

## Fixed-Point Arithmetic: An Introduction

*Randy Yates*

September 15, 2020



<http://www.digitalsignallabs.com>

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Fixed-Point Binary Representations</b>	<b>3</b>
2.1	Unsigned Fixed-Point Rationals . . . . .	4
<b>3</b>	<b>The Operations of One's Complement and Two's Complement</b>	<b>5</b>
<b>4</b>	<b>Signed Two's Complement Fixed-Point Rationals</b>	<b>5</b>
<b>5</b>	<b>Fundamental Rules of Fixed-Point Arithmetic</b>	<b>6</b>
5.1	Unsigned Wordlength . . . . .	6
5.2	Signed Wordlength . . . . .	6
5.3	Unsigned Range . . . . .	6
5.4	Signed Range . . . . .	6
5.5	Addition Operands . . . . .	6
5.6	Addition Result . . . . .	6
5.7	Unsigned Multiplication . . . . .	7
5.8	Signed Multiplication . . . . .	7
5.9	Unsigned Division . . . . .	7
5.10	Signed Division . . . . .	8
5.11	Wordlength Reduction . . . . .	8
5.12	Shifting . . . . .	8
5.12.1	Literal Shift . . . . .	9
5.12.1.1	Multiplying/Dividing By A Power of Two . . . . .	9
5.12.1.2	Modifying Scaling . . . . .	9
5.12.2	Virtual Shift . . . . .	10
<b>6</b>	<b>Dimensional Analysis in Fixed-Point Arithmetic</b>	<b>10</b>
<b>7</b>	<b>Concepts of Finite Precision Math</b>	<b>12</b>
7.1	Precision . . . . .	12
7.2	Resolution . . . . .	12
7.3	Range . . . . .	12
7.4	Accuracy . . . . .	12
7.5	Dynamic Range . . . . .	12
<b>8</b>	<b>Fixed-Point Analysis—An Example</b>	<b>13</b>
<b>9</b>	<b>Acknowledgments</b>	<b>14</b>
<b>10</b>	<b>Terms and Abbreviations</b>	<b>14</b>
<b>11</b>	<b>Revision History</b>	<b>14</b>
<b>12</b>	<b>References</b>	<b>15</b>

## List of Figures

## List of Tables

1	Revision History . . . . .	14
---	----------------------------	----

## 1 Introduction

This document presents definitions of signed and unsigned fixed-point binary number representations and develops basic rules and guidelines for the manipulation of these number representations using the common arithmetic and logical operations found in fixed-point DSPs and hardware components.

While there is nothing particularly difficult about this subject, I found little documentation either in hardcopy or on the web. What documentation I did find was disjointed, never putting together *all* of the aspects of fixed-point arithmetic that I think are important. I therefore decided to develop this material and to place it on the web not only for my own reference but for the benefit of others who, like myself, find themselves needing a complete understanding of the issues in implementing fixed-point algorithms on platforms utilizing integer arithmetic.

During the writing of this paper, I was developing assembly language code for the Texas Instruments TMS320C50 Digital Signal Processor, thus my approach to the subject is undoubtedly biased towards this processor in terms of the operation of the fundamental arithmetic operations. For example, the C50 performs adds and multiplies as if the numbers are simple signed two's complement integers. Contrast this against the Motorola 56k series which performs two's complement fractional arithmetic, with values always in the range  $-1 \leq x < +1$ .

It is my hope that this material is clear, accurate, and helpful. If you find any errors or inconsistencies, please email me at [yates@ieee.org](mailto:yates@ieee.org).

Finally, the reader may be interested in the author's related paper [1] on the application of fixed-point arithmetic to the implementation of FIR filters.

## 2 Fixed-Point Binary Representations

A collection of  $N$  ( $N$  a positive integer) binary digits (bits) has  $2^N$  possible states. This can be seen from elementary counting theory, which tells us that there are two possibilities for the first bit, two possibilities for the next bit, and so on until the last bit, resulting in

$$\underbrace{2 \times 2 \times \dots \times 2}_{N \text{ times}} = 2^N$$

possibilities.

