



STRUMENTA

About us ▾

Services

Products

Articles

Contacts

The tomassetti.me website has changed: it is now part of strumenta.com. You will continue to find all the news with the usual quality, but in a new layout.

# Parsing in Java: Tools and

L

W  
G  
in  
P

X

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe  
at any time.

**Table of contents**

# STRUMENTA

If you need to parse a language, or document, from

Articles

Products

In Java there are fundamentally three ways to solve the problem:

- use an existing library supporting that specific language: for example a library to parse XML
- building your own custom parser by hand
- a tool or library to generate a parser: for example ANTLR, that you can use to build parsers for any language

## Use An Existing Library ×

All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.

parser generators, or you have specific requirements



STRUMENTA

that you cannot satisfy using a typical parser generator. For instance, because you need the best possible performance or a deep integration between different components.

About us | Services | Products  
Articles | Contacts

## A Tool Or Library To Generate A Parser

In all other cases the third option should be the default one, because is the one that is most flexible and has the shorter development time. That is why on this article we concentrate on the tools and libraries that correspond to this option.

*Note: text in blockquote describing a program comes from the original documentation*

All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.





# Tools To Create Parsers

[About us ▾](#)[Services](#)[Products](#)[Articles](#)[Contacts](#)

We are going to see:

- tools that can generate parsers usable from Java (and possibly from other languages)
- Java libraries to build parsers

Tools that can be used to generate the code for a parser are called **parser generators** or **compiler compiler**. Libraries that create parsers are known as **parser combinators**.

Parser generators (or parser combinators) are not trivial: you need some time to learn how to use them

All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

Email Address

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe at any time.



# Useful Things To Know About Parsers

[About us](#)[Services](#)[Products](#)[Articles](#)[Contacts](#)

To make sure that these list is accessible to all programmers we have prepared a short explanation for terms and concepts that you may encounter searching for a parser. We are not trying to give you formal explanations, but practical ones.

## Structure Of A Parser

A parser is usually composed of two parts: a *lexer*, also known as *scanner* or *tokenizer*, and the proper parser. Not all parsers adopt this two-steps schema:

All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe at any time.



**STRUMENTA**

The lexer scans the text and find ‘4’, ‘3’, ‘7’ and then the space ‘ ‘. The job of the lexer is to ~~about the~~ recognize that **Products** the first characters constitute one token of type **Articles** **NUM**. Then the lexer finds a ‘+’ symbol, which corresponds to a second token of type **PLUS**, and lastly it finds another token of type **NUM**.

[Contacts](#)

The parser will typically combine the tokens produced by the lexer and group them.

The definitions used by lexers or parser are called *rules* or *productions*. A lexer rule will specify that a sequence of digits correspond to a token of type **NUM**, while a parser rule will specify that a sequence of tokens of type **NUM**, **PLUS**, **NUM** corresponds to an expression.



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe at any time.





STRUMENTA

There are two terms that are related and sometimes

they are used inter**changeable**: parse**tree** and services

Products

Abstract SyntaxTree (AST).

Articles

Contacts

Conceptually they are very similar:

- they are both **trees**: there is a root representing the whole piece of code parsed. Then there are smaller subtrees representing portions of code that become smaller until single tokens appear in the tree
- the difference is the level of abstraction: the parse tree contains all the tokens which appeared in the program and possibly a set of intermediate rules. The AST instead is a polished version of the parse tree where the

