



Writing your own programming language and compiler with Python

Marcelo Andrade [Follow](#)

Jun 29 · 7 min read



```
mov    rax, 0x0a687f
push   rax
xor    rax, rax
mov    rsi, rsp
mov    rdi, 1
mov    rdx, 8
call   write
```

* Top highlight

Introduction

After studying compilers and programming languages, I felt like internet tutorials and guides are way too complex for beginners or are missing some important parts about these topics.

My goal with this post is to help people that are seeking a way to start developing their first programming language/compiler.

Requirements

On this guide, I'll be using [PLY](#) as lexer and parser, and [LLVMlite](#) as low level intermediate language to do code generation with optimizations (if you don't know what I'm talking about, don't worry, I'll explain it later).

So, the requirements for this project are:

- [Anaconda](#) (way simpler to install LLVMlite through *conda* than *pip*)
- [LLVMlite](#)

```
$ conda install --channel=numba llvmlite
```

- [RPLY](#) (same as PLY but with a better API)

```
$ conda install -c conda-forge rply
```

- [LLC](#) (LLVM static compiler)
- [GCC](#) (or other linking tool)

Getting Started

"Where to begin?". This is the most common question when trying to create your programming language.

I'll start by defining my own language. Let's call it *TOY*, here's a simple example of a *TOY* program:

```
var x;
x := 4 + 4 * 2;
if (x > 3) {
    x := 2;
} else {
    x := 5;
}
print(x);
```

Although this example is really simple, it is not so easy to be implemented as a programming language. So let's start with a simpler example:

```
print(4 + 4 - 2);
```

But how do you formally describe a language grammar? It is way to hard to create all possible examples of a language to show everything it can do.

To do this, we do what is called an [ERNF](#). It is a metalanguage to define with

To do this, we do what is called an EBNF. It is a metalinguage to define within one document, all possible grammar structures of a language. You can find most programming languages EBNFs easily.

To understand better how a EBNF grammar works, I recommend reading [this post](#).

EBNF

Let's create a EBNF that describes the minimal possible functionality of *TOY*, only a sum operation. It will describe the following example:

```
4 + 2;
```

It's EBNF can be described as:

```
expression = number, "+", number, ";";
number = digit+;
digit = [0-9];
```

This example is way too simple to be useful as a programming language, so let's add some functionalities. The first one, is to be able to add as many numbers as you want, and the second one, is to be able to subtract numbers as well.

Here's an example of our new programming language:

```
4 + 4 - 2;
```

And it's EBNF can be described as:

```
expression = number, { ("+"|"-"), number }, ";";
number = digit+;
digit = [0-9];
```

And finally, adding *print* to our programming language:

```
print(4 + 4 - 2);
```

We get to this EBNF:

```
program = "print", "(", expression, ")"; ;
expression = number, { ("+"|"-"), number } ;
number = digit+;
digit = [0-9];
```

Now that we've defined our grammar, how do we translate it to code? So we can validate and understand a program? And after that, how can we compile it to a binary executable?

Compiler

A compiler is a program that turns a programming language into machine language or other languages. In this guide, I'm going to compile our programming language into [LLVM IR](#) and then into machine language.

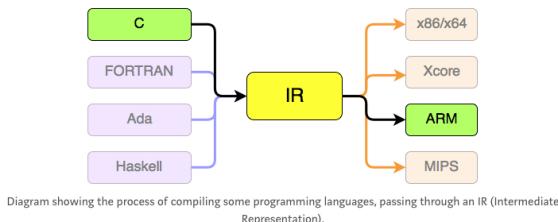


Diagram showing the process of compiling some programming languages, passing through an IR (Intermediate Representation).

Using LLVM, it is possible to optimize your compilation without learning compiling optimization, and LLVM has a really good library to work with compilers.

