



HAL
open science

Leveraging Power Spectral Density for Scalable System-Level Accuracy Evaluation

Benjamin Barrois, Karthick Parashar, Olivier Sentieys

► **To cite this version:**

Benjamin Barrois, Karthick Parashar, Olivier Sentieys. Leveraging Power Spectral Density for Scalable System-Level Accuracy Evaluation. IEEE/ACM Conference on Design Automation and Test in Europe (DATE), Mar 2016, Dresden, Germany. pp.6. hal-01253494

HAL Id: hal-01253494

<https://inria.hal.science/hal-01253494v1>

Submitted on 10 Jan 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Leveraging Power Spectral Density for Scalable System-Level Accuracy Evaluation

Benjamin Barrois

University of Rennes 1, INRIA/IRISA

6 rue de Kerampont

F-22300 Lannion

Email: benjamin.barrois@irisa.fr

Karthick Parashar

IMEC vzw

Kapeldreef 75

B-3001 Leuven

Email: karthick.parashar@imec.be

Olivier Sentieys

INRIA/IRISA, University of Rennes 1

6 rue de Kerampont

F-22300 Lannion

Email: olivier.sentieys@irisa.fr

Abstract—The choice of fixed-point word-lengths critically impacts the system performance by impacting the quality of computation, its energy, speed and area. Making a good choice of fixed-point word-length generally requires solving an NP-hard problem by exploring a vast search space. Therefore, the entire fixed-point refinement process becomes critically dependent on evaluating the effects of accuracy degradation. In this paper, a novel technique for the system-level evaluation of fixed-point systems, which is more scalable and that renders better accuracy, is proposed. This technique makes use of the information hidden in the power-spectral density of quantization noises. It is shown to be very effective in systems consisting of more than one frequency sensitive components. Compared to state-of-the-art hierarchical methods that are agnostic to the quantization noise spectrum, we show that the proposed approach is $5\times$ to $500\times$ more accurate on some representative signal processing kernels.

I. INTRODUCTION

Signal processing applications popularly use fixed-point data types for implementation. The choice of fixed-point data types is driven usually by cost constraints such as power, area and timing. The objective of fixed-point refinement during the design process is to make sure that chosen data types are precise enough to achieve the expected quality of computation while minimizing the cost constraint. In the recent years, several application scenarios resilient to errors of various kinds are being explored in the context of approximate computing [1], [2]. The errors with fixed-point data types are classified into two types arising from finite precision on one hand and finite dynamic range on the other. Although the impact of errors due to violation of finite dynamic range is more pronounced, these errors can be mitigated by techniques such as range analysis using affine arithmetic, interval arithmetic or more complex statistical techniques such as [3]. In spite of allowing for good dynamic range, the lack of precision causes errors that are perceived as bad quality of computation. In case of wireless applications, this can be measured as bit error rate (BER), in image and signal processing as signal to quantization noise ratio (SQNR), and, in general, as quantization noise power. Measuring the impact of finite precision on the output quality of computation is discussed in this paper.

Commercial tools for performing the fixed-point accuracy analysis are primarily based on facilitating fixed-point simula-

tion with user-defined word-lengths using software fixed-point constructs and libraries. Although very useful, evaluation by simulation can be very time consuming. The time required for fixed-point evaluation grows in proportion with the number of fixed-point variables and also the number of input sample size.

Using the analytical approach for accuracy evaluation, the noise power is obtained by evaluating a closed-form expression as a function of the number of bits assigned to various signals in the system. This approach requires a one-time effort for arriving at the closed-form expression for a given system. The analytical technique evaluates the first two moments of the quantization noise sources and propagates it through the signal-flow graph from all noise sources to the system output. On relatively small systems, the evaluation of path functions can be accomplished manually. As the system complexity grows, it would require automation support and eventually for very large systems, the automation could also prove painstakingly slow. Therefore, several divide and conquer approaches have been proposed such as [4], [5] to overcome the apparent complexity of large systems which respectively suffer from loss of information or enumerating all paths in the graph.

