



**How To
Create
Your Own
Freaking
Awesome
Programming
Language**

TABLE OF CONTENT

1. [Table of Content](#)
2. [Introduction](#)
 1. [Summary](#)
 2. [About The Author](#)
 3. [Before We Begin](#)
3. [Overview](#)
 1. [The Four Parts of a Language](#)
 2. [Meet Awesome: Our Toy Language](#)
4. [Lexer](#)
 1. [Lex \(Flex\)](#)
 2. [Ragel](#)
 3. [Python Style Indentation For Awesome](#)
 4. [Do It Yourself I](#)
5. [Parser](#)
 1. [Bison \(Yacc\)](#)
 2. [Lemon](#)
 3. [ANTLR](#)
 4. [PEGs](#)
 5. [Operator Precedence](#)
 6. [Connecting The Lexer and Parser in Awesome](#)
 7. [Do It Yourself II](#)
6. [Runtime Model](#)
 1. [Procedural](#)
 2. [Class-based](#)
 3. [Prototype-based](#)
 4. [Functional](#)
 5. [Our Awesome Runtime](#)
 6. [Do It Yourself III](#)

7. [Interpreter](#)
 1. [Evaluating The Nodes in Awesome](#)
 2. [Do It Yourself IV](#)
8. [Virtual Machine](#)
 1. [Byte-code](#)
 2. [The Stack](#)
 3. [Prototyping a VM in Ruby](#)
9. [Compilation](#)
 1. [Compiling to Byte-code](#)
 2. [Compiling to Machine Code](#)
10. [Mio, a minimalist homoiconic language](#)
 1. [Homoicowhat?](#)
 2. [Messages all the way down](#)
 3. [The Runtime](#)
 4. [Implementing Mio in Mio](#)
 5. [But it's ugly](#)
11. [Going Further](#)
 1. [Homoiconicity](#)
 2. [Self-Hosting](#)
 3. [What's Missing?](#)
12. [Resources](#)
 1. [Books & Papers](#)
 2. [Events](#)
 3. [Forums and Blogs](#)
 4. [Classes](#)
 5. [Interesting Languages](#)
13. [Farewell!](#)
14. [Solutions to Do It Yourself](#)
 1. [Solutions to Do It Yourself I](#)
 2. [Solutions to Do It Yourself II](#)

3. [Solutions to Do It Yourself III](#)
4. [Solutions to Do It Yourself IV](#)

Revision #5, Published June 2013.

Cover background image © [Asja Boros](#)

Content of this book is © Marc-André Cournoyer. All right reserved. This eBook copy is for a single user. You may not share it in any way unless you have written permission of the author.

This is a sample chapter.
Buy the full book online at
<http://createyourproglang.com>

LEXER

The lexer, or scanner, or tokenizer is the part of a language that converts the input, the code you want to execute, into tokens the parser can understand. If you think of your program as a series of sentences, tokens would be the words in those sentences. The job of the lexer is to extract those words (tokens) and tag them with a type (is it a string, number, operator?).

Let's say you have the following code:

```
1 print(  
2     "I ate",  
3     3,  
4     pies  
5     )
```

Once this code goes through the lexer, it will look something like this:

```
1 [IDENTIFIER print] ["("]  
2           [STRING "I ate"] [","]  
3           [NUMBER 3] [","]  
4           [IDENTIFIER pies]  
5           [")"]
```

What the lexer does is split the code into atomic units (tokens) and tag each one with the type of token it contains. This job can be done by some parsers, as we'll see in the next chapter, but separating it into two distinct processes makes it simpler for us developers and easier to understand.

Lexers can be implemented using regular expressions, but more appropriate tools exist.

Each of these tools take a grammar that will be compiled into the actual lexer. The format of these grammars are all alike. Regular expressions on the left hand side are repeatedly matched, in order, against the next portion of the input code string. When a match is found, the action on the right is taken.

LEX (FLEX)

[Flex](#) is a modern version of [Lex](#) (that was coded by Eric Schmidt, Ex-CEO of Google, by the way) for generating C lexers. Along with Yacc, Lex is the most commonly used lexer for parsing.

It has been ported to several target languages.

- [Rexical for Ruby](#)
- [JFlex for Java](#)

Lex and friends are not lexers per se. They are lexer compilers. You supply it a grammar and it will output a lexer. Here's what that grammar looks like:

```
1 %%
2 // Whitespace
3 [ \t\n]+      /* ignore */
4
5 // Literals
6 [0-9]+       yylval = atoi(yytext); return T_NUMBER;
7
8 // Keywords
9 "end"        yylval = yytext; return T_END;
10 // ...
```

On the left side a regular expression defines how the token is matched. On the right side, the action to take. The value of the token is stored in `yylval` and the type of token is returned. The `yy` prefix in the variable names is an heritage from Yacc, a parser compiler which we'll talk about in the next chapter.

More details in the [Flex manual](#).

