# Bits, Bytes, and Integers

15-213: Introduction to Computer Systems 2<sup>nd</sup> and 3<sup>rd</sup> Lectures, Jan 17 and Jan 22, 2013

#### Instructors:

Seth Copen Goldstein, Anthony Rowe, Greg Kesden

#### **MLK** recitations

- No recitations after 12:30, so ...
- The TAs have been kind enough to create some temporary sections:

■ GHC 4215: 10:30 & 11:30

**GHC 4102: 11:30** 

■ GHC 4101: 9:30 & 10:30



Carnegie Mellon

#### Waitlist

- Please be patient.
- If you register for autolab, get the work done → you will be ready when you get into the class

# **Today: Bits, Bytes, and Integers**

- Representing information as bits
- **Bit-level manipulations**
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings

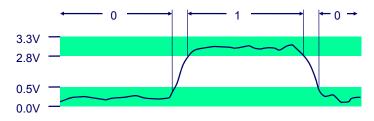
## **Binary Representations**

#### ■ Base 2 Number Representation

- Represent 15213<sub>10</sub> as 11101101101101<sub>2</sub>
- Represent 1.20<sub>10</sub> as 1.001100110011[0011]...<sub>2</sub>
- Represent 1.5213 X 10<sup>4</sup> as 1.1101101101101<sub>2</sub> X 2<sup>13</sup>

#### ■ Electronic Implementation

- Easy to store with bistable elements
- Reliably transmitted on noisy and inaccurate wires



# **Encoding Byte Values**

- Byte = 8 bits
  - Binary 000000002 to 111111112
  - Decimal: 0<sub>10</sub> to 255<sub>10</sub>
  - Hexadecimal 00<sub>16</sub> to FF<sub>16</sub>
    - Base 16 number representation
    - Use characters '0' to '9' and 'A' to 'F'
    - Write FA1D37B<sub>16</sub> in C as
      - 0xFA1D37B
      - 0xfa1d37b

He	+ 0e	zimal Binary
0	0	0000
0 1 2 3 4	0 1 2 3 4 5 6	0001
2	2	0010
3	3	0011
4	4	0100
5	5	0101
6 7	6	0110
7	7	0111
8	8	1000
9	9	1001
Α	10	1010
В	11	1011
C	12	1100
D	13	1101
Е	14	1110
F	15	1111

Carnegie Mellon

# **Data Representations**

C Data Type	Typical 32-bit	Intel IA32	x86-64
char	1	1	1
short	2	2	2
int	4	4	4
long	4	4	8
long long	8	8	8
float	4	4	4
double	8	8	8
long double	8	10/12	10/16
pointer	4	4	8

**Today: Bits, Bytes, and Integers** 

- Representing information as bits
- **Bit-level manipulations**
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings

# **Boolean Algebra**

#### Developed by George Boole in 19th Century

- Algebraic representation of logic
  - Encode "True" as 1 and "False" as 0

#### And

#### Or

#### Not

#### **Exclusive-Or (Xor)**

# **General Boolean Algebras**

- Operate on Bit Vectors
  - Operations applied bitwise

All of the Properties of Boolean Algebra Apply

#### Carnegie Mellon

11

# **Example: Representing & Manipulating Sets**

#### Representation

- Width w bit vector represents subsets of {0, ..., w-1}
- $a_i = 1$  if  $j \in A$ 
  - 01101001 { 0, 3, 5, 6 }
  - **76543210**
  - 01010101 { 0, 2, 4, 6 }
  - **76543210**

#### Operations

<b>.</b> &	Intersection	01000001	{ 0, 6 }
•	Union	01111101	{ 0, 2, 3, 4, 5, 6 }
■ ^	Symmetric difference	00111100	{ 2, 3, 4, 5 }
■ ~	Complement	10101010	{ 1, 3, 5, 7 }

# **Bit-Level Operations in C**

- Operations &, |, ~, ^ Available in C
- Apply to any "integral" data type
  - long, int, short, char, unsigned
- View arguments as bit vectors
- Arguments applied bit-wise

#### ■ Examples (Char data type)

- ~0x41 → 0xBE
  - $\bullet$  ~01000001<sub>2</sub> → 10111110<sub>2</sub>
- $\sim 0x00 \rightarrow 0xFF$ 
  - $\sim 00000000_2 \rightarrow 11111111_2$
- $0x69 \& 0x55 \rightarrow 0x41$ 
  - $01101001_2 \& 01010101_2 \rightarrow 01000001_2$
- 0x69 | 0x55 → 0x7D
  - $01101001_2 \mid 01010101_2 \rightarrow 01111101_2$

