

Lecture 2 : ~~2/20/04~~

S.A. Treter, "Comar System Design using DSP Algorithms,"
Kluwer 2003, New York

1st Lab due Tues, Jan Feb 18 of ~~Thurs~~, Jan 7/9

Finite # Systems

Finite Number Systems

integer = 32-bit word

$$\longleftrightarrow (d_{31} \dots d_1 d_0)$$

MSB LSB

$$0 \leq V = \sum_{n=0}^{31} d_n 2^n \leq 2^{32} - 1$$

(2's complement) signed integer : MSB is sign bit :

$$: d_{31} = 0 \quad +$$

$$: d_{31} = 1 \quad -$$

$$-2^{31} \leq V = -d_{31} 2^{31} + \sum_{n=0}^{30} d_n 2^n \leq 2^{31} - 1$$

Q_L number : move radix pt. from
rt of d₀ to any location
: programmer keeps track

$$Q_L \# \longleftrightarrow (d_{31} \dots d_L \cdot d_{L-1} \dots d_0)$$

$$V = 2^{-L} \left(-d_{31} 2^{31} + \sum_{n=0}^{30} d_n 2^n \right)$$

↑ shift left by L

NB : "integer" if this type occupy
one 32-bit register

IEEE Floating Pt. (single & dbl precision)

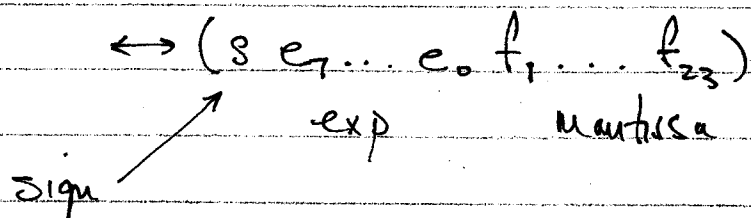
normal #s

Subnormal (< smallest pos normal)

NaNs

∞ #s

fp (single) = 32-bit word



~~Let's~~
Let's define some intermediate vars:

$$e = \sum_{k=0}^7 e_k 2^k \quad f = \sum_{k=1}^{23} f_k 2^{-k} \quad s \in \{0, 1\}$$

$$0 \leq e \leq 255 \quad \& \quad 0 \leq f \leq 1 - 2^{-23}$$

(Convention or code) 1. if $e = 255$ (11...1) & $f \neq 0$, then $x = \text{NaN}$, indep of s

(Code or code for $\pm \infty$) 2. if $e = 255$ (11...1) & $f = 0$ (0...0), then $x = (-1)^s \infty$ [$\pm \infty$]

3. if $0 < e < 255$, then x is a normal # with value

NB: $2^{-23} \approx 0.0000001192$ (3)

$$X = (-1)^s 2^{e-127} (1+f)$$

s sign
 $e-127$ exponent
 f fraction
 $\sum_{k=0}^{23} f_k 2^{-k}$
 $f_k = 1$

radix
 rt by 127 places or
 lft by 126 places

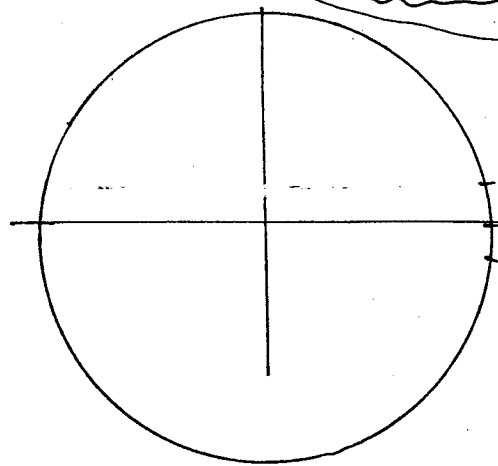
Insert Here

NB: The subnormal #'s transition btwn 0 and the smallest minimum zone value in the same increments as the resolution of the smallest normal # zone. Systems that do not support subnormal numbers round these values to zero. This is known as a flush to zero operation.

