

## Integer Arithmetic and Floating Point

Kai Shen

1

## Integer Arithmetic

- We (programmers) understand the semantics supported by current system (primarily hardware, may also involve compiler)
  - Main issue: limited by the data type size
  - Can there be other semantics? Yes, in your own hardware design and implementation.
- We do not discuss hardware implementation (leave that to the hardware engineers), but it doesn't hurt to know some implementation issues
  - Unsigned and signed (2's complement) integers are often computed in the same way, so their hardware implementation can share some components

2

## Unsigned Addition

Operands:  $w$  bits

True sum:  $w+1$  bits

Discard carry:  $w$  bits

$$\begin{array}{r}
 u \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}} \\
 + v \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}} \\
 \hline
 u+v \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}} \\
 \hline
 \text{UAdd}_w(u, v) \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}}
 \end{array}$$

- Semantics: standard addition, but ignore overflowed carry
  - Still commutative and associative
- Implements modular arithmetic
 
$$s = \text{UAdd}_w(u, v) = u + v \bmod 2^w$$

$$\text{UAdd}_w(u, v) = \begin{cases} u + v & u + v < 2^w \\ u + v - 2^w & u + v \geq 2^w \end{cases}$$

3

## Two's Complement Addition

Operands:  $w$  bits

True sum:  $w+1$  bits

Discard carry:  $w$  bits

$$\begin{array}{r}
 u \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}} \\
 + v \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}} \\
 \hline
 u+v \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}} \\
 \hline
 \text{TAdd}_w(u, v) \quad \boxed{\phantom{0}} \boxed{\phantom{0}} \dots \boxed{\phantom{0}} \boxed{\phantom{0}}
 \end{array}$$

- TAdd and UAdd have identical bit-level computation
  - Signed vs. unsigned addition in C:
 

```
int s, t, u, v;
s = (int) ((unsigned) u + (unsigned) v);
t = u + v;
Will give s == t
```

4

### TAdd Overflow

- Functionality
  - True sum requires  $w+1$  bits
  - Drop off MSB
  - Treat remaining bits as 2's comp. integer

True Sum

$0111\dots1$   $2^w-1$

$0100\dots0$   $2^{w-1}$

$0000\dots0$   $0$

$1011\dots1$   $-2^{w-1}-1$

$1000\dots0$   $-2^w$

PosOver

NegOver

TAdd Result

$011\dots1$

$000\dots0$

$100\dots0$

$$TAdd_w(u,v) = \begin{cases} u+v+2^w & u+v < TMin_w \text{ (NegOver)} \\ u+v & TMin_w \leq u+v \leq TMax_w \\ u+v-2^w & TMax_w < u+v \text{ (PosOver)} \end{cases}$$

5

### Negation: Complement & Increment

- Semantics of negation ( $-x$ )
  - Only meaningful for signed integer
  - TMIN** is smallest negative integer, what is  $-TMIN$ ?
- Claim: for 2's complement
  - $\sim x + 1 = -x$
- Complement
  - Observation:  $\sim x + x = 1111\dots111 = -1$

$$\begin{array}{r} x \quad 100011101 \\ + \sim x \quad 01100010 \\ \hline -1 \quad 11111111 \end{array}$$

6

### Unsigned Multiplication

Operands:  $w$  bits

$u$ 

$\square \square \dots \square$

$*$ 
 $v$ 

$\square \square \dots \square$

---

True product:  $2*w$  bits

$u \cdot v$ 

$\square \square \dots \square \square \dots \square$

---

Discard  $w$  bits:  $w$  bits

$UMult_w(u, v)$ 

$\square \square \dots \square$

- Semantics: standard multiplication, but ignore high order  $w$  bits
- Implements modular arithmetic
  - $UMult_w(u, v) = u \cdot v \text{ mod } 2^w$

7

### Signed Multiplication

- Under 2's complement, same bit-level computation as in unsigned case

8

### Code Security Example

```

void* copy_elements(void *ele_src[], int ele_cnt, size_t ele_size) {
    /*
     * Allocate buffer for ele_cnt objects, each of ele_size bytes
     * and copy from locations designated by ele_src
     */
    void *result = malloc(ele_cnt * ele_size);
    if (result == NULL)
        /* malloc failed */
        return NULL;
    void *next = result;
    int i;
    for (i = 0; i < ele_cnt; i++) {
        /* Copy object i to destination */
        memcpy(next, ele_src[i], ele_size);
        /* Move pointer to next memory region */
        next += ele_size;
    }
    return result;
}
    
```

9

### Power-of-2 Multiply with Shift

- Multiply is slow on most machines
- Operation
  - $u \ll k$  gives  $u * 2^k$
  - Both signed and unsigned