# **Contrast: Logic Operations in C**

- Contrast to Logical Operators
  - **&** &&, ||, !
    - View 0 as "False"
    - Anything nonzero as "True"
    - Always return 0 or 1
    - Early termination
- Examples (char data type)
  - !0x41 → 0x00
  - $!0x00 \rightarrow 0x01$
  - $!!0x41 \rightarrow 0x01$
  - $0x69 \&\& 0x55 \rightarrow 0x01$
  - 0x69 || 0x55 → 0x01
  - p && \*p (avoids null pointer access)

## **Contrast: Logic Operations in C**

- **Contrast to Logical Operators**
- **&** &&, ||, !
  - View 0 as "Fal
  - Anything ponze
  - Alway
  - Early Watch out for && vs. & (and || vs. |)...
- Example one of the more common oopsies in ■ !0x41 -
  - **C** programming

• !0x00

!!0x41

- $0x69 \&\& 0x55 \rightarrow 0x01$ ■ 0x69 || 0x55 → 0x01
- p && \*p (avoids null pointer access)

13

Carnegie Mellon

Carnegie Mello

# **Shift Operations**

- Left Shift: x << y
  - Shift bit-vector x left y positions
    - Throw away extra bits on left
    - Fill with 0's on right
- Right Shift: x >> y
  - Shift bit-vector x right y positions
    - Throw away extra bits on right
  - Logical shift
    - Fill with 0's on left
  - Arithmetic shift
    - Replicate most significant bit on left
- **Undefined Behavior** 
  - Shift amount < 0 or ≥ word size</p>

Argument x	01100010	
<< 3	00010 <i>000</i>	
Log. >> 2	00011000	
Arith. >> 2	00011000	

Argument x	10100010
<< 3	00010 <i>000</i>
Log. >> 2	00101000
Arith. >> 2	<b>11</b> 101000

# **Today: Bits, Bytes, and Integers**

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings
- Summary

# **Encoding Integers**

#### Unsigned

#### **Two's Complement**

$$B2U(X) = \sum_{i=0}^{w-1} x_i \cdot 2^i$$

$$B2T(X) = -x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-2} x_i \cdot 2^i$$

short int x = 15213; short int y = -15213;

Sign Bit

#### ■ C short 2 bytes long

1		Decimal	Hex	Binary
	x	15213	3B 6D	00111011 01101101
	У	-15213	C4 93	11000100 10010011

- Sign Bit
  - For 2's complement, most significant bit indicates sign
    - 0 for nonnegative
    - 1 for negative

17

Carnegie Mellon

# **Two-complement Encoding Example (Cont.)**

x = 15213: 00111011 01101101y = -15213: 11000100 10010011

Weight	152	13	-152	213
1	1	1	1	1
2	0	0	1	2
4	1	4	0	0
8	1	8	0	0
16	0	0	1	16
32	1	32	0	0
64	1	64	0	0
128	0	0	1	128
256	1	256	0	0
512	1	512	0	0
1024	0	0	1	1024
2048	1	2048	0	0
4096	1	4096	0	0
8192	1	8192	0	0
16384	0	0	1	16384
-32768	0	0	1	-32768
Sum	•	15213		-15213

Carnegie Mellor

# **Numeric Ranges**

#### Unsigned Values

• 
$$UMax = 2^w - 1$$
111...1

#### ■ Two's Complement Values

■ 
$$TMin = -2^{w-1}$$

$$TMax = 2^{w-1} - 1$$

#### Other Values

Minus 1

111...1

#### Values for W = 16

	Decimal	Hex	Binary
UMax	65535	FF FF	11111111 11111111
TMax	32767	7F FF	01111111 11111111
TMin	-32768	80 00	10000000 00000000
-1	-1	FF FF	11111111 11111111
0	0	00 00	00000000 00000000

## **Values for Different Word Sizes**

	W				
	8	16	32	64	
UMax	255	65,535	4,294,967,295	18,446,744,073,709,551,615	
TMax	127	32,767	2,147,483,647	9,223,372,036,854,775,807	
TMin	-128	-32,768	-2,147,483,648	-9,223,372,036,854,775,808	

#### Observations

- |*TMin* | = *TMax* + 1
  - Asymmetric range
- UMax = 2 \* TMax + 1

#### C Programming

- #include limits.h>
- Declares constants, e.g.,
  - ULONG\_MAX
  - LONG MAX
  - LONG\_MIN
- Values platform specific

# **Unsigned & Signed Numeric Values**

Χ	B2U( <i>X</i> )	B2T(X)
0000	0	0
0001	1	1
0010	2	2
0011	3	3
0100	4	4
0101	5	5
0110	6	6
0111	7	7
1000	8	-8
1001	9	<b>-</b> 7
1010	10	-6
1011	11	<b>-</b> 5
1100	12	-4
1101	13	-3
1110	14	-2
1111	15	-1

