

C Programming Language Overview

Author: Eric Laroche
Copyright © 2004 Eric Laroche

This paper may serve as a short overview of the C programming language.

C is a general-purpose programming language which features economy of expression, modern control flow and data structures, and a rich set of operators. C wears well as one's experience with it grows. -- *K&R2* [*The C Programming Language* by Brian W. Kernighan and Dennis M. Ritchie] ¹

C *stands* for effectiveness of language, good style, sound design. ¹ C typically uses a *compiler* ². C is *case-sensitive* [in its keywords and identifiers ³].

It is recommended to *skip* this text's *footnote texts* on the *first sweep*.

¹ [*The C Programming Language*, Brian W. Kernighan, Dennis M. Ritchie, Prentice Hall, 2nd Ed. 1988, ISBN 0-13-110362-8] ⁴

² A *compiler* is the tool [program] to translate a [higher-level] [programming] language [to a lower-level language, often into object files ⁵].

³ *Identifier case-sensitiveness* is not guaranteed for the link ⁶ phase.

⁴ For other *C standards* [ISO, ANSI], see [links](#).

⁵ An *object file* is the output of a *compiler* ², an *assembler* or a similar tool, usually [machine [language]] code that is *native* to a processor [or alternatively *byte-code* for some *abstract processor architecture* [that is usually *interpreted* on software that is called a *virtual machine*]].

⁶ *Linking* is the process of generating a [binary] executable, a shared library, a re-linked object [file] ⁵ or similar, from *object files* ⁵, *libraries* ⁷, start-up object code [and possibly additional resources] ⁸.

⁷ A *library* [in the narrower sense] is a collection of compiled ² source code files [so called *object files* ⁵].

⁸ The *link process* is usually *opaque* to a programmer, controlled by the *compiler driver program* ⁹.

⁹ The *compiler driver* [program] steers all phases ¹⁰ of compilation, from source code to binary.

¹⁰ The original *compile phases* were *preprocessing*, *compiling*, *optimizing*, *assembling*, *linking*.

Source code representation

C programming language source code is typically represented in *files* ¹¹. These *C code source files* typically carry a *file suffix* of `.c` ^{12 13 14}. A C file is called a *module* or *compilation unit* ¹⁵.

C source code files do not have [mandatory] header fields [such as a preamble e.g. in the first source code line] ¹⁶.

Interfaces are often represented in *header files*, that typically carry a file suffix of `.h`. Such *interfaces* consist of stuff that is used *inter-modular*, i.e. *between* different C source code files ^{17 18 19}.

¹¹ *Development environments* may chose to represent sources ²⁰ in another way, e.g. in some kind of *repository* ²¹.

¹² Other file suffixes are not usual ²², although compilers often support alternatives by offering *language-specifying compiler options*.

¹³ The *C compiler* [*driver program*] needs *input type information*, to know what to *do* (compile ²³, assemble ^{24 25}, or link).

¹⁴ *Build systems* (such as `make`) typically expect well-known file suffixes too, to deduct file types without file contents lookup ²⁶.

¹⁵ *Compilation units* are *compiled* ² *independently* from other compilation units.

¹⁶ There are no `/etc/magic` entries for C; the `file` tool deducts C code by heuristics only.

¹⁷ The header file (interface) inclusion mechanism is a general one, i.e. one can include any file (text, code).

¹⁸ These *interface files* can be used to specify *build dependencies* in software build environments. Large systems may compile faster if the interfaces are sufficiently fine granular ²⁷.

¹⁹ *Header files* may represent interfaces of a whole library ⁷ [facade pattern].

²⁰ Not just C sources.

²¹ *Repositories* may allow a more *fine granular code resolution* (e.g. function-granular) with more *meta information* (such as modification-author and comments) or just be *faster* than with file system based data accesses. *Repositories* may be implemented as a *file system layer abstraction* ²⁸, which makes it easy to use generic editors and build tools.

