

CSE251: System Programming

3. Bits, Bytes, Integers, and Floats

Seongil Wi

Motivation: Great Reality!

- **Question: Is $x^2 \geq 0$?**
 - In an ideal mathematical setting: Yes!
 - But in real computing systems? Let's see.

Motivation: Great Reality!

- In case of int's
 - $40000 * 40000 \rightarrow 1600000000$
 - $50000 * 50000 \rightarrow ?$

Motivation: Great Reality!

- In case of int's
 - $40000 * 40000 \rightarrow 1600000000$
 - $50000 * 50000 \rightarrow ?$

It is important to understand how modern computer systems store , represent, and manipulate data

Today's topic

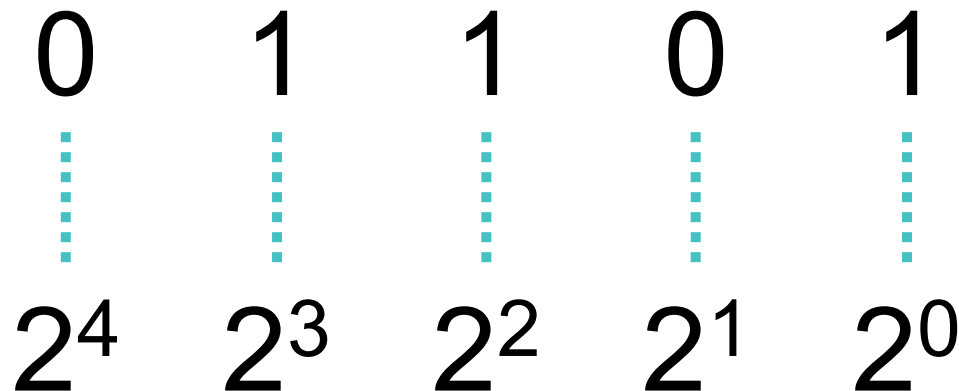
Bits

*Modern computers store and process information represented as **two-valued signals**. These lowly binary digits (i.e., **bits**) form the basis of the digital revolution.*

Binary Representation



- Leads to a simple binary, i.e., base-2 numbering system
 - 0 represents 0
 - 1 represents 1
 - Each “place” represents a power of two
- By encoding/interpreting sets of bits in various ways, we can represent different things:
 - Numbers, characters, operations, ...



Binary Representation



- Leads to a simple binary, i.e., base-2 numbering system
 - 0 represents 0
 - 1 represents 1
 - Each “place” represents a power of two
- By encoding/interpreting sets of bits in various ways, we can represent different things:
 - Numbers, characters, operations, ...

$$\begin{array}{ccccccccc} 0 & 1 & 1 & 0 & 1 & = & 13 & \text{(in decimal)} \\ & & + & & + & & & \\ 2^4 & 2^3 & & 2^2 & 2^1 & & 2^0 & \end{array}$$

Binary Representation: Simple Numbers



- For example, we can count in binary, a base-2 numbering system (e.g., 000, 001, 010, 011, 100, 101, 110, 111, ...)
 - $000 = 0 \cdot 2^2 + 0 \cdot 2^1 + 0 \cdot 2^0 = 0$ (in decimal)
 - $001 = 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 = 1$ (in decimal)
 - $010 = 0 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 = 2$ (in decimal)
 - $011 = 0 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = 3$ (in decimal)
 - Etc.
- (FYI) Consider some base-10 examples:
 - $000 = 0 \cdot 10^2 + 0 \cdot 10^1 + 0 \cdot 10^0$
 - $001 = 0 \cdot 10^2 + 0 \cdot 10^1 + 1 \cdot 10^0$
 - $357 = 3 \cdot 10^2 + 5 \cdot 10^1 + 7 \cdot 10^0$

Hexadecimal (Hex, Base 16 Number System)

- Writing out numbers in binary **takes too many digits** 😞
- We want a way to represent numbers more densely such that fewer digits are required
 - Any power-of-two base provides this property 😊
 - E.g, base-8, base-10, base-16
 - Each has been used by “computer people” over time
 - Hexadecimal (i.e., base-16) is often preferred because it is denser

Hexadecimal (Hex, Base 16 Number System)

- Hexadecimal 00_{16} to FF_{16}
 - Base 16 number representation
 - Use characters '0' to '9' and 'A' to 'F'
- Consider 1A2B in Hexadecimal:
$$1*16^3 + A*16^2 + 2*16^1 + B*16^0$$
$$= 1*16^3 + 10*16^2 + 2*16^1 + 11*16^0 = 6699 \text{ (in decimal)}$$
- The C language prefixes hex with "0x" so they aren't confused with decimal numbers
 - 0xFA1D37B
 - 0Xfa1d37b

Hex	Decimal	Binary
0	0	0000
1	1	0001
2	2	0010
3	3	0011
4	4	0100
5	5	0101
6	6	0110
7	7	0111
8	8	1000
9	9	1001
A	10	1010
B	11	1011
C	12	1100
D	13	1101
E	14	1110
F	15	1111

Boolean Algebra



- Developed by George Boole in 19th century
 - Algebraic representation of logic
 - 1: “True” / 2: “False”

Boolean Algebra: AND (&), OR (|)

And (&)

- $A \& B = 1$
when both $A=1$ and $B=1$

Input		Output
a	b	
0	0	0
0	1	0
1	0	0
1	1	1

Or (|)

- $A | B = 1$
when either $A=1$ or $B=1$

Input		Output
a	b	
0	0	0
0	1	1
1	0	1
1	1	1

Boolean Algebra: NOT (\sim), XOR (\wedge)

Not (\sim)

- $\sim A = 1$
when $A=0$

Input	Output
0	1
1	0

eXclusive-OR (\wedge)

- $A \wedge B = 1$
when either $A=1$ or $B=1$, but not both

Input		Output
a	b	
0	0	0
0	1	1
1	0	1
1	1	0

General Boolean Algebras



- Operate on Bit Vectors
 - Operations applied bitwise

$$\begin{array}{r} 01101001 \\ \& 01010101 \\ \hline 01000001 \end{array}$$

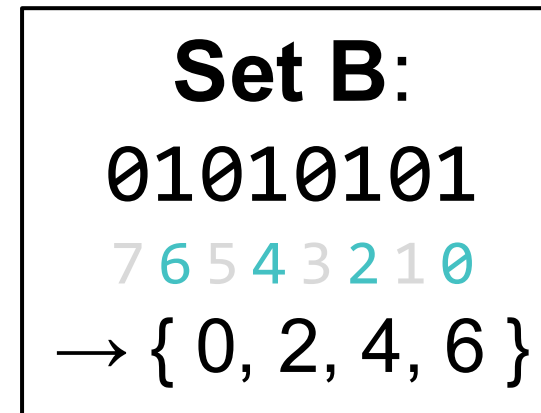
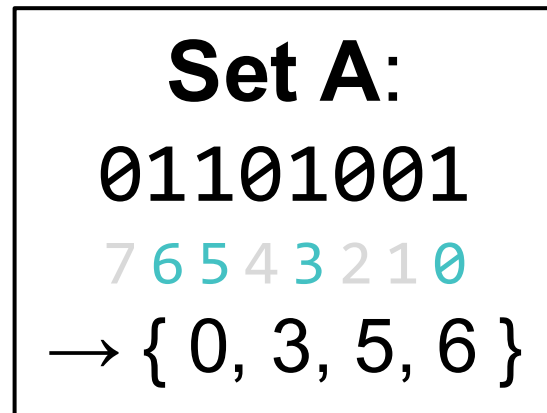
$$\begin{array}{r} 01101001 \\ | 01010101 \\ \hline 01111101 \end{array}$$

$$\begin{array}{r} 01101001 \\ \wedge 01010101 \\ \hline 00111100 \end{array}$$

$$\begin{array}{r} \sim 01010101 \\ \hline 10101010 \end{array}$$

Usage Example: Set Operations

- Representation: $a_j = 1$ if $j \in A$



- Operations

A & B: Intersection

01000001 { 0, 6 }

A | B: Union

01111101 { 0, 2, 3, 4, 5, 6 }

A ^ B: Symmetric difference

00111100 { 2, 3, 4, 5 }

$\sim B$: Complement

10101010 { 1, 3, 5, 7 }

Bit-Level Operations in C



- Operations `&` , `|` , `~` , `^` available in C
 - Apply to any “integral” data type
 - `char` (1 byte)
 - `short` (2 bytes)
 - `int` (4 bytes)
 - ...

(Background) Data Type in C

- Each variable in C has an associated *data type*
- It specifies **the type of data that the variable can store** like integer, character, floating, double, etc

C declaration		Bytes	
Signed	Unsigned	32-bit	64-bit machine
[signed] char	unsigned char	1	1
short	unsigned short	2	2
int	unsigned	4	4
long	unsigned long	4	8
int32_t	uint32_t	4	4
int64_t	uint64_t	8	8
char *		4	8
float		4	4
double		8	8

(Background) Data Type in C *

```
int age = 20;           // integer
float height = 5.7;    // floating-point
double pi = 3.14159;  // double-precision floating-point
char grade = 'A';     // character

printf("Age: %d\n", age);
printf("Height: %.1f\n", height);
printf("Pi: %.5lf\n", pi);
printf("Grade: %c\n", grade);
```

Print "Grade: "

Format string

(Background) Format String in C

- An argument right before “...” (variable-length arguments) that is used to convert C data types into a string (e.g., printf, sprintf, sscanf, syslog, ...)

```
int printf(const char *format, ...);
```

(Background) Format String in C

- An argument right before “...” (variable-length arguments) that is used to convert C data types into a string (e.g., printf, sprintf, sscanf, syslog, ...)

Format string

```
int printf(const char *format, ...);
```

(Background) Format String in C



Format	Meaning
%c	Character output
%d	Decimal output
%x	Hexadecimal output
%u	Unsigned decimal output
%s	String (terminated by \0) output
%f	Float/double output
...	...