In the most general sense, we can allow these states to represent anything conceivable. In the case of an  $N$ -bit binary word, some examples are up to  $2^N$ :

1. students at a university;
2. species of plants;
3. atomic elements;
4. integers;
5. voltage levels.

Drawing from set theory and elementary abstract algebra, one could view a representation as an onto mapping between the binary states and the elements in the representation set (in the case of unassigned binary states, we assume there is an "unassigned" element in the representation set to which all such states are mapped).

The salient point is that there is no meaning inherent in a binary word, although most people are tempted to think of them (at first glance, anyway) as positive integers (i.e., the *natural binary* representation, defined in the next section). However, **the meaning of an N-bit binary word depends entirely on its interpretation**, i.e., on the representation set and the mapping we choose to use.

In this section, we consider representations in which the representation set is a particular subset of the rational numbers. Recall that the rational numbers are the set of numbers expressible as  $a/b$ , where  $a, b \in \mathbb{Z}, b \neq 0$ . ( $\mathbb{Z}$  is the set of integers.) The subset to which we refer are those rationals for which  $b = 2^n$ . We also further constrain the representation sets to be those in which every element in the set has the same number of binary digits and in which every element in the set has the binary point at the same position, i.e., the binary point is fixed. Thus these representations are called “fixed-point.”

The following sections explain four common binary representations: unsigned integers, unsigned fixed-point rationals, signed two’s complement integers, and signed two’s complement fixed-point rationals. We view the integer representations as special cases of the fixed-point rational representations, therefore we begin by defining the fixed-point rational representations and then subsequently show how these can simplify to the integer representations. We begin with the unsigned representations since they require nothing more than basic algebra. Section 2.2 defines the notion of a “two’s complement” so that we may proceed well-grounded to the discussion of signed two’s complement rationals in section 2.3.

## 2.1 Unsigned Fixed-Point Rationals

An N-bit binary word, when interpreted as an unsigned fixed-point rational, can take on values from a subset  $P$  of the non-negative rationals given by

$$P = \{p/2^b \mid 0 \leq p \leq 2^N - 1, p \in \mathbb{Z}\}.$$

Note that  $P$  contains  $2^N$  elements. We denote such a representation  $U(a, b)$ , where  $a = N - b$ .

In the  $U(a, b)$  representation, the  $n$ th bit, counting from right to left and beginning at 0, has a weight of  $2^n/2^b = 2^{n-b}$ . Note that when  $n = b$  the weight is exactly 1. Similar to normal everyday base-10 decimal notation, the binary point is between this bit and the bit to the right. This is sometimes referred to as the **implied binary point**. A  $U(a, b)$  representation has  $a$  integer bits and  $b$  fractional bits.

The value of a particular N-bit binary number  $x$  in a  $U(a, b)$  representation is given by the expression

$$x = (1/2^b) \sum_{n=0}^{N-1} 2^n x_n$$

where  $x_n$  represents bit  $n$  of  $x$ . The range of a  $U(a, b)$  representation is from 0 to  $(2^N - 1)/2^b = 2^a - 2^{-b}$ .

For example, the 8-bit unsigned fixed-point rational representation  $U(6, 2)$  has the form

$$b_5 b_4 b_3 b_2 b_1 b_0 . b_{-1} b_{-2},$$

where bit  $b_k$  has a weight of  $2^k$ . Note that since  $b = 2$  the binary point is to the left of the second bit from the right (counting from zero), and thus the number has six integer bits and two fractional bits. This representation has a range of from 0 to  $2^6 - 2^{-2} = 64 - 1/4 = 63\ 3/4$ .

The unsigned integer representation can be viewed as a special case of the unsigned fixed-point rational representation where  $b = 0$ . Specifically, an N-bit unsigned integer is identical to a  $U(N, 0)$  unsigned fixed-point rational. Thus the range of an N-bit unsigned integer is

$$0 \leq U(N, 0) \leq 2^N - 1.$$

and it has  $N$  integer bits and 0 fractional bits. The unsigned integer representation is sometimes referred to as “natural binary.”