Operands:  $w$  bits

True product:  $w+k$  bits

Discard  $k$  bits:  $w$  bits

$UMult_u(u, 2^k)$   
 $TMult_u(u, 2^k)$

- Examples
  - $u \ll 3$  ==  $u * 8$
  - $u \ll 5 - u \ll 3$  ==  $u * 24$

10

### Compiled Multiplication Code

C Function

```

int mull2(int x)
{
    return x*12;
}
    
```

Compiled Arithmetic Operations	Explanation
<pre> leal (%eax,%eax,2), %eax sall \$2, %eax                     </pre>	<pre> t = x+x*2 return t &lt;&lt; 2;                     </pre>

- C compiler automatically generates shift/add code when multiplying by constant

11

### Division

- Integer division: divide one integer over another, output an integer
- Semantics:
  - Like standard division
  - No overflow problem (except divide by zero)
  - Round toward zero (round down on positive side, round up on negative side)
- Implementation for signed/unsigned division is very different
- Division is slower than multiply, so converting to shift etc. will help even more

12

### Unsigned Power-of-2 Divide with Shift

- Quotient of unsigned by power of 2
  - $u \gg k$  gives  $\lfloor u / 2^k \rfloor$
  - Uses logical shift

Operands:

/  $2^k$

Division:

Result:

$u$

$u / 2^k$

$\lfloor u / 2^k \rfloor$

Binary Point

	Division	Computed	Hex	Binary
$x$	15213	15213	3B 6D	00111011 01101101
$x \gg 1$	7606.5	7606	1D B6	00011101 10110110
$x \gg 4$	950.8125	950	03 B6	00000011 10110110
$x \gg 8$	59.4257813	59	00 3B	00000000 00111011

13

### Signed Power-of-2 Divide with Shift

- Quotient of signed by power of 2
  - $x \gg k$  gives  $\lfloor x / 2^k \rfloor$
  - Uses arithmetic shift
  - Rounds wrong direction when  $u < 0$

Operands:

/  $2^k$

Division:

Result:

$x$

$x / 2^k$

RoundDown( $x / 2^k$ )

Binary Point

	Division	Computed	Hex	Binary
$y$	-15213	-15213	C4 93	11000100 10010011
$y \gg 1$	-7606.5	-7607	E2 49	11100010 01001001
$y \gg 4$	-950.8125	-951	FC 49	11111100 01001001
$y \gg 8$	-59.4257813	-60	FF C4	11111111 11000100

14

### Correct Power-of-2 Divide

- Quotient of negative number by power of 2
  - Want  $\lceil x / 2^k \rceil$  (Round Toward 0)
  - Compute as  $\lfloor (x+2^k-1) / 2^k \rfloor$ 
    - In C:  $(x + (1 \ll k) - 1) \gg k$

15

### Integer C Puzzles

```
Initialization
int x = foo();
int y = bar();
unsigned ux = x;
unsigned uy = y;
```

- $ux \geq 0$
- $ux > -1$
- $x > 0 \ \&\& \ y > 0 \Rightarrow x + y > 0$
- $ux \gg 3 == ux/8$
- $x \gg 3 == x/8$

16

### Floating Point

- Background: fractional binary numbers
- Bit representation: IEEE floating point standard
- Rounding, addition, multiplication

17

### Fractional Binary Numbers

What is  $1011.101_2$ ?

- Representation
  - Bits to right of "binary point" represent fractional powers of 2
  - Represents rational number:  $\sum_{k=j}^i b_k \times 2^k$

18

### Fractional Binary Numbers: Examples

Value	Representation
5 & 3/4	$101.11_2$
2 & 7/8	$10.111_2$
1 & 7/16	$1.0111_2$

- Observations
  - Divide by 2 by shifting right
  - Multiply by 2 by shifting left
  - Numbers of form  $0.111111\dots_2$  are just below 1.0
    - $1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0$
    - Use notation  $1.0 - \epsilon$

19

### Representable Numbers

- Limitation
  - Can only exactly represent numbers of the form  $x/2^k$
  - Other rational numbers have repeating bit representations
- Value Representation
  - 1/3       $0.0101010101[01]\dots_2$
  - 1/5       $0.001100110011[0011]\dots_2$
  - 1/10      $0.0001100110011[0011]\dots_2$

20

### Bits Representation of Fractional Binary Numbers

- How?
  - Think about the limited number of bits we have

21


### IEEE Floating Point

- IEEE Standard 754
  - Established in 1985 as uniform standard for floating point representation
    - Before that, many idiosyncratic formats
  - Supported by all major CPUs
- Driven by numerical concerns
  - Nice standards for rounding, overflow, underflow
  - Hard to make fast in hardware
    - Numerical analysts predominated over hardware designers in defining standard