(Background) Data Type in C

```
int age = 20;           // integer
float height = 5.7;    // floating-point
double pi = 3.14159;  // double-precision floating-point
char grade = 'A';     // character

printf("Age: %d\n", age);
printf("Height: %.1f\n", height);
printf("Pi: %.5lf\n", pi);
printf("Grade: %c\n", grade);
```

Print "Grade: "

Format string

New line

Output

```
Age: 20
Height: 5.7
Pi: 3.14159
Grade: A
```

Bit-Level Operations in C (Example)



```
char a = 0x69;  
char b = 0x55;  
printf("0x%x\n", a&b);
```

Bit-Level Operations in C (Example)

```
char a = 0x69;  
char b = 0x55;  
printf("0x%x\n", a&b);
```

$0x69 \ \& \ 0x55 \rightarrow 01101001_2 \ \& \ 01010101_2$
 $\rightarrow 01000001_2$
 $\rightarrow 0x41$

Output:

0x41

Logic Operations



- Logic operations: `&&` , `||` , `!`
 - View `0` as “False”
 - Anything `nonzero` as “True”
 - Always return “True” or “False”
 - Example:
 - `!0x41` → `0x00`
 - `!0x00` → `0x01`
 - `0x69 && 0x55` → `0x01`

- Bit-level Operations: `&` , `|` , `~` , `^`

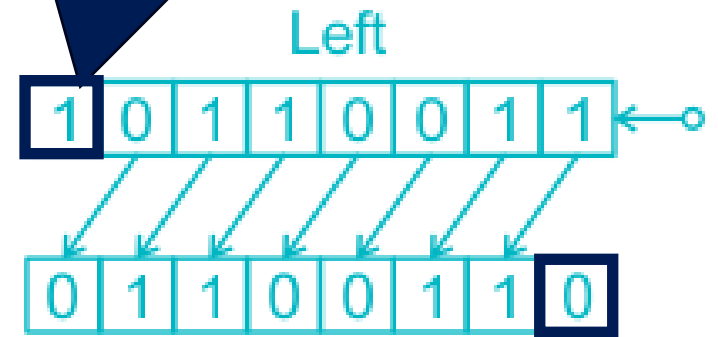
Different from `&&` , `||` , ...

Shift Operations in C

- Left shift: $x \ll y$
 - Shift bit-vector x left y positions
 - Throw away extra bits on left
 - Fill with 0's on right

- Right shift: $x \gg y$
 - Shift bit-vector x right y positions
 - Throw away extra bits on left
 - Logical shift
 - Fill with 0's on left
 - Arithmetic shift
 - Replicate most significant bit (MSB) on left

Throw away extra bits on left

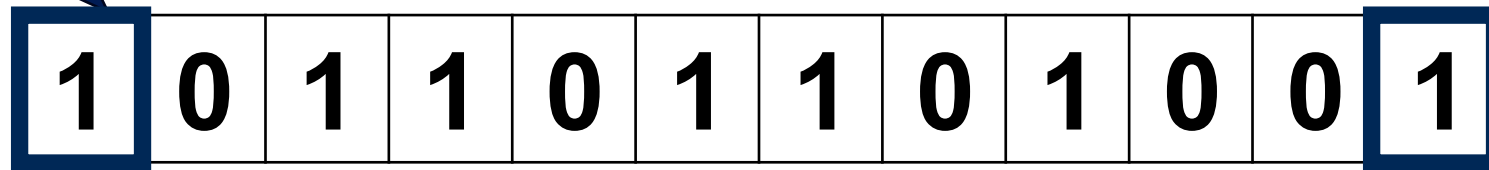


Fill with 0's on right

FYI: LSB and MSB



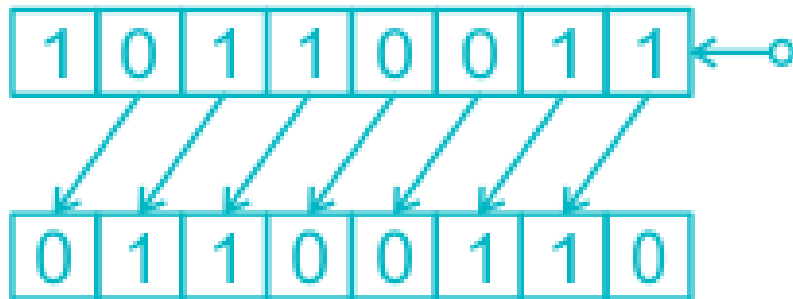
Most Significant Bit (MSB)



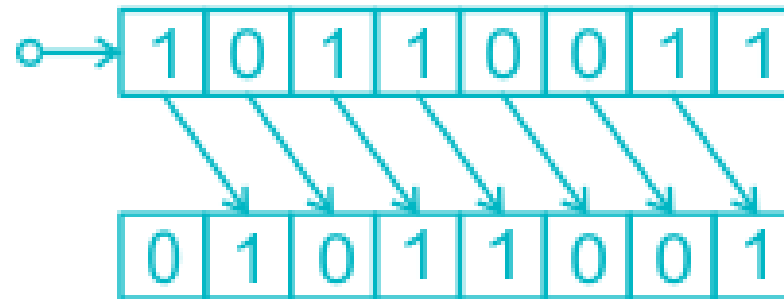
Least Significant Bit (LSB)

Two Types of Shift Operations

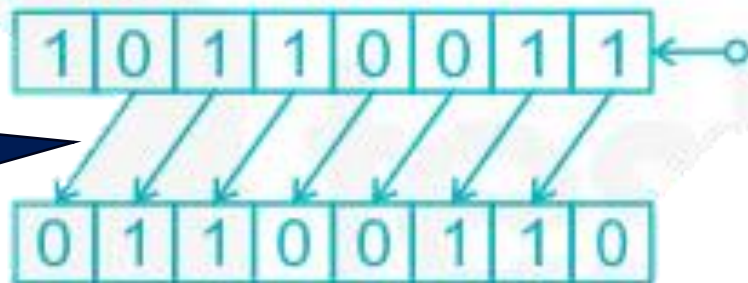
Logical Shift Left



Right



Arithmetic Shift Left



Same as shift left logical

- Shift left by 1 bit \rightarrow multiply by 2
- Shift left by 2 bits \rightarrow multiply by 2^2
- ...

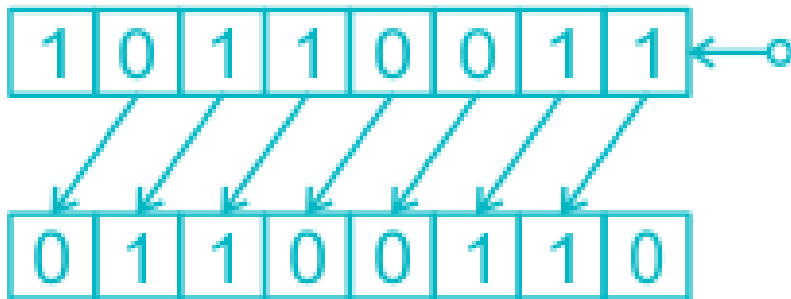
$$00001011_2 = 11_{10}$$

Arithmetic shift
left by 1 bit

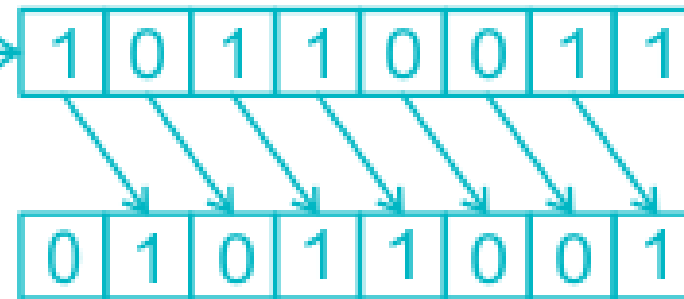
$$00010110_2 = 22_{10}$$

FYI: Logical Shift and Arithmetic Shift

Logical Shift Left



Right



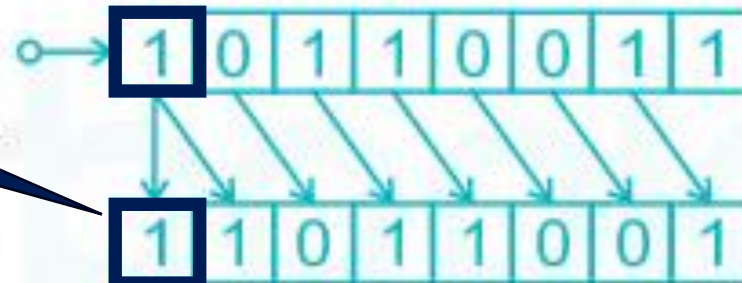
Arithmetic Shift Left



Same as shift left logical

Preserve the sign bit

Right



- Shift left by 1 bit \rightarrow multiply by 2
- Shift left by 2 bits \rightarrow multiply by 2^2
- ...

- Shift right by 1 bit \rightarrow divide by 2
- Shift right by 2 bits \rightarrow divide by 2^2
- ...



Arithmetic Right Shifts vs. Logical Right Shifts

- The C standards do not precisely define which type of right shift should be used for the >> operator
- However, almost all compiler/machine behaves as follows:
 - For signed data, an arithmetic right shift is performed
 - For unsigned data, a logical right shift is performed

Shift Operations in C



```
int a = 0xdeadbeef;  
printf("0x%x\n", a << 4);  
printf("0x%x\n", a >> 4);
```

Output:

```
0xeadbeef0  
0xfdeadbee
```

Bytes



1 Byte = 8 Bits

Show Bytes in C



```
void show_bytes(unsigned char* start, int len) {  
    int i;  
    for (i = 0; i < len; i++)  
        printf("%02x ", start[i]);  
    printf("\n");  
}
```

(Example) Show Bytes in C

```
void show_bytes(unsigned char* start, int len) {
    int i;
    for (i = 0; i < len; i++)
        printf("%02x ", start[i]);
    printf("\n");
}

int main() {
    char b[] = "Hello";
    show_bytes(b, 5);
    return 0;
}
```

0x10e	
0x10d	\0 (0x00)
0x10c	o (0x6f)
0x10b	l (0x6c)
0x10a	l (0x6c)
0x109	e (0x65)
0x108	H (0x48)
0x107	
0x106	

