

# COS 217: Introduction to Programming Systems

## Interfaces and Implementations\* Using the Smallest C Data Type: Character

\*Note: We will have more lectures on this topic in a more general context later



**PRINCETON UNIVERSITY**

# Agenda



Interface, implementation and design decisions for characters

Interfaces and implementations for the human reader

Another difference from Java

- Variable declarations in C89

# To Use Characters in Programs, What Do We Need?



- A representation for characters
- Ways to input and output characters
- Ways to manipulate characters
  - Convert from lowercase to uppercase, etc.



# Character Representation: The ASCII Standard

Mapping from integer values to characters on pretty much all machines:  
**ASCII** (American Standard Code for Information Interchange) (/ 'æski/)

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	NUL									HT	LF					
16																
32	SP	!	"	#	\$	%	&	'	(	)	*	+	,	-	.	/
48	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
64	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
80	P	Q	R	S	T	U	V	W	X	Y	Z	[	\	]	^	_
96	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
112	p	q	r	s	t	u	v	w	x	y	z	{		}	~	

Notes: Many non-printing characters left blank in table above

UPPER-CASE and lower-case letters are 32 apart

... but they're internally contiguous. So are digits 0 through 9.



## Converting to Uppercase in upper Program

```
if ((c >= 97) && (c <= 122))  
    c -= 32;
```

What's wrong?

# A Different Representation/Implementation: EBCDIC



Extended Binary Coded Decimal Interchange Code (/ ' εbsɪdɪk/)

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	NUL					HT										
16																
32						LF										
48																
64	SP											.	<	(	+	
80	&										!	\$	*	)	;	
96	-	/										,	%	_	>	?
112										`	:	#	@	'	=	"
128		a	b	c	d	e	f	g	h	i			{			
144		j	k	l	m	n	o	p	q	r			}			
160		~	s	t	u	v	w	x	y	z						
176																
192		A	B	C	D	E	F	G	H	I						
208		J	K	L	M	N	O	P	Q	R						
224	\		S	T	U	V	W	X	Y	Z						
240	0	1	2	3	4	5	6	7	8	9						

Partial map



# C Provides Character Literals

Translate to different integer values in different encodings

Single quote syntax: 'a' is a value of type char with the value 97 in ASCII

Use backslash to write special characters

- Examples (with numeric equivalents in ASCII, EBCDIC):

'a'	the a character (97, 129)
'A'	the A character (65, 193)
'0'	the zero character (48, 240)
'\0'	the NUL (nullbyte) character (0, 0)
'\n'	the newline character (10, 37)
'\t'	the horizontal tab character (9, 5)
'\\'	the backslash character (92, 224)
'\''	the single quote character (39, 125)
'\"'	the double quote character (34, 127)

## Converting to Uppercase: Version 2



```
if ((c >= 'a') && (c <= 'z'))  
    c += 'A' - 'a';
```

Arithmetic  
on chars?

What's wrong now?





# Recall EBCDIC

## Extended Binary Coded Decimal Interchange Code

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	NUL					HT										
16																
32						LF										
48																
64	SP											.	<	(	+	
80	&										!	\$	*	)	;	
96	-	/										,	%	_	>	?
112										`	:	#	@	'	=	"
128		a	b	c	d	e	f	g	h	i			{			
144		j	k	l	m	n	o	p	q	r			}			
160		~	s	t	u	v	w	x	y	z						
176																
192		A	B	C	D	E	F	G	H	I						
208		J	K	L	M	N	O	P	Q	R						
224	\		S	T	U	V	W	X	Y	Z						
240	0	1	2	3	4	5	6	7	8	9						

Partial map

Note: UPPER CASE not contiguous; same for lower case.



## Converting to Uppercase: Version 3

- Provide a real **interface**:
  - Character data type
  - API for operations on characters
    - Works on all machines and representations
    - **Implemented** differently for different machines and representations (ASCII, EBCDIC, etc)

```
#include <ctype.h>

if (islower(c))
    c = toupper(c);
```



# What about Input/Output?

C does not provide I/O facilities in the language :

- They're provided in the standard library, declared in `stdio.h`
  - Constant: `EOF`
  - Data type: `FILE` (described later in course)
  - Variables: `stdin`, `stdout`, and `stderr`
  - Functions: (numerous)

## Reading characters

- `getchar()` function with return type wider than `char` (specifically, `int`)
- Returns `EOF` (a special non-character `int`) to indicate failure
- **Reminder: there is no such thing as "the EOF character"**

## Writing characters

- `putchar()` function accepting one parameter
- For symmetry with `getchar()`, parameter is an `int`

## Aside: Unicode

Back in 1970s, the language only cared about English characters:

ASCII:

# American Standard Code for Information Interchange

## Other languages?

## Diacritics?

[illegible]

결	경	결	국	국
ACB1	ACD1	ACE1	ACE2	ACF2
결	계	겸	곤	국
ACB4	ACD4	ACE4	ACE4	ACF4
결	객	겹	곡	공
ACB5	ACD5	ACE5	ACE5	ACF5
결	객	겹	공	국
ACB9	ACD9	ACE9	ACE9	ACF9
결	객	겹	곧	국
ACB7	ACD7	ACE7	ACE7	ACF7
결	겐	겹	곧	곡
ACB8	ACD8	ACE8	ACE8	ACF8



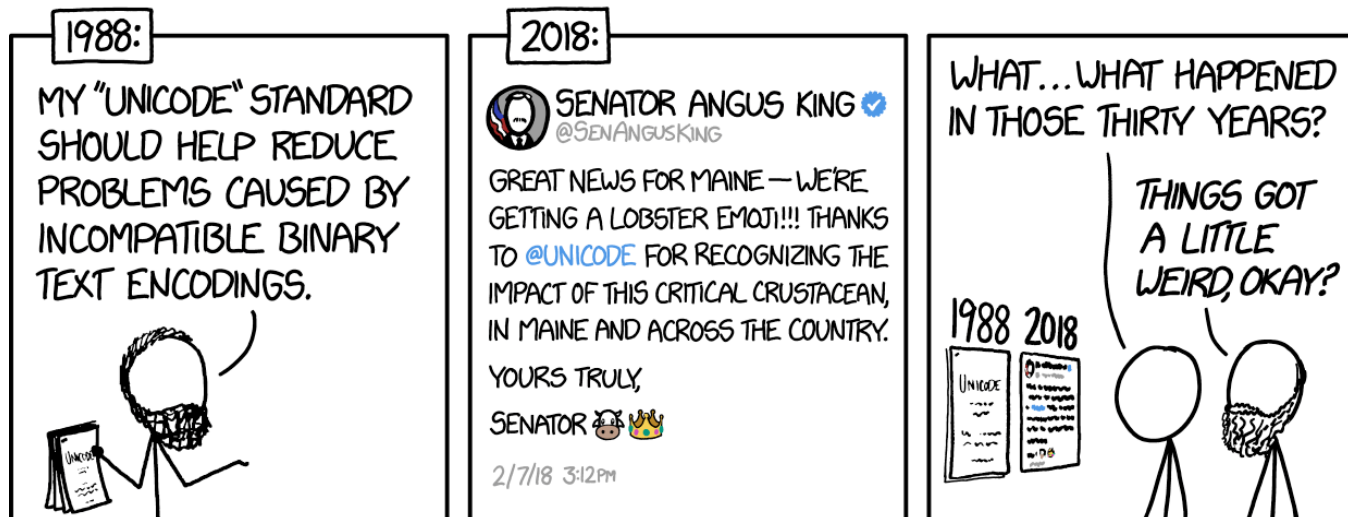
·	2	3	4
1001	1012	1023	1034
·	2	3	4
1001	1013	1023	1033
:	0	1	2
1004	1014	1024	1034
·	1	2	3
1005	1015	1025	1035
-	6	7	8
1006	1016	1026	1036
+	0	1	2
1007	1017	1027	1037
0	1	2	3
1010	1020	1030	1040
1	2	3	4
1011	1021	1031	1041
0	1	2	3
1012	1022	1032	1042
0	1	2	3
1013	1023	1033	1043
0	1	2	3
1014	1024	1034	1044
0	1	2	3
1015	1025	1035	1045
0	1	2	3
1016	1026	1036	1046
0	1	2	3
1017	1027	1037	1047
0	1	2	3
1018	1028	1038	1048
0	1	2	3
1019	1029	1039	1049



# Modern Unicode

When C was designed, characters fit in 8 (really 7) bits, so C's chars are 8 bits long.

When Java was designed, Unicode fit in 16 bits, so Java's chars are 16 bits. Then ...



<https://xkcd.com/1953/>

Result: modern systems use *variable length* (UTF-8/16/32) encoding for Unicode.