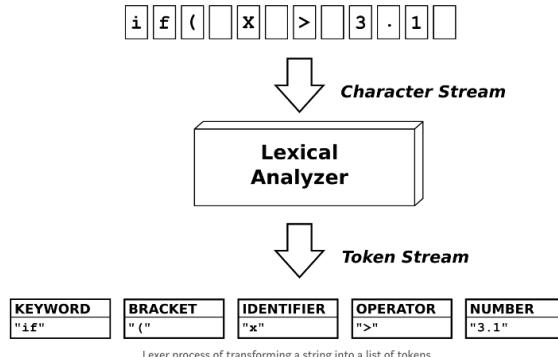
Our compiler can be divided into three components:

- Lexer
- Parser
- Code Generator

For the **Lexer** and **Parser** we'll be using RPLY, really similar to PLY: a Python library with lexical and parsing tools, but with a better API. And for the **Code Generator**, we'll use LLVMLite, a Python library for binding LLVM components.

Lexer

The first component of our compiler is the **Lexer**. It's role is to take the program as input and divide it into *Tokens*.



We use the minimal structures from our EBNF to define our tokens. For example, with the following input:

```
print(4 + 4 - 2);
```

Our Lexer would divide this string into this list of tokens:

```

Token('PRINT', 'print')
Token('OPEN_PAREN', '(')
Token('NUMBER', '4')
Token('SUM', '+')
Token('NUMBER', '4')
Token('SUB', '-')
Token('NUMBER', '2')
Token('CLOSE_PAREN', ')')
Token('SEMI_COLON', ';')

```

So, let's start coding our compiler. First, create a file named `lexer.py`. We'll define our tokens on this file. We'll only use `LexerGenerator` class from RPLY to create our Lexer.

```

1  from rply import LexerGenerator
2
3
4  class Lexer():
5      def __init__(self):
6          self.lexer = LexerGenerator()
7
8      def _add_tokens(self):
9          # Print
10         self.lexer.add('PRINT', r'print')
11         # Parenthesis
12         self.lexer.add('OPEN_PAREN', r'\(')
13         self.lexer.add('CLOSE_PAREN', r'\)')
14         # Semicolon
15         self.lexer.add('SEMI_COLON', r';')
16         # Operators
17         self.lexer.add('SUM', r'\+')
18         self.lexer.add('SUB', r'\-')
19         # Number
20         self.lexer.add('NUMBER', r'\d+')
21         # Ignore spaces
22         self.lexer.ignore('\s+')
23
24     def get_lexer(self):
25         self._add_tokens()
26         return self.lexer.build()

```

lexer.py hosted with ❤ by GitHub [view raw](#)

After this, create your main file named `main.py`. We'll combine all three compiler components on this file.

```

1  from lexer import Lexer
2
3  text_input = """

```

```

1  -----
2  print(4 + 4 - 2);
3  ***
4
5
6
7  lexer = Lexer().get_lexer()
8  tokens = lexer.lex(text_input)
9
10 for token in tokens:
11     print(token)

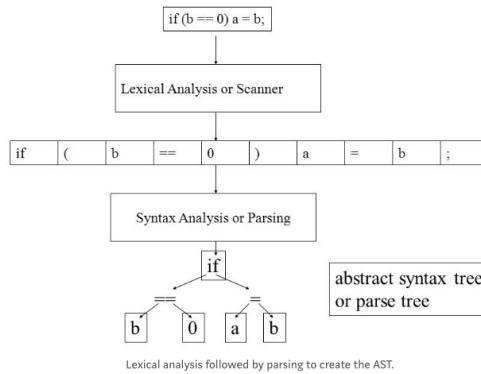
```

main1.py hosted with ❤ by GitHub [view raw](#)

If you run `$ python main.py`, the output of tokens will be the same as described above. You can change the name of your tokens if you want, but I recommend keeping the same to keep consistency with the **Parser**.

Parser

The second component in our compiler is the **Parser**. It's role is to do a syntax check of the program. It takes the list of tokens as input and create an AST as output. This concept is more complex than a list of tokens, so I highly recommend a little bit of research about Parsers and ASTs.