(byte-addressable) Memory

(FYI) ASCII Table

Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char
0	0	[NULL]	32	20	[SPACE]	64	40	@	96	60	`
1	1	[START OF HEADING]	33	21	!	65	41	A	97	61	a
2	2	[START OF TEXT]	34	22	"	66	42	B	98	62	b
3	3	[END OF TEXT]	35	23	#	67	43	C	99	63	c
4	4	[END OF TRANSMISSION]	36	24	\$	68	44	D	100	64	d
5	5	[ENQUIRY]	37	25	%	69	45	E	101	65	e
6	6	[ACKNOWLEDGE]	38	26	&	70	46	F	102	66	f
7	7	[BELL]	39	27	'	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	H	104	68	h
9	9	[HORIZONTAL TAB]	41	29)	73	49	I	105	69	i
10	A	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	B	[VERTICAL TAB]	43	2B	+	75	4B	K	107	6B	k
12	C	[FORM FEED]	44	2C	,	76	4C	L	108	6C	l
13	D	[CARRIAGE RETURN]	45	2D	-	77	4D	M	109	6D	m
14	E	[SHIFT OUT]	46	2E	.	78	4E	N	110	6E	n
15	F	[SHIFT IN]	47	2F	/	79	4F	O	111	6F	o
16	10	[DATA LINK ESCAPE]	48	30	0	80	50	P	112	70	p
17	11	[DEVICE CONTROL 1]	49	31	1	81	51	Q	113	71	q
18	12	[DEVICE CONTROL 2]	50	32	2	82	52	R	114	72	r
19	13	[DEVICE CONTROL 3]	51	33	3	83	53	S	115	73	s
20	14	[DEVICE CONTROL 4]	52	34	4	84	54	T	116	74	t
21	15	[NEGATIVE ACKNOWLEDGE]	53	35	5	85	55	U	117	75	u
22	16	[SYNCHRONOUS IDLE]	54	36	6	86	56	V	118	76	v
23	17	[ENG OF TRANS. BLOCK]	55	37	7	87	57	W	119	77	w
24	18	[CANCEL]	56	38	8	88	58	X	120	78	x
25	19	[END OF MEDIUM]	57	39	9	89	59	Y	121	79	y
26	1A	[SUBSTITUTE]	58	3A	:	90	5A	Z	122	7A	z
27	1B	[ESCAPE]	59	3B	;	91	5B	[123	7B	{
28	1C	[FILE SEPARATOR]	60	3C	<	92	5C	\	124	7C	
29	1D	[GROUP SEPARATOR]	61	3D	=	93	5D]	125	7D	}
30	1E	[RECORD SEPARATOR]	62	3E	>	94	5E	^	126	7E	~
31	1F	[UNIT SEPARATOR]	63	3F	?	95	5F	_	127	7F	[DEL]

(Example) Show Bytes in C

```
void show_bytes(unsigned char* start, int len) {
    int i;
    for (i = 0; i < len; i++)
        printf("%02x ", start[i]);
    printf("\n");
}
```

```
int main() {
    char b[] = "Hello";
    show_bytes(b, 5);
    return 0;
}
```

Output

48

start 0x107
 b 0x107

0x10e	
0x10d	\0 (0x00)
0x10c	o (0x6f)
0x10b	l (0x6c)
0x10a	l (0x6c)
0x109	e (0x65)
0x108	H (0x48)
0x107	
0x106	

(byte-addressable) Memory

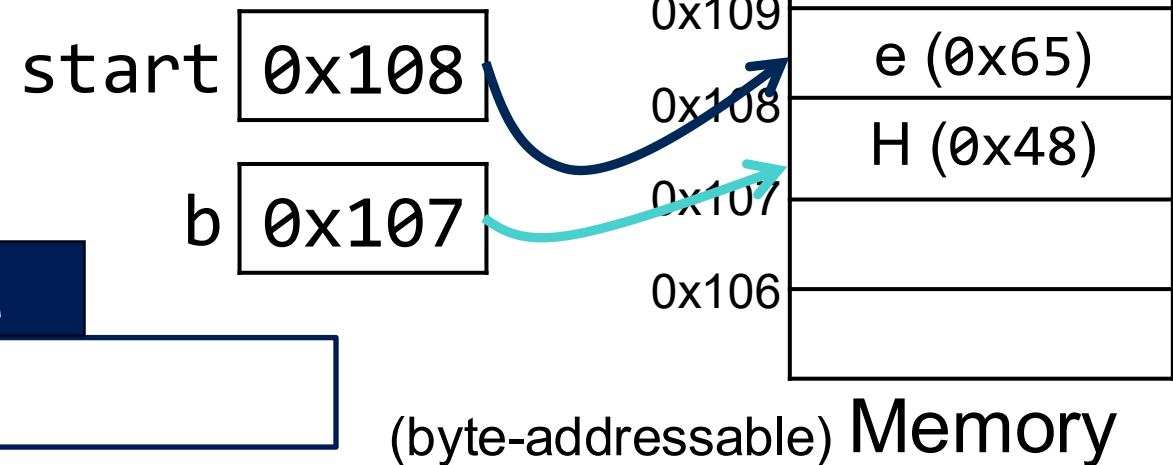
(Example) Show Bytes in C

```
void show_bytes(unsigned char* start, int len) {
    int i;
    for (i = 0; i < len; i++)
        printf("%02x ", start[i]);
    printf("\n");
}
```

```
int main() {
    char b[] = "Hello";
    show_bytes(b, 5);
    return 0;
}
```

Output

48 65



(Example) Show Bytes in C

```
void show_bytes(unsigned char* start, int len) {
    int i;
    for (i = 0; i < len; i++)
        printf("%02x ", start[i]);
    printf("\n");
}
```

```
int main() {
    char b[] = "Hello";
    show_bytes(b, 5);
    return 0;
}
```

Output

48 65 6c

start 0x109

b 0x107

0x10e	
0x10d	\0 (0x00)
0x10c	o (0x6f)
0x10b	l (0x6c)
0x10a	l (0x6c)
0x109	e (0x65)
0x108	H (0x48)
0x107	
0x106	

(byte-addressable) Memory

(Example) Show Bytes in C

```
void show_bytes(unsigned char* start, int len) {
    int i;
    for (i = 0; i < len; i++)
        printf("%02x ", start[i]);
    printf("\n");
}
```

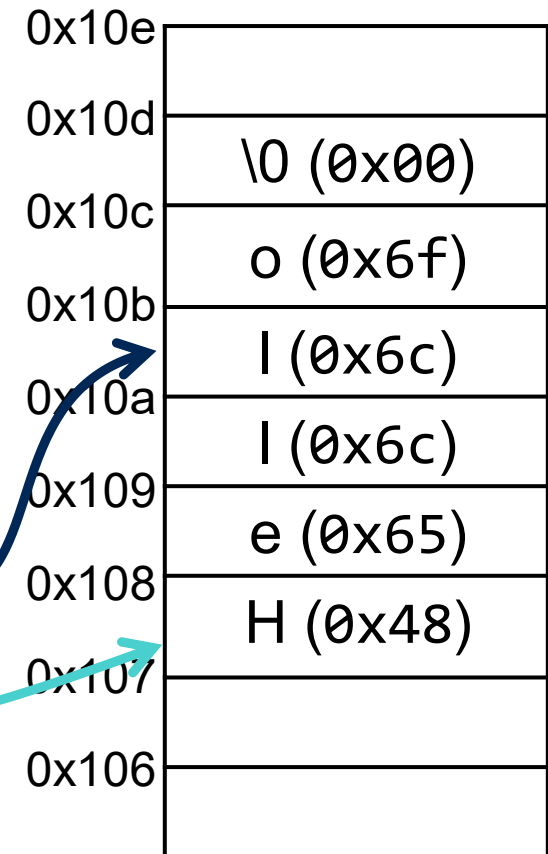
```
int main() {
    char b[] = "Hello";
    show_bytes(b, 5);
    return 0;
}
```

Output

48 65 6c 6c

start 0x10a

b 0x107



(byte-addressable) Memory

(Example) Show Bytes in C

```
void show_bytes(unsigned char* start, int len) {
    int i;
    for (i = 0; i < len; i++)
        printf("%02x ", start[i]);
    printf("\n");
}
```

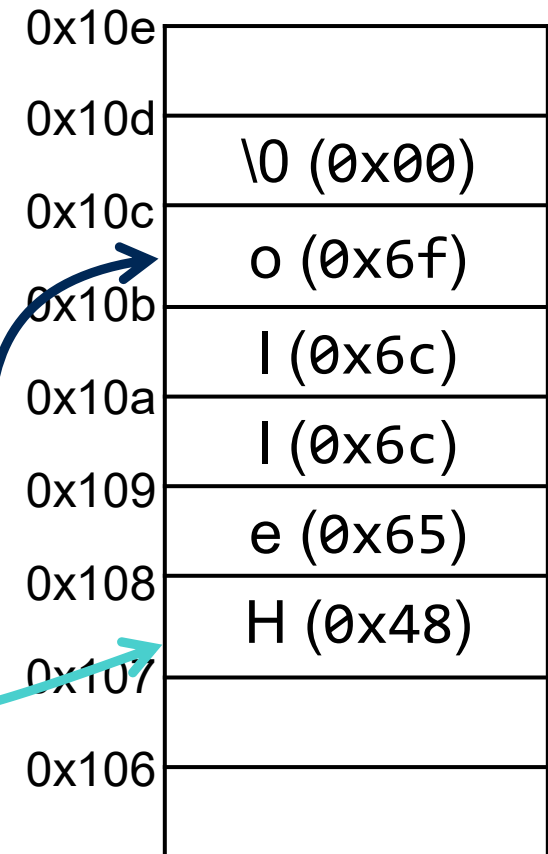
```
int main() {
    char b[] = "Hello";
    show_bytes(b, 5);
    return 0;
}
```

Output

48 65 6c 6c 6f

start 0x10b

b 0x107



(byte-addressable) Memory

How About This One?



```
void show_bytes(unsigned char* start, int len) {  
    int i;  
    for (i = 0; i < len; i++)  
        printf("%02x ", start[i]);  
    printf("\n");  
}  
  
int main() {  
    int a = 98303; // 0x00017fff  
    show_bytes(a, 4);  
    return 0;  
}
```

How About This One?



```
void show_bytes(unsigned char* start, int len) {  
    int i;  
    for (i = 0; i < len; i++)  
        printf("%02x ", start[i]);  
    printf("\n");  
}
```

```
int main() {  
    int a = 98303; // 0x00017fff  
    show_bytes(a, 4);  
    return 0;  
}
```

Output

ff 7f 01 00

0x10e	
0x10d	
0x10c	
0x10b	
0x10a	0x00
0x109	0x01
0x108	0x7f
0x107	0xff
0x106	

(byte-addressable) Memory

How About This One?



```
void show_bytes(unsigned char* start, int len) {
    int i;
    for (i = 0; i < len; i++)
```