In C, this is supported not in the language but via libraries

# Agenda



Interfaces, implementations and design decisions for characters

Interfaces and implementations for the human reader

Another difference from Java

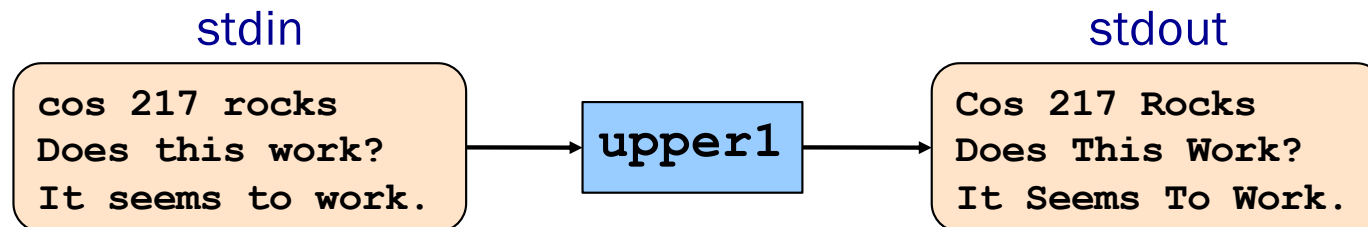
- Variable declarations in C89



# Recall: The upper1 Product Specification

## Functionality

- Read all chars from stdin
- Capitalize the first letter of each word
  - “cos 217 rocks”  $\Rightarrow$  “Cos 217 Rocks”
- Write result to stdout



# upper1 Version 3



```
#include <stdio.h>
#include <ctype.h>
enum Statetype {NORMAL, INWORD};

enum Statetype handleNormalState(int c)
{
    enum Statetype state;
    if (isalpha(c)) {
        putchar(toupper(c));
        state = INWORD;
    } else {
        putchar(c);
        state = NORMAL;
    }
    return state;
}

enum Statetype handleInwordState(int c)
{
    enum Statetype state;
    if (!isalpha(c)) {
        putchar(c);
        state = NORMAL;
    } else {
        putchar(c);
        state = INWORD;
    }
    return state;
}
```

```
int main(void)
{
    int c;
    enum Statetype state = NORMAL;
    while ((c = getchar()) != EOF) {
        switch (state) {
            case NORMAL:
                state = handleNormalState(c);
                break;
            case INWORD:
                state = handleInwordState(c);
                break;
        }
    }
    return 0;
}
```

That's an A-, at best.  
No comments!





# Upper1: Improving the Interface

## Problem:

- The program works, but...
- It's too hard for the human reader

## Solution:

- Add function-level comments as part of the interface



# Function Comments: A Key Part of Interfaces

Function comment should describe

***what the function does*** (from the caller's viewpoint)

- Data coming into the function
  - Parameters, input streams
- What it does to those data (at a high level)
- Data going out from the function
  - Return value, output streams, call-by-reference parameters

Function comment should **not** describe

***how the function works***



# Function Comment Examples

**Bad** main() function comment: Describes how the function works

```
Read a character from stdin using getchar.  
Depending upon the current DFA state, pass the  
character to an appropriate state-handling  
function. The value returned by the state-  
handling function is the next DFA state. Repeat  
until end-of-file. Return 0.
```

**Good** main() comment: Describes what the function does (from caller's perspective)

```
Read text from stdin. Convert the first character  
of each "word" to uppercase, where a word is a  
sequence of uppercase or lowercase letters. Write  
the result to stdout. Return 0.
```



## Upper1: Other Interface Comments

```
/* defines constants representing each state in the DFA */  
enum Statetype {NORMAL, INWORD};
```

```
/* Implement the NORMAL state of the DFA. c is the current DFA  
character. Write c's uppercase equivalent, if it has one, or  
otherwise c itself, to stdout. Return the next state specified  
by the DFA. */
```

```
enum Statetype handleNormalState(int c) {
```

```
/* Implement the INWORD state of the DFA. c is the current  
DFA character. Write c to stdout. Return the next DFA state. */
```

```
enum Statetype handleInwordState(int c) {
```

```
/* Read text from stdin. Convert the first character of each  
"word" to uppercase, where a word is a sequence of  
letters. Write the result to stdout. Return 0. */
```

```
int main(void) {
```

```
/* Use a DFA approach. state is the current DFA state. */  
enum Statetype state = NORMAL;
```



## Other Good Things for the Human Reader

- Comments in the implementation
- Readable code
  - Often a tradeoff with efficiency and “showing coding prowess”
  - Lack of language support for other data types can make it worse



## iClicker Question



Q: Is the `if` statement in `upper` really necessary?

A. Gee, I don't know.  
Let me check  
the man page  
(again)

```
#include <stdio.h>
#include <ctype.h>
int main(void)
{
    int c;
    while ((c = getchar()) != EOF) {
        if (islower(c))
            c = toupper(c);
        putchar(c);
    }
    return 0;
}
```



# ctype.h Functions

```
$ man toupper
```

## NAME

`toupper`, `tolower` - convert letter to upper or lower case

## SYNOPSIS

```
#include <ctype.h>
int toupper(int c);
int tolower(int c);
```

## DESCRIPTION

`toupper()` converts the letter `c` to upper case, if possible.  
`tolower()` converts the letter `c` to lower case, if possible.

