

Syntactic Analysis (Parsing)

José Miguel Rivero

`rivero@cs.upc.edu`

Barcelona School of Informatics (FIB)

Universitat Politècnica de Catalunya BarcelonaTech (UPC)

Summary

- Objectives of Syntax Analysis
- Context Free Grammars. Applications
- Parsing in Compilers / Interpreters
- Syntax vs. Semantics
- Derivation. Parse Tree
- Cocke-Younger-Kasami Parsing Algorithm
- Parse Tree and Abstract Syntax Tree (AST)
- Ambiguous Grammars
- Linear Parsing Algorithms
 - Top-down $LL(1)$ parsers
 - Bottom-up $LR(1)$ parsers

Objectives of Syntax Analysis

- Given a CFG (Context Free Grammar) G , with start symbol S , and a word w (a sequence of tokens), can w be generated by G ?

$$w \in \mathcal{L}(G) ? \quad S \Rightarrow^* w ?$$

- Analyze the sequence of tokens to determine their grammatical structure with respect to G .
Compute the *parse tree* corresponding to the input
- Detect, diagnose, and recovery from *syntax errors*
- Accepts some invalid constructs, filtered out by the semantic analysis

Context Free Grammars

● Expressive power

- $\mathcal{L}_1 = \{ a^n b^n \mid n \geq 0 \}$

- $\mathcal{L}_2 = \{ w c w^{-1} \mid w \in (a|b)^* \}$

- $\mathcal{L}_3 = \{ w c w \mid w \in (a|b)^* \}$

- $\mathcal{L}_4 = \{ a^n b^m c^n d^m \mid n, m \geq 0 \}$

- Algebraic expressions involving numbers, operations $+$ and $*$, and left and right parentheses

- Constructions of programming languages such as declarations, statements (assignment, `if`, `while`, ...), expressions, etc.

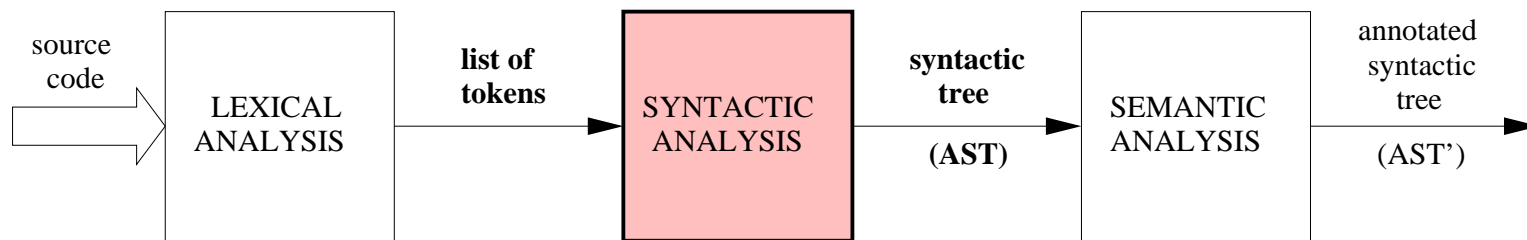
Context Free Grammars

● Applications

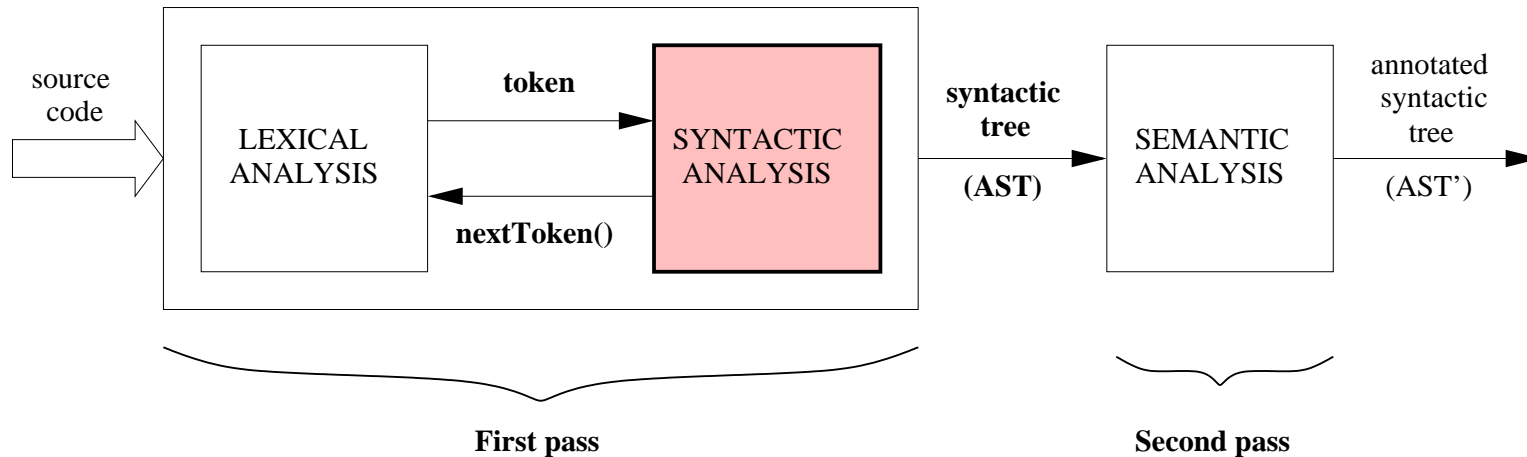
- Natural language processing (NLP)
- Type-setting languages (nroff, postscript), document and extensible markup languages (\LaTeX , SGML, XML, ...)
- Query for databases and information systems (SQL, DataLog, LDAP, ...)
- Logical synthesis and simulation of electronic circuits (VHDL)
- Graphical or modelling languages (UML, Energy Systems Language)
- **Many applications come with their built-in language.**
For example, the scripting language for Adobe Flash (ActionScript)

Parsing in Compilers / Interpreters

Conceptual structure:

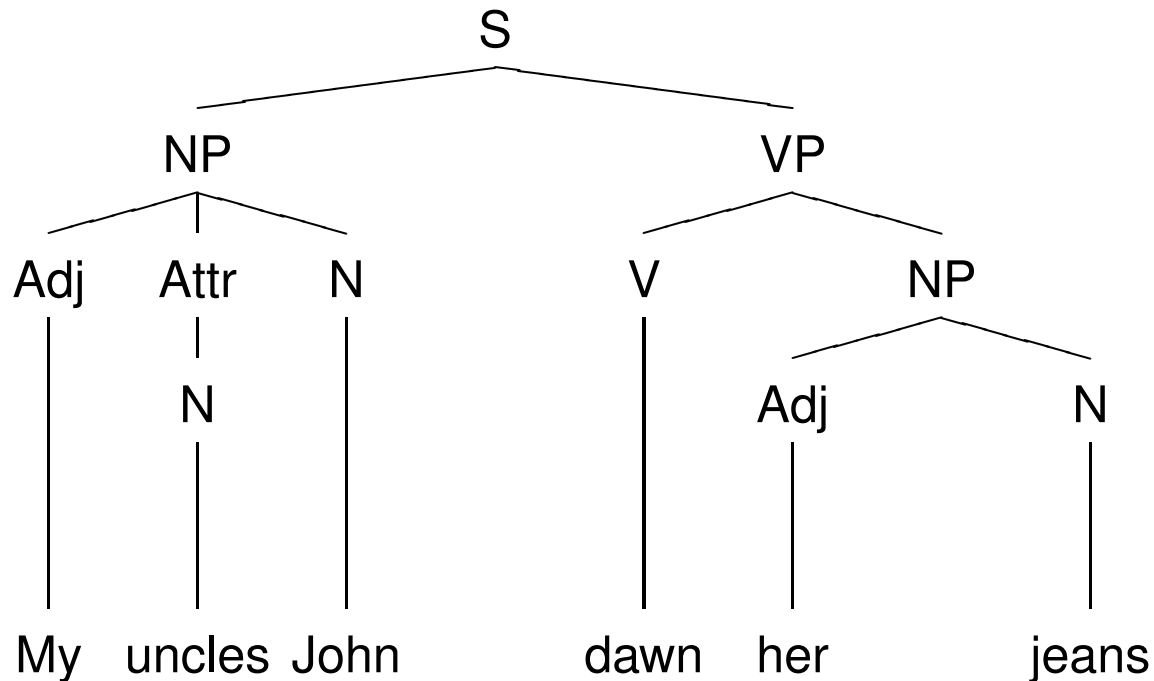


Usual structure:



Syntax vs. Semantics

- “My uncles John dawn her jeans”



- “Marketplace for : input the a .”

- “They are hunting dogs”

Syntax vs. Semantics

- Programming languages.

Syntax errors:

```
public class {
    int static i;
    boolean j;
    public double i(UnknownClass k {
        i..x = "hello";
        if (i / 2)
            i + 1 = j;
        return i > ;
    }
}
```


Syntax vs. Semantics

- Programming languages.

Semantic errors:

```
public class MyFirstClass {
    static int i;
    boolean j;
    public double i(UnknownClass k) {
        i.x[10] = "hello";
        if (i / 2)
            i + 1 = j;
        return i > 3.14;
    }
}
```

Syntax vs. Semantics

- Programming languages.

Ambiguity:

```
class Bar {  
    public float foo( ) {  
        return 1;  
    }  
  
    public int foo( ) {  
        return 2;  
    }  
};
```

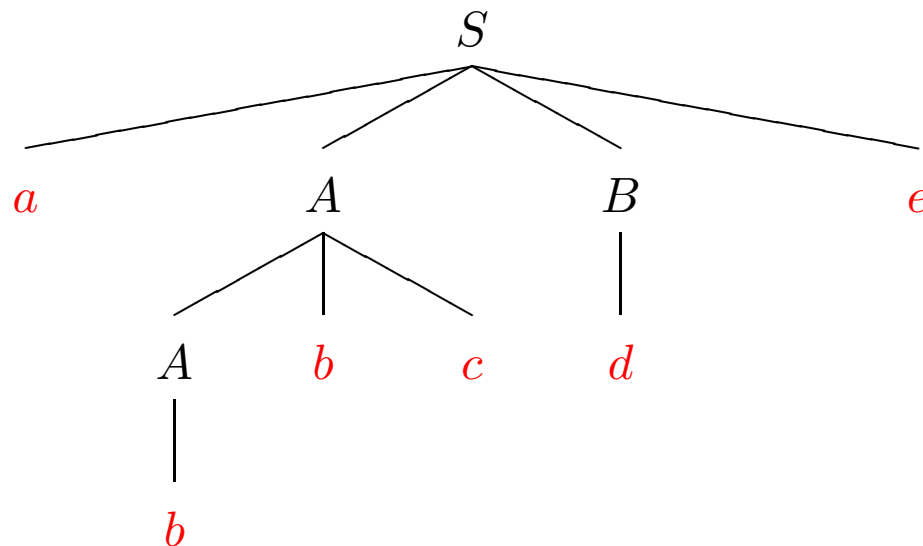
```
float x = Bar::foo( );
```

Derivation. Parse Tree

$$\text{Grammar } G: \begin{cases} S \rightarrow aABe \\ A \rightarrow Abc \mid b \\ B \rightarrow d \end{cases} \quad w = abcde$$

$$S \Rightarrow aABe \Rightarrow aAbcBe \Rightarrow abcBe \Rightarrow abcde = w$$

Parse tree:



Cocke-Younger-Kasami Algorithm

Input : A CFG in *Chomsky Normal Form* G , and a word $w = a_1 \dots a_n$

Only rules of the form $A \rightarrow a$ or $A \rightarrow BC$

Output : An $n \times n$ table T , where $T[i, l]$ is the set of nonterminals A that generate the substring $a_i \dots a_{i+l-1}$

Example : $w = a_1 a_2 \dots a_8$

	1	2	3	4	1=5	6	7	8
1								
2								.
i = 3		B			A		.	.
4						.	.	.
5			C	
6			
7		
8	