Examples:

1.  $U(6, 2)$ . This number has  $6 + 2 = 8$  bits and the range is from 0 to  $2^6 - 1/2^2 = 63.75$ . The value 8Ah (1000,1010b) is

$$(1/2^2)(2^1 + 2^3 + 2^7) = 34.5.$$

2.  $U(-2, 18)$ . This number has  $-2 + 18 = 16$  bits and the range is from 0 to  $2^{-2} - 1/2^{18} = 0.2499961853027$ . The value 04BCh (0000,0100,1011,1100b) is

$$(1/2^{18})(2^2 + 2^3 + 2^4 + 2^5 + 2^7 + 2^{10}) = 1212/2^{18} = 0.004623413085938.$$

3.  $U(16, 0)$ . This number has  $16 + 0 = 16$  bits and the range is from 0 to  $2^{16} - 1 = 65,535$ . The value 04BCh (0000,0100,1011,1100b) is

$$(1/2^0)(2^2 + 2^3 + 2^4 + 2^5 + 2^7 + 2^{10}) = 1212/2^0 = 1212.$$

### 3 The Operations of One’s Complement and Two’s Complement

Consider an  $N$ -bit binary word  $x$  interpreted as if in the  $N$ -bit natural binary representation (i.e.,  $U(N, 0)$ ). The *one’s complement* of  $x$  is defined to be an operation that inverts every bit of the original value  $x$ . This can be performed arithmetically in the  $U(N, 0)$  representation by subtracting  $x$  from  $2^N - 1$ . That is, if we denote the one’s complement of  $x$  as  $\tilde{x}$ , then

$$\tilde{x} = 2^N - 1 - x.$$

The two’s complement of  $x$ , denoted  $\hat{x}$ , is determined by taking the one’s complement of  $x$  and then adding one to it:

$$\begin{aligned} \hat{x} &= \tilde{x} + 1 \\ &= 2^N - x. \end{aligned} \tag{1}$$

Examples:

1. The one’s complement of the  $U(8, 0)$  number 03h (0000,0011b) is FCh (1111,1100b).
2. The two’s complement of the  $U(8, 0)$  number 03h (0000,0011b) is FDh (1111,1101b).

### 4 Signed Two’s Complement Fixed-Point Rationals

An  $N$ -bit binary word, when interpreted as a signed two’s complement fixed-point rational, can take on values from a subset  $P$  of the rationals given by

$$P = \{p/2^b \mid -2^{N-1} \leq p \leq 2^{N-1} - 1, p \in \mathcal{Z}\}.$$

Note that  $P$  contains  $2^N$  elements. We denote such a representation  $A(a, b)$ , where  $a = N - b - 1$ .

The value of a specific  $N$ -bit binary number  $x$  in an  $A(a,b)$  representation is given by the expression

$$x = (1/2^b) \left[ -2^{N-1} x_{N-1} + \sum_{n=0}^{N-2} 2^n x_n \right],$$

where  $x_n$  represents bit  $n$  of  $x$ . The range of an  $A(a,b)$  representation is

$$-2^{N-1-b} \leq x \leq +2^{N-1-b} - 1/2^b.$$

Note that the number of bits in the magnitude term of the sum above (the summation, that is) has one less bit than the equivalent prior unsigned fixed-point rational representation. Further note that these bits are the  $N - 1$  least significant bits. It is for these reasons that the most-significant bit in a signed two's complement number is usually referred to as the *sign bit*.

Example:

$A(13,2)$ . This number has  $13+2+1=16$  bits and the range is from  $-2^{13} = -8192$  to  $+2^{13} - 1/4 = 8191.75$ .

## 5 Fundamental Rules of Fixed-Point Arithmetic

The following are practical rules of fixed-point arithmetic. For these rules we note that when a scaling can be either signed ( $A(a,b)$ ) or unsigned ( $U(a,b)$ ), we use the notation  $X(a,b)$ .

### 5.1 Unsigned Wordlength

The number of bits required to represent  $U(a,b)$  is  $a + b$ .

### 5.2 Signed Wordlength

The number of bits required to represent  $A(a,b)$  is  $a + b + 1$ .

### 5.3 Unsigned Range

The range of  $U(a,b)$  is  $0 \leq x \leq 2^a - 2^{-b}$ .

### 5.4 Signed Range

The range of  $A(a,b)$  is  $-2^a \leq x \leq 2^a - 2^{-b}$ .