If `c` is not an unsigned char value, or EOF, the behavior of these functions is undefined.

## RETURN VALUE

The value returned is that of the converted letter,  
or `c` if the conversion was not possible.



## iClicker Question



Q: Is the if statement really necessary?

- A. Yes, necessary for correctness.
- B. Not necessary, but I'd leave it in.
- C. Not necessary, and I'd get rid of it.

```
#include <stdio.h>
#include <ctype.h>
int main(void)
{
    int c;
    while ((c = getchar()) != EOF) {
        if (islower(c))
            c = toupper(c);
        putchar(c);
    }
    return 0;
}
```





## iClicker Question: How to Structure a Loop?



Q: There are several ways to structure a loop – which is best?

A. `for (c = getchar(); c != EOF; c = getchar())  
 putchar(toupper(c));`

B. `while ((c = getchar()) != EOF)  
 putchar(toupper(c));`

C. `for (;;)   
{ c = getchar();  
 if (c == EOF)  
 break;  
 putchar(toupper(c));  
}`

D. `c = getchar();  
while (c != EOF)  
{putchar(toupper(c));  
 c = getchar();  
}`



## Other Good Things for the Human Reader

- Comments in the implementation
- **Readable code**
  - Often a tradeoff with efficiency and “showing coding prowess”
  - Lack of language support for other data types can make it worse



# Unlike Java, C has No Boolean (Logical) Data Type

- Represent logical data using type char or int
  - Or any primitive type! 🤖
- Conventions:
  - Statements (if, while, etc.) use  $0 \Rightarrow \text{FALSE}$ ,  $\neq 0 \Rightarrow \text{TRUE}$
  - Relational operators ( $<$ ,  $>$ , etc.) and logical operators ( $!$ ,  $\&\&$ ,  $||$ ) produce the result 0 or 1
- Would have been nice to have a Boolean type



# Issues with Lack of Logical Data Type

Imagine this type of code in Java code and in C.  
What happens in each case?

```
...  
int i;  
...  
i = 0;  
...  
if (i = 5)  
    statement1;  
...
```

What happens  
in C?

What happens  
in Java?

DALL-E 2

prompt: impressionist painting of a  
computer programmer with a lack  
of sleep debugging late at night

# Sample Exam Question (Spring 2016, Exam 1)



Indicate what value this expression evaluates to:

What happens  
in C?



```
...  
-10 < i < -1  
...
```



What happens  
in Java?



DALL-E 2

prompt: impressionist painting of a  
computer programmer with a lack  
of sleep debugging late at night



# Doing Logic with Integers

Using integers to represent logical data permits shortcuts

```
...  
int i;  
...  
if (i) /* same as (i != 0) */  
    statement1;  
else  
    statement2;  
...
```

It also permits really bad code...

```
i = (1 != 2) + (3 > 4);
```



## iClicker? More like iBrainteaser!



Q: What is `int i` set to in the following code?

```
i = (i < (i < 0)) + (i >= (i > 0)) + ((i-i) < (i == i));
```

A. Depends on the initial value of i

B. 0

C. 1

D. 2

E. 3

D.

If i is negative, this will be  $1 + 0 + 1$

If i is non-negative, this will be  $0 + 1 + 1$



# Agenda

Interfaces, implementations and design decisions for characters

Interfaces and implementations for the human reader

Another difference from Java

- Variable declarations in C89





# Declaring Variables

C requires variable declarations (some languages don't: awk, bash)

Declaring variables requires more from the programmer

- Extra verbiage
- Thinking ahead about how it will be used

But it has many benefits

- Allows compiler to check “spelling”
- Allows compiler to allocate memory more efficiently
- Declaring variables' types produces fewer surprises at runtime



# Declaring Variables

C requires variable declarations.

- Declaration statement specifies type of variable (and other attributes too)

Examples:

```
int i;  
int i, j;  
int i = 5;  
const int i = 5; /* value of i cannot change */  
static int i; /* covered later in course */  
extern int i; /* covered later in course */
```



# Declaring Variables

C requires variable declarations.

- Declaration statement specifies type of variable (and other attributes too)
- Unlike Java (and later versions of C), declaration statements in C89 must appear **before** any other kind of statement in any compound statement.
  - Note this **doesn't** mean that all declarations must be at the top of a *function* exclusively!

```
{  
    int i;  
  
    /* Non-declaration  
       statements, even if  
       they don't use j */  
  
    int j;  
}
```

Illegal in C89

```
{  
    int i;  
    int j;  
  
    /* Non-declaration  
       statements that use  
       i and/or j. */  
}
```

Legal in C89

Next time ... numbers! (Bigger than 127.)

