

# LINGI2132

Languages and Translators

**Emilyen Laffineur**

Teacher :

Nicolas Laurent

Assistant:

Mathieu Jadin & Alexandre Dubray



Ecole Polytechnique de Louvain  
Universtié catholique de Louvain  
Belgium

Scholar year : 2020-2021

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Compilation Pipeline</b>	<b>2</b>
2.1	Lexing . . . . .	2
2.2	Parsing . . . . .	2
2.3	Semantic Analysis . . . . .	3
2.3.1	Type checking . . . . .	3
2.3.2	Name Binding . . . . .	3
2.3.3	Rest . . . . .	3
2.4	Optimization . . . . .	3
2.4.1	Tree-Walk Interpreter . . . . .	4
2.5	Code Generation . . . . .	4
2.6	Optimizations . . . . .	5
2.6.1	Inlining . . . . .	6
<b>3</b>	<b>Formal Grammars</b>	<b>7</b>
3.1	Language vs Grammar . . . . .	7
3.2	Grammar notation . . . . .	7
3.2.1	PEG . . . . .	7
3.2.2	EBNF and BNF . . . . .	8
3.2.3	Summary . . . . .	8
<b>4</b>	<b>Parsers</b>	<b>9</b>
4.1	Parsing Tools . . . . .	9
4.1.1	Parsing Generator . . . . .	9
4.1.2	Parsing Library . . . . .	10
4.2	(Abstract) Syntax Tree . . . . .	10
4.2.1	Syntax Tree . . . . .	10
4.2.2	Abstract Syntax Tree . . . . .	11
4.3	Summary 1 . . . . .	11
4.4	Coding a parser . . . . .	12
4.5	Parser Combinator . . . . .	12
<b>5</b>	<b>PEG &amp; CFG Semantics</b>	<b>13</b>

5.1	Recap . . . . .	13
5.2	Context-Free Grammars . . . . .	13
5.2.1	Semantics Derivation . . . . .	14
5.3	PEG Semantics . . . . .	14
5.4	PEG vs CFG . . . . .	14
5.4.1	Semantics . . . . .	14
5.4.2	The difference . . . . .	14
5.4.3	The big difference . . . . .	14

# Chapter 1

## Introduction

This class is about compiler. The goal of compiler is to implement a programming language. It take source code and make it executable (or pass it to a program that execute it). We may think that language like Python or JavaScript does not use compilers however they use it as well! Even if they are dynamical. Note that this class is not about programming paradigms, it's truly about implementing languages in general.

We will learn about compilers and their parts, allowing us to understand how languages works. We will also learn some useful patterns (Domain Specific Language (DSL), trees, interpreter pattern).

## Chapter 2

# Compilation Pipeline

Here is the classical compiler pipeline :

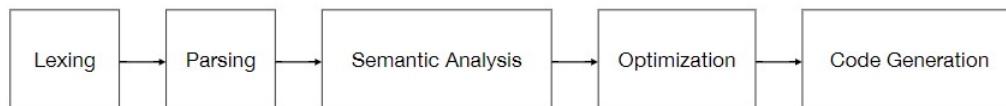


Figure 2.1: Compiler pipeline

We will go through each step.

### 2.1 Lexing

**Definition 2.1.1** (Lexing). It's a lexical analysis, tokenization, scanning. It take the text and produce tokens (= "words"). e.g : `foo = bar + 42` will produce `|foo| = |bar| + |42|`.

### 2.2 Parsing

**Definition 2.2.1** (Parsing). It take the output of lexing (tokens) and extract it's structure (by producing an Abstract Syntax Tree (AST)). It also catches syntax errors (like missing parenthesis around if in Java). Other kind of errors (like something that is not defined by the language) is not a syntax error, thus, the parser cannot catch it.

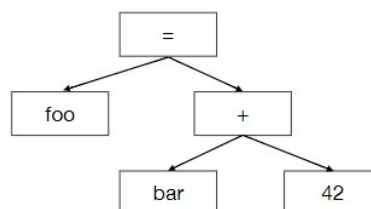


Figure 2.2: Abstract Syntax Tree

Note that lexing is optional while parsing (called scanner-less parser) for others parsing is actually lexing AND parsing.

## 2.3 Semantic Analysis

It consists of a series of check that is perform on the AST, two mains are Name binding et type checking.

