

Compilers

module of the course
“Professional English”

Yulia Burkatovskaya

Department of Computer Engineering

Associate professor

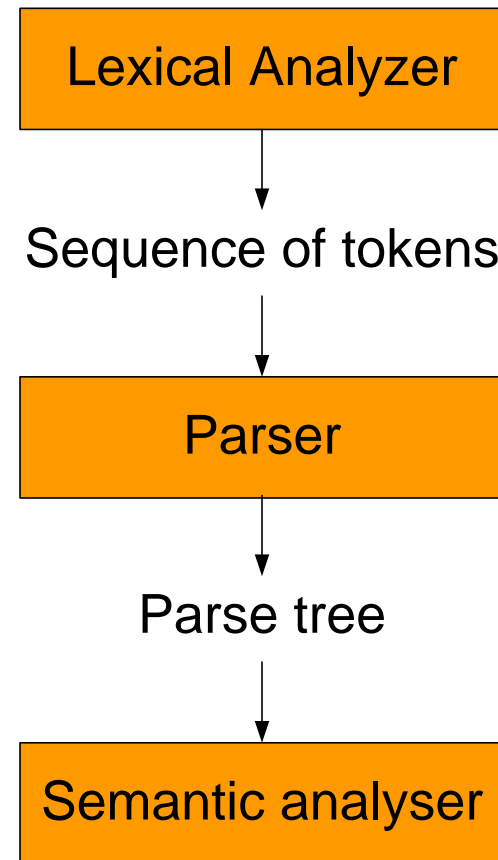
[3. Parsing]

- Basics of parsing
- Context-free grammar
- Parse tree
- Ambiguity
- Associativity and precedence
- Left-recursion
- Error-recovery

3.1. Basics of parsing

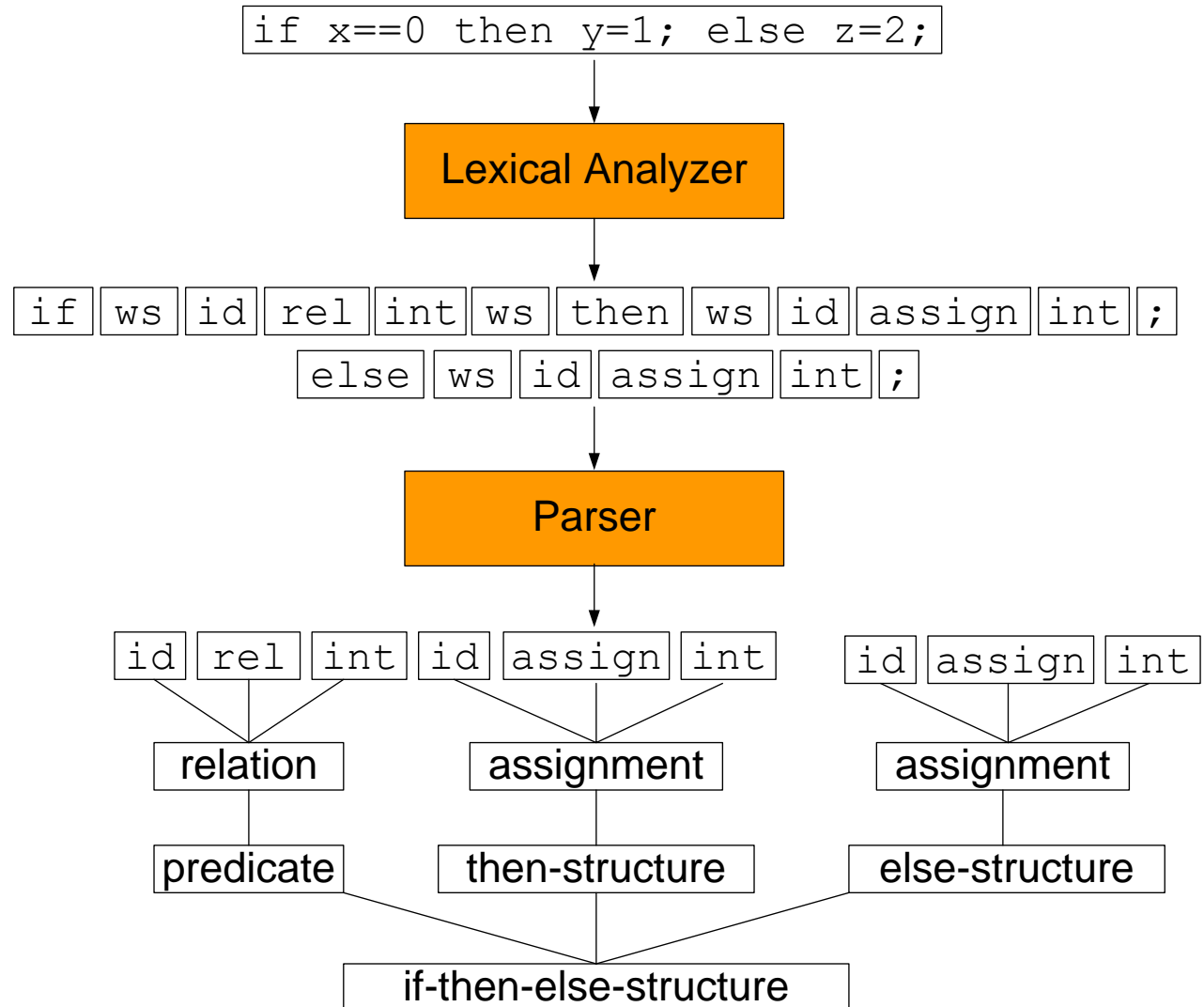
Syntax analysis tasks:

- To understand the meaning of the sequence of tokens
- To distinguish between valid and invalid strings of tokens
- To construct a parse tree for the semantic analyzer



