

# Context-Free Grammar (CFG)

A *context-free grammar* looks like this bunch of rules:

$$E \rightarrow E + E$$

$$E \rightarrow M$$

$$M \rightarrow M \times M$$

$$M \rightarrow A$$

$$A \rightarrow 0$$

$$A \rightarrow 1$$

$$A \rightarrow (E)$$

Main idea:

- ▶  $E, M, A$  are *non-terminal symbols* aka *variables*. When you see them, you apply rules to expand. One of them is designated as the *start symbol*. You always start from it. I will designate  $E$  as the start symbol.
- ▶  $+, \times, 0, 1, (, )$  are *terminal symbols*. They are the characters you want in your language.

## Derivation (aka Generation)

*Derivation* is a finite sequence of applying the rules until all non-terminal symbols are gone. Often aim for a specific final string.

$$\begin{array}{ll} E \rightarrow M & \rightarrow 1 \times (M + E) \\ \rightarrow M \times M & \rightarrow 1 \times (A + E) \\ \rightarrow A \times M & \rightarrow 1 \times (0 + E) \\ \rightarrow 1 \times M & \rightarrow 1 \times (0 + M) \\ \rightarrow 1 \times A & \rightarrow 1 \times (0 + M \times M) \\ \rightarrow 1 \times (E) & \rightarrow 1 \times (0 + A \times M) \\ \rightarrow 1 \times (E + E) & \rightarrow 1 \times (0 + 1 \times A) \\ & \rightarrow 1 \times (0 + 1 \times 1) \end{array}$$

Context-free grammars can support: matching parentheses, **unlimited nesting**.

## Backus-Naur Form (BNF)

Backus-Naur Form is a computerized, practical notation for CFGs.

- ▶ Surround non-terminal symbols by  $\langle \rangle$ ; allow multi-letter names.
- ▶ Merge rules with the same LHS.
- ▶ (Some versions.) Surround terminal strings by single or double quotes.
- ▶ Use  $::=$  for  $\rightarrow$ .

Our example grammar in BNF:

```
 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle "+" \langle \text{expr} \rangle \mid \langle \text{mul} \rangle$   
 $\langle \text{mul} \rangle ::= \langle \text{mul} \rangle "*" \langle \text{mul} \rangle \mid \langle \text{atom} \rangle$   
 $\langle \text{atom} \rangle ::= "0" \mid "1" \mid "(" \langle \text{expr} \rangle ")"$ 
```

## Extended Backus-Naur Form (EBNF)

- ▶ {...} for 0 or more occurrences.
- ▶ [...] for 0 or 1 occurrences.
- ▶ Some versions: No <> needed around non-terminal symbols.

Example: Lisp/Scheme S-expression<sup>1</sup> grammar (basic):

In BNF:

```
<s-expr> ::= <identifier> | "(" <s-exprs> ")"  
<s-exprs> ::= <s-expr> | <s-expr> <s-exprs>
```

In EBNF:

```
<s-expr> ::= <identifier>  
           | "(" <s-expr> { <s-expr> } ")"
```

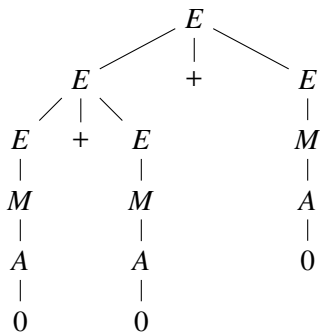
So you need fewer artificial non-terminals and rules that merely mean “at least 0 of this”, “at least 1 of that”, etc.

---

<sup>1</sup>“symbolic expression”

## Parse Tree aka Derivation Tree

A *parse tree* aka *derivation tree* presents a derivation with more structure (tree), less repetition.



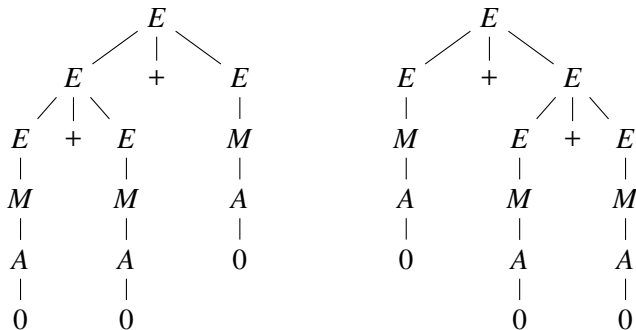
This example generates  $0 + 0 + 0$ .