### 5.5 Addition Operands

Two binary numbers must be scaled the same in order to be added. That is,  $X(c,d) + Y(e,f)$  is only valid if  $X = Y$  (either both  $A$  or both  $U$ ) and  $c = e$  and  $d = f$ .

### 5.6 Addition Result

The scale of the sum of two binary numbers scaled  $X(e,f)$  is  $X(e+1,f)$ , i.e., the sum of two  $M$ -bit numbers requires  $M + 1$  bits.

## 5.7 Unsigned Multiplication

$$U(a_1, b_1) \times U(a_2, b_2) = U(a_1 + a_2, b_1 + b_2).$$

## 5.8 Signed Multiplication

$$A(a_1, b_1) \times A(a_2, b_2) = A(a_1 + a_2 + 1, b_1 + b_2).$$

## 5.9 Unsigned Division

Let  $U(a_3, b_3) = \frac{U(a_1, b_1)}{U(a_2, b_2)}$  and consider the largest possible result:

$$\begin{aligned} \text{largest result} &= \frac{\text{largest dividend}}{\text{smallest divisor}} \\ &= \frac{2^{a_1} - 2^{-b_1}}{2^{-b_2}} \\ &= 2^{a_1+b_2} - 2^{b_2-b_1}. \end{aligned} \quad (2)$$

Thus we require

$$2^{a_3} - 2^{-b_3} \geq 2^{a_1+b_2} - 2^{b_2-b_1}. \quad (3)$$

It is natural to let  $a_3 = a_1 + b_2$ , in which case the inequalities below result:

$$\begin{aligned} 2^{a_3} - 2^{-b_3} &\geq 2^{a_3} - 2^{b_2-b_1} \\ -2^{-b_3} &\geq -2^{b_2-b_1} \\ 2^{-b_3} &\leq 2^{b_2-b_1} \\ -b_3 &\leq b_2 - b_1 \\ b_3 &\geq b_1 - b_2. \end{aligned} \quad (4)$$

Thus we have a constraint on  $b_3$  due to  $b_1$  and  $b_2$ .

Now consider the smallest possible result:

$$\begin{aligned} \text{smallest result} &= \frac{\text{smallest dividend}}{\text{largest divisor}} \\ &= \frac{2^{-b_1}}{2^{a_2} - 2^{-b_2}}. \end{aligned} \quad (5)$$

This then requires  $b_3$  to obey the following constraint:

$$\begin{aligned} 2^{-b_3} &\leq \frac{2^{-b_1}}{2^{a_2} - 2^{-b_2}} \\ b_3 &\geq b_1 + \log_2(2^{a_2} - 2^{-b_2}) \end{aligned} \quad (6)$$

If we assume  $b_2$  is positive, (6) is the more stringent of the two constraints (4) and (6) on  $b_3$ . We then express (6) in a slightly simpler form:

$$b_3 \geq \log_2(2^{a_2+b_1} - 2^{b_1-b_2}). \quad (7)$$

The final result is then

$$U(a_1, b_1)/U(a_2, b_2) = U(a_1 + b_2, \lceil \log_2(2^{a_2+b_1} - 2^{b_1-b_2}) \rceil). \quad (8)$$

## 5.10 Signed Division

Let  $r = n/d$  where  $n$  is scaled  $A(a_n, b_n)$  and  $d$  is scaled  $A(a_d, b_d)$ . What is the scaling of  $r$  ( $A(a_r, b_r)$ )?

$$|r_M| = \frac{|n_M|}{|d_m|} \quad (9)$$

$$= \frac{2^{a_n}}{2^{-b_d}} \quad (10)$$

$$= 2^{a_n+b_d}. \quad (11)$$

Since this maximum can be positive (when the numerator and denominator are both negative),  $a_r = a_n + b_d + 1$ .

Similarly,

$$|r_m| = \frac{|n_m|}{|d_M|} \quad (12)$$

$$= \frac{2^{-b_n}}{2^{a_d}} \quad (13)$$

$$= 2^{-(a_d+b_n)}. \quad (14)$$

This implies  $b_r = a_d + b_n$ .