Basics of parsing

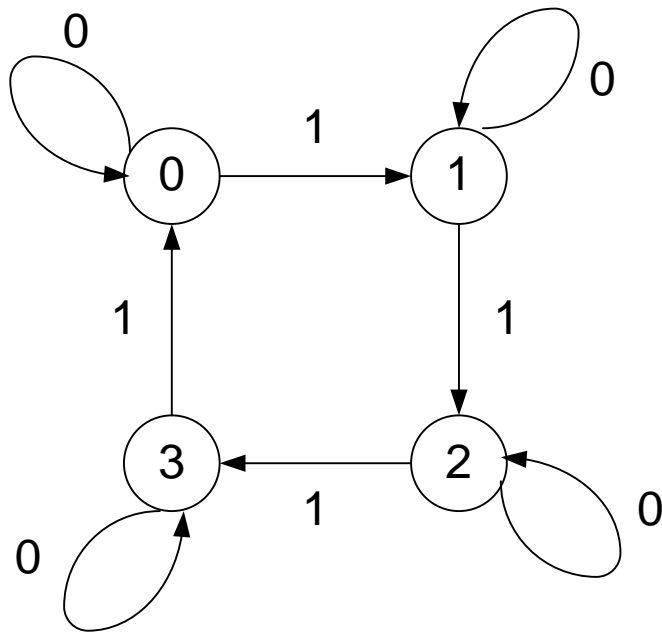
Example



Basics of parsing

- **Regular languages**

Counting “mod k”.



- **Context-free languages**

Recursive structures.

```
if ...  
  then ...  
    if ...  
      then ...  
      ...  
    else ...  
  else ...
```

[Basics of parsing]

Not every string of tokens is a program!

- **Regular expressions** describe tokens.
- **Context-free grammars** describe valid strings of tokens.

3.2. Context-free grammars

Context-free grammar $G = \langle A, T, N, P \rangle$

- A – a *start symbol*;
- T – a set of *terminals*;
- N – a set of *non-terminals*;
- P – a set of *productions*.

Production: $X \rightarrow Y_1 Y_2 \dots Y_k$;

$X \in N$;

$Y_i \in T \cup N \cup \{\epsilon\}$.

A set of productions: $X \rightarrow S_1 \mid S_2 \mid \dots \mid S_j$.

Context-free grammars

- Example (if-else structure)

EXP \rightarrow if EXP then ST;

EXP \rightarrow if EXP then ST; else ST;

EXP \rightarrow ID

EXP \rightarrow ID COM ID

EXP \rightarrow ID COM INT

ST \rightarrow ID = ID | ID = INT

COM \rightarrow == | < | <= | > | >=

ID \rightarrow LET | ID LET | ID DIG

INT \rightarrow DIG | INT DIG

LET \rightarrow a | ... | z | A | ... | Z

DIG \rightarrow 0 | ... | 9

- EXP=EXPRESSION
- ST=STATEMENT
- ID=IDENTIFIER
- COM=COMPARISON
- INT=INTEGER
- LET=LETTER
- DIG=DIGIT

Context-free grammars

Derivation

Let $G = \langle A, T, N, P \rangle$ be a CFG and:

- $S_1 S_2 \dots S_k \dots S_n \in (T \cup N \cup \{\epsilon\})^*$;
- $S_k \rightarrow Y_1 \dots Y_j \in P$,

then $S_1 S_2 \dots S_k \dots S_n \rightarrow S_1 S_2 \dots Y_1 \dots Y_j \dots S_n$ is a *step of derivation*.

$\alpha_0 \rightarrow \alpha_1 \rightarrow \dots \rightarrow \alpha_n$: α_n derives from α_0 in n steps ($\alpha_0 \rightarrow^* \alpha_n$).

Context-free grammars

EXP \rightarrow if EXP then ST;

EXP \rightarrow if EXP then ST; else ST;

EXP \rightarrow ID

EXP \rightarrow ID COM ID

EXP \rightarrow ID COM INT

ST \rightarrow ID = ID | ID = INT

COM \rightarrow == | < | <= | > | >=

ID \rightarrow LET | IDLET | IDDIG

INT \rightarrow DIG | INTDIG

LET \rightarrow a | ... | z | A | ... | Z

DIG \rightarrow 0 | ... | 9

EXP \rightarrow if EXP then ST; \rightarrow if if EXP then ST; else ST; then ST; \rightarrow

if if ID COM INT then ST; else ST; then ST; \rightarrow

if if ID COM INT then ST; else ST; then ID = ID; \rightarrow ...

if if $x_1 \leq 12$ then $x_2 = 23$; else $y = z$; then $z = x_2$;

Context-free grammars

- Leftmost derivation: the leftmost non-terminal is always chosen to replace.

EXP \rightarrow if EXP then ST; \rightarrow if if EXP then ST; else ST; then ST; \rightarrow
if if ID COM INT then ST; else ST; then ST; \rightarrow

- Rightmost derivation: the rightmost non-terminal is always chosen to replace.

EXP \rightarrow if EXP then ST; \rightarrow if EXP then ID = ID; \rightarrow
if EXP then ID = IDDIG; \rightarrow if EXP then ID = ID2;

Context-free grammars

$L(G) = \{\alpha: \alpha \in T^* \cup \{\epsilon\}, A \rightarrow^* \alpha\}$ – the **language generated by G**
(context-free language).

$G(A, T, N, P)$ and $G'(A', T, N', P')$ are **equivalent** if they generate the same language, i.e. $L(G) = L(G')$.

Implementing tools are sensitive to grammar.

