# On Finite Wordlength Properties of Block-Floating-Point Arithmetic

Abhijit Mitra

*Abstract*— A special case of floating point data representation is block floating point format where a block of operands are forced to have a joint exponent term. This paper deals with the finite wordlength properties of this data format. The theoretical errors associated with the error model for block floating point quantization process is investigated with the help of error distribution functions. A fast and easy approximation formula for calculating signal-to-noise ratio in quantization to block floating point format is derived. This representation is found to be a useful compromise between fixed point and floating point format due to its acceptable numerical error properties over a wide dynamic range.

*Keywords*— Block floating point, Roundoff error, Block exponent distribution fuction, Signal factor.

## I. INTRODUCTION

THE so-called *block floating point* (BFP) data representation [1]- [2] has come forth recently in several digital audio signal broadcasting standards as a *near instantaneous* data companding (compression + expansion) technique to reduce the bitflow of data per sample. The standards using such a companding method include for example NICAM (digital two-channel sound system for PAL and MAC TV standards), MUSE (Japanese HDTV standard) and DSR (German digital satellite radio) [3]. The BFP format has become the natural choice in the above standards after exhibiting the overly preciseness of the IEEE standard on floating point arithmetic for certain DSP applications and also for achieving a significant hardware simplicity compared to its floating point counterpart [4]-[5]. Morever, it has recently been shown [1] that some properties of BFP format can make it a better alternative than fixed point or floating point arithmetic for implementing recursive linear systems.

In BFP representation, the incoming data is arranged in nonoverlapping blocks and depending on the relative magnitudes of the data samples in each block, a common exponent is assigned. This arrangement, in fact, combines two widely utilized number representation formats, fixed point (FxP) and floating point (FP), to exploit some benefits of FP like the wide dynamic range on one hand, with the simplicity of FxP like reduced computational complexity on the other. The effects of exploiting BFP arithmetic in recursive and non-recursive digital filters was first investigated by Oppenheim [6]. The recent work of Kalliojärvi [7] has shown that BFP formats

are generally very bit efficient schemes. More recent work of Mitra [8]- [9] has proved the effectiveness of this representation in the complicated area of adaptive filtering. However, to the best of our knowledge, no extensive work on the finite-precision properties of this data format has so far been reported in the literature.

In this paper, the finite wordlength properties of BFP is studied from a deterministic viewpoint. In particular, theoretical analysis of a significant quantization error model associated with the BFP format is carried out by deriving the distribution of block exponents. The signal to quantization noise ratio of the BFP format is found to be dependent on this block exponent distribution, which requires precise information about the signal statistics. Thus, an easy-to-use formula is derived for fast and simple calculations and a signal factor, dependent on the block length, is introduced for this specific purpose. A comparison with the FP and the FxP systems, based on several simulations, is also discussed to show the efficiency of the BFP representation over the FP and FxP formats. The outline of this paper is as follows: Section 2 deals with block floating point representation. Section 3 analyses the quantization errors involved in BFP arithmetic within several subsections by taking help of the scaled additive roundoff error model, deriving the block exponent distribution and then introducing a fast and easy formula for calculating signal to noise ratio. Conclusions regarding the efficiency of BFP arithmetic, based on the above analysis and simulation results, are drawn in Section 4.

In the sequel, we shall denote $\mathbb{R}$ as the set of real numbers, $\mathbb{Z}$ as the set of integers, $\mathbb{R}^n$ as the set of length $n$ real vectors and $\mathbb{R}^{n \times n}$ as the set of $(n \times n)$ real matrices.