### 2.3.1 Type checking

**Definition 2.3.1** (Type checking). Verify that the type passed is what is expected.

```
1 void print(String str)
2 print(42)
```

**Definition 2.3.2** (Type inference). Let the compiler guess the type of the variable given the context.

Note that type checking is inference. When :

```
1 String x = "foobar"
```

is used, it first need to guess what is "foobar", then check the consistency with the type.

Sometimes, we need to know the return value's type of a function to perform the check. That's what Name Binding is.

### 2.3.2 Name Binding

**Definition 2.3.3** (Name binding). Allow us to infer the type of a variable given the return type of a function.

Note that it becomes much more complex with OOP.

### 2.3.3 Rest

It also performs flow check (missing return for example), access control (public, private, etc.).

## 2.4 Optimization

Actually, the pipeline we saw in figure 2.1 it's only the old school C compiler. New C compiler, do optimization on the produced machine code (LLVM, Low Level Virtual Machine)! Do not get confuse, it does not use a virtual machine at all... LLVM generate IR (intermediate representation), optimize it and then generate machine code.

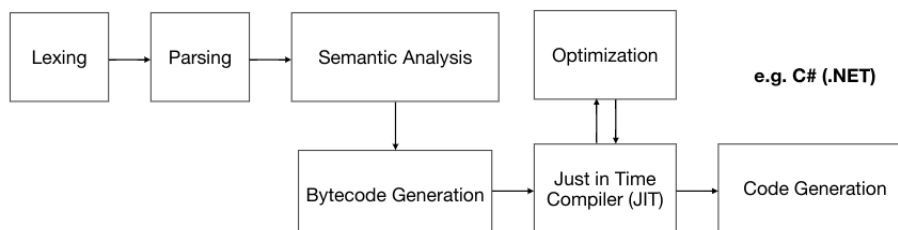


Figure 2.3: C# pipeline

This one generate bytecode (Java name), that use usable but that won't be directly use by the machine. After the generation of the bytecode, it is passed to a Just In Time Compiler (JIT).

What is done until JIT is done at compilation time, JIT and the following is done while running the program!

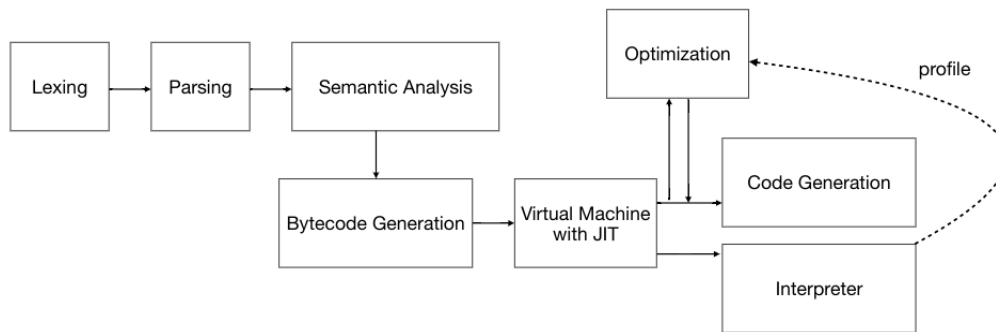


Figure 2.4: Java and Python pipeline

That's the architecture that most recent languages uses. The bytecode is passed to a VM executing a JIT. This VM can generate machine code (code generation) or pass it to in interpreter which will run a runtime profile that is much better for optimization! That the reason why Java can be as fast (or faster) the C. Pypy (Python), Java, TruffleRuby use this kind of architecture!

We could also not optimize it! We could directly use the Tree-Walk interpreter, Ruby before 1.9 did that. CPython (standard one), and Ruby generate the bytecode and run it in a VM.