0x10e

**Why is 98303, which is 0x17fff, represented by ff 7f 01 00?
The answer is endianness.**

```
int main() {
    int a = 98303; // 0x00017fff
    show_bytes(a, 4);
    return 0;
}
```

Output

ff 7f 01 00

0x108	0x7f
0x107	0xff
0x106	

(byte-addressable) Memory

Endianness



- The order in which a sequence of bytes are stored in memory
- Big Endian = The MSB goes to the lowest address
- Little Endian = The LSB goes to the lowest address

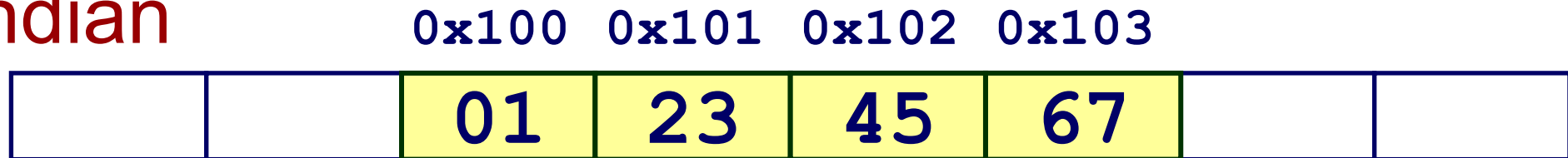
***X86 (and X86_64) uses
Little Endian***

Byte Ordering Example

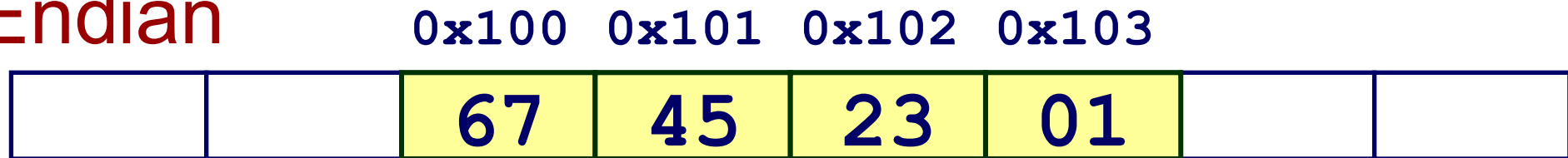


- Variable x has 4-byte value of 0x01234567
- Address given by &x is 0x100

Big Endian



Little Endian



Byte Ordering



- Big Endian: Sun (Oracle SPARC), PPC Mac, Internet
 - Least significant byte has highest address
- Little Endian: x86, ARM processors running Android, iOS, and Linux
 - Least significant byte has lowest address
- Important notes:
 - Bits are not reversed, as the low order bit is the reference point.
 - Doesn't affect chars, or strings (arrays of chars), as chars are only one byte

Integers

Unsigned Binary Integers

- Given an n-bit number

$$x = x_{n-1}2^{n-1} + x_{n-2}2^{n-2} + \dots + x_12^1 + x_02^0$$

- Range: 0 to $2^n - 1$
- Example

$$\begin{aligned} &0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1011_2 \\ &= 0 + \dots + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 \\ &= 0 + \dots + 8 + 0 + 2 + 1 = 11_{10} \end{aligned}$$



How to represent negative numbers?

- Range of 32 bits ($n=32$): 0 to +4,294,967,295
- Range of 64 bits ($n=64$): 0 to +18,446,744,073,709,551,615

Approach #1: Sign Magnitude

$$001 = 1$$

$$101 = -1$$

Approach #1: Sign Magnitude

$$001 = 1$$

$$101 = -1$$

MSB represents the sign

- 0: +
- 1: -

Approach #1: Sign Magnitude

$$\begin{array}{l} 001 = 1 \\ 101 = -1 \end{array}$$

MSB represents the sign

- 0: +
- 1: -

The rest of the bits represent the magnitude



Problems?

Problems of Sign Magnitude



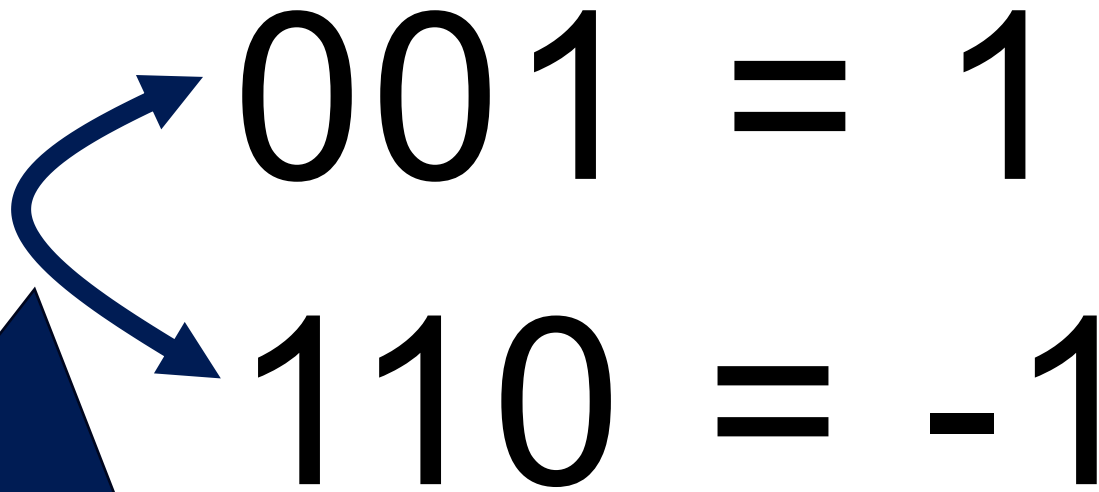
$$001 = 1$$

$$+ 110 = -2$$

$$111 \neq -1$$

Inconsistencies in operations

Approach #2: One's Complement


$$001 = 1$$
$$110 = -1$$

To represent a negative number, all the bits of the corresponding positive number are inverted

Approach #2: One's Complement

$$001 = 1$$

$$+ 101 = -2$$

$$110 = -1$$

Provides ease of operations in terms of computing

Approach #2: One's Complement

$$\begin{array}{r} 001 = 1 \\ + 101 = -2 \\ \hline \end{array}$$

$$110 = -1$$



Problems?

Provides ease of operations in terms of computing

Problems of One's Complement (+ Sign Magnitude)

59

Sign Magnitude

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

$$100 = -0$$

$$101 = -1$$

$$110 = -2$$

$$111 = -3$$

One's Complement

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

$$100 = -3$$

$$101 = -2$$

$$110 = -1$$

$$111 = -0$$

Problems of One's Complement (+ Sign Magnitude)

Sign Magnitude

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

$$100 = -0$$

$$101 = -1$$

$$110 = -2$$

$$111 = -3$$

One's Complement

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

$$100 = -3$$

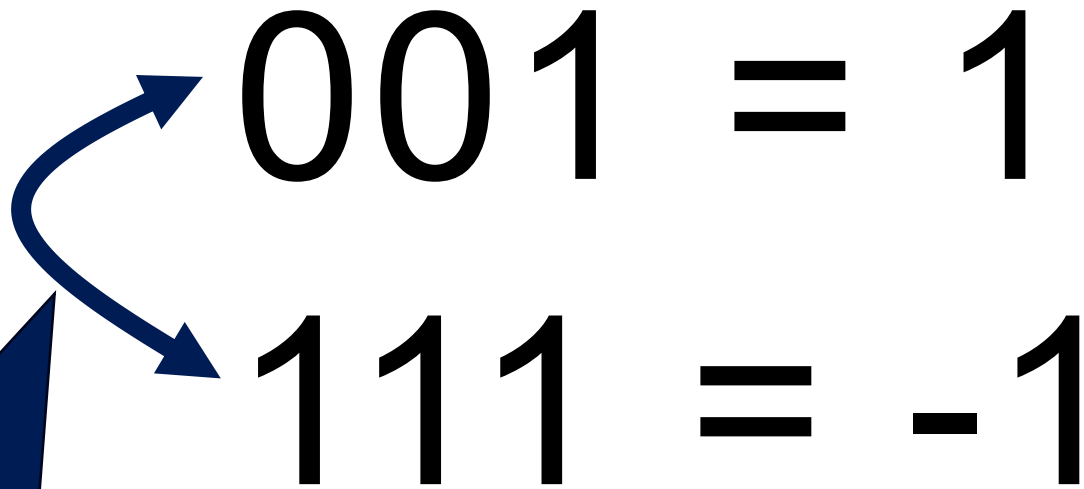
$$101 = -2$$

$$110 = -1$$

$$111 = -0$$

Two representations
for zero

Approach #3: Two's Complement


$$001 = 1$$
$$111 = -1$$

= One's Complement + 1

Approach #3: Two's Complement

Sign Magnitude

000 = +0
001 = +1
010 = +2
011 = +3
100 = -0
101 = -1
110 = -2
111 = -3

One's Complement

000 = +0
001 = +1
010 = +2
011 = +3
100 = -3
101 = -2
110 = -1
111 = -0

Two's Complement

000 = +0
001 = +1
010 = +2
011 = +3
100 = -4
101 = -3
110 = -2
111 = -1

Approach #3: Two's Complement

Sign Magnitude

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

Range: -3 ~ +3

$$100 = -0$$

$$101 = -1$$

$$110 = -2$$

$$111 = -3$$

One's Complement

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

Range: -3 ~ +3

$$100 = -3$$

$$101 = -2$$

$$110 = -1$$

$$111 = -0$$

Two's Complement

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

Range: -4 ~ +3

$$100 = -4$$

$$101 = -3$$

$$110 = -2$$

$$111 = -1$$

Approach #3: Two's Complement

Sign Magnitude

$$000 = +0$$

$$001 = +1$$

One's Complement

$$000 = +0$$

$$001 = +1$$

Two's Complement

$$000 = +0$$

$$001 = +1$$

Provides ease of operations and
a unique representation for zero

$$100 = -0$$

$$101 = -1$$

$$110 = -2$$

$$111 = -3$$

$$100 = -3$$

$$101 = -2$$

$$110 = -1$$

$$111 = -0$$

$$100 = -4$$

$$101 = -3$$

$$110 = -2$$

$$111 = -1$$

2's complement code is the most widely used representation of signed numbers in computer systems 😊