## II. BFP REPRESENTATION

In binary FP format, any number $x \in \mathbb{R}$ is written as

$$x = sign(x).m.2^{\psi} \tag{1}$$

where the exponent $\psi \in \mathbb{Z}$ is chosen so that the mantissa $m \in \mathbb{R}$ is within the range $[\frac{1}{2}, 1)$, i.e., the mantissas are normalized [10]. The BFP representation can be considered as a special case of FP format, where the incoming data is grouped into nonoverlapping blocks of N consecutive samples and each block has a joint scaling factor corresponding to the data samples with the largest magnitude in the block. In other words, given a data vector $\mathbf{x} (\in \mathbb{R}^N) = [x_1, ..., x_N]$, we represent it as [8]

$$\mathbf{x} = [\overline{x}_1, ..., \overline{x}_N].2^{\gamma} = \overline{\mathbf{x}}.2^{\gamma} \tag{2}$$

where $\overline{x}_i (= x_i.2^{-\gamma})$ represent the signed mantissas for $i =$

TABLE I

DIFFERENT BFP FORMATS USED IN DIGITAL AUDIO TRANSMISSION STANDARDS.

| Standards | N | $B_d + 1$ | $B_\gamma + 1$ | Equivalent bits/sample |
|---|---|---|---|---|
| NICAM | 32 | 10 | 3 | 10.09375 |
| MUSE: A-mode | 32 | 8 | 3 | 8.09375 |
| MUSE: B-mode | 48 | 11 | 3 | 11.0625 |
| DSR | 64 | 14 | 3 | 14.046875 |

$1, 2, ..., N$ and the block exponent $\gamma \in \mathbb{Z}$ is defined as

$$\gamma = \lfloor log_2 Max \rfloor + 1 \qquad (3)$$

with $Max = max(|x_1|, ..., |x_N|)$ and '$\lfloor . \rfloor$' being the so-called floor function, i.e., $\lfloor a \rfloor \in \mathbb{Z}$ is closest but not exceeding $a \in \mathbb{R}$. For determining $\gamma$ in such a way, the mantissas are always block normalized, i.e.,

$$\frac{1}{2} \leq max(|\overline{x}_i|) < 1. \qquad (4)$$

The above BFP representation can also be extended for the data matrix $\mathbf{X}(\in \mathbb{R}^{N \times N}) = [x]_{ij}$ by assigning the new block exponent $\gamma_c \in \mathbb{Z}$ as

$$\gamma_c = \lfloor log_2(max\{|x_{ij}|\}) \rfloor + 1. \qquad (5)$$

It should be noted that if ($B_d$ + one sign) bits are used to represent each mantissa within the block and if ($B_\gamma$ + one sign) bits are used to account for the block exponent, then effectively, under BFP system, each sample can be equivalently represented with $(1 + B_d) + (1 + B_\gamma)/N$ bits because the block exponent is taken only once for the whole block. This particular strength makes this format more considerable than FxP or FP systems. Table 1 shows the different BFP formats used in digital audio transmission standards utilizing such an adavantage.

### III. QUANTIZATION ERROR ANALYSIS

#### A. Quantization Error Model

The roundoff error models used with FxP and FP format can not be applied directly for the case of BFP representation. The additive roundoff error model [10] of the fixed point format, which represents the absolute error in rounding process, should not be used as BFP is a scaled number representation system with different exponents for different blocks. In FP arithmetic on the other hand, relative error is more important than absolute error as the quantization affects only the mantissa, meaning that FP errors are multiplicative rather than additive. Thus, a relative roundoff error model is used in the FP analysis and the relative error $\epsilon$ is defined as

$$Q[x] = x(1 + \epsilon) = x + x\epsilon$$
$$\Rightarrow \epsilon = \frac{Q[x] - x}{x} = \frac{Q[m] - m}{m} = \frac{\xi}{m} \qquad (6)$$

where $Q[.]$ denotes the respective quantity after quantization and $\xi$ is the mantissa quantization error. This relative roundoff error model, however, can not be used for investigating the BFP roundoff errors because the relative errors will not be bounded in this case. This is due to the fact that the mantissas in BFP arithmetic can be arbitrarily close to zero, unlike the normalized FP mantissas which cannot take a value beyond the lower limit $\frac{1}{2}$. Therefore, the BFP quantization error is modeled with a *scaled additive roundoff error model* [2] or *mantissa additive roundoff error model*, defined as follows

$$\alpha = Q[x_i] - x_i$$
$$= (Q[\overline{x}_i] - \overline{x}_i).2^\gamma = e_m.2^\gamma \qquad (7)$$

where $Q[.] \in \mathbb{R}$ denotes the quantized value of a quantity and $e_m$ is the mantissa quantization error.

