Data Representation



CS/COE 0449 Introduction to Systems Software

wilkie

(with content borrowed from Vinicius Petrucci and Jarrett Billingsley)

Spring 2019/2020

BIT MANIPULATION

Flippin' Switches

2

CS/COE 0449 - Spring 2019/2020

What are "bitwise" operations?

- The "numbers" we use on computers aren't *really* numbers right?
- It's often useful to treat them instead as a pattern of bits.
- Bitwise operations treat a value as a pattern of bits.



The simplest operation: NOT (logical negation)



- We can summarize this in a truth table.
- We write NOT as $\sim A$, or $\neg A$, or \overline{A}
- In C, the NOT operation is the "!" operator

Applying NOT to a whole bunch of bits

• If we use the not instruction (~ in C), this is what happens:



That's it.

only 8 bits shown cause 32 bits on a slide is too much

Let's add some switches

- There are two switches in a row connecting the light to the battery.
- How do we make it light up?



AND (Logical product)

- AND is a binary (two-operand) operation.
- It can be written a number of ways: A&B $A \land B$ $A \cdot B$ AB
- If we use the and instruction (& in C):

&

we did 8 independent AND operations

"Switching" things up ;))))))))))))))))))))))))))))

• NOW how can we make it light up?



OR ("Logical" sum...?)

- We might say "and/or" in English.
- It can be written a number of ways:
 A B AVB A+B
- If we use the or instruction (| in C):



B Α 0 0 1 1 0 1 1 0

We did 8 independent OR operations.

Bit shifting

• Besides AND, OR, and NOT, we can move bits around, too.



Left-shifting in C/Java

C (and Java) use the << operator for left shift

B = A << 4; // B = A shifted left 4 bits

If the bottom 4 bits of the result are now 0s...

...what happened to the top 4 bits?

0011 0000 0000 1111 1100 1101 1100 1111

the bit bucket is not a real place

it's a programmer joke ok

in the UK they might say the "Bit Bin"

bc that's their word for trash



(animated)

Bit Bucket

> > > > (\cdot)

• We can shift right, too

0 1 1 0 0 0 0 0 0 0 1 1 1 1 1 1 0 0 1 1 0 1 1 1 0 0 1 1 1 1 0 1 1 0 0 0 0 0 0 0 1 1 1 1 1 1 0 0 1 1 0 1 1 1 0 0 1 1 1 1 1

• C/Java use >>, this in MIPS is the srl (Shift Right Logical) instruction

see what I mean about 32 bits on a slide

Q: What happens when we shift a negative number to the right? 12

Shift Right (Logical)

• We can shift right, too (srl in MIPS)



if we shift these bits **right by 1...** we stick a **0** at the top

again!

AGAIN!

Wait... what if this was a negative number?

Shift Right (Arithmetic)

• We can shift right with sign-extension, too (MIPS: sra)



if we shift these bits **right by 1...**

we copy the **1** at the top (or 0, if MSB was a 0)

again!

AGAIN!

AGAIN!!!!! (It's still negative!)

Huh... that's weird

• Let's start with a value like 5 and shift left and see what happens:

Binary	Decimal
101	5
1010	10
10100	20
101000	40
1010000	80

Why is this happening Well uh... what if I gave you **49018853** How do you multiply that by 10? by 100?

by 100000?

Something **very similar** is happening here

a << n == a * 2ⁿ

- Shifting left by n is the same as multiplying by 2ⁿ
 - You probably learned this as "moving the decimal point"
 - And moving the decimal point *right* is like shifting the digits *left*
- Shifting is fast and easy on most CPUs.
 - Way faster than multiplication in any case.
 - It's not a great reason to do it when you're using C though
- Hey... if shifting *left* is the same as multiplying...

a >> n == a / 2ⁿ, ish

- You got it
- Shifting right by n is like dividing by 2ⁿ
 sort of.
- What's 101₂ shifted right by 1?
 - 10₂, which is 2...
 - It's like doing integer (or flooring) division
- Generally, compilers are smart enough that you just multiply/divide
 - It's confusing to shift just to optimize performance.
 - It's good to not be clever until it is proven that you need to be.