### 2.4.1 Tree-Walk Interpreter

Let's take the following program :

```

1  var foo = 42 + 52
2  print(foo)

```

The AST is the following :

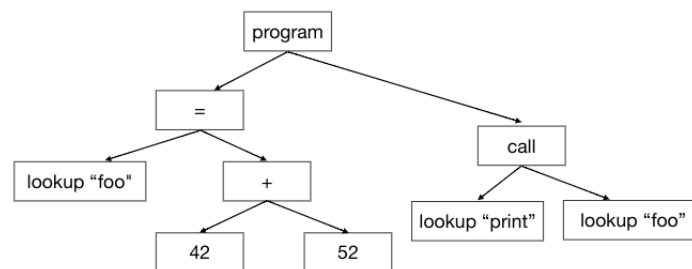


Figure 2.5: Resulting AST

(Note that print has its own AST). We have the scope, containing the value of foo saved in the store and the definition of print. The procedure in order to execute is recursive, it will lookup for foo, perform the addition, store it. Lookup at the print definition, lookup for foo and foo will be passed as a parameter to print.

## 2.5 Code Generation

It could be machine or bytecode generation. Example of Machine Code (e.g. x64) :

```

1  int square(int num) {
2      return num*num
3  }

```

Becomes :

```

1  square(int):
2      push rbp
3      mov rbp, rsp
4      mov DWORD PTR [rbp-4], edi
5      mov eax, DWORD PTR [rbp-4]
6      imul eax, eax
7      pop rbp
8      ret
9  With optimization it becomes :
10     mov eax, edi
11     imul eax, edi
12     ret

```

While bytecode generation is:

```

1  public class Hello {
2      int square(int num) {
3          return num*num
4      }
5  }

```

Which becomes (stacked base instead of registry based as before):

```

1  public class Hello {
2      public <init>()V
3          L0 LINENUMBER 1 L0
4          ALOAD 0
5          INVOKEVIRTUAL java/lang/Object.<init> ()V
6          RETURN
7      L1
8      LOCALVARIABLE this LHello; L0 L1 0
9      MAXSTACK = 1
10     MAXLOCALS = 1
11     square(I)I
12     L0
13         LINENUMBER 3 L0
14         ILOAD 1
15         ILOAD 1
16         IMUL
17         IRETURN
18     L1
19     LOCALVARIABLE this LHello; L0 L1 0
20     LOCALVARIABLE num I L0 L1 1
21     MAXSTACK = 2
22     MAXLOCALS = 2}

```

## 2.6 Optimizations

As we have seen, we can optimize on tree or on target code. They are two types of optimization that are the base of everything :

- Inlining : pulling a function in another one
- Partial evaluation : propagating known information (e.g constant-folding)

Some other exists like loop unrolling, etc.



### 2.6.1 Inlining

Take the following program :

```
1  var foo = false
2  int a(){
3      if (foo) return somethingLongAndBoring()
4      else return b(true)+1
5  }
6  int b(boolean bar) {
7      return bar ? 42 : 52
8  }
```

If we optimize the function a. If we inline somethingLongAndBoring and b. However, if we do this and do not enter the condition, this could be bad for cache locality. Indeed, it will lead to a cache miss, with a lot of code, the time to retrieve the actual code will be huge. However, we know that foo is always false so we can get rid of the first condition! Then we can inline b. Thanks to the inlining of b, then constant folding it will lead on only this :

```
1  int a(){
2      return 43;
3  }
```

Which is much more efficient! Thanks to that we can see that there is a complementary between constant folding and inlining.

## Chapter 3

# Formal Grammars

### 3.1 Language vs Grammar

**Definition 3.1.1** (Grammar). The (formal) definition (*description*) of a language

**Definition 3.1.2** (Language). A (potentially infinite) *set* of sentences, in programming language, a sentence is a source file, REPL expression, etc.

**Definition 3.1.3** (Alphabet). Composed of tokens, lexemes.

Let's take an example which is JSON. The language is all valid JSON expression. Sentences : e.g.

```
{
  "version": 17,
  "bundles": [
    { "name" : "org.graalvm.component.installer.Bundle" },
    { "name" : "org.graalvm.component.installer.commands.Bundle" },
    { "name" : "org.graalvm.component.installer.remote.Bundle" },
    { "name" : "org.graalvm.component.installer.os.Bundle" }
  ]
}
```

Figure 3.1: JSON Sentence

Grammar :

```
VALUE  ::= STRINGLIT / NUMBER / OBJECT / ARRAY
OBJECT ::= "{" (PAIR (" " PAIR)* )? "}"
PAIR   ::= STRINGLIT ":" VALUE
ARRAY  ::= "[" (VALUE (" " VALUE)* )? "]"
```

**Alphabet  
(Tokens)**

Figure 3.2: JSON Grammar

### 3.2 Grammar notation

#### 3.2.1 PEG

Parsing Expression Grammar. The "?" denotes optional, "\*" means 0 or more times.

```

VALUE    ::= STRINGLIT / NUMBER / OBJECT / ARRAY
OBJECT   ::= "{" (PAIR ("," PAIR)* )2 "}"
PAIR     ::= STRINGLIT ":" VALUE
ARRAY    ::= "{" (VALUE ("," VALUE)* )2 "}"

```

Figure 3.3: PEG notation

PEG is a grammar formalism and a notation.