$A \rightarrow BC \in G$

$B \Rightarrow^* a_3 a_4$

$C \Rightarrow^* a_5 a_6 a_7$

$A \Rightarrow BC \Rightarrow^* a_3 a_4 a_5 a_6 a_7$

Cocke-Younger-Kasami Algorithm

Input : A CFG in *Chomsky Normal Form* G , and a word $w = a_1 \dots a_n$

Only rules of the form $A \rightarrow a$ or $A \rightarrow BC$

Output : An $n \times n$ table T , where $T[i, l]$ is the set of nonterminals A that generate the substring $a_i \dots a_{i+l-1}$

Algorithm :

for all $i \in [1..n]$: $T[i, 1] := \{A \mid A \rightarrow a_i \in G\}$

for all $l \in [2..n]$:

for all $i \in [1..n-l]$:

for every $A \rightarrow BC \in G$:

if there is a $l_1 \in [1..l-1]$ such that

$B \in T[i, l_1]$ and $C \in T[i+l_1, l-l_1]$

then

$T[i, l] := T[i, l] \cup \{A\}$

Using dynamic programming, each table entry can be filled in $O(n)$ time.

The algorithm runs in $O(n^3)$ time.

Parsing expressions

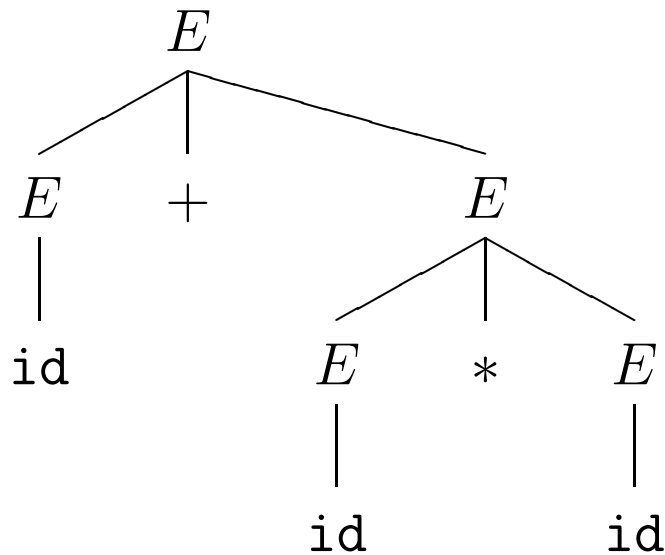
$G : E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$

$w = \text{id}_1 + \text{id}_2 * \text{id}_3$

(Leftmost) Derivation #1 :

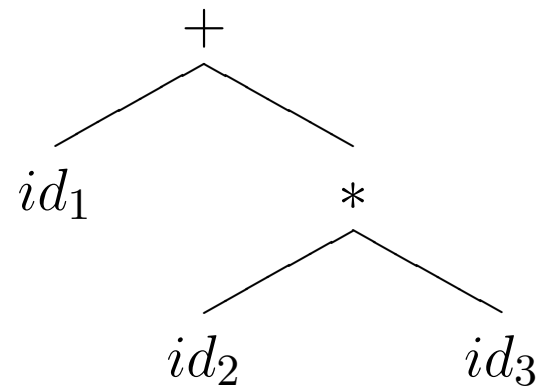
$E \Rightarrow_{ld} E + E \Rightarrow_{ld} \text{id} + E \Rightarrow_{ld} \text{id} + E * E \Rightarrow_{ld} \text{id} + \text{id} * E \Rightarrow_{ld} \text{id} + \text{id} * \text{id}$

Parse Tree:



$\text{id}_1 + (\text{id}_2 * \text{id}_3)$

Abstract Syntax Tree (AST):



Parsing expressions

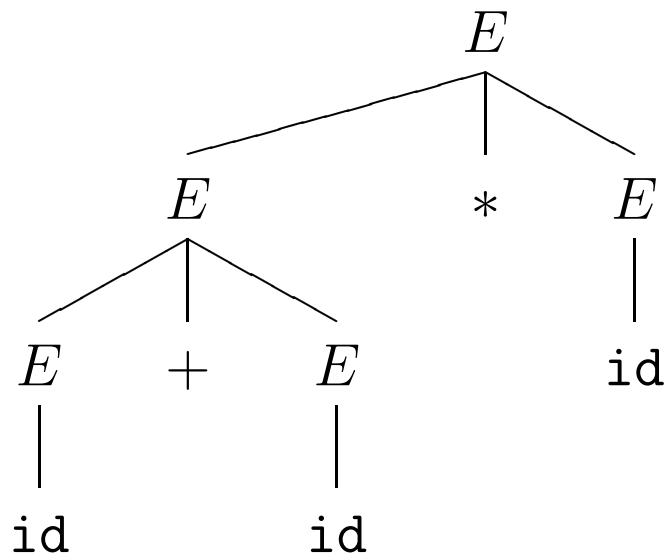
$G : E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$

$w = \text{id}_1 + \text{id}_2 * \text{id}_3$

(Leftmost) Derivation #2 :

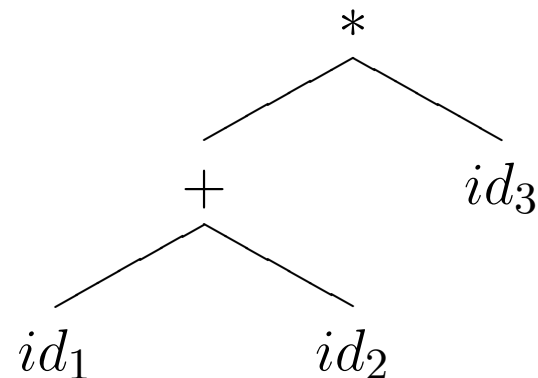
$E \Rightarrow_{ld} E * E \Rightarrow_{ld} E + E * E \Rightarrow_{ld} \text{id} + E * E \Rightarrow_{ld} \text{id} + \text{id} * E \Rightarrow_{ld} \text{id} + \text{id} * \text{id}$

Parse Tree:



$(\text{id}_1 + \text{id}_2) * \text{id}_3$

Abstract Syntax Tree (AST):



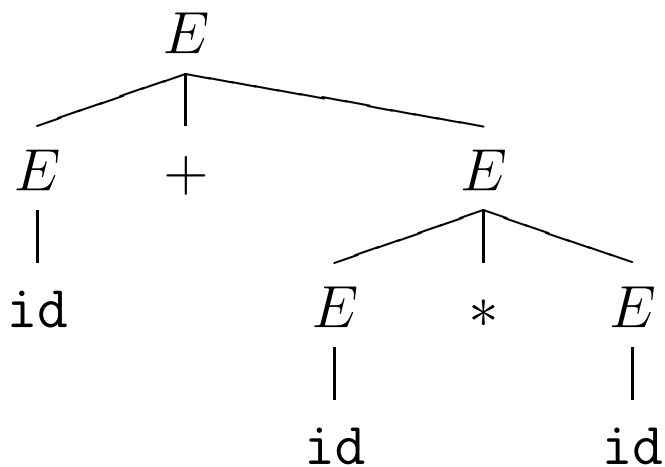
Ambiguous grammar

A grammar G is said to be *ambiguous* if there is some string w that has more than one parse tree or more than one leftmost derivation.

Ambiguous grammar

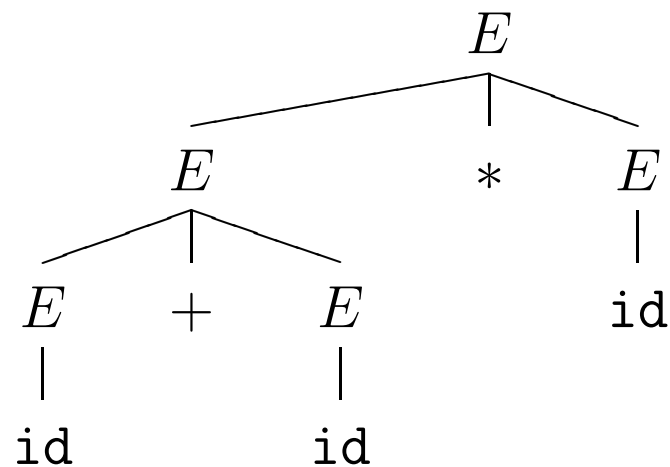
$G : E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$
 $w = \text{id}_1 + \text{id}_2 * \text{id}_3$

Parse Tree A:



$\text{id}_1 + (\text{id}_2 * \text{id}_3)$

Parse Tree B:



$(\text{id}_1 + \text{id}_2) * \text{id}_3$

Operator's precedence

• $id_1 + (id_2 * id_3 * id_4) + id_5$

• Precedence: $+ \prec_p *$

• This other grammar G' reflects the operator's precedence:

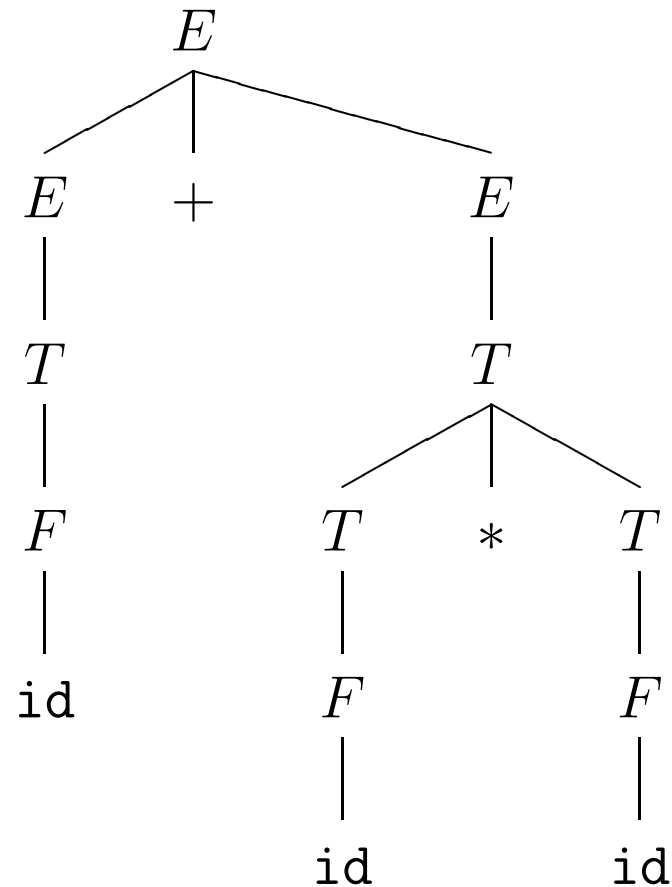
$$E \rightarrow E + E \mid T$$

$$T \rightarrow T * T \mid F$$

$$F \rightarrow (E) \mid id$$

• Now $w = id_1 + id_2 * id_3$ only has this parse tree:

$$id_1 + (id_2 * id_3)$$



Operator's precedence

- $id_1 + (id_2 * id_3 * id_4) + id_5$

- Precedence: $+ \prec_p *$

- This other grammar G' reflects the operator's precedence:

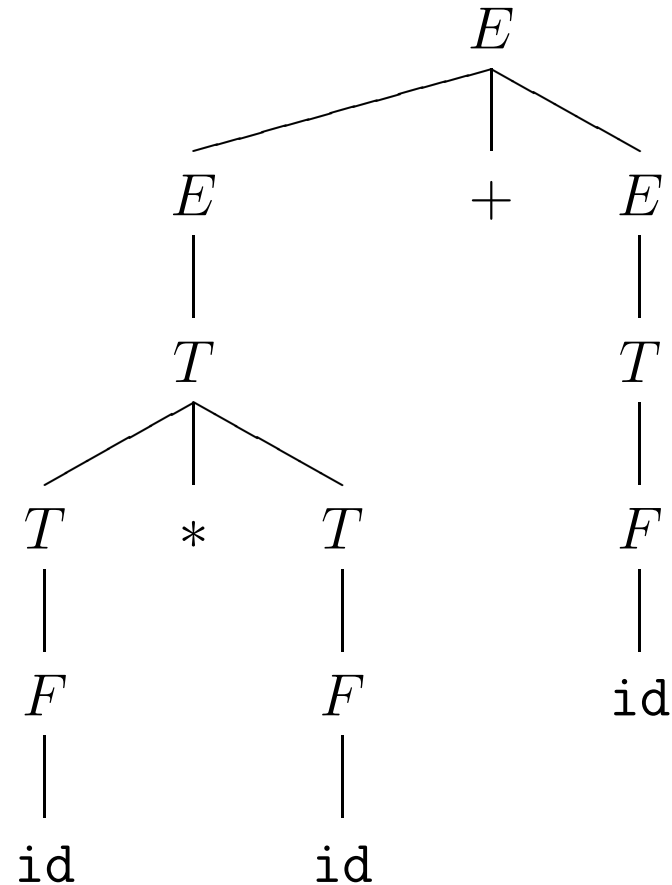
$$E \rightarrow E + E \mid T$$

$$T \rightarrow T * T \mid F$$

$$F \rightarrow (E) \mid id$$

- Now $w = id_1 * id_2 + id_3$ only has this parse tree:

$$(id_1 * id_2) + id_3$$



Operator's precedence

- $id_1 + (id_2 * id_3 * id_4) + id_5$

- Precedence: $+ \prec_p *$

- This other grammar G' reflects the operator's precedence:

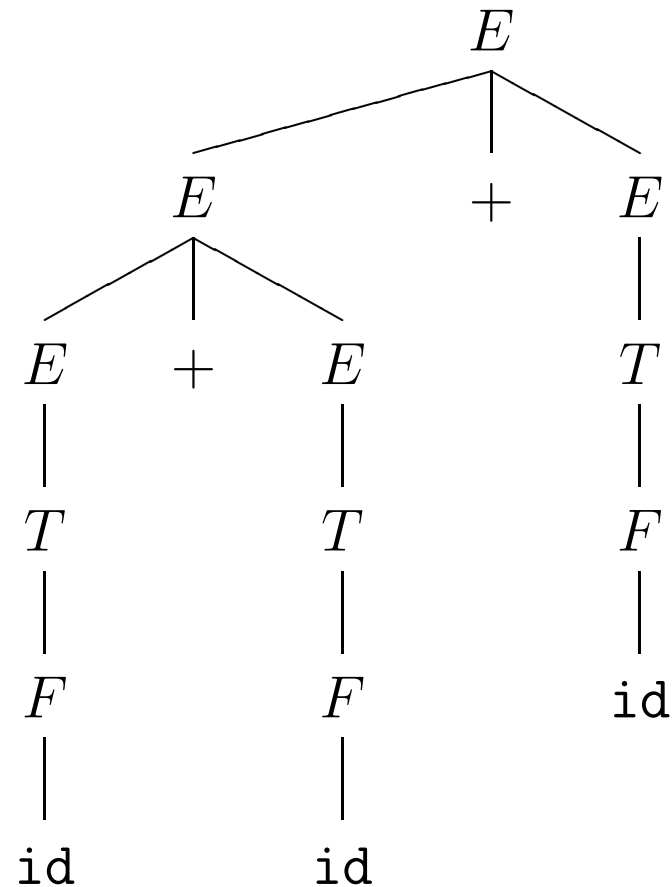
$$E \rightarrow E + E \mid T$$

$$T \rightarrow T * T \mid F$$

$$F \rightarrow (E) \mid id$$

- But $w = id_1 + id_2 + id_3$ still has two parse trees:

$$(id_1 + id_2) + id_3$$



Operator's precedence

- $id_1 + (id_2 * id_3 * id_4) + id_5$

- Precedence: $+ \prec_p *$

- This other grammar G' reflects the operator's precedence:

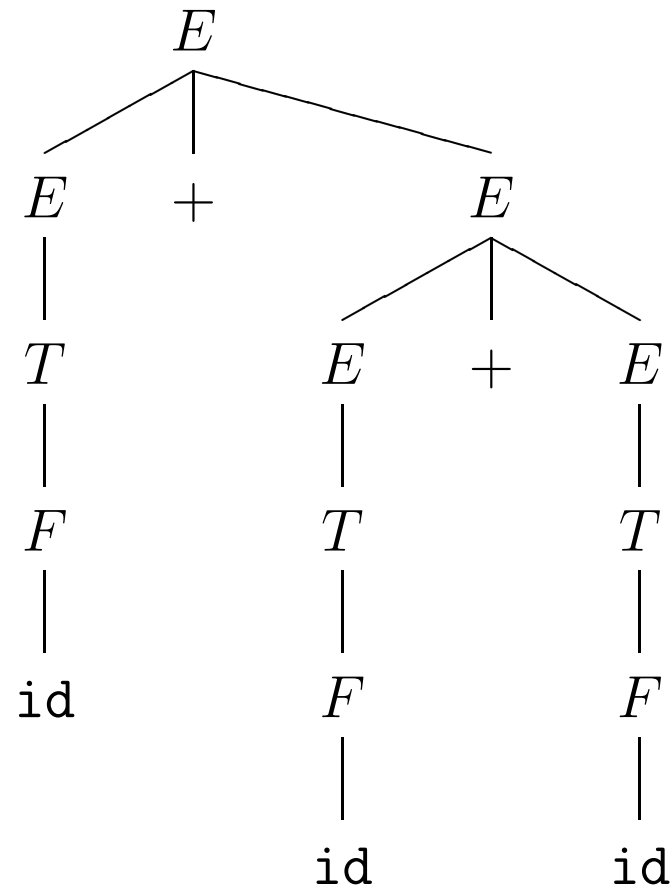
$$E \rightarrow E + E \mid T$$

$$T \rightarrow T * T \mid F$$

$$F \rightarrow (E) \mid id$$

- But $w = id_1 + id_2 + id_3$ still has two parse trees:

$$id_1 + (id_2 + id_3)$$



Operator's associativity

Left-associative:

$$e_1 \otimes e_2 \otimes e_3 \text{ means } (e_1 \otimes e_2) \otimes e_3$$

Right-associative:

$$e_1 \otimes e_2 \otimes e_3 \text{ means } e_1 \otimes (e_2 \otimes e_3)$$

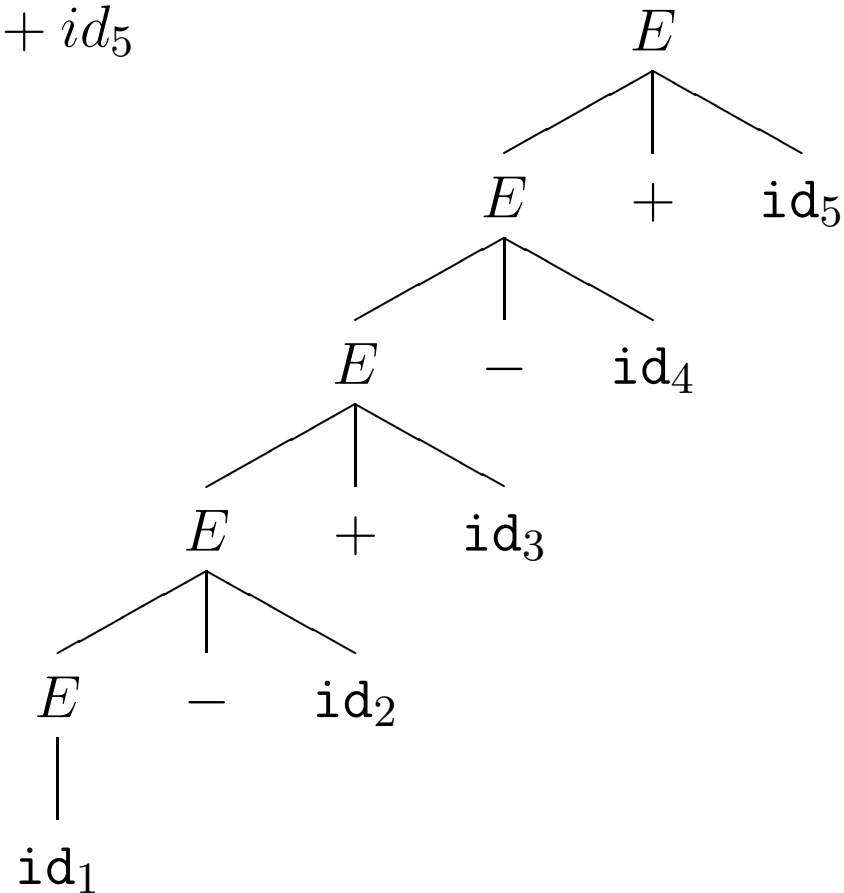
Non-associative:

$$e_1 \otimes e_2 \otimes e_3 \text{ is not correct}$$

Usually arithmetic (+, -, *, /) and logical operators (&&, ||) are left-associative, comparison operators (>, =, <, ...) are non-associative, and exponentiation (^) is right-associative.

Operator's associativity

$$\underbrace{\underbrace{\underbrace{\underbrace{(id_1 - id_2)}_E + id_3}_E - id_4}_E + id_5}_E$$



Operator's associativity

Left-associativity (+, -):

$$\begin{aligned} E &\rightarrow E + \text{id} \\ &\quad | E - \text{id} \\ &\quad | \text{id} \end{aligned}$$

Right-associativity (+, -):

$$\begin{aligned} E &\rightarrow \text{id} + E \\ &\quad | \text{id} - E \\ &\quad | \text{id} \end{aligned}$$

Non-associativity (>):

$$\begin{aligned} E &\rightarrow T > T \quad | \quad T \\ T &\rightarrow \text{id} \quad | \quad \text{num} \end{aligned}$$

Left-associativity (+, *):

$$\begin{aligned} E &\rightarrow E + T \\ &\quad | T \\ T &\rightarrow T * F \\ &\quad | F \\ F &\rightarrow (E) \\ &\quad | \text{id} \end{aligned}$$

Grammar for expressions

Grammar G :