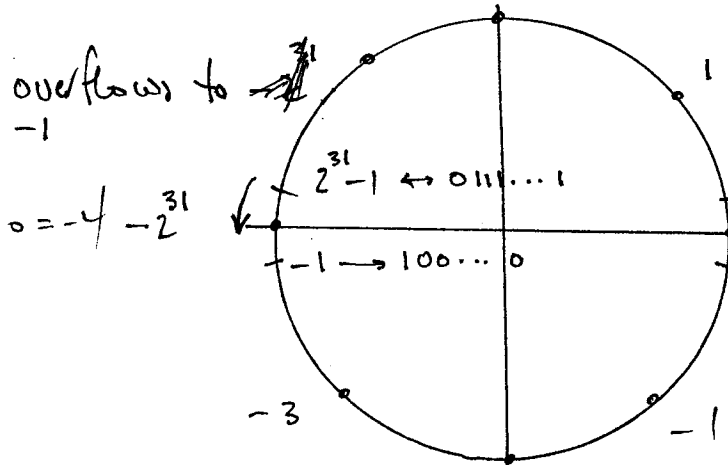
4. $e = 0 \ \& \ f \neq 0$, then $X = (-1)^s 2^{-126} f$
 (Subnormal) \Rightarrow smallest positive number $(2^{-126})(2^{-23}) = 2^{-149} \approx 1.04 \times 10^{-45}$

5. $e = 0 \ \& \ f = 0$, then $X = (-1)^s \cdot 0 = \pm 0$

Overflow



$0, 2^{32}$
 $2^{32} - 1$
 $11\dots1 + 0\dots01$
 overflows to zero



overflows to -1

$0 = -4 - 2^{31}$

$2^{31} - 1 \leftrightarrow 011\dots1$
 $-1 \leftrightarrow 100\dots0$

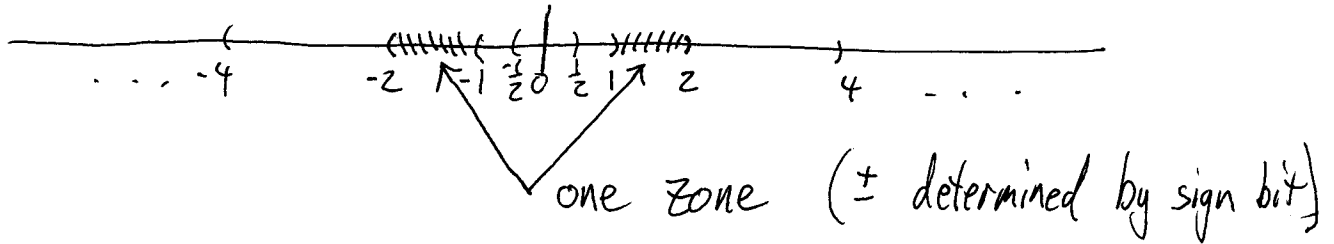
$-2 \leftrightarrow 11\dots1 = 2 - 1$

$$\begin{array}{r}
 011 \\
 + 100 \\
 \hline
 111
 \end{array}
 \quad
 \begin{array}{r}
 3 \\
 - 4 \\
 \hline
 -1 = -2 + 3
 \end{array}$$

$11\dots1 = 2^{32} - 1$

IEEE Floating-Point (Insert)

We can think of the set of real numbers being split into zones



$$0 < e < 255 \quad (\text{NB: } e=0, 255 \text{ are reserved})$$

$$0 \leq f \leq 1 - 2^{-23} \Rightarrow \text{zone resolution } 2^{-23} \approx 0.0000001192$$

If z^a is the minimum zone ^{value} resolution (associated with $f=0$), then the zone resolution is

$$(2^{-23}) z^a \quad (\text{NB: } a = e - 127)$$

(Comment: $e=127$ for the example above and the minimum zone value is one in the example above)

The smallest normal number ^(and minimum resolution value) is $2^{-126} \approx 1.175 \times 10^{-38}$ and the zone resolution is $(2^{-23})(2^{-126}) \approx 1.04 \times 10^{-45}$

The largest minimum resolution value is $2^{127} \approx 1.7 \times 10^{38}$ and the zone resolution is $(2^{-23})(2^{127}) = 2^{104} \approx 2.03 \times 10^{31}$

There are always 2^{23} pts of resolution in a given zone and the larger the minimum zone value is, the larger the jumps b/w number is. Hence, worse resolution.