Thus

$$\frac{A(a_n, b_n)}{A(a_d, b_d)} = A(a_n + b_d + 1, a_d + b_n). \quad (15)$$

## 5.11 Wordlength Reduction

Define the operation  $HIn(X(a, b))$  to be the extraction of the  $n$  most-significant bits of  $X(a, b)$ . Similarly, define the operation  $LOn(X(a, b))$  to be the extraction of the  $n$  least-significant bits of  $X(a, b)$ . For signed values,

$$\begin{aligned} HIn(A(a, b)) &= A(a, n - a - 1) \text{ and} \\ LOn(A(a, b)) &= A(n - b - 1, b). \end{aligned} \quad (16)$$

Similarly, for unsigned values,

$$\begin{aligned} HIn(U(a, b)) &= U(a, n - a) \text{ and} \\ LOn(U(a, b)) &= U(n - b, b). \end{aligned} \quad (17)$$

## 5.12 Shifting

We define two types of shift operations below, *literal* and *virtual*, and describe the scaling results of each.

Note that shifts are expressed in terms of right shifts by integer  $n$ . Shifting left is accomplished when  $n$  is negative.



### 5.12.1 Literal Shift

A *literal shift* occurs when the bit positions in a register move left or right. A literal shift can be performed for two possible reasons, to divide or multiply by a power of two, or to change the scaling.

In both cases note that this will possibly result in a loss of precision or overflow assuming the output register width is the same as the input register width.

**5.12.1.1 Multiplying/Dividing By A Power of Two** A literal shift that is done to multiply or divide by a power of two shifts the bit positions but keeps the output scaling the same as the input scaling.

Thus we have the following scaling:

$$X(a, b) \gg n = X(a, b). \quad (18)$$

Example:

Let's say  $X$  is a 16-bit signed two's complement integer that is scaled A(14,1), or Q1. Let's set that integer equal to 128:  $X = +128 = 0x0080$ , and thus its scaled value is

$$x = X/2^1 \quad (19)$$

$$= 128/2 \quad (20)$$

$$= 64.0 \quad (21)$$

Now we want to divide that by 4, so we shift it right by 2,  $X = X \gg 2$ , so that the new value of  $X$  is 32. This shift didn't change the scaling since we are dividing by actually shifting, so it's still scaled Q1 after the shift. So now  $X = 32$  and  $x = X/2 = 16.0$ . Since the original value of  $x$  was 64.0, we see that we have indeed divided that value by 4, which was the objective.

Note that this is probably a bad way to multiply or divide a fixed-point value since, if you're multiplying, you run the risk of overflowing, and if you're dividing, you run the risk of losing precision. It would be much better to perform the multiplication or division using the "virtual shift" method described in section 5.12.2.

**5.12.1.2 Modifying Scaling** A literal shift that is done to modify the scaling shifts the bit positions and makes the output scaling different than the input scaling. Thus we have the following scaling:

$$X(a, b) \gg n = X(a + n, b - n). \quad (22)$$

Example:

Again let's say  $X$  is a 16-bit signed two's complement integer that is scaled A(14,1), or Q1, and let's set that integer equal to 128:  $X = +128 = 0x0080$ , and thus its scaled value is

$$x = X/2^1 \quad (23)$$

$$= 128/2 \quad (24)$$

$$= 64.0 \quad (25)$$

Now let's say we want to change the scaling from Q1 to Q3 (or equivalently, from A(14,1) to A(12,3)). So we shift the integer left by two bits:  $X = X \ll 2$ . So our new integer value is 512, but we've now also changed our scaling as in equation 22 so that it's A(12,3) ( $n = -2$  here since we're shifting left).

So in this case our final fixed-point value is still  $512 / 8 = 64.0$ .

