

## VARIABLE LENGTH DATA FORMATS

Brakefield, James C., Technology Incorporated  
Houston, Texas 77058  
Quinn, Matthew J., NASA-Johnson Space Center  
Houston, Texas 77058

### ABSTRACT

From Information Theory, we learn that variable word length data formats can be more efficient than their constant word length counterparts. However, in the past, constant word length computers were designed because the hardware complexity demanded by variable word length machines offset any possible savings by completely filling up the memory. With the advent of LSI, many old concepts are no longer valid. The purpose of this paper is to discuss a number of variable length floating point and integer formats and to give the various advantages and disadvantages of their use. Often one knows in advance that a given integer will not exceed a certain magnitude or that a particular floating point number is only accurate to only "n" places of accuracy. Faced with this, it is good engineering to choose variable length floating point and integer formats which require the least amount of hardware or the minimum amount of software or which have some other dominant advantage. The formats discussed have the advantage that length change algorithms are invariant with respect to data types (unsigned, signed, floating point, integers, and complex numbers). Finally one particular computer, the STARAN associative array processor, uses a completely variable fixed point and floating point formats. Its hardware and software are described to show how a variable length data format is a practical goal.

## VARIABLE LENGTH DATA FORMATS

This paper describes the different fixed length formats in use throughout the computer industry. Integer formats are covered in one's complement, two's complement, and binary coded decimal. Then floating point numbers are discussed for a number of the popular computers. Finally, the hardware and software of the STARAN computer, one of the few variable word length computers in the industry, is described in some detail. Specific integer and floating point formats are recommended as most suitable for mathematical operations.

## VARIABLE LENGTH FLOATING POINT FORMATS

The following develops a number of "theorems" about equivalent number representation and formats. The notation is a variation on that of reference 1.

$r$  is the base, typically 2, 4, 8, 10 or 16

$\ell$  is the number of bits used in a base  $r$  number and includes any sign or exponent.

$US(r,\ell)$  is an unsigned number to base  $r$  using  $\ell$  bits

$ISM(r,\ell)$  is a sign and magnitude integer

$IOC(r,\ell)$  is a one's complement integer

$ITC(r,\ell)$  is a two's complement integer

$F(r,\ell)$  is a fractional number to base  $r$  using  $\ell$  bits

US, SM, OC, and TC will also be used as suffixes

$FP(r,\ell,e,m,a)$  is a floating point number of  $\ell$  bits with the mantissa normalized via base  $r$ .  $e$  is the exponent representation and  $m$  the mantissa representation.  $a$  is U for unnormalized and N for normalized.

$\approx$  is an associative and commutative relation expressing approximate equality of two number representations.

$\approx$  is like  $\approx$  except that the two operands need not be equally dense.

For Unsigned Numbers:

$$1.1 \quad 0 \leq US(2,\ell) \leq 2^{\ell-1}$$

$$1.2 \quad 0 \leq US(10,\ell) \leq 10^{\ell-1}$$

$$\text{Since } 2^m = 10^{m/\log_2 10}$$

$$1.3 \quad US(2,m) \approx US(10,4m/\log_2 10) \approx US(10,1.2m)$$

Which shows that base 10 requires 20% more bits to represent the same range of values.

The following will be used later:

$$2.1 \text{ US}(2, \ell) \approx 2^* \text{US}(2, \ell - \log_2 2) \approx 3^* \text{US}(2, \ell - \log_2 3) \approx 4^* \text{US}(2, \ell - \log_2 4)$$

$$2.2 \text{ So: } 2^{\text{US}(2, \ell)} \approx 4^{\text{US}(2, \ell - 1)} \approx 8^{\text{US}(2, \ell - 1.585)} \approx 16^{\text{US}(2, \ell - 2)}$$

For Integers:

$$3.1 -2^{\ell-1} + 1 \leq \text{ISM}(2, \ell) \leq 2^{\ell-1} - 1$$

$$3.2 -2^{\ell-1} + 1 \leq \text{IOC}(2, \ell) \leq 2^{\ell-1} - 1$$

$$3.3 -2^{\ell-1} \leq \text{ITC}(2, \ell) \leq 2^{\ell-1} - 1$$

$$3.4 \text{ Thus: } \text{US}(2, \ell) - 2^{\ell-1} \approx \text{ISM}(2, \ell) \approx \text{IOC}(2, \ell) \approx \text{ITC}(2, \ell)$$

$$3.5 \text{ Similarly: } I 'x' (2, \ell) \approx F 'x' (2, \ell) * 2^{\ell-1}$$

Note that ISM and IOC are symmetrical and thus contain a "-0" code.

ITC is not symmetrical, but the " $-2^{\ell-1}$ " code is not generally useful as " $-(-2^{\ell-1})$ " results in an overflow.

