

C: Part 13

Language of the 'C'

Following on from last month's article, Steven Goodwin, in this, the final part of our C tutorial, looks at how C can be unreadable, and why it becomes like that. [BY STEVEN GOODWIN](#)



A good friend of mine from University knows all the bad parts of town. She knows where the fights will be, and who'll be dealing in what, and where. Thing is – she's a nice girl! What is she doing knowing about the dodgy parts of town? Her answer was my inspiration; knowing where not to go, stops you from doing it. So in this article, I will tell you why bad code is written, how to understand it, and how to stop yourself from doing it.

Purpose In Life

The easiest case to understand as to why code is unreadable is where it was written so intentionally. This could be because it was written to demonstrate an interesting (mis)use of the language, or intended for a programming competition such as the IOCCC (see boxout). Some source code will be obfuscated on purpose to hide it's meaning along with any clever, novel, or interesting technology contained within.

This is sometimes referred to as 'shrouded source' where, although the code is available to the end user (enabling it to be distributed openly, needing only a recompile), it is

impossible to read and understand, since the meaning of the code has been perverted. This can happen by using obtuse (or even wrong) variable and function names (perhaps of a single letter), the removal of white space, an over-use of macros or any number of other techniques. Such code even makes Perl look readable!

Understanding such code is a considerable task, and not to be undertaken lightly. Only in exceptional cases (i.e. you're paid to, or the code is a puzzle you "just have to work out") is it worth trying to understand such code. Your time is better spent solving the problems yourself, and re-writing it in a sensible (preferably open), fashion.

In My Defence

For code that is unintentionally obfuscated, the most common cause is casual. Code is written in a particular style 'just because that's how a particular programmer writes' – the geek equivalent of Finnegans' Wake! Over years of programming, people drop into various habits. Some good. Some bad. All of them are completely natural to the person in question, but require more

thought by everyone else. Let us take a simple case:

```
if (fp = $$
    fopen("/etc/convert.conf", "r"))
    { /* process the file */ }
```

This is something we've seen before, and is quite a common structure for opening a file and handling its contents, should it exist. We've seen it before, so we're used to it. If we had not, it might be a different story. So, what if the expression was something with which we are unfamiliar? Here, the use of language is identical, but the situation is not.

```
if (n = CountItems())
    { /* Is this supposed to check
       the integrity of 'n'? */ }
```

This unintentional obfuscation can show its roots in a number of places, but because they are all quirks of the original programmer (which you are unlikely to know on a first hand basis) and so it gives you two things to think about, not one. For example, a programmer may have come from a different language to C, and was forcing his ideas into 'C' and

of expressions that use (or even rely on) precedence. Stepping back to the example I gave when discussing precedence, notice how ill formatting would confuse the issue.

```
ans = 10*x / 5*y;
```

Clever Trevor

Another case of obfuscation is where the code is cleverer than it needs to be. This can manifest itself in a couple of ways. Consider a loop to compute the sum of every number between 1 and 100.

```
int iTotal = 0;
for(i=1;i<=100;i++)
    iTotal += i;
```

This is simple, understandable and very straightforward. However, if the programmer knew about Gauss, he might have written:

```
for(i=1;i<=50;i++)
    iTotal += 101;
```

or even,

```
iTotal = 101 * 50;
```

This works because a mathematician named Gauss (1777-1855) deduced that working the arithmetic from both ends at once reduces the sums complexity. It becomes a list of 50 sums, all of them equal to 101, because $1 + 100 = 2 + 99 = 3 + 98$ and so on. Very simple to understand when you're writing it, but much more difficult to read; especially in the general case. If each step of the process is given a comment then there is some salvation – but there rarely is! This type of obfuscation requires you to understand the language used to implement the problem...and the method used to solve the problem. Using mathematical identities is often necessary to improve performance of software, but you should always document those methods so others can understand the code to make enhancements (and bug fixes) easily.

Code can also be too clever outside the field of mathematics, and may rely on assumptions that are implicit in the code or data.

```
j = 0;
for(i=1;i<100;i++)
    if (/*some condition*/)
        j=j?j:i;
```

This loop finds the first case where the condition is true, and produces the position in the variable 'j'. This a good example of bad coding because:

- 1) The variables are not meaningful
- 2) It relies on 'i' never starting at 0
- 3) It over-condenses the code without a comment

Again, splitting the conditional ?: will help understand this.

```
if (j)
    j = j;
else
    j = i;
```

So, if 'j' is non-zero, nothing will happen and j will remain unchanged. Otherwise (j = 0) it will be assigned to the value of 'i'. At this point it is no longer zero (because of the assumption that i always starts at 1), and so is unaffected for the rest of the loop.

Bright Lights, Big City

Trying to out-compile the compiler can also make code unreadable. This is where the writer has learnt / discovered / worked out how the compiler (or even CPU) will handle the code, and so has written parts of the program in such a way to produce code that gives better performance. This has the same symptoms as code that has been intentionally obfuscated, but suffers from the fact that not even the original author knows why it had to be done in that particular way. The irony is that any performance gained from compiler A is not valid on B (or even between versions of the same compiler)! One good example here is the code:

```
tmp = *ptr;
for(i=0;i<1000;i++)
    ptr[i] = 0;
```

Here the 'tmp' variable is never used, and appears to be redundant, so one might be tempted to remove it. However, on processors such as the Pentium, reading the memory location at 'ptr' may cause that section of memory (8K or so)

to be brought into the cache. Then, when the loop writes to memory, it does so to the (fast) cached version, and not main memory, as it might have done without the 'tmp' line.