²² Unlike with the *suboptimally* ²⁹ chosen `.C` for the *C++ programming language*, that led to [more compatible] alternatives such as `.cc`, `.cxx`, `.cpp`, etc.

²³ Compiling possibly *different* languages.

²⁴ Assembling with or without prior [C-] preprocessing.

²⁵ The original compile *phases* ¹⁰ had assembly sources as an *intermediate form*.

²⁶ [C] *file contents interpretation* is *not an easy thing* to do ¹⁶; make wouldn't be able to deduce a C file's contents.

²⁷ The focus lies however on providing *optimal interface layer abstractions*, not compilation speed.

²⁸ Most environments (typically the operation systems) allow the use of *generalized file systems*, typically by means of a *network file system* interface.

²⁹ The `.C` file suffix [note the upper case] is suboptimal with *case-insensitive filesystems*.

Types

The *C programming language* is a quite ³⁰ ³¹ *strongly typed* ³² ³³ language. This means that *[data] types, variables and functions must* ³⁴ *be declared* ³⁵ before their use.

³⁰ For historical reason, the use of undeclared functions is still supported ³⁶.

³¹ *Type aliases* (as defined by typedefs) may be used interchangeably; boolean and integer values are treated interchangeably ³⁷.

³² C is a *typed* language, but not so much a *type-centered* language (such as C++, which implements means to *protect* data accesses).

³³ *Strongly typed* implies *type safe* ³⁸.

³⁴ Functions *should* be declared before used ³⁰.

³⁵ Note the difference between *declaration* and *definition* ³⁹. A data *definition* actually reserves space for the data, whereas a data *declaration* ⁴⁰ introduces *variable name and type* only. A function *definition* provides the function's implementation, whereas a function *declaration* ⁴¹ introduces a *name and type* only. *Type definitions* do not per se generate code or data ⁴², so a distinction to a 'type declaration' (such as with functions and variables) is not needed; however an *incomplete data type* ⁴³ ⁴⁴, may be considered 'less than a definition'. Since a *definition* is an *implicit declaration* too, the programmer must ensure the explicit declaration is seen at definition location, to ensure *interface consistency*.

³⁶ The use of undeclared functions generally generates *compiler warnings* though.

³⁷ There is no *explicit* boolean data type, `int` is used. This may lead to possibly subtle bugs.

³⁸ *Type safe* means that a lot of problems with types are caught in the compilation stage.

³⁹ A *definition* is an *implicit declaration*, i.e. it is a *declaration in the broader sense*.

⁴⁰ *Data declarations (variable declarations)* are done using the `extern` keyword ⁴¹.

⁴¹ *Function declarations* do not require ⁴⁵ the `extern` keyword as syntactical means, since definition and declaration (implementation and prototype) are distinguished by the former's *code block*.

⁴² Unlike with the [more complex] C++ language, whose type (class) definition *may* implicitly generate code and/or data, e.g. a virtual function table, a default constructor, runtime type information [helper data], etc.

⁴³ *Incomplete data types* are introduced by typedefs on undefined structs ⁴⁶. Such a construct introduces a type *name* only, that can be used as pointer (reference) only ⁴⁴.

⁴⁴ *Incomplete data types* are used as an advanced decoupling (modularization) feature [bridge pattern].

⁴⁵ Function declarations (prototypes) *can use* a `extern` keyword though.

⁴⁶ E.g. `typedef struct T_st T;`

Code and data

C source code mainly consists of *code and data, and type definitions* [mentioned in the previous chapter].

Type definitions appear *early* in a source code file or even separated in a header file [47](#), since they glue parts of code together and are needed by *both parts*.

Code is bound to *function names*, which are globally or module-wide [48](#) visible [49](#).

Data is [at least temporarily [50](#)] bound to *variables*. Data may be bound to fields of some *other data structure* [which itself must be bound to such fields or a variable [51](#)].

Lexical stuff is considered *comments* and *spacing* [52](#). *Comments* are made of *freely formatted text* inside `/*` ... `*/` [53](#).

A *simple preprocessor* supports *macros* [54](#), *file inclusion* [55](#) and *conditional compiling* [56](#).