This means that for

really small #'s 0.00000000 X XXXXXX

small #'s X.XXXXXX

large #'s XXXXX.XX

really large #'s XXXXXXX 00000000.00

not necessarily zero, but will be in quantized levels and 9

NB: $2^{-23} \approx 0.000001192$

Have ^{accurate} resolution up to "6 decimal places" (or seven digits) ... see small #'s above (short and int have 5 decimal place resolution, respectively)
max short = 32768
max int = 2147483648

An aside

A number ~~where~~ where the number of digits after the decimal point is fixed is known a fixed-point #.

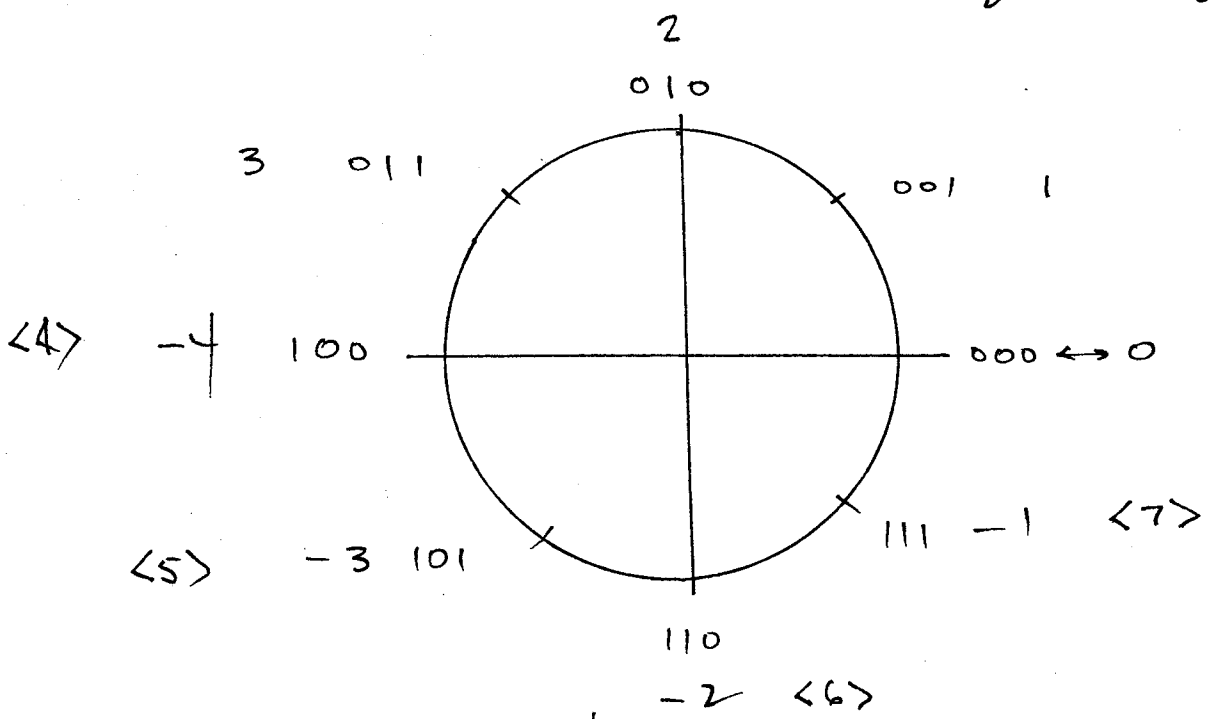
It is important to understand that fixed-point arithmetic is not limited to integer values (remember Qc formats). However, a fixed-point number may be scaled to an integer. In fact, this is required for implementation processing. Upon process completion, the fixed-point number may be unscaled. This point will be clarified in the lectures for Lab 4: FIR Filter Design and Implementation.