### 5.12.2 Virtual Shift

A *virtual shift* shifts the virtual binary point<sup>1</sup> without modifying the underlying integer value. It can be used as an alternate method of performing a multiplication or division by a power of two. However, unlike the literal shift case, the virtual shift method loses no precision and avoids overflow. This is because the bit positions don't actually move—the operation is simply a reinterpretation of the scaling.

$$X(a, b) \gg n = X(a - n, b + n). \quad (26)$$

## 6 Dimensional Analysis in Fixed-Point Arithmetic

Consider a fixed-point variable  $x$  that is scaled  $A(a_x, b_x)$ . Denote the scaled value of the variable as lowercase  $x$  and the unscaled value as uppercase  $X$  so that

$$x = X/2^{b_x}.$$

Units, such as inches, seconds, furlongs/fortnight, etc., may be associated with a fixed-point variable by assigning a weight to the variable. Denote the scaled weight as  $w$  and the unscaled weight as  $W$ , so that the value *and dimension* of a quantity  $\alpha_x$  that is to be represented by  $x$  can be expressed as

$$\begin{aligned} \alpha_x &= x \times w \\ &= X \times W. \end{aligned} \quad (27)$$

Since  $x = X/2^{b_x}$ ,

$$x \times w = X/2^{b_x} \times w, \quad (28)$$

and since  $x \times w = X \times W$ ,

$$X/2^{b_x} \times w = X \times W \Rightarrow w = 2^{b_x} W. \quad (29)$$

Example 1:

An inertial sensor provides a linear acceleration signal to a 16-bit signed two's complement A/D converter with a reference voltage of 2V peak. The analog sensor signal is related to acceleration in meters per second squared through the conversion  $m/(s^2 - volt)$ , i.e., the actual acceleration  $\alpha(t)$  in meters/second<sup>2</sup> can be determined from the sensor voltage  $v(t)$  as

$$\alpha(t) = v(t) [\text{volts}] \times \left[ \frac{\text{m}}{(s^2)(\text{volt})} \right]. \quad (30)$$

<sup>1</sup>The virtual binary point is called virtual since it doesn't actually exist anywhere except in the programmer's mind.

If we consider the incoming A/D samples to be scaled  $A(16, -1)$ , what is the corresponding scaled and unscaled weights?

Solution:

We know that -32768 corresponds to -2 volts. Using the equation

$$\begin{aligned}\alpha_x &= x \times w \\ &= X/2^{b_x} \times w\end{aligned}\quad (31)$$

we simply plug in the values to obtain

$$-2[\text{volts}] = -32768/2^{-1} \times w_v. \quad (32)$$

Solving for  $w_v$  provides the intermediate result

$$w_v = \left[ \frac{\text{volt}}{32768} \right]. \quad (33)$$

Let's check: An unscaled value of 32767 corresponds to a scaled value of  $v = 32767/2^{-1} = 65534$ , and thus the physical quantity  $\alpha_v$  to which this corresponds is

$$\begin{aligned}\alpha_v &= v \times w_v \\ &= 65534 \times \left[ \frac{\text{volt}}{32768} \right] \\ &= 1.999939[\text{volts}].\end{aligned}\quad (34)$$

Now simply multiply  $w_v$  by the original analog conversion factor  $\left[ \frac{\text{m}}{(\text{s}^2)(\text{volt})} \right]$  to obtain the acceleration weighting  $w_a$  directly:

$$\begin{aligned}w_a &= w_v \times \left[ \frac{\text{m}}{(\text{s}^2)(\text{volt})} \right] \\ &= \left[ \frac{\text{volt}}{32768} \right] \times \left[ \frac{\text{m}}{(\text{s}^2)(\text{volt})} \right] \\ &= \left[ \frac{\text{m}}{32768\text{s}^2} \right].\end{aligned}\quad (35)$$

The unscaled weight is then determined from the scaled weight and the scaling as

$$\begin{aligned}W_a &= w_a/2^{b_a} \\ &= \left[ \frac{\text{m}}{32768\text{s}^2} \right] / 2^{-1} \\ &= \left[ \frac{\text{m}}{16384\text{s}^2} \right].\end{aligned}\quad (36)$$

Example 2:

Bias is an important error in inertial measurement systems. An average scaled value of 29 was measured from the inertial measurement system in example 1 with the system at rest. What is the bias  $\beta$ ?

