

An Introduction to C Programming for Java Programmers

Mark Handley

M.Handley@cs.ucl.ac.uk

CONTENTS

1	Introduction	2
2	Basics of C	3
2.1	C Source Files	3
2.2	main()	4
2.3	Built-in Commands	5
2.4	Variables and Basic Types	8
2.5	Declaring Variables	10
2.6	Operators	11
2.6.1	Arithmetic Operators	11
2.6.2	Bitwise Operators	11
2.6.3	Relational and Logical Operators	11
2.6.4	Assignment	11
2.6.5	Increment and Decrement	12
2.7	Overflows, Assignments, and Other Trouble	13
2.8	Arrays	15
2.9	Pointers and Addresses	16
2.10	Dynamic Memory Allocation	17
2.11	Text Strings	19
2.12	Data Structures	22
2.13	Functions and Procedures	25
2.14	Input, Output and File Handling	30
3	The C Preprocessor	33
3.1	Including Header Files	33
3.2	Preprocessor Constants and Macros	33
3.3	Conditional Compilation	34

1 INTRODUCTION

This short tutorial is an introduction to programming in C. It is intended for students who already have some programming experience in Java, so know the basics of how to program and are familiar with the basic Java syntax which is shared with C. This is not a reference manual, nor a complete guide to the C language.

To learn more about C than is covered here, I recommend “The C Programming Language” by Kernighan and Ritchie (<http://cm.bell-labs.com/cm/cs/cbook/>).

In this tutorial I discuss many common mistakes made by C programmers, with particular emphasis on the sort of issues of which a Java programmer might be unaware. The intent is not to scare the reader away from C, but to highlight specific differences with Java, and avoid making the most common errors.

C is a fairly low-level language. Features and limitations of the underlying hardware are often exposed in C, whereas in Java they are largely hidden. The result is that good C programs will usually run quickly, with a small memory footprint. For operating system programming, C allows direct access to the hardware, which might simply not be possible in Java.

The downside of C compared to Java is that C does not protect the programmer from his or her own errors. It is not a strongly typed language, so the compiler provides little warning of many errors that would simply fail to compile in Java. This makes debugging harder. Memory management in C is manual, and manipulating pointers to memory is essential in any non-trivial C program. Getting this right is difficult.

But comparing C against Java on a good/bad scale is not productive; each serves a useful role for a certain set of tasks. For example, these days I would not recommend writing simple non-performance-critical network servers in C (although many are), simply because it is too hard to write really secure code. But writing a general-purpose operating system entirely in Java might not even be possible, and would certainly be rather inefficient on today’s common processors. A competent programmer should be aware of the advantages and limitations of a range of languages, and use the right tool for the job.

2 BASICS OF C

2.1 C Source Files

Sometimes a C program can be included in a single source file, but this is infeasible with non-trivial sized programs. There are three types of file that you need to be concerned with:

- C program files. The convention is that these have a `.c` extension such as `foo.c`. C program files contain the actual executable code of your program.
- C header files. The convention is that these have a `.h` extension such as `foo.h`. C header files contain interface definitions, preprocessor macros, and so on. When you need to call a function that is defined in one C program file from a function that is defined in another C program file, you will need to create a header file with a *prototype* for the external function, so that the compiler knows how to compile your function.
- library files. There are several different types of library file, depending on how they are linked to your program, but they all contain pre-compiled functions that you can use in your program.

In addition, it is common to automate the process of figuring out which source files need to be compiled together, using the utility `make`. Thus you'll often see a file called `Makefile` in a directory containing C source files. The `Makefile` contains information about which programs should be built, and how to compile and link them.

We'll discuss separate compilation and the use of libraries later.

2.2 main()

Every C program must contain a procedure called `main()`. This is the first procedure in your program to be executed. The simple *hello world* program in C then looks like:

```
#include <stdio.h>

main()
{
    printf("hello world!\n");
}
```

To compile this using `gcc` on a Unix or MacOS X system, if the source file is called `helloworld.c`, you'd type (in bold):

```
vulture.xorp.org: gcc -o helloworld helloworld.c
vulture.xorp.org:
```

The program `gcc` is the GNU C compiler, which is the most common C compiler found on most Unix and Linux systems. On some systems the C compiler might be called `cc`.

The `"-o helloworld"` argument indicates that the compiled program should be called `"helloworld"`. If no filename is given for the compiled program, the default on Unix is usually `a.out`. The last argument `"helloworld.c"` tells the compiler which source file or files to compile.

On Unix, see the `gcc` manual page for details of other command line options to `gcc` (run `man gcc`). A full `gcc` book can be found online at <http://www.network-theory.co.uk/gcc/gccintro/>.

Then to run your program, you'd type:

```
vulture.xorp.org: ./helloworld
hello world!
vulture.xorp.org:
```

In addition to showing how to compile code, this trivial example illustrates several things:

- The use of the preprocessor macro `#include` to include the `stdio` header file. Without this, the compiler wouldn't know how to compile the `printf` command which is in the standard C library.
- The use of `main()` as the first procedure to be executed in any C program.
- The use of `printf` to print to the standard output channel (in this case your terminal window). On Unix, see the `printf` manual page for more details (run `man printf`)
- The use of `\n` to indicate a newline in a string.

We'll see examples of separate compilation and linking of multiple program files in Sections 2.13 and 3.2.

