

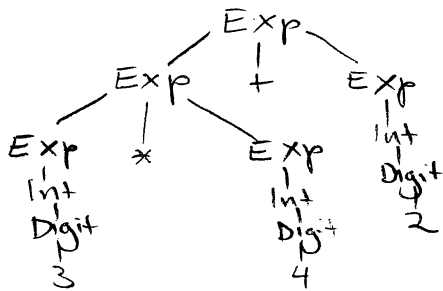
# CFG Ambiguity [H-M-U 5.4]

Example:

Integer  $\rightarrow$  Digit Integer | Digit

Digit  $\rightarrow$  0 | 1 | 2 | ... | 9

Exp(ression)  $\rightarrow$  Exp + Exp | Integer | Exp \* Exp



Thus  $\text{Exp} \xRightarrow{*} 3 * 4 + 2$   
 i.e.,  $3 * 4 + 2 \in L(\text{Exp})$

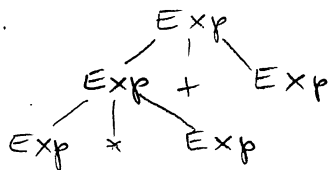
Observe: The same string can have multiple derivations, eg.

leftmost  $\rightarrow$   $\text{Exp} \Rightarrow \text{Exp} + \text{Exp} \Rightarrow \text{Exp} * \text{Exp} + \text{Exp} \xRightarrow{*} 3 * \text{Exp} + \text{Exp}$   
 $\xRightarrow{*} 3 * 4 + \text{Exp} \xRightarrow{*} 3 * 4 + 5$ , and

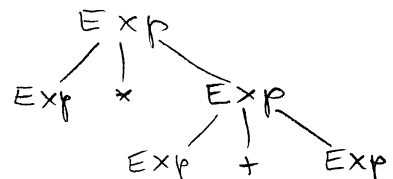
rightmost  $\rightarrow$   $\text{Exp} \Rightarrow \text{Exp} + \text{Exp} \xRightarrow{*} \text{Exp} + 2 \Rightarrow \text{Exp} * \text{Exp} + 2$   
 $\xRightarrow{*} \text{Exp} * 4 + 2 \xRightarrow{*} 3 * 4 + 2$

This is okay. It is usually not okay if there are multiple parse trees yielding the same string.

Eg.



versus



same yield,  $\text{Exp} * \text{Exp} + \text{Exp}$ , but very different meanings!

Def: A CFG is unambiguous if every terminal string has at most one parse tree.

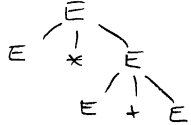
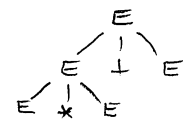
Equivalently, unambiguous  $\Leftrightarrow$  every terminal string has at most one leftmost derivation

ambiguous  $\Leftrightarrow$  some terminal string has two or more parse trees (or leftmost derivations)

## Removing ambiguity (leaving $L(G)$ the same)

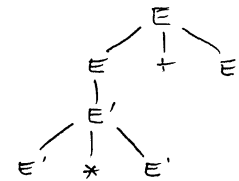
### Idea 1 Operator precedence

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



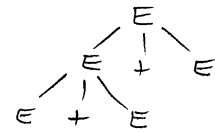
$$G': E \rightarrow E + E \mid E'$$

$$E' \rightarrow \text{Num} \mid E' * E' \mid (E)$$

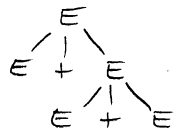


"+" should have strictly lower precedence, i.e. be higher in the parse tree, than "\*" - but  $G'$  is still ambiguous!

$E + E + E$  has parse trees



and

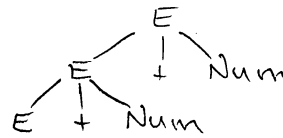


### Idea 2 Force left associativity

$$H: E \rightarrow \text{Num} \mid E + (E)$$

we don't want to allow this  $E$  to be expanded further

$$H': E \rightarrow \text{Num} \mid E + \text{Num}$$



$$G'': E \rightarrow E + E' \mid E'$$

$$E' \rightarrow E' * E'' \mid E''$$

$$E'' \rightarrow (E) \mid \text{Num}$$

Idea:  $E$  can be split into a "+" but  $E'$  can only expand as a "\*" or  $(E)$  and  $E''$  can only be  $(E)$  or Number.

