

Computer System Support for Scientific and Engineering Computation

Lecture 5, May 17, 1988 (revision date June 24, 1988)

Copyright ©1988 by W. Kahan and K.C. Ng
All rights reserved.

1 Conventional Floating-Point Formats

Let's remind ourselves what we have to deal with. We have to deal with radix β , which has to be drawn from a set of numbers

$$\text{Radix } \beta \in \{2, 4, 8, 10, 16\}$$

to correspond to machines actually in existence ($\beta = 4$ may be an exception, but we might just include it in although there are no machines that I know of today that use radix 4). I also talked about digits d_j

$$\text{Digits } d_j \in \{0, 1, 2, \dots, \rho - 1, \rho\} \text{ where } \rho = \beta - 1.$$

Now, conventional floating-point representation has to pass somehow three fields:

Sign	\pm
Significand, Mantissa, or Coefficient	$[d_1 d_2 \dots d_{p-1} d_p]$, p "significant" digits of Precision
Exponent	e , stored as $e + \text{Bias}$, an integer.

The *sign* field normally takes up one bit in binary but there are other representations that may take up a whole character. The second field has a number of names: you could call it *Coefficient*, or *Mantissa* ("Mantissa" is a very bad word: originally it represented the fraction part of a logarithm), or *Significand*. It seems that "Significand" is widely used now. Finally the *Exponent* field is simply an integer, often stored with bias (will be explained in a moment). Once you got these three fields, the interpreted value is as follows.

$\pm \underbrace{[d_1 d_2 \dots d_{p-1} d_p]}_{\text{Integer Coefficient}} \times \beta^{e_p}$	CDC Cyber 17× Burroughs B65××
$\pm [d_1 d_2 \dots d_{p-1} d_p] \times \beta^{e_0}$	CDC Cyber 18× IBM 370 DEC VAX "C" Language
$\pm [d_1 . d_2 \dots d_{p-1} d_p] \times \beta^{e_1}$	"Scientific" notation Calculators IEEE 754/854
... others

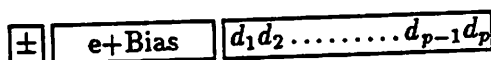
All values are the same if exponents are correlated:

$$e_0 = e_p + p = e_1 + 1.$$

Thus if you like to put the β -point in another place, you can accommodate that simply by changing the bias. What it really means is that the position in which you put the point has very little to do with the hardware (hardware is going to use integer arithmetic anyway). The trouble is that the representation is non-unique:

$$0.0123 = 0.123 \times 10^{-1} = 1.23 \times 10^{-2}$$

For hardware reasons it is a good idea to have a unique representation. So that is the purpose of normalization. A representation of a number will be called normalized if its first digit $d_1 \neq 0$.



$d_1 \neq 0$ for normalized nonzero number

An *unnormalized* nonzero number has $d_1 = 0$ but $e > \text{minimum}$. It is just another way to represent the value of a normalized number. Now whether zero is a *normalized* number or not is uninteresting and I won't argue about that, except that there are some unnormalized zeros:

Normalized zero : $d_1 = \dots = d_p = 0, e = \text{standard value, usually minimum.}$
Unnormalized zero : $d_1 = \dots = d_p = 0, e \neq \text{standard value.}$

There are reasons for unnormalized zeros on some machines. They can be used for converting a number to an integer represented in floating-point (cf. Fortran's AINT()). Finally there are something call subnormal:

Subnormal(denormalized)number: $d_1 = 0, e = \text{minimum value.}$

Note that because of the minimum exponent, subnormal numbers are *unique* representations for their numerical values, so they are *not unnormalized* numbers.

Examples: on a calculator with format $\pm d.ddd \times 10^{\pm ee}$

0.0123×10^0 is unnormalized
 1.230×10^{-2} is normalized
 0.012×10^{-99} is subnormal
 0.000×10^0 is normal zero by convention.

Note that the zero exponent in the normal zero is for display purposes. Internally it could be 0.000×10^{-99} . A more humane way is simply to display a single digit 0. That way you know you really got a zero rather than something that's hidden in the right hand side of the display window.

2 Why are most floating-point formats intended to be normalized?

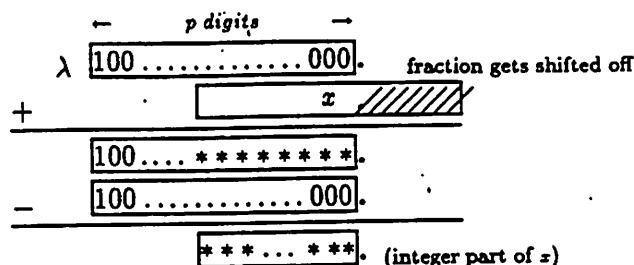
The reason for the normalization is to make the hardware simpler to design, especially for addition and subtraction.

Example: Add/Subtract $3.0 + 0.012 \Rightarrow 3.01$ to 3 significant decimals.

$\begin{array}{r} 3.00 \times 10^0 \\ + 1.20 \times 10^{-2} \\ \hline 3.00 \times 10^0 \\ 0.012 \times 10^0 \\ \hline 3.012 \times 10^0 \\ \downarrow \\ 3.01 \times 10^0 \end{array}$ <p>As usual</p>	$\begin{array}{r} 0.03 \times 10^2 \\ + 1.20 \times 10^{-2} \\ \hline 0.03 \times 10^2 \\ 0.00012 \times 10^2 \\ \hline 0.03012 \times 10^2 \\ \downarrow \\ 0.03 \times 10^2 \end{array}$ <p>just 3. !! in Unnormalized arithmetic IBM 370 ...</p>	$\begin{array}{r} 0.03 \times 10^2 \\ + 1.20 \times 10^{-2} \\ \hline 0.03 \times 10^2 \\ 0.00012 \times 10^2 \\ \hline 0.03012 \times 10^2 \\ \swarrow \text{normalization} \\ 3.012 \times 10^0 \\ \downarrow \\ 3.01 \end{array}$ <p>on Burroughs B65xx using double-width accumulator. (cf ALGOL)</p>
--	---	---

In the first column, when adding 0.012 and 3.0, the hardware would first discover that the operands' exponents are different, so it would initially do a pre-alignment. The last digit "2" has to go away. An interesting question is *when* should it go, and machines differ. Some machines will throw away the digit "2" before the addition, which doesn't do too much damage to the sum here, but you will see later that it can do a great deal of damage to subtraction. The point is, registers don't normally hold enough digits to keep everything that has been shifted right.

Now, the second column demonstrates that if an unnormal representation of 3.00 is allowed, the hardware will need four extra positions to hold all the digits. That seems to require a lot of extra digits. Most machines have floating-point add built in such a way that they only carry one or two extra digits. If that is the case the sum would end up with an incorrect result. That is why people say that unnormalized arithmetic is less accurate. It is not the unnormalized number that is less accurate; rather the hardware is designed in such a way that if you do arithmetic with unnormalized operands, you may get inaccurate results. That could be intentional because the shifting process provide a way to take the integer part of a number. Recall (see earlier notes for ceil and floor) that to get the integer part of a number, you add it to $\lambda = \beta^{p-1}$ and then subtract λ from the sum. The fraction part of x will then be rounded off.