⁴⁷ *Header files* need to be included sufficiently early [57](#).

⁴⁸ *Module* means a single C source code file, so module-wide means file-wide.

⁴⁹ C does not support *anonymous* functions [lambda expressions].

⁵⁰ E.g. temporarily bound to a `auto` variable and then returned as value.

⁵¹ If this reference *chain* breaks [usually by a software bug], *memory* will *leak* [58](#).

⁵² *Spacing* is usually discussed by *coding rules*.

⁵³ Comments cannot be *nested* [59](#).

⁵⁴ `#define`, with or without macro parameters, `#undef`.

⁵⁵ `#include`

⁵⁶ `#if`, `#endif`, `#elif`, `#else`, `#ifdef`, `#ifndef`.

⁵⁷ Function or data definitions *must* have seen their declarations too, to ensure *consistency*.

⁵⁸ C usually does not include *garbage collecting* that could clean up such *memory leaks*.

⁵⁹ `#if` *preprocessor directives* can be *nested* and can be used to *comment-out code* [60](#).

⁶⁰ E.g. inside `#if 0 ... #endif`.

Functions

A C function (procedure, routine, method) has one entry point and possibly several exit points [61](#). A function establishes a *code block*, which is delimited by a pair of *braces* (`{ }`) [62](#).

Any code block can have its own set of *local variables* [63](#) (which can overwrite (hide) other occurrences of the same name). A function can have argument variables [64](#) (*parameters*), which are always passed *by-value* [65](#); *by-reference* can be implemented by using *pointers* to variables [66](#).

Function definitions cannot [and need not] be nested [67](#).

The *user function* that is called at application startup is `main` [68](#).

⁶¹ *Several function exit points* are implemented by *multiple* `return` statements.

⁶² Sample: `int isqr(int i) {return i * i;}`

⁶³ Such *local variables* can be *automatic* (stack based; exclusive to a stack frame) or *static* (data segment based; typically process-wide shared) [69](#) [70](#).

⁶⁴ C also supports *variable number of arguments* [71](#) [72](#).

⁶⁵ *By-value* enable the programmer to use *any expressions* as *function arguments* [not just variable references [73](#)].

⁶⁶ *'Calls' by-name* can be implemented by using *macros*.

⁶⁷ *Function nesting* would be used to further confine function scope [however, the existing file scope should be narrow enough, especially with small source files], and allow outer-level variable access [which however would weaken data encapsulation].

⁶⁸ `int main(void) {...}` or `int main(int argc, char** argv) {...}` or `int main(int argc, char** argv, char** envp) {...}`.

⁶⁹ Such a *static variable* is also known as *singleton [pattern]*.

⁷⁰ *Static variables* are not per-se *thread-safe*.

⁷¹ As in `int printf(const char*, ...);`, where the first [format-] parameter specifies what follows.

⁷² The *variable arguments* however are not *type-checked*.

⁷³ Not allowing *calls-by-value* probably would imply lots of not-so-elegant [temporary] variables.

Program flow

C is designed to be a *terse programming language* (i.e. able to express *much* on a few lines or pages), so the number of *program flow keywords* is small ^{74 75}.

Program flow in a function is from top to bottom, executing statement by statement (a statement is terminated by `;` or it is a *compound statement* (a block) in braces), unless one of the following constructs is encountered.

Loop constructs are done with `while` [continuation test at the beginning of the loop], `do` [continuation test at the end of the loop] ⁷⁶ and `for` [continuation test at the beginning of the loop; additional initializer code and per-loop code ⁷⁷]. *Loops* can be left ⁷⁸ or short-cut by `break` and `continue` (and of course too by `return` and `goto`).

Alternative code flows are done by `if` and `else` ⁷⁹.

`switch/case/default` (and `break`) can be used when comparing to *compile-time constants* ^{80 81}.

`return` is used to *leave a function context*; `goto` allows *arbitrary jumps* inside *function code* ⁸².

