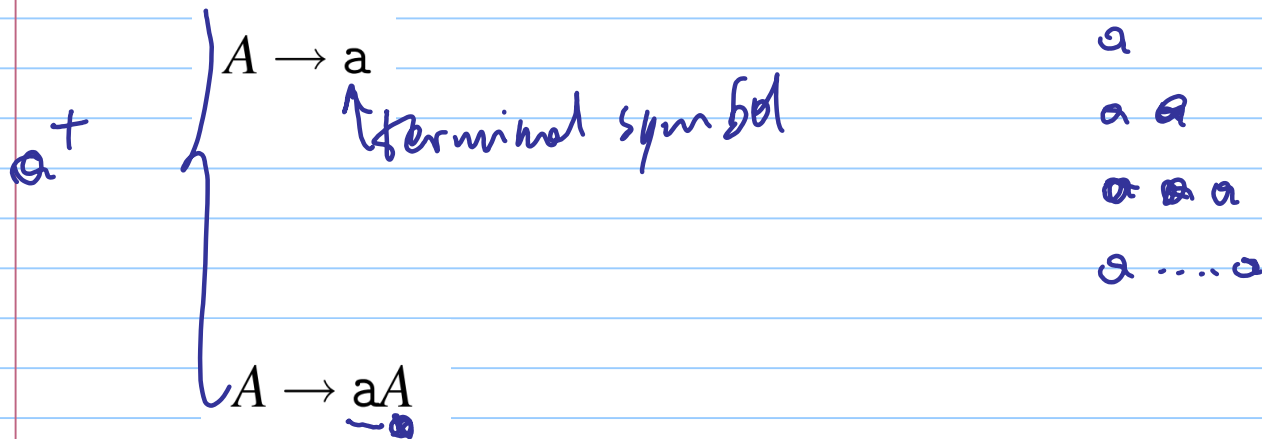Syntax analysis (Ch. 2 [M])

Takes a list of tokens, gives a syntax tree

context-free grammars

$$N \rightarrow X_1 \ldots X_n$$     production (rule, or production rule)

↑
one nonterminal symbol
on the LHS

$a^+$
{
$A \rightarrow$ a                       a

↑ terminal symbol           a a

                                     a a a

                                     a ....a

$A \rightarrow \underline{a}A$

$$a^* \quad \begin{cases} B \rightarrow \varepsilon \\ B \rightarrow aB \end{cases}$$

optionally

$\varepsilon, \ a, \ aa, \ aaa, \ \ldots$

| | | where : |
|---|---|---|
| $s^*$ | $\{``"\} \cup \{vw \mid v \in L(s), w \in L(s^*)\}$ | Each string in the language is a concatenation of any number of strings in the language of $s$. |

Up to this point, we have given examples of c.f. languages
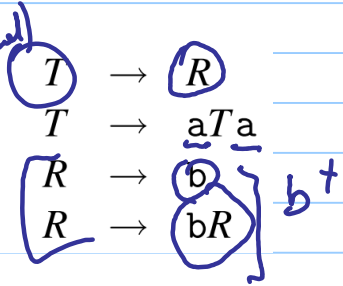that are also regular languages, for example:

$S \rightarrow \varepsilon$

$S \rightarrow aSb$

$\{a^n b^n \mid n \geq 0\} = \{\varepsilon, \ ab, \ aabb, \ aaabbb, \ldots\}$

(2)