2.3 Built-in Commands

C supports the usual commands to control execution flow that you are used to from Java:

if statements

```
if ( pi == 3 ) {  
    location = "indiana";  
}
```

The braces ({ and }) are only needed if more than one statement is conditional. In the example below, the `printf` statement will be executed irrespective of the value of `pi`.

```
location = "unknown";  
if (pi == 3)  
    location = "indiana";  
printf("Location = %s\n", location);
```

if/else statements

```
if (pi == 3) {  
    location = "indiana";  
} else {  
    location = "unknown";  
}
```

C also supports a shorthand if/else notation. The single statement below has exactly the same effect as the if/else example above.

```
location = (pi == 3 ? "indiana" : "unknown");
```

This notation is very useful in some circumstances such as preprocessor macro definitions, but it should be used sparingly as it is not always as readable as a regular if/else construction.

while statements

```
while (money > 0) {  
    money--;  
    drink_another_beer();  
}
```

do/while statements

```
do {  
    sleep++;  
} while (sleep < 8)
```

A `do/while` loop is similar to a `while` loop, except that the contents of the loop are always executed at least once, even if the condition is false. The regular `while` loop turns out to be more commonly used in practice.

for loops

```
int i;  
for (i = 1; i <= 10; i++) {  
    printf("%d squared is %d\n", i, i*i);  
}
```

The `for` command consists of three control statements and a body that may get executed multiple times. The three control statements in the above example are:

- `i = 1`. The first control statement is executed only once before the loop is entered for the first time.
- `i <= 10`. The second control statement is the test for the loop. It is executed before the loop is entered each time, and if the value is true, the body of the loop is executed. If the test is false on the first try, the body of the loop will never be executed.
- `i++`. The third control statement is executed after each pass through the loop body, before the test is run. The notation `i++` increments the value of `i` by one.

Thus in the above example, the order of execution is:

1. `i = 1`
2. `i <= 10`. Value is true, so loop body is executed.
3. `printf()`.
4. `i++`. Value of `i` is now 2.
5. `i <= 10`. Value is true, so loop body is executed.
6. `printf()`.
7. `i++`. Value of `i` is now 3.
- ...
...
8. `printf()`.
9. `i++`. Value of `i` is now 11.
10. `i <= 10`. Value is false, so execution continues with the next statement after the loop.

Note the use of `%d` to indicate that an integer argument is to be printed as a decimal value.

switch statements

```
char c;
c = getchar();
switch (c) {
    case 'c':
        cont = 1;
        break;
    case 'd':
        done = 1;
        break;
    case 'Q':
    case 'q':
        exit();
    default:
        printf("valid inputs and 'q', 'c' and 'd'\n");
}
```

The `case` keyword is used to indicate each of the possible values that the variable in the `switch` statement might take.

The `default` keyword is used to indicate a catchall case that matches everything not matched by a prior case statement.

The `break` keyword is used to leave the `switch` statement. If this is omitted, execution drops through to the body of the `case` statement below. This is shown in the `case 'Q'` line, where the inputs `'q'` and `'Q'` both cause the program to exit.

goto statements

```
char c;
int characters = 0;
int lines = 0;
while (1) {
    c = getchar();
    if ( c == EOF ) {
        goto Finished;
    }
    characters++;
    if ( c == '\n' ) {
        lines++;
    }
}

Finished:
printf( "lines: %d, characters: %d\n", lines, characters);
```

This simple and rather inefficient program counts the number of lines and characters in an input file. At the end of the file, it uses a `goto` statement to jump out of the `while` loop, continuing execution after the label `Finished`. In this particular case, a `break` statement could have been used instead, but if there are many

nested loops, `break` cannot trivially be used.

`goto` should be used sparingly. It's best used for handling exceptions and error cases. Unlike java, C has no built-in exception handling; `goto` can sometimes be used as an effective alternative. Too much use of `goto` can make your code unreadable.

2.4 Variables and Basic Types

The basic types available to you in C are:

- characters
- integer numbers
- floating point numbers

Various modifiers are available to change the size of these, and to specify whether they can hold signed or unsigned values. Unlike with Java, the precise amount of storage (and hence range of values available) for a variable depends on the particular system the program is running on.

On an Intel 32 bit x86 machine, the basic variable types are:

- `char` - 8 bit signed integer, also used to store characters for strings. Numeric value range: -128 to 127
- `unsigned char` - 8 bit unsigned integer. Value range: 0 to 255
- `short` - 16 bit signed integer. Value range: -32768 to 32767
- `unsigned short` - 16 bit signed integer. Value range: 0 to 65535
- `int` - 32 bit signed integer. Value range: -2,147,483,648 to +2,147,483,647
- `long` - same as `int`.
On some other systems, `int` might be 16 bit and `long` might be 32 bit.
- `unsigned int` - 32 bit unsigned integer. Value range: 0 to 4,294,967,295
- `unsigned long` - same as `unsigned int`.
On some other systems, `unsigned int` might be 16 bit and `unsigned long` might be 32 bit.
- `long long` - 64 bit signed integer. Value range: -2^{63} to $2^{63} - 1$.
- `float` - 32 bit floating point number.
- `double` - 64 bit floating point number.
- `long double` - 96 bit floating point number.

You might be surprised to learn that a `char` is an 8-bit signed integer. The reality is that if we look at the actual CPU's Arithmetic and Logic Unit (ALU), there are really only two basic types: floating point numbers and integers. C is not strongly typed, so assigning an integer to a `char` is perfectly OK, so long as the value fits in the appropriate data range for a `char`. Thus the following is perfectly legal C:

```
char c;
int i;
i = 67;
c = i;
printf("The ASCII character code for '%c' is %d\n", c, c);
```

And the output from this code fragment is:

```
The ASCII character code for 'C' is 67
```

Note that the `printf` statement prints the value of `c` twice, once as a character (`%c`) and once as an integer (`%d`).

The `sizeof()` command can be used to find out how many bytes a particular type or variable uses on the local system. This can be useful when trying to write code that is portable between 16-bit, 32-bit and 64-bit systems. For example:

```
short s;
int i;
long l;
long long ll;
float f;
double d;
long double ld;
printf("short: %d, int: %d, long: %d long long: %d\n",
      sizeof(s), sizeof(i), sizeof(l), sizeof(ll));
printf("float: %d, double: %d long double: %d\n",
      sizeof(f), sizeof(d), sizeof(ld));
```

On a 32-bit x86 machine, this prints out:

```
short: 2, int: 4, long: 4 long long: 8
float: 4, double: 4 long double: 12
```

On a 64-bit Solaris/Ultraspac machine, this prints out:

```
short: 2, int: 4, long: 4 long long: 8
float: 4, double: 4 long double: 16
```

The only difference here is in the size of a `long double`. However, on embedded processors or some other 64 bit machines, the some of the other values may vary.