A Ruby equivalent, using the `rexical` gem (a port of Lex to Ruby), would be:

```
1 rule
2   # Whitespace
3   [\ \t]+      # ignore
4
5   # Literals
6   [0-9]+      { [:NUMBER, text.to_i] }
7
8   # Keywords
9   end        { [:END, text] }
```

`Rexical` follows a similar grammar as `Lex`. Regular expression on the left and action on the right. However, an array of two items is used to return the type and value of the matched token.

More details on the [rexical project page](#).

RAGEL

A powerful tool for creating a scanner is [Ragel](#). It is very flexible, and can handle grammars of varying complexities and output lexers in several languages.

Here's what a `Ragel` grammar looks like:

```
1 %%{
2 machine lexer;
3
4 # Machine
5 number      = [0-9]+;
6 whitespace  = " ";
7 keyword     = "end" | "def" | "class" | "if" | "else" | "true" | "false" | "nil";
8
9 # Actions
```

```

10 main := |*
11   whitespace; # ignore
12   number      => { tokens << [:NUMBER, data[ts..te].to_i] };
13   keyword     => { tokens << [data[ts...te].upcase.to_sym, data[ts...te]] };
14 *|;
15
16 class Lexer
17   def initialize
18     %% write data;
19   end
20
21   def run(data)
22     eof = data.size
23     line = 1
24     tokens = []
25     %% write init;
26     %% write exec;
27     tokens
28   end
29 end
30 }%%

```

More details in the [Ragel manual \(PDF\)](#).

Here are a few real-world examples of Ragel grammars used as language lexers:

- [Min's lexer](#) (in Java)
- [Potion's lexer](#) (in C)

PYTHON STYLE INDENTATION FOR AWESOME

If you intend to build a fully-functioning language, you should use one of the previously mentioned tools. Since Awesome is a simplistic language and we want to illustrate the basic concepts of a scanner, we will build the lexer from scratch using regular expressions.

To make things more interesting, we'll use indentation to delimit blocks in our toy language, as in Python. All of indentation magic takes place within the lexer. Parsing blocks of code delimited with { ... } is no different from parsing indentation when you know how to do it.

Tokenizing the following Python code:

```
1 if tasty == True:
2     print "Delicious!"
```

will yield these tokens:

```
1 [IF] [IDENTIFIER tasty] [EQUAL] [IDENTIFIER True]
2   [INDENT] [IDENTIFIER print] [STRING "Delicious!"]
3 [DEDENT]
```

The block is wrapped in `INDENT` and `DEDENT` tokens instead of { and }.

The indentation-parsing algorithm is simple. You need to track two things: the current indentation level and the stack of indentation levels. When you encounter a line break followed by spaces, you update the indentation level. Here's our lexer for the Awesome language:

In file `code/lexer.rb`

Our lexer will be used like so: `Lexer.new.tokenize("code")`, and will return an array of tokens (a token being a tuple of `[TOKEN_TYPE, TOKEN_VALUE]`).

```
3 class Lexer
```

First we define the special keywords of our language in a constant. It will be used later on in the tokenizing process to disambiguate an identifier (method name, local variable, etc.) from a keyword.

```
7  KEYWORDS = ["def", "class", "if", "true", "false", "nil"]
8
9  def tokenize(code)
10     code.chomp! # Remove extra line breaks
11     tokens = [] # This will hold the generated tokens
12
```

We need to know how deep we are in the indentation so we keep track of the current indentation level we are in, and previous ones in the stack so that when we dedent, we can check if we're on the correct level.

```
16     current_indent = 0 # number of spaces in the last indent
17     indent_stack = []
18
```

Here is how to implement a very simple scanner. Advance one character at the time until you find something to parse. We'll use regular expressions to scan from the current position (*i*) up to the end of the code.

```
23     i = 0 # Current character position
24     while i < code.size
25         chunk = code[i..-1]
26
```

Each of the following `if/elsif`s will test the current code chunk with a regular expression. The order is important as we want to match `if` as a keyword, and not a method name, we'll need to apply it first.

First, we'll scan for names: method names and variable names, which we'll call identifiers. Also scanning for special reserved keywords such as `if`, `def` and `true`.

```
34     if identifier = chunk[/\A([a-z]\w*)/, 1]
35         if KEYWORDS.include?(identifier) # keywords will generate [:IF, "if"]
36             tokens << [identifier.upcase.to_sym, identifier]
37         else
38             tokens << [:IDENTIFIER, identifier]
39         end
40         i += identifier.size # skip what we just parsed
41
```

Now scanning for constants, names starting with a capital letter. Which means, class names are constants in our language.

```
44     elsif constant = chunk[/\A([A-Z]\w*)/, 1]
45         tokens << [:CONSTANT, constant]
46         i += constant.size
47
```

Next, matching numbers. Our language will only support integers. But to add support for floats, you'd simply need to add a similar rule and adapt the regular expression accordingly.

```
50     elsif number = chunk[/\A([0-9]+)/, 1]
51         tokens << [:NUMBER, number.to_i]
52         i += number.size
53
```