#### Equivalence

 Same encodings for nonnegative values

#### Uniqueness

- Every bit pattern represents a unique integer value
- Each representable integer has a unique bit encoding

#### ■ ⇒ Can Invert Mappings

- U2B(x) = B2U<sup>-1</sup>(x)
  - Bit pattern for unsigned integer
- T2B(x) = B2T $^{-1}$ (x)
  - Bit pattern for two's comp integer

# **Today: Bits, Bytes, and Integers**

- Representing information as bits
- **Bit-level manipulations**

#### Integers

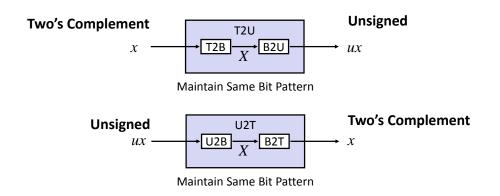
- Representation: unsigned and signed
- Conversion, casting
- Expanding, truncating
- Addition, negation, multiplication, shifting
- Summarv
- Representations in memory, pointers, strings

21

23

Carnegie Mellon

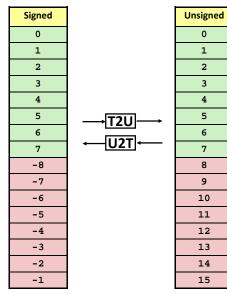
# **Mapping Between Signed & Unsigned**



■ Mappings between unsigned and two's complement numbers: keep bit representations and reinterpret

# Mapping Signed ↔ Unsigned

Bits
0000
0001
0010
0011
0100
0101
0110
0111
1000
1001
1010
1011
1100
1101
1110
1111



Carnegie Mello

1

2

3

4

5

6

7

8

9

10

11

12

13

14

Unsigned 0

1

2

3

6

10 11 12

13

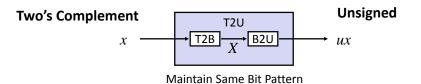
14

15

# Mapping Signed ↔ Unsigned

Bits	Signed	
0000	0	
0001	1	
0010	2	
0011	3	. = .
0100	4	$\leftarrow$
0101	5	
0110	6	
0111	7	
1000	-8	
1001	-7	
1010	-6	. / 46
1011	-5	+/- 16
1100	-4	
1101	-3	
1110	-2	
1111	-1	

**Relation between Signed & Unsigned** 





Large negative weight becomes

Large positive weight

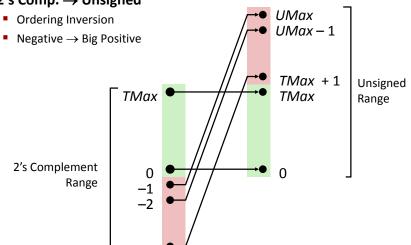
26

Carnegie Mello

Carnegie Mellon

#### **Conversion Visualized**

■ 2's Comp. → Unsigned



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 same as U2T and T2U

int tx, ty;
unsigned ux, uy;
tx = (int) ux;
uy = (unsigned) ty;

Implicit casting also occurs via assignments and procedure calls

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

27

# **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 <sub>1</sub>	Constant <sub>2</sub>	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

# Summary Casting Signed ↔ Unsigned: Basic Rules

- Bit pattern is maintained
- But reinterpreted
- Can have unexpected effects: adding or subtracting 2<sup>w</sup>
- Expression containing signed and unsigned int
  - int is cast to unsigned!!

30

Carnegie Mellon

29

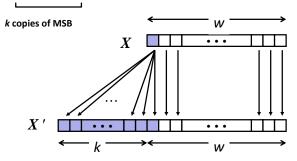
# **Today: Bits, Bytes, and Integers**

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings

Carnegie Mellon

# **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:
  - $X' = X_{w-1}, ..., X_{w-1}, X_{w-1}, X_{w-2}, ..., X_0$



# **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	1111111 11111111 11000100 10010011

- Converting from smaller to larger integer data type
- C automatically performs sign extension

33

# **Summary: Expanding, Truncating: Basic Rules**

- Expanding (e.g., short int to int)
  - Unsigned: zeros added
  - Signed: sign extension
  - Both yield expected result
- Truncating (e.g., unsigned to unsigned short)
  - Unsigned/signed: bits are truncated
  - Result reinterpreted
  - Unsigned: mod operation
  - Signed: similar to mod
  - For small numbers yields expected behavior

Carnegie Mello

Carnegie Mellon

# Fake real world example

- Acme, Inc. has developed a state of the art voltmeter they are connecting to a pc. It is precise to the millivolt and does not drain the unit under test.
- Your job is to develop the driver software.



printf("%d\n", getValue());