In this paper, we provide an alternative analytical accuracy evaluation approach for use with hierarchical techniques to be applied on linear and time invariant (LTI) systems. This technique captures the information associated with the frequency spread of quantization noise power by sampling its power spectral density (PSD). We show in this paper, how such information can be used for breaking the complexity of evaluating quantization noise at the output of large signal processing systems. Contributions of this paper are as follows:

- quantifying the accuracy of the proposed technique based on PSD propagation and
- demonstrating its high scalability at system level resulting from linear time complexity.

The rest of this paper is organized as follows. Section II reviews analytical methods for accuracy analysis at the algorithm level of errors due to finite arithmetic effects in systems using fixed-point arithmetic. In Section III, the proposed estimation method based on PSD is introduced and developed for general systems. Finally, in Section IV, two representative signal-processing benchmarks are chosen to showcase the efficiency

of the proposed method.

II. RELATED WORK ON ACCURACY ANALYSIS

The loss in accuracy due to finite precision imposed by fixed-point numbering format has been evaluated using several metrics. The most common among them are the error bounds and the mean-square error (*MSE*). While the first metric is used to determine the worst case impact, the *MSE* is an average case metric very useful in tuning the average performance of the system under consideration in terms of its energy and timing. Although the finite precision accuracy must be compared with infinite precision (or arbitrary precision) numbers, it is impossible to do so while simulating using a computer. So, the IEEE double-precision floating-point format, whose dynamic range and precision are several orders of magnitude higher compared to typical fixed-point word lengths is considered as the reference for all comparison purposes.

In the literature, the mean square value of the differences between computations by fixed-point system and the reference system implementation is also referred to as quantization noise power. This is a scalar quantity and it changes as a function of the fixed-point word-length. Evaluation of quantization noise power at the output of a fixed-point system is either performed by simulation-based technique or using analytical techniques. Simulation-based techniques are universal and can be made use of as long as there are enough computational resources. By the nature of it, simulation-based techniques take longer time and are subjected to the input stimulus bias. Analytical techniques, on the other hand, provide a closed-form expression for calculating the quantization noise power as a function of fixed-point word-lengths. However, they are limited due to their dependence upon the following properties [6]:

- 1) Quantization noise and the signal are uncorrelated.
- 2) Quantization noise at its source is spectrally white.
- 3) Effect of a small perturbation at the input of the operation generates a linearly proportional perturbation at the output of the operation.

The first two properties pertain to the quantization noise source under conditions defined in the pseudo-quantization-noise (PQN) model, the statistics of the noise and signal are uncorrelated and even though the signal itself may be correlated in time, the noise signal is uncorrelated in time [6].

The third property relates to the application of “perturbation theory” [7]. It is possible to propagate quantization noise through as long as the function defined by the operation can be linearized. Consider a binary operator whose inputs are x and y and the output is z . If the input signals be perturbed by b_x and b_y to obtain \bar{x} and \bar{y} respectively, the output is perturbed by the quantity b_z to obtain \bar{z} . In other words, as long as the fixed-point operator is smooth, the impact of small perturbations at the input translates to perturbation at the output of the operator without any change in its macroscopic behavior. In the realm of perturbation theory, the output noise b_z is a linear combination of the two input noises b_x and b_y such as

$$b_z = \nu_1 b_x + \nu_2 b_y \quad (1)$$

where ν_1, ν_2 are obtained from a first-order Taylor approximation [7] of the continuous and differentiable function f :

$$\begin{aligned} \bar{z} &= f(\bar{x}, \bar{y}) \\ &\simeq f(x, y) + \frac{\partial f}{\partial x}(x, y) \cdot (\bar{x} - x) + \frac{\partial f}{\partial y}(x, y) \cdot (\bar{y} - y). \end{aligned} \quad (2)$$

Therefore, the expression of the terms ν_1 and ν_2 are given as

$$\nu_1 = \frac{\partial f}{\partial x}(x, y) \quad \nu_2 = \frac{\partial f}{\partial y}(x, y). \quad (3)$$

Following the third property of quantization noise enumerated above, a further assumption for Eq. 1 to hold true is that the noise terms b_x and b_y are uncorrelated with one another. It has to be noted here that the terms ν_1 and ν_2 can be time varying. This method is not limited to binary operations only. In fact, this method can be applied at the functional level with any number of inputs and outputs and to all operators on a given data path in order to propagate the quantization noise from all error sources to the output.