That is a pain in the neck to get the integer part this way. There are people who say as long as you allow unnormalized zero in hardware, the same thing can be done faster by simply adding x to an aptly chosen unnormalized zero. The picture is the same as above except that the leading 1 of λ now is 0, and you only do the addition once. However, one should be careful about the exponent of the unnormalized zero (not necessarily p). For the hardware may carry extra digits (i.e., the fraction part of x may not be thrown away entirely if there are only p zeros ahead of the point) and you have to compensate by putting more digits in the unnormalized zero. I suppose it would make more sense to have an explicit instruction that converts to integer. That is what the IEEE Standard specified.

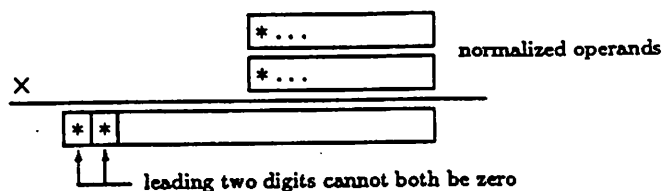
Finally, the third column on the previous page is what the Burroughs B65 machine does. Burroughs originally designed its B5000 machine to suit the language ALGOL. One of the characteristic of ALGOL is that integer operations may result in floating-point numbers. This avoids the conundrum in Fortran which would return 0 when computing $1/3$ because Fortran forces the result to be an integer. So ALGOL has both integer and floating-point "REAL" numbers in the sense that you cannot tell them apart. This was interpreted unnecessarily by the hardware people to mean that integers should look like floating-point numbers, and that's what they did. They encoded the floating-point numbers in such a way that the floating-point encoding for at least an unnormalized representation of an integral value would coincide bitwise with the encoding for an integer. Thus an unnormalized number could be as legitimate as others and in principle, any arithmetic you do with an unnormalized number should be the same as with its normalized counterpart. Burroughs accomplished that by keeping all the digits that other folks might throw away. In fact, the arithmetic was done in a double-width register (plus a digit), and then the result would be shifted to the left until either the left hand register filled or the bottom register is cleared. That explains why on Burroughs machines it is perfectly reasonable to have unnormalized numbers.

There have been various policies toward unnormalized operands. To summarize,

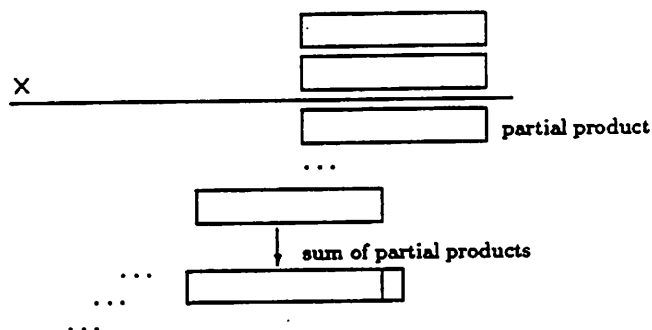
- IBM 370 Prenormalizes operands before \times , \div ; not before \pm (for the sake of $\text{AINT}(x)$).
- DEC VAX has no unnormalized operands
- IEEE 754/854 allows no unnormalized operands, but includes subnormal operands for gradual underflow.

Remark 1. Prenormalization before \times and \div is built in IBM 370 hardware. It is done that way to avoid unnecessary loss of accuracy: otherwise you would end up with a zero product if half of the leading bits in both operands were zeros. Or, if you decide to normalize the result, then you have to figure out in advance the amount of post-shifting, and you may end up with something that is harder to build. (But CYDOME's CYDRA5 does work this way.)

Remark 2. Another reason for prenormalization is cheaper hardware. It is possible to do the multiplication with a register that is just a little bit wider than the leading word if all you want is the leading word of the product. All you need are normalized numbers: for the product of two normalized numbers (length p) must have length either $2p$ or $2p - 1$; consequently you wouldn't have to carry more than one extra digit (called the *guard digit*) in order to know that there might be a post-shift.



Thus it's enough to use a register with some extra bits to hold the sum of the partial products:



3 Exponent Bias and Range Asymmetry

Why is there a Bias for the exponent? Part of the reason is to compensate for where you may have put the "point" in the documentation. However, the bias does turn up when you ask how balanced or symmetrical is the range of numbers. If you multiply the biggest number with the smallest normal number, then you get

$$(\text{Max Normal}) \times (\text{Min Normal nonzero}) \approx (\beta \times \beta^{e_{\max}})(1 \times \beta^{e_{\min}}) = \beta^{1+e_{\max}+e_{\min}}$$

If this product is too far from 1, the range is regarded as too asymmetrical, and causes troubles with expressions like

$$q := (\text{small integer})/x \quad \dots \text{can over/underflow too easily?}$$

In other words, the reciprocal of the smallest number may be much bigger than the biggest number (or the other way around). In that case over/underflow may come to you as a surprise. Here is another problem:

$$q := x \cdot y/z \quad \dots \text{what order avoids spurious over/underflow?}$$

We'll come back to the above problem later. Now let's see how to choose a bias. Recall if

$$X : \boxed{\pm} \boxed{e + \text{Bias}} \boxed{d_1 d_2 \dots d_{p-1} d_p}$$

then the interpreted value is

$$x = \pm[d_1 \cdot d_2 \dots d_{p-1} d_p] \times \beta^e, \quad e_{\min} \leq e \leq e_{\max}. \quad (1)$$

Normally the exponent field $e + \text{Bias}$ would be stored in a field l bits wide; thus

$$0 \leq e + \text{Bias} \leq 2^l - 1.$$

It follows that

$$e_{\min} = -\text{Bias}, \quad e_{\max} = 2^l - 1 - \text{Bias}.$$

Therefore $1 + e_{\max} + e_{\min} = 2^l - 2 \cdot \text{Bias}$. That explain how the choice of Bias can influence the symmetry of the exponent range. If you want not to have too asymmetrical a range, you should choose Bias near 2^{l-1} . If overflow is more to be avoided than underflow (e.g. if underflow is gradual), choose $\text{Bias} = 2^{l-1} - 1$. (DEC VAX¹: $\text{Bias} = 2^{l-1} + 1$ predisposes to overflow!)

Range Asymmetry is Severe on Some Machines:

- CDC Cyber 17x:

$$\begin{aligned} \text{Max. magnitude} &\approx 2^{1022+48} \\ \text{Min. normal magnitude} &= 2^{-1022+47} \end{aligned}$$

$$(\text{Max}) \cdot (\text{Min}) \approx 2^{95}$$

Hence $2^{94}/\text{Max. underflows!}$

- Similar trouble on Burroughs B65x.

IEEE 754, single precision:

$$\begin{aligned} \text{Max.} &\approx 2^{128} \\ \text{Min. Normal} &= 2^{-126} \\ \text{Min. Subnormal} &= 2^{-149} \end{aligned}$$

$$\begin{aligned} (\text{Max}) \cdot (\text{Min. Normal}) &\approx 4 \\ (\text{Max}) \cdot (\text{Min. Subnormal}) &\approx 2^{-21} \end{aligned}$$

Hence $2^{-20}/(\text{Min. Subnormal})$ overflows.

Remark. Subnormal numbers are intended to cope with certain situations in addition and subtraction (see Lecture 7b, section 4), not multiplication and division. Dividing a number by a subnormal number and getting an overflow simply tells the user the bad news of the earlier loss of significant digits. On machines that flush subnormal numbers to zero you will encounter division by zero anyway. So we choose our bias by looking at $(\text{Max}) \cdot (\text{Min. Normal})$ rather than $(\text{Max}) \cdot (\text{Min. Subnormal})$.