You really get a sense of the low-level nature of C here. Java provides a nice OS-independent set of types. C gives you what the raw hardware provides, and as a result is usually significantly faster. The downside is that it's harder to write portable code.

In general, you want to be careful to use variables that have enough space for the range of values you want to hold, and you want to check untrusted input to be sure it's in the right range before you store it. **Overflows will not be automatically detected by the compiler or at runtime!**

For floating point numbers, you almost always want to use `double`. The type `float` usually doesn't have enough precision. Also be careful about exact comparisons with floating point numbers - hardware rounding errors can often result in very slight variations from the value you're comparing against.

2.5 Declaring Variables

Variables must be declared at the start of a block, and they remain in scope until the end of the block. A block can be a procedure, if statement, while loop, or practically anything enclosed by { and }. For example:

```
#include <stdio.h>
main()
{
    int i = 1, j = 1;
    while (i == 1) {
        int j, k;
        j = 2;
        printf("In here, i is %d, j is %d\n", i, j);
        i--;
    }
    printf("Out here, i is %d, j is %d\n", i, j);
}
```

In this example `i` and `j` are two integer variables declared at the start of `main`.

Two more variables (`j` and `k`) are then declared inside the `while` loop.

But `j` was already declared outside the while loop. What happened here is that a new variable is created, overriding (or “shadowing”) the first version of `j` for the duration of the while loop, and then disappearing after termination of the loop.

When you run this program you get:

```
vulture.xorp.org: gcc -o foo foo.c
vulture.xorp.org: ./foo
In here, i is 1, j is 2
Out here, i is 0, j is 1
```

Thus the assignment `j = 2` does not change the value of the original `j`, but instead sets the value of the new version of `j` in the while loop.

Shadowing of this form is usually a mistake - it’s too confusing to do intentionally on a regular basis. You can get the compiler to help you spot such errors:

```
vulture.xorp.org: gcc -Wshadow -o foo foo.c
foo.c: In function 'main':
foo.c:6: warning: declaration of 'j' shadows previous local
```

In general, getting the compiler to tell you about possible errors is a good thing. Often it will spot problems you might have missed. See the `gcc` man page for lots of other `-w` flags you can set.

2.6 Operators

C provides all the usual operators you would expect.

2.6.1 Arithmetic Operators

+	addition
-	subtraction
*	multiplication
/	division
%	modulus (remainder). For example the value $10 \% 3$ is 1.

2.6.2 Bitwise Operators

~	one's complement
	bitwise or
&	bitwise and
^	bitwise exclusive or (xor)
>> <i>n</i>	shift right <i>n</i> bits
<< <i>n</i>	shift left <i>n</i> bits

Bitwise operators perform a logical operation on all the bits in a value.

For example, the value of $(1 | 4)$ is 5, and the value of $(3 \& 2)$ is 2.

2.6.3 Relational and Logical Operators

==	equal to
<	less than
<=	less than or equal to
>	greater than
>=	greater than or equal to
!=	not equal to
!	logical not
	logical or
&&	logical and

In C, there is no *boolean* type - the *int* type is simply overloaded. Where a boolean type is expected, an int with a value of zero is treated as *false*, and an int with a non-zero value is treated as *true*.

For example, the value of $(0 || 4)$ is 1 (true), and the value of $(3 \&\& 2)$ is also 1 (true).

Confusing bitwise and logical operators is a common error.

2.6.4 Assignment

C assigns variables in the usual way:

`a = 2;` sets the value of `a` to be 2.

C also supports many shorthand operator/assignments to change the value of a variable:

```
a += 2    same as a = a + 2
a -= 2    same as a = a - 2
a *= 2    same as a = a * 2
a /= 2    same as a = a / 2
a %= 2    same as a = a % 2
```

```
a &= b    same as a = a & b
a |= b    same as a = a | b
a ^= b    same as a = a ^ b
```

```
a <<= 2   same as a = a << 2
a >>= 2   same as a = a >> 2
```

A common error is to use = instead of == in a comparison.

For example:

```
WRONG:
    if (a = 2) {
        /*do something*/
    }
```

In the above example, a will be *assigned* the value of 2. The assignment itself returns the value of 2, which is interpreted as *true*, and the conditional branch will always be executed. **The compiler will not notice this error, because it's valid C**, but it's rarely what you intended.

```
CORRECT:
    if (a == 2) {
        /*do something*/
    }
```

In the version above, the correct comparison operator has been used.

```
CORRECT AND SAFER:
    if (2 == a) {
        /*do something*/
    }
```

In this version, the correct comparison operator has been used, but also the terms in the comparison have been reversed. This is good practice, because if you accidentally type = instead of ==, the compiler will report an error because you can't assign a value to the constant 2.

2.6.5 Increment and Decrement

```
++x    increment x, then use new value
x++    use old value of x, then increment x
--x    decrement x, then use new value
x--    use old value of x, then decrement x
```

For example:

```
int x, y;

x = 1;
y = ++x;
/* value of x is 2, value of y is 2 */

x = 1;
y = x++;
/* value of x is 2, value of y is 1 */

x = 1;
y = x++ * 10;
/* value of x is 2, value of y is 10 */
```

2.7 Overflows, Assignments, and Other Trouble

C allows you to assign pretty much anything to pretty much any type. It expects you know what you're doing. This can allow you to write very efficient and fast code, but it will also get you into trouble unless you're very careful.