To implement our parser, we'll use the structure created with our EBNF as model. Luckily, RPLY's parser uses a format really similar to the EBNF to create its parser, so it is really straightforward.

The most challenging is to attach the Parser with the AST, but when you get the idea, it becomes really mechanical.

First, create a new file named `ast.py`. It will contain all classes that are going to be called on the parser and create the AST.

```

1  class Number():
2      def __init__(self, value):
3          self.value = value
4
5      def eval(self):
6          return int(self.value)
7
8
9  class BinaryOp():
10     def __init__(self, left, right):
11         self.left = left
12         self.right = right
13
14
15  class Sum(BinaryOp):
16      def eval(self):
17          return self.left.eval() + self.right.eval()
18
19
20  class Sub(BinaryOp):
21      def eval(self):
22          return self.left.eval() - self.right.eval()
23
24
25  class Print():
26      def __init__(self, value):
27          self.value = value
28
29      def eval(self):
30          print(self.value.eval())

```

ast.py hosted with ❤ by GitHub [view raw](#)

Second, we need to create the parser. For that, we'll use `ParserGenerator` from RPLY. Create a file name `parser.py`:

```

1  from rply import ParserGenerator
2  from ast import Number, Sum, Sub, Print
3
4
5  class Parser():
6      def __init__(self):
7          self.pg = ParserGenerator(
8              # A list of all token names accepted by the parser.

```

```

9         ['NUMBER', 'PRINT', 'OPEN_PAREN', 'CLOSE_PAREN',
10        'SEMI_COLON', 'SUM', 'SUB']
11    )
12
13    def parse(self):
14      @self.pg.production('program : PRINT OPEN_PAREN expression CLOSE_PAREN SEMI_COLON')
15      def program(p):
16          return Print(p[2])
17
18      @self.pg.production('expression : expression SUM expression')
19      @self.pg.production('expression : expression SUB expression')
20      def expression(p):
21          left = p[0]
22          right = p[2]
23          operator = p[1]
24          if operator.gettokentype() == 'SUM':
25              return Sum(left, right)
26          elif operator.gettokentype() == 'SUB':
27              return Sub(left, right)
28
29      @self.pg.production('expression : NUMBER')
30      def number(p):
31          return Number(p[0].value)
32
33      @self.pg.error
34      def error_handle(token):
35          raise ValueError(token)
36
37      def get_parser(self):
38          return self.pg.build()

```

[parser1.py](#) hosted with ❤ by GitHub [view raw](#)

And finally, we'll update our file `main.py` to combine Parser with Lexer.

```

1  from lexer import Lexer
2  from parser import Parser
3
4  text_input = """
5  print(4 + 4 - 2);
6 """
7
8  lexer = Lexer().get_lexer()
9  tokens = lexer.lex(text_input)
10
11 pg = Parser()
12 pg.parse()
13 parser = pg.get_parser()
14 parser.parse(tokens).eval()

```

[main2.py](#) hosted with ❤ by GitHub [view raw](#)

Now, if you run `$ python main.py`, you'll see the output being the result of `print(4 + 4 - 2)`, which is equal to printing 6.

With these two components, we have a functional compiler that interprets *TOV* language with Python. However, it still doesn't create a machine language code and is not well optimized. To do this, we'll enter the most complex part of the guide, code generation with LLVM.

Code Generator

The third and last component of our compiler is the **Code Generator**. Its role is to transform the AST created from the parser into machine language or an IR. In this case, it's going to transform the AST into LLVM IR.

This component is the main reason why I'm writing this post. There aren't good guides on how to implement code generation with LLVM on Python.

LLVM can be really complex to understand, so if you wish to fully understand what is going on, I recommend reading [LLVMlite docs](#).

LLVMlite doesn't have a implementation to a `print` function, so you have to define your own.