Solution:

$$\begin{aligned}
 \beta &= x \times w \\
 &= 29 \times \left[ \frac{\text{m}}{32768\text{s}^2} \right] \\
 &= 0.88501 \times 10^{-3} \left[ \frac{\text{m}}{\text{s}^2} \right].
 \end{aligned} \tag{37}$$


---

## 7 Concepts of Finite Precision Math

### 7.1 Precision

*Precision* is the maximum number of non-zero bits representable. For example, an A(13,2) number has a precision of 16 bits. For fixed-point representations, precision is equal to the wordlength.

### 7.2 Resolution

*Resolution* is the smallest non-zero magnitude representable. For example, an A(13,2) has a resolution of  $1/2^2 = 0.25$ .

### 7.3 Range

*Range* is the difference between the most negative number representable and the most positive number representable,

$$X_R = X_{MAX+} - X_{MAX-}. \tag{38}$$

For example, an A(13,2) number has a range from -8192 to +8191.75, i.e., 16383.75.

### 7.4 Accuracy

*Accuracy* is the magnitude of the maximum difference between a real value and its representation. For example, the accuracy of an A(13,2) number is 1/8. Note that accuracy and resolution are related as follows:

$$A(x) = R(x)/2, \tag{39}$$

where  $A(x)$  is the accuracy of  $x$  and  $R(x)$  is the resolution of  $x$ .

### 7.5 Dynamic Range

*Dynamic range* is the ratio of the maximum absolute value representable and the minimum positive (i.e., non-zero) absolute value representable. For a signed fixed-point rational representation  $A(a, b)$ , dynamic range is

$$\times 2^a / 2^{-b} = 2^{a+b} = 2^{N-1}. \tag{40}$$

For an unsigned fixed-point rational representation  $U(a, b)$ , dynamic range is

$$(2^a - 2^{-b}) / 2^{-b} = 2^{a+b} - 1 = 2^N - 1. \tag{41}$$

For  $N$  of any significant size, the “-1” is negligible.

## 8 Fixed-Point Analysis—An Example

An algorithm is usually defined and developed using an algebraically complete number system such as the real or complex numbers. To be more precise, the operations of addition, subtraction, multiplication, and division are performed (for example) over the field  $(\mathbb{R}, +, \times)$ , where subtraction is equivalent to adding the additive inverse and division is equivalent to multiplying by the multiplicative inverse.

As an example, consider the algorithm for calculating the average of the square of a digital signal  $x(n)$  over the interval  $N$  (here, the signal is considered to be quantized in time but not in amplitude):

$$y(n) = \frac{1}{N} \sum_{k=0}^{N-1} x^2(n-k). \quad (42)$$

In this form, the algorithm implicitly assumes  $x(n) \in \mathbb{R}$ , and the operations of addition and multiplication are performed over the field  $(\mathbb{R}, +, \times)$ . In this case, the numerical representations have infinite precision.

This state of affairs is perfectly acceptable when working with pencil and paper or higher-level floating-point computing environments such as Matlab or MathCad. However, when the algorithm is to be implemented in fixed-point hardware or software, it must necessarily utilize a finite number of binary digits to represent  $x(n)$ , the intermediate products and sums, and the output  $y(n)$ .

Thus the basic task of converting such an algorithm into fixed-point arithmetic is that of determining the wordlength, accuracy, and range required for each of the arithmetic operations involved in the algorithm. In the terms of the fundamentals given in section 2, we need to determine a) whether the value should be signed ( $A(a, b)$ ) or unsigned ( $U(a, b)$ ), b) the value of  $N$  (the wordlength), and c) the values for  $a$  and  $b$  (the accuracy and range). Any two of wordlength, accuracy, and range determine the third. For example, given wordlength and accuracy, range is determined. In other words, we cannot independently specify all of wordlength, accuracy, and range.

Continuing with our example, assume the input  $x(n)$  is scaled  $A(15, 0)$ , i.e., plain old 16-bit signed two's complement samples. The first operation to be performed is to compute the square. According to the rules of fixed-point arithmetic,  $A(15, 0) \times A(15, 0) = A(31, 0)$ . In other words, we require 32 bits for the result of the square in order to guarantee that we will avoid overflow and maintain precision. It is at this point that design tradeoffs and other information begin to affect how we implement our algorithm.

For example, in one possible scenario, we may know *a-priori* that the input data  $x(n)$  do not span the full dynamic range of the  $A(15, 0)$  representation, thus it may be possible to reduce the 32-bit requirement for the result and still guarantee that the square operation does not overflow.