Of course, matching strings too. Anything between `"..."`.

```
55     elsif string = chunk[/\A("[^"]*)"/, 1]
56         tokens << [:STRING, string]
57         i += string.size + 2 # skip two more to exclude the `"`
58
```

And here's the indentation magic! We have to take care of 3 cases:

```
if true: # 1) The block is created.
    line 1
    line 2 # 2) New line inside a block, at the same level.
continue # 3) Dedent.
```

This `elsif` takes care of the first case. The number of spaces will determine the indent level.

```
68     elsif indent = chunk[/\A\n( +)/m, 1] # Matches "<newline> <spaces>"
69         if indent.size <= current_indent # indent should go up when creating a block
70             raise "Bad indent level, got #{indent.size} indents, " +
71                 "expected > #{current_indent}"
72         end
73         current_indent = indent.size
74         indent_stack.push(current_indent)
75         tokens << [:INDENT, indent.size]
76         i += indent.size + 2
77
```

The next `elsif` takes care of the two last cases:

- Case 2: We stay in the same block if the indent level (number of spaces) is the same as `current_indent`.
- Case 3: Close the current block, if indent level is lower than `current_indent`.

```
83     elsif indent = chunk[/\A\n( *) /m, 1] # Matches "<newline> <spaces>"
84         if indent.size == current_indent # Case 2
85             tokens << [:NEWLINE, "\n"] # Nothing to do, we're still in the same block
86         elsif indent.size < current_indent # Case 3
87             while indent.size < current_indent
88                 indent_stack.pop
89                 current_indent = indent_stack.last || 0
90             tokens << [:DEDENT, indent.size]
```

```

91         end
92         tokens << [:\NEWLINE, "\n"]
93     else # indent.size > current_indent, error!
94         raise "Missing ':'" # Cannot increase indent level without using ":"
95     end
96     i += indent.size + 1
97

```

Long operators such as `||`, `&&`, `==`, etc. will be matched by the following block. One character long operators are matched by the catch all `else` at the bottom.

```

101     elsif operator = chunk[/\A(\|\|&&|==|!=|<=|>=)/, 1]
102         tokens << [operator, operator]
103         i += operator.size
104

```

We're ignoring spaces. Contrary to line breaks, spaces are meaningless in our language. That's why we don't create tokens for them. They are only used to separate other tokens.

```

107     elsif chunk.match(/\A /)
108         i += 1
109

```

Finally, catch all single characters, mainly operators. We treat all other single characters as a token. Eg.: `() , . ! + - <`.

```

112     else
113         value = chunk[0,1]
114         tokens << [value, value]
115         i += 1
116
117     end
118
119 end
120

```

Close all open blocks. If the code ends without dedenting, this will take care of balancing the INDENT...DEDENTS.

```
123     while indent = indent_stack.pop
124         tokens << [:DEDENT, indent_stack.first || 0]
125     end
126
127     tokens
128 end
129 end
```

You can test the lexer yourself by running the test file included with the book. Run `ruby -Itest test/lexer_test.rb` from the code directory and it should output 0 failures, 0 errors. Here's an excerpt from that test file.

In file `code/test/lexer_test.rb`

```
1 code = <<-CODE
2 if 1:
3   if 2:
4     print("...")
5     if false:
6       pass
7     print("done!")
8   2
9
10 print "The End"
11 CODE
12 tokens = [
13   [:IF, "if"], [:NUMBER, 1],           # if 1:
14   [:INDENT, 2],
15   [:IF, "if"], [:NUMBER, 2],         # if 2:
16   [:INDENT, 4],
17   [:IDENTIFIER, "print"], ["(", "("], # print("...")
18   [:STRING, "..."],
19   [")", ")"],
20   [:NEWLINE, "\n"],
21   [:IF, "if"], [:FALSE, "false"],    # if false:
```

```

22     [:INDENT, 6],
23     [:IDENTIFIER, "pass"],           #     pass
24     [:DEDENT, 4], [:NEWLINE, "\n"],
25     [:IDENTIFIER, "print"], [("(" , "("),           #     print("done!")
26                               [:STRING, "done!"],
27                               [")", ")"],
28     [:DEDENT, 2], [:NEWLINE, "\n"],
29     [:NUMBER, 2],                   #     2
30     [:DEDENT, 0], [:NEWLINE, "\n"],
31     [:NEWLINE, "\n"],              #
32     [:IDENTIFIER, "print"], [:STRING, "The End"]   # print "The End"
33 ]
34 assert_equal tokens, Lexer.new.tokenize(code)

```

Some parsers take care of both lexing and parsing in their grammar. We'll see more about those in the next section.

DO IT YOURSELF I

- a. Modify the lexer to parse: `while condition: ...` control structures.
- b. Modify the lexer to delimit blocks with `{ ... }` instead of indentation.

[Solutions to Do It Yourself I.](#)

This is a sample chapter.
Buy the full book online at
<http://createyourproglang.com>