

Achmad Imam Kistijantoro
Anggrahita Bayu Sasmita
Yudistira Dwi Wardhana Asnar
Rahmat Mulyawan
Infall Syafalni

IF2130/II2130 – Organisasi dan Arsitektur Komputer

sumber: Greg Kesden, CMU 15-213, 2012

Representasi Informasi – Floating Point Operations

Floating Point Operations: Basic Idea

▶ $\mathbf{x} +_{\mathbf{f}} \mathbf{y} = \mathbf{Round}(\mathbf{x} + \mathbf{y})$

▶ $\mathbf{x} \times_{\mathbf{f}} \mathbf{y} = \mathbf{Round}(\mathbf{x} \times \mathbf{y})$

▶ Basic idea

- ▶ First **compute exact result**
- ▶ Make it fit into desired precision
 - ▶ Possibly overflow if exponent too large
 - ▶ Possibly **round to fit into frac**



Rounding

- ▶ Rounding Modes (illustrate with \$ rounding)

▶		\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
▶	Towards zero	\$1	\$1	\$1	\$2	-\$1
▶	Round down ($-\infty$)	\$1	\$1	\$1	\$2	-\$2
▶	Round up ($+\infty$)	\$2	\$2	\$2	\$3	-\$1
▶	Nearest Even(default)	\$1	\$2	\$2	\$2	-\$2



Closer Look at Round-To-Even

▶ Default Rounding Mode

- ▶ Hard to get any other kind without dropping into assembly
- ▶ All others are statistically biased
 - ▶ Sum of set of positive numbers will consistently be over- or under-estimated

▶ Applying to Other Decimal Places / Bit Positions

- ▶ When exactly halfway between two possible values
 - ▶ Round so that least significant digit is even
- ▶ E.g., round to nearest hundredth

1.2349999	1.23	(Less than half way)
1.2350001	1.24	(Greater than half way)
1.2350000	1.24	(Half way—round up)
1.2450000	1.24	(Half way—round down)



Rounding Binary Numbers

▶ Binary Fractional Numbers

- ▶ “Even” when least significant bit is 0
- ▶ “Half way” when bits to right of rounding position = $100..._2$

▶ Examples

- ▶ Round to nearest $1/4$ (2 bits right of binary point)

Value Value	Binary	Rounded	Action	Rounded
$2 \frac{3}{32}$	$10.000\textcolor{red}{11}_2$	10.00_2	($< 1/2$ —down)	2
$2 \frac{3}{16}$	$10.00\textcolor{red}{110}_2$	10.01_2	($> 1/2$ —up)	$2 \frac{1}{4}$
$2 \frac{7}{8}$	$10.11\textcolor{red}{100}_2$	11.00_2	($= 1/2$ —up)	3
$2 \frac{5}{8}$	$10.10\textcolor{red}{100}_2$	10.10_2	($= 1/2$ —down)	$2 \frac{1}{2}$



FP Multiplication

- ▶ $(-1)^{s_1} M_1 2^{E_1} \times (-1)^{s_2} M_2 2^{E_2}$
- ▶ Exact Result: $(-1)^s M 2^E$
 - ▶ Sign s : $s_1 \wedge s_2$
 - ▶ Significand M : $M_1 \times M_2$
 - ▶ Exponent E : $E_1 + E_2$
- ▶ Fixing
 - ▶ If $M \geq 2$, shift M right, increment E
 - ▶ If E out of range, overflow
 - ▶ Round M to fit **frac** precision
- ▶ Implementation
 - ▶ Biggest chore is multiplying significands



Floating Point Addition

▶ $(-1)^{s1} M1 2^{E1} + (-1)^{s2} M2 2^{E2}$

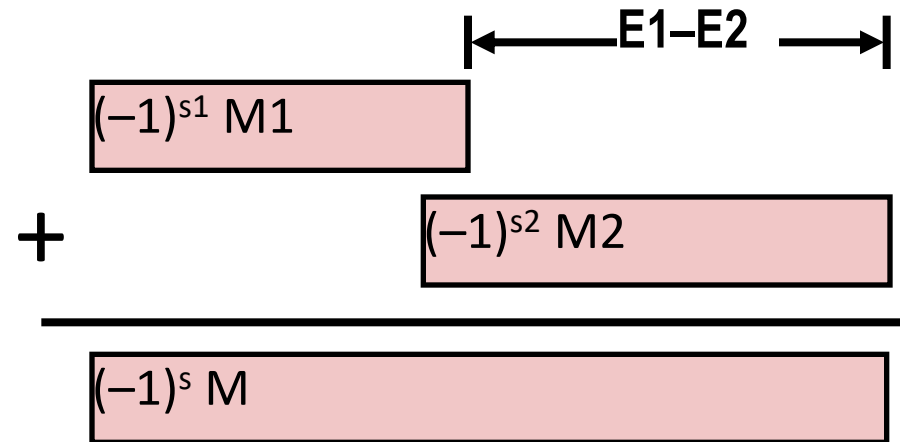
▶ Assume $E1 > E2$

▶ Exact Result: $(-1)^s M 2^E$

▶ Sign s , significand M :

▶ Result of signed align & add

▶ Exponent E : $E1$



▶ Fixing

▶ If $M \geq 2$, shift M right, increment E

▶ if $M < 1$, shift M left k positions, decrement E by k

▶ Overflow if E out of range

▶ Round M to fit **frac** precision

Floating Point in C

- ▶ **C Guarantees Two Levels**

- ▶ **float** single precision
- ▶ **double** double precision

- ▶ **Conversions/Casting**

- ▶ Casting between **int**, **float**, and **double** changes bit representation
- ▶ **double/float** → **int**
 - ▶ Truncates fractional part
 - ▶ Like rounding toward zero
 - ▶ Not defined when out of range or NaN: Generally sets to TMin
- ▶ **int** → **double**
 - ▶ Exact conversion, as long as **int** has ≤ 53 bit word size
- ▶ **int** → **float**
 - ▶ Will round according to rounding mode



Summary

- ▶ IEEE Floating Point has clear mathematical properties
- ▶ Represents numbers of form $M \times 2^E$
- ▶ One can reason about operations independent of implementation
 - ▶ As if computed with perfect precision and then rounded
- ▶ Not the same as real arithmetic
 - ▶ Violates associativity/distributivity
 - ▶ Makes life difficult for compilers & serious numerical applications programmers



End of Segment

.