several nonterminals : T and R

Start (nonterminal) symbol

$T \rightarrow R$
$T \rightarrow aTa$
$R \rightarrow b$
$R \rightarrow bR$

$b^{+}$

disjunction

$T \rightarrow R \mid aTa$
$R \rightarrow b \mid bR$

a shorthand for this

In EBNF:
$T \rightarrow b+ \mid aTa$

context-free grammar, using regexp notation

You may use $+$, $*$, $?$

| Form of $s_i$ | Productions for $N_i$ |
| --- | --- |
| $\varepsilon$ | $N_i \rightarrow$ |
| a | $N_i \rightarrow$ a |
| $s_j s_k$ | $N_i \rightarrow N_j N_k$ |
| $s_j \| s_k$ | $N_i \rightarrow N_j$<br>$N_i \rightarrow N_k$ |
| $s_j*$ | $N_i \rightarrow N_j N_i$<br>$N_i \rightarrow$ |
| $s_j+$ | $N_i \rightarrow N_j N_i$<br>$N_i \rightarrow N_j$ |
| $s_j?$ | $N_i \rightarrow N_j$<br>$N_i \rightarrow$ |

Fig. From regular expression to
2.1     context-free grammars

Syntactic category: constructs of a (programming) language that differ in meaning. Three typical ones:

Expressions: are evaluated to yield a value

Commands: are executed to change memory or perform I/O

Declarations: define properties of names used in other parts of the program

expressions without (and procedures) parentheses; can be described by a regular expression

$$\textbf{num}((+|-|*|/)\textbf{num})*$$

an ambiguous grammar

$$
\begin{aligned}
Exp &\rightarrow Exp + Exp \\
Exp &\rightarrow Exp - Exp \\
Exp &\rightarrow Exp * Exp \\
Exp &\rightarrow Exp / Exp \\
Exp &\rightarrow \textbf{num} \\
Exp &\rightarrow (Exp)
\end{aligned}
$$

$2 \ , \quad 5 + 2 - 3 \ , \quad \underline{3 - 6 * 5} \quad \underline{4 / 5 * 2}$

✓ $\underline{(3-6) * 5}$ ✗

Each syntactic category is denoted by a nonterminal, such as Exp in the grammar above.

*Stat -> **id** := Exp*          assignment

*Stat -> Stat ; Stat*          sequence / (list) of statements

*Stat -> **if** Exp **then** Stat **else** Stat*   two-way conditional

*Stat -> **if** Exp **then** Stat*      one-way conditional

A grammar for (simple) statements

Ex. 2, 3 [M]  (first part)

$P \rightarrow \varepsilon$

$P \twoheadrightarrow (P)$

$P \rightarrow P P$

$\varepsilon, \ (), \ ((\ )), \ ()(), \ (())(),$

$()(()), \ldots$

✗ $((\ )$

✗ $()\ ())$

1

Derivation $\Rightarrow$ or, sometimes $::=$

"may be rewritten as"

↓ (rewrite)

1. $\alpha N \beta \Rightarrow \alpha \gamma \beta$    if there is a production $N \to \gamma$
2. $\alpha \Rightarrow \alpha$                           (reflexivity)
3. $\alpha \Rightarrow \gamma$    if there is a $\beta$ such that $\alpha \Rightarrow \beta$ and $\beta \Rightarrow \gamma$    (transitivity)

**Definition 3.1** *Given a context-free grammar G with start symbol S, terminal symbols T and productions P, the language L(G) that G generates is defined to be the set of strings of terminal symbols that can be obtained by derivation from S using the productions P, i.e., the set $\{w \in T^* \mid S \Rightarrow w\}$.*

sentential form

$\Rightarrow$ for one step in a derivation

$\overset{*}{\Rightarrow}$ for what Mogensen calls a derivation.

