

acwj

Trganda

May 30, 2021

## 1 Part 0: Introduction

I've decided to go on a compiler writing journey. In the past I've written some [assemblers](#), and I've written a [simaple compiler](#) for a typeless language. But I've never written a compiler that can compile itself. So that's where I'm headed on this journey.

As part of the process, I'm going to write up my work so that others can follow along. This will also help me to clarify my thoughts and ideas. Hopefully you, and I, will find this useful!

### 1.1 Goals of the journey

Here are my goals, and non-goals, for the journey:

- To write a self-compiling compiler. I think that if the compiler can compile itself, it gets to call itself a real compiler.
- To target at least one real hardware platform. I've seen a few compilers that generate code for hypothetical machines. I want my compiler to work on real hardware. Also, if possible, I want to write the compiler so that it can support multiple backends for different hardware platforms.
- Practical before research. There's a whole lot of research in the area of compilers. I want to start from absolute zero on this journey, so I'll tend to go for a practical approach and not a theory-heavy approach. That said, there will be times when I'll need to introduce (and implement) some theory-based stuff.
- Follow the KISS principle: keep it simple, stupid! I'm definitely going to be using Ken Thompson's principle here: "When in doubt, use brute force."
- Take a lot of small steps to reach the final goal. I'll break the journey up into a lot of simple steps instead of taking large leaps. This will make each new addition to the compiler a bite-sized and easily digestible thing.

### 1.2 Target Language

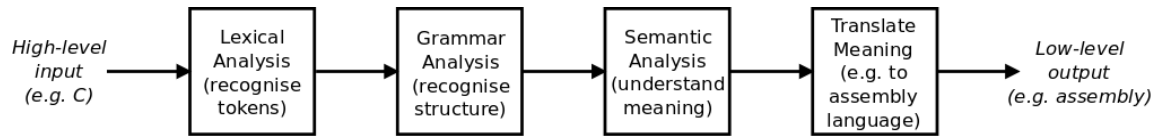
The choice of a target language is difficult. If I choose a high-level language like **Python**, **Go** etc., then I'll have to implement a whole pile of libraries and classes as they are built-in to the language.

I could write a compiler for a language like Lisp, but these can be [done easily](#).

Instead, I've fallen back on the old standby and I'm going to write a compiler for a subset of C, enough to allow the compiler to compile itself. C is just a step up from assembly language (for some subset of C, not [C18](#)), and this will help make the task of compiling the C code down to assembly somewhat easier. Oh, and I also like C.

## 1.3 The Basic of a Compiler's Job

The job of a compiler is to translate input in one language (usually a high-level language) into a different output language (usually a lower-level language than the input). The main steps are:



- Do [lexical analysis](#) to recognise the lexical elements. In several languages, `=` is different to `==`, so you can't just read a single `=`. We call these lexical elements *tokens*.
- [Parse](#) the input, i.e. recognise the syntax and structural elements of the input and ensure that they conform to the *grammar* of the language. For example, your language might have this decision-making structure:

---

```
if (x < 23) {  
    print("x is smaller than 23\n");  
}
```

---

but in another language you might write:

---

```
if (x < 23)  
    print("x is smaller than 23\n");
```

---

This is also the place where the compiler can detect syntax errors, like if the semicolon was missing on the end of first *print* statement.

- Do [semantic analysis](#) of the input, i.e. understand the meaning of the input. This is actually different from recognising the syntax and structure. For example, in English, a sentence might have the form `<subject> <verb> <adjective> <object>`.

---

```
David ate lovely bananas.  
Jennifer hates green tomatoes.
```

---

- [Translate](#) the meaning of the input into a different language. Here we convert the input, parts at a time, into a lower-level language.

## 1.4 Resources

There's a lot of compiler resources out on the Internet. Here are the ones I'll be looking at.

