

Introduction to Compilers and Language Design

Copyright © 2018 Douglas Thain.

Hardcover ISBN: 978-0-359-13804-3

Paperback ISBN: 978-0-359-14283-5

First edition.

Anyone is free to download and print the PDF edition of this book for personal use. Commercial distribution, printing, or reproduction without the author's consent is expressly prohibited. All other rights are reserved.

You can find the latest version of the PDF edition, and purchase inexpensive hardcover copies at <http://compilerbook.org>

Revision Date: May 2, 2019

Appendix C – Coding Conventions

C has been the language of choice for implementing low level systems like compilers, operating systems, and drivers since the 1980s. However, it is fair to say that C does not enforce a wide variety of good programming practices, in comparison to other languages. To write solid C, you need to exercise a high degree of self-discipline.¹

For many students, a compilers course in college is the first place where you are asked to create a good sized piece of software, refining it through several development cycles until the final project is reached. This is a good opportunity for you to pick up some good habits that will make you more productive.

To that end, here are the coding conventions that I ask my students to observe when writing C code. Each of these recommendations requires a little more work up front, but will save you headaches in the long run.

Use a version control system. There are a variety of nice open source systems for keeping track of your source code. Today, Git, Mercurial, and Subversion are quite popular, and I’m sure next year will bring a new one. Pick one, learn the basic features, and shape your code gradually by making small commits.²

Go from working to working. Never leave your code in a broken state. Begin by checking in the simplest possible sketch of your program that compiles and works, even if it only prints out “hello world”. Then, add to the program in a minor way, make sure that it compiles and runs, and check it in.³

Eliminate dead code. Students often pick up the habit of commenting out one bit of code while they attempt to change it and test it. While this is

¹Why not use C++ to address some of these disciplines? Although C++ is a common part of many computer science curricula, I generally discourage the use of C++ by students. Although it has many features that are attractive at first glance, they are not powerful enough to allow you to dispense with the basic C mechanisms. (For example, even if you use the C++ string class, you still need to understand basic character arrays and pointers.) Further, the language is so complex that very few people really understand the complete set of features and how they interact. If you stick with C, what you see is what you get.

²Some people like to spend endless hours arguing about the proper way to use arcane features of these tools. Don’t be one of those people: learn the basic operations and spend your mental energy on your code instead.

³This advice is often attributed as one of Jim Gray’s “Laws of Data Engineering” in slide presentations, but I haven’t been able to find an authoritative reference.

a reasonable tactic to use for a quick test, don't allow this dead code to pile up in your program, otherwise your source code will quickly become incomprehensible. Remove unused code, data, comments, files, or anything else that is unnecessary to the program, so that you can clearly see what it does now. Trust your version control system to allow you to go back to a previously working version, if needed.

Use tools to handle indenting. Don't waste your time arguing about indenting style; find a tool that does it for you automatically, and then forget about it. Your editor probably has a mode to indent automatically. If not, use the standard Unix tool `indent`.

Name things consistently. In this book, you will see that every function consists of a noun and a verb: `expr_typecheck`, `decl_codegen`, etc. Each one is used consistently: `expr` is always used for expressions, `codegen` is always used for code generation. Every function dealing with expressions is in the `expr` module. It may be tempting to take shortcuts or make abbreviations in the heat of battle, but this will come back to bite you. Do it right the first time.

Put only the interface in a header file. In C, a header file (like `expr.h`) is used to describe the elements needed to call a function: function prototypes and the types and constants necessary to invoke those functions. If a function is only used within one module, it should *not* be mentioned in the header file, because nobody outside the module needs that information.

Put only the implementation in a source file. In C, a source file (like `expr.c`) is used to provide the definitions of functions. In the source file, you should include the corresponding header (`expr.h`) so that the compiler can check that your function definitions match the prototypes. Any function or variable that is private to the module should be declared `static`.