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

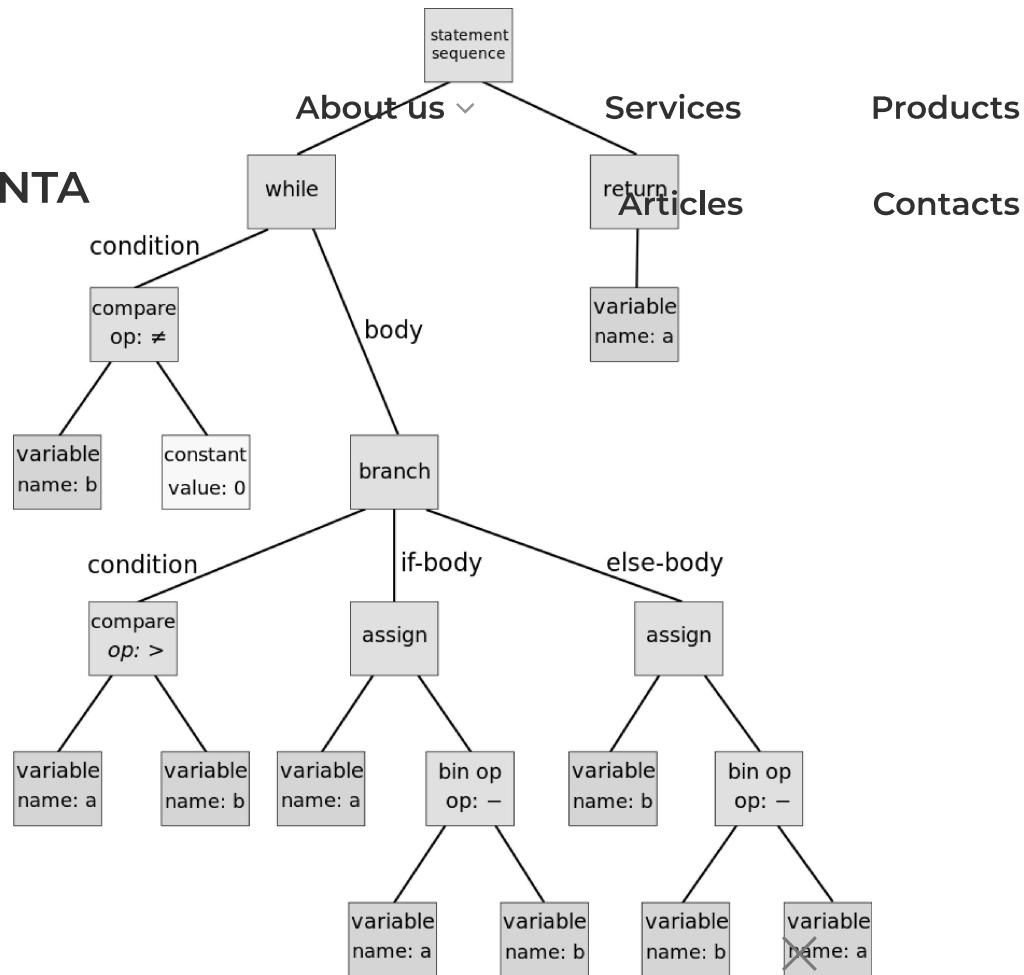
Email Address

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe  
at any time.



STRUMENTA



All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

---

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



## STRUMENTA

In simple terms is a list of rules that define how each construct can be composed. For example, a rule for **Products** an if statement could specify that it must starts with the “if” keyword, followed by a left parenthesis, an expression, a right parenthesis and a statement.

A rule could reference other rules or token types. In the example of the if statement, the keyword “if”, the left and the right parenthesis were token types, while expression and statement were references to other rules.

The most used format to describe grammars is the **Backus-Naur Form (BNF)**, which also has many variants, including the **Extended Backus-Naur Form**.

The Extended variant has the advantage of including a simple way to denote repetitions. A typical rule in a

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

# Left-recursive Rules


[About us](#) ▾

[Services](#)
[Products](#)

In the context of parsers an important feature is the support for left-recursive rules. This means that a rule could start with a reference to itself. This reference could be also indirect.

[Articles](#)
[Contacts](#)

Consider for example arithmetic operations. An addition could be described as two expression(s) separated by the plus (+) symbol, but an expression could also contain other additions.

- ```

1. addition      ::= expression '+'  

   expression  

2. multiplication ::= expression '*'  

   expression  

3. // an expression could be an addition or

```



**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.





# Types Of Languages And

## Grammars

STRUMENTA

About us ▾

Services

Products

Articles

Contacts

We care mostly about two types of languages that can be parsed with a parser generator: *regular languages* and *context-free languages*. We could give you the formal definition according to the Chomsky hierarchy of languages, but it would not be that useful. Let's look at some practical aspects instead.

