

Lecture 4: Unicode & bitwise operations

CSE 29: Systems Programming and Software Tools

Olivia Weng

Announcements

- Problem set I released
- TA office hours today during discussion

What about non-English characters?

What about non-English characters?

- Thousands more characters used in languages around the world
- ASCII does not define:
 - Spanish: é
 - Chinese: 中
 - Emoji: 🐫
- char datatype of 1 byte only encodes 256 possible bit patterns
- Challenge: Millions of lines of code written that assumed 1 byte ASCII chars

UTF-8: Unicode encoding

- Use more bits to encode more characters!
- Code point: an integer representing a character (e.g., `'A' == 65`)
- Normal ASCII code point: Highest order bit of byte is 0xxxxxx
 - UTF-8 is backwards compatible with ASCII! ←
- Multi-byte code point: Highest order bit of byte is 1xxxxxx

How many bytes do you need?

- Multi-byte code point: Highest order bit of byte is $1xxxxxx$
- Bit flags indicate code point length

- $110xxxxx = 2$ bytes
- $1110xxxx = 3$ bytes
- $11110xxx = 4$ bytes

$1110xxxx$ $10xxxxxx$ $10xxxxxx$
1st byte 2nd byte 3rd byte

- Bytes after the first byte start with

- $10xxxxxx$

2 bytes: $110xxxxx$
flag code point part
1st byte

$10xxxxxx$
flag code point part
2nd byte

Code point construction

$0xC3$
↓ ↓
1100 0011

$0xA9$
↓ ↓
1010 1001

- é = 1100011 101001
 - bit flags : only encodes metadata (i.e., only provides information about)
 - the code point = 000011 101001 = 233

$$\begin{array}{r} 128 + 64 + 32 + 8 + 1 \\ \downarrow \\ 11101001 = 233 \\ 2^7 2^6 2^5 2^4 2^3 2^2 2^1 2^0 \end{array}$$

$$\begin{array}{r} 192 \\ + 32 \\ \hline 224 \\ + 8 \\ \hline 232 \end{array}$$

Demo

- printing codepoints

How many bytes do you need?

1 byte char 0xxxxxx

- Multi-byte code point: Highest order bit of byte is 1xxxxxx
- Bit flags indicate code point length
 - 110xxxxx = 2 bytes
 - 1110xxxx = 3 bytes
 - 11110xxx = 4 bytes
- Bytes after the first byte start with
 - 10xxxxxx
- How do we check for these specific bit flags to know how many bytes we have?
 - Need a way of inspecting individual bits!

Bit operations

- Special mathematical operations in C for manipulating bits

- &: AND
- |: OR
- ~: NOT
- ^: XOR
- >>: shift right
- <<: shift left



AND: &

$$T \& T = T$$

0 = "false"
1 = "true"

- AND each bit together

Truth table

input_a	input_b	Result
0	0	0
0	1	0
1	0	0
1	1	1

$$\begin{array}{r} \downarrow \\ 0\ 0\ 0\ 0\ 1\ 0\ 1\ 0 \\ \& 0\ 0\ 0\ 0\ 1\ 1\ 1 \\ \hline 0\ 0\ 0\ 0\ 1\ 0\ 1\ 0 \end{array}$$

What is the result?

$$\begin{array}{r} 00110011 \\ \& 10100101 \\ \hline 00100001 \end{array}$$

OR: |

$$\begin{array}{l} T \mid F = T \\ F \mid T = T \\ T \mid T = 1 \end{array}$$

0 = "false"
1 = "true"

- OR each bit together

Truth table

input_a	input_b	Result
0	0	0
0	1	1
1	0	1
1	1	1

$$\begin{array}{r} 00010001 \\ | 11001010 \\ \hline 11011011 \end{array}$$

What is the result?

$$\begin{array}{r} 00100111 \\ | \quad 10101100 \\ \hline 10101111 \\ \hline \end{array}$$

A F

NOT: ~ *operates on 1 operand*

- NOT each bit together

Truth table

input_a	Result
0	1
1	0

$$\begin{array}{r} \sim 1010\ 1100 \\ \hline 0101\ 0011 \end{array}$$

Shift right: >> Ex: 4 bits two's complement

- Shift the bits "off a cliff" to the right

$$\begin{array}{r} 0100 \gg 2 = \underline{\quad 0\ 0\ 0\ 1} \cancel{\otimes} \\ \backslash 4 \cdot 1.2 \cdot 1.2 \end{array} = \begin{array}{l} 0001 \\ \underline{1} \end{array} \leftarrow \begin{array}{l} \text{sign extending} \\ \text{by 0} \end{array}$$

$$\begin{array}{r} 0100 \gg 1 = \underline{\quad 0\ 0\ 1\ 0} \cancel{\otimes} \\ \backslash 4 \cdot 1.2 \end{array} = \begin{array}{l} 0010 \\ \underline{2} \end{array}$$

$$\begin{array}{r} 1000 \gg 1 = \underline{\quad 1\ 1\ 0\ 0} \cancel{\otimes} \\ \backslash -8 \cdot 1.2 \end{array} = \begin{array}{l} [100 \\ -4] \end{array} \leftarrow \begin{array}{l} \text{sign extending} \\ \text{by 1} \end{array}$$

4 bits 2's complement

What is the result in binary and decimal?

- ~~1010~~ $\gg 3 = \underline{\underline{111}}$ -1

- 1000 $\gg 2 = \underline{\underline{1110}} -2$

- $1011 \gg 4 \approx \underline{\underline{1111}} -1$

$-5 // 2 = -2 // 2 = -1 // 2 = 0 \neq -1$

division pattern break

- $1101 \gg 2 = \underline{\underline{111}} -1$

Shift left: <<

4 bits unsigned

- Shift the bits "off a cliff" to the left

$$0001 \ll 1 = \cancel{0001}0 = 0010$$

| *2

 2

$$0001 \ll 3 = \cancel{0001}000 = 1000$$

| *2*2*2

 = 8

What is the result in binary and decimal (unsigned)?

- $0010 \ll 2 = 1000$ 8

- $0001 \ll 3 = 0000$ 8

- $1000 \ll 4 = 0000$ 0

- $0110 \ll 2 = 1000$ 8

Demo: practice with bitwise operators

```
int bitwise_is_even(int8_t num);
```



```
int count_1_bits(int32_t num);
```

How to use bit operators to select specific bits?

How to use bit operators to select specific bits?

- Bit masking: select specific individual bits out of a binary representation of a number
 - Examples:

- `lowest_four_bits(0b01010101) =` 0000 1111

- highest_four_bits(`0b10110011`) = 1111 0000

- `lowest_four_bits(192) =`
 $128 + 64 = 192$ 

- highest four bits(lowest four bits(200)) =

$$128 + 64 + 8 = 200 \rightarrow \begin{array}{r} 110001000 \\ 000000111 \\ \hline 1100001000 \end{array} \Rightarrow \underline{\underline{11110000}} \quad \underline{\underline{00000000}}$$

How to implement masking?

```
char lowest_four_bits(char c) {  
    return c & 0b00001111;  
}
```

```
char highest_four_bits(char c) {  
    return c & 0b11110000;  
}
```

The bitwise **&** operator selects **specific bit positions** based on the **pattern of 1's** that the variable is **&'d** with.

& truth table	Result
0 & 0	0
0 & 1	0
1 & 0	0
1 & 1	1

Why can't we use the bitwise | operator?

How do we check the codepoint bit flags?

- Bit masking!
- codepoint_size(char string[])

2 byte codepoint