## Parse Trees: General Points

- ▶ Internal nodes are non-terminal symbols.
- ▶ Both operators and operands are terminal symbols at leaves.
- ▶ Whole string recorded, just scattered.
- ▶ Purpose: Present derivation, help visualize derivation and grammar.

# Ambiguous Grammar

Two different trees generate the same  $0 + 0 + 0$ :



If this happens, the grammar is *ambiguous*.

We try to design unambiguous grammars.

(Bad news: CFG ambiguity is undecidable.)

## Unambiguous Grammar Example

An unambiguous grammar that generates the same language as our ambiguous grammar example:

```
<expr> ::= <expr> "+" <mul> | <mul>  
<mul> ::= <mul> "*" <atom> | <atom>  
<atom> ::= "0" | "1" | "(" <expr> ")"
```

Exercise: Find the parse trees for  $0 + 0 + 0$  and  $0 \times 0 \times 0$ . Observe that you are forced only one answer, and it's left-leaning.

(Bad news: Equivalence of two CFGs is also undecidable.)



## Left Recursive vs Right Recursive

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle "+" \langle \text{mul} \rangle$

That is a *left recursive* rule. The recursion is at the beginning (left).

$\langle \text{expr} \rangle ::= \langle \text{mul} \rangle "+" \langle \text{expr} \rangle$

That is a *right recursive* rule. The recursion is at the end (right).

Sometimes they convey intentions of left association or right association. But not always.

They affect some parsing algorithms.

# Recursive Descent Parsing

*Recursive descent parsing* is a simple strategy for writing a parser.

- ▶ Write a procedure for each rule.
- ▶ Non-terminals on RHS become procedure calls, possibly recursive calls. (Thus “recursive descent”, also “top-down”.)
- ▶ Left recursion needs special treatment to avoid infinite loops.
- ▶ Terminal symbols: Consume input and check.
- ▶ Alternatives require lookahead or backtracking.

Some options for handling left recursion:

- ▶ Re-design grammar to not have left recursion.
- ▶ Many left recursive rules just express left-associating operators. Can be done without left recursive code.

## Recursive Descent Parser Example

Example grammar suitable for recursive descent parsing:

```
<sub> ::= <atom> "-" <sub> | <atom>  
<atom> ::= "0" | "1" | "(" <sub> ")"
```

Pseudo-code of recursive descent parser:

sub:

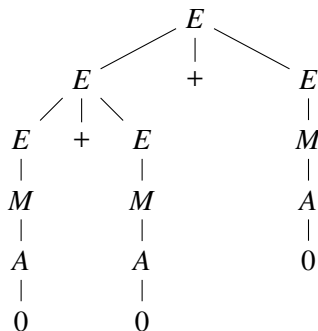
```
  try (atom;  
      read; if not "-" then fail;  
      sub;)  
  if that failed: atom;
```

atom:

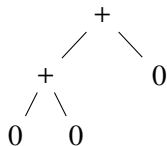
```
  read;  
  if "0" or "1": success;  
  if "(": sub;  
      read; if not ")" then fail;  
  else: fail;
```

# Abstract Syntax Tree (AST) (vs Parse Tree)

Parse tree:



*Abstract syntax tree:*



## Abstract Syntax Tree: General Points

- ▶ Internal nodes are operators/constructs.  
Example construct: if-then-else.
- ▶ Non-terminal symbols gone or replaced by constructs.
- ▶ Many terminal symbols gone too if they play no role other than nice syntax (e.g., spaces, parentheses, punctuations).  
Those bearing content, replaced by appropriate representations, not stay as characters.  
E.g., Character '+' replaced by a data constructor, character '0' replaced by number 0.
- ▶ Purpose: Present essential structure and content, ready for interpreting, compiling, analyses.
- ▶ Parsers usually output abstract syntax trees when successful.

# Lexical Analysis aka Tokenization

In principle: Grammar and parser can work on characters directly.  
But sometimes messy.

In practice: two stages:

1. Chop into chunks and classify into *lexemes* aka *tokens*, discard spaces, e.g.,

```
"(xa * xb)**25"
```

↳

```
[Open, Var "xa", Op Mul, Var "xb", Close,  
Op Exp, NumLiteral 25]
```

*Lexical Analysis, tokenization. Lexers/tokenizers need only regular expressions.*

2. Parsing based on CFG, but terminal symbols are tokens, not characters.