Another possible scenario is that we do not require all of the precision in the result, and this also will reduce the required wordlength.

In yet a third scenario, we may look ahead to the summation to be performed and realize that if we don't scale back the result of each square we will overflow the sum that is to subsequently be performed (assuming we have a 32-bit accumulator). On the other hand, we may be using a fixed-point processor such as the TI TMS320C54x which has a 40-bit accumulator, thus we have 8 "guard bits" past the 32-bit result which may be used in the accumulations to prevent overflow for up to 256 ( $8 = \log_2(256)$ ) sums.

To complete our example, let's further assume that a) we keep all 32 bits of the result of the squaring operation, b) the averaging "time,"  $N$ , does not exceed  $2^4 = 16$  samples, c) we are using a fixed-point processor with an accumulator of  $32 + 4 = 36$  bits or greater, and d) the output wordlength for  $y(n)$  is 16 bits ( $A(15, 0)$ ). The final decision that must be made is to determine which method we will use to form a 16-bit value from our 36-bit sum.

It is clear that we should take the 16 bits from bits 20 to 35 of the accumulator (where bit 0 is the LSB) in order to avoid overflowing the output, but shall we truncate or round? Shall we utilize some type of dithering or noise-shaping? These are all questions that relate to the process of *quantization* since we are quantizing a 36-bit word to a 16-bit word. The theory of quantization and the tradeoffs to be made are outside the scope of this topic.

## 9 Acknowledgments

I wish to thank my colleague John Storbeck, a fellow DSP programmer and full-time employee of GEC-Marconi Hazeltine (Wayne, NJ), for his guidance during my early encounters with the C50 and fixed-point programming, and for his stealy critiques of my errant thoughts on assorted programming solutions for the C50. I also wish to thank Dr. Donald Heckathorn for nurturing my knowledge of fixed-point arithmetic during the same time-frame. Finally, I wish to thank my friend and fellow DSP engineering colleague Robert Bristow-Johnson for his encouragement and guidance during my continuing journey through the world of fixed-point programming and digital signal processing.

## 10 Terms and Abbreviations

**DSP** Digital Signal Processor, or Digital Signal Processing.

**FIR** Finite Impulse Response. A type of digital filter that does not require a recursive architecture (i.e., the use of *feedback*), that is inherently stable, and generally easier to design and implement than IIR filters. However, an FIR filter requires more computational resources than an equivalent IIR filter.

**IIR** Infinite Impulse Response. A type of digital filter that requires a recursive architecture, is potentially unstable, and substantially more difficult to design and implement than FIR filters. However, an IIR filter requires less computational resources than an equivalent FIR filter.

## 11 Revision History

Table 1 lists the revision history for this document.

Rev.	Date/Time	Person	Changes
PA1	circa 1997	Randy Yates	Initial Version
PA2	unknown	Randy Yates	1. Converted to Digital Signal Labs document format. 2. Corrected signed dynamic range in equation 40.
PA3	29-Mar-2007	Randy Yates	1. Updated document design.
PA4	17-Apr-2007	Randy Yates	1. Removed "this is a test" from introductory paragraph (!). 2. Added justification of signed division scaling.
PA5	23-Aug-2007	Randy Yates	Updated shift section.
PA6	07-Jul-2009	Randy Yates	1. Added derivation of unsigned division result. 2. Changed units notation to remove what looked like a subtractions (-).
PA7	01-Jan-2013	Randy Yates	1. Changed "right" to "left" in description of binary point position in last paragraph of page four (section 2.1). Thanks to Rick Lyons for finding this error! 2. Added examples to sections 5.12.1.1 and 5.12.1.2.
PA8	02-Jan-2013	Randy Yates	1. Changed "will lose precision" to "run the risk of losing precision" in section 5.12.1.1.

Table 1: Revision History

Author  
Randy YatesDate  
September 15, 2023Time  
13:34Rev  
PA8No.  
n/aReference  
fp.tex

## 12 References

### References

- [1] R. Yates, "Practical Considerations in Fixed-Point Arithmetic: FIR Filter Implementations."