Note: Yacc/bison can automatically deal with precedence/associativity

## Inherent ambiguity [H-M-U ch. 5.4.4]

- $L = \{a^m b^m c^n d^n \mid m, n \geq 0\} \cup \{a^m b^n c^n d^m \mid m, n \geq 0\}$  is a context-free language (CFL)

Theorem: Every grammar  $G$  for  $L$  is ambiguous!

-essentially strings like  $a^m b^m c^m d^m$  must have two parse trees, for  $m$  large enough.

## Simplifying the presentation of Context-free Grammars

- ① Eliminate "useless" symbols
- ② Convert to "Chomsky normal form" (CNF)

Def: A variable or terminal  $X$  is useful if

- Ⓐ  $X$  derives a terminal string, and
- Ⓑ  $X$  is reachable from the start, i.e.,  
 $S \xRightarrow{*} \alpha X \beta$

### Eliminating useless symbols

- Ⓐ Eliminate symbols that cannot generate a terminal string
- Ⓑ Eliminate symbols that are not reachable

Ⓐ Example  $G: S \rightarrow AB|C$   
 $A \rightarrow OB|C$   
 $B \rightarrow 1|AO$   
 $C \rightarrow AC|C$   
 $G': S \rightarrow AB$   
 $A \rightarrow OB, B \rightarrow 1|AO$

	symbols					
	0	1	A	B	C	S
round	✓	✓				
2	✓	✓		✓		
3	✓	✓	✓	✓		
4	✓	✓	✓	✓		✓
5	no changes $\Rightarrow$ eliminate C					

Algorithm: Loop through the rules

- every terminal is generating
- if  $A \rightarrow X_1 X_2 \dots X_k$  and each  $X_j$  is generating, then so is  $A$

updating the symbols in each round; until no changes are made in a loop.

Ⓑ Use a graph traversal algorithm

•  $S$  is reachable

• If  $A$  is reachable and  $A \rightarrow X_1 X_2 \dots X_k$ , then so are the  $X_j$

round	0	1	S	A	B
1			✓		
2			✓	✓	✓
3	✓	✓	✓	✓	✓

Important: Remember to eliminate non-generating symbols first.

Example:  $S \rightarrow AB|a$   
 $A \rightarrow b$

should simplify to  $S \rightarrow a$

In general, if  $S \xRightarrow{*} CD$  and  $C$  is not generating, then  $D$  is reachable, but this derivation can't go anywhere

## Chomsky Normal Form (CNF)

Definition: A context-free grammar is in Chomsky normal form if all its productions are of the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a,$$

where  $A, B, C$  are variables and  $a$  is a terminal.

Theorem: For any context-free language  $L$ , there exists a context-free grammar  $G$  with

- $L(G) = L \setminus \{\epsilon\}$
- $G$ 's symbols are all useful
- $G$  is in Chomsky Normal Form.

Example:  $E \rightarrow E + E \mid E * E \mid \text{Digit}$   
 $\text{Digit} \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

in CNF:  $E \rightarrow E \text{ Eplus} \mid E \text{ Etimes} \mid 0 \mid 1 \mid \dots \mid 9$   
 $E_{\text{plus}} \rightarrow \text{PLUS } E$   
 $E_{\text{times}} \rightarrow \text{TIMES } E$   
 $\text{PLUS} \rightarrow +$   
 $\text{TIMES} \rightarrow *$

Why? - simpler because every production has one of two forms  
- for proving the Pumping Lemma for context-free languages

- Proof idea:
- ① Eliminate useless symbols ✓
  - ② Eliminate  $\epsilon$  productions ( $L \rightarrow L \setminus \{\epsilon\}$ )  
- Find "nullable" symbols  $A \xrightarrow{*} \epsilon$ , and replace them  
if  $A \rightarrow BC$  are  $B, C$  are nullable,  
then add  $A \rightarrow BC \mid B \mid C$   
throw away  $A \rightarrow \epsilon$  productions
  - ③ Eliminate unit productions
  - ④ Convert to Chomsky Normal Form by splitting larger productions, adding new variables