#### B. Error Distribution Functions

The quantized mantissas $Q[\overline{x}_i]$ are assumed to be uniformly distributed in the interval $[0, (1 - \triangle)]$ and the mantissa error $e_m$ is bounded by

$$-\triangle/2 \leq e_m \leq \triangle/2 \qquad (8)$$

for *rounding-to-nearest*, and by

$$-\triangle \leq e_m \leq 0 \qquad (9)$$

for *truncation*, where $\triangle = 2^{-B_d}$ is the quantization step size of the mantissas. The above said quantization error $\alpha$ and the mantissa quantization error $e_m$ can be considered as white noise sequences (with zero mean when rounding-to-nearest method is used) as well as uncorrelated with the data sequence $x(n)$. The mantissa roundoff errors are also uncorrelated with the block exponents $\gamma$. These important properties have been verified by quantizing data to BFP format and simulating the respective autocorrelation and crosscorrelation functions. As an example, auto- and cross-correlation functions in quantizing 0 dB Gaussian data (2000 samples) to 1+7+(1+3)/8-bit BFP format are plotted in Figures 1, 2 and 3, which clearly show that $\alpha$ and $e_m$ can be considered as white noise sequences, and $\alpha$ is uncorrelated with the data sequence $x(n)$.

From these Figures, it is clear that certain joint distributions can be split up due to the statistical independence property [12], e.g., if $p_1(e_m)$ and $p_2(m)$ are the probability densities of $e_m$ and $Q[\overline{x}_i]$ respectively, and $p_3(e_m, m)$ is the joint probability density of these variables, then we can write
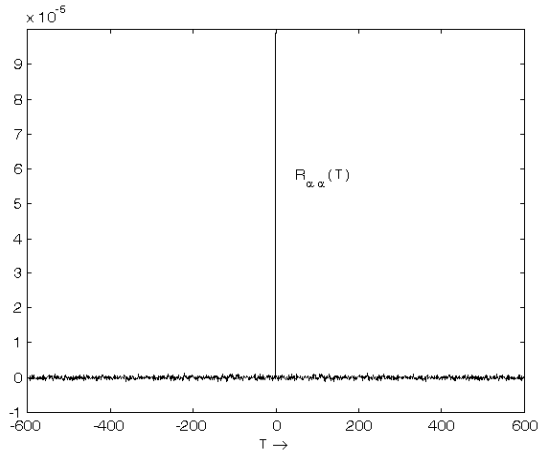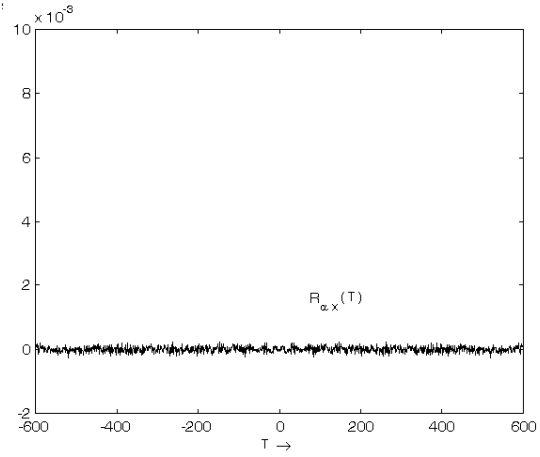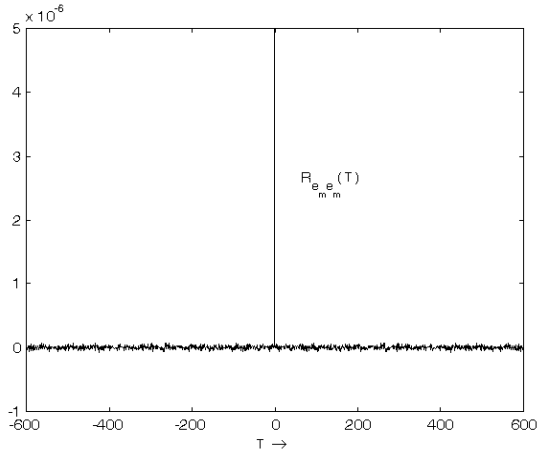
$$p_3(e_m, m) = p_1(e_m).p_2(m) \qquad (10)$$

