

```

1 #include <stdio.h>
2 #include <string.h>
3
4 void inspect(char s[]) {
5     int index = 0;
6     printf("%s, length %ld:\n", s, strlen(s));
7     while(s[index] != 0) {
8         char current = s[index];
9         printf(" %c (%03hhu 0b%08hbb)\n", current, current, current);
10        index += 1;
11    }
12    printf("\n");
13 }
14
15 int main() {
16     char s1[] = "José";
17     char s2[] = "ピカチュウ";
18     char s3[] = "é";
19     inspect(s1);
20     inspect(s2);
21     inspect(s3);
22 }

```

```
$ gcc inspect.c -o inspect
```

```
$ ./inspect
```

```
José, length 5:
```

```
J (074 0b01001010)
```

```
o (111 0b01101111)
```

```
s (115 0b01110011)
```

```
é (195 0b11000011)
```

```
 (169 0b10101001)
```

```
$ ./inspect
```

```
ピカチュウ length 15:
```

```
 (227 0b11100011)
```

```
 (131 0b10000011)
```

```
 (148 0b10010100)
```

```
// ... continues
```

```
$ ./inspect
```

```
é, length 4:
```

```
 (240 0b11110000)
```

```
 (159 0b10011111)
```

```
 (166 0b10100110)
```

```
 (128 0b10000000)
```

110 prefix for 2-byte character

10 prefix on "continuation" bytes

1110 prefix for 3-byte character

10 prefix on "continuation" bytes

11110 prefix for 4-byte character

10 prefix on "continuation" bytes

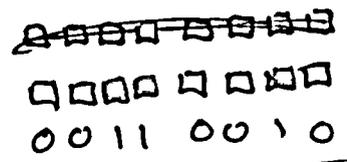
é

11000011 10101001

00011101001

233 "unicode code point 233"

```
char s1[] = "cse29";
printf("%c %d", s[3], s[3]);
```



**A1**: What does this print? 2 50

```
char s2[] = "cse30_next";
printf("%d", s2[10]); // fill in to print 0
printf("%c", s2[4]); // fill in to print 0
```

**(A2)**  
**(A3)**

	unsigned	signed
0000 0000	0	0
0000 0001	1	1
0000 0010	2	2
0000 0011	3	3
...	...	...
0111 1111	127	127
1000 0000	128	-128
1000 0001	129	-127
1000 0010	130	-126
...	...	...
1111 1111	255	-1

ASCII

$$\begin{array}{r}
 '0'iii' \quad 'iiii' \\
 + 0000 \quad 0001 \\
 \hline
 10000000
 \end{array}$$

**2's complement**

- most significant digit counts as negative
- addition "just works"
- no wasted/ambiguous reps
- can do "by hand"

# UTF8

~~Continuation~~

0xxxxxxx 1-byte UTF8 = ASCII

~~10~~

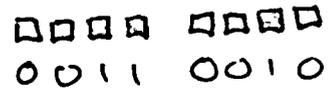
110xxxxx 10yyyyyy 2-byte UTF8 "code point"

xxxxxyyyyyy ← this is a number in the  
Unicode standard for  
a "glyph" or character  
we see

1110xxxx 10yyyyyy 10zzzzzz - 3-byte

11110xxx 10yyyyyy 10zzzzzz 10wwwww - 4-byte

```
char s[] = "cse29";  
printf("%c %d", s[3], s[3]);
```



A1: What does this print?      2    50

```
char s2[] = "cse30 next";  
printf("%d", s2[ 10 ]);  
printf("%c", s2[ 4 ]);
```

// fill in to print 0  
// fill in to print ~~0~~ 0

A2  
A3

	unsigned	signed
0000 0000	0	0
0000 0001	1	1
0000 0010	2	2
0000 0011	3	3
0000 0100	4	4
⋮		
0111 1111	127	127
1000 0000	128	-128
1000 0001	129	-127
1000 0010	130	-126
⋮		
1111 1111	255	-1

$$\begin{array}{r}
 0111\ 1111 \\
 + \quad \quad \quad 1 \\
 \hline
 10000000
 \end{array}$$

2's complement

Idea: most significant bit is interpreted as negative

- addition "just works"
- unique/unambiguous 0
- adding leading 1's to a negative # does not change its meaning

# UTF-8

1-byte	0xxxxxxx	ASCII		
2-byte	110xxxxx	10yyyyyy		
3-byte	1110xxxx	10yyyyyy	10zzzzzz	
4-byte	11110xxx	10yyyyyy	10zzzzzz	10wwwww

xxxxxyyyyyy → this is binary # of a code point in Unicode