Some examples of unexpected consequences:

```
unsigned char c;
c = 255;
c += 2;
```

Adding 2 to the unsigned char `c`, which holds a value of 255 as shown above, results in the value of `c` wrapping. The end result is that `c` has a value of 1.

```
int i;
short s;
i = 32769;
s = i;
```

The least significant 16 bits of the 32-bit `int` `i` are copied into the 16-bit `short` `s`. Unfortunately as a 16-bit signed number, 32769 becomes -32767 which is then the new value of `s`.

```
int i;
float f;
i = 1234567890;
f = i;
```

The integer `i` is copied to the float `f`. The compiler does the right thing, and adds code to convert from integer representation to floating point representation, but this value is too large to be held in a 32-bit float without losing precision. `f` ends up holding the value 1234567936.000000. Using a `double` here would have

avoided this particular rounding error.

```
int i = 1;
int j = 0;
int k;
k = i/j;
```

This is an arithmetic error. In Java, this would throw an exception. In C, it causes the process to crash, and to dump a core file:

```
vulture.xorp.org: gcc -g -o foo foo.c
vulture.xorp.org: ./foo
Floating exception (core dumped)
vulture.xorp.org:
```

The CPU raised a divide-by-zero exception, which the operating system trapped, and aborted execution of your program. There's no indication as to what exactly you did wrong, or where in your program the error is. But be grateful - if the OS hadn't handled this, the system would have crashed - it did the best it could, and preserved the evidence. Your best hope now is to use a debugger such as *gdb* to examine that core file. Note the use of the `-g` flag above - without this the executable will not include debugging symbols that *gdb* needs.

2.8 Arrays

An array in C is declared and accessed as follows:

```
int a[10];
int i;
for (i = 0; i < 10; i++) {
    a[i] = 0;
}
```

The statement `int a[10];` declares an array of ten integers, and allocates memory on the stack for this array. As it's on the stack, the memory will be automatically freed when `a` goes out of scope.

The `for` loop here simply sets the values of the elements of `a` to be zero. You cannot expect C to zero memory for you. Note that the first element of `a` is `a[0]` and the last element is `a[9]`.

C does not perform any array bounds checking - this makes array use in C very fast, but it also makes it a little dangerous. A common error is to index beyond the end of an array. For example:

```
int x=0;
int a[10];
int i;
for (i = 0; i <= 10; i++) {
    a[i]=27;
}
printf("The value of x is %d\n", x);
```

This code fragment produces the following output:

```
aardvark.cs.ucl.ac.uk: ./foo
The value of x is 27
```

But the value of `x` was set to 0, so what happened here? The `for` loop set the values of `a[0]` to `a[10]` to all be 27. But the declaration of `a` only allocated enough memory for `a[0]` to `a[9]`. When we wrote the value of 27 into `a[10]`, this overwrote the next four bytes of memory. As the memory on the stack on this system is allocated top-down, the next four bytes of memory happen to correspond to the memory used by variable `x`, so the value of `x` is overwritten.

The moral here is to be very careful with array subscripting - the compiler won't help you out. Out-by-one errors are notoriously common, and in C then can have subtle adverse side effects without necessarily causing the program to immediately crash.

2.9 Pointers and Addresses

In Java you have no direct access to memory - all accesses to memory have to go through the API for one variable or another, and all accesses are checked. C gives you direct access to memory, as you can see from the array example above.

C also gives you access to *memory addresses*. For example, if you want to find out the address in memory where a variable is stored, you can do the following:

```
int x=0;
printf("The address of x is %x\n", &x);
```

The notation `&x` means “the address of x”. This code fragment gives the following output (the address is printed in hex):

```
The address of x is bfbff3ec
```

In C, an address is itself a data type, known as a *pointer*. Pointers can be assigned and copied just like any other variable, and they can be *dereferenced* so that the memory being pointed to can be read or written. For example:

```
int x=0;

int *p;
p = &x;

printf("The address of x is %x\n", p);
printf("The value of x is %d\n\n", *p);

*p = 27;
printf("The value of x is %d\n", x);
```

In the example, a regular variable `x` is declared. A pointer variable is also declared: `int *p;` In a declaration, the “*” used like this declares a pointer to an type, rather than the type itself, so this declares that the variable called `p` is a variable that holds a pointer to an integer.

The line `p = &x;` takes the address of `x` and stores this in pointer `p`.

The code prints out the value of the pointer `p` (ie the address of `x`), and then it prints out the value of the memory pointed to by `p`. The notation `*p`, used like this outside of a declaration, means to dereference the pointer, and hence in this context it gives the value of `x` because `p` points at `x`.

Then we have the line: `*p = 27;`

This dereferences `p` and sets the value of the memory pointed to by `p` to be 27. So in this case it sets the value of `x` to be 27 because `p` points at `x`.

Thus the output of the code fragment above is:

```
vulture.xorp.org: ./foo
The address of x is bfbff3ec
The value of x is 0

The value of x is 27
```

Pointers are fundamental to the use of C for any non-trivial program. They're used to pass around references to data, they're used for all dynamic memory allocation, and they're a fundamental part of handling text strings in C.

2.10 Dynamic Memory Allocation

So far, we have seen variables and arrays allocated on the *stack* - that is the storage for the variable will be allocated when the variable is declared, and will be removed when the program leaves the block where the variable is in scope. Often we need to store data in one part of a program, leave that part of the program, and without copying it, pass the stored data to another part of the program. To do this in C we need to dynamically allocate memory on the *heap*.

The main way to dynamically allocate memory in C is to use the `malloc()` function, or one of a family of closely related functions. For example:

```
#include <stdlib.h>

main()
{
    int *buffer;
    int i;

    buffer = malloc( 10 * sizeof(int) );

    for (i = 0; i < 10; i++) {
        buffer[i] = i;
    }

    for (i = 0; i < 10; i++) {
        printf("p[%d] is %d\n", i, buffer[i]);
    }

    free(buffer);
}
```

Running this program gives:

```
aardvark.cs.ucl.ac.uk: ./foo p[0] is 0
p[1] is 1
p[2] is 2
p[3] is 3
p[4] is 4
p[5] is 5
p[6] is 6
p[7] is 7
p[8] is 8
p[9] is 9
```

In the example above, we declared a pointer called `buffer`, and then used `malloc()` to allocate enough memory to store ten integers.

After this we then accessed this memory just as if it were an array. This might seem surprising at first, but in C, an array is essentially just a block of memory and a pointer to the base of that memory. Thus array indexing operations can be used to access any block of memory, even one created using `malloc`.