where, in the case of *rounding*, due to eq.(8)

$$p_1(e_m) = \begin{cases} 1/\triangle & \text{for } |e_m| \leq \triangle/2 \\ 0 & \text{otherwise} \end{cases} \qquad (11)$$

and, in case of *truncation*, due to eq.(9)

$$p_1(e_{b_m}) = \begin{cases} 1/\triangle & \text{for } -\triangle \leq e_m \leq 0 \\ 0 & \text{otherwise.} \end{cases} \qquad (12)$$

Fig. 1 Autocorrelation function of scaled additive roundoff error $\alpha$ ($R_{\alpha\alpha}(T)$).



Fig. 3 Crosscorrelation function of $\alpha$ and quantized data $Q[x]$ ($R_{\alpha x}(T)$).



Fig. 2 Autocorrelation function of mantissa quantization error $e_m$ ($R_{e_m e_m}(T)$).

Also note that, if $p_4(\gamma)$ denotes the probability mass function of integer valued block exponents, then we can write $p_5(e_m, \gamma) = p_1(e_m).p_4(\gamma)$. In general, finding out the value of $p_4(\gamma)$ is a difficult task. A simplified approach to find out the same is given next.

### C. Distribution of Block Exponents

The derivation of the *probability mass function* (PMF)[1] of block exponents $p_\gamma(\gamma_k)$ ($k = 1, ..., N_\gamma$, where $N_\gamma$ is the available distinct block exponent levels) is a tedious task in the

[1] For any quantized random variable, the probability distribution is defined as the sum of different probability mass functions where the dummy variable can assume the values greater than or equal to the random variable. A detailed study can be found in Oppenheim and Schafer [14].

analysis of BFP roundoff errors because $p_\gamma(\gamma_k)$ can not be selected as a uniform distribution for its direct relation with the input data. However, it can be derived in the following way: The probability of the largest magnitude block variable $M$ to be at most $t$ is

$$\Phi(t) = P\{M \leq t\}$$
$$= \int_{-t}^{t} ... \int_{-t}^{t} \psi_{x_1,...,x_N}(x_1, ..., x_N)\, dx_1...dx_N \quad (13)$$

where $\psi_{x_1,...,x_N}(x_1, ..., x_N)$ is the *joint probability density function* (JPDF) [12] of the block variables. The block exponent is $\gamma_k$ if the largest magnitude of block variables M is in the interval $2^{\gamma_k - 1} \leq M < 2^{\gamma_k}$. Then we can write [2]

$$p_\gamma(\gamma_k) = \Phi(2^{\gamma_k}) - \Phi(2^{\gamma_k - 1}). \quad (14)$$

In general, the derivation of this PMF requires evaluation of N dimensional integrals, which is difficult. Therefore it is often practical to approximate the JPDF of the block variables with the corresponding marginal distributions, i.e., assuming the block variables are uncorrelated and statistically independent. If the block variables are also identically distributed, then the probability function becomes

$$\Phi_{iid}(t) = [\int_{-t}^{t} \psi_x(x)\, dx]^N. \quad (15)$$

When the quantized signal $x$ is white noise, this assumption is quite reasonable. For example, if the signal is Gaussian distributed with variance $\sigma_x^2$, then the distribution of block exponents becomes

$$p_\gamma(\gamma_k) = [erf(\frac{2^{\gamma_k}}{\sqrt{2}\sigma_x})]^N - [erf(\frac{2^{\gamma_k - 1}}{\sqrt{2}\sigma_x})]^N \quad (16)$$

where $erf(x)$ is the *error function*, i.e.,

$$erf(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2}\, dt. \quad (17)$$

The theoretical distribution function as defined by eq. (16) and also the simulated distribution function of the block exponents by quantizing a 0 dB Gaussian signal to BFP format with 2000 samples and taking the block length equal to 4, are plotted in Figures 4 and 5. It is seen that the plots exhibit a very good similarity between them.