# Fake real world example

- Acme, Inc. has developed a state of the art voltmeter they are connecting to a pc. It is precise to the millivolt and does not drain the unit under test.
- Your job is to develop the driver software.





printf("%d\n", getValue());

#### Lets run some tests

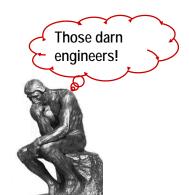
printf("%d\n", getValue());

- **50652**
- 1500
- 9692
- **26076**
- **17884**
- **42460**
- **34268**
- **50652**

Lets run some tests

int x=getValue(); printf("%d %08x\n",x, x);

- 0000c5dc **50652**
- 1500 000005dc
- 9692 000025dc
- 26076 000065dc
- 17884 000045dc
- **42460** 0000a5dc
- 34268 000085dc
- 50652 0000c5dc



Carnegie Mellon

Carnegie Mellon

# Only care about least significant 12 bits

```
int x=getValue();
x=(x & 0x0fff);
printf("%d\n",x);
```



# Only care about least significant 12 bits

```
int x=getValue();
x=x(\&0x0fff);
printf("%d\n",x);
```



printf("%x\n", x);

# Must sign extend

```
int x=getValue();
x=(x&0x00fff) | (x&0x0800?0xfffff000:0);
printf("%d\n",x);
```



There is a better way.

41

# **Because you graduated from 213**

int x=getValue();
x=(x&0x00fff) | (x&0x0800?0xfffff000:0);
printf("%d\n",x);



Carnegie Mellon

# Lets be really thorough



**Today: Bits, Bytes, and Integers** 

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
- Representations in memory, pointers, strings
- Summary

# **Unsigned Addition**

Operands: w bits

True Sum: w+1 bits

Discard Carry: w bits

 $UAdd_{w}(u, v)$ 

#### Standard Addition Function

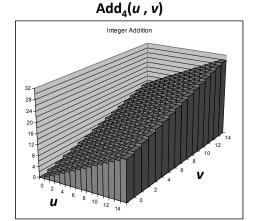
- Ignores carry output
- Implements Modular Arithmetic

$$s = UAdd_w(u, v) = u + v \mod 2^w$$

# **Visualizing (Mathematical) Integer Addition**

#### ■ Integer Addition

- 4-bit integers u, v
- Compute true sum Add₄(u, v)
- Values increase linearly with u and v
- Forms planar surface



45

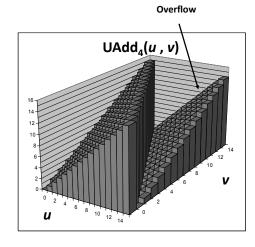
Carnegie Mellon

# **Visualizing Unsigned Addition**

#### Wraps Around

- If true sum  $\ge 2^w$
- At most once

# True Sum $\begin{array}{ccc} 2^{w+1} & & & \\ 2^{w} & & & & \\ 0 & & & & \\ \end{array}$ Modular Sum



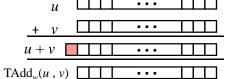
# **Two's Complement Addition**

Operands: w bits

True Sum: w+1 bits

Discard Carry: w bits

ts



#### TAdd and UAdd have Identical Bit-Level Behavior

Signed vs. unsigned addition in C:

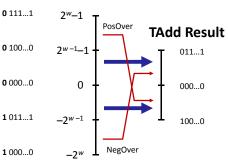
Will give s == t

#### **TAdd Overflow**

#### Functionality

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

#### **True Sum**



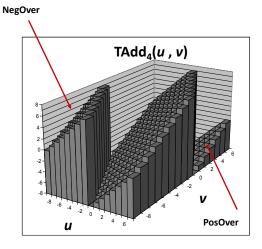
# **Visualizing 2's Complement Addition**

#### Values

- 4-bit two's comp.
- Range from -8 to +7

#### Wraps Around

- If sum  $\geq 2^{w-1}$ 
  - Becomes negative
  - At most once
- If sum  $< -2^{w-1}$ 
  - Becomes positive
  - At most once



Carnegie Mello

Carnegie Mellon

## Multiplication

- Goal: Computing Product of w-bit numbers x, y
  - Either signed or unsigned
- But, exact results can be bigger than w bits
  - Unsigned: up to 2w bits
    - Result range:  $0 \le x * y \le (2^w 1)^2 = 2^{2w} 2^{w+1} + 1$
  - Two's complement min (negative): Up to 2w-1 bits
    - Result range:  $x * y \ge (-2^{w-1})*(2^{w-1}-1) = -2^{2w-2} + 2^{w-1}$

1 011...1

- Two's complement max (positive): Up to 2w bits, but only for (TMin,,,)2
  - Result range:  $x * y \le (-2^{w-1})^2 = 2^{2w-2}$
- So, maintaining exact results...
  - would need to keep expanding word size with each product computed
  - is done in software, if needed
    - e.g., by "arbitrary precision" arithmetic packages