A regular language can be defined by a series of regular expressions, while a context-free one need something more. A simple rule of thumb is that if a grammar of a language has recursive elements it is not a regular language. For instance, as we said 

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



# The Differences Between

## PEG and CFG

ABOUT US ▾

Services

Products

STRUMENTA

The main difference between PEG and CFG is that the

ordering of choices is meaningful in PEG, but not in CFG. If there are many possible valid ways to parse an input, a CFG will be ambiguous and thus wrong.

Instead with PEG the first applicable choice will be chosen, and this automatically solve some ambiguities.

Another difference is that PEG use scannerless parsers: they do not need a separate lexer, or lexical analysis phase.

Traditionally both PEG and some CFG have been unable to deal with left-recursive rules, but some

All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.



**STRUMENTA**

The basic workflow of a parser generator tool is quite

simple: you write a grammar that defines the

language, or document, and you run the tool to

generate a parser usable from your Java code.

[Services](#)[Products](#)[Articles](#)[Contacts](#)

The parser might produce the AST, that you may have to traverse yourself or you can traverse with additional ready-to-use classes, such as **Listeners** or **Visitors**. Some tools instead offer the chance to embed code inside the grammar to be executed every time the specific rule is matched.

Usually you need a runtime library and/or program to use the generated parser.

## Regular (I ever)



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



STRUMENTA

AnnoFlex is an annotation-based tool, but it does not use proper Java annotations. Instead it relies on custom Javadoc tags, that are interpreted by AnnoFlex. The reason for this design choice is to avoid the need to double escape regular expressions inside strings, which would have been necessary with strings inside annotations.

This is an interesting compromise: on one hand is easier to create regular expressions, but it looks a bit inelegant. It is pragmatism over formality and we think that most people would like it.

The following example shows a simple AnnoFlex program.



All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.





STRUMENTA

```

20.          // this function and the
21.          following will be generated by AnnoFlex
22.          About us Services
23.          scanner.setString("Hello
24.          WorldGoodbye World");
25.          String curToken =
26.          scanner.getNextToken();
27.          while (curToken != null) {
28.              // it prints I am arriving or
29.              // leaving: "Hello World"
30.              // and I am arriving or
31.              // leaving: "Goodbye World"
32.              System.out.println(curToken+":
33.              "+scanner.getMatchText());
34.              curToken =
35.              scanner.getNextToken();
36.          }
37.      }
38.  }
39. }
```

**Products****Articles****Contacts**

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



AnnoFlex requires Java 7 or later.



About us ▾

Services

Products

Articles

Contacts

JFlex is a lexical analyzer (lexer) generator based upon deterministic finite automata (DFA). A JFlex lexer matches the input according to the defined grammar (called spec) and executes the corresponding action (embedded in the grammar).

It can be used as a standalone tool, but being a lexer generator is designed to work with parser generators: typically it is used with CUP or BYacc/J. It can also work with ANTLR.

The typical grammar (spec) is divided three parts,

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



**STRUMENTA**

```

13.
14. LineTerminator = r|n|rn
15.     About us ▾
16.     %% Services
17. // third section Products
18.
19. /* keywords */
20. <YYINITIAL> "abstract" { return
21.     symbol(sym.ABSTRACT); }
22. <YYINITIAL> "boolean" { return
23.     symbol(sym.BOOLEAN); }
24. <STRING> {
25.     " yybegin(YYINITIAL);
26.     return
27.     symbol(sym.STRING_LITERAL,
28.     string.toString()); }
29. [...]

```

Services

Articles

Contacts



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

developed by the author and described in this paper.



STRUMENTA

Adaptive LL(\*) Parsing: The Power of Dynamic

Analysis (PDF)

About us ▾

Services

Products

Articles

Contacts

It can output parsers in many languages. But the real added value of a vast community it is the large amount of grammars available. The version 4 supports direct left-recursive rules.

It provides two ways to walk the AST, instead of embedding actions in the grammar: visitors and listeners. The first one is suited when you have to manipulate or interact with the elements of the tree, while the second is useful when you just have to do something when a rule is matched.