Negative Number in C



```
int i = -1;  
printf("0x%x\n", i);
```

Output:

```
0xffffffff
```

2's-Complement Signed Integers

- Given an n-bit number

$$X = -x_{n-1}2^{n-1} + x_{n-2}2^{n-2} + \dots + x_12^1 + x_02^0$$

- Range: -2^{n-1} to $+2^n - 1$
- Example

$$\begin{aligned} & 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1100_2 \\ &= -1 \times 2^{31} + 1 \times 2^{30} + \dots + 1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0 \\ &= -2,147,483,648 + 2,147,483,644 = -4_{10} \end{aligned}$$

- Range of 32 bits ($n=32$): $-2,147,483,648$ to $+2,147,483,647$
- Range of 64 bits ($n=64$): $-9,223,372,036,854,775,808$ to $+9,223,372,036,854,775,807$

2's-Complement: Number Range

- 32-bit signed numbers:

0000	0000	0000	0000	0000	0000	0000	0000	$_{two} = 0_{ten}$
0000	0000	0000	0000	0000	0000	0000	0001	$_{two} = + 1_{ten}$
0000	0000	0000	0000	0000	0000	0000	0002	$_{two} = + 2_{ten}$
...								
0111	1111	1111	1111	1111	1111	1111	1110	$_{two} = + 2,147,483,646_{ten}$
0111	1111	1111	1111	1111	1111	1111	1111	$_{two} = + 2,147,483,647_{ten}$
1000	0000	0000	0000	0000	0000	0000	0000	$_{two} = - 2,147,483,648_{ten}$
1000	0000	0000	0000	0000	0000	0000	0001	$_{two} = - 2,147,483,647_{ten}$
...								
1111	1111	1111	1111	1111	1111	1111	1110	$_{two} = - 2_{ten}$
1111	1111	1111	1111	1111	1111	1111	1111	$_{two} = - 1_{ten}$

maxint

minint

Integer Overflow

Example:

$maxint + 1 = minint$ ($\neq maxint + 1$ in real world)

$$\begin{array}{l}
 0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_{two} = +\ 2,147,483,647_{ten} \\
 1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_{two} = -\ 2,147,483,648_{ten}
 \end{array}$$

maxint

minint

Integer Overflow in C



```
int i = 2147483647;  
  
printf("%d\n", i);  
printf("%d\n", i + 1);
```

Output:

```
2147483647  
-2147483648
```

Zero-Day Alert: Google Chrome Under Active Attack, Exploiting New Vulnerability

📅 Nov 29, 2023 👤 Ravie Lakshmanan

Zero-Day / Web Browser

— Trending News

Google has rolled out security updates to fix seven security issues in its Chrome browser, including a zero-day that has come under active exploitation in the wild.

Tracked as CVE-2023-6345, the high-severity vulnerability has been described as an **integer overflow** bug in Skia, an open source 2D graphics library.



Another Benefit of 2's-Complement: Simplifies Overflow Detection *

Basic idea: If ...

Positive Number + Positive Number = Negative Number

Or...

Negative Number + Negative Number = Positive Number

Overflow is occurred!

Another Benefit of 2's-Complement: Simplifies Overflow Detection *

From a logic circuit perspective:

$$\begin{array}{r} 010 = 2 \\ + 011 = 3 \\ \hline 101 = -3 \quad (\neq 5) \end{array}$$

Another Benefit of 2's-Complement: Simplifies Overflow Detection *

From a logic circuit **C_{out}** effective: **C_{in}**

Carry in (C _{in})	Carry out (C _{out})	Overflow occurred
0	0	X
0	1	O
1	0	O
1	1	X

0 1

010 = 2

+ 011 = 3

101 = -3 (≠5)

Another Benefit of 2's-Complement: Simplifies Overflow Detection *

From a logic circuit C_{out} effective: C_{in}

Carry in (C_{in})	Carry out (C_{out})	Overflow occurred
0	0	X
0	1	O
1	0	O
1	1	X

Exception:
Overflow has occurred

$$\begin{array}{r}
 \begin{array}{c} 0 \\ \downarrow \\ 010 = 2 \end{array} \\
 + \begin{array}{c} 1 \\ \downarrow \\ 011 = 3 \end{array} \\
 \hline
 101 = -3 (\neq 5)
 \end{array}$$

2's-Complement Signed Integers

- 31st-bit (MSB) is sign bit
 - 1 for negative numbers
 - 0 for non-negative numbers
- $-(-2^{n-1})$ can't be represented
- Non-negative numbers have the same unsigned and 2s-complement representation
- Some specific numbers
 - 0: 0000 0000 ... 0000
 - -1: 1111 1111 ... 1111
 - Most-negative: 1000 0000 ... 0000
 - Most-positive: 0111 1111 ... 1111

Signed Negation



- Complement and add 1
 - Complement means $1 \rightarrow 0, 0 \rightarrow 1$

$$x + \bar{x} = 1111\dots111_2 = -1$$

$$\bar{x} + 1 = -x$$

- Example: negate +2
 - +2 = 0000 0000 ... 0010₂
 - 2 = 1111 1111 ... 1101₂ + 1
 - = 1111 1111 ... 1110₂

Sign Extension



- Representing a number using more bits
 - Preserve the numeric value
- **Replicate the sign bit to the left**
 - c.f. unsigned extension: extend with 0s
- Examples: 8-bit to 16-bit
 - +2: 0000 0010 => 0000 0000 0000 0010
 - 2: 1111 1110 => 1111 1111 1111 1110

Sign Extension in C



```
char c = 0x80;    // equivalent to -128
int i = c;        // sign extension
printf("0x%x\n", i);
```

Output:

0xffffffff80

Float / Double

Floating-point Number: Motivation



We need a way to represent ...

- *Infinite decimal*
(e.g., 3.1415926535...)
- *Very small numbers*
- *Very large numbers*

from computer!

**Solution: Floating-point
Number Representation**

Floating-point Number: Motivation

We need a way to represent ...

- *Infinite decimal*
(e.g., 3.1415926535...)

→
Approximate value

3.1415

- *Very small numbers*

→
Floating decimal point

0.001 × 10⁻²⁰

- *Very large numbers*

→
Floating decimal point

3.15576 × 10¹⁹

Can be represented with a limited number of bits!

Solution: Floating-point Number Representation

Floating-point Number: Notations

- **Scientific notation:** renders numbers with a single digit to the *left* of the point
 - Example: **7**.15576 x 10⁴, **0**.314 x 10¹
- **Normalized scientific notation:** scientific notation that has **no leading 0s**
 - Example: **7**.15576 x 10⁴, **3**.14 x 10⁰

Normalized



Floating-point Number: Notations

- **Scientific notation:** renders numbers with a single digit to the *left* of the point

$$1 \leq \textit{Significand} < 10$$

- **Normalized scientific notation:** scientific notation that has *no leading 0s*

–Example: 7.15576 x 10⁴, 3.14 x 10⁰

Significand

Base

Exponent

(sign) x significand x base^{exponent}

Floating-point Number: Binary

- **Scientific notation:** renders numbers with a single digit to the *left* of the point

$$1 \leq \textit{Significand} < 2$$

- **Normalized scientific notation:** scientific notation that has *no leading 0s* (= *Always 1 to the left of the point*)

– Example: 1.1×2^{-1} ($0.75_{\text{ten}} \rightarrow 0.5 + 0.25$
 $\rightarrow 0.11_{\text{two}} \rightarrow 1.1 \times 2^{-1}$)

2

(sign) x significand x base^{exponent}

Floating-point Number: Binary

- **Scientific notation:** renders numbers with a single digit to the *left* of the point

$$1 \leq \textit{Significand} < 2$$

How are floating-point numbers represented in computers?

$$\rightarrow 0.11_{\text{two}} \rightarrow 1.1 \times 2^{-1}$$

2

$$(\textit{sign}) \times \textit{significand} \times \textit{base}^{\textit{exponent}}$$

IEEE 754 Floating-point Standard

- Developed in response to *divergence of representations*

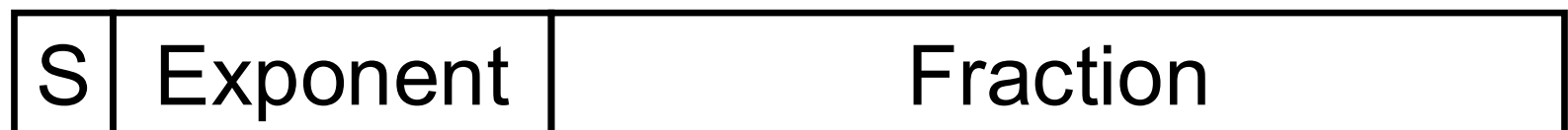
Divergence of representations

$$0.11_{\text{two}} = 1.1_{\text{two}} \times 2^{-1} = 11_{\text{two}} \times 2^{-2}$$

Normalized representation

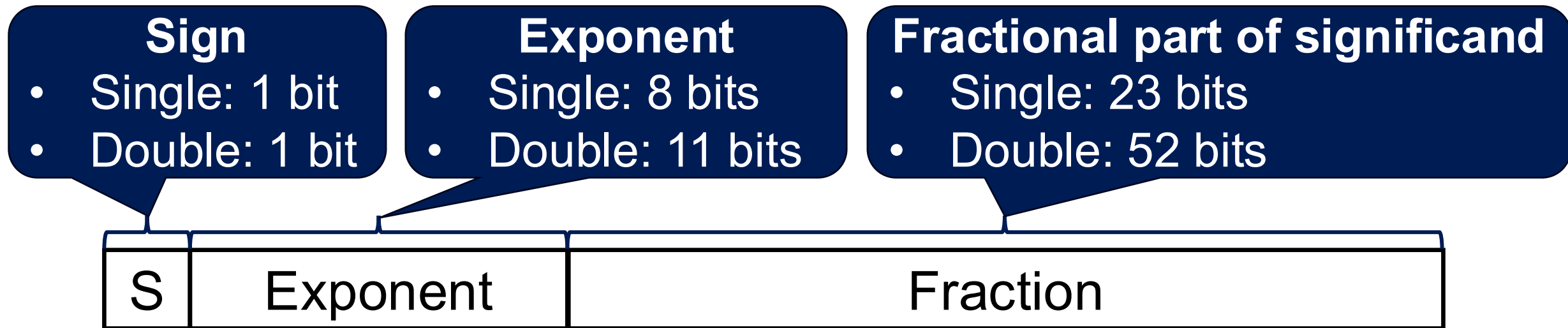
$$1.1_{\text{two}} \times 2^{-1}$$

IEEE 754 representation



IEEE 754 Floating-point Standard

- Two representations
 - **Single precision (32-bit)**: type `float` in C
 - **Double precision (64-bit)**: type `double` in C