⁷⁴ See [links](#) for a list of the C programming language keywords, with a few comments.

⁷⁵ *Keywords* cannot be *reused* as identifiers [they're truly *reserved* words].

⁷⁶ `do {...}; while (...);` is about the same as `for (;;) {...; if (!(...)) break;}`.

⁷⁷ `for`'s *per-loop code* is usually used for *incrementers*, e.g. in `for (i = 0; i < n; i++) {...;}` [being roughly ⁸³ equivalent to `i = 0; while (i < n) {...; i++;}`].

⁷⁸ However, there's no *multilevel-break* ⁸⁴.

⁷⁹ *Specially nested if/else lists* [*multi-way decisions*] are typically '*linearized*' in C code, e.g. `if (...) {...;} else if (...) {...;} else {...;}` instead of `if (...) {...;} else {if (...) {...;} else {...;}}` [all but terminal else-case do *not* establish nested blocks].

⁸⁰ The `switch` statement historically ⁸⁵ allowed a faster dispatch ⁸⁶.

⁸¹ An annoyance with the `switch` statement is the *reuse* of the `break` keyword, for which reason one cannot [by `break`] *leave an enclosing loop*.

⁸² A `goto` may not be used to jump from one function to another ⁸⁷.

⁸³ Apart from behavior if `continue` is used.

⁸⁴ Can be done by `goto`.

⁸⁵ More recent analyzers/optimizers can quite reasonably handle and take advantage of *expression constantness* in `if/else` statements [and elsewhere] too.

⁸⁶ E.g. *jump table driven* dispatch.

⁸⁷ C is versatile enough to allow *library code* ⁸⁸ being able to implement out-of-context jumps, e.g. with `long jmp`.

⁸⁸ `long jmp` and similar functions are found in [more general] library code ⁸⁹ rather than in C source code, since they e.g. modify *special* ⁹⁰ CPU registers, which can't be done in C ⁹¹.

⁸⁹ *Libraries* ⁷ can be of *mixed-language*; the *standard C library* e.g. often contains some [platform-specific] assembly code compilations ⁹².

⁹⁰ In `long jmp` most notably the *stack pointer register* ⁹³, to *adjust stack frames*.

⁹¹ But in [platform-specific] *assembly* code.

⁹² Either to implement things that can't be done in C, such as `longjmp` or `long long` number type multiplication [helper functions], or to provide performance-optimized implementations, e.g. for `memcpy`.

⁹³ In *stack-based* environments; [for very tiny environments [e.g. on small embedded systems]] a C environment can be implemented *without stack*.

Primitive data types

C defines *integer numbers* of different sizes: `char` ⁹⁴, `short` [int], `int` ⁹⁵, `long` [int], `long long` [int] ^{96 97}, which can e.g. be found to be of sizes 1, 2, 4, 4, 8 bytes ⁹⁸ [depending on *platform* and *compilation model* ⁹⁹]. Those integer types are (with exception of `char` ¹⁰⁰) implicitly signed, i.e. *operations* consider the numbers to carry a sign bit ¹⁰¹). The `unsigned` modifier changes this behavior.

C defines *floating point numbers* of different sizes: `float`, `double`, `long double` ^{96 97}, which can e.g. be found to be of sizes 4, 8, 16 bytes. Floating point number types do not support `unsigned`.

Implicit (where applicable) and *explicit* conversion rules between these number types are defined.

`void` is used to specify an un-specified reference type ¹⁰², an empty function parameter list ^{103 104} or a cast on a unused expression ¹⁰⁵.

⁹⁴ `char` is mainly used for [ASCII ¹⁰⁶] strings, i.e. [readable] *text* used in programs.

⁹⁵ `int` was assumed to fit the target CPU's native *register size*.

⁹⁶ The pattern here is to *re-apply* the `long` attribute.

⁹⁷ Not defined on every platform.

⁹⁸ C doesn't define the *exact* integer number sizes. However it defines that `long` is at least as large as `int`, which is at least as large as `short`.