C Bitwise Operations: Summary

C code	Description	MIPS instruction
х у	or	or x, x, y
х&у	and	and x, x, y
х ^ у	xor	xor x, x, y
! x	not	seq x, x, \$0 ("seqz")
~x	complement (negate)	nor x, x, \$0 ("not")
x << y	left-shift logical	sll x, x, y
x >> y	right-shift logical	srl x, x, y

When x is signed (most of the time...): x >> y right-shift arithmetic sra x, x, y

FRACTIONAL ENCODING

Every Time I Teach Floats I Want Some Root Beer

19

CS/COE 0449 - Spring 2019/2020

Fractional numbers

- Up to this point we have been working with integer numbers.
 - Unsigned and signed!

2019

• However, Real world numbers are... Real numbers. Like so:

2019.320

That creates new challenges!

Let's start by taking a look at them.

Just a fraction of a number

- The numbers we use are written positionally: the position of a digit within the number has a meaning.
- What about when the numbers go over the decimal point?

2019.320 1000s 100s 10s 1s 10ths 100ths 1000ths 10³ 10² 10¹ 10⁰ 10⁽⁻¹⁾ 10⁽⁻²⁾ 10⁽⁻³⁾

A fraction of a bit?

- Binary is the same!
- Just replace 10s with 2s.

2⁽⁻¹⁾ 2⁽⁻²⁾ 2⁽⁻³⁾ 2⁽⁻⁴⁾ **2**³ **2**² 2¹ 20 4s 1s 8ths 2s 2ths² 4ths 16ths **8**s

To convert into decimal, just add stuff

 $0 \times 8 +$ $1 \times 4 +$ $1 \times 2 +$ $0 \times 1 +$ $1 \times .5 +$ $1 \times .25 +$ $0 \times .125 +$ $1 \times .0625$

$$= 6.8125_{10}$$

From decimal to binary? Tricky?



So, it's easy right? Well...



So, it's easy right? Well.....



So, it's easy right? Well.....





0. 1₁₀

How much is it worth?

•Well, it depends on where you stop!

 $0.0001_2 = 0.0625$ $0.00011001_2 = 0.0976...$

$0.000110011001_2 = 0.0998...$

Fixing the point

• If we want to represent decimal places, one way of doing so is by assuming that the lowest *n* digits are the decimal places.

\$12.34 1234 +\$10.81 +1081 this is \$23.15 2315

this is called *fixed-point representation*

A rising tide

• Maybe half-and-half? 16.16 number looks like this:

the largest (signed) value we can the smallest fraction we can represent is +32767.9999ish represent is 1/65536

...but if we place the binary point to the right... **0011 0000 0101 1010 1000 0000.1111 1111** ...then we trade off accuracy for *range* further away from 0.

Move the point

- What if we could float the point around?
 - Enter scientific notation: The number **-0.0039** can be represented:

-0.39 × 10⁻² -3.9 × 10⁻³

- These both represent the same number, but we need to move the decimal point according to the power of ten represented.
- The bottom example is in normalized scientific notation.
 - There is only one non-zero digit to the left of the point.
- Because the decimal point can be moved, we call this representation:



IEEE 754

- Established in 1985, updated as recently as 2008.
- Standard for floating-point representation and arithmetic that virtually every CPU now uses.
- Floating-point representation is based around scientific notation:

$$1348 = +1.348 \times 10^{+3}$$

-0.0039 = -3.9 $\times 10^{-3}$
-1440000 = -1.44 $\times 10^{+6}$
sign significand exponen

t

Binary Scientific Notation

- Scientific notation works equally well in any other base!
 - (below uses base-10 exponents for clarity)
- $+1001 0101 = +1.001 0101 \times 2^{+7}$ × 2⁻³ $-0.001\ 010 = -1.010$ × 2⁺¹⁵ -1001 0000 0000 0000 = -1.001What do you notice about the digit before the **binary** point? f - fraction $(+/-)1.f \times 2^{exp}$ 1.f - significand
 - exp exponent