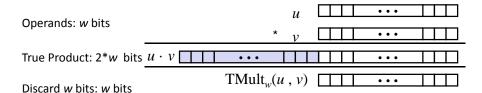
## **Unsigned Multiplication in C**

Operands: w bits True Product:  $2^*w$  bits  $u \cdot v$  $UMult_{w}(u, v)$ Discard w bits: w bits

- Standard Multiplication Function
  - Ignores high order w bits
- Implements Modular Arithmetic

$$UMult_{w}(u, v) = u \cdot v \mod 2^{w}$$

# **Signed Multiplication in C**

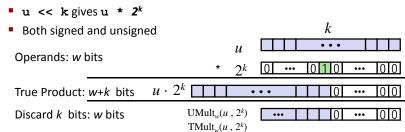


#### Standard Multiplication Function

- Ignores high order w bits
- Some of which are different for signed vs. unsigned multiplication
- Lower bits are the same

#### **Power-of-2 Multiply with Shift**

#### Operation



#### Examples

- u << 3 == u \* 8
- u << 5 u << 3 == u \* 2</pre>
- Most machines shift and add faster than multiply
  - Compiler generates this code automatically

53

J4

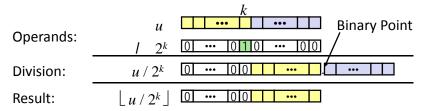
Carnegie Mello

Carnegie Mellon

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

#### Quotient of Unsigned by Power of 2

- $\mathbf{u} >> \mathbf{k}$  gives  $\left[\mathbf{u} / \mathbf{2}^{k}\right]$
- Uses logical shift

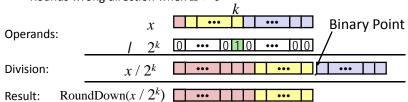


	Division	Computed	Hex	Binary
x	15213	15213	3B 6D	00111011 01101101
x >> 1	7606.5	7606	1D B6	00011101 10110110
x >> 4	950.8125	950	03 B6	00000011 10110110
x >> 8	59.4257813	59	00 3B	00000000 00111011

# **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 x< 0</li>



	Division	Computed	Hex	Binary
Y	-15213	-15213	C4 93	11000100 10010011
y >> 1	-7606.5	-7607	E2 49	11100010 01001001
y >> 4	-950.8125	-951	FC 49	<b>1111</b> 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  $\lceil \mathbf{x} / 2^k \rceil$  (Round Toward 0)
- Compute as  $\lfloor (x+2^k-1)/2^k \rfloor$ 
  - In C: (x + (1 << k) -1) >> k
  - Biases dividend toward 0

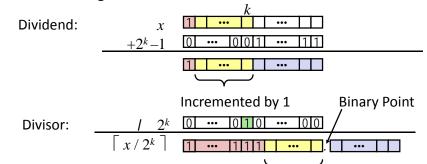
Case 1: No rounding

Biasing has no effect

57

# **Correct Power-of-2 Divide (Cont.)**

Case 2: Rounding



Incremented by 1

Biasing adds 1 to final result

Carnegie Mellon

# **Today: Bits, Bytes, and Integers**

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings

**Arithmetic: Basic Rules** 

#### Addition:

- Unsigned/signed: Normal addition followed by truncate, same operation on bit level
- Unsigned: addition mod 2<sup>w</sup>
  - Mathematical addition + possible subtraction of 2<sup>w</sup>
- Signed: modified addition mod 2<sup>w</sup> (result in proper range)
  - Mathematical addition + possible addition or subtraction of 2<sup>w</sup>

#### Multiplication:

- Unsigned/signed: Normal multiplication followed by truncate, same operation on bit level
- Unsigned: multiplication mod 2<sup>w</sup>
- Signed: modified multiplication mod 2<sup>w</sup> (result in proper range)

# 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 Performing Modular Arithmetic
  - Multiprecision arithmetic
- Do Use When Using Bits to Represent Sets
  - Logical right shift, no sign extension

61

# **Integer C Puzzles**

- Assume 32-bit word size, two's complement integers
- For each of the following C expressions: true or false? Why?
  - x < 0</li>
- $\Rightarrow$  ((x\*2) < 0)
- ux >= 0
- x & 7 == 7
- $\Rightarrow$  (x<<30) < 0
- ux > -1
- X > Y
- $\Rightarrow$  -x < -y

Initialization

unsigned uy = y;