22

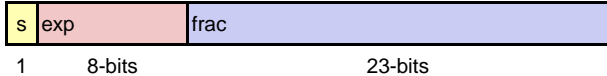

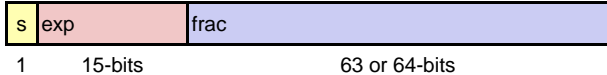
### Floating Point Representation

- Numerical form:  $(-1)^s M 2^E$ 
  - Sign bit  $s$**  determines whether number is negative or positive
  - Significand  $M$**  normally a fractional value in range [1.0,2.0).
  - Exponent  $E$**  weights value by power of two
- Encoding
  - MSB  $s$  is sign bit  $s$
  - exp field encodes  $E$  (but is not equal to  $E$ )
  - frac field encodes  $M$  (but is not equal to  $M$ )



23

### Precisions

- Single precision: 32 bits
 
- Double precision: 64 bits
 
- Extended precision: 80 bits (Intel only)
 

24

### Normalized Values

A diagram showing a floating-point format with three fields: a single-bit sign field 's', an exponent field 'exp', and a fraction field 'frac'.

- Condition:  $\text{exp} \neq 000\dots 0$  and  $\text{exp} \neq 111\dots 1$
- Exponent coded as **biased** value:  $E = \text{Exp} - \text{Bias}$ 
  - $\text{Exp}$ : unsigned value exp
  - $\text{Bias} = 2^{k-1} - 1$ , where  $k$  is number of exponent bits
    - Single precision: 127 (Exp: 1...254, E: -126...127)
    - Double precision: 1023 (Exp: 1...2046, E: -1022...1023)
- Significand coded with implied leading 1:  $M = 1.\text{xxx}\dots\text{x}_2$ 
  - xxx...x: bits of frac
  - Minimum when  $000\dots 0$  ( $M = 1.0$ )
  - Maximum when  $111\dots 1$  ( $M = 2.0 - \epsilon$ )
  - Get extra leading bit for "free"

25

### Normalized Encoding Example

- float  $F = 15213.0$ ;
  - $15213_{10} = 11101101101101_2$   
 $= 1.1101101101101_2 \times 2^{13}$
- Significand
  - $M = 1.1101101101101_2$
  - frac = 11011011011010000000000<sub>2</sub>
- Exponent
  - $E = 13$
  - $\text{Bias} = 127$
  - $\text{Exp} = 140 = 10001100_2$
- Result:
 

0	10001100	11011011011010000000000
<b>s</b>	<b>exp</b>	<b>frac</b>

26

### Denormalized Values

A diagram showing a floating-point format with three fields: a single-bit sign field 's', an exponent field 'exp', and a fraction field 'frac'.

- Condition:  $\text{exp} = 000\dots 0$
- Significand coded with implied leading 0:  $M = 0.\text{xxx}\dots\text{x}_2$ 
  - Equispaced leading to 0
- Exponent value:  $E = 1 - \text{Bias}$  (instead of  $E = 0 - \text{Bias}$ )
  - Smooth transition between largest value of  $\text{exp} = 00\dots 00$  and smallest value of  $\text{exp} = 00\dots 01$
- Special case:  $\text{exp} = 000\dots 0$ ,  $\text{frac} = 000\dots 0$ 
  - Represents zero value (distinct +0 and -0)

27

### Special Values

A diagram showing a floating-point format with three fields: a single-bit sign field 's', an exponent field 'exp', and a fraction field 'frac'.

- Condition:  $\text{exp} = 111\dots 1$
- Case:  $\text{exp} = 111\dots 1$ ,  $\text{frac} = 000\dots 0$ 
  - Represents value  $\infty$  (infinity)
  - Operation that overflows
  - Both positive and negative
  - E.g.,  $1.0/0.0 = -1.0/-0.0 = +\infty$ ,  $1.0/-0.0 = -\infty$
- Case:  $\text{exp} = 111\dots 1$ ,  $\text{frac} \neq 000\dots 0$ 
  - Not-a-Number (NaN)
  - Represents case when no numeric value can be determined
  - E.g.,  $\text{sqrt}(-1)$ ,  $\infty - \infty$ ,  $\infty \times 0$

28

### Tiny Floating Point Example

s	exp	frac
1	4-bits	3-bits

- 8-bit floating point representation
  - the sign bit is in the most significant bit
  - the next four bits are the exponent, with a bias of 7
  - the last three bits are the **frac**
- Same general form as IEEE Format
  - normalized, denormalized
  - representation of 0, NaN, infinity