Hopefully this explains the term "fixed-point approximation" to π , namely PI = 3.14159265359

2's Complement & Overflow

Illustrate for 3-bit word ($d_2 d_1 d_0$)

$$V = -d_2 \cdot 2^2 + d_0 + d_1 \cdot 2$$



Let's do binary addition in usual way

100 -4
 011 3
 100 -4

 011 3
 (OVERFLOW)

011 3
 010 2

 100 -4

 001 1 (OK)

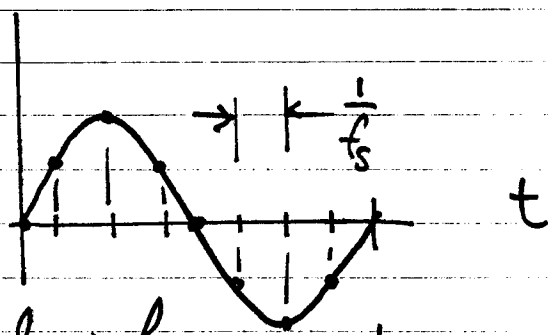
intermediate overflow to 101 = -3

Adds move you CCW thru virtual < > \neq
 Subs move you CW thru virtual < >. Intermediate
 reads can be bad but final reads OK.

Remarks on Sampling

Consider a sine wave

$$X(t) = \sin 2\pi f_0 t$$



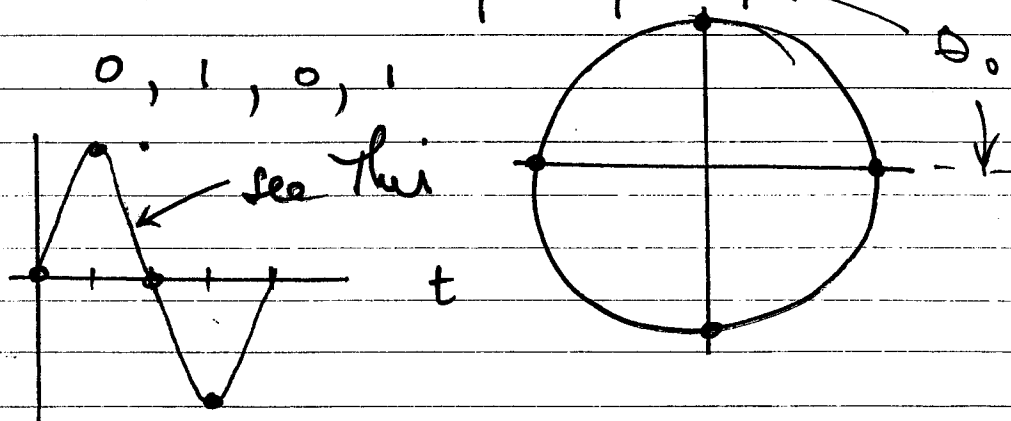
If we sample @ rate f_s , then $\frac{1}{f_0}$

$$X[n] = \sin 2\pi \frac{f_0}{f_s} n = X(t = \frac{n}{f_s}) = X(t = nt_s)$$

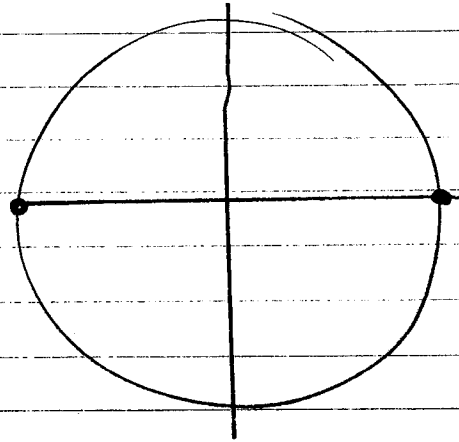
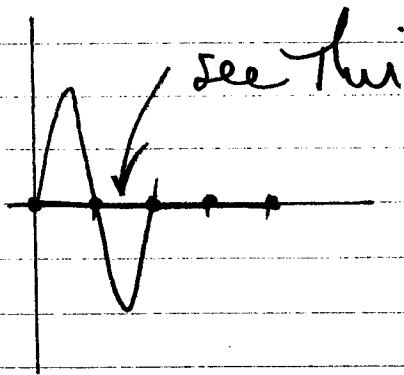
If $f_0 = \frac{1}{4} f_s$, then $(\frac{1}{f_0} = \frac{4}{f_s})$