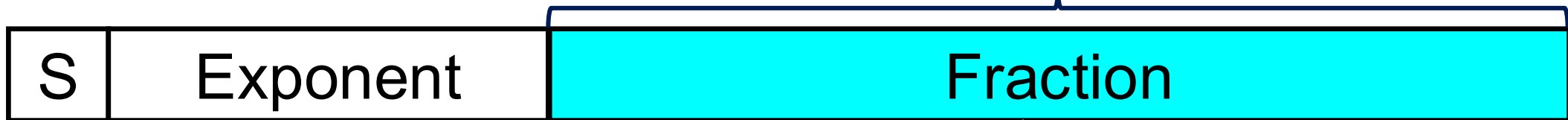
$$+1.1 \times 2^{-1} \text{ two}$$

IEEE 754 Floating-point Standard

- Two representations
 - **Single precision (32-bit)**: type `float` in C
 - **Double precision (64-bit)**: type `double` in C

Fractional part of significand

- Single: 23 bits
- Double: 52 bits

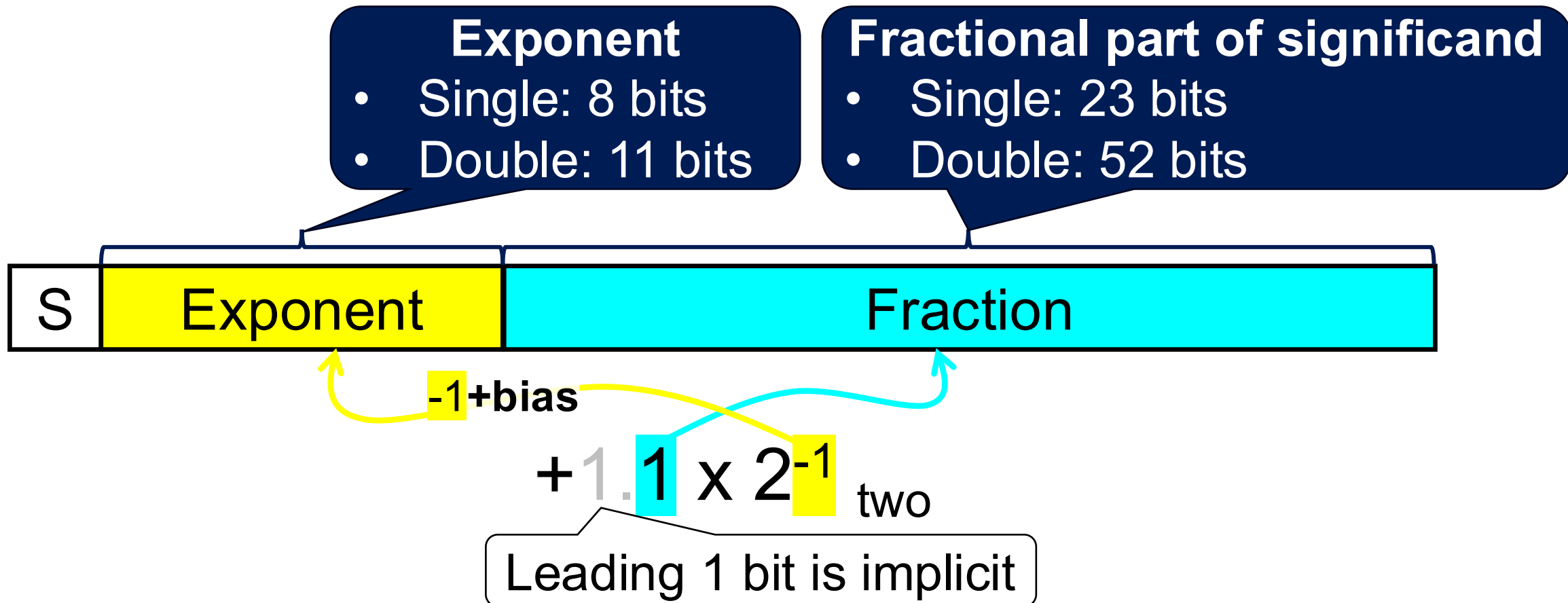


$$+1.\mathbf{1} \times 2^{-1} \text{ two}$$

Leading 1 bit is implicit

IEEE 754 Floating-point Standard

- Two representations
 - **Single precision (32-bit)**: type `float` in C
 - **Double precision (64-bit)**: type `double` in C



Exponent: Why Biased?



- To make sorting easier
- **Bias:** *127 for single precision and 1023 for double precision*

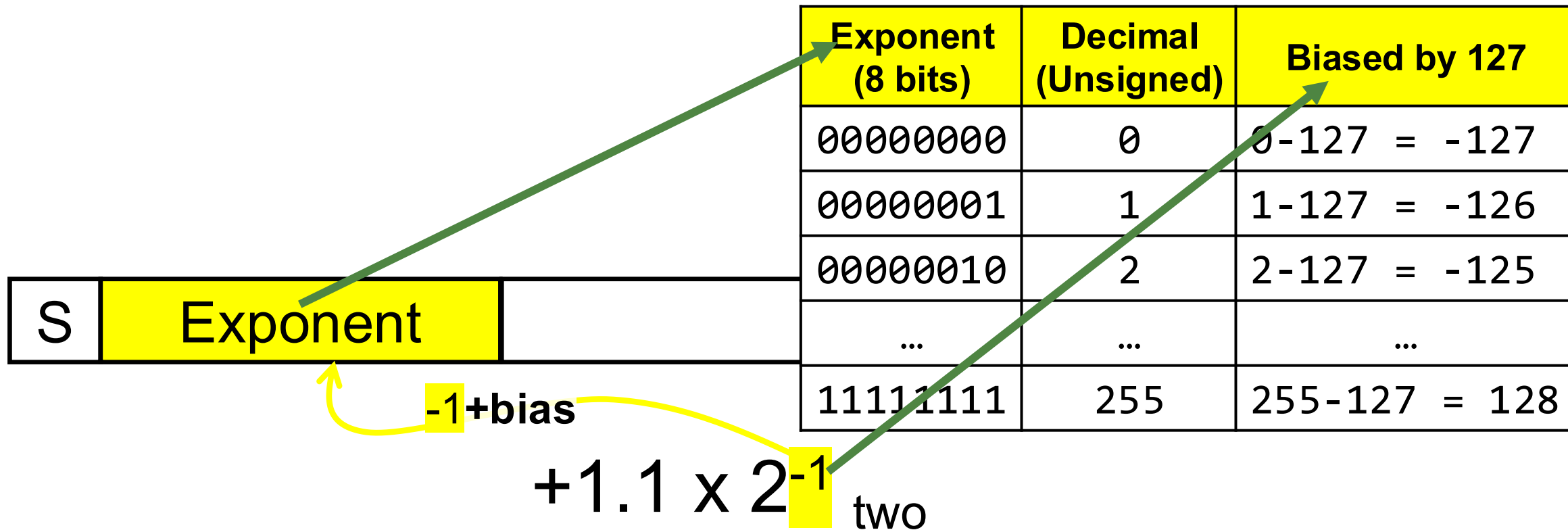


$-1 + \text{bias}$

$+1.1 \times 2^{-1}$ two

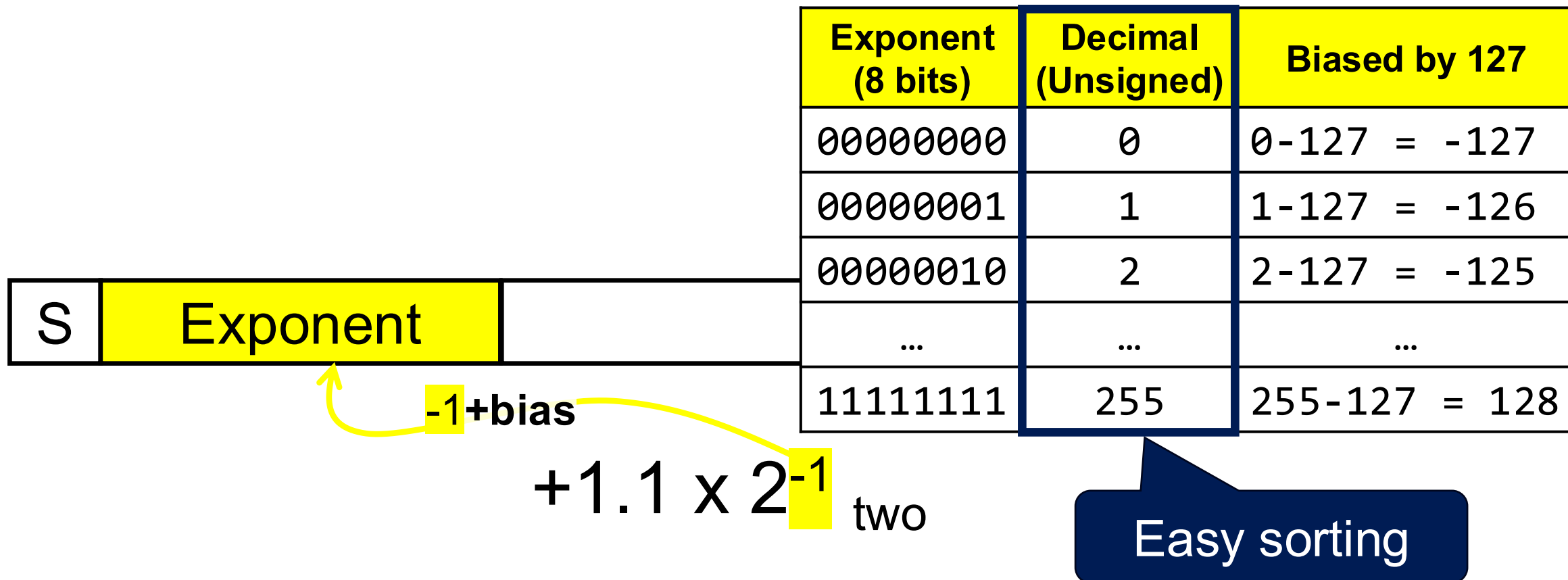
Exponent: Why Biased?