$E \rightarrow (E)$

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow ID$

$E \rightarrow T$

$T \rightarrow E + T$

$T \rightarrow E * T$

$T \rightarrow ID$

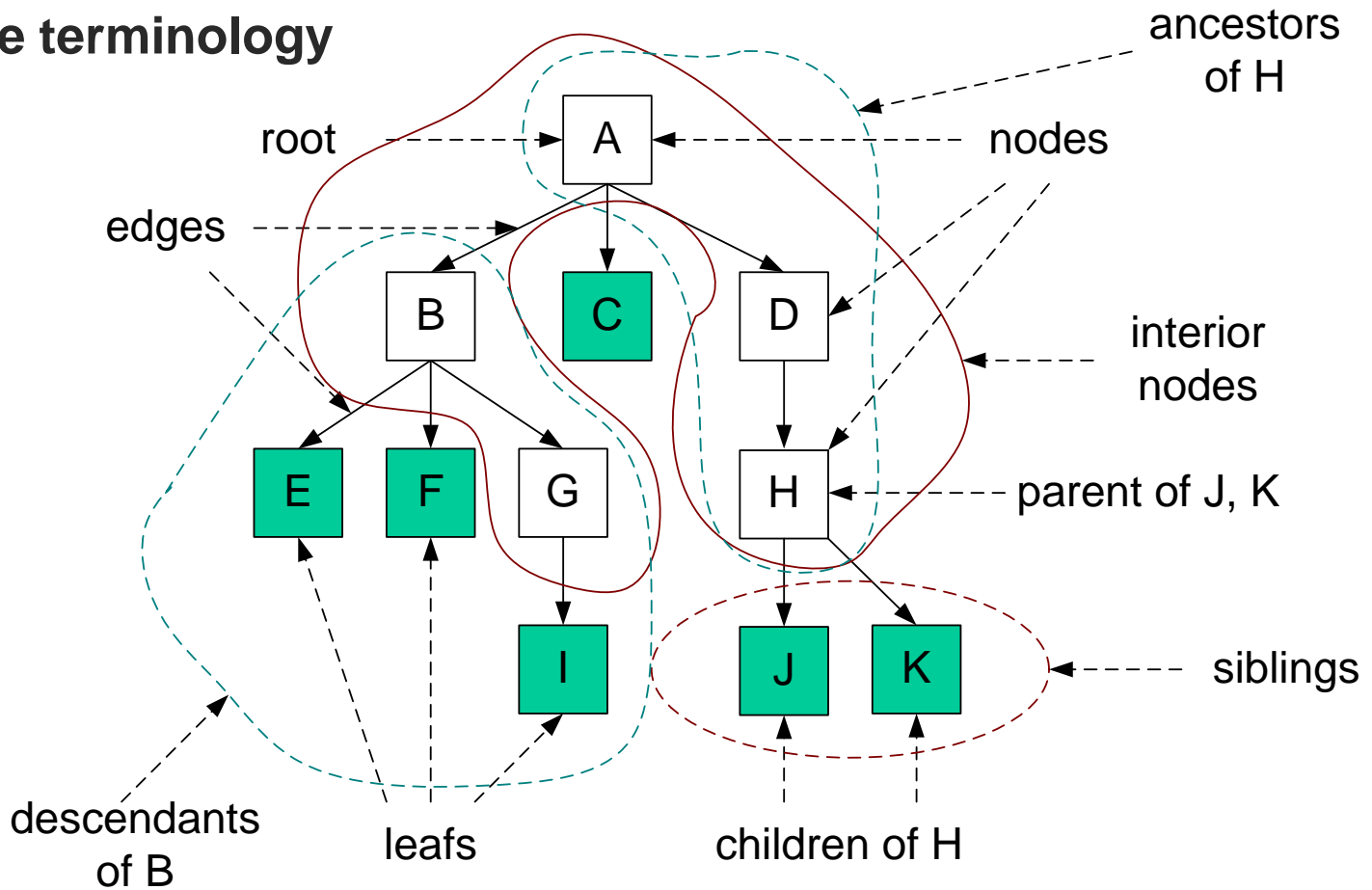
$T \rightarrow (E)$

[CFG vs Rexp]

- Separating the syntactic structure of a language into lexical and nonlexical parts provides a convenient way of modularizing the front end of a compiler into two manageable-sized components.
- The lexical rules of a language are frequently quite simple, and to describe them we do not need a notation as powerful as grammars.
- Regular expressions generally provide a more concise and easier-to-understand notation for tokens than grammars.
- More efficient lexical analyzers can be constructed automatically from regular expressions than from arbitrary grammars.