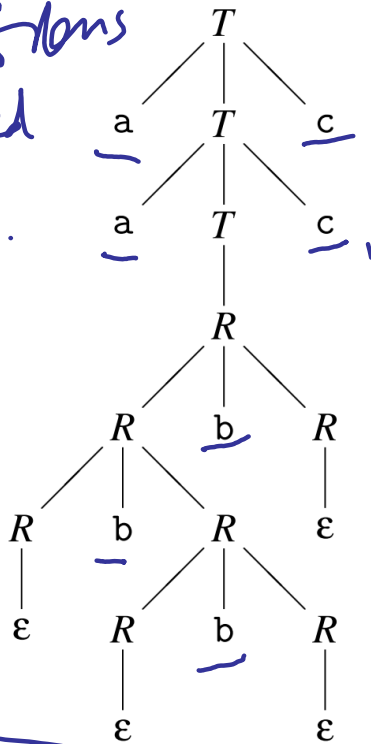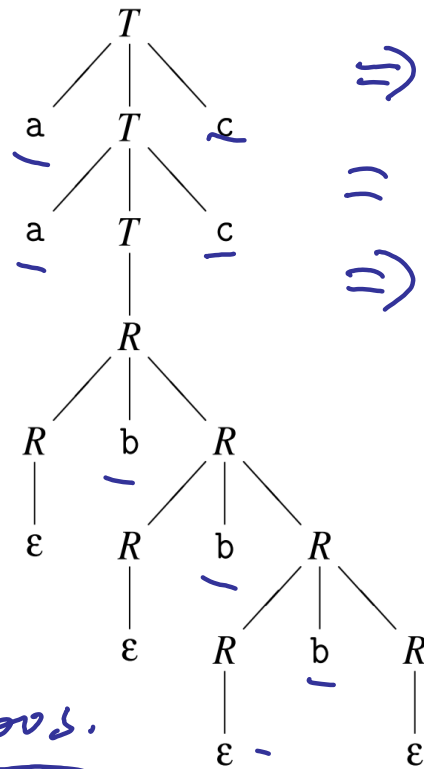$P \Rightarrow (P) \Rightarrow ((P)) \Rightarrow (())$    sentence

$T \rightarrow R$

$T \rightarrow aTc$

$R \rightarrow$

$R \rightarrow RbR$

ac

$T \Rightarrow aTc \Rightarrow aRc \Rightarrow ac$



$\underline{T}$

$\Rightarrow a\underline{T}c$

$\Rightarrow aa\underline{T}cc$

$\Rightarrow aa\underline{R}cc$

$\Rightarrow aaRb\underline{R}cc$

$\Rightarrow aa\underline{R}bcc$

$\Rightarrow aaRb\underline{R}bcc$

$\Rightarrow aaRb\underline{R}bRbcc$

$\Rightarrow aa\underline{R}bbRbcc$

$\Rightarrow aabb\underline{R}bcc$

$\Rightarrow aabbbcc$

neither a leftmost nor a rightmost derivation

$\underline{T}$

$\Rightarrow$ a$\underline{T}$c

$\Rightarrow$ aa$\underline{T}$cc

$\Rightarrow$ aa$\underline{R}$cc ,

$\Rightarrow$ aa$\underline{R}$b$R$cc

$\Rightarrow$ aa$\underline{R}$b$R$b$R$cc

$\Rightarrow$ aab$\underline{R}$b$R$cc

$\Rightarrow$ aab$\underline{R}$b$R$b$R$cc

$\Rightarrow$ aabb$\underline{R}$b$R$cc

$\Rightarrow$ aabbb$\underline{R}$cc

$\Rightarrow$ aabbbcc

a leftmost derivation

## Syntax tree

This tree corresponds to the two different two derivations described before.



$$\underline{a}\ a\ b\ b\ c$$

A grammar such that there is exists a string in its language that has two distinct syntax trees is ambiguous.

$$T \Rightarrow a\underline{T}c \Rightarrow a a \underline{T}cc \Rightarrow a a \underline{R}cc \Rightarrow$$

$$\Rightarrow a a \underline{R}bRcc \Rightarrow a a b \underline{R}cc \Rightarrow$$

$$\Rightarrow a a b \underline{R}bRcc \Rightarrow a a b b \underline{R}cc \Rightarrow$$

$$\Rightarrow a b b \underline{R}bRcc \Rightarrow$$

$$= a b b b \underline{R}cc \Rightarrow$$

$$\Rightarrow \boxed{a\ b\ b\ b\ c c}$$

(equivalent) grammar

This grammar is a non-ambiguous version of

$$T \to R$$
$$T \to aTc$$
$$R \to$$
$$R \to bR$$

$$T \to R$$
$$T \to aTc \quad , \text{which is ambiguous}$$
$$R \to$$
$$R \to RbR$$