### D. Roundoff Error Variance Calculation

As mentioned earlier, if rounding-to-nearest is used as the rounding method, the roundoff error $\alpha$ in eq. (7) has zero mean and variance

$$\sigma_\alpha^2 = \sigma_{e_m}^2 . E[2^{2\gamma}] = \frac{2^{-2B_d}}{12} . \sum_{k=1}^{N_\gamma} p_\gamma(\gamma_k) 2^{2\gamma_k} \qquad (18)$$

where $p_\gamma(\gamma_k)$, $k = 1, ..., N_\gamma$, is the PMF of the block exponents and $N_\gamma$ is the available distinct block exponent levels. A quantitative measurement of $p_\gamma(\gamma_k)$ has already been given in the previous subsection. With the help of eq. (18) above, the signal to noise ratio (SNR) in BFP quantization of a zero mean signal with variance $\sigma_x^2$ can now be expressed as

$$
\begin{aligned}
SNR_{BFP} &= 10log_{10} \frac{\sigma_x^2}{\sigma_\alpha^2} \qquad (19)\\
&= 6.02B_d + 10.79 + 10log_{10}\left(\frac{\sigma_x^2}{\sum_k p_\gamma(\gamma_k)2^{2\gamma_k}}\right)
\end{aligned}
$$

which, apart from the number of bits used to represent each mantissa, i.e., $B_d$, also depends on the ratio of signal variance and the mean-square of the block exponents.
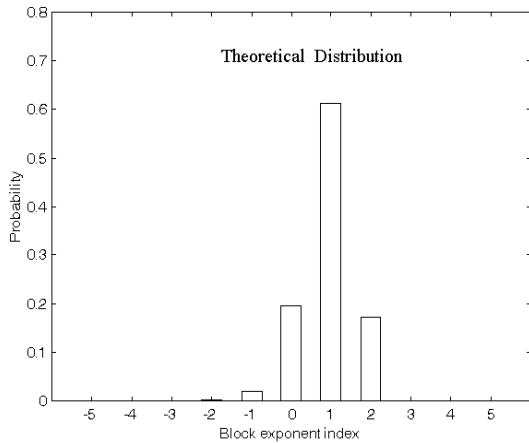


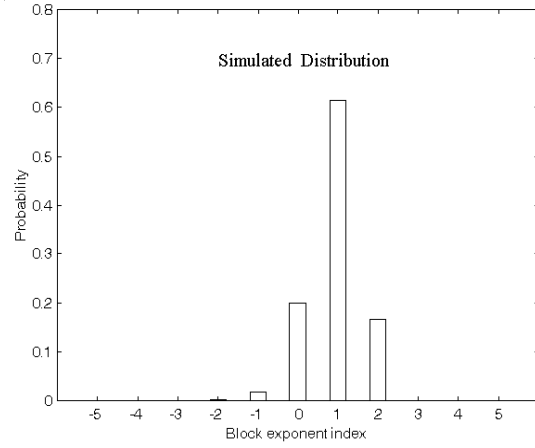Fig. 4 Theoretical block exponent distribution for quantizing 0 dB Gaussian data in BFP format.



Fig. 5 Block exponent distribution from simulation for quantizing 0 dB Gaussian data in BFP format.

### E. Proposed Approximation of SNR Calculation

It should be noted that the theoretical error investigation requires precise knowledge of the distribution of quantized data for calculating block exponent distribution. This information, however, is not always available, which makes the theoretical error calculations only a little useful. Therefore the derivation of an approximate, easy-to-handle formula for calculating SNR in BFP quantization is needed for the analytical purpose.