- To make sorting easier
- **Bias:** *127 for single precision and 1023 for double precision*



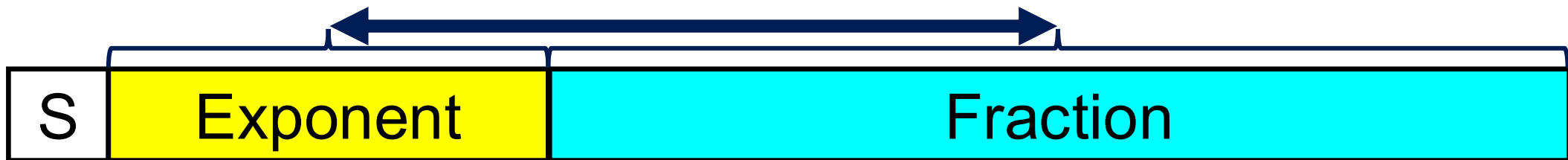
Exponent: Why Biased?

- To make sorting easier
- **Bias:** *127 for single precision and 1023 for double precision*



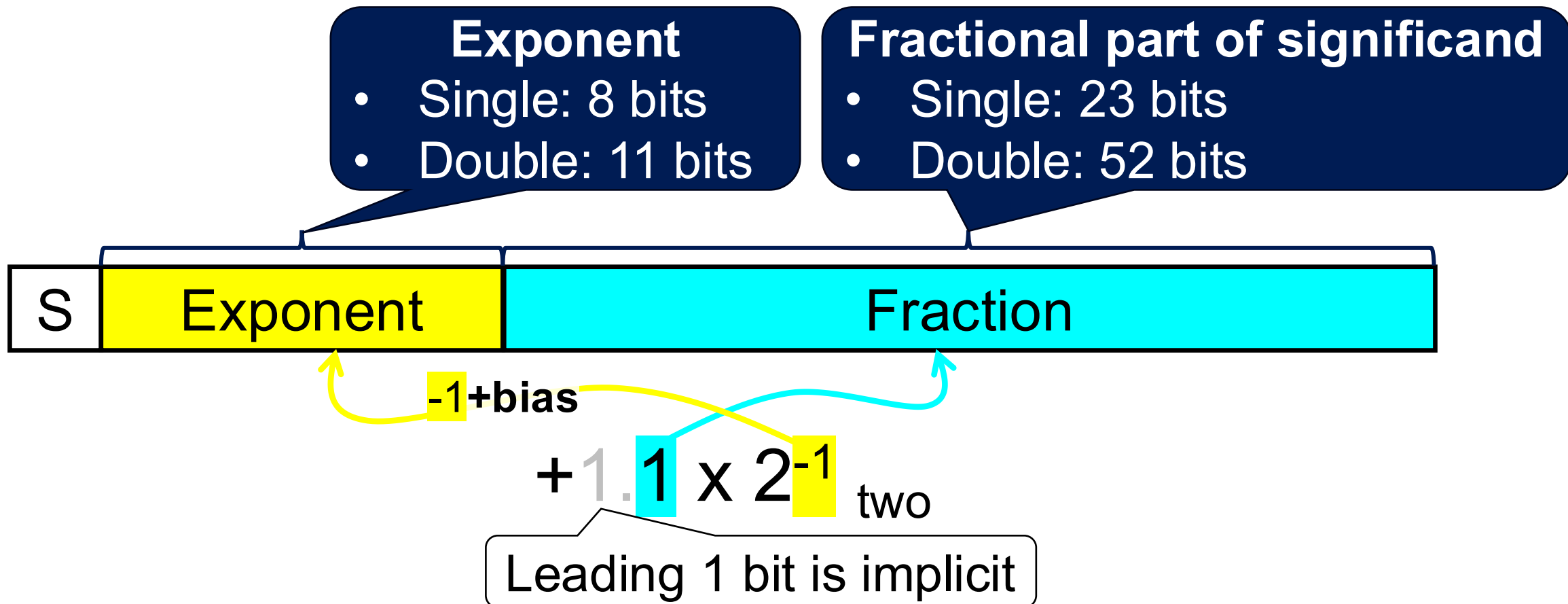
Tradeoff between Precision and Range

- Increasing *the size of the fraction* enhances *the precision*
 - Shorter length: $1.\mathbf{1} \times 2^{-1}$
 - Longer length: $1.\mathbf{10110} \times 2^{-1}$
- Increasing *the size of the exponent* increases *the range*
 - Shorter length: $1.1 \times 2^{\mathbf{2}}$
 - Longer length : $1.1 \times 2^{\mathbf{23}}$



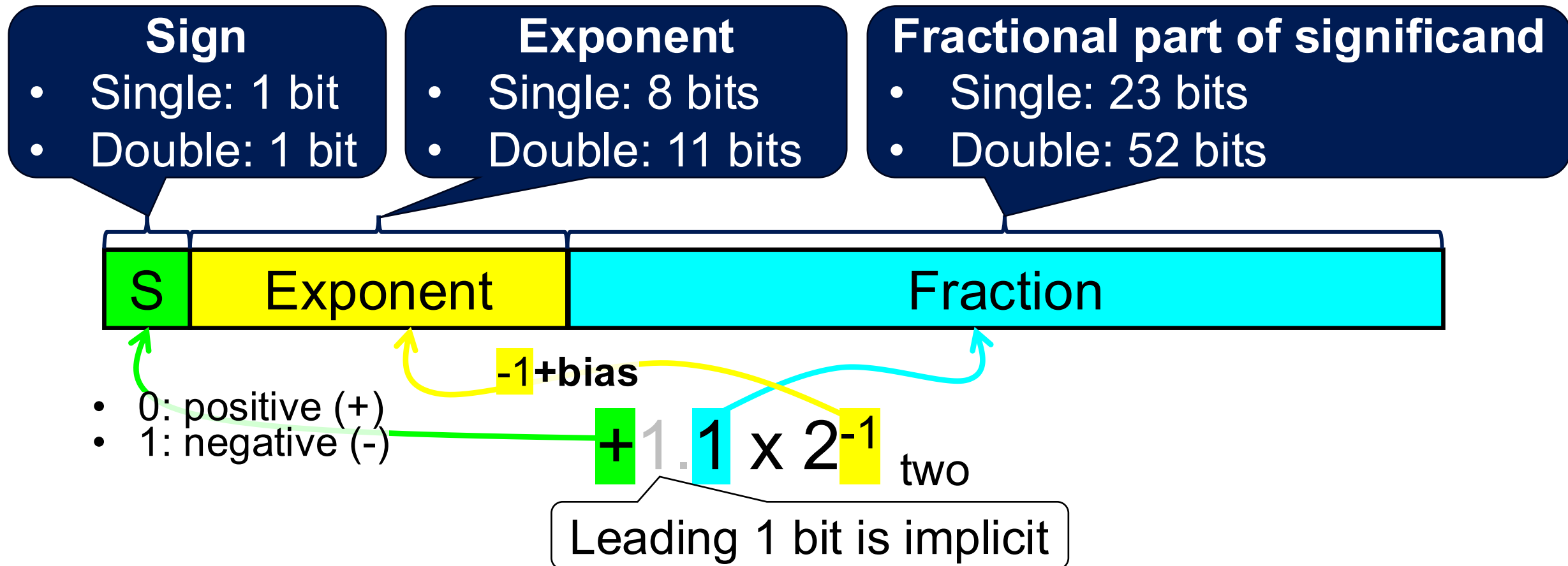
IEEE 754 Floating-point Standard

- Two representations
 - **Single precision (32-bit)**: type `float` in C
 - **Double precision (64-bit)**: type `double` in C



IEEE 754 Floating-point Standard

- Two representations
 - **Single precision (32-bit)**: type `float` in C
 - **Double precision (64-bit)**: type `double` in C



IEEE 754 Floating-point Standard: Summary

$$(-1)^{\text{sign}} \times (1 + \text{fraction}) \times \text{base}^{\text{exponent} - \text{bias}}$$



Single-precision Range



- Exponents 00000000 and 11111111 reserved
- **Smallest value**
 - Exponent: 00000001
⇒ actual exponent = $1 - 127 = -126$
 - Fraction: 000...00 ⇒ significand = 1.0
 - $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$
- **Largest value**
 - Exponent: 11111110
⇒ actual exponent = $254 - 127 = +127$
 - Fraction: 111...11 ⇒ significand ≈ 2.0
 - $\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$

Used for
special cases

Double-precision Range



- Exponents 0000...00 and 1111...11 reserved

- **Smallest value**

- Exponent: 000000000001

- \Rightarrow actual exponent = $1 - 1023 = -1022$

- Fraction: 000...00 \Rightarrow significand = 1.0

- $\Rightarrow \pm 1.0 \times 2^{-1022} \approx \pm 2.2 \times 10^{-308}$

- **Largest value**

- Exponent: 11111111110

- \Rightarrow actual exponent = $2046 - 1023 = +1023$

- Fraction: 111...11 \Rightarrow significand ≈ 2.0

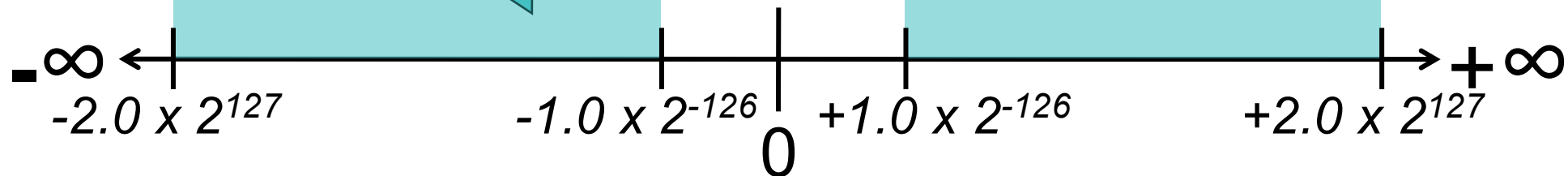
- $\Rightarrow \pm 2.0 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$