$$E \rightarrow E + T$$

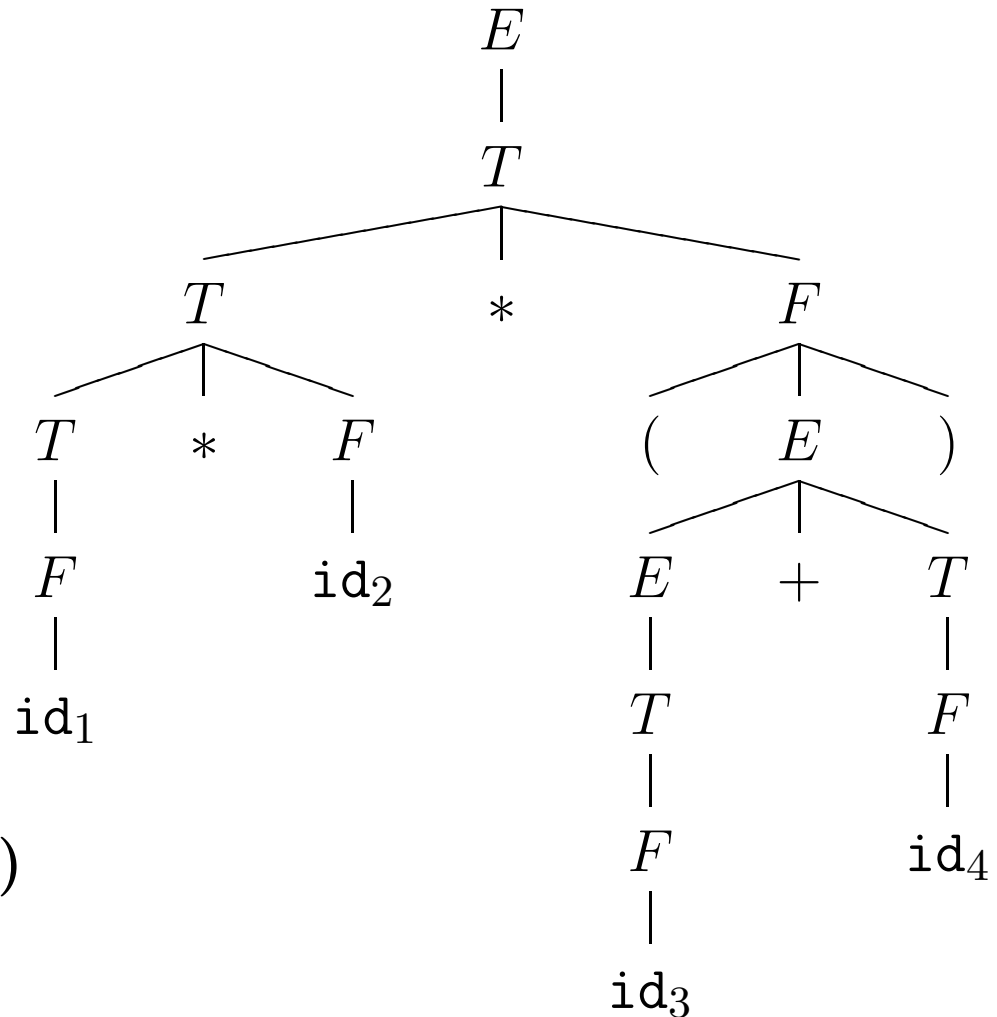
$$| T$$

$$T \rightarrow T * F$$

$$| F$$

$$F \rightarrow (E)$$

$$| \text{id}$$



$$w = \underbrace{id_1 * id_2} * \underbrace{(id_3 + id_4)}$$

Grammar for expressions

Grammar G :

$$E \rightarrow E + T$$

$$| T$$

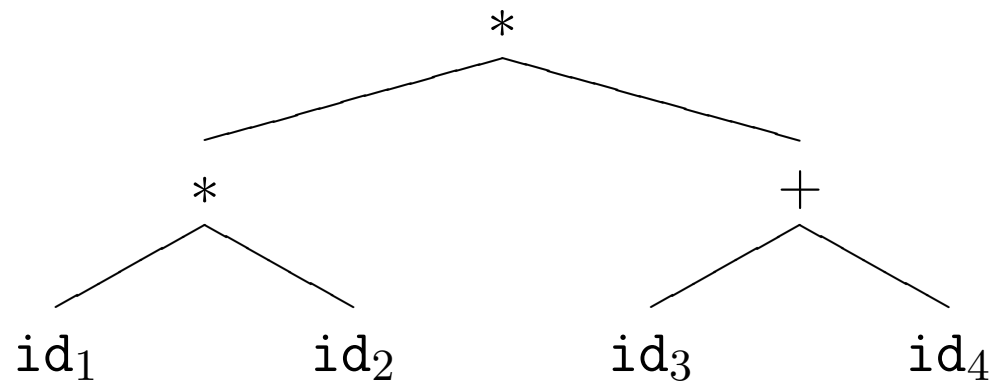
$$T \rightarrow T * F$$

$$| F$$

$$F \rightarrow (E)$$

$$| \text{id}$$

Abstract Syntax Tree
(AST) :



$$w = \underbrace{id_1 * id_2} * \underbrace{(id_3 + id_4)}$$

Exercises

- Extend the grammar for expressions with the following operators and precedence:

$$\{>, <, =\} \prec_p \{+, -\} \prec_p \{*, /\} \prec_p \{-_u\}$$

where $-_u$ denotes unary minus ($-e_1$).

All the binary operators are left-associative.

- Extend the previous grammar incorporating the access-to-array operator $[]$ (with the usual syntax $(e_1[e_2])$), the access-to-struct operator $.$, and the postfix access-to-pointed operator $^$ ($e_1^$).

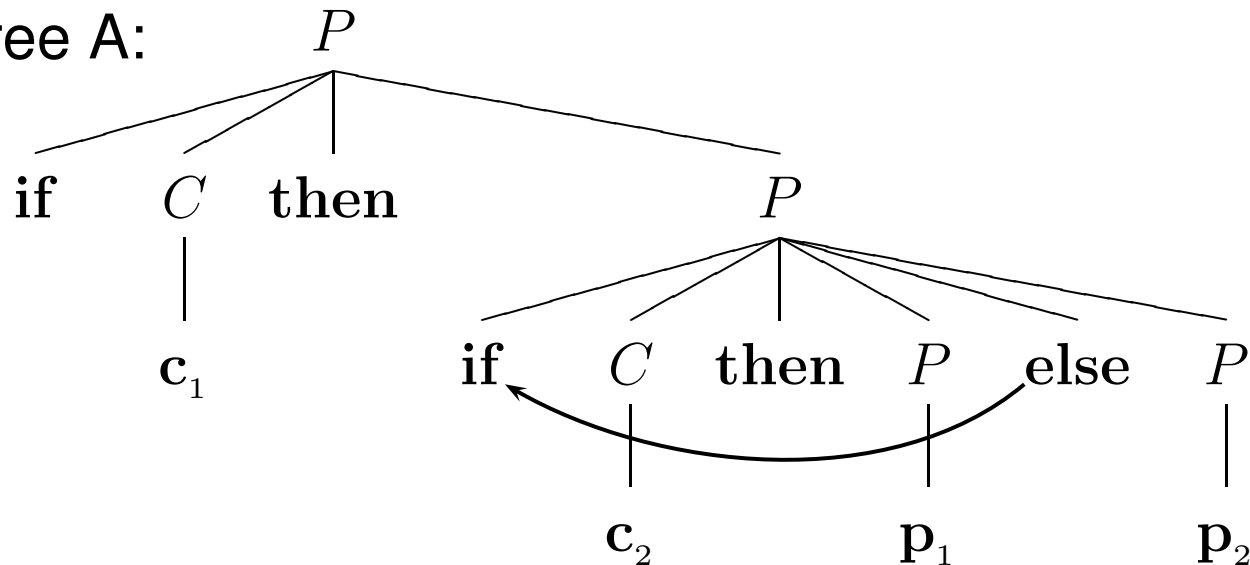
All of them have maximal precedence.

