

Parsing and Ambiguity (Partly overlaps CSE 305)

Def: A parse tree for a string $x \in L(G)$ or a sentential form $X \in (ZUV)^*$

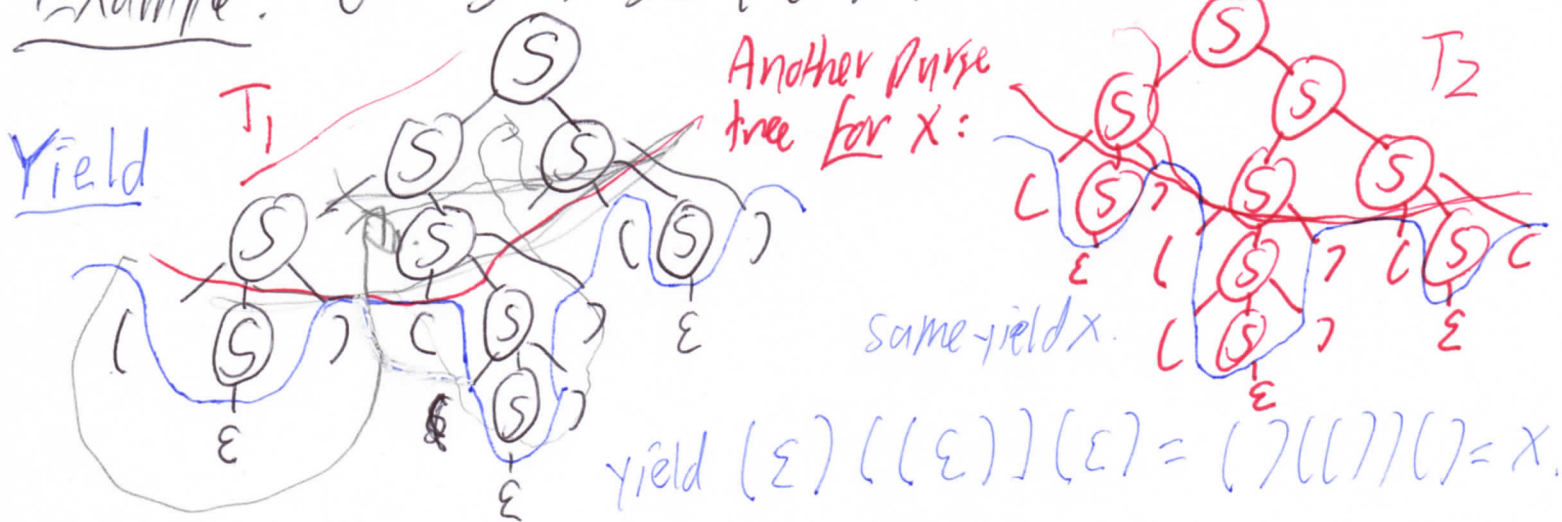
But for "LA" on the next page we'd want the root label to be A not S.

(*) is a tree with internal nodes labeled by variables root label is S leaves labeled by terminals or ϵ

such that for each internal node with label $A \in V$, its children have labels U_1, \dots, U_m such that

$A \rightarrow U_1 \dots U_m$ is a rule in the CFG G and finally the leaves in left-to-right order form x (or X).

Example: $G = S \rightarrow SS \mid (S) \mid \epsilon$ $x = () (()) ()$



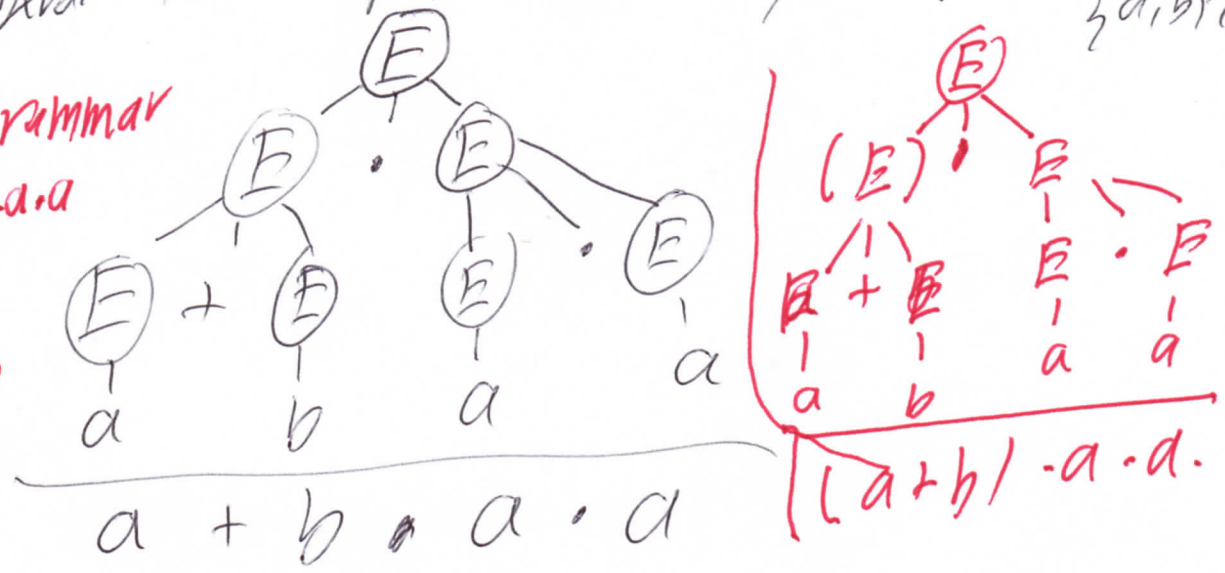
A Typical Expression Grammar
 Define $V = \{E\}$, $S = E$, and the rules