The approximate analysis comes up from the fundamental idea of quantization. When the number of bits used in each mantissa with FxP, FP and BFP format are same, the SNR in BFP quantization is not as good as that of FP but better than the FxP case. This follows from the simple fact that the BFP quantized mantissas have less number of significant bits compared to FP mantissas, which are always normalized. Again, BFP mantissas are block normalized, i.e., at least the maximum amplitude mantissa is normalized, which gives them the provision to have more significant bits than the FxP mantissas (except the case of 0 dB signal level, which is discussed later on). Hence, the $SNR_{BFP}$ can be interpolated from the FP and FxP SNR formulas with the introduction of a *signal factor* $\tau \in \mathbb{R}$, as

$$SNR_{BFP} = \tau SNR_{FxP} + (1 - \tau)SNR_{FP} \qquad (20)$$

where the signal factor $\tau$, in our treatment, is distributed within the range $[0, 1]$ and is a function of the block length N and signal statistics. To derive the $SNR_{BFP}$, we now need to know only about $SNR_{FP}$ and $SNR_{FxP}$.

The relative error $\epsilon$ in FP arithmetic is assumed to be white noise, uncorrelated with the signal. The total quantization error $e = \epsilon x$ is then white noise with variance

$$\sigma_e^2 = \sigma_\epsilon^2 \sigma_x^2 \qquad (21)$$

with $\sigma_x^2$ being the variance of the signal and

$$\sigma_\epsilon^2 = \frac{\triangle^2}{8ln2} \approx (.18)\triangle^2 \qquad (22)$$

where $\triangle = 2^{-B_d}$. Thus the SNR for a FP quantizer comes as

$$\begin{aligned} SNR_{FP} &= 10log_{10}\frac{\sigma_x^2}{\sigma_e^2} \\ &= 6.02B_d + 7.44 \quad dB. \end{aligned} \qquad (23)$$

For FxP format, the variance of additive roundoff error $\beta$ becomes $\sigma_\beta^2 = \frac{\triangle^2}{12}$, leading to

$$\begin{aligned} SNR_{FxP} &= 10log_{10}\frac{\sigma_x^2}{\sigma_\beta^2} \\ &= 6.02B_d + 10.79 + 10log_{10}\sigma_x^2 \quad dB. \end{aligned} \qquad (24)$$

For a 0 dB Gaussian signal, $SNR_{FxP}$ is even higher than the $SNR_{FP}$ as the FxP mantissas cover the entire dynamic range to have more number of significant bits. Simulations have also supported these theoretical results regarding the SNRs, which have been shown in Fig. 6. Eq. (24) now can be modified by dividing the input signal by a scaling factor $\chi$ to reduce the input amplitude so that the signal, for any level, does not exceed the dynamic range of the FxP quantization process. Thus the modified SNR becomes

$$SNR_{FxP} = 6.02B_d + 10.79 - 20log_{10}\Lambda \quad dB \qquad (25)$$

where $\Lambda = \frac{\chi}{\sigma_x}$. Substituting eq. (23) and (25) in eq. (20), we get

$$SNR_{BFP} = 6.02B_d + 7.44 - \tau(20log_{10}\Lambda - 3.35) \quad dB. \qquad (26)$$

Deriving the exact value of signal factor $\tau$ is rather difficult for the nonlinear effects of block length $N$ and quantizer range limit for different types of signals. Nevertheless, for our purpose, the obtained signal factor with $\Lambda = 1$ was

$$\tau = \frac{(1 - N^{-1})}{(1 + N^{-1})} \qquad (27)$$

which has shown exact values for the terminal cases but failed to interpolate good results on few points in between the entire range. A better formula for $\tau$ was obtained through fitting this curve to the data accumulated by theoretical calculations as