So, to start, let's create a file named `codegen.py` that will contain the class `CodeGen`. This class is responsible to configure LLVM and create and save the IR code. We also declare the `Print` function on it.

```

1  from llvm import ir, binding
2
3
4  class CodeGen():
5      def __init__(self):
6          self.binding = binding
7          self.binding.initialize()
8          self.binding.initialize_native_target()
9          self.binding.initialize_native_asmprinter()
10         self._config_llvm()
11         self._create_execution_engine()
12         self._declare_print_function()
13
14     def _config_llvm(self):
15         # Config LLVM
16         self.module = ir.Module(name=__file__)
17         self.module.triple = self.binding.get_default_triple()
18         func_type = ir.FunctionType(ir.VoidType(), [], False)
19         base_func = ir.Function(self.module, func_type, name="main")
20         base_func.append_basic_block(name="entry")

```

```

21         self.builder = ir.IRBuilder(block)
22
23     def __create_execution_engine(self):
24         """
25             Create an ExecutionEngine suitable for JIT code generation on
26             the host CPU. The engine is reusable for an arbitrary number of
27             modules.
28         """
29         target = self.binding.Target.from_default_triple()
30         target_machine = target.create_target_machine()
31         # And an execution engine with an empty backing module
32         backing_mod = binding.parse_assembly("")
33         engine = binding.create_mcjit_compiler(backing_mod, target_machine)
34         self.engine = engine
35
36     def __declare_print_function(self):
37         # Declare Printf function
38         voidptr_ty = ir.IntType(8).as_pointer()
39         printf_ty = ir.FunctionType(ir.IntType(32), [voidptr_ty], var_arg=True)
40         printf = ir.Function(self.module, printf_ty, name="printf")
41         self.printf = printf
42
43     def __compile_ir(self):
44         """
45             Compile the LLVM IR string with the given engine.
46             The compiled module object is returned.
47         """
48         # Create a LLVM module object from the IR
49         self.builder.ret_void()
50         llvm_ir = str(self.module)
51         mod = self.binding.parse_assembly(llvm_ir)
52         mod.verify()
53         # Now add the module and make sure it is ready for execution
54         self.engine.add_module(mod)
55         self.engine.finalize_object()
56         self.engine.run_static_constructors()
57         return mod
58
59     def create_ir(self):
60         self.__compile_ir()
61
62     def save_ir(self, filename):
63         with open(filename, 'w') as output_file:
64             output_file.write(str(self.module))

```

codegen.py hosted with ❤ by GitHub

[view raw](#)

After this, let's update our `main.py` file to call `CodeGen` methods:

```

1  from lexer import Lexer
2  from parser import Parser
3  from codegen import CodeGen
4
5  fname = "input.toy"
6  with open(fname) as f:
7      text_input = f.read()
8
9  lexer = Lexer().get_lexer()
10 tokens = lexer.lex(text_input)
11
12 codegen = CodeGen()
13
14 module = codegen.module
15 builder = codegen.builder
16 printf = codegen.printf
17
18 pg = Parser(module, builder, printf)
19 pg.parse()
20 parser = pg.get_parser()
21 parser.parse(tokens).eval()
22
23 codegen.create_ir()
24 codegen.save_ir("output.ll")

```

main3.py hosted with ❤ by GitHub

[view raw](#)

As you can see, I removed the input program from this file and created a new file called `input.toy` to simulate a external program. It's content is the same as the input described.

```
print(4 + 4 - 2);
```

Another change that was made is passing `module`, `builder` and `printf` objects to the Parser. This was made so we could pass this objects to the AST, where the LLVM AST is created. So, we change `parser.py` to receive these objects and pass them to the AST.

```

1  from rply import ParserGenerator
2  from ast import Number, Sum, Sub, Print
3
4
5  class Parser():
6      def __init__(self, module, builder, printf):
7          self.pg = ParserGenerator(
8              # A list of all token names accepted by the parser.
9              ["NUMBER", "PRINT", "OPEN_PAREN", "CLOSE_PAREN",
10               "SEMI_COLON", "SUM", "SUB"])
11
12          self.module = module
13          self.builder = builder
14          self.printf = printf
15
16      def parse(self):
17          #self.pg.production('program : PRINT OPEN_PAREN expression CLOSE_PAREN SEMI_COLON')