The typical grammar is divided in two parts: lexer

**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.

written. If you are ready to become a professional



STRUMENTA

ANTLR developer, you can buy our video course to

**Build professional parsers and languages using**

Products

ANTLR

Articles

Contacts

*. Note: the course is taught using Python.*

## APG

APG is a recursive-descent parser using a variation of **Augmented BNF**, that they call Superset Augmented BNF. ABNF is a particular variant of BNF designed to better support bidirectional communications protocol. APG also support additional operators, like syntactic predicates and custom user defined matching functions.



It can generate parsers in C/C++, Java & JavaScript

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



# STRUMENTA

```

1. // example from a tutorial of the author
   of the tool available here
2. // https://www.sitepoint.com/alternative-
   to-regular-expressions/
3. phone-number = ["("] area-code sep
   office-code sep subscriber
4. area-code      = 3digit
   ; 3 digits
5. office-code   = 3digit
   ; 3 digits
6. subscriber    = 4digit
   ; 4 digits
7. sep           = *3(%d32-47 / %d58-126 /
   %d9) ; 0-3 ASCII non-digits
8. digit         = %d48-57
   ; 0-9

```

[About us](#)
[Services](#)
[Products](#)
[Articles](#)
[Contacts](#)

## BYACC/J



All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe  
at any time.

separated by '%%': DECLARATIONS, ACTIONS and



STRUMENTA

CODE. The second one contains the grammar rules

and the third one ~~the custom user code services~~

Products

Articles

Contacts

```

1. // from the documentation
2. %
3. import java.lang.Math;
4. import java.io.*;
5. import java.util.StringTokenizer;
6. %
7.
8. /* YACC Declarations */
9. %token NUM
10. %left '-' '+'
11. %left '*' '/'
12. %left NEG /* negation--unary minus */
13. %right '^' /* exponentiation */
14.
15. /* Grammar follows */
16. %%

```



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

<https://tomassetti.me/parsing-in-java/>

**STRUMENTA**

recursive descent parser. Attributed grammar means that the rules, that are written in an EBNF variant, can be annotated in several ways to change the methods of the generated parser.

The scanner includes support for dealing with things like compiler directives, called pragmas. They can be ignored by the parser and handled by custom code.

The scanner can also be suppressed and substituted with one built by hand.

Technically all the grammars must be LL(1), that is to say the parser must be able to choose the correct rule only looking one symbol ahead. But Coco/R provides several methods to bypass this limitation, including semantic checks, which are basically custom

~~functions that must return a boolean value. The~~

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



STRUMENTA

```

15.
16. TOKENS
17. ident About us letter { letter | digit
18. }.
19. [...]
20. PRODUCTIONS
21. // just a rule is shown
22. IdentList =
23. ident <out int x> (. int n = 1; .)
24. {', ' ident (. n++; .)
25. } (.  

26. Console.WriteLine("n = " + n); .)
27. .
28. "END" ident '!'

```

[Products](#)[Articles](#)[Contacts](#)

Coco/R has a good documentation, with several examples grammars. It supports several languages

including Java, C# and C++

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

compatibility.



It requires Java 7 to generate the parser, but it can run

on earlier versions

[About us](#) ▾

[Services](#)

[Products](#)



**STRUMENTA**

A typical parser defined with annotations will look

[Examples](#)

[Contacts](#)

like this.

```
1. // required import
2. import org.yuanheng.cookcc.*;
3.
4. @CookCCOption (lexerTable = "compressed",
   parserTable = "compressed")
5. // the generated parser class will be a
   parent of the one you define
6. // in this case it will be "Parser"
7. public class Calculator extends Parser
8.
9.     // code
10.
11.    // a lexer rule
```



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



**STRUMENTA**

For the standard of parser generators, using Java annotations it is a ~~peculiar~~ choice. Compared to an alternative like ANTLR there is certainly a less clear division between the grammar and the actions. This could make the parser harder to maintain for complex languages. Also porting to another language could require a complete rewrite.

On the other hand this approach permit to mix grammar rules with the actions to perform when you match them. Furthermore it has the advantage of being integrated in the IDE of your choice, since it is just Java code.