Ambiguous If-then-else (dangling else)

Grammar: $w = \text{if } c_1 \text{ then if } c_2 \text{ then } p_1 \text{ else } p_2$

$P \rightarrow \text{if } C \text{ then } P$
 $\quad \quad \quad | \text{if } C \text{ then } P \text{ else } P$
 $\quad \quad \quad | p$
 $C \rightarrow c$

Parse Tree A:

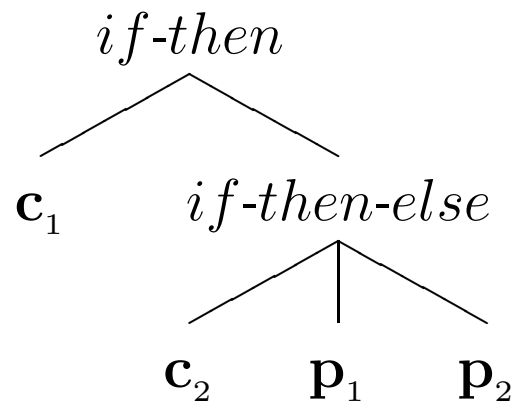


Ambiguous If-then-else (dangling else)

Grammar: $w = \text{if } c_1 \text{ then if } c_2 \text{ then } p_1 \text{ else } p_2$

$$\begin{aligned} P &\rightarrow \text{if } C \text{ then } P \\ &\quad | \text{if } C \text{ then } P \text{ else } P \\ &\quad | p \\ C &\rightarrow c \end{aligned}$$

Abstract Syntax Tree (AST) A:

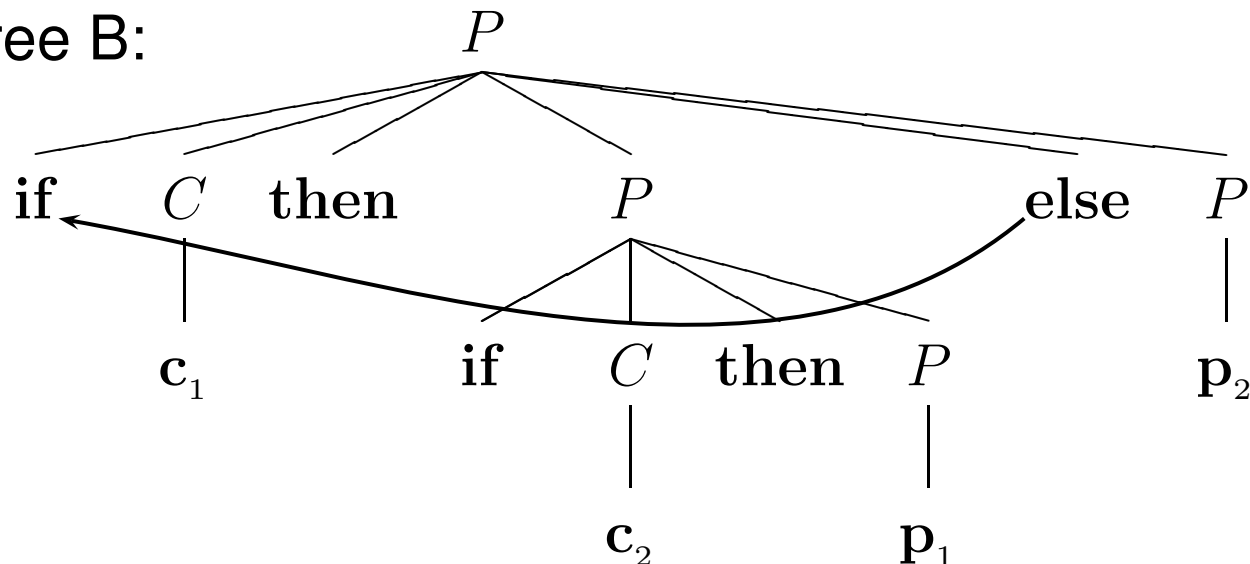


Ambiguous If-then-else (dangling else)

Grammar: $w = \text{if } c_1 \text{ then if } c_2 \text{ then } p_1 \text{ else } p_2$

$P \rightarrow \text{if } C \text{ then } P$
 $\quad \quad \quad | \text{if } C \text{ then } P \text{ else } P$
 $\quad \quad \quad | p$
 $C \rightarrow c$

Parse Tree B:

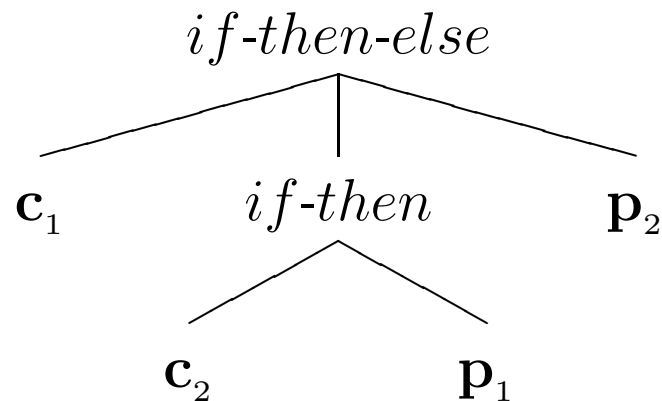


Ambiguous If-then-else (dangling else)

Grammar: $w = \text{if } c_1 \text{ then if } c_2 \text{ then } p_1 \text{ else } p_2$

$$\begin{aligned} P &\rightarrow \text{if } C \text{ then } P \\ &\quad | \text{if } C \text{ then } P \text{ else } P \\ &\quad | p \\ C &\rightarrow c \end{aligned}$$

Abstract Syntax Tree (AST) B:



Exercise

- Find an alternative non-ambiguous grammar for *this* language.

- Give a clue?

There are two kind of propositions:

- *closed*, where a following `else` cannot correspond to them.
- *open*, where a following `else` will correspond to them (to their `if`).