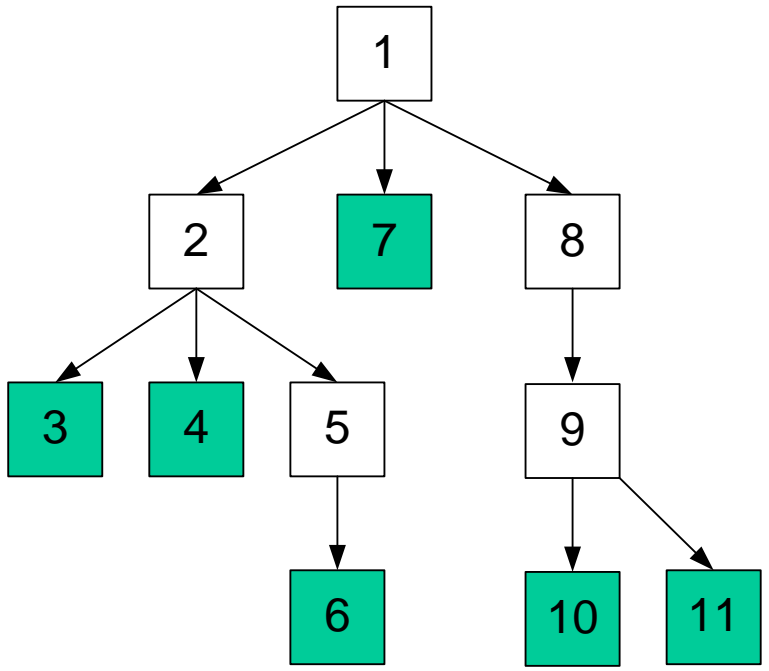
3.3. Parse tree

Tree terminology

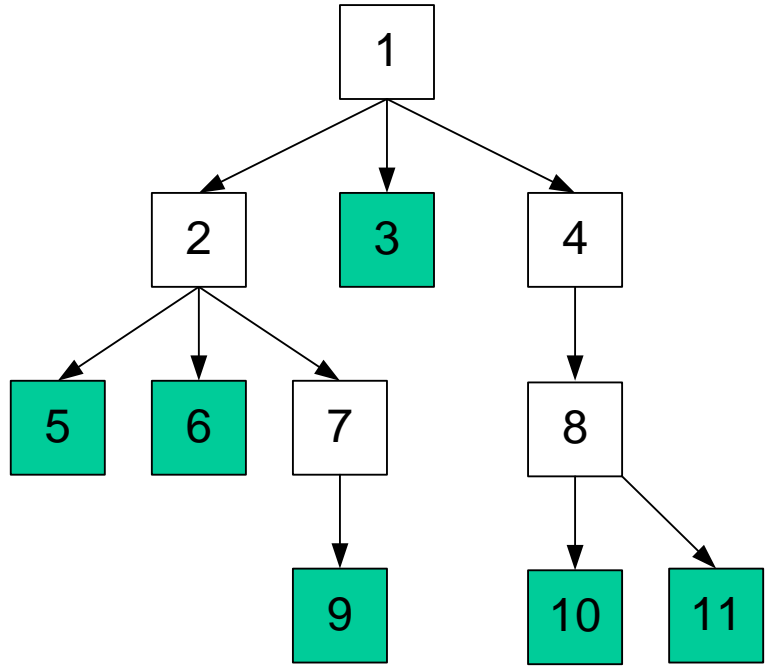


[Parse tree]

depth-first search

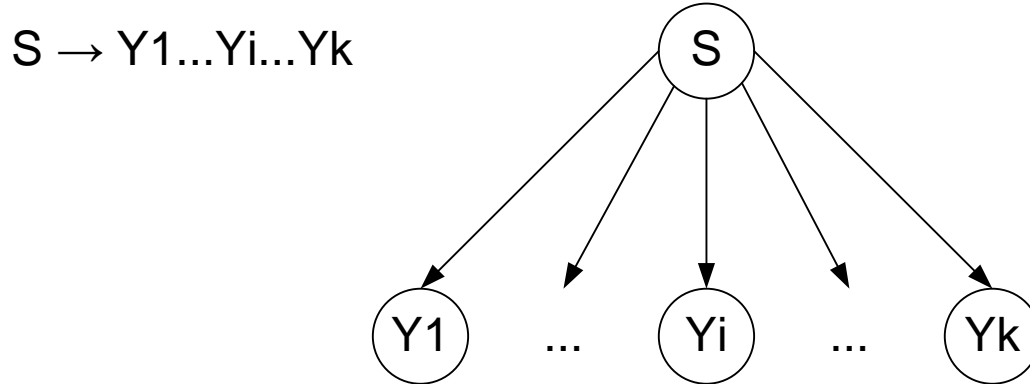


breadth-first search



Parse tree

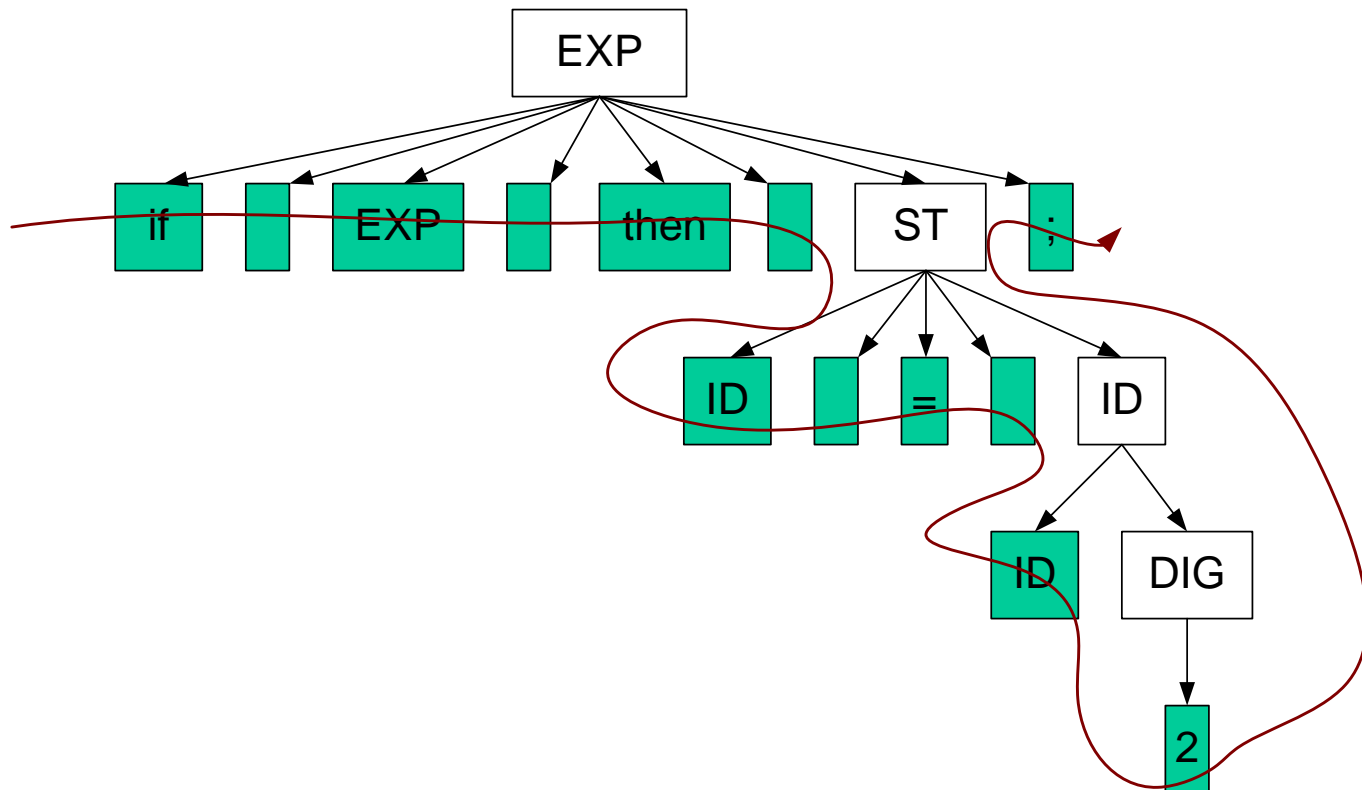
Parse tree is a graphical representation of a derivation.