Used for
special cases

Overflow and Underflow



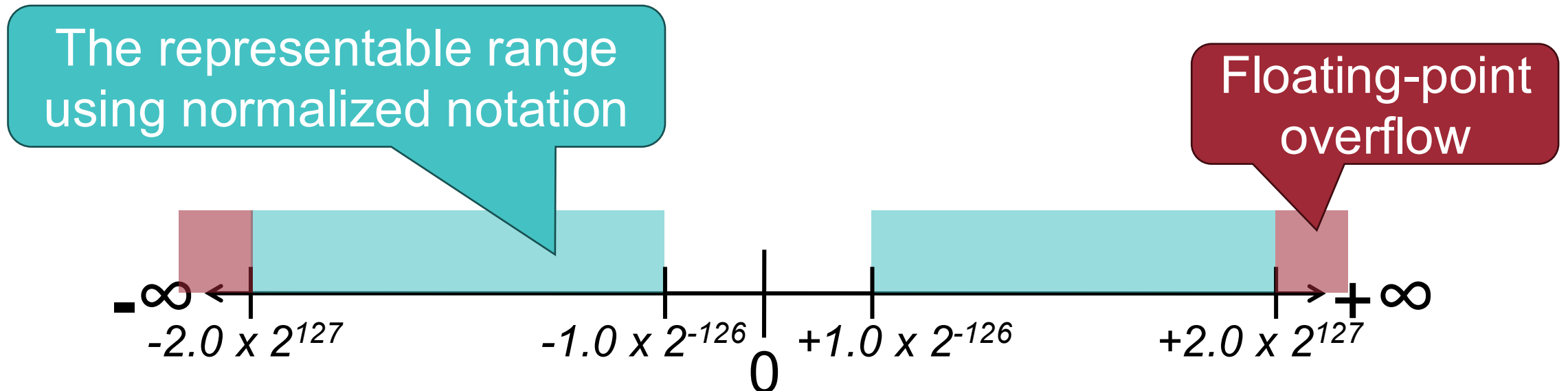
The representable range using normalized notation



Overflow and Underflow



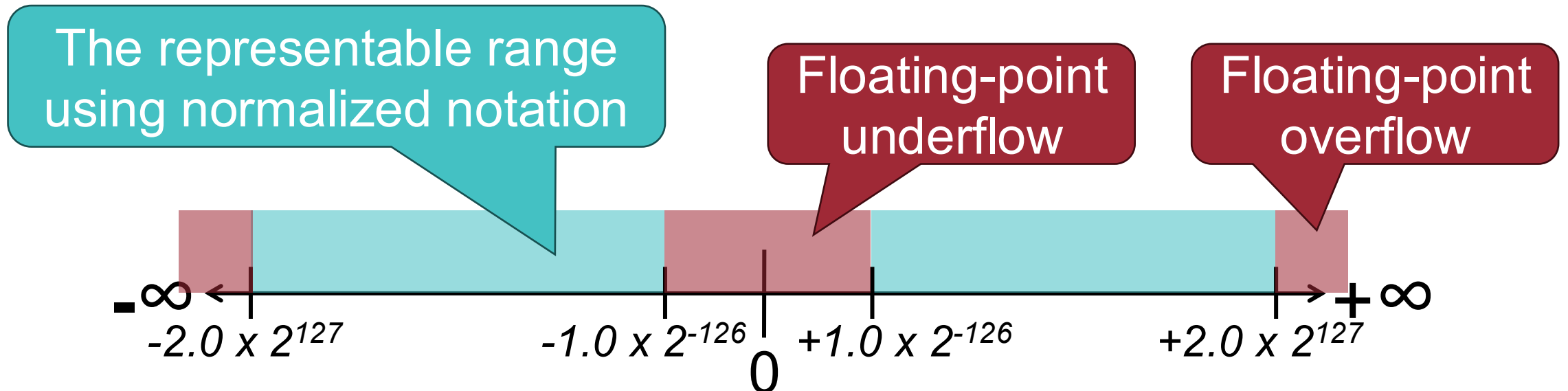
- **Overflow:** occurs when a result has a magnitude *too big* to be represented



Overflow and Underflow



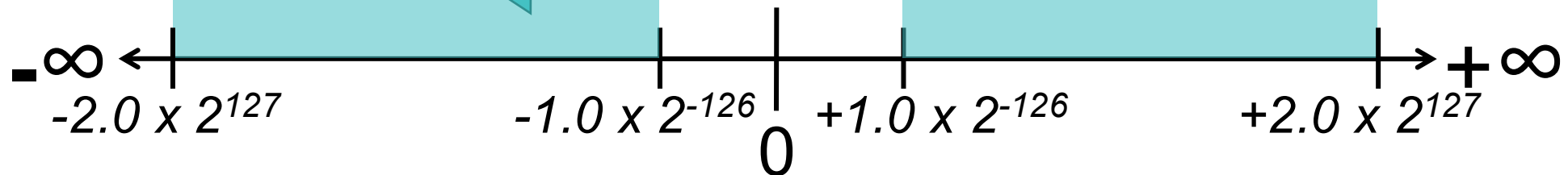
- **Overflow:** occurs when a result has a magnitude *too big* to be represented
- **Underflow:** occurs when a result has a magnitude *too small* to be represented



IEEE 754: Special Cases



The representable range using normalized notation

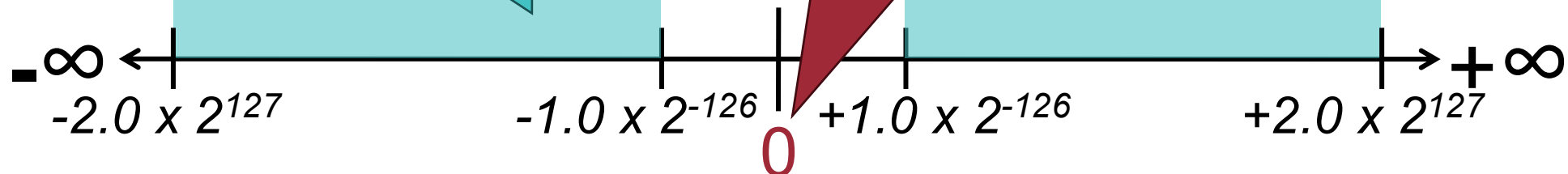


IEEE 754: Special Cases

- Exponent = 00...0, Fraction = 00...0
→ Not 1.0×2^{-127} but **0**

The representable range using normalized notation

How to represent 0?

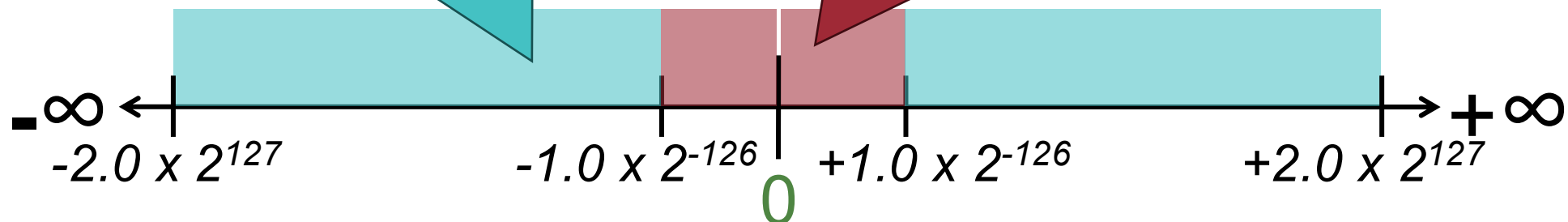


IEEE 754: Special Cases

- Exponent = 00...0, Fraction = 00...0
→ Not 1.0×2^{-127} but **0**
- Exponent = 00...0, Fraction \neq 00...0
→ Not $(1 + \text{fraction}) \times 2^{-127}$ but **$(0 + \text{fraction}) \times 2^{-126}$**
→ Denormalized real numbers (to represent very small numbers)

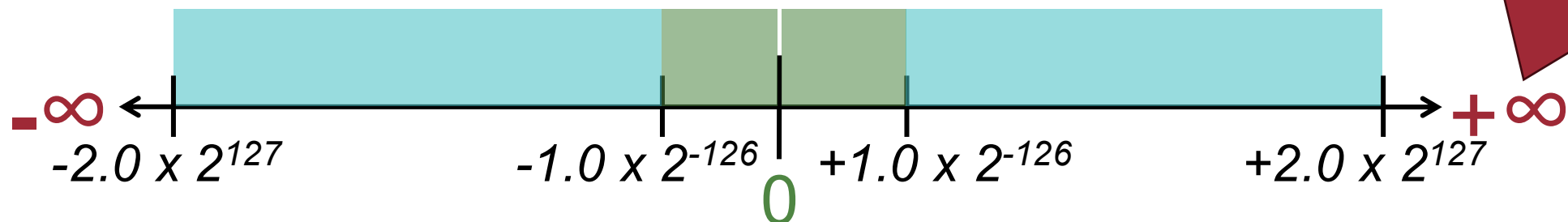
The representable range using normalized notation

How to represent very small numbers?
→ ***Denormalized numbers***



IEEE 754: Special Cases

- Exponent = 00...0, Fraction = 00...0
→ Not 1.0×2^{-127} but **0**
- Exponent = 00...0, Fraction \neq 00...0
→ Not $(1 + \text{fraction}) \times 2^{-127}$ but **$(0 + \text{fraction}) \times 2^{-126}$**
→ Denormalized real numbers (to represent very small numbers)
- Exponent = 11...1, Fraction = 00...0
→ **\pm infinity**



IEEE 754: Special Cases



- Exponent = 00...0, Fraction = 00...0
→ Not 1.0×2^{-127} but **0**
- Exponent = 00...0, Fraction \neq 00...0
→ Not $(1 + \text{fraction}) \times 2^{-127}$ but **$(0 + \text{fraction}) \times 2^{-126}$**
→ Denormalized real numbers (to represent very small numbers)
- Exponent = 11...1, Fraction = 00...0
→ **\pm infinity**
- Exponent = 11...1, Fraction \neq 00...0
→ **Not-a-Number (NaN)**
→ Indicates illegal or undefined result

How to represent the result of invalid operations (e.g., 0/0)?

IEEE 754 Encoding of Floating-point Numbers

11



Single precision		Double precision		Object represented
Exponent	Fraction	Exponent	Fraction	
0	0	0	0	0
0	Nonzero	0	Nonzero	\pm denormalized number
1–254	Anything	1–2046	Anything	\pm floating-point number
255	0	2047	0	\pm infinity
255	Nonzero	2047	Nonzero	NaN (Not a Number)

Question?