Billericay Dickie

The flip side of code trying to be too clever is code that is not clever at all. This could be because it uses the wrong method, but gets the right answer, or uses the right methods but in the wrong cases. Consider this example I found in some code on a Windows machine.

```
len = strlen(szFilename);
szFilename[len - 4]=0;
```

It gets the right answer (most of the time) but uses the wrong method! The intention was to remove the file extension from szFilename, before concatenating a different one onto the end. This is confusing because that's not what it actually does: imagine if the filename did not have an extension!

A similar case happens with comments. Comments are good unless they disagree with the code. Or the variables used within the code are given names that do not apply to their job. Neither happens when writing a program, but as it changed and new features are added, the comments are not updated, and a 'LastItem' variable is now used (or even re-used) as a count of the number items – and slowly the code becomes less clear than it once was.

Both situations should be rectified by making the code do what it is supposed to, in the way that it is supposed to do it! Removing a filename extension means taking all characters after the dot – so look for the last dot (I've even seen code that removed the first dot, causing other problems!) and remove those characters. If you're writing an interactive application and you notice there are 12 characters after the dot (i.e. it probably does not have an extension, just a dot in the name) you can always report the error to the user and let them confirm the action.

The 'write as you mean' rule is the best way to code, ensuring both man and machine understand what's going on. Consider the 1 to 100 summing loop above. If we'd written it as,

End Note Sidebar

As this marks the final part of 'Language of the C', I would like to take the opportunity to thank a few people. Notable John Hearn for encouraging me to write it, John Southern & Colin Murphy for letting me write it, Alan Troth for reading it (and forcing me to re-write it), and TULS for the beer and curry!

```
for(i=0;i<100;i++)
    iTotals += i+1;
```

This may fit in with C's zero-indexing policy, but it does not make (as much) sense because the meaning is lost. The number zero is not part of the question, so it should not be integral to the finding of the solution. And do not use the excuse of 'code will run slow' when cutting corners, either. As Knuth once said, "Premature optimisation is the root of all evil".

Old Before I Die

Some programs are difficult to read because of their over-reliance on the C pre-processor. Especially by beginners who have come from Pascal, say, and would still rather type 'begin' instead of '{' to start each code block. It is not unknown for them to start each file with:

```
#define begin {
#define end }
```

This, although quaint, is ultimately confusing to the reader (since the word begin looks like it should be a function or a variable), and prevents the author from moving away from Pascal.

They will never have to think in C and so are likely to implement substandard solutions (that are by their nature more difficult to read) because they are not considering (and working to) the strengths of the language. In these cases, you need to pre-process the source files to expand the macros into something that looks more like C.

In extreme cases people may be working with programs that have been converted, line-by-line, from another language into C.

These conversions are often the technical equivalent of badly translated Kung Fu movies – the words may be

correct, but they do not make sense in context! Depending on the importance of the software (and the salary involved!) it may be worth re-writing them, not from the source, but from the original algorithms.

Old Red Eyes is Back

One case of obfuscation that happens (but rarely) involves old code. When software has been ported from an old system, or you are working on an old Unix system, there may be some historical features that can be confusing.

The programmer might have used functions that no longer exist in the standard libraries, or those that have since been replaced or renamed. One example is strchr.

This used to be called index, and might exist in some code. Now, since this function has not been documented for many years one might be tempted to look for it outside of the standard libraries, and not find it.

In extreme cases, you might be working on a compiler that supports features of the old K&R style of C. On these systems, the language was much younger than it is now, and supports strange syntax such as:

```
int x 1;
```

Which is actually a simple declaration and assignment that we know as:

```
int x=1;
```

Similarly, code like

```
x -= 1;
```

Would (on an ANSI C compiler) assign minus one to x. However, in the 'olden days', it would decrement x by 1, because -= was the original form of -=.

With Linux being comparatively new, this should be a rare case, especially as gcc does not support it.

The End Of the World

Despite the fact that this series has taught the C language and its many (varied) uses, it is still possible to construct legitimate code that looks wrong, strange, or confusing. My favourite example of this is Duffs Device.

```
register n = (count + 7) / 8;
/* count > 0 assumed */
switch (count % 8)
{
case 0: do { *to = *from++;
case 7:      *to = *from++;
case 6:      *to = *from++;
case 5:      *to = *from++;
case 4:      *to = *from++;
case 3:      *to = *from++;
case 2:      *to = *from++;
case 1:      *to = *from++;
            } while (--n > 0);
}
```

(Copyright 1984, 1988, Tom Duff)

(The 'to' address is mapped to a device, and therefore it does not need to be incremented within the program).

Any language powerful enough to produce original code, is also (by its very nature) powerful enough to produce oddities or quirks of use that were not considered when originally designing the language.

No language course could ever hope to cover every single obtuse case of syntax in existence – and there's always one programmer who will find more evil ways of abusing the language. In these cases, you have little choice but to work through the code, line by line, function by function, understanding what the compiler would do in these situations and mimic it. This technique (called dry-running) is carried out by language lawyers to understand and demonstrate vagrancies of a particular language. And you should to. ■

IOCCC

The International Obfuscated C Code Contest. A yearly competition to (ab)use C in the most esoteric manner possible. The winning entries are somewhat scarier than the 'simple' examples given here. www.ioccc.org

THE AUTHOR

The language of 'C' has been brought to you today by Steven Goodwin and the pages 68–72. Steven is a lead programmer, who has just finished off a game for the Nintendo GameCube console. When not working, he can often be found relaxing at London LONIX meetings.