```

```

18     def program(p):
19         return Print(self.builder, self.module, self.printf, p[2])
20
21     @self.pg.production("expression : expression SUM expression")
22     @self.pg.production("expression : expression SUB expression")
23     def expression(p):
24         left = p[0]
25         right = p[2]
26         operator = p[1]
27         if operator.gettokentype() == 'SUM':
28             return Sum(self.builder, self.module, left, right)
29         elif operator.gettokentype() == 'SUB':
30             return Sub(self.builder, self.module, left, right)
31
32     @self.pg.production("expression : NUMBER")
33     def number(p):
34         return Number(self.builder, self.module, p[0].value)
35
36     @self.pg.error
37     def error_handle(token):
38         raise ValueError(token)
39
40     def get_parser(self):
41         return self.pg.build()

```

parser2.py hosted with ❤ by GitHub [view raw](#)

And finally, and most important, we change the `ast.py` to receive these objects and create the LLVM AST using LLVMlite methods.

```

1  from llvm import ir
2
3
4  class Number():
5      def __init__(self, builder, module, value):
6          self.builder = builder
7          self.module = module
8          self.value = value
9
10     def eval(self):
11         i = ir.Constant(ir.IntType(8), int(self.value))
12         return i
13
14
15     class BinaryOp():
16         def __init__(self, builder, module, left, right):
17             self.builder = builder
18             self.module = module
19             self.left = left
20             self.right = right
21
22
23         class Sum(BinaryOp):
24             def eval(self):
25                 i = self.builder.add(self.left.eval(), self.right.eval())
26                 return i
27
28
29         class Sub(BinaryOp):
30             def eval(self):
31                 i = self.builder.sub(self.left.eval(), self.right.eval())
32                 return i
33
34
35         class Print():
36             def __init__(self, builder, module, printf, value):
37                 self.builder = builder
38                 self.module = module
39                 self.printf = printf
40                 self.value = value
41
42             def eval(self):
43                 value = self.value.eval()
44
45                 # Declare argument list
46                 voidptr_ty = ir.IntType(8).as_pointer()
47                 fmt = "%i\n"
48                 c_fmt = ir.Constant(ir.ArrayType(ir.IntType(8), len(fmt)),
49                                     bytarray(fmt.encode("utf8")))
50                 global_fmt = ir.GlobalVariable(self.module, c_fmt.type, name="fmt")
51                 global_fmt.linkage = "internal"
52                 global_fmt.global_constant = True
53                 global_fmt.initializer = c_fmt
54                 fmt_arg = self.builder.bitcast(global_fmt, voidptr_ty)
55
56                 # Call Print Function
57                 self.builder.call(self.printf, [fmt_arg, value])

```

ast2.py hosted with ❤ by GitHub [view raw](#)

With these changes, our compiler is ready to transform a TOY program into an LLVM IR file `output.ll`. To compile this `.ll` file into a executable, we'll use LLVM to create a object file `output.o`, and finally, GCC (you could use other linking programs) to create the final executable file.

```
$ llc -filetype=obj output.ll
$ gcc output.o -o output
```

And you can finally run your executable file compiled from the initial program.

```
$ ./output
```

Next Steps

After this guide, I hope you can understand an EBNF and the three basic concepts of a compiler. With this knowledge, you now can create your own programming language and write a optimized compiler to it with Python. I encourage you to go further and add new elements to your language and compiler, here are some ideas:

- Statements
- Variables
- New Binary Operators (Multiplication, Division)
- Unary Operators
- If Statement
- While Statement

Feel free to send me any compiler projects. I'll be happy to help you with anything.