Any grammar with productions like these:

$$\begin{cases} N \to N\alpha N \\ N \to \beta \end{cases}$$ is ambiguous

both left & right recursion

$\beta \overset{a}{\underset{}{}} \beta \overset{a}{\underset{}{}} \beta$ as a sentential form
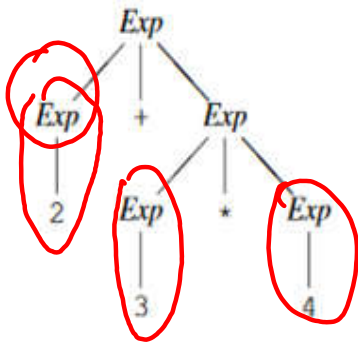
**Fig. 2.10** Fully reduced tree for the syntax tree in Fig. 2.7

**Fig. 2.11** Preferred syntax tree for 2+3*4 using Grammar 2.2, and the corresponding fully reduced tree



*✳ associates (binds) more tightly than +*

$$Exp \rightarrow Exp + Exp$$
$$Exp \rightarrow num$$

We need to get rid of productions, that are both left and right recursive

# Operator precedence & associativity

ambiguous $\begin{cases} E & \to & E \oplus E \\ E & \to & \textbf{num} \end{cases}$ — prod. is both left and right recursive

non-ambiguous assuming $\oplus$ is left-associative $\begin{cases} E & \to & E \oplus E' \\ E & \to & E' \\ E' & \to & \textbf{num} \end{cases}$ = production is only left recursive

$5 + 2 - 3 = \begin{cases} (5+2)-3 & \text{?} \\ 5 + (2-3) \end{cases}$

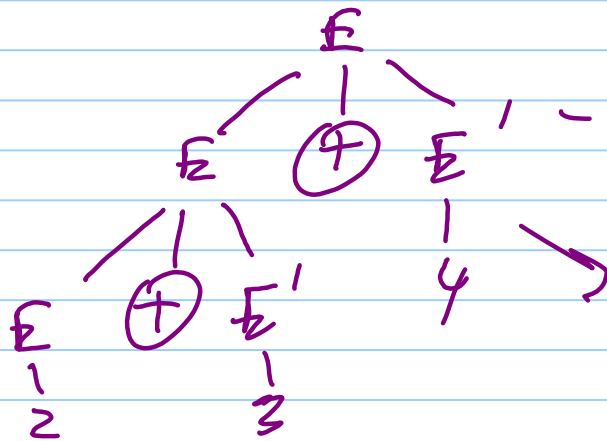$5 * 2 / 3 = \begin{cases} (5*2)/3 & \text{?} \\ 5 * (2/3) \end{cases}$

$2 - 3 - 4 = \begin{cases} 2 - (3-4) = 2-(-1)=3 & \text{?} \\ (2-3)-4 = -1-4=-5 \end{cases}$

By convention, $-$ and $/$ are left-associative, so

$2 - 3 - 4 = (2-3) - 4$

$$E \rightarrow E' \oplus E$$
$$E \rightarrow E'$$
$$E' \rightarrow \textbf{num}$$

$\oplus$ is a right associable operator

For example, the list constructor : in Haskell!

$$1 : 2 : [] = 1 : (2 : []) = 1 : [2] = [1,2]$$

$$E \rightarrow E' \oplus E'$$
$$E \rightarrow E'$$
$$E' \rightarrow \textbf{num}$$