However, if the memory is allocated using `malloc`, it will not automatically be freed when the pointer to it is removed. Thus when you have finished using the memory, you need to explicitly free it using a call to `free()`, as shown above.

Failure to free memory will result in a memory leak - your program will get bigger and bigger until it crashes because it can't get any more memory.

Equally bad, but more subtle, is when you do free the memory, but leave behind a pointer to that memory, and later access the memory using that pointer. It is likely that before long the memory you freed will be reallocated to your program when you call `malloc` again. Your old pointer will now point at memory being used elsewhere, and accessing this memory will have seemingly random results.

In general, dynamically allocating memory, freeing it correctly, and making sure it is never referenced again after being freed, are among the hardest things to get consistently right in a C program. But dynamic memory allocation is also essential in any non-trivial program, so great care must be taken to make it very clear who is responsible for freeing memory.

2.11 Text Strings

Most programming languages have support for text strings as first class types. However, most CPUs have no special support for strings, so all this entails quite a bit of work by the compiler. C generally closely reflects the hardware capabilities, and so C doesn't have much in the way special support for text strings in the base language. Thus strings in C are implemented using arrays and pointers, as shown below:

```
#include <stdlib.h>

main()
{
    char *s1;
    char s2[80];
    char *s3;

    s1 = "Hello";
    s3 = malloc(80);

    printf("Enter your first name:\n");
    fgets(s2, 80, stdin);

    printf("Enter your last name:\n");
    fgets(s3, 80, stdin);

    printf("%s %s %s\n", s1, s2, s3);
    free(s3);
}
```

Some sample output is shown below:

```
vulture.xorp.org: ./foo
Enter your first name:
Donald
Enter your last name:
Duck
Hello Donald Duck
```

In this example, three variables are declared, and all three can be regarded as being strings.

First, `s1` is declared to be of type “char *”. Then `s1` is set to point at the literal string “Hello”. This literal string is part of the program code; it can't be changed. The assignment simply sets `s1` to point to this constant string. You can read from this memory, hence print out the string, but not write to it.

Second, `s2` is declared to be an array of 80 chars.

Third, `s3` is declared to be of type “char *”, but then we `malloc()` 80 bytes of memory, and set `s3` to point to it.

The only practical difference between `s2` and `s3` is that `s2`'s memory is on the stack, whereas `s3`'s memory is on the heap. If we wanted to, we could subsequently re-assign pointer `s3` to point to some other memory, but we can't do this with `s2` because of the way it's declared.

In any event, at this stage both `s2` and `s3` are effectively strings. We can read text into them - in this case we read it in using a standard library call `fgets()` which reads from standard input (usually the keyboard).

And we can print out the values stored, in this case using the `%s` formatting code to `printf`.

In C, strings are *null-terminated*. For example, `s2` above has 80 bytes of storage allocated, but we only read in 6 characters (“Donald”). A 7th character known as a NULL (it has character code 0) is then appended to the end of the string. When `printf` prints the string, it prints each character until it reaches the NULL, and then it stops. All the string handling functions in C expect a string to be null-terminated in this way.

It is important to understand that a variable of type `char *` is only ever a pointer from the compiler’s point of view. It’s easy to forget this. For example, if you want to compare two strings, you might incorrectly type:

```
WRONG:
char s1[80];
char s2[80];

printf("Enter first word:\n");
fgets(s1, 80, stdin);

printf("Enter second word:\n");
fgets(s2, 80, stdin);

if (s1 == s2) {
    printf("Words are the same\n");
} else {
    printf("Words are different\n");
}
```

No matter what you enter, this will always say that the strings are different. This is because `s1` and `s2` are really pointers, and in this case they point to different memory. The comparison `if (s1 == s2)` merely checks if the pointers point to the same memory - it does not compare the contents of that memory.

To actually compare the contents of strings, you need to use the `strcmp()` function from the standard library:

```
CORRECT:
#include <string.h>

...

char s1[80];
char s2[80];

printf("Enter first word:\n");
fgets(s1, 80, stdin);

printf("Enter second word:\n");
fgets(s2, 80, stdin);

if ( strcmp(s1,s2) == 0 ) {
    printf("Words are the same\n");
} else {
    printf("Words are different\n");
}
```

To use string functions, you need to include the `string.h` system header file, as shown above. We’ll discuss

the C preprocessor more in Section 3.

Other useful string functions include:

- `strcmp()` - compare two strings.
- `strncmp()` - compare the first n characters of two strings.
- `strdup()` - save a copy of a string.
- `strncpy()` - copy a string.
- `strncat()` - concatenate strings.
- `snprintf()` - formatted print into a string.

More details of these functions can be found in the man pages.

You may also discover that variants of these functions such as `strcpy()` (as opposed to `strncpy()`) exist. The difference is that you specify the amount of memory available with `strncpy()` whereas you don't with `strcpy()`. Good code **always** uses `strncpy()`, `strncat()`, and `snprintf()` rather than `strcpy()`, `strcat()` and `sprintf()`. The versions that don't require a size to be specified are just too easy to accidentally overflow the available memory. In networked programs, such simple buffer overflows frequently lead to serious security problems.

2.12 Data Structures

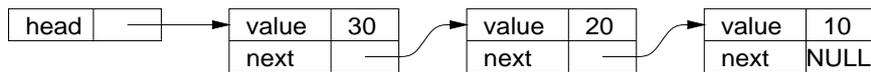
It is a very common requirement to need to group a set of information together and handle it as one entity. In Java, this is done using classes, but C is not object-oriented, and so has no concept of a class. What C does have is `struct`, which has the same ability as a class to group data, but without the accompanying class access methods.

For example, suppose you want to read in a series of integers from the keyboard, and store them in a linked list, and print them out in reverse order. Each element in the list needs to hold an integer and a pointer to the next element. You can use a `struct` for each list element, as shown in the program on the next page.

