

```

1 char firstfour(char c) {
2     return c & 0b11110000;
3 }

```

```

1 char lastfour(char c) {
2     return c & 0b00001111;
3 }

```

```
1 char startswith1110(char c) {
```

```
2
```

```
3
```

```
4
```

```
5
```

```
6
```

A ✓ return ~~c~~ firstfour(c) == Ob 11100000

B ✗ return c > Ob11100000

C ✗ return (c & Ob11100000) == Ob11100000
 ↑
 change to 1
 to correct

very
wrong
(monkey)

```

1 #include <stdio.h>
2 #include <string.h>
3
4 void inspect(char s[]) {
5     int index = 0;
6     printf("%s, length %d: ", s, strlen(s));
7     while(s[index] != 0) {
8         char current = s[index];
9         printf("%c (%hu %b", current, current);
10        for(int place = 128; place > 0; place /= 2) {
11            if((current & place) == 0) printf("0");
12            else printf("1");
13        }
14        printf(" ");
15        index += 1;
16    }
17    printf("\n");
18 }

```

```

20 int main() {
21     char input[100];
22     fgets(input, 100, stdin);
23     printf("\n");
24     inspect(input);
25 }

```

firstfour(Ob01010101) =	<u>Ob01010000</u>
lastfour(Ob11110011) =	<u>Ob00000011</u>
firstfour(192) =	<u>Ob11000000</u>
firstfour(lastfour(x)) =	<u>0</u>

can you find

a counterexample

to where

c doesn't start
1110, but returns true

Ob11110000

A Ob11110000

= =

Ob11100000

(false, correct)

B Ob11110000 >

Ob11100000

true (wrong)

C Ob11110000

& Ob11100000

Ob11100000 == Ob11100000

true (wrong)

printing the
binary rep of
a character

single & (bitwise &)

```
$ gcc inspect.c -o inspect
$ ./inspect
José, length 5: J (74 0b01001010) o (111 0b01101111) s (115 0b01110011) é (195 0b11000011) ô (169 0b10101001)
$ ./inspect
ピカチュウ
ピカチュウ length 15: ô (227 0b11100011) (131 0b10000011) (148 0b10010100) // ... continues
$ ./inspect
🦀, length 4: ô (240 0b11110000) ♦ (159 0b10011111) ♦ (166 0b10100110) ♦ (128 0b10000000)
```

ASCII does not define é or ô or 🦀. There are 100s of thousands of chars

There are millions of lines of code that use ASCII

1993 USENIX Rob Pike San Diego

UTF-8

Bytes like 0xxxxxxx

are ASCII (code points)

Bytes like 1xxxxxxx

are part of multi-byte characters (code points)

é is 2 bytes

0b110xxxxxx

0b10xxxxxx

code points up to z"

ô is 3 bytes

0b1110xxxx

0b10xxxxxx

0b10xxxxxx

cps up to 2^{16}

🦀 is 4 bytes

0b11110xxx

0b10xxxxxx (3 more)

cps up to 2^{21}

"Unicode" is a map from code points to glyph, character