### 1.4.1 Learning Resources

If you want to start with some books, papers and tools on compilers, I'd highly recommend this list:

- [Curated list of awesome resources on Compilers, Interpreters and Runtimes](#) by Ahmad Alhour

### 1.4.2 Existing Compilers

While I'm going to build my own compiler, I plan on looking at other compilers for ideas and probably also borrow some of their code. Here are the ones I'm looking at:

- [SubC](#) by Nils M Holm
- [Swieros C Compiler](#) by Robert Swierczek
- [fbcc](#) by Fabrice Bellard
- [tcc](#), also by Fabrice Bellard and others
- [catc](#) by Yuichiro Nakada
- [amacc](#) by Jim Huang
- [Small C](#) by Ron Cain, James E. Hendrix, derivatives by others

In particular, I'll be using a lot of the ideas, and some of the code, from the SubC compiler.

## 1.5 Setting Up the Development Enviroment

Assuming that you want to come along on this journey, here's what you'll need. I'm going to use a Linux development enviroment, so download and set up your favourite Linux system: I'm using Lubuntu 18.04.

I'm going to target two hardware platforms: Intel x86-64 and 32-bit ARM. I'll use a PC running Lubuntu 18.04 as the Intel target, and a Raspberry Pi running Raspbian as the ARM target.

On the Intel platform, we are going to need an existing C compiler. So, install this package. If there are any more tools required for a vanilla Linux system let me know.

Finally, clone a copy of this Github repository.

## 1.6 The Next Step

In the next part of our compiler writing journey, we will start with the code to scan our input file and find the *tokens* that are the lexical elements of our language.

## 2 Part 1: Introduction to Lexical Scanning

We start our compiler writing journey with a simple lexical scanner. As I mentioned in the previous part, the job of the scanner is to identify the lexical elements, or *tokens*, in the input language.

We will start with a language that has only five lexical elements:

- the four basic maths operators: +, /, \* and -
- decimal whole numbers which have 1 or more digits 0 .. 9

Each token that we scan is going to be stored in this structure (from **defs.h**)

---

```
// Token structure
struct token {
    int token;
    int intvalue;
};
```

---

where the **token** field can be one of these values (from **defs.h**):

---

```
// Tokens
enum {
    T_PLUS, T_MINUS, T_STAR, T_SLASH, T_INTLIT
};
```

---

When the token is a *T\_INTLIT* (i.e. an integer literal), the *intvalue* field will hold the value of the integer that we scanned in.

## 2.1 Functions in *scan.c*

The *scan.c* file holds the functions of our lexical scanner. We are going to read in one character at a time from our input file. However, there will be times when we need to "put back" a character if we have read too far ahead in the input stream. We also want to track what line we are currently on so that we can print the line number in our debug messages. All of this is done by the *next()* function:

---

```
static int next(void) {
    int c;

    if (Putback) {           // Use the character put
        c = Putback;         // back if there is one
        Putback = 0;
        return c;
    }

    c = fgetc(Infile);       // Read from input file
    if ('\n' == c)           // Increment line count
        Line++;
    return c;
}
```

---

The *Putback* and *Line* variables are defined in *data.h* along with our input file pointer:

---

```
extern_ int    Line;
extern_ int    Putback;
extern_ FILE   *Infile;
```

---

All C files will include this where *extern\_* is replaced with *extern*. But *main.c* will remove the *extern\_*; hence, these variables will "belong" to *main.c*.

Finally, how do we put a character back into the input stream? Thus:

---

```
// Put back an unwanted character
static void putback(int c) {
    Putback = c;
}
```

---

## 2.2 Ignoring Whitespace

We need a function that reads and silently skips whitespace characters until it gets a non-whitespace character, and returns it. Thus:

---

```
// Skip past input that we don't need to deal with,
// i.e. whitespace, newlines. Return the first
// character we do need to deal with.
```

---

```

static int skip(void) {
    int c;

    c = next();
    while ( ' ' == c || '\t' == c || '\n' == c || '\r' == c || '\f' == c ) {
        c = next();
    }
    return (c);
}