The output from running this program looks like:

```
aardvark.cs.ucl.ac.uk: ./foo
Enter an integer: 10
Enter an integer: 20
Enter an integer: 30
Enter an integer: Finished
Value: 30
Value: 20
Value: 10
```

In this program we declared a `struct` called `element`, which contains two fields: an integer called `value` and a pointer to another `element` called `next`. We then use this `struct` to build our linked list like this:



The variable `head` is used to keep track of the start of the list.

Each time through the while loop, we read a line of text from standard input, and if the end of file has not been reached, then we convert this text into an integer using `atoi()`. We create a new list element to store the integer in, using `malloc()` to allocate the memory.

Next we set then `value` field of the new element to be the integer we just read in. This is done in the line:

```
new_element->value = number;
```

The notation `new_element->value` means the field called `value` pointed to by the pointer called `new_element`.

We then need to link the element into the existing list. As we want to print out the numbers in reverse order, we'll put the new element at the start of the list. To do this, we set the `next` field of the new element to point to what used to be the first element of the list. Then we change the `head` pointer to point to the new element we just added.

Finally, when there's no more input, we loop through the list printing out the values, until we find an element whose `next` pointer is `NULL`. The value `NULL` is used to denote a pointer that doesn't point to anything; it has the numeric value of zero.

```

#include <stdio.h>
#include <stdlib.h>
#define MAXLEN 80

/* struct to hold an element of the linked list */
struct element {
    int value;
    struct element* next;
};

main()
{
    /* variable to hold the head of the list*/
    struct element *head = NULL;

    /* temporary variable for printing out the list */
    struct element *p;

    while (1) {
        int number;
        char buffer[MAXLEN];
        struct element *new_element;

        /* read in an integer from the keyboard */
        printf("Enter an integer: ");
        fgets(buffer, MAXLEN, stdin);
        if ( feof(stdin) ) {
            printf("Finished \n");
            break;
        }
        number = atoi(buffer);

        /* create a new list element */
        new_element = malloc( sizeof(struct element) );
        new_element->value = number;

        /* add element to the list */
        new_element->next = head;
        head = new_element;
    }

    /* print out the numbers we stored */
    p = head;
    while ( p != NULL ) {
        printf("Value: %d\n", p->value);
        p = p->next;
    }
}

```

Declaring Structs on the Stack

In the linked list example, we allocated heap memory for a struct using malloc, and accessed its fields using “->” to follow a pointer. Structs can also be allocated on the stack, as you might with a regular variable. For example, in the system header file `sys/time.h`, a struct called `timeval` is defined:

```
struct timeval {
    long tv_sec;    /* seconds */
    long tv_usec;  /* and microseconds */
};
```

To find out the current time, we might use this as follows:

```
#include <stdio.h>
#include <sys/time.h>

main()
{
    /* a struct timeval to hold the current time. */
    /* timeval is defined in sys/time.h */
    struct timeval time;

    double fractionaltime;

    /* get the current time and store it in "time" */
    gettimeofday(&time, NULL);

    /* convert the time to a fraction */
    fractionaltime = time.tv_sec + (time.tv_usec/1000000.0);

    printf("%f seconds since 1st January 1970\n", fractionaltime);
}
```

Which might produce output such as:

```
vulture.xorp.org: ./foo
1096801178.236713 seconds since 1st January 1970
```

In this example we declared the variable `time` to be a struct `timeval` just like we might declare any other local variable on the stack.

We then called the standard system call `gettimeofday()` to get the current time, passing in a pointer to `time` (we get a pointer by using “&time”). The call to `gettimeofday()` reads the system clock to find out the current time, and fills in the fields of the struct `time` using the pointer we gave it.

In this example we want to print out the time as a decimal number, so we convert it from the struct `timeval` representation by reading the fields of `time` individually. The notation `time.tv_sec` means the `tv_sec` field of the struct called `time`. Finally we print out the value.

The notation from this example of `time.tv_sec` is functionally similar to the `new_element->value` notation in the previous example. The difference in notation is because `new_element` was a pointer to a struct, whereas `time` actually *is* a struct.

2.13 Functions and Procedures

C is a procedural language, which means that the main code-structuring principle is the *procedure* or *function*. In contrast, Java is object-oriented, and its main code-structuring principle is that of the object or class. In this document, we mostly use the terms *procedure* and *function* interchangeably.

A function is simply a named section of code that you can define. Parameters can be passed into the function, and a single result is returned at the end. Functions allow the same code to be called from multiple places in your program without copying the code. But just as importantly, used right they make your code much more readable and easier to maintain than it would otherwise be.

Java methods within classes are very similar in concept to C functions, but the details differ somewhat.

A simple example of a function is shown in the example below.

```
#include <stdio.h>

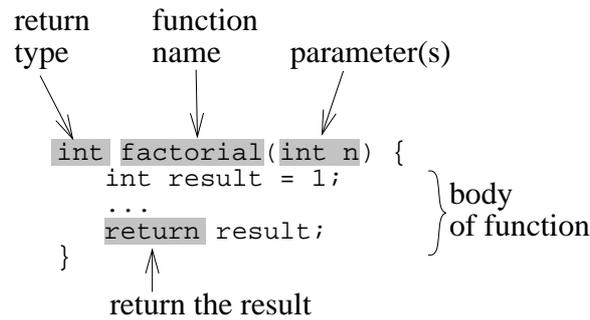
int factorial(int n) {
    int result = 1;
    int ctr;
    for ( ctr = 2; ctr <= n; ctr++) {
        result *= ctr;
    }
    return result;
}

main()
{
    int i;
    int fact_i;
    for (i = 1; i <= 10 ; i++) {
        fact_i = factorial(i);
        printf("factorial(%d) is %d\n", i, fact_i);
    }
}
```

The output of this program is:

```
vulture.xorp.org: ./foo
factorial(1) is 1
factorial(2) is 2
factorial(3) is 6
factorial(4) is 24
factorial(5) is 120
factorial(6) is 720
factorial(7) is 5040
factorial(8) is 40320
factorial(9) is 362880
factorial(10) is 3628800
```

In this example, we define the function `factorial()`, which takes a single parameter “int n”, and also returns a value of type `int`:



We then call this function multiple times from `main()` to calculate the factorials of the integers between one and ten.