¹VAX's manual said $\text{Bias} = 2^{l-1}$, but here we interpreted x 's value as in (1), whereas VAX puts the "point" before d_1 .

Exercise: Exhibit a program that starts from any three given positive numbers x, y, z and computes $p := x \cdot y \cdot z$ in some order that avoids undeserved over/underflow. Do likewise for $q := x \cdot y/z$.

The first one is easy. Here is how you do it: sort x, y, z so that, say,

$$x \geq y \geq z > 0,$$

and then compute

$$p := (x \cdot z) \cdot y.$$

That is, you multiply two extreme numbers first. You should be able to verify, if p is computed that way, over/underflow occurs only if it has to. The second one is not that easy. We could sort x, y , and $1/z$ except $1/z$ might over/underflow.

Can you find a way?

I will leave to you the indistinct possibility that you might be able to find a portable program that will work on all machines. This is a problem that does arise in "real life" from time to time; cf. program SVD2 \times 2.

4 Lexicographic Order

The format we have chosen for floating-point numbers preserves "lexicographic order". That is the reason we chose a biased exponent rather than an exponent with an independent sign (like Burroughs). Consider the string X

$$\text{String } X: \quad \boxed{\pm} \underbrace{\boxed{e + \text{Bias}}}_{l \text{ bits}} \boxed{d_1 d_2 \dots d_{p-1} d_p}$$

with its floating-point value

$$x = \pm [d_1 \cdot d_2 \dots d_{p-1} d_p] \times \beta^e, \quad 0 \leq e + \text{Bias} \leq 2^l - 1.$$

Assume X is Normalized (i.e., either $d_1 > 0$ or $(e + \text{Bias}) = [d_1 d_2 \dots d_{p-1} d_p] = 0$). Then the order of x 's as floating-point values is the same as the order of X 's as integers with the same sign-convention. This permits *quick* floating-point comparison without doing floating-point subtraction. Note that we must have the correct sign-convention in order to make the lexicographic ordering work. Let's examine three sign-conventions that come to mind:

$$\begin{array}{ll} X & \boxed{+} \boxed{e + \text{Bias}} \boxed{d_1 d_2 \dots d_{p-1} d_p} \\ -X & \boxed{-} \boxed{\text{What } \dots} \boxed{\dots \text{ goes in here?}} \end{array}$$

1. Sign-magnitude, used by IBM, DEC, IEEE 754. Only the sign bit changes.
2. 1's complement, used by CDC Cyber 17x. The exponent and significand are 1's-complemented as well as the sign.
3. 2's complement, used by HP-3000, GE/Honeywell. The entire word is 2's-complemented.

Here are the crucial considerations:

- Sign-Magnitude \equiv 1's complement has ± 0 , allows $\pm\infty$.
- 2's complement has one *unsigned* 0, must have a negative value that either *overflows* when negated, or is an outlaw (e.g. unsigned ∞).

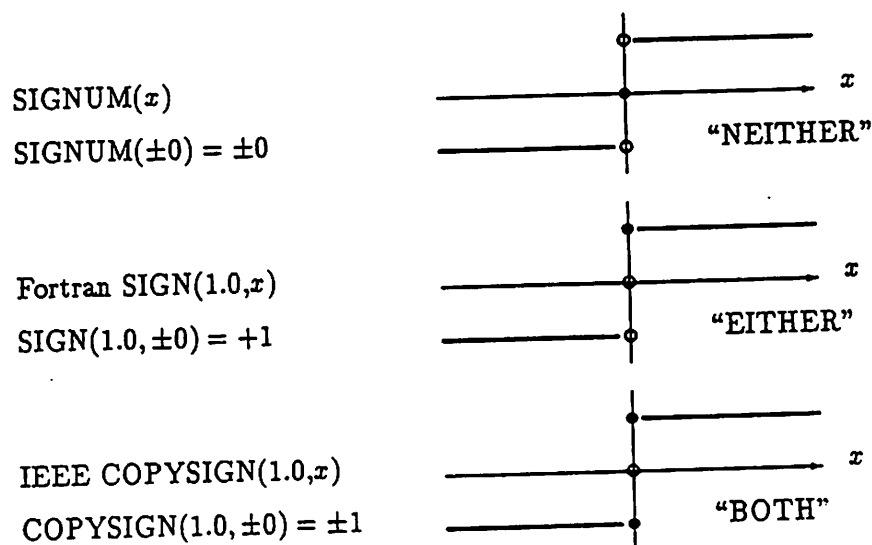
It turns out that Sign-Magnitude is best because it allows ± 0 . There are valuable uses for a signed zero in complex arithmetic, especially in conformal maps of slitted domains (cf. Powell & Iserles (ed.) (1987) "The state of the Art in Numerical Analysis" ch.7 Oxford U.P.).

People may get uncomfortable about signed zero because it is very difficult to understand how to distinguish a $+0$ and a -0 . In fact on a well designed machine they are equal:

$$+0 = -0.$$

But, if $f(x)$ is *discontinuous* at $x = 0$, then it *may* distinguish $f(+0)$ from $f(-0)$. Consider to which side of a discontinuity the discontinuity itself should be attached:

e.g.



The discontinuity of SIGNUM function is attached to neither of the continuous pieces. The Fortran SIGN function attaches the discontinuity to one of its two continuous components. IEEE COPYSIGN function attaches the discontinuity at both sides, making the graph of the function two *closed* components.

So the ambiguity of signed zero has to be resolved only when the function is actually discontinuous at zero; otherwise you can't care. The reciprocal function is an example: if you take the reciprocal of zero, its sign shows up

$$\frac{1.0}{\pm 0.0} = \pm\infty \Rightarrow \frac{1}{1/x} \approx x \text{ even for } x = \pm\infty.$$

It's true that most people don't care about these uses, and therefore it is vitally important that signed zeros be implemented in such a way that for every operation that is continuous at zero, they cannot be told apart. That is true in the IEEE Standards: $+0 = -0$ unless you have a discontinuous operation in mind. The signed zero propagates properly in the IEEE Standards, e.g., $3.0 \times (-0.0) = -0.0$ (on an IBM 370, $3.0 \times (-0.0) = +0.0$).

There are some conundrums about signed zero. For example, thanks to minus zero, we have to ask compiler writers not to replace $a - b$ by $-(b - a)$ because they have opposite signs when they are zero. It OK to do $(-b) + a$ though. So minus zero does pose some technical problems in the handling of details. But those problems in all situations that I know of do not introduce any performance disadvantage. That's why IEEE Standards have minus zero: it is useful, and doesn't cost anything except to the implementors.

5 Precision = Resolution (not Accuracy)

Precision is a measure of the resolving power of the arithmetic; accuracy is something that tells how right or wrong the result is. If I write $\pi = 3.14159023487628374$ then the precision is the number of digits (18) I wrote but the accuracy is only six figures. I can utter my inaccurate approximation of π to a *very high precision*. The distinction is necessary because in general people believe that if they carry more precision they will get more accuracy. That is usually true but beautiful examples (which we will come to) show that carrying more precision doesn't improve accuracy in every instance.