- The **root** is labeled by the start symbol.
- **Leaves** are labeled by terminals. These labels read from left to right constitute a sequential form called the **yield** or **frontier** of the tree.

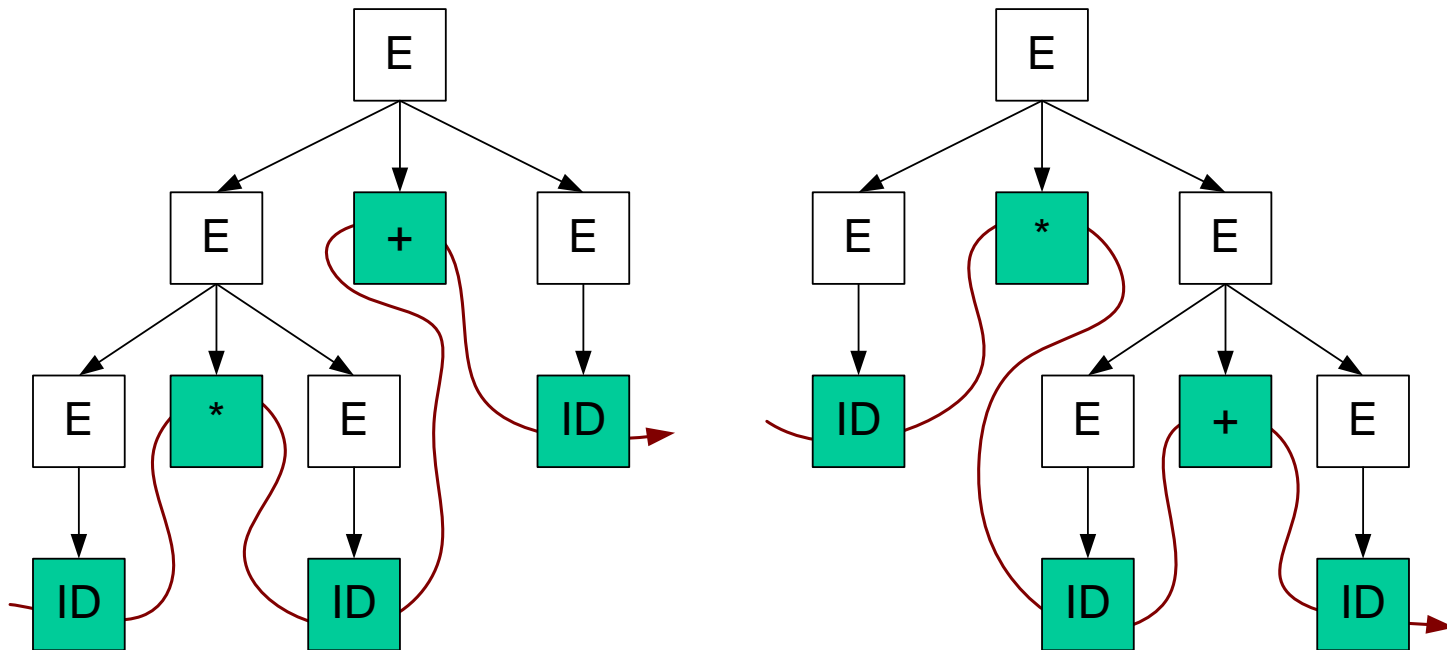
[Parse tree]

EXP \rightarrow if EXP then ST; \rightarrow if EXP then ID = ID; \rightarrow
if EXP then ID = IDDIG; \rightarrow if EXP then ID = ID2;



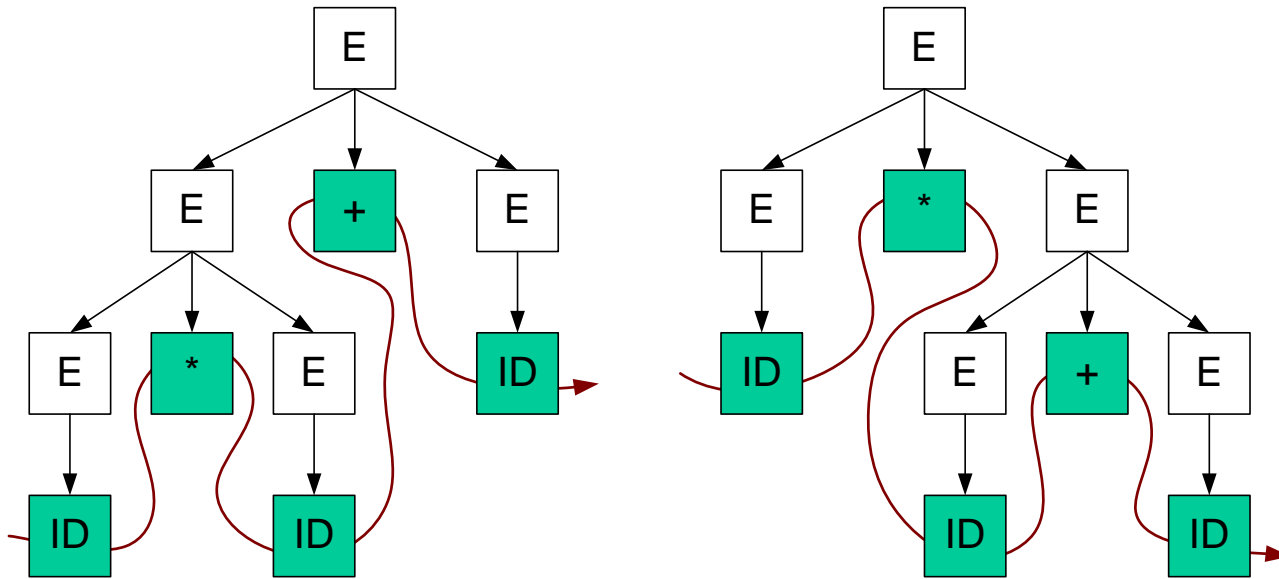
3.4. Ambiguity

- Grammar: $E \rightarrow E+E|E*E|(E)|ID$
- String: $ID*ID+ID$



[Ambiguity]

- A grammar is **ambiguous** if it has more than one parse tree for some string.
- It can cause different interpretations of a program!