It is important to realize that parameters in C are passed by *value*, whereas in Java they are passed by *reference*. When a parameter is passed by value, a *copy* is made, and is passed into the function. Here is a simple example to illustrate this:

```

#include <stdio.h>

void silly_example(int x)
{
    x = 2;
    printf("During call, x is changed to %d\n", x);
}

main()
{
    int i = 1;
    printf("Before call, i is %d\n", i);
    silly_example(i);
    printf("After call, i is %d\n", i);
}
  
```

```

vulture.xorp.org: ./foo
Before call, i is 1
During call, x is changed to 2
After call, i is 1
  
```

As you can see, within a function you can change the value of parameters just like you can with any other variable. Such changes have no effect on the original value of the variable passed into the function. By the way, the use of `void` here as the return type indicates that the function doesn’t return any value.

What if you *did* want to change the *original* value of the variable passed into the function? Then you need to pass in a pointer, and the function needs to change the variable using that pointer:

```
#include <stdio.h>

void another_silly_example(int *x, int *y)
{
    *x = 2; *y = 3;
    printf("During call, *x is changed to %d, *y is changed to %d\n", *x, *y);
}

main()
{
    int a = 1, b = 1;
    printf("Before call, a is %d, b is %d\n", a, b);
    another_silly_example(&a, &b);
    printf("After call, a is %d, b is %d\n", a, b);
}
```

```
vulture.xorp.org: ./foo
Before call, a is 1, b is 1
During call, *x is changed to 2, *y is changed to 3
After call, a is 2, b is 3
```

In this case we passed in pointers to `a` and `b`. Within the function, because of the ordering of parameters, the pointer to `a` was called `x` and the pointer to `b` was called `y`.

We then changed the values of the memory pointed to by these pointers (`*x = 2; *y = 3;`). Thus when we return from the function, the values of `a` and `b` are seen to have changed.

String Parameters

As we saw in section 2.11, C does not have a first class string datatype. Thus when strings are passed as parameters, they're usually passed as type `char*`:

```
#include <stdio.h>

int string_length(char *s) {
    int i = 0;
    char *p = s;
    while (*p++ != '\0')
        i++;
    return i;
}

main()
{
    char* hw = "Hello World!";
    int len;

    len = string_length(hw);
    printf("The length of '%s' is %d\n", hw, len);
}
```

```
vulture.xorp.org: ./foo
The length of 'Hello World!' is 12
```

As a string was already really just a pointer to a region of memory, passing it into a function doesn't really change anything - it's still a pointer to the same memory. But of course if you change the contents of that memory from within the function, that memory will remain changed after the function.

Function Prototypes

C programs can be compiled in stages if you split your code into multiple files. If you do this though, the compiler will be confused if you try to reference a function in one file from within another file, because it won't know what types your functions return. To avoid this problem, you can declare a *function prototype* to help the compiler out. For example, suppose we wanted to split the string example above into two files called `strlen.c` and `main.c`, then we could do the following:

```
strlen.c:

int string_length(char *s) {
    int i = 0;
    char *p = s;
    while (*p++ != '\0')
        i++;
    return i;
}
```

```
main.c:

#include <stdio.h>

int string_length(char *s);

main()
{
    char* hw = "Hello World!";
    int len;

    len = string_length(hw);
    printf("The length of '%s' is %d\n", hw, len);
}
```

The actual `string_length()` function is defined in `strlen.c`, but before we use this function in `main.c` we declare a prototype:

```
int string_length(char *s);
```

This tells the compiler what the parameters and return value are for the `string_length()` function when it's separately compiling `main.c`.

To actually compile the code, you could use the `-c` flag to `gcc`, which indicates to compile the module specified to a `.o` object file, but not to link it. For example `strlen.c` gets compiled to `strlen.o`. Then the different object files are all linked together using `gcc` to produce the final executable program. Thus:

```
vulture.xorp.org: gcc -c strlen.c
vulture.xorp.org: gcc -c main.c
vulture.xorp.org: gcc -o foo strlen.o main.o
```

2.14 Input, Output and File Handling

We have already seen some examples of input and output, including `printf` for formatted printing and `fgets` to read a string. In this section we'll cover input, output, and file handling in a little more detail.

By default, a C program on Unix has access to three open *file descriptors* or *streams*:

- `stdin` - the standard input, which is read-only.
- `stdout` - the standard output, which is write-only.
- `stderr` - the standard error, which is write-only.

If the program was started from the command line, then by default, `stdin` will receive input from the keyboard, and both `stdout` and `stderr` will print to the terminal or terminal window.

The `printf` command sends formatted output to `stdout`. It takes a variable number of parameters. The first parameter is always a formatting string which includes formatting flags dictating how the remaining parameters should be interpreted. The formatting flags begin with a “%” and take forms like “%d” or “%5.2f”. Each flag consists of a basic code, and various modifiers indicating how that value is to be formatted or otherwise modified. The most common basic codes are:

- `d` - print an integer as a signed decimal number.
- `u` - print an integer as a unsigned decimal number.
- `o` - print an integer as a unsigned octal number.
- `x` - print an integer as a unsigned hexadecimal number.
- `f` - print a double as a decimal number.
- `c` - print a single character.
- `s` - print a `char*` argument as a null-terminated string.

Additionally a field width and a precision may be specified before the basic code. For example:

- `%3d` - print an integer, padding the left with spaces if necessary to make at least three characters.
- `%5.2f` - print a double, padding the left with spaces to make five characters if necessary, and printing two digits to the right of the decimal point.

Many more flags and modifiers than these are available - see the `printf` man page for more details.

The `fprintf` command is similar to `printf`, but takes an additional parameter specifying which file descriptor to print to. For example:

```
fprintf(stderr, "Something bad happened!\n");
```

would print to the `stderr` stream.

To write a raw block of bytes to a stream, `fwrite` can be used. For example, to write 80 lots of 4 bytes to `stdout` from memory pointed to by `int *buffer`, you could write:

```
fwrite(buffer, 4, 80, stdout);
```

There are many different ways to read from a stream, including:

- `fread` - read a raw block of bytes from a stream.
- `fgets` - read a string until a newline is found.
- `fgetc` - read a character from a stream.