**Definition 3.2.1** (Formalism). Mathematical system to define the language (set of sentences)

**Definition 3.2.2** (Notation). A way to denote a grammar to be interpreted by the formalism.

### 3.2.2 EBNF and BNF

Note that (E)BNF is a notation for CFG (Context-Free Grammars).

#### EBNF

Extended Backus-Naur Form. "[" means optional, "{" means 1 or more times, so in order to make "0 or more" with have to combine symbols.

```

VALUE    ::= STRINGLIT | NUMBER | OBJECT | ARRAY
OBJECT   ::= "{" [ PAIR [ "{" " , " PAIR ] ] "}"
PAIR     ::= STRINGLIT ":" VALUE
ARRAY    ::= "{" (VALUE [ "{" " , " VALUE ] ] "}"

```

Figure 3.4: EBNF notation

#### BNF

Backus-Naur Form. It is at the base of EBNF, it does not contain "{" or "[", so it uses recursion and a lot more rules. Note that  $\epsilon$  is the mathematical notation for "nothing"

```

VALUE    ::= STRINGLIT | NUMBER | OBJECT | ARRAY
OBJECT   ::= "{" PAIRS? "}"
PAIRS?   ::=  $\epsilon$  | PAIRS
PAIRS    ::= PAIR TPAIRS
TPAIRS   ::=  $\epsilon$  | " , " PAIRS
PAIR     ::= STRINGLIT ":" VALUE
ARRAY    ::= "[" VALUES? "]"
VALUES?  ::=  $\epsilon$  | VALUES
VALUES   ::= VALUE TVALUES
TVALUES  ::=  $\epsilon$  | " , " VALUES

```

Figure 3.5: BNF notation

### 3.2.3 Summary

CFG and PEG are different formalism. In some case, PEG and CFG can be the same (like in JSON) because it is a very simple definition. By it is not true in general. Sometimes, we can use PEG notations like +, \*, ? in CFGs (coming from REGEx, which are related to CFGs.). Actually we can use any notation that we want, we just need to define it!

# Chapter 4

## Parsers

**Definition 4.0.1** (Parser). A parser is a program that :

- Accepts/rejects input a sentence in the language (recognizer)
- Extracts a(n) (abstract) syntax tree

### 4.1 Parsing Tools

A parsing tool lets you create parsers. We can denote two kinds : Generator and Library. Note that is just a generalization. The main distinction to retain is the generation vs interpretation.

Parsing tools are often called "parsers", it is acceptable, but incorrect, parsing tools create and/or run parsers.

Parsing tools are compilers!

- Language = grammar notation
- Code generation (parser generators) or interpretation (parsing libraries) (depending on the used tool)
- Domain Specific Language (DSL) (depending on how we define it. DSL is opposed to Turing-Complete language like Java, Ruby, Python, etc. A DSL could be XML, JSON, Grammar, SQL, etc.)

#### 4.1.1 Parsing Generator

Note : in the following figure, all solid boxes denote programs.

**Definition 4.1.1** (Parsing Generator). A Parsing Generator is :

- Outputs a parser program (typically source code)
- Typically a command line tool
- Typically uses own grammar language

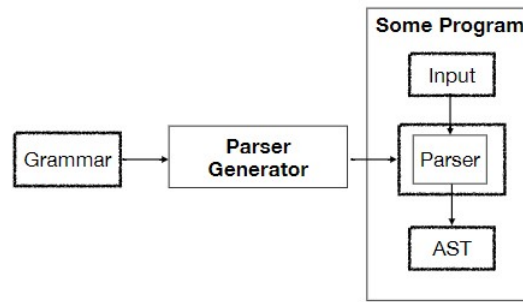


Figure 4.1: Parser Generator

### 4.1.2 Parsing Library

**Definition 4.1.2** (Parsing Library). A Parsing Library is :

- Let us "interpret" a grammar
- Grammar typically defined with a DSL.