$$\tau = (1 - N^{-.72})(1 + N^{-.72})^{-1} \qquad (28)$$

with the values $\tau = 0$ for $N = 1$ and $\tau = 1$ for $N \to \infty$, as expected. The signal factor is plotted as a function of $N$ in Fig. 7. Note that the above formula can be applied only when the quantized signal occupies the whole amplitude range of the quantizer, i.e., for a 0 dB signal.
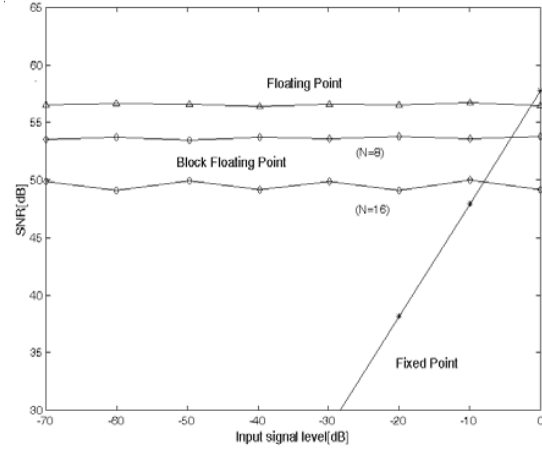


Fig. 6 SNR diagram for quantizing the uncorrelated Gaussian data with 1+7 bit FxP, 1+7+(1+3) bit FP and 1+7+(1+3)/N bit BFP format, with N=8 and N=16.
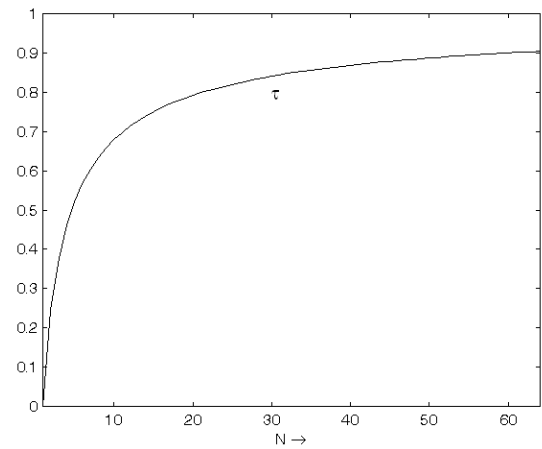


Fig. 7 The signal factor $\tau$ plotted as a function of block length N for a 0 dB Gaussian signal.

### F. A Note on BFP Format Accuracy

The accuracy given by eq. (28) is quite acceptable for small and moderate block lengths. For larger block lengths, the SNR shows some ripple due to dependence on signal. The approximate SNR formula in eq. (26) is most reliable within a certain dynamic range where the quantized signal utilizes the whole amplitude range. Beyond this range, the $SNR_{BFP}$ actually falls off rapidly and can be increased again by adding more number of block exponent bits.

In general, the BFP system gives the best performance for small block lengths. Increasing the block length $N$ reduces the computational complexity as well as the effective number of bits per sample to a little extent on one hand, but also decreases the SNR on the other. Hence, one must be choosy about finding out the lowest possible $N$ for fulfilling the operational

requirements to work with this format properly. For example, if the length of a direct-form FIR digital filter is $L$, then it is sufficient to choose the block length $N \geq L - 1$, in order to minimize the inter-block adjustments.

## IV. CONCLUSIONS

The finite wordlength properties of BFP representation system have been studied in this paper. BFP format is a better choice as an alternative of FP systems due to its simplicity of the associated hardware with sufficient SNR over a wide dynamic range. The theoretical calculations have been done with the help of a significant quantization model for this format and a probability distribution function for the block exponents has also been derived step-by-step for this purpose. The theoretical SNR of BFP system is found to be dependent on the distribution function of the block exponents, which, in turn, requires a precise knowledge on signal properties to be evalueted. Thus an easy usable formula for calculating the BFP format SNR has been derived with the introduction of a signal factor, which has been proved to be quite acceptable for a good dynamic range. Such an excellent format with high potential is expected to be investigated by researchers in several challenging signal processing application areas in wireless communications, where numerical error properties of chosen data format play an important role.

## REFERENCES

[1] K. R. Ralev and P. H. Bauer, "Realization of Block Floating Point Digital Filters and Application to Block Implementations," *IEEE Trans. Signal Processing*, vol. 47, no. 4, pp. 1076-1086, April 1999.