When above conditions hold true, the output quantization noise power of the system is obtained by linear propagation of all quantization noise sources [8] as

$$E[b_y^2] = \sum_{i=1}^{N_e} K_i \sigma_i^2 + \sum_{i=1}^{N_e} \sum_{j=1}^{N_e} L_{ij} \mu_i \mu_j \quad (4)$$

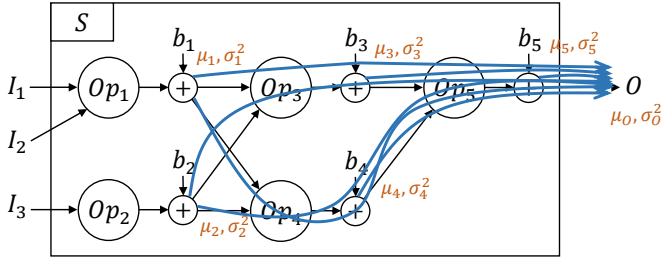
where $E[\cdot]$ is the expectation function, b_y is the error signal associated with its corresponding system output signal y . The system under consideration consists of N_e fixed-point operations and the i^{th} operation is generating quantization noise b_i with mean and standard deviation μ_i and σ_i . Fig. 1.a illustrates this noise propagation. The terms K_i and L_{ij} are constants and depend on the path function h_i from the i^{th} source to the output y and are calculated as

$$K_i = \sum_{k=-\infty}^{\infty} E[h_i^2(k)], \quad (5)$$

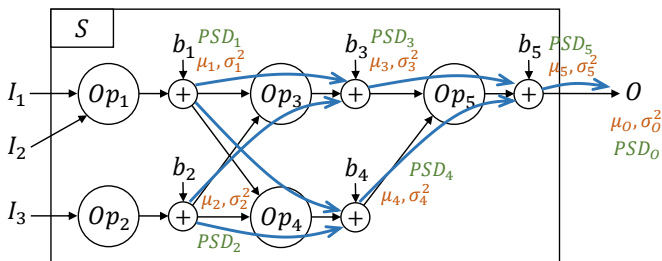
$$L_{ij} = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} E[h_i(k)h_j(l)]. \quad (6)$$

Hierarchical techniques for evaluation of quantization noise power have been proposed [9], [4] to overcome the scalability concerns associated with fixed-point systems. In this approach, the system components are evaluated one at a time and then combined by superposition at the output (Fig. 1.b, blind propagation of μ_i, σ_i^2). If simulation-based technique is used for evaluation of quantization noise power at the output, the hierarchical evaluation process helps parallelize simulation of each of the components. When employing analytical technique such as the one in Eq. 4, the number of paths required to be evaluated is reduced dramatically. This reduction is very interesting from the design automation perspective. The paths are broken around the system component boundaries and each component can be evaluated separately thereby reducing the burden of semantic analysis. However, it has to be borne in mind that the application of the technique in Eq. 4 requires that

the quantization noise satisfies the three properties enumerated above and also that the noise signals are always uncorrelated, which is often false and can cause severe errors. This paper addresses this problem and suggests a technique that exploits the information hidden in the power spectral density (PSD) of the quantization noise [6] [10] signal to achieve very accurate estimates.



(a) Traditional flat method: propagation of μ_i, σ_i^2 from each noise to output



(b) PSD agnostic: blind propagation of μ_i, σ_i^2 . Proposed PSD method: propagation of μ_i, σ_i^2 , and PSD_i

Fig. 1: Comparison of noise parameters propagation using traditional flat, PSD agnostic and proposed PSD methods