$$X[n] = \sin 2\pi \frac{1}{4} n \quad ; \quad \theta_0 = \frac{2\pi}{4}$$

And we get 4 samples per period

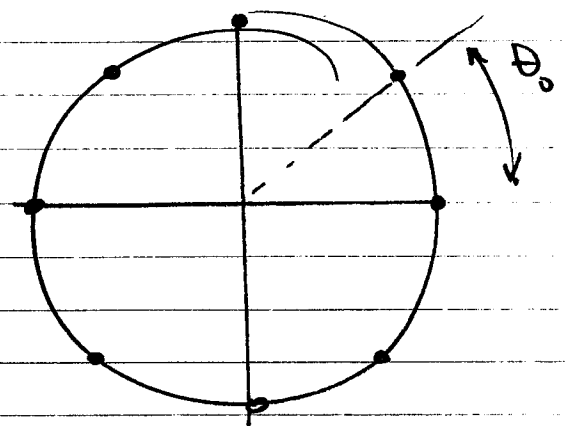
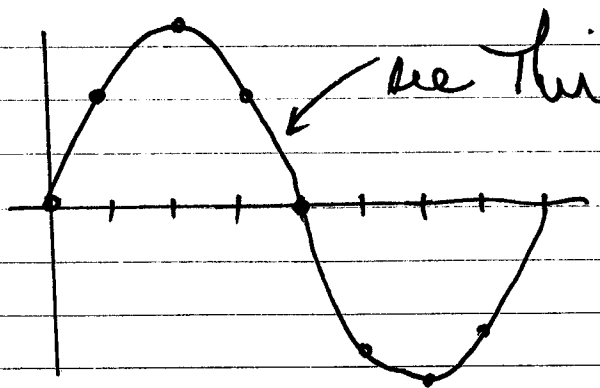


If $f_0 = f_s / 2$, Then $\frac{1}{f_0} = \frac{2}{f_s}$ And we get $x[n] = \sin 2\pi \frac{1}{2} n$ for 2 samples per period ; $\theta_0 = \frac{2\pi}{2}$



If $f_0 = f_s / 8$, Then $\frac{1}{f_0} = \frac{8}{f_s}$ And

$x[n] = \sin 2\pi \frac{1}{8} n$ And 8 samples per



$$\theta_0 = \frac{2\pi}{8}$$

You will see dots in your code
↳ contin curve on your oscilloscope

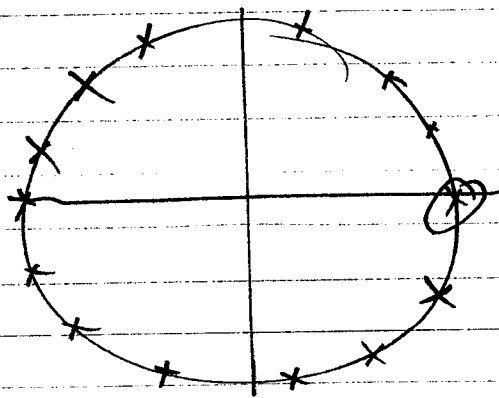
If $f_0 = \frac{p}{q} f_s$, then

$$\theta_0 = \frac{p}{q}$$

and we go q samples to get $p 2\pi$

$$\sin 2\pi \frac{p}{q} n$$

$$p/q = \frac{7}{16}$$



If $f_0 = \text{irr } f_s$, then samples never
overwrite.

With "PI = 3.14159..." we don't actually
overwrite the way the theory says.