- x \* x >= 0
- x > 0 & y > 0
- $\Rightarrow x + y > 0$
- x >= 0x <= 0</li>
- $\Rightarrow -x <= 0$  $\Rightarrow -x >= 0$
- (x|-x)>>31==-1
- ux >> 3 == ux/8
- x >> 3 == x/8
- x & (x-1) != 0

Carnegie Mello

Carnegie Mellon

# **Today: Bits, Bytes, and Integers**

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings

**Byte-Oriented Memory Organization** 

# 00.0

- Programs refer to data by address
  - Conceptually, envision it as a very large array of bytes
    - In reality, it's not, but can think of it that way
  - An address is like an index into that array
    - and, a pointer variable stores an address
- Note: system provides private address spaces to each "process"
- Think of a process as a program being executed
- So, a program can clobber its own data, but not that of others

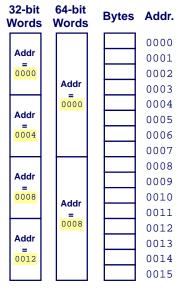
#### **Machine Words**

- Any given computer has a "Word Size"
  - Nominal size of integer-valued data
    - and of addresses
  - Most current machines use 32 bits (4 bytes) as word size
    - Limits addresses to 4GB (2<sup>32</sup> bytes)
    - Becoming too small for memory-intensive applications
      - leading to emergence of computers with 64-bit word size
  - Machines still support multiple data formats
    - Fractions or multiples of word size
    - Always integral number of bytes

65

# **Word-Oriented Memory Organization**

- Addresses Specify Byte Locations
  - Address of first byte in word
  - Addresses of successive words differ by 4 (32-bit) or 8 (64-bit)



Carnegie Mellon

# For other data representations too ...

C Data Type	Typical 32-bit	Intel IA32	x86-64	
char	1	1	1	
short	2	2	2	
int	4	4	4	
long	4	4	8	
long long	8	8	8	
float	4	4	4	
double	8	8	8	
long double	8	10/12	10/16	
pointer	4	4	8	

Carnegie Mellor

# **Byte Ordering**

- So, how are the bytes within a multi-byte word ordered in memory?
- Conventions
  - Big Endian: Sun, PPC Mac, Internet
    - Least significant byte has highest address
  - Little Endian: x86
    - Least significant byte has lowest address

# **Byte Ordering Example**

#### Example

- Variable x has 4-byte value of 0x01234567
- Address given by &x is 0x100

Big Endian	l	0x100	0x101	0x102	0x103	
		01	23	45	67	
Little Endia	an	0x100	0x101	0x102	0x103	
		67	45	23	01	

69

# **Representing Integers**

Decimal: 15213

Binary: 0011 1011 0110 1101

Hex: 3 B 6 D

Carnegie Mellon

```
int A = 15213;
                               long int C = 15213;
   IA32, x86-64
                   Sun
                                  IA32
                                               x86-64
                                                               Sun
                                                                00
                                    6D
                                                  6D
                    00
      3B
                                    3B
                                                  3B
                                                                00
                                    00
                                                  00
                                                                3B
      00
                                    00
int B = -15213:
                                                  00
                                                  00
   IA32, x86-64
                   Sun
                                                  00
                    FF
                    FF
      C4
                    C4
                            Two's complement representation
      FF
                    93
```

Carnegie Mellon

Carnegie Mellon

# **Examining Data Representations**

- Code to Print Byte Representation of Data
  - Casting pointer to unsigned char \* allows treatment as a byte array

```
typedef unsigned char *pointer;

void show_bytes(pointer start, int len){
  int i;
  for (i = 0; i < len; i++)
     printf("%p\t0x%.2x\n",start+i, start[i]);
  printf("\n");
}</pre>
```

#### **Printf directives:**

%p: Print pointer %x: Print Hexadecimal

# show\_bytes Execution Example

```
int a = 15213;
printf("int a = 15213;\n");
show_bytes((pointer) &a, sizeof(int));
```

#### Result (Linux):

```
int a = 15213;
0x11ffffcb8 0x6d
0x11ffffcb9 0x3b
0x11ffffcba 0x00
0x11ffffcbb 0x00
```

# **Reading Byte-Reversed Listings**

#### Disassembly

- Text representation of binary machine code
- Generated by program that reads the machine code

#### **■** Example Fragment

Address	Instruction Code	Assembly Rendition
8048365:	5b	pop %ebx
8048366:	81 c3 ab 12 00 00	add \$0x12ab,%ebx
804836c:	83 bb 28 00 00 00 00	cmpl \$0x0,0x28(%ebx)

#### Deciphering Numbers

- Value:
- Pad to 32 bits:
- Split into bytes:
- Reverse:

0x12ab

0x000012ab

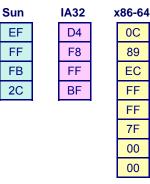
00 00 12 ab

ab 12 00 00

73

# **Representing Pointers**

int 
$$B = -15213;$$
  
int \*P = &B



Different compilers & machines assign different locations to objects

74

Carnegie Mello

Carnegie Mellon

# **Representing Strings**

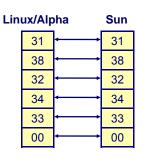
char S[6] = "18243";

#### Strings in C

- Represented by array of characters
- Each character encoded in ASCII format
  - Standard 7-bit encoding of character set
  - Character "0" has code 0x30
    - Digit i has code 0x30+i
- String should be null-terminated
  - Final character = 0

#### Compatibility

Byte ordering not an issue

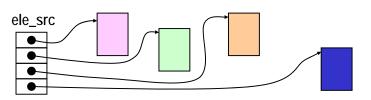


# **Code Security Example**

#### SUN XDR library

Widely used library for transferring data between machines

void\* copy\_elements(void \*ele\_src[], int ele\_cnt, size\_t ele\_size);



malloc(ele\_cnt \* ele\_size)



Carnegie Mellor

#### arnegie Mellon

Carnegie Mello

#### **XDR Code**

```
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;
}</pre>
```

# **XDR Vulnerability**

malloc(ele\_cnt \* ele\_size)

- What if:
  - ele\_cnt = 2<sup>20</sup> + 1
  - ele\_size = 4096 = 2<sup>12</sup>
  - Allocation = ??
- How can I make this function secure?

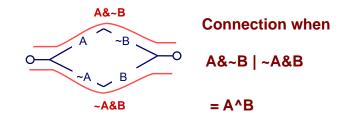
77

Carnegie Mellon

#### **Bonus extras**

# **Application of Boolean Algebra**

- Applied to Digital Systems by Claude Shannon
  - 1937 MIT Master's Thesis
  - Reason about networks of relay switches
    - Encode closed switch as 1, open switch as 0



79

## **Code Security Example**

```
/* Kernel memory region holding user-accessible data */
#define KSIZE 1024
char kbuf[KSIZE];

/* Copy at most maxlen bytes from kernel region to user buffer */
int copy_from_kernel(void *user_dest, int maxlen) {
    /* Byte count len is minimum of buffer size and maxlen */
    int len = KSIZE < maxlen ? KSIZE : maxlen;
    memcpy(user_dest, kbuf, len);
    return len;
}</pre>
```

- Similar to code found in FreeBSD's implementation of getpeername
- There are legions of smart people trying to find vulnerabilities in programs

# **Typical Usage**

```
/* Kernel memory region holding user-accessible data */
#define KSIZE 1024
char kbuf[KSIZE];

/* Copy at most maxlen bytes from kernel region to user buffer */
int copy_from_kernel(void *user_dest, int maxlen) {
    /* Byte count len is minimum of buffer size and maxlen */
    int len = KSIZE < maxlen ? KSIZE : maxlen;
    memcpy(user_dest, kbuf, len);
    return len;
}</pre>
```

```
#define MSIZE 528

void getstuff() {
    char mybuf[MSIZE];
    copy_from_kernel(mybuf, MSIZE);
    printf("%s\n", mybuf);
}
```

Carnegie Mellon

# Malicious Usage /\* Declaration of library function memcpy \*/

```
/* Declaration of library function memcpy */
void *memcpy(void *dest, void *src, size_t n);
```

```
/* Kernel memory region holding user-accessible data */
#define KSIZE 1024
char kbuf[KSIZE];

/* Copy at most maxlen bytes from kernel region to user buffer */
int copy_from_kernel(void *user_dest, int maxlen) {
    /* Byte count len is minimum of buffer size and maxlen */
    int len = KSIZE < maxlen ? KSIZE : maxlen;
    memcpy(user_dest, kbuf, len);
    return len;
}</pre>
```

```
#define MSIZE 528

void getstuff() {
    char mybuf[MSIZE];
    copy_from_kernel(mybuf, -MSIZE);
    . . .
}
```

#### **Mathematical Properties**

- Modular Addition Forms an Abelian Group
  - Closed under addition

```
0 \leq \mathsf{UAdd}_{w}(u, v) \leq 2^{w} - 1
```

Commutative

```
UAdd_{u}(u, v) = UAdd_{u}(v, u)
```

Associative

```
UAdd_{w}(t, UAdd_{w}(u, v)) = UAdd_{w}(UAdd_{w}(t, u), v)
```

0 is additive identity

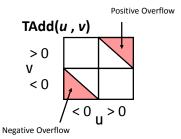
```
UAdd_{u}(u,0) = u
```

- Every element has additive inverse
  - Let  $UComp_w(u) = 2^w u$  $UAdd_w(u, UComp_w(u)) = 0$