29

### Dynamic Range (Positive Only)

	s	exp	frac	E	Value	
Denormalized numbers	0	0000	000	-6	0	
	0	0000	001	-6	$1/8 * 1/64 = 1/512$	closest to zero
	0	0000	010	-6	$2/8 * 1/64 = 2/512$	
	...					
	0	0000	110	-6	$6/8 * 1/64 = 6/512$	
	0	0000	111	-6	$7/8 * 1/64 = 7/512$	largest denorm
Normalized numbers	0	0001	000	-6	$8/8 * 1/64 = 8/512$	smallest norm
	0	0001	001	-6	$9/8 * 1/64 = 9/512$	
	...					
	0	0110	110	-1	$14/8 * 1/2 = 14/16$	
	0	0110	111	-1	$15/8 * 1/2 = 15/16$	closest to 1 below
	0	0111	000	0	$8/8 * 1 = 1$	
	0	0111	001	0	$9/8 * 1 = 9/8$	closest to 1 above
	0	0111	010	0	$10/8 * 1 = 10/8$	
	...					
	0	1110	110	7	$14/8 * 128 = 224$	
0	1110	111	7	$15/8 * 128 = 240$	largest norm	
	0	1111	000	n/a	inf	

30

### Interesting Numbers

Description			{single, double}
	exp	frac	Numeric Value
■ Zero	00...00	00...00	0.0
■ Smallest Pos. Denorm.	00...00	00...01	$2^{-(23,52)} \times 2^{-(126,1022)}$
<ul style="list-style-type: none"> <li>■ Single <math>\approx 1.4 \times 10^{-45}</math></li> <li>■ Double <math>\approx 4.9 \times 10^{-324}</math></li> </ul>			
■ One	01...11	00...00	1.0
■ Largest Normalized	11...10	11...11	$(2.0 - \epsilon) \times 2^{(127,1023)}$
<ul style="list-style-type: none"> <li>■ Single <math>\approx 3.4 \times 10^{38}</math></li> <li>■ Double <math>\approx 1.8 \times 10^{308}</math></li> </ul>			

31

### Special Properties of Encoding

- Floating point zero same as integer zero
  - All bits = 0
- Can (almost) use unsigned integer comparison
  - Must first compare sign bits
  - Must consider -0 = 0
  - NaNs problematic
    - Will be greater than any other values
    - What should comparison yield?
  - Otherwise OK
    - Denorm vs. normalized
    - Normalized vs. infinity

32

### Floating Point Operations: Basic Idea

- $x \pm_f y = \text{Round}(x + y)$
- $x \times_f y = \text{Round}(x \times y)$
- Basic idea
  - First **compute exact result**
  - Make it fit into desired precision
    - Possibly overflow if exponent too large
    - Possibly **round to fit into frac**

33

### Rounding

- Rounding modes (illustrate with \$ rounding)
  - Towards zero
 

	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
■	\$1	\$1	\$1	\$2	-\$1
  - Round down ( $-\infty$ )
 

	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
■	\$1	\$1	\$1	\$2	-\$2
  - Round up ( $+\infty$ )
 

	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
■	\$2	\$2	\$2	\$3	-\$1
  - Nearest even (default)
 

	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
■	\$1	\$2	\$2	\$2	-\$2
- Nearest even
  - Round to nearest acceptable point
  - When exactly halfway between two adjacent points, round so that least significant digit is even
- Why?
  - Statistically unbiased, even digit is easier to represent

34

### 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

35

### Properties of FP Add

- Commutative?
- Associative?
  - Overflow and inexactness of rounding
  - $(1e20 + -1e20) + 3.14 = 3.14$
  - $1e20 + (-1e20 + 3.14) = ??$
- Monotonicity:  $a \geq b \Rightarrow a+c \geq b+c?$ 
  - Except for infinities & NaNs

36

## Floating Point Multiplication

- $(-1)^{s1} M1 2^{E1} \times (-1)^{s2} M2 2^{E2}$
- Exact result:  $(-1)^s M 2^E$ 
  - Sign  $s$ :  $s1 \wedge s2$
  - Significand  $M$ :  $M1 \times M2$
  - Exponent  $E$ :  $E1 + E2$
- 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

37

## Mathematical Properties of FP Mult

- Multiplication is commutative?
- Multiplication is associative?
  - Possibility of overflow, inexactness of rounding
- Monotonicity:  $a \geq b \ \& \ c \geq 0 \Rightarrow a*c \geq b*c$ ?
  - Except for infinities & NaNs

38

## 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**  $\rightarrow$  **int**
    - Truncates fractional part
    - Like rounding toward zero
    - Not defined when out of range or NaN: Generally sets to TMin
  - **int**  $\rightarrow$  **double**
    - Exact conversion, as long as **int** has  $\leq 53$  bit word size
  - **int**  $\rightarrow$  **float**
    - Will round according to rounding mode

39

## Disclaimer

These slides were adapted from the CMU course slides provided along with the textbook of "Computer Systems: A programmer's Perspective" by Bryant and O'Hallaron. The slides are intended for the sole purpose of teaching the computer organization course at the University of Rochester.

40