III. PSD-BASED ACCURACY EVALUATION

It is clear from the state of the art that there exists two types of limitations to the existing accuracy evaluation techniques. While the analytical technique reduces the simulation time greatly, its preprocessing time can grow exponentially requiring to employ hierarchical techniques such as [4]. However, these techniques introduce the problem of inaccuracy by approximating error quantities with just mean and variance. This is especially true in cases when large systems are broken down to smaller sub-systems for analysis. For illustration, consider the system shown in Fig. 1.b. The system S consists of several sub-systems marked as $Op_{1...5}$. The noise generated at the output of each system, correspondingly marked as $b_{1...5}$, is propagated (blue arrows) through several parts of the system for calculation of the moments of error at the system output. Suppose there are memory elements in Op_1 and Op_2 , propagation of noise b_1 and b_2 (say) through Op_3 by just using the first two moments of the quantization noise (as described in the previous section) can lead to errors in estimates at the output of Op_3 which can further be amplified by Op_5 all the way to output O . Similarly, the path through Op_4 also influences the error of the estimate through Op_5 leading to

very large error margins for O . In order to analytically arrive at the moments of the system output, additional information pertaining to quantization noise at points of convergence of two or more noise paths is required. We refer to the methods that do not consider PSD information (such as [9]) as PSD-agnostic methods. In this section, we propose a technique which efficiently makes use of power-spectral-density (PSD) of the quantization noise for evaluating the error output and which is scalable both in terms of accuracy and system size.

A. PSD of a quantization noise

A large signal processing system can be divided into a number of sub-systems, each characterized by its transfer function. The transfer function defines the magnitude and phase relationship of the path for input signals of different frequencies. Since our interest is only the noise power, we ignore the phase spectrum and consider only its magnitude spectrum or the PSD. With the knowledge of the PSD distribution of the input and the system PSD profile it is possible to calculate the PSD of the output. The PSD $S_{xx}(F)$ of a signal x at any normalized frequency F is defined as the Fourier transform ($\mathcal{F}\{\cdot\}$) of the autocorrelation function of x as

$$S_{xx}(F) = \mathcal{F}\{x(n) \cdot x^*(n+m)\}, \quad (7)$$

$$S_{xx}(F) = \mathcal{F}\{x\} \cdot \mathcal{F}\{x\}^* = |\mathcal{F}\{x\}|^2. \quad (8)$$

With the knowledge of the PSD of x , the MSE and the mean of x is obtained by summing up the power in each frequency component as

$$E[x^2] = \int_{-1}^1 S_{xx}(F) dF = \mu^2 + \sigma^2$$

$$S_{xx}(0) = \mu^2. \quad (9)$$

The PSD of the quantization noise generated by a fixed-point data type with d fractional bits is (as discussed in Section II) white except for $F = 0$, which depends on mean. By discretizing the PSD into N_{PSD} regular bins including the DC component, the PSD of a generated quantization noise b_x is given by

$$S_{b_x}(F) = \begin{cases} \frac{1}{N_{PSD}} \sigma^2 & \text{if } F \neq 0, \\ \mu^2 & \text{if } F = 0. \end{cases} \quad (10)$$

where mean and variance μ and σ^2 for both truncation and rounding modes with d bits is as given in [6].

B. PSD propagation across a fixed-point LTI system

In this paper, we will focus on linear and time-invariant (LTI) systems, which constitute the major part of signal processing systems. An LTI system can be represented by a signal flow graph (SFG) composed of boxes corresponding to sub-systems defined by their impulse response and delimited by additive quantization noise sources such as the one described in Section II. The proposed PSD evaluation method then consists of three steps:

- 1) Detect cycles in SFG and break them to obtain an equivalent acyclic SFG that can be used for noise propagation using classical SFG transformations [11].

- 2) The discrete PSD of each signal processing block and the additive noise associated with the input signal is calculated on N_{PSD} points.
- 3) The noise PSD parameters are propagated from inputs to outputs, using Equations 11 and 14.

Let x be the input of a system of impulse response h . Then the output y is obtained by the convolution operation ($*$) of x and h as $y = x * h$. In the Fourier transform domain it can be written as $Y = X \cdot H$ where $Y = \mathcal{F}\{y\}$. Following this, the output PSD $S_{yy}(F)$ is obtained as [10]

$$S_{yy}(F) = S_{xx}(F) \cdot \|H(F)\|^2. \quad (11)$$

where $\|H(F)\|$ is the magnitude response of the system h .