**CUP**

**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.

1. // example from the documentation



```

2. // CUP specification for a simple
   expression evaluator (w/ actions)
3. About us ▾ Services Products
4. import java_cup.runtime.*;
5. Articles Contacts
6. /* Preliminaries to set up and use the
   scanner. */
7. init with {: scanner.init();
8. :};
9. scan with {: return scanner.next_token();
10. :};
11. terminal SEMI, PLUS, MINUS,
   TIMES, DIVIDE, MOD;
12. terminal UMINUS, LPAREN,
   RPAREN;
13. terminal Integer NUMBER;
14.
15. /* Non-terminals */
16. non terminal expr_list, expr_part;

```

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

**Grammatica**



STRUMENTA

“

*Grammatica is a C# and Java parser generator (compiler). It reads a grammar file (in an EBNF format) and creates well-commented and readable C# or Java source code for the parser. It supports LL(k) grammars, automatic error recovery, readable error messages and a clean separation between the grammar and the source code.*

[About us](#)[Services](#)[Products](#)[Articles](#)[Contacts](#)

The description on the Grammatica website is itself a good representation of Grammatica: simple to use, well-documented, with a good amount of features. You can build a listener by subclassing the generated classes, but not a visitor. There is a good reference,

**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.





STRUMENTA

```

13.    WHITESPACE          = <<[
      tnr]>> %ignore%
14.          About us ▾
15.    %productions%
16.          Services
17.    Expression = Term [ExpressionTail] ;
18.
19.    ExpressionTail = "+" Expression
20.          | "-" Expression ;
21.
22.    Term = Factor [TermTail] ;
23.
24.    [..]
25.
26.    Atom = NUMBER
27.          | IDENTIFIER ;

```

Services

Products

Articles

Contacts

**Jacc**

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.





STRUMENTA

Compared to ANTLR the grammar file is much less

clean and include ~~a lot of Java source code~~

Products

```

1.  javacc_options
2.  // "PARSER_BEGIN" "(" <IDENTIFIER> ")"
3.  PARSER_BEGIN(SimpleParser)
4.  public final class SimpleParser { //
   Standard parser class setup...
5.
6.  public static void main(String args[])
{
7.      SimpleParser parser;
8.      java.io.InputStream input;
9.
10. }
11. PARSER_END(SimpleParser)
12.
13. // the rules of the grammar
14. // token rules
15. TOKEN .

```



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

projects, like JavaParser. This has left some quirks in



STRUMENTA

the documentation and usage. For instance,  
technically JavaCC ~~itself does not build an AST, but it Products~~  
comes with a tool that does it, JTree, so for practical  
purposes it does.

Services

Articles

Contacts

There is a grammar repository, but it does not have many grammars in it. It requires Java 5 or later.

## ModelCC

“

*ModelCC is a model-based parser generator that decouples language specification from language processing [...]. ModelCC receives a conceptual model as input, along with constraints that*

×

**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.

authors of the tools, but a limited documentation.

Nonetheless there are examples available, including the following modal:



The modal window displays a calculator interface with a numeric keypad, arithmetic operators (+, -, ×, ÷), and a decimal point (.). The text "Calculator" is visible at the top left of the calculator area.

# S the fo here **STRUMENTA**

## Articles

## Contacts

```
1. public abstract class Expression
2.     implements IModel {
3.     public abstract double eval();
4. }
5. [...]
6.
7. public abstract class UnaryOperator
8.     implements IModel {
9.     public abstract double eval(Expression
10. e);
11. }
12. [...]
```



# All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

Email Address

## **Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

29. }

30

31. [...]



STRUMENTA

About us ▾

Services

Products

Articles

Contacts

## SableCC

SableCC is a parser generator created for a thesis and with the aim to be easy to use and to offer a clean separation between grammar and Java code. Version 3 should also offer an included a ready-to-use way to walk the AST using a visitor. But that is all in theory because there is virtually no documentation and we have no idea how to use any of these things.

Also, a version 4 was started in 2015 and apparently lies abandoned.



All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

3.      Digits ::= '0'..'9';





```

4.    }
5.
6.    token { About us ▾      Services      Products
7.    Space ::= [' ', #8, #9]*;
8.    EOLN ::= [#10, #13];
9.    EOF ::= [#65535];
10.
11.   [...]
12.   Identifier ::= [Letters] [Letters,
13.   Digits]*;
14.
15.   rules {
16.     Variable ::= "var", Identifier;
17.
18.     Element ::= Number | Identifier;
19.
20.     PlusExpression ::= Element, '+',
21.         Expression;
22.     [...]
23.   }

```



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.