Floating Point:

Most computer floating point is with binary numbers, the only variations being the choice of exponent and mantissa representations and the radix used in normalizing the mantissa.

Examining the different exponent and mantissa representations:

$$4.1 \text{ FP}(r, \ell, \text{US}(s, k), -2^{k-1}, m, a) \approx$$

$$\text{FP}(r, \ell, \text{ISM}(s, k), m, a) \approx$$

$$\text{FP}(r, \ell, \text{IOC}(s, k), m, a) \approx$$

$$\text{FP}(r, \ell, \text{ITC}(s, k), m, a)$$

4.2 and

$$\text{FP}(r, \ell, e, \text{ISM}(s, k), a) \approx$$

$$\text{FP}(r, \ell, e, \text{IOC}(s, k), a) \approx$$

$$\text{FP}(r, \ell, e, \text{ITC}(s, k), a) \approx$$

$FP(r, l, e + k-1, FSM(s, k), a) \approx$

$FP(r, l, e + k-1, FOC(s, k), a) \approx$

$FP(r, l, e + k-1, FTC(s, k), a)$

Which means that one can use either a biased unsigned exponent or a signed integer exponent of any type with little difference in the range of floating point values which can be represented. Also the mantissa can be of any of the integer types. If the mantissa is of any of the fraction types the effect is to offset the exponent by the mantissa length (without sign).

Now for different  $r$  using unnormalized FP with specific  $e$  and  $m$  for illustration, (using 2.2):

5.1  $FP(16, l + 1, ISM(2, k-2), ISM(16, l-k+3), U) \approx$

$FP(8, l + .415, ISM(2, k-1.585), ISM(8, l-k+2), U) \approx$

$FP(4, l, ISM(2, k-1), IAM(4, l-k+1), U) \approx$

$FP(2, l, ISM(2, k), ISM(2, l-k), U)$

The summary of the above is that base 8 and 16 waste .415 and 1.0 bits respectively compared to base 2. This happens because although the exponent requires fewer bits to represent the same range, the mantissa requires even more bits to keep the same accuracy. Thus convenience in normalization costs in either more bits or less accuracy.

The final point is that for  $FP(2, l, e, m, N)$  an additional bit can be saved since the leading digit of the mantissa must be 1 bit and thus need not be stored. Various fractions of a bit can in theory be saved in the other radicies. (using 5.1).

$$6.1 \quad FP(2, \ell, e, m, U) \approx FP(2, \ell-1, e, m, N) \approx$$

$$FP(4, \ell, e, m, U) \approx FP(4, \ell - \log_2 4/3, e, m, N) \approx$$

$$FP(8, \ell + .415, e, m, U) \approx FP(8, \ell - \log_2 8/7 + .415, e, m, N) \approx$$

$$FP(16, \ell + 1, e, m, U) \approx FP(16, \ell - \log_2 16/15 + 1, e, m, N)$$

Thus  $FP(2, \ell-1, e, m, U)$  is optimum from the storage cost/fixed accuracy point of view.

Now we ask the question: What is the best format for variable length quantities? This is in theory a hardware or software engineering question which would have some relation to the "gate costs" of a particular technology. We will submit that truncation is the simplest method of shortening a quantity and the zero extension is the simplest method of lengthening a quantity. The following formats have the property that only truncation on the left is used for shortening and that only zero extension on the left is used for lengthening.

$$7.1 \quad US(2, \ell) = m_{\ell-1} 2^{\ell-1} + m_{\ell-2} 2^{\ell-2} \dots + m_1 2^1 + m_0 2^0$$

$$= m_{\ell-1}, m_{\ell-2}, \dots, m_1, m_0$$

$$7.2 \quad ISM(2, \ell) = (1-2*s)(m_{\ell-2} 2^{\ell-2} + m_{\ell-3} 2^{\ell-3} \dots + m_1 2^1 + m_0 2^0)$$

$$= m_{\ell-2}, m_{\ell-3}, \dots, m_1, m_0, s$$

$$7.3 \quad FP(2, 4k, ISM(2, k), FSM(2, 3k) + 1, N) =$$

$$(1-2*m_s)(1+m_{-1} 2^{-1} + m_{-2} 2^{-2} + \dots + m_{-3k+2} 2^{-3k+2} + m_{-3k+1} 2^{-3k+1}) * 2^{\uparrow}$$

$$(e_{k-2} 2^{k-2} + e_{k-3} 2^{k-3} + \dots + e_1 2^1 + e_0 2^0) * (1-2*e_s)$$

$$= m_{-3k+1}, m_{-3k+2}, m_{-3k+3}, e_{k-2}, \dots, e_1, m_{-5}, m_{-4}, m_{-3}, e_0, m_{-2}, m_{-1}, m_s, e_s$$

I.e., sign and magnitude with sign on the left, floating point with exponent and mantissa interleaved, and with the mantissa bit reversed.

