Saudi Arabia since Sept. 1984, where he is currently an Associate Professor. He was a Visiting Researcher in DIAS (UMIST) during summer 1996. His research interests are mainly focused on DSP (theory and applications) and its interactions with control and estimation theory. Besides teaching, he is also involved in research, supervision of students, reviewing papers for international conferences and journals (including IEEE Tr. SP), and research proposals for grant-awarding agencies. He is a Senior Member of the IEEE, a Member of the IEEE, and EURASIP. He served on the editorial board of the International Journal of Electrical Engineering Education during 1994–2005.



**S. Akhtar** received his B.S. degree in electrical (communication) engineering in 1993 from the University of Engineering and Technology, Lahore, Pakistan. He received his M.S. systems (control) engineering degree in 1998 from King Fahd University of Petroleum and Minerals, Saudi Arabia. He worked as Instrument/Control Engineer at a chemical plant for three years from 1993 to 1996. He is now a Lecturer at King Fahd University of Petroleum and Minerals, Saudi Arabia. His research interests include pattern recognition, signal processing, and neural and fuzzy techniques for control and measurement.