STRUMENTA

generated parsers have no runtime

dependency on Canopy itself

[About us](#)[Services](#)[Products](#)It also provides easy access to the parse tree nodes. [Articles](#) [Contacts](#)

A Canopy grammar has the neat feature of using actions annotation to use custom code in the parser.

In practical terms, you just write the name of a function next to a rule and then you implement the function in your source code.

The Java file containing the action code.

```
1.  [...]
2.  import maps.Actions;
3.  [...]
4.
```



**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.





STRUMENTA

generator with support for runtime

grammar rules.

About us ▾

Services

Products

Laja is a code generator and a parser generator and it

Articles

Contacts

is mainly designed to create external DSLs. This means that it have some peculiar features. With Laja you must specify not just the structure of the data, but also how the data should be mapped into Java structures. These structures are usually objects in a hierarchy or flat organization. In short, it makes very easy to parse data files, but it is less suitable for a generic programming language.

Laja options, like output directory or input file, are set in a configuration file.

A Laja grammar is divided in a rules section and the

All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

Email Address

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe at any time.

# Mouse


[About us](#) ▾

[Services](#)
[Products](#)
[Articles](#)
[Contacts](#)

*Mouse is a tool to transcribe PEG into an executable parser written in Java.*

It does not use packrat and thus it uses less memory than the typical PEG parser (the manual explicitly compares Mouse to Rats!).

It does not have a grammar repository, but there are grammars for Java 6-8 and C.

A Mouse grammar is quite clean. To include custom code, a feature called semantic predicates, you do something similar to what you do in Canopy. You include a name in the grammar and then later, in a

**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.

Compiler). It is based on PEG, but it uses “additional



STRUMENTA

expressions and operators necessary for generating actual parsers". It ~~supports~~ left-recursive productions. It can automatically generate an AST.

[Articles](#)[Contacts](#)

It requires Java 6 or later.

The grammar can be quite clean, but you can embed custom code after each production.

```

1. // example from Introduction to the Rats!
   Parser Generator
2. // http://cs.nyu.edu/courses/fall11/CSCI-
   GA.2130-001/rats-intro.pdf
3. /* module intro */
4. module Simple;
5. option parser(SimpleParser);
6.
7. /* productions for syntax analysis */ X
8. public String program = e:expr EOF {
```

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



STRUMENTA

“

*Rekex is a new PEG parser generator for Java 17. It unifies grammar definition and AST construction in the most natural and intuitive way, leading to the simplest approach to writing parsers.*

[About us](#)[Services](#)[Products](#)[Articles](#)[Contacts](#)

Rekex is a new parser generator with a novel approach that flips writing a parser on its head. With traditional parser generators you write a grammar and then the generated parser produces a parse tree. One issue with this approach is that the parse tree is rarely what you want. So, you need to post-process the parse tree to create a data structure that fits your program. This can be a long process in itself.



Particularly if you are dealing with a large grammar

**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.

even a few examples in their repository and





# STRUMENTA

explanation of subprojects used by the library. This is crucial given the ~~new approach of this library~~, so Products that every user can understand if it is a good fit for them.

[Articles](#)
[Contacts](#)

Okay, now we can see what a Rekex parser looks like.

```

1.  [...]
2.
3.  public interface
   ExampleParser_Calculator1
4.  {
5.    // in this example:
6.    // - whitespaces: optional
   Space/Tab between tokens
7.    // - token datatypes are extracted
8.    // - eval() defined inside
   datatypes
9.

