

04834580 Software Engineering (Honor Track) 2024-25

Parsing

Sergey Mechtaev

mechtaev@pku.edu.cn

School of Computer Science, Peking University



Definition ([1])

Parsing is a process of analyzing a string of symbols, either in programming languages or data structures, conforming to the rules of a formal grammar by breaking it into parts.

The context-free grammar [2] describing arithmetic expressions with addition, subtraction, and parentheses is as follows:

$$\begin{array}{l} E \rightarrow E + T \\ \quad | E - T \\ \quad | T \\ T \rightarrow (E) \\ \quad | \text{id} \end{array}$$

Here:

- ▶ E : Represents an expression.
- ▶ T : Represents a term.
- ▶ id : Represents an identifier, which can be a number or variable.
- ▶ The operators $+$ and $-$ denote addition and subtraction respectively.
- ▶ Parentheses (E) group sub-expressions for precedence.

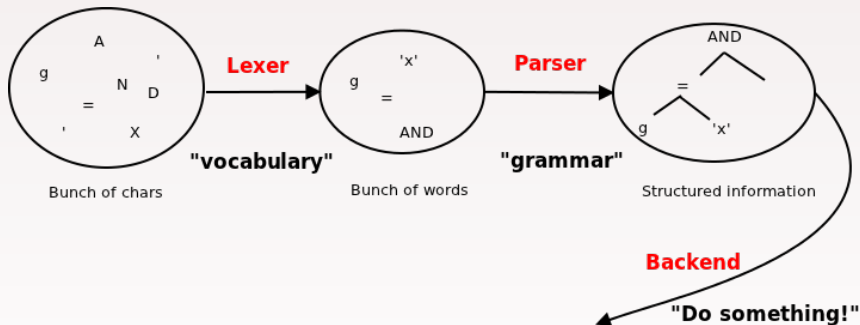
| |
|--|
| $\langle \text{Stmt} \rangle \rightarrow \langle \text{Id} \rangle = \langle \text{RExpr} \rangle ;$ |
| $\langle \text{Stmt} \rangle \rightarrow \{ \langle \text{StmtList} \rangle \}$ |
| $\langle \text{Stmt} \rangle \rightarrow \text{if} (\langle \text{RExpr} \rangle) \langle \text{Stmt} \rangle$ |
| $\langle \text{StmtList} \rangle \rightarrow \langle \text{Stmt} \rangle$ |
| $\langle \text{StmtList} \rangle \rightarrow \langle \text{StmtList} \rangle \langle \text{Stmt} \rangle$ |
| $\langle \text{RExpr} \rangle \rightarrow \langle \text{RExpr} \rangle > \langle \text{AExpr} \rangle$ |
| $\langle \text{RExpr} \rangle \rightarrow \langle \text{RExpr} \rangle < \langle \text{AExpr} \rangle$ |
| $\langle \text{RExpr} \rangle \rightarrow \langle \text{RExpr} \rangle \geq \langle \text{AExpr} \rangle$ |
| $\langle \text{RExpr} \rangle \rightarrow \langle \text{RExpr} \rangle \leq \langle \text{AExpr} \rangle$ |
| $\langle \text{RExpr} \rangle \rightarrow \langle \text{AExpr} \rangle$ |
| $\langle \text{AExpr} \rangle \rightarrow \langle \text{AExpr} \rangle + \langle \text{PExpr} \rangle$ |
| $\langle \text{AExpr} \rangle \rightarrow \langle \text{AExpr} \rangle - \langle \text{PExpr} \rangle$ |
| $\langle \text{AExpr} \rangle \rightarrow \langle \text{PExpr} \rangle$ |
| $\langle \text{PExpr} \rangle \rightarrow \langle \text{Id} \rangle$ |
| $\langle \text{PExpr} \rangle \rightarrow \langle \text{Num} \rangle$ |
| $\langle \text{Id} \rangle \rightarrow x$ |
| $\langle \text{Id} \rangle \rightarrow y$ |
| $\langle \text{Num} \rangle \rightarrow 0$ |
| $\langle \text{Num} \rangle \rightarrow 1$ |
| $\langle \text{Num} \rangle \rightarrow 9$ |

| $\langle \text{Stmt} \rangle$ | |
|---|--|
| $\text{if} (\langle \text{RExpr} \rangle)$ | $\langle \text{Stmt} \rangle$ |
| $\text{if} (\langle \text{RExpr} \rangle > \langle \text{AExpr} \rangle)$ | $\langle \text{Stmt} \rangle$ |
| $\text{if} (\langle \text{AExpr} \rangle > \langle \text{AExpr} \rangle)$ | $\langle \text{Stmt} \rangle$ |
| $\text{if} (\langle \text{PExpr} \rangle > \langle \text{PExpr} \rangle)$ | $\langle \text{Stmt} \rangle$ |
| $\text{if} (\langle \text{Id} \rangle > \langle \text{Num} \rangle)$ | $\langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ \langle \text{StmtList} \rangle \}$ |
| $\text{if} (x > 9)$ | $\{ \langle \text{StmtList} \rangle \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ \langle \text{Stmt} \rangle \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ \langle \text{Id} \rangle = \langle \text{RExpr} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = \langle \text{AExpr} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = \langle \text{PExpr} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = \langle \text{Num} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = 0 ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = 0 ; \langle \text{Id} \rangle = \langle \text{RExpr} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = 0 ; y = \langle \text{AExpr} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = 0 ; y = \langle \text{AExpr} \rangle + \langle \text{PExpr} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = 0 ; y = \langle \text{PExpr} \rangle + \langle \text{PExpr} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = 0 ; y = \langle \text{Id} \rangle + \langle \text{Num} \rangle ; \langle \text{Stmt} \rangle$ |
| $\text{if} (x > 9)$ | $\{ x = 0 ; y = y + 1 ; \}$ |