Non-associable operators

$$add :: Int \rightarrow (Int \rightarrow Int)$$
$$add \; x \; y = x + y$$

"Increment by x"

e.g.  $<$ in Pascal is non-associative

In C,  $(3 < 4) < 5 = 1 < 5 = $ true

is left associative

$$E \rightarrow E + E'$$
$$E \rightarrow E - E'$$
$$E \rightarrow E'$$
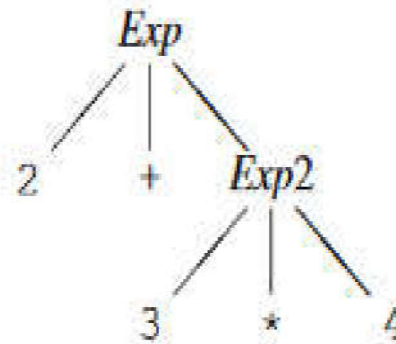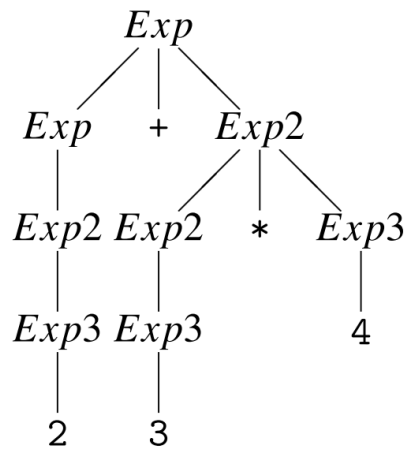$$E' \rightarrow \textbf{num}$$

$$E \rightarrow E + E'$$
$$E \rightarrow E' \oplus E$$
$$E \rightarrow E'$$
$$E' \rightarrow \textbf{num}$$

$$Exp \rightarrow Exp + Exp2$$
$$Exp \rightarrow Exp - Exp2$$
$$Exp \rightarrow Exp2$$
$$Exp2 \rightarrow Exp2 * Exp3$$
$$Exp2 \rightarrow Exp2 / Exp3$$
$$Exp2 \rightarrow Exp3$$
$$Exp3 \rightarrow \mathbf{num}$$
$$Exp3 \rightarrow （Exp）$$

```
            Exp
          /  |  \
       Exp   +   Exp2
        |        / | \
      Exp2    Exp2 * Exp3
        |       |       |
      Exp3    Exp3      4
        |       |
        2       3
```

```
        Exp
       / | \
      2  +  Exp2
            / | \
           3  *  4
```

$$Stat \quad \rightarrow \quad Stat2 \; ; \; Stat$$

$$Stat \quad \rightarrow \quad Stat2$$

$$Stat2 \quad \rightarrow \quad Matched$$

$$Stat2 \quad \rightarrow \quad Unmatched$$

$$Matched \quad \rightarrow \quad \texttt{if } Exp \texttt{ then } Matched \texttt{ else } Matched$$