```

---

## 2.3 Scanning Tokens: *scan()*

So now we can read characters in while skipping whitespace; we can also put back a character if we read one character too far ahead. We can now write our first lexical scanner:

---

```

// Scan and return the next token found in the input.
// Return 1 if token valid, 0 if no tokens left.
int scan(struct token *t) {
    int c;

    // Skip whitespace
    c = skip();

    // Determine the token based on
    // the input character
    switch (c) {
    case EOF:
        return (0);
    case '+':
        t->token = T_PLUS;
        break;
    case '-':
        t->token = T_MINUS;
        break;
    case '*':
        t->token = T_STAR;
        break;
    case '/':
        t->token = T_SLASH;
        break;
    default:
        // More here soon
    }

    // We found a token
    return (1);
}

```

---

That's it for the simple one-character tokens: for each recognised character, turn it into a token. You may ask: why not just put the recognised character into the *structtoken*? The answer is that later we will need to recognise multi-character tokens such as `==` and keywords like *if* and *while*. So it will make life easier to have an enumerated list of token values.

## 2.4 Integer Literal Values

In fact, we already have to face this situation as we also need to recognise integer literal values like 3827 and 87731. Here is the missing *default* code from the *switch* statement:

---

```
default:
    // If it's a digit, scan the
    // literal integer value in
    if (isdigit(c)) {
        t->intvalue = scanint(c);
        t->token = T_INTLIT;
        break;
    }

    printf("Unrecognised character %c on line %d\n", c, Line);
    exit(1);
```

---

Once we hit a decimal digit character, we call the helper function *scanint()* with this first character. It will return the scanned integer value. To do this, it has to read each character in turn, check that it's a legitimate digit, and build up the final number. Here is the code:

---

```
// Scan and return an integer literal
// value from the input file. Store
// the value as a string in Text.
static int scanint(int c) {
    int k, val = 0;

    // Convert each character into an int value
    while ((k = chrpos("0123456789", c)) >= 0) {
        val = val * 10 + k;
        c = next();
    }

    // We hit a non-integer character, put it back.
    putback(c);
    return val;
}
```

---

We start with a zero *val* value. Each time we get a character in the set 0 to 9, we convert this to an *int* value with *chrpos()*. We make *val* 10 times bigger and then add this new digit to it.

For example, if we have the character 3, 2, 8, we do:

- $val = 0 * 10 + 3$ , i.e. 3
- $val = 3 * 10 + 2$ , i.e. 32
- $val = 32 * 10 + 8$ , i.e. 328

Right at the end, did you notice the call to *putback(c)*? We found a character that's not a decimal digit at this point. We can't simply discard it, but luckily we can put it back in the input stream to be consumed later.

You may also ask at this point: why not simply subtract the ASCII value of 0 from *c* to make it an integer? The answer is that, later on, we will be able to do *chrpos("0123456789abcdef")* to convert hexadecimal digits as well.

Here's the code for *chrpos()*:

---

```
// Return the position of character c
// in string s, or -1 if c not found
static int chrpos(char *s, int c) {
    char *p;

    p = strchr(s, c);
    return (p ? p - s : -1);
}
```

---

And that's it for the lexical scanner code in *scan.c* for now.

## 2.5 Putting the Scanner to Work

The code in *main.c* puts the above scanner to work. The *main()* function opens up a file and then scans it for tokens:

---

```
void main(int argc, char *argv[]) {
    ...
    init();
    ...
    Infile = fopen(argv[1], "r");
    ...
    scanfile();
    exit(0);
}
```

---

And *scanfile()* loops while there is new token and prints out the details of the token:

---

```
// List of printable tokens
char *tokstr[] = { "+", "-", "*", "/", "intlitt" };

// Loop scanning in all the tokens in the input file.
// Print out details of each token found.
static void scanfile() {
    struct token T;

    while (scan(&T)) {
        printf("Token %s", tokstr[T.token]);
        if (T.token == T_INTLIT)
            printf(", value %d", T.intvalue);
        printf("\n");
    }
}
```

---

## 2.6 Some Example Input Files

I've provided some example input files so you can see what tokens the scanner finds in each file, and what input files the scanner rejects.