In any signal processing system, the quantization noise sources from various inputs converge in at either an adder or a multiplier. Considering the LTI subset, multiplications are nothing but multiplication with constants and hence correspond to linear scaling factors for noise powers. In the case of adders, if the sum of two quantities x and y is obtained as $z = x + y$, then $S_{zz}(F)$ is given by

$$S_{zz}(F) = S_{xx}(F) + S_{yy}(F) + S_{xy}(F) + S_{yx}(F), \quad (12)$$

where $S_{xy}(F)$ is obtained using the cross-correlation spectrum of x and y and is obtained as

$$S_{xy}(F) = \mathcal{F}\{x(n) \cdot y^*(n+m)\}. \quad (13)$$

Also, $S_{yx}(F)$ is obtained as the complex conjugate of $S_{xy}(F)$. Indeed, if x and y are uncorrelated, the cross-correlation is rendered zero and $S_{zz}(F)$ is simply the sum of $S_{xx}(F)$ and $S_{yy}(F)$.

$$S_{zz}(F) = S_{xx}(F) + S_{yy}(F). \quad (14)$$

The complexity of propagating PSD parameters through the system essentially depends on the number of discrete points N_{PSD} . The total time for evaluation of the PSD parameters can be split into two parts: one, τ_{pp} corresponding to the preprocessing stage which involves evaluating the N_{PSD} -point Fourier transform of transfer function of the sub-systems with complexity $\mathcal{O}\{N \log(N)\}$; two, the actual time required for evaluation τ_{eval} which is $\mathcal{O}\{N\}$ from Equations 11 and 14. τ_{eval} is required for evaluating the accuracy for various inputs and can be repeatedly performed without any preprocessing say N_{eval} times. Since the time spent on pre-processing is a one time investment, the actual evaluation time is dominated by the τ_{eval} which is linear with N_{PSD} .

IV. RESULTS

In this section, the proposed method is evaluated using a three step approach. First, we show experimentally that the estimates obtained by proposed PSD technique are close to simulation. Then, we present the impact of choosing the number N_{PSD} to capture PSD information on the accuracy as well as the execution times of the proposed approach. Finally, we also discuss the improved accuracy in estimation

and compare it with the result obtained by PSD agnostic method.

All experiments are performed using Matlab R2014b. The MSE deviation E_d is chosen as the metric for comparison in all these experiments. It is calculated as

$$E_d = \frac{E[err_{sim}^2] - E[err_{est}^2]}{E[err_{sim}^2]}, \quad (15)$$

where $E[err_{sim}^2]$ is the output error power obtained by simulation and $E[err_{est}^2]$ is obtained by proposed analytical estimation. From this metric, an accuracy equivalent to less than one bit corresponds to the range $E_d \in (-75\%, 300\%)$, which can be trivially proven considering the error power relative to two successive word-lengths. In the following sections, we first present the experiments and provide a discussion of the results obtained.

A. Experiment Setup

1) *FIR, IIR filters*: The first experiment consists of evaluating the PSD of a single FIR (finite impulse response) and IIR (infinite impulse response) filter blocks as described in Section III. The quantized input signal is propagated through the chosen filter and the output quantization noise power is measured by simulation and by the proposed PSD method. The error in estimates of the noise power E_d is obtained on a total of 147 FIR and 147 IIR filters obtained by attributing different functionalities (bandpass, low-pass and hi-pass), various taps involving memory elements between 16 and 128 taps for FIR filters and from 2 to 10 taps for IIR filter. Simulation is run on 10^6 inputs and PSD estimation is performed on 1024 samples.

2) *Frequency Domain Filtering*: The system described in Figure 2 is a frequency domain band-pass filter. It consists in a 16-tap low-pass FIR filter H_{hp} followed by a frequency-domain filter, composed of a 16-point FFT block, a multiplication by the 16 coefficients of a high-pass FIR filter H_{lp} and an inverse FFT. The frequency domain filter applies the filter using the popular overlap save method. Simulations are carried out on a set of 10^7 input samples.