```



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



**STRUMENTA**<sup>1</sup>

```

28. -----
29.          About us ▾      Services      Products
30.      enum Op1
31.      {
32.          @Word("+) plus,
33.          @Word("-") minus;
34.
35.          public int eval(int x, int y){
36.              return this==plus ? x+y : x-
37.                  y;
38.          }
39.      enum Op2
40.      {
41.          @Word({"/", "\u00f7"}) div,
42.          @Word({"*", "\u00d7"}) mul,
43.          @Word("") mul_implicit; // a b ==
44.                  a * b
45.          public int eval(int x, int y){
46.              return this==div ? x/y : x*yX
47.          }

```

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

```

66.          // [...] other records with eval
               methods

```



# STRUMENTA<sup>9</sup>

```

67.
68.      // teAbout us Services Products
-----  

STRUMENTA9.Articles Contacts
70.      public static PegParser<Input>
    parser()
71.      {
72.          return PegParser.of(Input.class);
73.      }
74.      public static Function<Input, Integer>
    eval()
75.      {
76.          return (input) -> input.e0.eval();
77.      }
78.  }

```

In practice, you write the parser using a well-defined structure and conventions (e.g., eval classes) in Java.  
In other words, the grammar is Java (17) code but a 

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



STRUMENTA

In practice this means that they are very useful for all the little parsing problems you find. If the typical developer encounters a problem, that is too complex for a simple regular expression, these libraries are usually the solution. In short, if you need to build a parser, but you do not actually want to, a parser combinator may be your best option.

## Jparsec

Jparsec is the port of the parsec library of Haskell.

Parser combinators are usually used in one phase, that is to say they are without lexer. This is simply because it can quickly become too complex to

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.



STRUMENTA

```

OPERATORS =
8.      Terminals.operators("+", "-", "*",
9.      "/", "(", About us ▾ Services Products
10.     [...] Articles Contacts
11.
12.     static final Parser<?> TOKENIZER =
13.
14.     Parsers.or(Terminals.DecimalLiteral.TOKENI
15.     ZER, OPERATORS.tokenizer());
16.
17.     static Parser<Double>
18.     calculator(Parser<Double> atom) {
19.         Parser.Reference<Double> ref =
20.             Parser.newReference();
21.         Parser<Double> unit =
22.             ref.lazy().between(term("("),
23.             term(")").or(atom);
24.         Parser<Double> parser = new
25.             OperatorTable<Double>()
26.                 .infixl(op("+", (l, r) -> l + r), X

```

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

beginning of 2017).



# Parboiled

About us ▾

Services

Products

*Parboiled provides a recursive descent PEG parser implementation that operates on PEG rules you specify.*

Articles

Contacts

The objective of parboiled is to provide an easy to use and understand way to create small DSLs in Java. It puts itself in the space between a simple bunch of regular expressions and an industrial-strength parser generator like ANTLR. A parboiled grammar can include actions with custom code, included directly into the grammar code or through an interface.



**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.



STRUMENTA

```

21.      [...]
22.      public Rule OperatorRule(Rule subRule,
23.          Rule operatorRule) {
24.          Var<Character> op = new
25.              Var<Character>();
26.          return Sequence(
27.              subRule,
28.              ZeroOrMore(
29.                  operatorRule,
30.                  op.set(matchedChar())),
31.              subRule,
32.              push(new
33.                  CalcNode(op.get(), pop(1), pop())))
34.          );
35.      }
36.      [...]
37.      public Rule Number() {
38.          return Sequence(

```

[Products](#)[Articles](#)[Contacts](#)

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.





Parboiled works a bit like a cross between a parser combinator and a parser generator. You create rules in code, with ready-to-use methods like Sequence or

Optional, just like a parser combinator. However, the end result is a parser class that you are supposed to use like a generated parser. Parboiled is not suited to create individually used rules, i.e., to parse bits and pieces, the way a parser combinator can. You use it to parse a coherent language.

It does not build an AST for you, but it provides a parse tree and some classes to make it easier to build it. That is because its authors maintain that the AST is heavily dependent on your exact project needs, so they prefer to offer an “open and flexible approach”. It

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.



STRUMENTA

*combinators, parsing expression**grammars and packrat parsers to*[About us](#)[Services](#)[Products](#)*model grammars and parsers as**objects that can be reconfigured*[Articles](#)[Contacts](#)*dynamically.*

PetitParser is also between a parser combinator and a traditional parser generator. All the information is written in the source code, but the source code is divided in two files. In one file you define the grammar, while in the other one you define the actions corresponding to the various elements. The idea is that it should allow you to dynamically redefine grammars. While it is smartly engineered, it is debatable if it is also smartly designed.



