Floating Point

Arithmetic

(The IEEE Standard)

Floating Point Arithmetic (and The IEEE Standard)

* Floating Point Arithmetic

- Representations
- Issues
- Normalized, Unnormalized, Subnormal
- Precision
- Wobble

★ The IEEE Standard

- Why
- What it contains, what it doesn't contain
- Formats
- Rounding
- Operations
- Infinities, NANs
- Exceptions

- Traps

Several Issues Come Up:

- How many bits for range, how many bits for precision?
- ★ What to do with numbers too small to represent with this scheme?
- What to do with numbers that do not correspond to exact representations?
- ★ What to do with numbers too large to be represented?
- Shall we distinguish numbers too large with true infinities?
- What about nonsense numbers? (Examples:

Arcsin 2, $\frac{0}{0}$, $\infty - \infty$)

First, An Example



In DEC format: (-1)^s * 0.1 fra * 2^{EXP-4}



First, Some General Stuff:

A number can be represented as

$$\pm d_0 \cdot d_1 d_2 \dots \beta^e$$

These numbers correspond to points on the real line. If we insist that all representations be normalized, then the points are shown (normalized can mean: $d_0 = \emptyset$, $d_1 = 1$)



(We can, incidentally, store the number in signed-magnitude format:)

s e + BIAS
$$d_2 d_3 d_4 \dots$$

SIGN +, -

Normalized, Unnormalized, Subnormal Again, we are looking at $\pm d_0.d_1d_2...*\beta^e$ 1. If it is <u>normalized</u>, it is:

$$\pm$$
 0.1 d₂ d₃ ··· * β^{e}

2. <u>Unnormalized</u> (after a subtract of like signs, for example)

 $\pm 0.0001 d_2 d_3 \dots * \beta^{e+3}$

- 3. <u>Subnormal</u> means it can't be represented in the machine in normalized format
 - Recall the format ± e+BIAS d2 d3 ...

Corresponds to $\pm 0.1 d_2 d_3 \dots \star \beta^e$

 Suppose we successively divide by β. We can do this <u>until</u> e+BIAS = 1. Below that we can't represent numbers (except 0). Why? Suppose we let e+BIAS = 0. How do we now represent 0?

Precision



The IEEE Standard

Reasons:

1. Direct Support for:

- Execution-time diagnosis of anomalies
- Smoother handling of exceptions
- Interval arithmetic at reasonable cost
- 2. Provide for development of:
 - Standard elementary functions
 - Very high precision arithmetic
 - Coupling of numeric & symbolic computation

The IEEE Standard (Continued)

What does it contain:

- Formats: single, double, extended
- Operations: +,-,*, +, , , REM,CMP
- Rounding modes
- Conversion: Int/FI., Dec/FI., FI/FI
- Exceptions: Underflow, Overflow, Div Ø, Inexact, Invalid

What it does not contain:

- Requirements for implementation in HDWR or SFWR
- Interpretation of NaNs
- Formats for Integers, BCD
- Conversions other than above

The Formats

There are four; we start with one as an example.

Single

Representable Numbers:

* Normalized

1.d₁d₂d₃...d₂₃ * 2^e

where $-126 \le e \le +127$



Note: The range of exponents

-126 <u><</u> e <u><</u> +127

Coupled with the BIAS (127) which is added to the exponent yields an 8 bit string from 00000001

11111110

Two strings remain: 0000000, 11111111

★ Subnormal numbers (Exp field = 0000000)

$$0.d_1d_2...d_{23}$$
 2^{-126}

* Infinities (Exp field = 11111111) s 11111111 000 ... 0

Formats (Continued)

That still leaves those strings characterized:

s 11111111 Not Zero

These are defined as NaNs.

They result from invalid operations (Like,: $\frac{0}{0}$, $\frac{\infty}{\infty}$, $\infty - \infty$)

Generalizing to the other formats <u>Single Single-X Double Double-X</u>

Precision	24 bits	≥32	53	≥64
Exponent	8 bits	≥11	11	≥15
Word Length	32 bits	≥ 4.3	64	≥79
Exp BIAS	+127		+1023	
e _{max}	+127	≥1023	+1023	≥16382
e _{min}	-126	≤ -1022	-1022	<mark>≤ -16382</mark>

Rounding

- 1st We perform the <u>operation</u> & produce the infinitely precise result
- 2nd We round to fit it into the destination format

Four Rounding Modes

- <u>Default:</u> To nearest. If equally near, then to the one having A Ø in LSB
- 2. Directed roundings
 - Toward + ∞
 - Toward ∞
 - Toward Ø (Chop)

Operations

★ Arithmetic: +, -, *, +, REM
When y ≠ Ø, r = x REM y, is defined:
r = x-y*n, where n is the integer
nearest X/y

whenever $\left| \begin{array}{c} n - \frac{x}{y} \right| = \frac{1}{2}$, then n is EVEN

∴ Remainder is always exact

- ★ Square root: Result defined if ARG ≥ Ø.
- Conversion from one format to another
 - To fewer bits: rounded
 - To more bits: exact

Operations(Continued)

- Conversion FI. Pt. <---> Integers
 Binary <---> Decimal
- ★ Comparison
 - Always exact
 - Never underflow, overflow
 - Four relations are possible {>, =, <, unordered}
- Note: Invalid is signaled if unordered operands are compared and unordered is not the basis but > or < is the basis.

Examples:

Predicate	>	<	=	?	unordered
=	F	F	Τ	F	No
?≠	Т	Т	F	Т	No
>	Т	F	F	F	Yes
?<≠	F	Т	Т	Ţ	No

Infinities, NaNs, ±Ø

- ∞:
 - ★ ∞ < (finite) < + ∞
 - ***** Arithmetic on ∞ is <u>exact</u>
 - $\star \infty$ is created by
 - Overflow
 - "Divide by zero"
- NaN:
 - ★ Signaling & Quiet

Signaling - Reserved operand that signals the invalid Op. Exception for all operations in the standard. If no trap occurs, a quiet NaN is delivered

<u>Quiet</u> - Operations on quiet NaNs produce quiet NaNs. They provide hooks to retrospective diagnostic information.

Exceptions

When detected: Take Trap, or Set Flag, or Both

Flag can be reset <u>only</u> under program control

- * Invalid
 - Operation on a signaling NaN.
 - · ∞ ∞ 0/0
 - **0** ∗ ∞ ∞/∞
 - x REM y, where y=0 or x= ∞
 - **VNEG**
 - Conversion from FI. to int. or decimal, when overflow, infinity, or NaN prevents the conversion
 - Comparison via predicates involving > or <, and Not?, when the operands are unordered

Exceptions (Continued)

* Divide by zero

When f(finite) --> Infinite and exact

★ <u>Overflow</u>

When the destinations largest finite number is exceeded by what would have been the <u>rounded</u> floating point result if the exponent range were unbounded





Why?

- ***** Underflow
 - Tiny value (which could cause subsequent overflow)
 - Loss of precision

Delivered result may be zero, subnormal No., or ± 2^{min-exp}

Exceptions (Continued)

* Underflow (continued)

Trapped underflows! [All operations except conversions]

1 <u>st</u>, Multiply infinitely precise Result by 2^a

* Inexact

When the result of an operation is not exact, or on non-trapped overflow.

Traps

For any of the five exceptions, a user should be able to:

- ***** Specify a handler
- Request that an existing handler be disabled, saved, restored.

When a system traps, the trap handler should be able to determine:

- Which exception occurred on this operation
- * The kind of operation being performed
- ***** The destination format
- In overflow, underflow, & inexact, the correctly rounded result
- In invalid & divide by zero, the operand values