More details can be found in the relevant man pages.

The standard I/O streams are not the only places that you can read input and write output. You can open new streams using `fopen()` to open a stream to or from a file.

For example, to write to a file:

```
FILE *file;
char *filename = "/tmp/foo";

/* open the file for writing */
file = fopen(filename, "w");
if (file == NULL) {
    fprintf(stderr, "File %s could not be opened\n", filename);
    exit(1);
}

/* write to the file */
fprintf(file, "Hello World!\n");

/* close the file */
fclose(file);
```

In this example, we `fopen` the file “/tmp/foo” for writing. This should return us a `FILE` data structure, which we can use as a file handle in subsequent calls. If we are unsuccessful in opening the file, `file` will be `NULL`, and then we print an error message on `stderr`.

If we opened the file correctly, we then use `fprintf` to write to the file, and finally close the file using `fclose`.

It is important not to forget to `fclose()` a file you’ve written to. File operations on files you’ve opened with `fopen` are *buffered*, meaning that the bytes may not immediately reach the file. This improves performance, but it also means that if your program exits without calling `fclose`, the file may be empty or truncated.

We can use `fopen` to open a file for reading in a similar manner. In the example below, we open a file, and then loop through reading a line at a time from the file and printing them out.

```
FILE *file;
char *filename = "/tmp/foo";

/* open the file for writing */
file = fopen(filename, "r");
if (file == NULL) {
    fprintf(stderr, "File %s could not be opened\n", filename);
    exit(1);
}

/* loop while reading a line at a time from the file and printing */
while (1) {
    char buffer[80];
    fgets(buffer, 80, file);

    /* if it's the end of file, break out of this loop */
    if (feof(file))
        break;

    printf("%s", buffer);
}

/* close the file */
fclose(file);
```

In addition to buffered I/O, where the file handle is typically a pointer to a `FILE` data structure, C also supports unbuffered I/O. The file handle for unbuffered I/O is an integer *file descriptor*. Such a file descriptor is returned by the `open()` system call (for file operations) or by the `socket()` system call (for networking operations). File descriptors are closed using the `close()` system call, and reading and writing to file descriptors makes use of the `read()` and `write()` system calls. The basic semantics of these calls are the same as for buffered I/O, but the order and types of the parameters differ. Again more details can be found in the relevant man pages.

These unbuffered operations are a slightly lower level interface than the buffered interface. Most often you will use buffered I/O for file operations and unbuffered I/O for networking operations.

3 THE C PREPROCESSOR

The process of C compilation has three main passes - the preprocessor pass, the compilation pass, and program linking.

The preprocessor is practically a language in its own right. It allows you to include header files into your code, define constants and macros, and conditionally compile different parts of your code.

3.1 Including Header Files

In many of the examples so far, you have already seen:

```
#include <stdio.h>
```

This includes the contents of the system header file `stdio.h` into your source code file before starting the compilation pass. The file `stdio.h` contains definitions for things like the `FILE` file handle seen in section 2.14, and function prototypes for standard library functions like `printf`. If you don't include `stdio.h` then the compiler won't know how to compile these functions.

You can define your own header files if you have definitions you want to share across multiple C source files. For example, if you have the header file `constants.h`, then you might include this file at the top of each of your source files using:

```
#include "constants.h"
```

The only difference between this and the `stdio.h` example is the use of quotes instead of angle-brackets - this affects where the compiler searches for the file to include. Typically you use quotes for your own header files and angle-brackets for the system header files.

3.2 Preprocessor Constants and Macros

The preprocessor also supports the definition of constants and macros. Constants are very simple:

```
#define PI 3.141592
#define UCL "University College London"
```

The preprocessor will simply do a text substitution, replacing every occurrence of the characters `PI` with the characters `3.141592`. These preprocessor statements are not regular C commands, so unless you want the substituted text to include a semicolon, do not put a semicolon at the end of a `#define` statement.

The preprocessor command `#define` can also be used to define macros to make your code easier to read. For example, if you wanted to define a macro to return the maximum of two numbers, you could define the macro:

```
#define max(X, Y) ((X) > (Y) ? (X) : (Y))
```

This defines a macro that takes two arguments. In this case it uses C's `?/:` notation to represent a simple if/else statement (see section 2.3). For example, you might use it like:

```
#define max(X, Y) ((X) > (Y) ? (X) : (Y))

a = max(c + d, e + f);
```

After preprocessing, the C compiler will see:

```
a = (( c + d ) > ( e + f ) ? ( c + d ) : ( e + f )) ;
```

3.3 Conditional Compilation

Sometimes you need to have two different fragments of code in the same source code file, and compile one or the other depending on circumstances. Common uses are to conditionally enable debugging, or to use different code depending on the operating system on which the code is being compiled. The C preprocessor can do this for you using the `#ifdef` and `#if` commands. For example:

```
#include <stdio.h>

/* #define DEBUG */

#ifdef DEBUG
#define debug(s) printf("%s\n", s);
#else
#define debug(s)
#endif

main()
{
    int a = 1;
#ifdef DEBUG
    printf("a is %d\n", a);
#endif

    debug("reached here!\n");
}
```

This program will print out no output. But all you need to do is uncomment the line `#define DEBUG` and all the debugging output will be enabled.

Note that the macro `debug(s)` has two possible definitions. In the debugging case, the program gets compiled including the relevant `printf()` statements. In the no-debugging case, all the debugging information is removed by the preprocessor before the C compilation pass, so all this debugging information doesn't slow down the final program.

In this tutorial we have only touched on a few common uses for the preprocessor. Some C programmers make very extensive use of the preprocessor to balance readability with performance. However, excessive use of the preprocessor can make C programs hard to understand, because it's not so obvious where to look for where something is defined. For an extreme example, take a look at this winning example from the International Obfuscated C Competition by Vern Paxson:

<http://www.icir.org/vern/ioccc92.c>