# **Characterizing TAdd**

#### Functionality

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



$$TAdd_{w}(u,v) \quad = \quad \begin{cases} u+v+2^{w} & u+v < TMin_{w} \text{ (NegOver)} \\ u+v & TMin_{w} \leq u+v \leq TMax_{w} \\ u+v-2^{w} & TMax_{w} < u+v \text{ (PosOver)} \end{cases}$$

# **Mathematical Properties of TAdd**

- Isomorphic Group to unsigneds with UAdd
  - TAdd<sub>w</sub>(u , v) = U2T(UAdd<sub>w</sub>(T2U(u ), T2U(v)))
    - Since both have identical bit patterns
- Two's Complement Under TAdd Forms a Group
  - Closed, Commutative, Associative, 0 is additive identity
  - Every element has additive inverse

$$TComp_w(u) = \begin{cases} -u & u \neq TMin_w \\ TMin_w & u = TMin_w \end{cases}$$

Carnegie Mello

Carnegie Mellon

# **Compiled Multiplication Code**

#### **C** Function

```
int mul12(int x)
{
   return x*12;
}
```

#### **Compiled Arithmetic Operations**

#### **Explanation**

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

# **Compiled Unsigned Division Code**

#### **C** Function

```
unsigned udiv8(unsigned x)
{
  return x/8;
}
```

#### **Compiled Arithmetic Operations**

```
shrl $3, %eax
```

#### **Explanation**

# Logical shift
return x >> 3;

- Uses logical shift for unsigned
- For Java Users
  - Logical shift written as >>>

# **Compiled Signed Division Code**

#### **C** Function

```
int idiv8(int x)
{
  return x/8;
}
```

#### **Compiled Arithmetic Operations**

```
test1 %eax, %eax
js L4
L3:
  sarl $3, %eax
  ret
L4:
  add1 $7, %eax
  jmp L3
```

#### **Explanation**

```
if x < 0
   x += 7;
# Arithmetic shift
return x >> 3;
```

- Uses arithmetic shift for int
- For Java Users
  - Arith. shift written as >>

89

#### **Arithmetic: Basic Rules**

- Unsigned ints, 2's complement ints are isomorphic rings: isomorphism = casting
- Left shift
  - Unsigned/signed: multiplication by 2<sup>k</sup>
  - Always logical shift
- Right shift
  - Unsigned: logical shift, div (division + round to zero) by 2<sup>k</sup>
  - Signed: arithmetic shift
    - Positive numbers: div (division + round to zero) by 2k
    - Negative numbers: div (division + round away from zero) by 2<sup>k</sup>
       Use biasing to fix

Carnegie Mellon

# **Negation: Complement & Increment**

■ Claim: Following Holds for 2's Complement

$$\sim x + 1 == -x$$

- Complement
- **■** Complete Proof?

# **Complement & Increment Examples**

x = 15213

	Decimal	Hex	Binary
x	15213	3B 6D	00111011 01101101
~x	-15214	C4 92	11000100 10010010
~x+1	-15213	C4 93	11000100 10010011
У	-15213	C4 93	11000100 10010011

x = 0

	Decimal	Hex	Binary
0	0	00 00	00000000 00000000
~0	-1	FF FF	11111111 11111111
~0+1	0	00 00	00000000 00000000

# **Properties of Unsigned Arithmetic**

- Unsigned Multiplication with Addition Forms Commutative Ring
  - Addition is commutative group
  - Closed under multiplication

$$0 \leq \mathsf{UMult}_w(u, v) \leq 2^w - 1$$

Multiplication Commutative

$$UMult_w(u, v) = UMult_w(v, u)$$

Multiplication is Associative

$$UMult_w(t, UMult_w(u, v)) = UMult_w(UMult_w(t, u), v)$$

1 is multiplicative identity

$$UMult_w(u, 1) = u$$

Multiplication distributes over addtion

$$UMult_w(t, UAdd_w(u, v)) = UAdd_w(UMult_w(t, u), UMult_w(t, v))$$

# Properties of Two's Comp. Arithmetic

- Isomorphic Algebras
  - Unsigned multiplication and addition
    - Truncating to w bits
  - Two's complement multiplication and addition
    - Truncating to w bits
- **■** Both Form Rings
  - Isomorphic to ring of integers mod 2<sup>w</sup>
- Comparison to (Mathematical) Integer Arithmetic
  - Both are rings
  - Integers obey ordering properties, e.g.,

$$u > 0$$
  $\Rightarrow u + v > v$   
 $u > 0, v > 0$   $\Rightarrow u \cdot v > 0$ 

• These properties are not obeyed by two's comp. arithmetic

```
TMax + 1 == TMin
15213 * 30426 == -10030 (16-bit words)
```