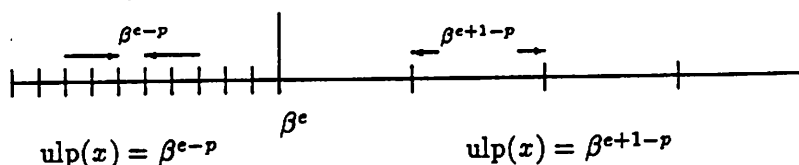
The *Relative Precision* to which you can specify a number x is the following quotient:

$$\frac{\text{ulp}(x)}{|x|}.$$

Here ulp stands for unit in the last place and if $x = \pm[d_1 \cdot d_2 \dots d_{p-1} d_p] \times \beta^e$, then $\text{ulp}(x) = [0.0 \dots 01] \times \beta^e$. In general, if $\beta^e < x < \beta^{e+1}$ then $\text{ulp}(x) = \beta^{e+1-p}$. But there is ambiguity:

$$\text{ulp}(\beta^e) = \beta^{e+1-p} \text{ or } \beta^{e-p}$$

depending on whether you look left or right.



Every time you cross a power of the radix the unit in the last place jumps by a factor of β . This is called *wobbling precision*. The larger the radix, the worse the wobble.

$$\underbrace{\beta^{-p} \leq \frac{\text{ulp}(x)}{|x|} \leq \beta^{1-p}}_{\text{"Wobble" by factor } \beta}$$

6 Error in Floating-point Operations

Let $\oplus \in \{+, -, \times, /\}$. Let A, B, C be the corresponding names of variables whose values are a, b, c . I'll use this notation to distinguish between the assignment statement in a program and the corresponding mathematically exact algebraic relation.

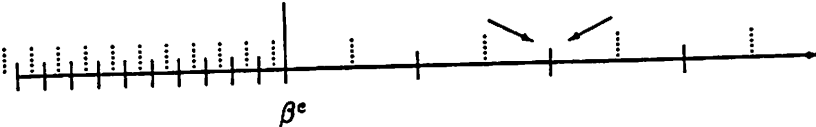
$$\begin{array}{ll} \text{Program statement} & "A := B \oplus C" \\ \text{produces value} & a = b \oplus c \text{ Rounded} \end{array}$$

The discrepancy between the exact result and the rounded one is easy to figure out if rounding is done in the best possible way: you will lose only half of the unit in the last place.

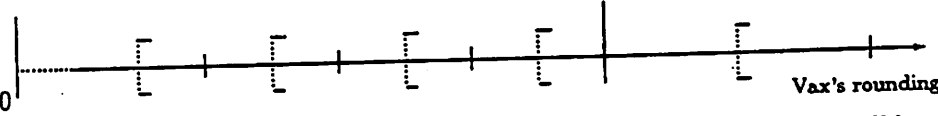
$$\begin{aligned} b \oplus c &= \pm x.xxx \dots xxx \, xxx \dots \\ a &= \pm x.xxx \dots xxx \\ &\quad \leftarrow \quad p \quad \rightarrow \end{aligned}$$

$$\text{correct rounding} \Rightarrow |a - b \oplus c| \leq \frac{1}{2} \text{ulp}(b \oplus c)$$

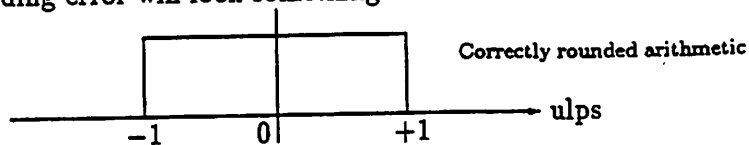
There are two special cases: (i) $a = b \oplus c$ if it is representable exactly, and (ii) $a = \text{something else}$ if over/underflow occurs. So values that fall between two dotted lines will be rounded to the center number (heavily marked):



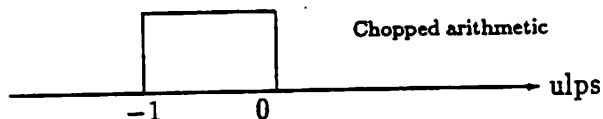
The only ambiguity now is what if $b \oplus c$ falls midway between two adjacent representable values (i.e., the value lies on the dotted line)? This ambiguity gets resolved in a variety of ways, of which only two deserve close attention. These two roundings are both called *correct rounding* because their error doesn't exceed half of an ulp. But one is slightly more correct than the other. First let's look at DEC VAX. DEC VAX rounds midway case *away* from zero, i.e., *up* in magnitude:



A fair question is why do we prefer correct rounding? The reason has been described in terms of the law of averages which the bias in rounding is smaller if you do round correctly. Here is a brief argument (by the way this is *not* the best argument): imagine that the values you round vary a lot and their digits behave randomly. For correct rounding, it looks as if as many numbers are being rounded up as are being rounded down. Roughly speaking the distribution of rounding error will look something like

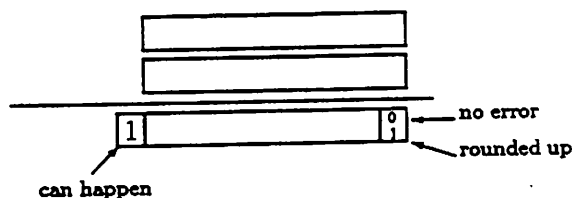


Here I am plotting the magnitude of rounding error in ulps. What this tells us is that we can calculate means and standard deviations. If all the rounding errors contribute their effect independently, then the means look like 0. That's a lot better than chopping (IBM 370, Cray, Cyber 170). The error for chopped arithmetic is biased (probabilistically). On the average it is $-\frac{1}{2}$ ulp:

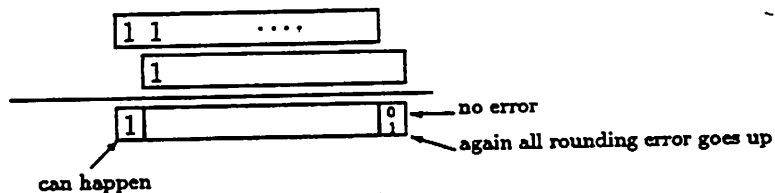


This bias shows up in many computation. For many years people have decried IBM for this fault because you can do computations (e.g. a long string of multiplies and adds) where you see your numbers steadily and systematically drifting downwards as you keep on computing with them.

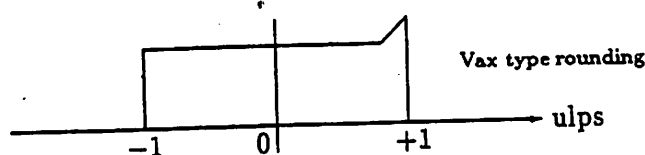
The VAX's type of rounding seems to have bias equal to 0. But that's not quite true. Here is the reason. A typical add might not involve any pre-shifting. It is often found when you do a floating-point add that the exponent of the operands are the same. The result may have a carry on the left, which means you have to shift one digit off:



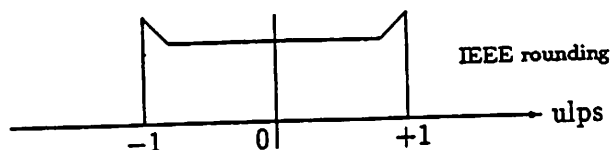
Consequently, all such rounding errors go up on a VAX. Another similar situation is a one-bit shifted subtract:



So there is a spike near +1 caused by the midway round up. And the bias is a bit off zero.



On the other hand, the IEEE rounds midway cases to nearest even, which means the rounding may go down as well as go up, depending on the least significant bit. So there are two spikes: one on -1 and one on $+1$. And the bias would again very near zero.