[2] K. Kalliojärvi and J. Astola, "Roundoff Errors in Block-Floating-Point Systems," *IEEE Trans. Signal Processing*, vol. 44, no. 4, pp. 783-790, April 1996.

[3] J. Kontro, K. Kalliojärvi and Y. Neuvo, "Floating-point arithmetic in signal processing," in *Proc. 1992 IEEE Int. Symp. Circuits, Syst.*, San Diego, CA, May 10-13, 1992, pp. 1784-1791.

[4] S. Sridharan and G. Dickman, "Block floating point implementation of digital filters using the DSP56000," *Microprocess. Microsyst.*, vol. 12, no. 6, pp. 299-308, July-Aug. 1988.

[5] P. H. Bauer, "Absolute Error Bounds for Block-Floating-Point Direct-Form Digital Filters," *IEEE Trans. Signal Processing*, vol. 43, no. 8, pp. 1994-1996, Aug. 1995.

[6] A. V. Oppenheim, "Realization of digital filters using block floating point arithmetic," *IEEE Trans. Audio Electroaccoust.*, vol. AE-18, no. 2, pp. 130-136, June 1970.

[7] K. Kalliojärvi, "Analysis of Block-Floating-Point Quantization Error," in *Proc. 11th Euro. Conf. Circuit Theo., Design*, Davos, Switzerland, Aug. 30- Sep. 3, 1993, pp. 791-796.

[8] A. Mitra, "A New Block-based NLMS Algorithm and Its Realization in Block Floating Point Format," *Int. J. Info. Tech.*, vol. 1, no. 4, pp. 244-248, 2004.

[9] A. Mitra, "Efficient Realization of Gradient Based Adaptive Filters using Block Floating Point Arithmetic," Ph.D. Dissertation, Indian Institute of Technology Kharagpur, India, Jan. 2004.

[10] A. V. Oppenheim and R. W. Schafer, *Discrete-Time Signal Processing*, Englewood Cliffs, NJ: Prentice-Hall, 1989.

[11] A. Fettweis, "On Properties of Floating-Point Roundoff Noise," *IEEE Trans. Accoust. Speech Signal Processing*, pp. 149-151, April 1974.

[12] A. Papoulis, *Probability, Random Variables and Stochastic Processes*, New York: McGraw-Hill, 1965.

[13] T. Kaneko and B. Liu, "On local roundoff errors in floating-point arithmetic," *Journal Ass. Comp. Mach.*, vol. 20, pp. 391-398, July 1973.

[14] A. V. Oppenheim and R. W. Schafer, *Digital Signal Processing*, Englewood Cliffs, NJ: Prentice-Hall, 1975.

[15] B. Liu, "Effect of finite wordlength on the accuracy of digital filters-A review," *IEEE Trans. Circuit Theory*, vol. CT-18, pp. 670-677, Nov 1971.

[16] J. H. Wilkinson, *Rounding Errors in Algebraic Processes*, Englewood Cliffs, NJ: Prentice-Hall, 1963.

[17] A. B. Sripad and D. L. Snyder, "Quantization Errors in Floating-Point Arithmetic," *IEEE Trans. Accoust. Speech Signal Processing*, vol. ASSP-26,pp. 149-151, Oct 1978.

**Abhijit Mitra** was born in Serampore, India, in 1975. He received the B.E.(Honors) degree from R. E. College, Durgapur, India, in 1997, M.E.Tel.E. degree from Jadavpur University, India, in 1999 and Ph.D. degree from Indian Institute of Technology, Kharagpur, India, in 2004, all in electronics and communication engineering. Since 2004, he is with Indian Institute of Technology, Guwahati, India, as an Assistant Professor. His research interests include finite wordlength digital signal processing, statistical signal processing, adaptive signal processing and wireless communications.

Dr. Mitra is a member of IEEE and also serves as a reviewer of IEEE Transactions on Signal Processing. He is the recipient of several national scholarships and research fellowships. He also has edited a proceedings volume of all India students' paper contest.