(id*id)+id

id*(id+id)

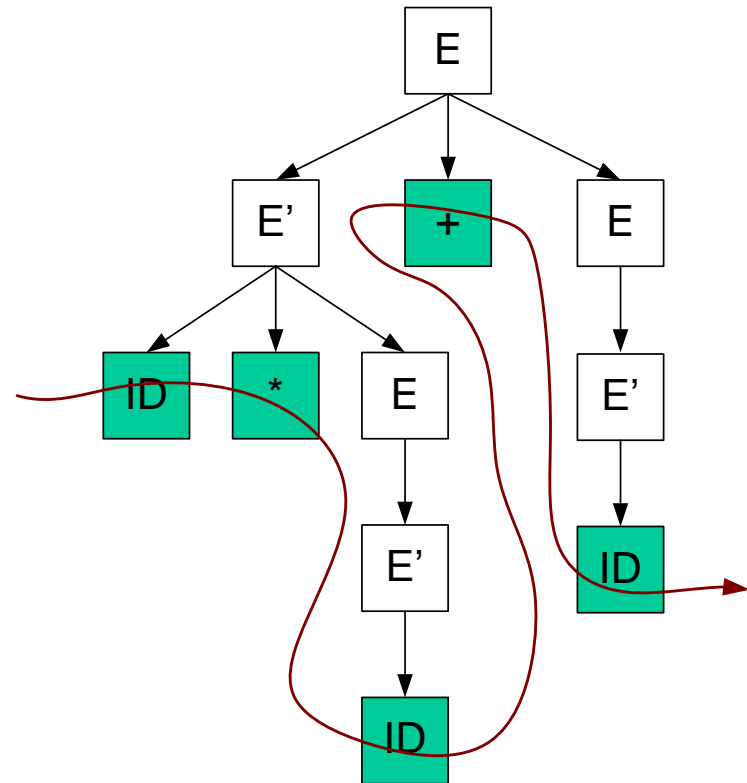
[Ambiguity]

To eliminate ambiguity:

- to rewrite a grammar;

$$E \rightarrow E'+E|E'$$

$$E' \rightarrow ID*E'|(E)*E'|(E)|D$$



$(id * id) + id$

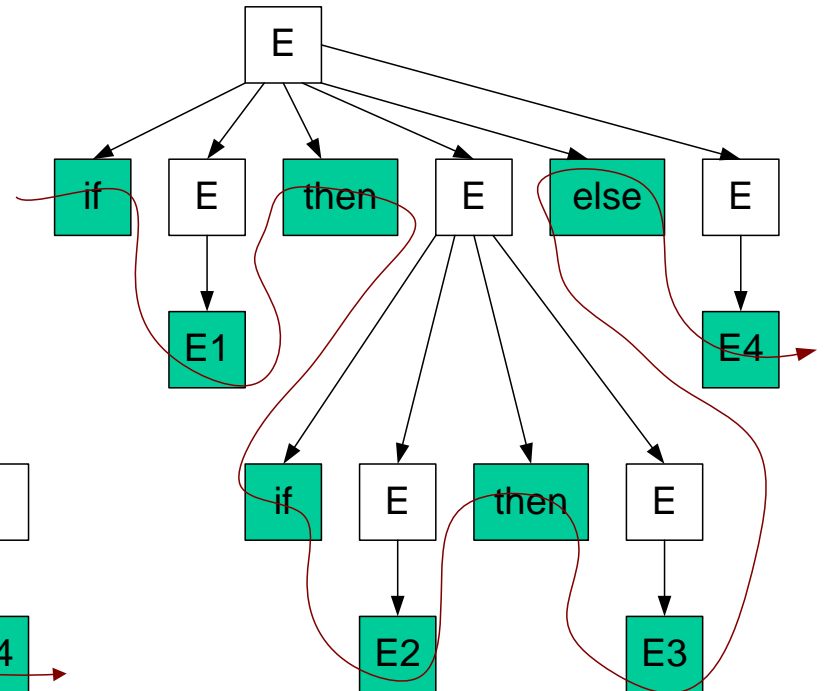
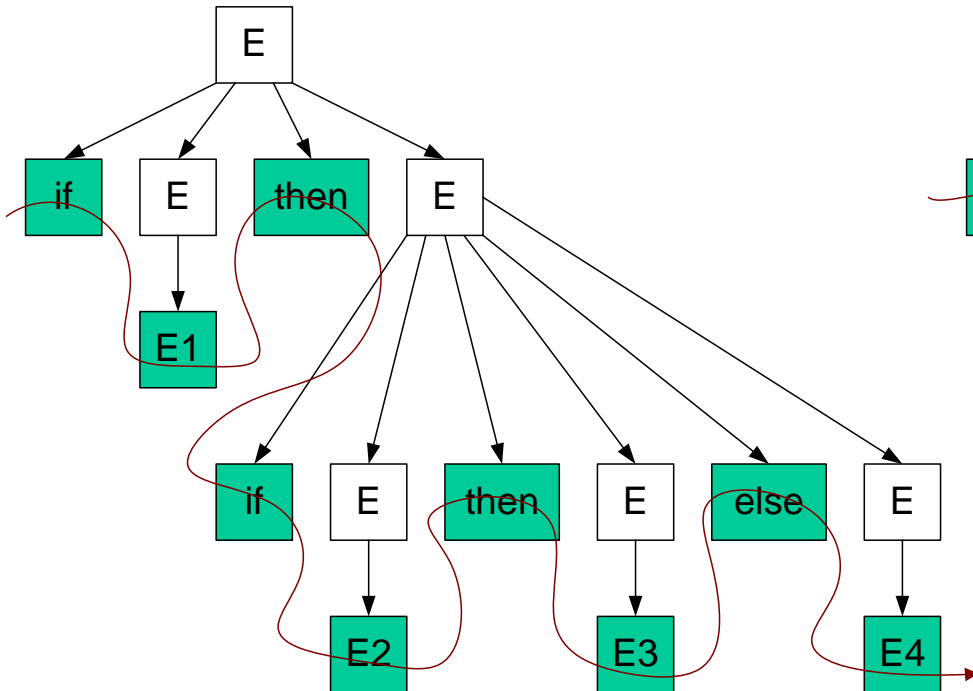
Ambiguity