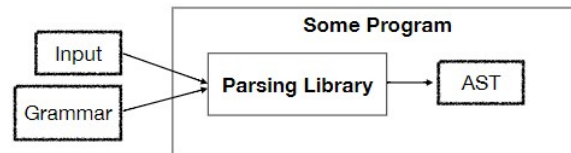


Figure 4.2: Parser Library

Note that input AND grammar could be in the "Some Program" box.

## 4.2 (Abstract) Syntax Tree

### 4.2.1 Syntax Tree

Remember the notation of figure 3.2 and 3.1. To make it easier for Syntax Trees we will use the figure 3.5. Doing so, we can extract the following ST :

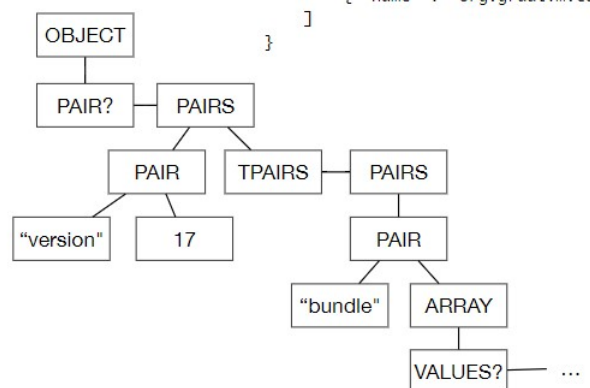


Figure 4.3: Syntax Tree

Syntax Tree is the following : for each phrase in the sentence we follow the grammar. It is the only thing we can do!

### 4.2.2 Abstract Syntax Tree

Here is the same Syntax Tree but abstracted :

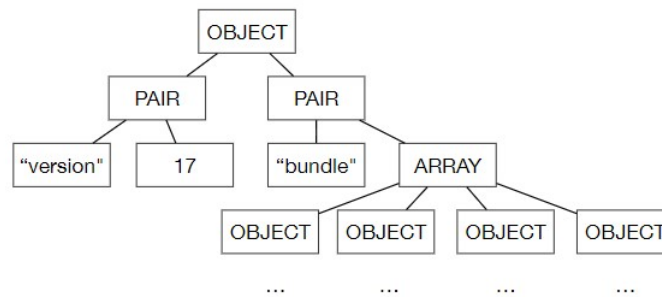


Figure 4.4: Abstract Syntax Tree

As we can see, it is much more clear than the classical Syntax Tree.

## 4.3 Summary 1

Parsing Tools that allow us to use grammar definition like PEG or EBNF (with more symbols that denotes more explicit stuffs) will generally produce better Syntax Tree than the others. However AST let us decide precisely what we want as node/leaves in the tree. Let's take an example using Java :

- 1 list.forEach(x -> System.out.println(x));
- 2 list.reduce((x,y) -> x + y);

As we can see, the difference is that with 1 identifier we do not need parenthesis. That must be take into account when generating the Syntax Tree. It would lead to the following Syntax Trees :

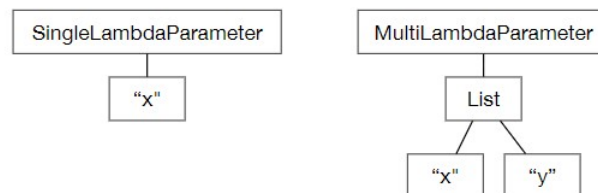


Figure 4.5: First ST

However, as we have seen, parsing is just the start of the pipeline! We have to do a lot of work with it (check for errors, generate bytecode, etc.). With this kind of Syntax Tree, we would have to deal with different kind of node for the same lambda principle. Using AST we can have the following :



Figure 4.6: Second AST

Which allow us to use the same code to deal with both!

## 4.4 Coding a parser

General principles :

- One function per symbol (non-terminals & terminals);
- Model in the input as globally accessible character array;
- The current input position is globally accessible;
- Calling a function = attempting to parse the symbol it denotes at the current input position;
- A parsing function returns :
  - true if the symbol was matched, updating the input position past the matched input.
  - false if they failed to match, the input remains unchanged.