Be lazy and recursive. Many language data structures are hierarchically nested. When designing an algorithm, take note of the nested data structures, and pass responsibility to other functions, even if you haven't written them yet. This technique generally results in code that is simple, compact, and readable. For example, to print out a variable declaration, break it down into printing the name, then the type, then the value, with some punctuation in between:

```
printf("%s:\n", d->name);
type_print(d->type);
printf(" = ");
expr_print(d->value);
printf(" ;\n");
```

Then proceed to writing `type_print` and `expr_print`, if you haven't done them already.

Use a Makefile to build everything automatically. Learn how to write a Makefile, if you haven't already. The basic syntax of Make is very simple.

The following rule says that `expr.o` depends upon `expr.c` and `expr.h`, and can be built by running the command `gcc`:

```
expr.o: expr.c expr.h
    gcc expr.c -c -o expr.o -Wall
```

There are many variations of Make that include wildcards, and pattern substitution, and all manner of other things that can be confusing to the uninitiated. Just start by writing plain old rules whose meaning is clear.

Null pointers are friends. When designing a data structure, use null pointers to indicate when nothing is present. You cannot dereference a null pointer, of course, and so you must check before using it. This can lead to code cluttered with null checks everywhere, like this:

```
void expr_codegen( struct expr *e, FILE *output )
{
    if(e->left) expr_codegen(e->left, output);
    if(e->right) expr_codegen(e->right, output);

    . . .
}
```

You can eliminate many of them by simply placing the check at the beginning of the function, and programming in a recursive style:

```
void expr_codegen( struct expr *e, FILE *output )
{
    if(!e) return;

    expr_codegen(e->left, output);
    expr_codegen(e->right, output);

    . . .
}
```

Automate regression testing. A compiler has to handle a large number of details, and it is all too easy for you to accidentally introduce a new bug when attempting to fix an old one. To handle this, create a simple test suite that consists of a set of sample programs, some correct and some incorrect. Write a little script that invokes your compiler on each sample program, and makes sure that it succeeds for the good tests, and fails on the bad tests. Make it a part of your Makefile, so that every time you touch the code, the tests are run, and you will know if things are still working.

Index

- a.out, 144
- absolute address, 141
- abstract syntax tree, 73
- abstract syntax tree (AST), 7, 8, 83
- accepting states, 16
- accepts, 16
- Acorn Archimedes, 165
- Acorn RISC Machine, 165
- address computation, 141
- Advanced RISC Machine (ARM), 165
- alternation, 14
- ambiguous grammar, 37
- ARM (Advanced RISC Machine), 165
- assembler, 6
- associativity, 15
- AST (abstract syntax tree), 7, 8, 83
- atomic types, 101

- backtracking, 13
- base pointer, 138
- base-relative, 153
- base-relative address, 141
- basic block, 123
- binary blob, 144
- bottom-up derivation, 37
- break, 134
- BSS size, 145
- bytecode, 1

- callee saved, 158
- caller saved, 158
- calling convention, 139

- canonical collection, 50
- CFG (context-free grammar), 36
- Chomsky hierarchy, 62
- CISC (Complex Instruction Set Computer), 165
- closure, 50
- code generator, 7
- code hoisting, 199
- code segment, 134
- comments, 11
- commutativity, 15
- compact finite state machine, 50
- compiler, 1, 5
- complex, 153
- Complex Instruction Set Computer (CISC), 165
- compound types, 102
- concatenation, 14
- conditional execution, 171
- conflict graph, 207
- constant folding, 122, 196
- context free languages, 62
- context sensitive languages, 63
- context-free grammar (CFG), 36
- control flow graph, 123, 200
- core, 61
- crystal ball interpretation, 18

- DAG (directed acyclic graph), 118
- data segment, 134
- data size, 145
- declaration, 83
- delete, 136
- derivation, 36
- deterministic finite automaton (DFA), 16

- directed acyclic graph (DAG), 118
- directives, 149
- distribution, 15
- domain specific languages, 2
- dot, 50
- dynamically typed language, 99

- entry point, 145
- enumerations, 102
- epsilon closure, 22
- evaluated, 83
- executable formats, 144
- explicitly typed language, 100
- expression, 83
- Extensible Linking Format (ELF), 145
- external format, 117