- to use additional rules to resolve the ambiguities;

If-else grammar:

- $E \rightarrow \text{if } E \text{ then } E \mid \text{if } E \text{ then } E \text{ else } E \mid \text{ID}$

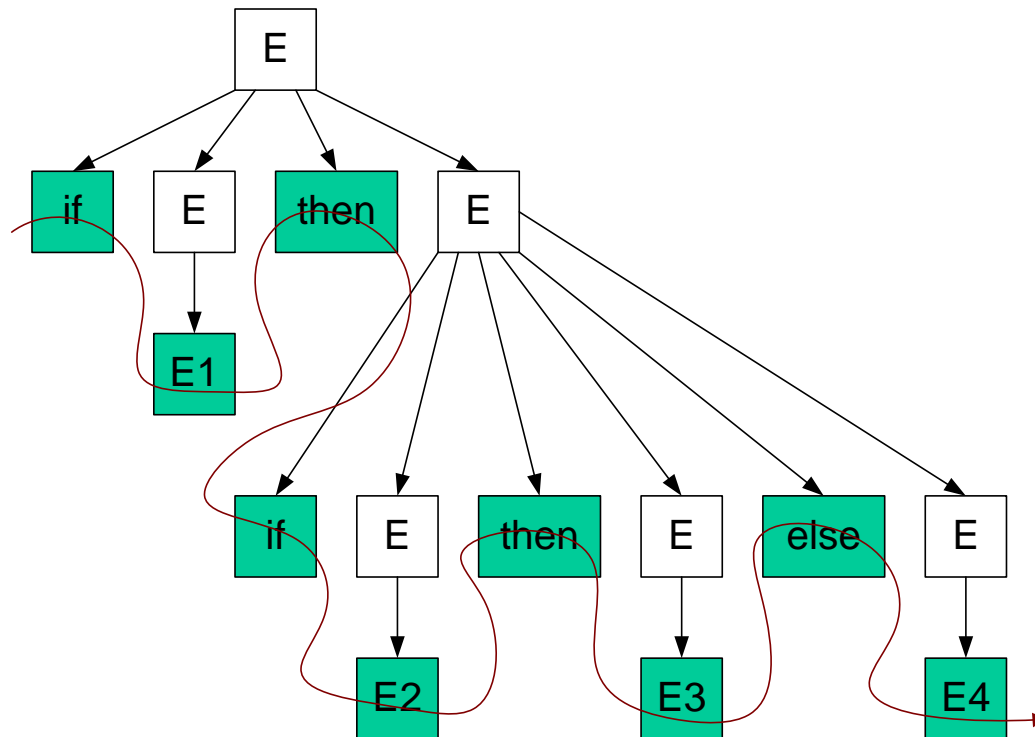
if E1 then if E2 then E3 else E4



Ambiguity

- **General rule:** “Match each **else** with the closest unmatched **then**”.

if E1 then if E2 **then** E3 **else** E4



[Ambiguity]

Unambiguous if-else grammar:

- $STMT \rightarrow MATCHEDSTMT \mid OPENSTMT$
- $MATCHEDSTMT \rightarrow \text{if } EXPR \text{ then } MATCHEDSTMT \text{ else } MATCHEDSTMT$
- $OPENSTMT \rightarrow \text{if } EXPR \text{ then } STMT \mid \text{if } EXPR \text{ then } MATCHEDSTMT \text{ else } OPENSTMT$

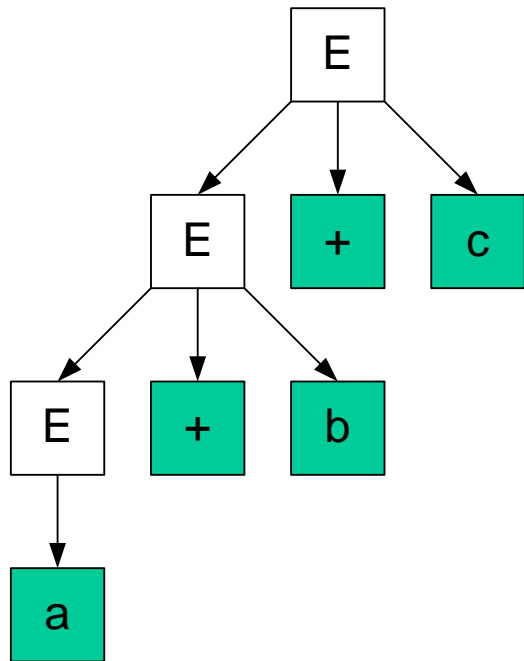
3.5. Associativity and precedence

Associativity of operators:

- left-associated operators ($-$, $+$, $*$, $/$, ...) – an operand with operator signs on both sides of it belongs to the operator to its left;
- right-associated operators ($=$, exponentiation, ...) – an operand with operator signs on both sides of it belongs to the operator to its right.