Let $\Sigma = \{a, b\}$ but let $T = \{a, b, +, \cdot, *, \epsilon, \emptyset, (,)\}$

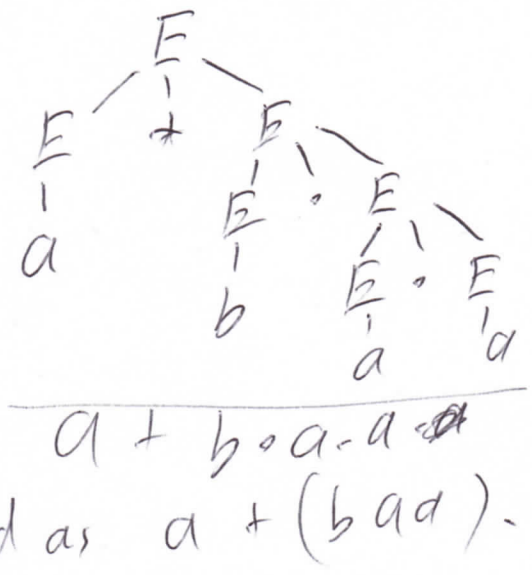
$$G_1: E \rightarrow (E) \mid E + E \mid E \cdot E \mid E^* \mid a \mid b \mid \epsilon \mid \emptyset$$

$L(G_1) = \{ \text{"liberal" textual representations of regular expressions over } \{a, b\} \}$

note: the grammar allows $(a+b) \cdot a \cdot a$ but this is a different string $x' \in L(G_1)$



We can also get x by a parse tree that respects the **actual** reading:



Hence G_1 is ambiguous (many times over), but happily, $L(G_1)$ is not inherently ambiguous...

An "overkill" = fix: CFG G_2

$$E \rightarrow (E + E) \mid (E \cdot E) \mid (E^*) \mid a \mid b \mid \epsilon \mid \emptyset$$

$L(G_2) = \{ \text{fully parenthesized regexps} \} \neq L(G_1)$. eg. $x'' = (a + (b \cdot a) \cdot a)$

An unambiguous grammar G_3 for regexps
 the way we write them, using three "categories":

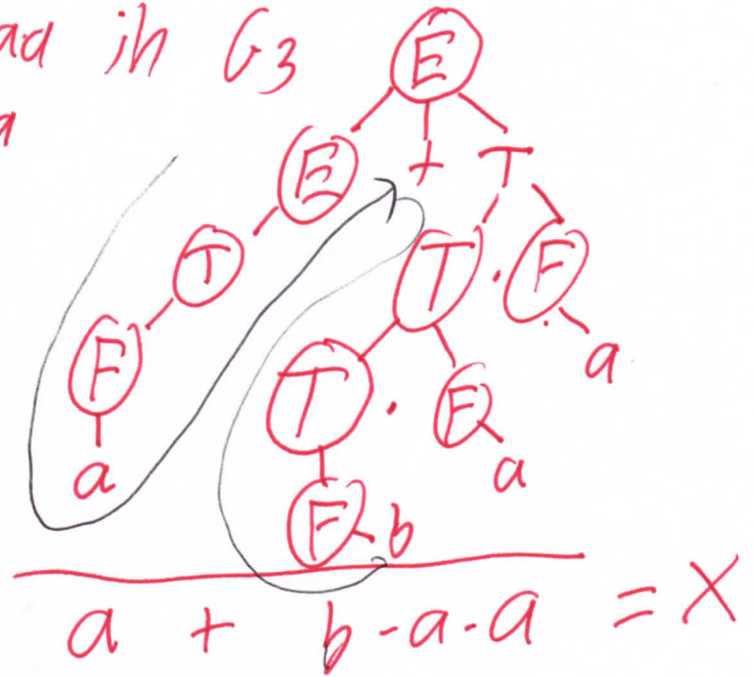
G_3 Expression $E \rightarrow T \mid E + T$ ~~$E \rightarrow T + E$~~
 Term $T \rightarrow F \mid T \cdot F$
 Factor $F \rightarrow a \mid b \mid \epsilon \mid \emptyset \mid (E) \mid F^*$

Let us try to derive $a + b \cdot a \cdot a$ in G_3

Now try to group $(a + b) \cdot a \cdot a$
 Can we do it without parens?



doesn't work



Added: The " $E \rightarrow T \cdot F$ "
 design pattern is used
 to define much wider
 precedence hierarchies
 of operators in many
 programming languages.
 My lecture stopped short
 of the "left-associative"
 vs "right-associative"
 issue, which comes into
 play when $-$ and $/$
 (etc.) are operators and which you should see in CS 4305.
 The big "gotcha" is that $a/b \cdot c$ parses as $(a/b) \cdot c$ when you probably intended $a/(b \cdot c)$ as if writing $\frac{a}{bc}$.



conclusion: to group $(a + b) \cdot a \cdot a$
 the grammar will force the
 parentheses. That's intuitively
 why it is unambiguous but
 gives the same language as G_1 .