Chapter 3 – Fixed-Point vs Floating-Point Processing

Fixed-point and floating-point refer to the two basic kinds of number systems used in computers. There are two concepts to consider:

- numerical representation, i.e., how the bits in a datum are interpreted
- arithmetic operations on the data

The c2000 family includes processors that are fixed-point and processors that are floating-point.



- integers are 16 bits on c2000
- long integers are 32 bits on c2000
- typically can achieve faster clock rate than floating-point
- typically processor has lower power consumption than floating-point
- typically processor is cheaper than floating-point processor
- programmer must be aware of limitations and may need to take steps to accommodate them:
 - o introduce scaling
 - o turn on "saturation mode"
 - implement "block floating-point"

2. Floating-Point Processing

- numbers are real
- floating-point representation
- arithmetic is done on real numbers
- regular floating-point numbers are <u>32 bits</u> on c2000¹
- floating-point has large dynamic range
- typically processor has higher power consumption than fixed-point
- typically processor is more expensive than fixed-point processor
- programmer has less concern about range limitations

¹ c2000 also supports "double" = 64 bits

Fixed-Point Applications

- good for "bit-exact" applications
- good for "control" parts of code:
 - o if, then, else decisions
 - o uses less memory than floating-point

Floating-Point Applications

• good for applications in which the data have a large dynamic range

Floating-Point Number in Memory Single-Precision 32-Bit **float** IEEE 754 Standard Format



Fixed-Point vs Floating-Point Processing

These are the bits in the number -1.00	A000 🔹							
10111111 1000000 0000000 0000000								
	-							



These a: 0100000	re t] 0 01(he 010	bits 101	in 0101	the 010	nur 1 01	ibe: .01	r 3. 0101	33:	333	33	*
📓 float.dat												
Offset(h) 00000000	00 01 55 55	02 55	03 04 40	05 0	6 07)	08 09	OA	OB OC	OD	OE	OF	UUU0

These are the bits in the number 1.875e-010 00101111 01001110 00101000 10001111



*

.





(Note: Neither of these are a "real" zero as per the standard S1.M x 2^{E-127} .)



"Not a Number" (NaN) e.g. 0/0 or x/0, i.e., result of divide-by-zero



Given the following 32-bit binary number written in standard floating-point format, what decimal number does it represent?

1100 0001 1010 0000 0000 0000 0000 0000

Given the following decimal number, write it as a binary number in standard floatingpoint format.

0.40625

Fixed-Point Saturation (assume 8-bit 2's complement numbers)



Note that fixed-point addition can potentially overflow, but fixed-point multiplication will not. Note that floating-point addition can also potentially overflow, but it is much more unlikely.

Example of saturation - high-voltage sine-wave output signal



Some Examples of Fixed- and Floating-Point Additions and Multiplications

Fixed-Point Add Operation – Integer Interpretation



, look up Q format in book or Internet

Fixed-Point Add Operation – Q7 Interpretation



Fixed-Point Multiply Operation – Integer Interpretation



Fixed-Point Multiply Operation – Q7 Interpretation



Floating-Point Add Operation



Floating-Point Multiply Operation