So that's the argument (not a very good one) for IEEE rounding in term of the law of averages. Unfortunately the law of averages is not something you should not depend on too much. Because when you get hurt by floating-point computation (which is rare), the pain is caused by the worst case, not the average. A much better argument comes about because certain computations are repetitive:

```

loop:  y = x+z
      ...
      ...      (y,z unchanged inside the loop)
      ...
      x = y-z
      goto loop

```

This may occur in solving a large system of coupled equations, when some unknowns converge faster than the others. Now the question is what does the loop do to the value of x ? A remarkable theorem says, if the value of x changes it will change once and never again, provided you use the IEEE rounding. If you round as the Vax does, you run the risk that any time you go around this cycle, the value of x may keep on drifting along one direction. That is the reason that we really want to round to nearest even; it prevents certain kinds of cyclical operations from drifting. This is a very general result and is often seen under the name of Harry Diamond's theorem². What the theorem says is roughly the following: if you compute the monotonic convex function $y = f(x)$ and then its inverse $x = f^{-1}(y)$, repeatedly, rounding off each time, then if you round the way IEEE Standard does, you could only make at most two changes in x .

7 Range/Precision Tradeoff for Radices $\beta = 2^k$

There are machines with various radices of the form $\beta = 2^k$:

²Harry Diamond didn't prove this theorem, but he proved something so near that it's worth attaching his name to it.

Name	β	k	Who?
Binary	2	1	IEEE 754, DEC VAX, CDC, CRAY, ...
Quaternary	4	2	...no more ...
Octal	8	3	Burroughs B65x
Hexadecimal	16	4	IBM 370, Amdahl, ...

The question is which radix is *best*? To answer this question, let's work out the wordsize, the range, and the precision of radix $\beta = 2^k$. Consider the floating-point word:

$$X = \boxed{\begin{array}{c} \overset{1 \leftarrow l \text{ bits} \rightarrow}{\pm} \overset{\leftarrow p \text{ sig. digits} = p \cdot k \text{ bits} \rightarrow}{e + \text{Bias}} \mid d_1 d_2 \dots d_{p-1} d_p \end{array}}$$

Here for convenience the interpreted value has no "point" after d_1 :

$$x = \pm \beta^e \times [d_1 d_2 \dots d_{p-1} d_p]$$

where $\beta = 2^k$ and

$$0 \leq e + \text{Bias} \leq 2^l - 1$$

$$0 \leq [d_1 d_2 \dots d_{p-1} d_p] \leq \beta^p - 1$$

Using the above notation, we have

$$\boxed{\text{Total wordsize: } w = 1 + l + p \cdot k \text{ bits}}$$

Let $\rho = \beta - 1$ so $[00 \dots 00] \leq [d_1 d_2 \dots d_{p-1} d_p] \leq [\rho \rho \dots \rho \rho]$. Normally X is normalized, i.e., $d_1 \geq 1$ unless $x = 0$. Now the range over which X varies has

$$\begin{aligned} \frac{\text{Max. } x}{\text{Min. } x > 0} &= \frac{\beta^{2^l-1-\text{Bias}} \times [\rho \rho \dots \rho \rho]}{\beta^{0-\text{Bias}} \times [10 \dots 00]} \\ &= \frac{\beta^{2^l-1} \times (\beta^p - 1)}{\beta^{p-1}} \\ &\approx \beta^{2^l} = 2^{k \cdot 2^l} \end{aligned}$$

Thus

$$\boxed{\text{Range: } 2^{k \cdot 2^l}}$$

Finally the worst-case precision is

$$\max_{x>0} \frac{(\text{successor of } x) - x}{x} = \frac{[100 \dots 001] - [100 \dots 000]}{[100 \dots 000]} = \frac{1}{\beta^{p-1}} = 2^{k \cdot (1-p)}.$$

Thus

$$\boxed{\text{Worst-Case Precision: } 2^{k \cdot (1-p)}}$$

Now we ask what is the word size of a binary format that has the same *range* and *worst-case precision* as $\beta = 2^k$? The result is interesting: you end up with smaller wordsize if you use binary. Here is the analysis. Let's say we use l' exponent bits and p' significant bits for binary. To achieve the same range and precision, we must have

$$2^{2^{l'}} = 2^{k \cdot 2^{l'}}, \quad 2^{1-p'} = 2^{k \cdot (1-p)}.$$

Thus $l' = l + \log_2 k$ and $p' = 1 + k \cdot (p - 1)$, and the wordsize is $w' = 1 + l + \log_2 k + k \cdot (p - 1) = w - (k - \log_2 k - 1)$. Hence $w - w' = k - \log_2 k - 1 \geq 0$ for all $k \geq 1$. The following table shows how many digits get lost for $k = 1, 2, 3, 4$:

Name	k	lost bits ($k - \log_2 k - 1$)
Binary	1	0 (-1 for Hidden bit!)
Quaternary	2	0
Octal	3	$(2 - \log_2 3) \approx 0.415$
Hexadecimal	4	1

Note that it seems binary and quaternary are the same. However in binary with normalized numbers the leading bit is 1 and you don't have to store it. You can't do that for other radices. So with a hidden bit,

Binary beats Quaternary by 1 bit
 Octal by 1.415 bits
 Hex. by 2 bits

That is an argument that has been used to prove that binary is best; the argument is not that one can save a few bits using binary radix, but that for a given memory wordsize it has greater precision than other radices. not much of an argument since memory is cheap. But a stronger argument is *WOBBLING PRECISION* (see section 5). In order to understand the importance of wobbling precision, you should read my note "Roundoff in Polynomial Evaluation, 18 October 1986". On the other hand the advantage of decimal is that people can understand it. So to conclude: use decimal when your machines will produce numbers that you expect people to read (people like tax preparers and collectors). Because they think decimal, run your machine in decimal so you do just what they naively imagine you do. Otherwise, use binary (for engineers, scientists, etc).

CONVENTIONAL FLOATING-POINT FORMATS

Radix $\beta = 2, 4, 8, 10 \text{ or } 16$.

Digits $d_j \in \{0, 1, 2, \dots, \beta-1, \beta\}$ where $\beta = \beta-1$.

Three FIELDS :

Sign $\dots \pm$

$\left\{ \begin{array}{l} \text{Significant} \\ \text{Mantissa} \\ \text{Coefficient} \end{array} \right\} \dots [d_1 d_2 \dots d_{p-1} d_p] \dots p \text{ "significant" digits of Precision}$

Exponent $\dots e$, stored as $e + \text{Bias}$, an integer.

Interpreted VALUE :

$\pm \underbrace{[d_1 d_2 \dots d_{p-1} d_p]}_{\text{Integer Coefficient}} \times \beta^{e_p} \dots \text{CDC Cyber 17x, Burroughs B65xx}$

OR $\pm [d_1 d_2 \dots d_{p-1} d_p] \times \beta^{e_0} \dots \text{IBM 370, DEC VAX, CDC Cyber 18x, "C" language}$

OR $\pm [d_1 \cdot d_2 d_3 \dots d_{p-1} d_p] \times \beta^{e_1} \dots \text{"Scientific" notation, Calculators, IEEE 754/854}$

OR $\dots \text{others} \dots$

All values are the same if exponents are correlated;

$$e_0 = e_p + p = e_1 + 1.$$

Shift of β -point without change in value \equiv change in Bias.

Kahan
17 May

NON-UNIQUENESS OF REPRESENTATION

$$0.0123 = 0.123 \times 10^{-1} = 1.23 \times 10^{-2}$$

NORMALIZATION \Rightarrow UNIQUE REPRESENTATION



$d_1 \neq 0$ for NORMALIZED NONZERO No's.