And! the hardware need only know the new and old " $\ell$ " and not whether the quantity is US, ISM, or FP. The best interleaving factor, on the basis of required exponent range appears to be one exponent bit to three mantissa bits.

## EXAMPLE OF A VARIABLE WORD LENGTH COMPUTER

The STARAN (Goodyear Aerospace Corporation) is an example of a variable word length computer. This computer can have thirty two array processors with 256 processing elements in each array. At the NASA-Johnson Space Center, the STARAN consists of two arrays allowing 512 operations to be performed simultaneously. Each processing element can perform any one of the basic sixteen Boolean functions of two logic variables. Each array consists of 256 by 256 bits of memory with 256 processing elements. Each processing element accepts the three inputs from the X response store register, the Y response store register, and the M response store register. (Refer to Figure 1). The X and Y registers are used to perform the arithmetic operations. The M register is used to indicate which words of the array participate in the masked operation. When the mask bit is set, the programmed operation between X and Y is performed. When the mask bit is zero, no operation is performed between X and Y. The combination of one X bit, one Y bit, and one M bit trigger the arithmetic unit that outputs a bit to the 256 by 256 bit memory.

Since the computer can perform any one of the sixteen basic Boolean functions, the programmer has great flexibility in what he can do. He can build the overall program from the elementary Boolean building blocks. While this gives him considerable flexibility, it requires longer and more complex programs. To ease the programmer's work, a higher order language, APPLE (Array Processor Programming Language) is available with a considerable number of macros to perform frequently used functions. By combining Boolean operations, the more complex mathematical operations can be performed.

Addition, subtraction, multiplication, and division can be performed to any degree of accuracy using machine language or APPLE instructions. Furthermore, they can be done with reasonable programs. The convenient limit of 256 bits for addition and subtraction and 128 bits for multiplication and division is currently available in the STARAN at this time. To double these lengths of operands is still a reasonable program. This flexibility makes the STARAN an ideal candidate to simulate computer accuracy.

By programming a number of different solutions to problems involved with different word lengths, the computer architect would be free to explore the limitations of his paper design before he builds the prototype. The normal method of analyzing accuracy of a matrix inversion using an  $n \times n$  Hilbert Matrix, can be done using variable word length numbers to attain the required accuracy. No longer does one need to analyze accuracy budgets to determine if a particular word length gives the desired accuracy. By programming a general problem with a variable word length and variable parameters, actual accuracy of the computer can be determined. In the case of an  $n \times n$  Hilbert Matrix, the variable  $n$  can be varied as well as the word length of the numerical elements of the matrix array. This gives the computer architect a two dimensional picture of the accuracy of his hardware and software system.

#### VARIABLE WORD LENGTH EXAMPLE

The STARAN computer uses a variable word length format called S.I.F. where S is the sign, I is the integer and F is the fractional element. For example, i.8.7 indicates one sign bit, 8 integer bits, and 7 fractional bits.

Another example, i.7.0, refers to one sign bit and seven integer bits, whereas 1.0.23 refers to one sign bit and 23 fractional bits. In these examples, the length was never over 24 bits. However, there is no limit on the word length. Basically 256 bits on addition and subtraction, and 128 bits on multiplication and division is a practical word length constraint. However, the length can be extended by clever manipulation of the STARAN to include even greater word length if it should be necessary.

In summary, a number of different formats has been presented. A method of obtaining these formats is available. By programming the STARAN computer for variable word length, essentially any format is available to the computer architect.

#### References:

Sterbenz, Pat H., Floating Point Computation, p. 12, 1974 Prentice-Hall.

IBM 1620      FP(10,5n+10,ISM(10,10),FSM(10,5n))  
  . M M M . . . M M E E

Burroughs 2500 FP(10,4n+16,ISM(10,12),ISM(10,4n+4))  
E<sub>s</sub> E E M<sub>s</sub> M M . . . M M.

CDC 1604 FPOC(2,48,US(2,11)-1024,FUS(2,36))  
S e e e e e e e e . m m m . . . . . . . . m m m m

Philco 213 FP(2,48,ITC(2,12),FTC(2,36))  
S.m.m.m . . . . . m.m.m.m e e e e e e e e e e e e

BM 7030 FP(2,60,ISM(2,11),FSM(2,49))  
e e e e e e e e . m m m m . . . . . . . . m m m m +

optimal FP(2,4n,ISM(2,n),FSM(2,3n)+1.)  
m m m m m e . . . . . . . m m m e m m m s e s