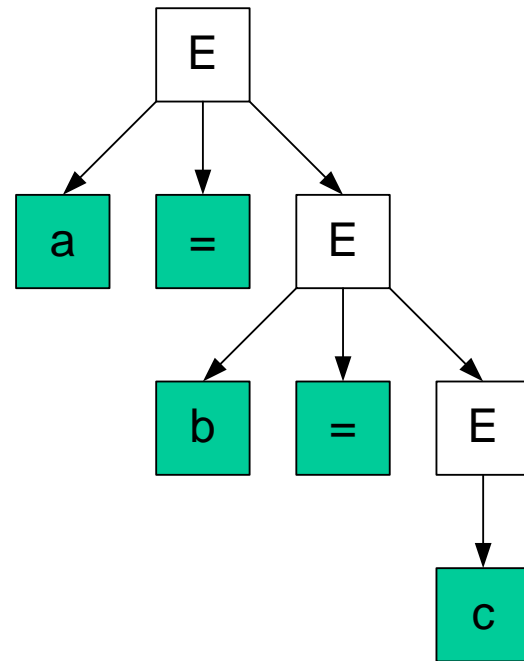
[Associativity and precedence]

$E \rightarrow E+a \mid \dots \mid E+z \mid a \mid \dots \mid z$



$(a+b)+c$

$E \rightarrow a=E \mid \dots \mid z=E \mid a \mid \dots \mid z$

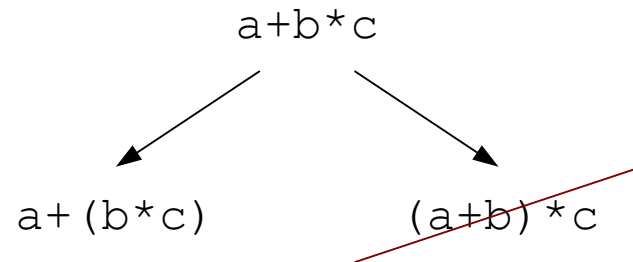


$a=(b=c)$

[Associativity and precedence]

We say that * has **higher precedence** than + if * takes its operands before + does.

Operators *, / have higher precedence than +, -



Associativity and precedence

Arithmetic grammar

- $E \rightarrow LET = E \mid LET = EXPR \mid LET$
- $EXP \rightarrow EXP + TERM \mid EXP - TERM \mid TERM$
- $TERM \rightarrow TERM * FACTOR \mid TERM / FACTOR \mid FACTOR$
- $FACTOR \rightarrow DIG \mid (EXP)$
- $DIG \rightarrow 0 \mid \dots \mid 9$
- $FACTOR \rightarrow LET \mid (EXP)$
- $LET \rightarrow a \mid \dots \mid z$

3.6. Left-recursion

- **Left-recursive grammar:** $A \rightarrow Aa$

Left-recursive arithmetic grammar:

- $E \rightarrow E+T \mid T$
- $T \rightarrow T*F \mid F$
- $F \rightarrow (E) \mid ID$

Non-left-recursive arithmetic grammar:

- $E \rightarrow TE' \mid T$
- $E' \rightarrow +TE' \mid +T$
- $T \rightarrow FT' \mid F$
- $T' \rightarrow *FT' \mid *F$
- $F \rightarrow (E) \mid ID$

Left-recursion

Left-recursive grammar:

- $A \rightarrow Aw_1 \mid \dots \mid Awn$
- $A \rightarrow v_1 \mid \dots \mid v_m$

Non-left-recursive grammar:

- $A \rightarrow v_1A' \mid \dots \mid v_mA'$
- $A \rightarrow v_1 \mid \dots \mid v_m$
- $A' \rightarrow w_1A' \mid \dots \mid w_nA'$
- $A' \rightarrow w_1 \mid \dots \mid w_n$

3.7. Left-factoring

- When the choice between two alternative productions for a non-terminal is not clear, we may be able to rewrite the productions to defer the decision until enough of the input has been seen that we can make the right choice.

If-else grammar:

- $S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid ID$

If-else grammar after left-factoring:

- $S \rightarrow \text{if } E \text{ then } S S' \mid ID$
- $S' \rightarrow \varepsilon \mid \text{else } S$

[Left-factoring]

Before:

- $A \rightarrow \sigma\alpha$
- $A \rightarrow \sigma\beta$

σ is the longest common prefix

After:

- $A \rightarrow \sigma A'$
- $A' \rightarrow \alpha$
- $A' \rightarrow \beta$

3.8. Error-recovery

Syntactic errors include

- misplaced semicolons or extra or missing braces; that is, "{" or " } . " ;
- in C or Java, the appearance of a case statement without an enclosing switch.

Error-recovery strategies:

- panic-mode;
- phrase-level;
- error-productions;
- global-correction.

Error-recovery

Panic-mode: the parser discards input symbols one at a time until one of a designated set of *synchronizing tokens* (semicolon ;, brace }, ...) is found.

- Panic-mode correction often skips a considerable amount of input without checking it for additional errors (-);
- it has the advantage of simplicity (+);
- it is guaranteed not to go into an infinite loop (+).

$(1++2) * 3$

Skip ahead to next integer (2) and then continue.

[Error-recovery]

Phrase-level: a parser may perform local correction on the remaining input; that is, it may replace a prefix of the remaining input by some string that allows the parser to continue (to replace a comma by a semicolon, delete an extraneous semicolon, or insert a missing semicolon).

- It can lead to infinite loops (for example, if we always inserted something on the input ahead of the current input symbol) (-).
- Phrase-level replacement has been used in several error-repairing compilers, as it can correct any input string (+).
- It has the difficulty in coping with situations in which the actual error has occurred before the point of detection (-).

Error-recovery

Error Productions: known common mistakes are specified in the grammar. A parser constructed from a grammar detects the anticipated errors when an error production is used during parsing.

- The parser can then generate appropriate error diagnostics about the erroneous construct that has been recognized in the input (+).
- It complicates the grammar (-).
- 5x instead of 5*x: add the production $E \rightarrow EE$

[Error-recovery]

Global Correction: find a correct “nearby” program trying token insertions and deletions and other changes in tokens.

- Hard to implement (-).
- Slows down parsing of correct programs (-).
- “Nearby” is not necessarily “the intended” program (-).

[Error-recovery]

Past

- Slow recompilation cycle (even once a day)
- Find as many errors in one cycle as possible

Present

- Quick recompilation cycle
- Users tend to correct one error/cycle
- Complex error recovery is less compelling