IEEE 754 Single-precision

- Known as float in C/C++/Java etc., 32-bit float format
- 1 bit for sign, 8 bits for the exponent, 23 bits for the *fraction*



Tradeoff:

- More accuracy = More fraction bits
- More range = More exponent bits

• Every design has tradeoffs ⁻_(ツ)_/⁻

Welcome to Systems!

illustration from user Stannered on Wikimedia Commons

IEEE 754 Single-precision

- Known as float in C/C++/Java etc., 32-bit float format
- 1 bit for sign, 8 bits for the exponent, 23 bits for the *fraction*

- The fraction field only stores the digits after the binary point
- The 1 before the binary point is implicit!
 - This is called normalized representation
 - In effect this gives us a 24-bit significand
 - The only number with a 0 before the binary point is 0!
- The significand of floating-point numbers is in sign-magnitude!
 - Do you remember the downside(s)?

illustration from user Stannered on Wikimedia Commons

The exponent field

- The exponent field is 8 bits, and can hold positive or negative exponents, but... it doesn't use S-M, 1's, or 2's complement.
- It uses something called biased notation.
 - biased representation = signed number + bias constant
 - single-precision floats use a bias constant of 127

- The exponent can range from -126 to +127 (1 to 254 biased)
 - 0 and 255 are reserved!
- Why'd they do this?
 - so you can sort floats with integer comparisons??

Binary Scientific Notation (revisited)

Our previous numbers are actually

+1.001 0101 ×
$$2^{+7} = (-1)^{0}$$
 × 1.001 0101 × $2^{134-127}$
-1.010 × $2^{-3} = (-1)^{1}$ × 1.010 × $2^{124-127}$
-1.001 × $2^{+15} = (-1)^{1}$ × 1.001 × $2^{142-127}$

Encoding an integer as a float

- You have an integer, like 2471 = 0000 1001 1010 0111₂
 - 1. put it in scientific notation
 - 1.001 1010 $0111_2 \times 2^{+11}$

2. get the exponent field by adding the bias constant

- $11 + 127 = 138 = 10001010_2$
- 3. copy the bits after the binary point into the fraction field

	S	exponent	fraction
	0	10001010	00110100111 000000000
	t		
ро	sitiv	ve	start at the left side! 39

Encoding a number as a float

You have a number, like -12.59375₁₀

- **1. Convert to binary:** Integer part: 1100₂ Fractional part: 0.10011₂
- **2.** Write it in scientific notation: $1100.10011_2 \times 2^0$
- **3.** Normalize it: $1.10010011_2 \times 2^3$
- 4. Calculate biased exponent $+3 + 127 = 130_{10} = 10000010_2$

S	exponent	fraction
1	10000010	10010011 0000000000000

0xC1498000

4()

while (computers don't do real math) { ... }

```
public class FloatTest {
   public static void main(String[] args) {
      double var = Double.MAX_VALUE;
```



```
while (var == Double.MAX_VALUE) {
    // Hang the program while the condition holds
    var = var + 0.1;
}
```

```
System.out.println("This never prints.");
```

}

Q: Consider and/or review the IEEE 754 standard. What is happening here? $_4$

Other formats

• The most common other format is double-precision (C/C++/Java double), which uses an 11-bit exponent and 52-bit fraction

This could be a whole unit itself...

- Floating-point arithmetic is COMPLEX STUFF.
- But it's not super useful to know unless you're either:
 - Doing lots of high-precision numerical programming, or
 - Implementing floating-point arithmetic yourself.
- However...
 - It's good to have an understanding of why limitations exist.
 - It's good to have an *appreciation* of how complex this is... and how much better things are now than they were in the 1970s and 1980s!
 - It's good to know things do not behave as expected when using float and double!!