⁹⁹ A compiler [environment] may allow e.g. different data pointer sizes, to allow larger or more compact applications. In a larger model ¹⁰⁷ either the full register size [instead of e.g. half] is used, or two registers are used for addressing ¹⁰⁸ or calculating.

¹⁰⁰ `char`'s signed/unsigned status is *implementation specific*. *Explicit casts* ¹⁰⁹ should be used where *sign matters*.

¹⁰¹ Usually the most significant bit, in *two's complement* representation.

¹⁰² `void*` ¹¹⁰

¹⁰³ A e.g. in `int main(void) ...`

¹⁰⁴ Leaving the function parameter list *empty*, as in `int main() ...`, as opposed to specifying it `void`, was historically used to leave the parameter list *unspecified* ¹¹¹.

¹⁰⁵ A e.g. in `(void)printf("hello, world\n");`, to indicate one is not interested in a function call's *return value* [but only in its *side-effect(s)*].

¹⁰⁶ The alternative *unicode* is supported by a larger, derived data type, `wchar_t` [often a [unsigned] `short` (size 2 bytes)].

¹⁰⁷ Such *models* can either be *mixed* [in an application or library ⁷, or different-model libraries can be used together], in which case often additional data modifier keywords ¹¹² are provided, or such memory models can't be *mixed* [and hence e.g. can only run in *different* processes].

¹⁰⁸ *Two register addressing* likely requires *segmented* memory layout.

¹⁰⁹ E.g. in `printf("%02x", (int)(unsigned char)c);`, to suppress unwanted sign extension.

¹¹⁰ Which of course cannot be dereferenced.

¹¹¹ A *unspecified* function parameter list was [historically] used to just specify the function's *return value type*, not its parameter types.

¹¹² Such *proprietary* [platform specific] keywords typically start with `_` or `__` ¹¹³.

¹¹³ Which makes them *system-internal* identifiers or keywords [by C definition ¹¹⁴].

¹¹⁴ To separate at least the *namespace* of the *C environment implementation* from the user namespace.

Composite data types

C allows *data structures* (`struct`), *overlapping* data structures (`union`), enumerated *integer constants* (`enum`) and *type name aliases* [115](#) (`typedef`). Composite data types are, together with pointers and arrays, called *derived data types*.

[Data] *definition modifiers* are `auto` [116](#) (non-static (stack) local variable), `const` (read-only variable or variable's read-only contents [117](#)) [118](#), `extern` (data *declaration* instead of *definition*), `register` (hint to the optimizer), `static` (non-stack data or non-globally visible), `volatile` (anti-hint to the optimizer).

¹¹⁵ Type aliases can be defined on primitive and composite types.

¹¹⁶ `auto` is the *default* storage class for *local variables*.

¹¹⁷ Depending on `const`'s position [119](#).

¹¹⁸ `const` pointers in *function parameters* indicate that contents will not be modified, i.e. *treated read-only* [interface contract]. `const` modifier keyword on *static data* will likely locate that data in a *read-only data segment*.

¹¹⁹ E.g. `char const* p` [120](#) vs. `char* const p` [vs. `char const* const p`] [121](#).

¹²⁰ `const char* p` is a synonym to `char const* p`, i.e. only `const`'s position relative to the indirection specifiers is relevant.

¹²¹ There is a possible `const` modifier per *indirection*.

Operators

arithmetic: `+` `-` `*` `/` `%` [addition, subtraction, multiplication, division, modulus] [122](#) [123](#)

bit: `&` `|` `^` `~` `<<` `>>` [and, or, exclusive-or, not, shifts]

logical: `&&` `||` `!` [and [124](#), or [124](#), not]

relational: `==` `!=` `<` `<=` `>` `>=` [equal, not equal, ...] [125](#)

assignment: `=` `*=` `/=` `%=` `+=` `-=` `<<=` `>>=` `&=` `|=` `^=` [126](#)

data selection: `.` `->` `[]` [data member selection, with prior referencing [127](#), array element selection]

referencing/dereferencing: `&` `*`

function call: `()`

explicit conversion: `()`

increment/decrement: `++` `--` [128](#)

size of an expression/type: `sizeof`

ternary operator: `?:`

sequence operator: `,` [129](#)

The operators have 15 levels of precedences [130](#) and left- or right-associativities.