UNNORMALIZED NONZERO NO. has $d_1 = 0$
but $e > \text{minimum}$.

\therefore UNNORMALIZED NO. is another way to
represent the value of a NORMALIZED NO.

NORMALIZED ZERO: $d_1 = d_2 = \dots = d_{p-1} = d_p = 0$
& $e = \text{standard value, usually minimum}$.

UNNORMALIZED ZERO: $d_1 = d_2 = \dots = d_{p-1} = d_p = 0$
but $e \neq \text{standard value}$.

SUBNORMAL (DENORMALIZED) NO. has $d_1 = 0$
and $e = \text{minimum value}$.

SUBNORMAL NO. is UNIQUE REPRESENTATION of its value,
so NOT UNNORMALIZED.

e.g. calculator $\pm d.ddd \times 10^{\pm ee}$

0.0123×10^0 is UNNORMALIZED

1.230×10^{-2} is NORMALIZED

0.012×10^{-99} is SUBNORMAL

0.000×10^0 is NORMAL ZERO by convention.

WHY ARE MOST FLOATING-POINT FORMATS
INTENDED TO BE NORMALIZED?

(NOT for better accuracy!)

Because it simplifies the hardware.

e.g. ADD/SUBTRACT $3.0 + 0.012 \Rightarrow 3.01$ to 3 sig dec

$$\begin{array}{r} 3.00 \times 10^0 \\ + 1.20 \times 10^{-2} \\ \hline \end{array}$$

$$\begin{array}{r} 3.00 \times 10^0 \\ 0.012 \times 10^0 \\ \hline \end{array}$$

$$3.012 \times 10^0$$

↓

$$3.01 \times 10^0$$

As usual

$$\begin{array}{r} 0.03 \times 10^2 \\ + 1.20 \times 10^{-2} \\ \hline \end{array}$$

$$0.03 \times 10^2$$

$$0.00012 \times 10^2$$

$$0.03012 \times 10^2$$

↓

$$0.03 \times 10^2$$

just 3. !!

in UNNORMALIZED
ARITHMETIC

IBM 370,

$$\begin{array}{r} 0.03 \times 10^2 \\ 1.20 \times 10^{-2} \\ \hline \end{array}$$

$$0.03 \times 10^2$$

$$0.00012 \times 10^2$$

$$0.03012 \times 10^2$$

↙

$$3.01 \times 10^0$$

↓

$$3.01$$

on Burroughs B55xx

using double-width
accumulator. (cf. 41604)

IBM 370 Prenormalizes operands before \times , \div .

not before \pm , ... $\text{AINT}(x)$

DEC VAX has NO UNNORMALIZED OPERANDS

IEEE 754/854 allows no UNNORMALIZED OPERANDS,

but includes SUBNORMAL OPERANDS

for GRADUAL UNDERFLOW, q.v.

EXPONENT BIAS and RANGE (A) SYMMETRY

$$x = + [d_1 d_2 d_3 \dots d_{p-1} d_p] \times \beta^e$$

$$e_{\min} \leq e \leq e_{\max}$$

$$X = \boxed{+} \boxed{e + \text{Bias}} \boxed{d_1 d_2 d_3 \dots d_{p-1} d_p}$$

$$(\text{Max. Normal } x) \cdot (\text{Min. Normal Nonzero } x)$$

$$\doteq (\beta \times \beta^{e_{\max}}) (1 \times \beta^{e_{\min}}) = \beta^{1+e_{\max}+e_{\min}}$$

If this product is too far from 1, the range is regarded as too asymmetrical, and causes troubles with expressions like

$$\dots q := (\text{small integer}) / x \quad \dots \text{can over/underflow too easily?}$$

$$\dots q := x \cdot y / z \quad \dots \text{what order avoids spurious over/underflow?}$$

$$\text{If } 0 \leq e + \text{Bias} \leq 2^l - 1 \quad \text{for } l \text{ bits}$$

$$\text{then } e_{\min} = -\text{Bias}, \quad e_{\max} = 2^l - 1 - \text{Bias},$$

$$1 + e_{\max} + e_{\min} = 2^l - 2 \times \text{Bias} \quad \dots \text{for balance, choose Bias near } 2^{l-1}.$$

If Overflow is more to be avoided than Underflow, (e.g. if Underflow is gradual), choose $\text{Bias} = 2^{l-1} - 1$.

(DEC VAX: $\text{Bias} = 2^{l-1} + 1$ predisposes to overflow)

RANGE ASYMMETRY is SEVERE
ON SOME MACHINES.

e.g. CDC Cyber 17x :

$$\begin{aligned}\text{Max. magnitude} &\doteq 2^{1022+48} \\ \text{Min. normal magnitude} &= 2^{-1022+47}\end{aligned}$$

$$(\text{Max.}) \cdot (\text{Min.}) \doteq 2^{95}$$

∴ $2^{94} / \text{Max.}$ UNDERFLOWS !

Similar trouble on Burroughs B65xx.

IEEE 754 , SINGLE PRECISION :

$$\text{Max.} \doteq 2^{128}$$

$$\text{MIN. NORMAL} \doteq 2^{-126}$$

$$\text{MIN. SUBNORMAL} = 2^{-149}$$

$$\circ \circ (\text{Max.}) \cdot (\text{MIN. NORMAL}) = 4$$

$$(\text{Max.}) \cdot (\text{MIN. SUBNORMAL}) \doteq 2^{-21}$$

∴ $2^{-20} / (\text{MIN. SUBNORMAL})$ OVERFLOWS.

PROBLEM: Exhibit a program that starts from any three given positive numbers x, y, z and computes $p := x \cdot y \cdot z$ in some order that avoids undeserved over/underflow. Do likewise for $q := x \cdot y / z$.

To compute p : Sort x, y, z so that, say, $x \geq y \geq z > 0$, and then compute $p := (x \cdot z) \cdot y$.

To compute q , we could sort x, y and $1/z$ except $1/z$ might over/underflow.

CAN YOU FIND A WAY?

This is a problem that does arise in "real life" from time to time; cf. program SVD2x2.

LEXICOGRAPHIC ORDER

$$\text{String } X = \boxed{\pm} \overset{\leftarrow l \text{ bits} \rightarrow}{\boxed{\text{exponent} + \text{Bias}}} \boxed{d_1 d_2 \dots d_{p-1} d_p}$$

$$\text{Floating-pt value } x = \pm [d_1 d_2 \dots d_{p-1} d_p] \times \beta^{\text{exponent}}$$

$$\text{for } 0 \leq \text{exponent} + \text{Bias} \leq 2^l - 1.$$

Assume X is Normalized;

i.e. either $d_1 > 0$

or $(\text{exponent} + \text{Bias}) = [d_1 d_2 \dots d_{p-1} d_p] = 0.$

Then order of x 's as floating-point values
is same as order of X 's as integers
with the same sign-convention

i.e. sign-magnitude, or

1's complement if $-X \equiv 1's \text{ compl.}(X)$, or

2's complement $-X \equiv 2's \text{ compl.}(X).$

This can permit QUICK FLOATING-PT. COMPARISON
without floating-pt. subtraction

$$X = \boxed{+} \boxed{\text{exponent} + \text{Bias}} \boxed{d_1 d_2 \dots d_p d_r}$$

$$-X = \boxed{-} \boxed{\phantom{\text{exponent} + \text{Bias}}} \boxed{}$$

Sign-Magnitude ?