$$Matched \quad \rightarrow \quad \mathbf{id} := Exp$$

$$Unmatched \quad \rightarrow \quad \texttt{if } Exp \texttt{ then } Matched \texttt{ else } Unmatched$$
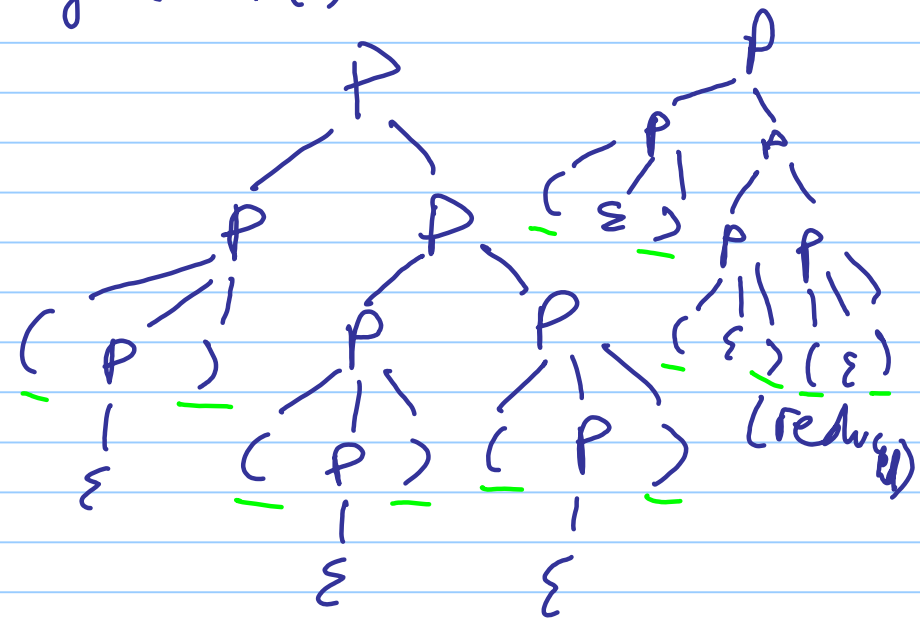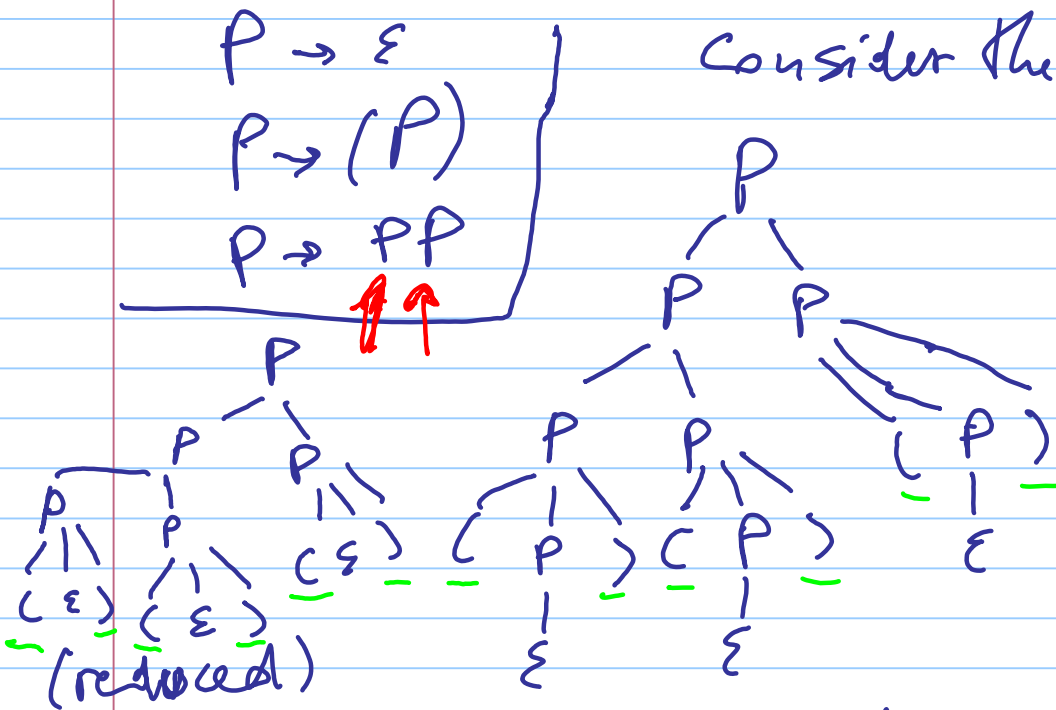
$$Unmatched \quad \rightarrow \quad \texttt{if } Exp \texttt{ then } Stat2$$

The following grammar (for the language of balanced parentheses) is ambiguous

$P \to \varepsilon$

$P \to (P)$

$P \to PP$

Consider the string ()()()

Note that the last production is both left and right recursive.

We can remove left recursion:

$P \to \varepsilon$

$P \to (P)P$

$()()()$

<span style="color:red">we removed the left recursion</span>