CS 33

Data Representation (Part 2)

CS33 Intro to Computer Systems

VIII-1 Copyright © 2022 Thomas W. Doeppner. All rights reserved.

4-Bit Computer Arithmetic



Signed vs. Unsigned in C

Constants

- by default are considered to be signed integers
- unsigned if have "U" as suffix
 - OU, 4294967259U

Casting

- explicit casting between signed & unsigned

```
int tx, ty;
unsigned ux, uy; // "unsigned" means "unsigned int"
tx = (int) ux;
uy = (unsigned int) ty;
```

- implicit casting also occurs via assignments and function calls

```
tx = ux;
uy = ty;
```

Casting Surprises

Expression evaluation

- if there is a mix of unsigned and signed in single expression, signed values implicitly cast to unsigned
- including comparison operations <, >, ==, <=, >=
- examples for W = 32: TMIN = -2,147,483,648, TMAX = 2,147,483,647

Constant ₁	Constant ₂	Relation	Evaluation
0	0U	==	unsigned
-1	0	<	signed
-1	0U	>	unsigned
2147483647	-2147483647-1	>	signed
2147483647U	-2147483647-1	<	unsigned
-1	-2	>	signed
(unsigned)-1	-2	>	unsigned
2147483647	2147483648U	<	unsigned
2147483647	(int)2147483648U	>	signed

Quiz 1

What is the value of (unsigned long) -1 - (long) ULONG_MAX ??? a) 0 b) -1 c) 1 d) ULONG_MAX

Sign Extension

- Task:
 - given w-bit signed integer x
 - convert it to w+k-bit integer with same value
- Rule:
 - make k copies of sign bit:

$$-X' = x_{w-1}, \dots, x_{w-1}, x_{w-1}, x_{w-2}, \dots, x_0$$



VIII–6

Sign Extension Example

short int x = 15213; int ix = (int) x; short int y = -15213; int iy = (int) y;

	Decimal	Hex	Binary
x	15213	3B 6D	00111011 01101101
ix	15213	00 00 3B 6D	00000000 00000000 00111011 01101101
У	-15213	C4 93	11000100 10010011
iy	-15213	FF FF C4 93	11111111 1111111 11000100 10010011

• Converting from smaller to larger integer data type

C automatically performs sign extension

Does it Work?

$$val_{w} = -2^{w-1} + \sum_{i=0}^{w-2} b_{i} \cdot 2^{i}$$

$$val_{w+1} = -2^{w} + 2^{w-1} + \sum_{i=0}^{w-2} b_{i} \cdot 2^{i}$$
$$= -2^{w-1} + \sum_{i=0}^{w-2} b_{i} \cdot 2^{i}$$

$$val_{w+2} = -2^{w+1} + 2^{w} + 2^{w-1} + \sum_{i=0}^{w-2} b_i \cdot 2^i$$
$$= -2^{w} + 2^{w-1} + \sum_{i=0}^{w-2} b_i \cdot 2^i$$
$$= -2^{w-1} + \sum_{i=0}^{w-2} b_i \cdot 2^i$$

CS33 Intro to Computer Systems

VIII-8 Copyright © 2022 Thomas W. Doeppner. All rights reserved.

Unsigned Multiplication



- Standard multiplication function
 - ignores high order w bits
- Implements modular arithmetic

 $UMult_w(u, v) = u \cdot v \mod 2^w$

Signed Multiplication



- Standard multiplication function
 - ignores high order w bits
 - some of which are different from those of unsigned multiplication
 - lower bits are the same
 - » but most-significant bit of TMULT determines sign

Power-of-2 Multiply with Shift

Operation

- $-u \ll k$ gives $u \times 2^k$ k both signed and unsigned U operands: w bits * 2^k 010 ... ••• $u * 2^{k}$ true product: *w*+*k* bits ... $\mathrm{UMult}_{w}(u, 2^{k})$ discard k bits: w bits ... ••• $\mathrm{TMult}_{w}(u, 2^{k})$ Examples u << 3 == u * 8
 - u << 5 u << 3 == u * 24
 - most machines shift and add faster than multiply
 - » compiler generates this code automatically

Unsigned Power-of-2 Divide with Shift

- Quotient of unsigned and power of 2
 - $-u \gg k \text{ gives } \lfloor u / 2^k \rfloor$
 - uses logical shift