IBM, DEC

1's-complement ?

CDC Cyber 17x

2's complement ?

HP, GE/Honeywell

Which is best ?

- Sign-Magnitude \equiv 1's complement ,
has ± 0 , allows $\pm \infty$.
- 2's complement has UNSIGNED 0 ,
must have a negative value
that either OVERFLOWS
when negated, or is an
outlaw (e.g. UNSIGNED ∞).

Sign-Magnitude is best because it
allows ± 0 . ! Valuable for
COMPLEX ARITHMETIC - CONFORMAL MAPS OF
SLITTED DOMAINS.

1. Powell & Iserles (ed.) (1987) "The State of the Art
in Numerical Analysis" ch.7 Oxford U.P.

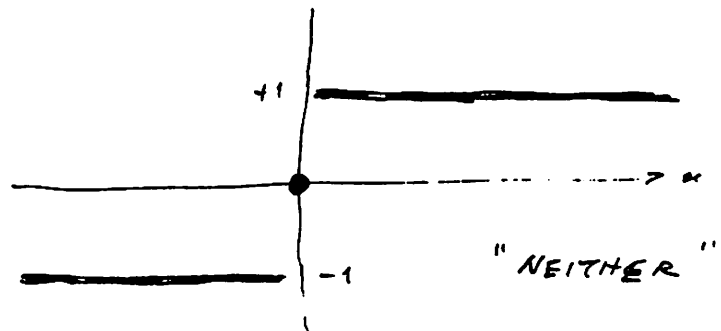
$$+0 = -0$$

But, if $f(x)$ is DISCONTINUOUS
at $x=0$, then it may

distinguish $f(+0)$ from $f(-0)$:

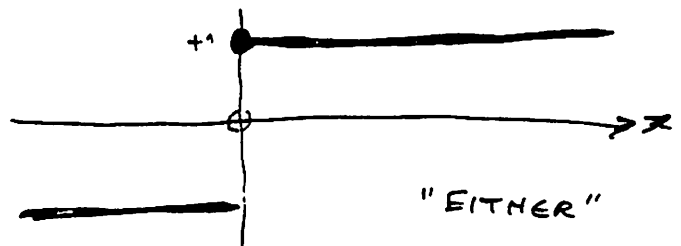
e.g.

SIGNUM (x)



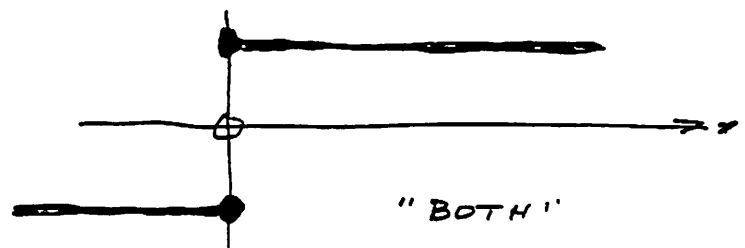
$$\text{SIGNUM}(\pm 0) = \pm 0$$

FORTRAN SIGN(1.0, x)



$$\text{SIGN}(1.0, \pm 0.0) = +1.0$$

IEEE COPYSIGN(1.0, x)



$$\text{COPYSIGN}(1.0, \pm 0.0) = \pm 1.0$$

$$1.0 / \pm 0.0 = \pm \infty \Rightarrow 1/(1/x) \doteq x.$$

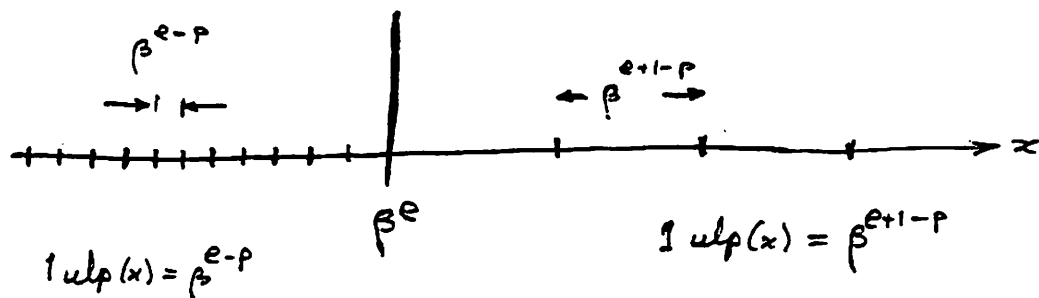
PRECISION = RESOLUTION
(not ACCURACY)

Relative Precision at x :

$$\frac{1 \text{ulp}(x)}{|x|}$$

where, if $x = \pm [d_1 d_2 \dots d_p] \times \beta^e$
then $\text{ulp}(x) = [0.0 \dots 0 1] \times \beta^e$

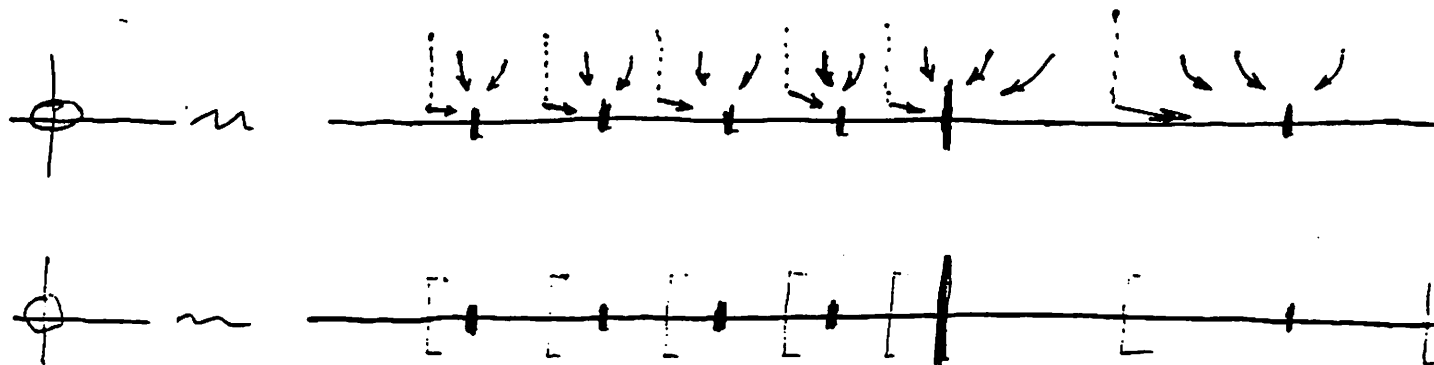
i.e. if $\beta^e < x < \beta^{e+1}$ then $\text{ulp}(x) = \beta^{e+1-p}$
but $\text{ulp}(\beta^e) = \beta^{e+1-p}$ or β^{e-p} ambiguously



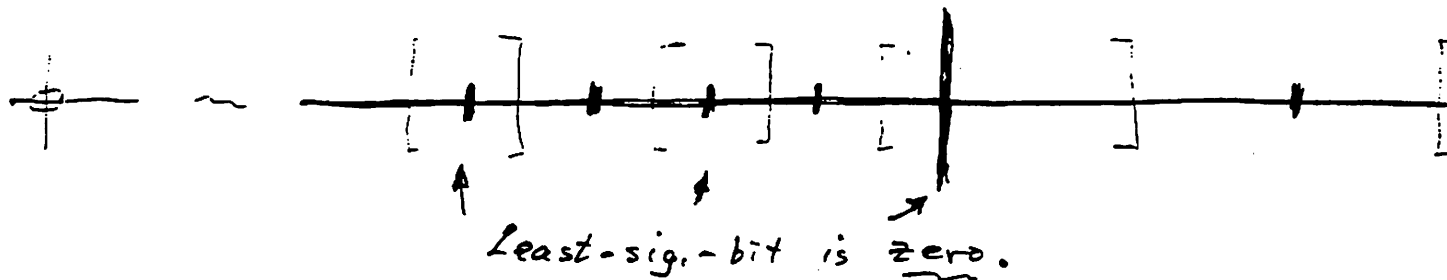
$\therefore \beta^{-p} \leq \frac{1 \text{ulp}(x)}{|x|} \leq \beta^{1-p}$
 \uparrow "WOBBLE"
 by factor β

"CORRECT" ROUNDING OF MIDWAY CASES

DEC VAX rounds MIDWAY CASES away from zero,
i.e. UP in magnitude.