- FA (finite automaton), 15
- finite automata, 13
- finite automaton (FA), 15
- frame pointer, 138
- free, 136
- function inlining, 199

- GIMPLE (GNU Simple Representation), 128
- Global data, 141
- global value, 152
- GNU Simple Representation (GIMPLE), 128
- grammar, 7, 8
- graph coloring, 208
- guard page, 135

- Heap data, 142
- heap segment, 134

- idempotency, 15
- identifiers, 11
- immediate value, 152
- implicitly typed language, 100
- indirect value, 153
- instruction selection, 7
- instructions, 149

- intermediate representation (IR), 7, 117
- interpreter, 1, 67
- IR (intermediate representation), 117
- items, 50

- Java Virtual Machine (JVM), 130
- JIT, 1
- just in time compiling, 1
- JVM (Java Virtual Machine), 130

- kernel, 50, 194
- keywords, 11
- Kleene closure, 14

- labels, 149
- LALR (Lookahead LR), 61
- language, 36
- leaf function, 141, 160, 173
- left recursion, 40
- lifetime, 126
- linker, 6
- literal pool, 168
- little languages, 2
- live ranges, 207
- LL(1) parse table, 46
- Local data, 142
- logical segments, 133
- lookahead, 59
- Lookahead LR (LALR), 61
- loop unrolling, 197
- LR(0) automaton, 50

- magic number, 145
- malloc, 136
- many-worlds interpretation, 18
- memory fragmentation, 137

- name resolution, 97, 109
- new, 136
- NFA (nondeterministic finite automaton), 17
- non-terminal, 36
- nondeterministic finite automaton (NFA), 17

- numbers, 11
- object code, 6
- optimization, global, 193
- optimization, interprocedural, 193
- optimization, local, 193
- optimization, peephole, 202
- optimizers, 7
- page fault, 135
- parser, 7
- parser generator, 67
- partial execution, 122
- PC-relative address, 142
- preprocessor, 5
- record type, 102
- recursive descent parser, 44
- recursively enumerable languages, 63
- reduce, 49
- reduce-reduce conflict, 53
- Reduced Instruction Set Computer (RISC), 165
- redundant load elimination, 202
- register allocation, 7, 206
- register value, 152
- regular expression, 14
- regular expressions, 13
- regular languages, 62
- rejects, 16
- RISC (Reduced Instruction Set Computer), 165
- rules, 36
- runtime library, 189
- safe programming language, 99
- scanner, 6
- scanner generator, 27
- scope, 106
- scratch registers, 179
- section table, 145
- segmentation fault, 135
- semantic actions, 72
- semantic routines, 7, 8
- semantic type, 77
- semantic values, 72
- semantics, 97
- sentence, 36
- sentential form, 36
- sentinel, 11
- shift, 49
- shift-reduce, 49
- shift-reduce conflict, 53
- side effect, 83
- side effects, 185
- Simple LR (SLR), 54
- SLR (Simple LR), 54
- SLR grammar, 54
- SLR parse tables, 54
- source language, 1
- SSA (static single assignment), 125
- stack, 138
- stack backtrace, 175
- stack frame, 138, 161
- stack machine, 127
- stack pointer, 138, 157, 171
- stack segment, 134
- start symbol, 36
- statement, 83
- static single assignment (SSA), 125
- statically typed language, 99
- strength reduction, 197
- strings, 11
- structure tag, 143
- structure type, 102
- subset construction, 22
- symbol table, 97, 105, 145
- System V ABI, 158
- target language, 1
- terminal, 35
- text segment, 134
- text size, 145
- tokens, 6, 7, 11
- toolchain, 5
- top-down derivation, 37
- translator, 67
- tree coverage, 202
- type, 98
- type checking, 97

typechecking, 8

union types, 102

unrolling factor, 197

unsafe programming language, 98

user-defined types, 101

validator, 67, 71

value, 83, 185

value-number method, 121

variant type, 103

virtual machine, 1

virtual registers, 126

virtual stack machine, 127

weak equivalence, 37

whitespace, 11

YYSTYPE, 77