	Division	Computed	Hex	Binary
x	15213	15213	3B 6D	00111011 01101101
x >> 1	7606.5	7606	1D B6	00011101 10110110
x >> 4	950.8125	950	03 B6	00000011 10110110
x >> 8	59.4257813	59	00 3B	0000000 00111011

Signed Power-of-2 Divide with Shift

- Quotient of signed and power of 2
 - $-\mathbf{x} \gg \mathbf{k}$ gives $\lfloor \mathbf{x} / 2^k \rfloor$
 - uses arithmetic shift
 - rounds wrong direction when x < 0



	Division	Computed	Hex	Binary	
У	-15213	-15213	C4 93	11000100 10010011	
y >> 1	-7606.5	-7607	E2 49	1 1100010 01001001	
y >> 4	-950.8125	-951	FC 49	1111 1100 01001001	
y >> 8	-59.4257813	-60	FF C4	1111111 11000100	

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 << k) 1) >> k
- » biases dividend toward 0

Case 1: no rounding



Biasing has no effect

Correct Power-of-2 Divide (Cont.)

Case 2: rounding



Biasing adds 1 to final result

Why Should I Use Unsigned?

- Don't use just because number nonnegative
 - easy to make mistakes

```
unsigned i;
for (i = cnt-2; i >= 0; i--)
a[i] += a[i+1];
```

– can be very subtle

```
#define DELTA sizeof(int)
int i;
for (i = CNT; i-DELTA >= 0; i-= DELTA)
....
```

- Do use when using bits to represent sets
 - logical right shift, no sign extension

Word Size

- (Mostly) obsolete term
 - old computers had items of one size: the word size
- Now used to express the number of bits necessary to hold an address
 - 16 bits (really old computers)
 - 32 bits (old computers)
 - 64 bits (most current computers)

Byte Ordering

- Four-byte integer
 - -0x76543210

Stored at location 0x100

- which byte is at 0x100?
- which byte is at 0x103?





CS33 Intro to Computer Systems

VIII-19 Copyright © 2022 Thomas W. Doeppner. All rights reserved.

Quiz 2

```
int main() {
    long x=1;
    func((int *)&x);
    return 0;
}
```

```
void func(int *arg) {
   printf("%d\n", *arg);
}
```

What value is printed on a big-endian 64-bit computer? a) 1 b) 0 c) 2³² d) 2³²-1

Which Byte Ordering Do We Use?

```
int main() {
    unsigned int x = 0x03020100;
    unsigned char *xarray = (unsigned char *)&x;
    for (int i=0; i<4; i++) {
        printf("%02x", xarray[i]);
    }
    printf("\n");
    return 0;
}
</pre>
```

03020100

CS33 Intro to Computer Systems

VIII-21 Copyright © 2022 Thomas W. Doeppner. All rights reserved.

Fractional binary numbers

• What is 1011.101₂?

Fractional Binary Numbers



Representable Numbers

• Limitation #1

- can exactly represent only numbers of the form n/2^k
 - » other rational numbers have repeating bit representations

Limitation #2

- just one setting of decimal point within the w bits
 - » limited range of numbers (very small values? very large?)

IEEE Floating Point

IEEE Standard 754

- established in 1985 as uniform standard for floating point arithmetic
 - » before that, many idiosyncratic formats
- supported on 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

Floating-Point Representation

Numerical Form:

- 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)

S	ехр	frac
---	-----	------

Precision options

• Single precision: 32 bits

S	ехр	frac
1	8-bits	23-bits

Double precision: 64 bits



• Extended precision: 80 bits (Intel only)

	s	exp	frac		
_	1	15-bits		64-bits	
CS33	Intro	to Computer Systems	VIII–27		

"Normalized" Values

- When: exp ≠ 000...0 and exp ≠ 111...1
- Exponent coded as biased value: E = Exp Bias
 - exp: unsigned value exp
 - 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.xxx...x2
 - xxx...x: bits of frac
 - minimum when frac=000...0 (M = 1.0)
 - maximum when frac=111...1 (M = 2.0ϵ)
 - get extra leading bit for "free"

Normalized Encoding Example

- Value: float F = 15213.0;
 - $-15213_{10} = 11101101101_{2}$
 - = 1.1101101101₂ x 2¹³

Significand

Μ	=	1.1101101101_2
frac	=	<u>1101101101101</u> 000000000 ₂

• Exponent

Ε	=	13		
bias	=	127		
exp	=	140	=	10001100 ₂

• Result:

0 10001100 1101101101000000000 s exp frac

Denormalized Values

- Condition: exp = 000...0
- Exponent value: E = –Bias + 1 (instead of E = 0 Bias)
- Significand coded with implied leading 0: M = 0.xxx...x2
 - xxx...x: bits of frac

Cases

- $\exp = 000...0, \operatorname{frac} = 000...0$
 - » represents zero value
 - » note distinct values: +0 and -0 (why?)
- $-\exp = 000...0, \operatorname{frac} \neq 000...0$
 - » numbers closest to 0.0
 - » equispaced

Special Values

- **Condition**: exp = 111...1
- Case: exp = 111...1, frac = 000...0
 - represents value ∞ (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: exp = 111...1, frac $\neq 000...0$
 - not-a-number (NaN)
 - represents case when no numeric value can be determined
 - e.g., sqrt(-1), $\infty \infty$, $\infty \times 0$

Visualization: Floating-Point Encodings



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

Dynamic Range (Positive Only)

	S	exp	frac	Е	Value		
	0	0000	000	-6	0		
	0	0000	001	-6	1/8*1/64 =	: 1/512	closest to zero
Denormalized	0	0000	010	-6	2/8*1/64 =	: 2/512	
numbers							
	0	0000	110	-6	6/8*1/64 =	6/512	
	0	0000	111	-6	7/8*1/64 =	: 7/512	largest denorm
	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
Normalized	0	0111	000	0	8/8*1 =	: 1	
numbers	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		
CS33 Intro to	Cor	nputer S	ystems		VIII–34		

Distribution of Values

- 6-bit IEEE-like format
 - e = 3 exponent bits
 - f = 2 fraction bits

• Notice how the distribution gets denser toward zero. 8 values



Distribution of Values (close-up view)

- 6-bit IEEE-like format
 - e = 3 exponent bits
 - f = 2 fraction bits

– bias is 3

	S	exp	frac
-	1	3-bits	2-bits



Quiz 3

6-bit IEEE-like format

- e = 3 exponent bitssexpfrac- f = 2 fraction bits13-bits2-bits

What number is represented by 0 010 10? a) 3 b) 1.5 c) .75 d) none of the above

Mapping Real Numbers to Float

- The real number 3 is represented as 0 100 10
- The real number 3.5 is represented as 0 100 11
- How is the real number 3.4 represented?
 0 100 11
- How is the real number π represented?
 0 100 10



Mapping Real Numbers to Float

- If R is a real number, it's mapped to the floating-point number whose value is closest to R
- What if it's midway between two values?
 - rounding rules coming up soon!

Floats are Sets of Values

- If A, B, and C are successive floating-point values
 - e.g., 010001, 010010, and 010011
- B represents all real numbers from midway between A and B through midway between B and C



CS33 Intro to Computer Systems

VIII-40 Copyright © 2022 Thomas W. Doeppner. All rights reserved.

Significance

- Normalized numbers
 - for a particular exponent value E and an S-bit significand, the range from 2^E up to 2^{E+1} is divided into 2^S equi-spaced floating-point values
 - » thus each floating-point value represents 1/2^s of the range of values with that exponent
 - » all bits of the signifcand are important
 - » we say that there are S significant bits for reasonably large S, each floating-point value covers a rather small part of the range
 - high accuracy
 - for S=23 (32-bit float), accurate to one in 2²³ (.0000119% accuracy)

Significance

- Unnormalized numbers
 - high-order zero bits of the significand aren't important
 - in 8-bit floating point, 0 0000 001 represents 2⁻⁹
 - » it is the only value with that exponent: 1 significant bit (either 2⁻⁹ or 0)
 - 0 0000 010 represents 2⁻⁸
 0 0000 011 represents 1.5*2⁻⁸
 - » only two values with exponent -8: 2 significant bits (encoding those two values, as well as 2⁻⁹ and 0)
 - fewer significant bits mean less accuracy
 - 0 0000 001 represents a range of values from .5*2-9 to 1.5*2-9
 - 50% accuracy

+/- Zero

- Only one zero for ints
 - an int is a single number, not a range of numbers, thus there can be only zero
- Floating-point zero
 - a range of numbers around the real 0
 - it really matters which side of 0 we're on!
 - » a very large negative number divided by a very small negative number should be positive

 $-\infty/-0 = +\infty$

» a very large positive number divided by a very small negative number should be negative

 $+\infty /-0 = -\infty$