Given that departs from the usual design of a parser

**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.



STRUMENTA

```

8.
9.     public JsonGrammarDefinition() {
10.       def("start", ref("value").end());
11.       Services
12.       def("array", of('[').trim());
13.           Articles
14.           .seq(ref("elements")).optional());
15.           .seq(of(']').trim()));
16.       def("elements",
17.           ref("value").separatedBy(of(',').trim())));
18.       def("members",
19.           ref("pair").separatedBy(of(',').trim()));
20.       def("trueToken",
21.           of("true").flatten().trim());
22.       def("falseToken",
23.           of("false").flatten().trim());
24.       def("nullToken",
25.           of("null").flatten().trim());
26.       def("stringToken",
27.           ref("stringPrimitive").flatten().trim()); X
28.       def("numberToken",

```

Products

Contacts

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

Functions.withoutSeparators());

9. action("members",





STRUMENTA

```

10.      action("array", new
Function<List<?>, List<?>> Services
11.      @Override
12.      public List<?> apply(List<?> Articles?>
input) {
13.          return input.get(1) != null ?
input.get(1) : new ArrayList<>();
14.      }
15.  });
16.
17.  [...]
18. }
19. }
```

Products

Contacts

There is a version written in Java, but there are also versions in Smalltalk, Dart, PHP, and TypeScript.

The documentation is lacking, but there are ~~example grammars available~~

**All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.**

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe at any time.



## STRUMENTA

wants to provide a simple internal Domain Specific

Language to express grammar languages

Products

as grammar languages

services

The library is simple, but offers everything necessary

to create simple parsers.

```
1. Parser<Chr, Integer> sum =  
    intr.andL(chr('+')).and(intr).map(Integer:  
    :sum);
```

In this example from the documentation you can see how it is possible to combine the parsers for integers (intr) and the one for characters (chr) to parse a simple sum expression. The expression is also evaluated using the map function to call the normal

All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

recommended by **Danny van Bruggen**, the maintainer

of JavaParser, so you know it is good.



STRUMENTA

About us ▾

Services

Products

Articles

Contacts

# Java Libraries That Parse Java: JavaParser

There is one special case that requires a specific comment: the case in which you want to parse Java code in Java. In this case we suggest using a library named JavaParser. Incidentally we heavily contribute to JavaParser, but this is not the only reason why we suggest it. The fact is that JavaParser is a project with tens of contributors and thousands of users, so it is pretty robust.

All the major parsing libraries and tools reviewed for Java, Python, C# and JavaScript.

[Sign Up and Get the Guide](#)

We respect your privacy. Unsubscribe at any time.

# Summary



STRUMENTA

About us ▾

Services

Products

Articles

Contacts

Parsing in Java is a broad topic and the world of parsers is a bit different from the usual world of programmers. You will find the best tools coming directly from academia, which is typically not the case with software. Some tools and libraries have been started for a thesis or a research project. The upside is that tools tend to be easily and freely available. The downside is that some authors prefer to have a good explanation of the theory behind what their tools do, rather than a good documentation on how to use them. Also, some tools end up being abandoned as the original authors finish their master or their PhD.



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

most technically correct solution might not be ideal in





STRUMENTA

real life with all its constraints. But we have searched

and tried many similar tools in our work and

something like this article would have helped us save

some time. So we wanted to share what we have

learned on the best options for parsing in Java.

*We would like to thank Stefan Czaska for having informed us of AnnoFlex.*

*We would like to thank Danny van Bruggen for having informed us of funcj.*

*We would like to thank Zhong Yu for having informed us of rekex.*



All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe  
at any time.

Spring



STRUMENTA

Parsing

Software Development

Boot

About us

Migrator

Services

14 September 2022

Products

27 September 2022

Articles

Contacts

**STRUMENTA**

Tools to solve complex problems

Privacy Policy

P.IVA 11817320010

Company Information

Strumenta

Strumenta

Federico Tomassetti

Federico Tomassetti



**All the major  
parsing libraries  
and tools  
reviewed for Java,  
Python, C# and  
JavaScript.**

Email Address

**Sign Up and Get the Guide**

We respect your privacy. Unsubscribe

at any time.