IEEE 754/254 rounds MIDWAY CASES to "NEAREST EVEN":



e.g. to 4 sig. dec.

IEEE 854 (HP-71 B)	{	3.142	←	3.1415	→	3.142	}	like DEC VAX HP-15C
		3.140	←	3.1405	→	3.141		
		3.140	←	3.1395	→	3.140		
		3.138	←	3.1385	→	3.139		
		3.138	←	3.1375	→	3.138		

Error in floating-point operations

Let $\odot \in \{+, -, \times, /\}$

Let A be name of variable whose value is a ,
 B b ,
 C c .

Program statement " $A := B \odot C$ "

produces value $a = b \odot c$ ROUNDED

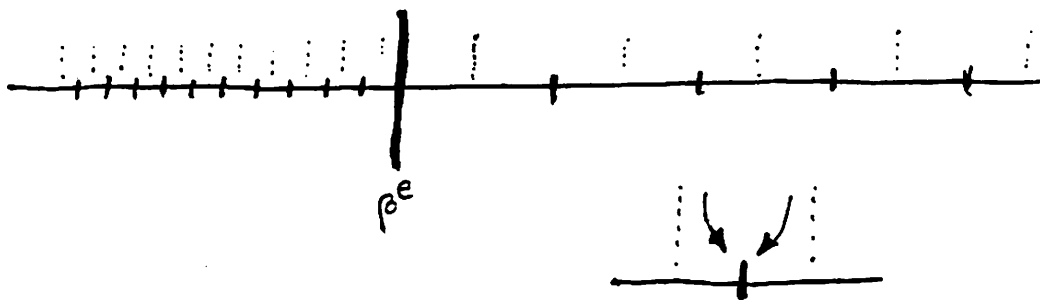
$$b \odot c = \pm x.x \times x \dots x \times x \times x \times x \dots$$

$$a = \pm x.x \times x \dots x \times x$$

$\longleftarrow p \longrightarrow$

"CORRECT" ROUNDING $|a - b \odot c| \leq \frac{1}{2} \text{ulp}(b \odot c)$

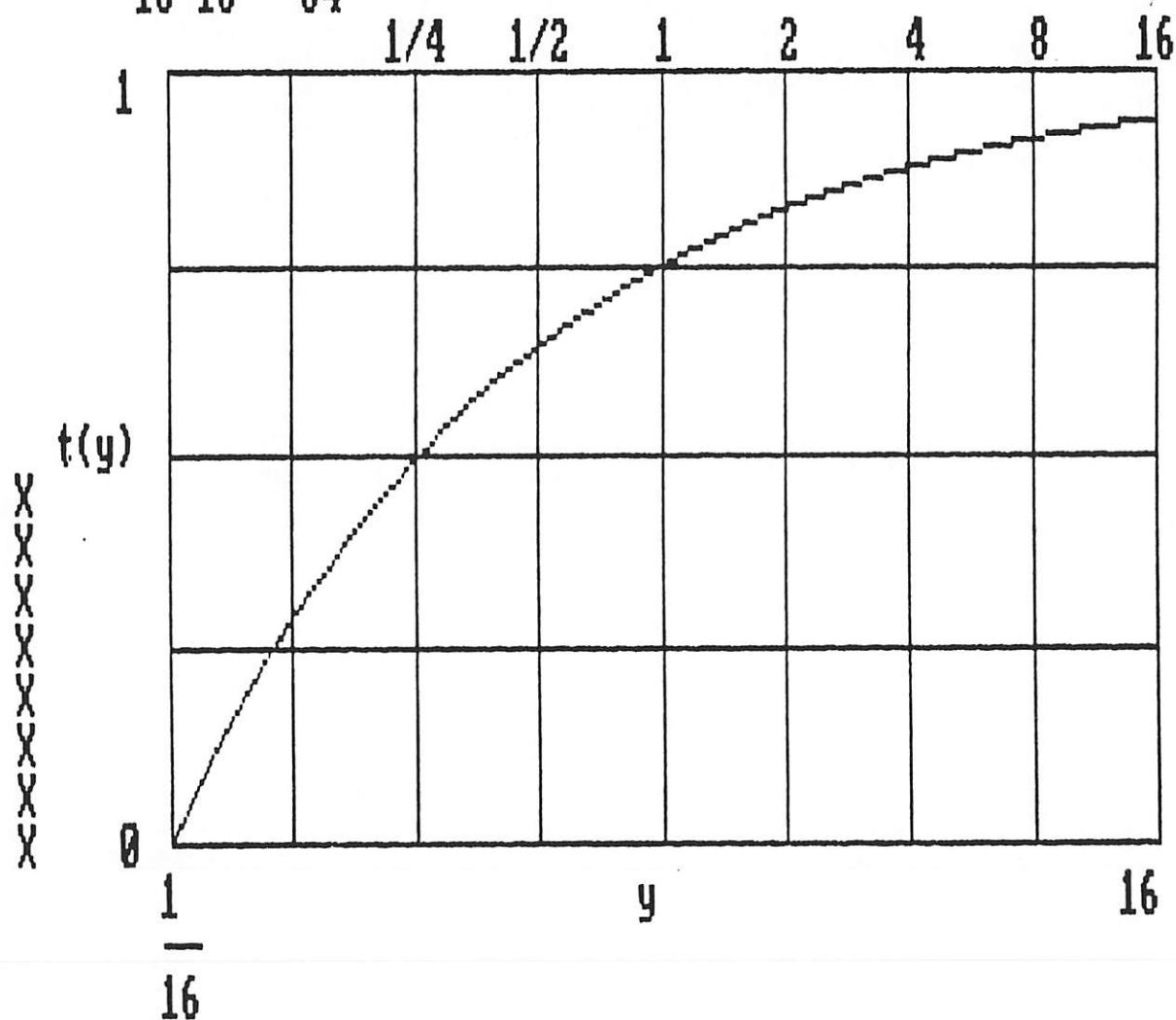
except $\begin{cases} a = b \odot c & \text{if it is representable exactly,} \\ a = \text{something else} & \text{if over/underflow occurs} \end{cases}$



What if $b \odot c$ falls midway between two adjacent representable values?

$$y := \frac{1}{16}, \frac{1}{16} + \frac{1}{64} \dots 16$$

$$t(y) := 1 - \frac{0.25}{\sqrt{y}}$$



--> y is on a
Logarithmic scale.