The Backus–Naur Form (BNF) [3] is a notation system for defining the syntax of programming languages:

$$\begin{aligned} E &::= T (('+' \mid '-') T)^* \\ T &::= '(' E ')' \mid \text{id} \end{aligned}$$

- ▶ $E ::= T (('+' \mid '-') T)^*$ an expression E consists of a term T optionally followed by zero or more sequences of addition “+” or subtraction “-” operators, each paired with another term T . The “*” indicates repetition (zero or more occurrences).
- ▶ $T ::= '(' E ')' \mid \text{id}$ means a term T can either be an expression E enclosed in parentheses to group subexpressions (E), or an identifier id , which represents variables or literals like numbers.



A lexer represents a stream of tokens:

```
import re

class Lexer:
    def __init__(self, input_text):
        self.tokens = re.findall(r'\d+|[(\)+\~]', input_text)
        self.position = 0

    def get_next_token(self):
        if self.position < len(self.tokens):
            token = self.tokens[self.position]
            self.position += 1
            return token
        return None

    def peek(self):
        if self.position < len(self.tokens):
            return self.tokens[self.position]
        return None
```

Definition ([4])

A kind of top-down parser built from a set of mutually recursive procedures where each such procedure implements one of the nonterminals of the grammar.


```
class Parser:
    def __init__(self, lexer):
        self.lexer = lexer
        self.current_token = self.lexer.get_next_token()

    def consume(self):
        self.current_token = self.lexer.get_next_token()

    def parse_E(self):
        node = self.parse_T()
        while self.current_token in ('+', '-'):
            op = self.current_token
            self.consume()
            node = (op, node, self.parse_T())
        return node

    ...
```

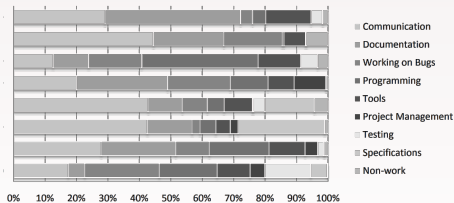
```
...

def parse_T(self):
    if self.current_token == '(':
        self.consume() # Consume '('
        node = self.parse_E()
        if self.current_token == ')':
            self.consume() # Consume ')'
            return node
        else:
            raise SyntaxError("Missing closing parenthesis")
    elif self.current_token.isdigit():
        node = int(self.current_token) # Convert id (number) to integer
        self.consume()
        return node
    else:
        raise SyntaxError(f"Unexpected token: {self.current_token}")
```

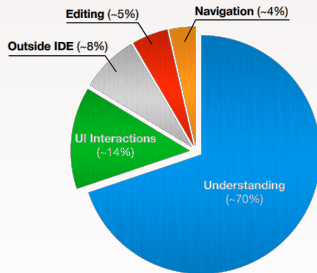
Client code:

```
if __name__ == "__main__":  
    input_text = "1 + (2 - 3)"  
    lexer = Lexer(input_text)  
    parser = Parser(lexer)  
    syntax_tree = parser.parse_E()  
    print(syntax_tree)
```

2008's study with 8 new Microsoft developers showed that programming occupies only 10–20% of their time. [5]



2015's study with 18 developers showed that they spend 70% of their time on understanding. [6]



Parser generators such as generate parsers based on grammar descriptions such as BNF. Popular generations include:

- ▶ Bison for C/C++
- ▶ ANTRL for Java and other languages
- ▶ Lark for Python

```
grammar ExpressionGrammar;
```

```
// parser
```

```
expr  : left=expr op=('*' | '/' ) right=expr #opExpr
      | left=expr op=('+' | '-' ) right=expr #opExpr
      | '(' expr ')'                        #parenExpr
      | atom=INT                            #atomExpr
      ;
```

```
// lexer
```

```
INT : ('0' .. '9') +;
```

```
WS : [ \r\n\t ] + -> skip ;
```

```
?start: product
    | start "+" product -> add
    | start "-" product -> sub
?product: atom
    | product "*" atom -> mul
    | product "/" atom -> div
?atom: NUMBER -> number
    | "-" atom -> neg
    | "(" start ")"
%import common.NUMBER
%import common.WS_INLINE
%ignore WS_INLINE
```

Definition

A parser combinator is a higher-order function that accepts several parsers as input and returns a new parser as its output.

- [1] Wikipedia Authors.
Parsing.
<https://en.wikipedia.org/wiki/Parsing>, 2025.
- [2] John E Hopcroft, Rajeev Motwani, and Jeffrey D Ullman.
Introduction to automata theory, languages, and computation.
Acm Sigact News, 32(1):60–65, 2001.
- [3] John W Backus.
The syntax and the semantics of the proposed international algebraic language of the zurich acm-gamm conference.
In *ICIP Proceedings*, pages 125–132, 1959.
- [4] Wikipedia Authors.
Recursive descent parser.
https://en.wikipedia.org/wiki/Recursive_descent_parser, 2025.

- [5] Andrew Begel and Beth Simon.
Novice software developers, all over again.
In International Workshop on Computing Education Research, pages 3–14, 2008.
- [6] Roberto Minelli, Andrea Mocci, and Michele Lanza.
I know what you did last summer-an investigation of how developers spend their time.
In International Conference on Program Comprehension, pages 25–35. IEEE, 2015.