---

```
$ make
cc -o scanner -g main.c scan.c

$ cat input01
2 + 3 * 5 - 8 / 3
```

```
$ ./scanner input01
Token intlit, value 2
Token +
Token intlit, value 3
Token *
Token intlit, value 5
Token -
Token intlit, value 8
Token /
Token intlit, value 3

$ cat input04
23 +
18 -
45.6 * 2
/ 18

$ ./scanner input04
Token intlit, value 23
Token +
Token intlit, value 18
Token -
Token intlit, value 45
Unrecognised character . on line 3
```

---

## 2.7 Conclusion and What's Next

We've started small and we have a simple lexical scanner that recognises the four main maths operators and also integer values. We saw that we needed to skip whitespace and put back characters if we read too far into the input.

Single character tokens are easy to scan, but multi-character tokens are a bit harder. But at the end, the *scan()* function returns the next token from the input file in a *structtoken* variable.

In the next part of our compiler writing journey, we will build a recursive descent parser to interpret the grammar of our input files, and calculate & print out the final value for each file.

## 3 Part 2: Introduction to Parsing

In this part of our compiler writing journey, I'm going to introduce the basics of a parser. As I mentioned in the first part, the job of the parser is to recognise the syntax and structural elements of the input and ensure that they conform to the *grammar* of the language.

We already have several language elements that we can scan in, i.e. our tokens:

- the four basic maths operators: \*, /, + and -
- decimal whole numbers which have 1 or more digits 0 .. 9

Now let's define a grammar for the language that our parser will recognise.

### 3.1 BNF: Backus-Naur Form

You will come across the use of [BNF](#) at some point if you get into dealing with computer languages. I will just introduce enough of the BNF syntax here to express the grammar we want to recognise.



We want a grammar to express maths expressions with whole numbers. Here is the BNF description of the grammar:

---

```
expression: number
| expression '*' expression
| expression '/' expression
| expression '+' expression
| expression '-' expression
;

number: T_INTLIT
;
```

---

The vertical bars separate options in the grammar, so the above says:

- An expression could be just a number, or
- An expression is two expressions separated by a  $*$  token, or
- An expression is two expressions separated by a  $/$  token, or
- An expression is two expressions separated by a  $+$  token, or
- An expression is two expressions separated by a  $-$  token, or
- A number is always a  $T\_INTLIT$  token

It should be pretty obvious that the BNF definition of the grammar is *recursive*: an expression is defined by referencing other expressions. But there is a way to **bottom-out** the recursion: when an expression turns out to be a number, this is always a  $T\_INTLIT$  token and thus not recursive.

In BNF, we say that "expression" and "number" are *non-terminal* symbols, as they are produced by rules in the grammar. However,  $T\_INTLIT$  is a *terminal* symbol as it is not defined by any rule. Instead, it is an already-recognised token in the language. Similarly, the four maths operator tokens are terminal symbols.

## 3.2 Recursive Descent Parsing

Given that the grammar for our language is recursive, it makes sense for us to try and parse it recursively. What we need to do is to read in a token, then *lookahead* to the next token. Based on what the next token is, we can then decide what path we need to take to parse the input. This may require us to recursively call a function that has already been called.

In our case, the first token in any expression will be a number and this may be followed by maths operator. After that there may only be a single number, or there may be the start of a whole new expression. How can we parse this recursively?

We can write pseudo-code that look like this:

---

```
function expression() {
    Scan and check the first token is a number. Error if it's not
    Get the next token
    If we have reached the end of the input, return, i.e. base case

    Otherwise, call expression()
}
```

---

Let's run this function on the input  $2+3-5T\_EOF$  where  $T\_EOF$  is a token that reflects the end of the input. I will number each call to *expression()*.