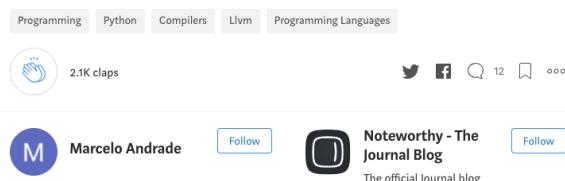
You can contact me at marceloga1@al.insper.edu.br

Hope you all enjoyed this post and got some love to programming languages and compilers!

You can see the [final code on GitHub](#).

This story is published in Noteworthy, where 10,000+ readers come every day to learn about the people & ideas shaping the products we love.

Follow our publication to see more product & design stories featured by the [Journal team](#).

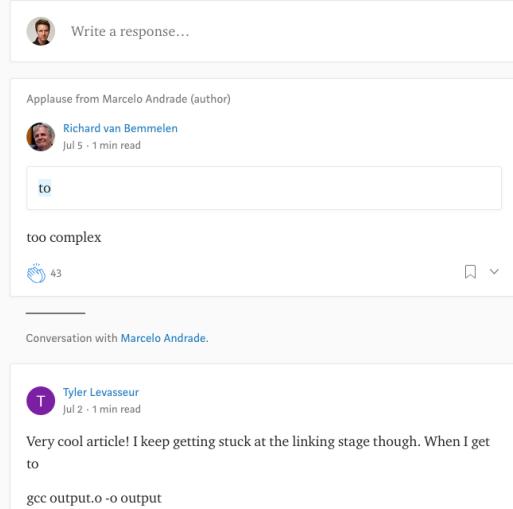


The image shows the Noteworthy publication profile on a platform like Medium. It includes a header with categories: Programming, Python, Compilers, LLVM, and Programming Languages. Below the header are social sharing icons for Twitter, Facebook, and LinkedIn, along with a clap count of 2.1K. The publication's logo is a blue circle with a white 'M'. The author's name, Marcelo Andrade, is listed with a 'Follow' button. The title of the publication is 'Noteworthy - The Journal Blog' with the subtitle 'The official Journal blog'.



The image displays three article cards from Noteworthy - The Journal Blog. The first card features a cartoon illustration of a face with hands holding it, titled 'Little known features of JavaScript' by Viral Shah, with a 10 min read duration. The second card shows a close-up of a coffee cup with the word 'HIGH' written on it, titled 'Hate Your Current Life? Start Waking Up Early' by Alisha Beotra, with a 6 min read duration. The third card is a collage of various images related to UX design, titled 'A year of building UX at Virgin Atlantic' by Martyn Reding, with a 12 min read duration.

Responses



The image shows the Noteworthy response interface. At the top, there is a text input field with a placeholder 'Write a response...'. Below it, a comment from Richard van Bemmelen is shown, dated Jul 5, with a 1 min read duration. The comment text is 'to too complex'. Below this, a response from Tyler Levasseur is shown, dated Jul 2, with a 1 min read duration. The response text is 'Very cool article! I keep getting stuck at the linking stage though. When I get to gcc output.o -o output'.

I get the error:

```
/usr/bin/ld: output.o: relocation R_X86_64_32 against `'.rodata' can not be  
used when making a shared object; recompile with -fPIC  
/usr/bin/ld: final...
```

[Read more...](#)



5

1 response

Marcelo Andrade
Jul 4 · 1 min read

I was looking into this problem, and it seems that it is a problem with LLVM 6 and GCC (correct me if I'm wrong and this is not wrong case).

I couldn't find what causes it on this version on LLVM, but a solution I found was changing the linker to Clang.

[Read more...](#)



6

1 response

Conversation with [Marcelo Andrade](#).

Shrek Wu
Jun 30 · 1 min read

llc command not found?



2 responses

Marcelo Andrade
Jun 30 · 1 min read

Have you installed LLVM? I listed it on requirements section.



1 response

[Show all responses](#)