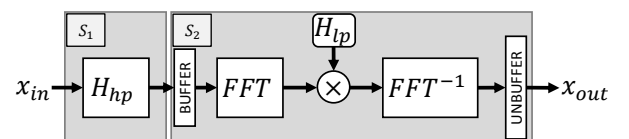


Fig. 2: Band-pass frequency filtering scheme

3) *Daubechies 9/7 Wavelet*: A 2-level *Daubechies 9/7 DWT* pair which forms the basis of many modern image and video codecs such as JPEG-2000, H.264 is shown in Figure3. For this experiment, 196 grayscale images extracted from USC-SIPI and RPI-CIPR image databases and from Brodatz texture images [12] used generally for evaluating JPEG2000 compression algorithms. Two levels of sub-band decomposition are performed on the sample images using the hierarchical signal flow graph. For the encoder, the first filtering and downsampling is applied on rows and the second

one on columns. Then, the second level coding is applied on the low-pass components (x_{ll}). Symmetrically, The decoder first performs upsampling and filtering is applied on columns followed by the second upsampling and filtering for the rows. For this experiment, fractional word-lengths d of all variables are set to the same value and are varied across 8 – 32 bits in steps of 4 and N_{PSD} is set to 1024 .

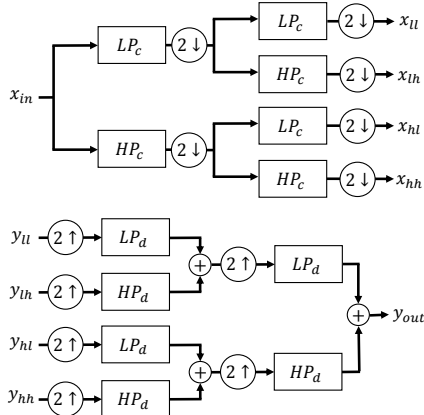


Fig. 3: 1-level DWT coder and decoder

B. Validation of the approach for LTI systems

The *min*, *max* and *absolute mean* of E_d for FIR and IIR filters is given in Table I. In the case of FIR filters, E_d is

	FIR filters	IIR filters
$\min(E_d)$	-0.37%	-19.4%
$\max(E_d)$	0.37%	31.2%
$\text{mean}(E_d)$	0.11%	9.44%

TABLE I: Relative error power estimation statistics E_d

contained within an absolute value of 0.37% in comparison with simulation. In the case of IIR filters, E_d bounds are higher because of their recursive nature and the high filter orders tested. FIR and IIR filters shows an average absolute E_d of respectively 0.11% and 9.44%, showing a generally very accurate estimation. For both, the accuracy is anyway largely less than one-bit equivalent. Moreover, classical flat estimation [8] applied to the same filters gives *exactly the same results* in terms of E_d , showing their strict equivalence on an elementary filtering block.

Figure 4 presents the results for the other two experiments as the number of fractional bits are changed between 8 and 32 bits with a maximum deviation in error of only about 10%. The maximum error in estimate is by far too small to make an impact on the final optimization.

C. Influence of the number of PSD samples

The proposed PSD estimation method achieves very good accuracy with a large number of sampling PSD samples. However, as discussed in Section III-B, larger number of N_{PSD} increases the evaluation time. Therefore, it would

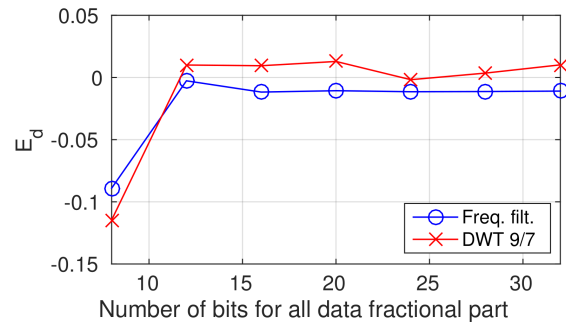


Fig. 4: E_d versus fractional bit-width d

be of interest to know the impact of finding out how this choice affects the estimation accuracy. To observe this, in both examples chosen in this paper the fixed-point error is obtained by both simulation and the proposed PSD method with different values of N_{PSD} in powers of 2 ranging from 16 to 1024. In this example, fractional bit-width d is uniformly set to 32 for all signals. Output error power deviation E_d value for this experiment is plotted Figure 5 versus N_{PSD} . As expected, increasing the number of PSD samples leads to an improvement of E_d . For $N_{PSD} = 16$, E_d is slightly inferior to -8% for the frequency filtering system, and slightly superior to 1% for the DWT system. Then, both curves tend to a value inside $\pm 1\%$. The accuracy obtained is better than the sub-one-bit objective. The accuracy of estimates obtained using the proposed method is a function of the system complexity.