See `parser.java` file. Note that, in order to be sure our implementation really work, we can use the debugger and check for the pos cursor at the end, or to go further we could also inspect the parse tree if we have build one.

As we can see, the parser we wrote is verbose and cumbersome. Also it does not have any bells or whistles, if it fails it just fails, it does not tell us where it fails or why. What we have done is implementing a PEG semantics for the grammar (hint, `a* a` is empty). Backtracking in PEG is to reset the counter before the start of the sequence. CFG also have a backtracking technique but it is not the same.

We can also add AST to our parser. In order to do that in Java, we just need to add a Dequeue and modify only a few methods. See `parser_ast.java` for that.

## 4.5 Parser Combinator

The parser we have implemented from PEG grammar. As we have seen, this parsers have a lot of duplication (each "expression" is handled the same way). The idea to solve that is to build an AST where each node is an "expression" and interpret it.

The idea, is that each combinator is a node/parsing expression in the AST. See file `Combinators.java`. Using combinators allow us to cut down on verbosity & code duplication (across all grammars and within a specific grammar). However it still have some drawbacks :

- Harder to debug
- It's slower because of megamorphism (we'll come back later on this)

Note that theses issues can be eliminated by using combinators for code generation. Good frameworks (like Autumn) will mitigate usability issues.

## Chapter 5

# PEG & CFG Semantics

### 5.1 Recap

**Definition 5.1.1** (Grammar). The (formal) definition of a language

**Definition 5.1.2** (Language). A (potentially infinite) set of sentences

**Definition 5.1.3** (Parser). Recognizes a sentence in the language + extract a syntax tree

**Definition 5.1.4** (Parsing Tool). Generate and/or runs parsers

Two types of formalism :

- CFG (Context-Free-Grammar)
- PEF (Parsing Expression Grammar)

Notations :

- (E)BNF (usually for CFG)
- PEG Notation (usually for PEG)

**Definition 5.1.5** (Non-Terminal). Things that we define

**Definition 5.1.6** (Terminals). Tokens, strings, etc. that cannot be extended further.

### 5.2 Context-Free Grammars

Usually a CFG is defined as a tuple of 4 components  $\text{Grammar} = (N, \Sigma, P, S)$ :

- $N$ : Non-Terminals
- $\Sigma$ : Alphabet (Terminals)
- $P$ : Production Rules ( $P : N \rightarrow (\Sigma \cup N)^*$ )
- Starting Symbol ( $S \in N$ )

In order to have a CFG from a "full" grammar, we first need to replace all syntactic sugars by the complete recursive rules, then eliminate all choices by introducing as many rules as we have choices.



### 5.2.1 Semantics Derivation

- The language defined by a CFG is the set of all sentences that can be derived from its rules
- Start from the start symbol, replace it by the right-hand side of one its production
- At each step, replace a non-terminal from the current string symbol by its definition (until no non-terminals are left in the string)
- Any terminal, the order does not matter

Note that by doing that, we define the language, not a parsing algorithm and we're doing it in a generative way grammar  $\rightarrow$  sentences.

Also, if we would have used characters derivation instead of tokens we would have need to continue to derive which would of potentially lead to an infinite derivation. So, we can say that for a sentence to be part of a language in need to be derivable however we cannot describe the all derivation table because in many case, it would be infinite.

## 5.3 PEG Semantics

We can see in the slide (9 of PEG - CFG semantics's pdf) the use of a top-down recursive descent parsers that produce production rules (with lookahead operator as an exception).

## 5.4 PEG vs CFG

### 5.4.1 Semantics

CFG are generative (their semantics is given by constructing the language set by the grammar through derivation). On the other hand PEG are recognition based : a sentence is in the language defined by a PEG grammar only if it is recognized by the language recognizer. In order to formalize PEG grammar we need to formalize the recognizer for the grammar.

These two approaches are very different in the mathematical point of view. Also, CFG is easier to defined in the mathematical language, on the other hand the PEG is very much related to the practice.

### 5.4.2 The difference

- CFG has unordered choices, while derivating we can pick any non-terminal we want
- PEG has ordered choices, the first matching is the "correct" one

### 5.4.3 The big difference