The order of sub-expression evaluation in binary operations is not defined; exceptions are the logical operators [124](#) and the sequence operator.

¹²² The arithmetic operators are defined for number types, whereas plus and minus are defined for *pointer arithmetic* [131](#) too.

¹²³ Mathematically associative operators are not treated *computationally associative* for reasons of [intermediate] *overflow* and *rounding*.

¹²⁴ *Logical and and or* establish *sequence points*, which makes their syntax more useful [132](#).

¹²⁵ The relational operators are defined for number types and pointer types [in which case the memory positions [addresses] are compared].

¹²⁶ The *assignment expressions* have a value too [133](#) [134](#).

¹²⁷ `p->m` can be written as `(*p).m` or `p[0].m`, `->` is more convenient.

¹²⁸ C defines *pre-* and *post* increment/decrement [to allow even terser expressions].

¹²⁹ The *sequence operator* is a *convenience*; to allow several sub-statements where syntactically one statement is expected. It allows *elegant* solutions.

¹³⁰ See [links](#) for a list of the C programming language operators ordered by precedences levels.

¹³¹ $\&(p[n])$ can be written as $p + n$.

¹³² *Useful syntax* by allowing expressions such as `if (p == NULL || *p == '\0') ...`

¹³³ E.g. `a = b = c;` being `b = c;` `a = b;` [left associative].

¹³⁴ This is part of C's *orthogonality*.

Runtime library

The runtime library required to implement selfcontained programs is tiny. ¹ Standard library functions are only called explicitly [135](#) [136](#).

Even some *fundamental* functionalities used in programming, such as *input/output* [137](#) [138](#) [`stdio.h`: `fopen`, `printf`, ...], *string handling* [139](#) [`string.h`: `strcpy`, `strcmp`, ...], *conversions* [`stdlib.h`: `strtol`, `strtod`, ...], *dynamic* [140](#) *memory* [141](#) [`stdlib.h`: `malloc`, `free`, ...], are *provided in the standard library*, not in the C language [in the narrower sense] itself.

Operating system interfaces are provided syntactically the same way library functions are [142](#).

¹³⁵ This *explicitness* allows *easy* control over C code's performance.

¹³⁶ E.g. *structure assignments* may trigger `memcpy`.

¹³⁷ I.e. there is no *built-in* input/output; it's just an API [application programming interface].

¹³⁸ *Standard input and output* are known as `stdin` and `stdout` [also declared in `stdio.h`]; there's an additional *output channel* known as the standard error output, `stderr` [usually initially mapping to the same as `stdout`].

¹³⁹ The *string handling* routines [143](#) are samples of functions that may be known to a *compiler* as *intrinsic*s [inlineable functions] and expanded [directly] into [optimized] code [144](#).

¹⁴⁰ *Dynamically acquired* [at run-time only; not reserved from application start] and *dynamically sized* [data/buffer size not known at compile-time].

¹⁴¹ Dynamic memory is bound to *reference (pointer) variables*; dereferencing allows access to the *dynamic storage* [the *heap memory*].

¹⁴² Interfaces such as `open` are e.g. implemented as *stubs* that e.g. map to an interrupt/trap with an appropriate *system call* number.

¹⁴³ Together with *memory copying* functions, such as `memcpy`.

¹⁴⁴ E.g. `strlen("hello, world\n")` could be translated into the *compile-time constant* 13 [i.e. function call left away].

Links

[Frequently asked C/C++ questions \[on lrdev\]](#)

Provides answers to ISO/IEC 9899 availability, useful entry points for C programming language information, C programming language keywords, etc.

[C programming language coding guidelines](#)

C's do's and don't's.

[C/C++ operator precedences](#)

Precedences of the C/C++ programming languages operators.