D. Comparison with PSD-agnostic methods

The deviation of the error estimates between the proposed and the PSD agnostic method is tabulated in Table II. The max error is obtained with $N_{PSD} = 16$ and min error is obtained with $N_{PSD} = 1024$. In all cases, it can be observed that the PSD agnostic method is much more erroneous than even the maximum error obtained using the proposed technique. It has to be noted that for the DWT example, the PSD agnostic method renders an error of 610%. The PSD agnostic method is $4.5\times$ worse off in its estimate for frequency filtering, and $554\times$ for DWT. For the best case, these values raise respectively to $3.5 \cdot 10^3\times$ and $6.7 \cdot 10^4\times$.

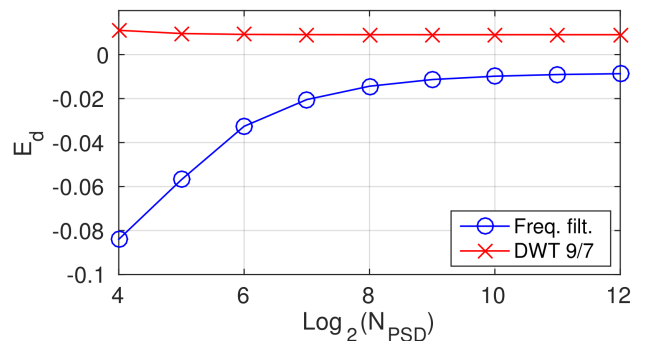


Fig. 5: E_d versus number of PSD samples N_{PSD}

	Proposed PSD method (max accuracy)	Proposed PSD method (min accuracy)	PSD agnostic method
Freq. Filt.	-8.40%	-0.87%	29.5%
DWT 9/7	1.10%	0.90%	610%

TABLE II: Comparison of E_d between PSD agnostic method and proposed PSD method

Time spent on this estimation is usually another critical resource. Figure 6 gives the time of output error estimation using the proposed PSD method versus N_{PSD} . With $N_{PSD} = 16$ the proposed method requires about one millisecond in case of both experiments. With more PSD samples, the time taken by frequency filtering example grows slower than Daubechies DWT example owing to its small size. A speed up factor of 3 – 5 orders of magnitude compared to simulation is obtained in both cases even for the highest value of N_{PSD} .

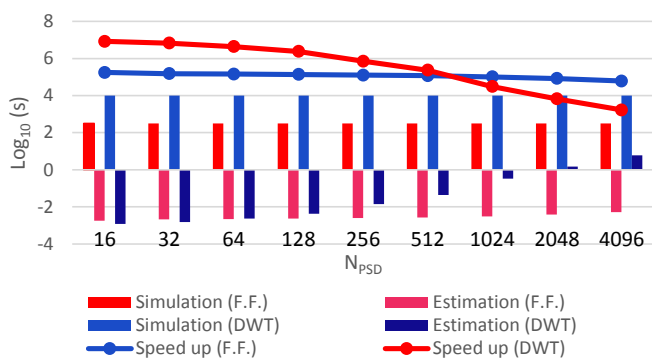


Fig. 6: Execution time in seconds and speed up for frequency filtering and DWT systems versus the number of PSD samples

E. Frequency repartition of output error

Another interesting feature inherent to the proposed estimation method is to know the frequency repartition of errors, which is relevant for refinement of fixed-point signal processing systems, and which is not estimated with conventional methods. Figure 7 gives a visual comparison between the PSDs of output error obtained by intensive simulation and PSD method on 1024 samples for a 2-level DWT encoding and decoding, with all data fractional parts set to 12 bits. Black to white values represent log-normalized low to high errors. The center of the image represents low frequencies, while the borders represent the high ones. These images show that proposed method achieves a very good estimation of frequency repartition of the output error, taking only a few milliseconds whereas simulation takes hours. Such a fast and accurate information can be used for refining the system word-lengths to reach a better output quality, basing the refinement not only on output error intensity but also on what frequency repartition is best for the application.

V. CONCLUSIONS

This paper has proposed the characterization and propagation of quantization noises in a fixed-point signal processing

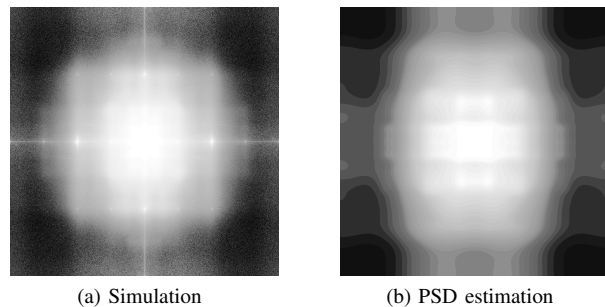


Fig. 7: Output frequency repartition of the fixed-point error after DWT encoding and decoding

system using its power spectral density. This method is applied at block level, which dramatically reduces the complexity of fixed-point system evaluation when compared to classical flat estimation method. It therefore leads to a significant speed up for accuracy evaluation, going from 3 to 5 orders of magnitude when compared to Monte-Carlo simulation in tested examples. Results demonstrate that the proposed estimation method leveraging PSD information achieves a less than one-bit accuracy with a large margin. They also show that complexity-equivalent PSD-agnostic techniques evaluate the accuracy with large errors. The proposed PSD technique also allows the observation of useful frequential properties of the output error that could not be achieved with conventional scalar methods.

REFERENCES

- [1] K. Alajel, W. Xiang, and J. Leis, "Error resilience performance evaluation of H.264 I-frame and JPWL for wireless image transmission," in *Int. Conf. Sig. Proc. and Com. Sys. (ICSPCS)*, 2010, pp. 1–7.
- [2] H. R. Wu, W. Lin, and L. Karam, "An overview of perceptual processing for digital pictures," in *Multimedia and Expo Workshops (ICMEW'12)*, July 2012, pp. 113–120.
- [3] E. Özer, A. P. Nisbet, and D. Gregg, "Stochastic bit-width approximation using extreme value theory for customizable processors," in *Compiler Construction*. Springer, 2004, pp. 250–264.
- [4] K. Parashar, R. Rocher, D. Menard, and O. Sentieys, "A hierarchical methodology for word-length optimization of signal processing systems," in *23rd Int. Conf. on VLSI Design (VLSID)*, 2010, pp. 318–323.
- [5] D. Novo, I. Tzimi, U. Ahmad, P. Ienne, and F. Catthoor, "Cracking the complexity of fixed-point refinement in complex wireless systems," in *IEEE Work. on Signal Processing Systems (SiPS'13)*, 2013, pp. 18–23.
- [6] B. Widrow and I. Kollár, *Quantization Noise: Roundoff Error in Digital Computation, Signal Processing, Control, and Communications*. Cambridge University Press, 2008.
- [7] G. Constantinides, "Perturbation analysis for word-length optimization," in *11th Annual IEEE Symposium on FCCM*, 2003, pp. 81–90.
- [8] D. Menard, R. Rocher, and O. Sentieys, "Analytical Fixed-Point Accuracy Evaluation in Linear Time-Invariant Systems," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 55, no. 1, Nov. 2008.
- [9] J.-W. Weijers, V. Derudder, S. Janssens, F. Petré, and A. Bourdoux, "From MIMO-OFDM algorithms to a real-time wireless prototype: a systematic Matlab-to-hardware design flow," *EURASIP J. Appl. Signal Process.*, vol. 2006, pp. 138–138, Jan. 2006.
- [10] L. B. Jackson, "Roundoff-noise analysis for fixed-point digital filters realized in cascade or parallel form," *IEEE Transactions on Audio and Electroacoustics*, vol. 18, no. 2, pp. 107–122, 1970.
- [11] K. Ogata, *Modern Control Engineering*, 4th ed. Upper Saddle River, NJ, USA: Prentice Hall PTR, 2001.
- [12] P. Brodatz, *Textures: A Photographic Album for Artists and Designers*. New-York: Dover Publications, 1966.