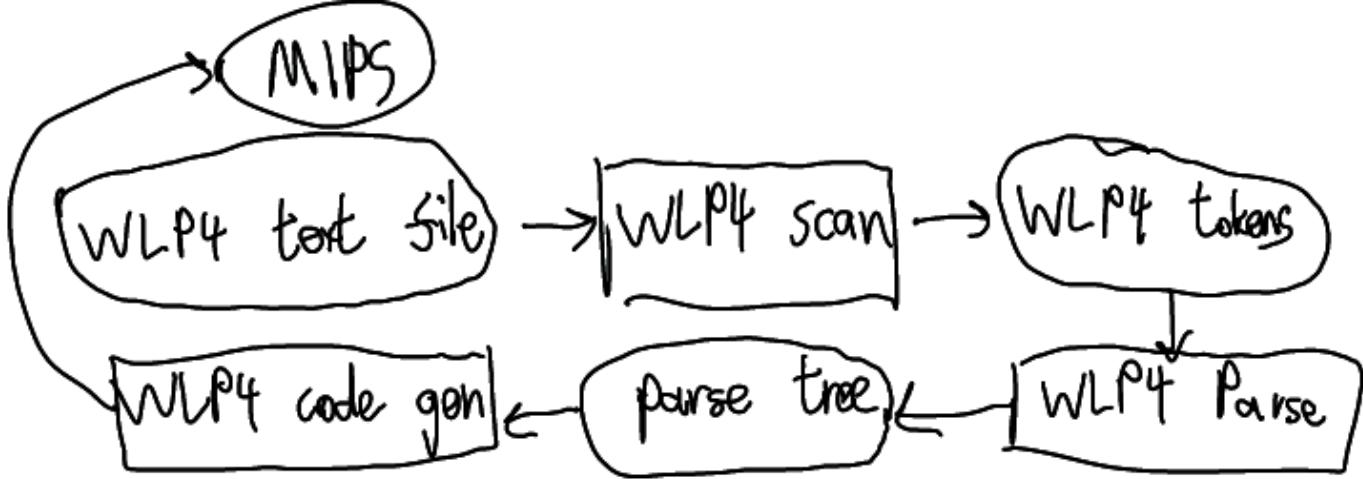


Lecture 13+14

Context-Free Languages

CS 241: Foundations of Sequential Programs
Fall 2014

Troy Vasiga et al
University of Waterloo



Big Picture of Compilation

Big picture:

- ▶ lexical analysis → scanning + tokenizing
- ▶ syntactic analysis
↳ parsing
- ▶ context-sensitive (semantic) analysis
↳ meaning
- ▶ synthesis (code generation)

Other points:

- ▶ staging can improve error messages
- ▶ staging uses the "right tool for the job"
- ▶ doing extra-work at an early stage is possible but over-complicating

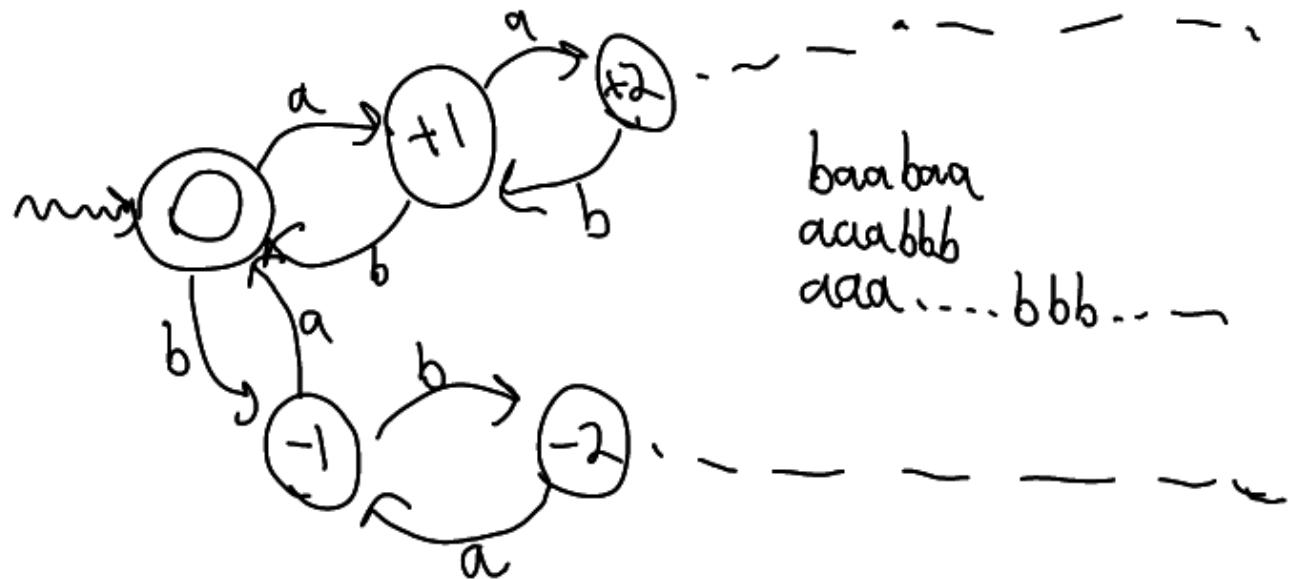
Non-regular languages

Give a DFA over $\Sigma = \{a, b\}$ for

finite

$L = \{w : \text{numbers of } a's \text{ in } w = \text{number of } b's \text{ in } w\}$

$\epsilon \in L$ abba $\in L$ aba $\notin L$



N states

3

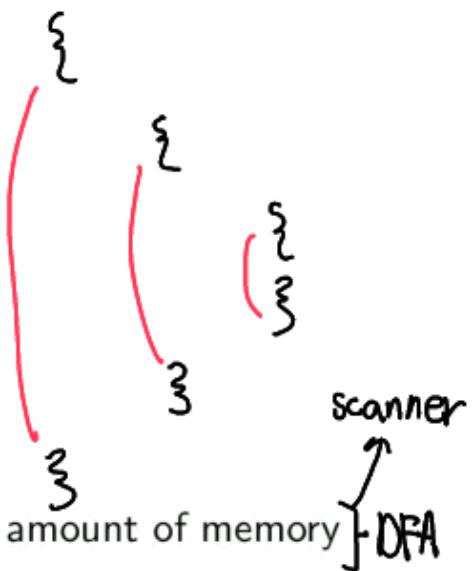
aa...a bb...b
 $N+3$ $N+3$

Proof: CS360

Context-Free Languages

- ▶ L is not a regular language, but it is *context-free*
- ▶ Note that a compiler needs to do this and then some:

int wahr(int a, int b)
if ((aa) (b bb))



- ▶ Context-free languages are built from:
 - ▶ finite sets
 - ▶ concatenation
 - ▶ union
 - ▶ recursion
- ▶ recognizers for regular languages use a finite amount of memory
- ▶ recognizers for context-free languages use a finite amount of memory plus one (unbounded) stack

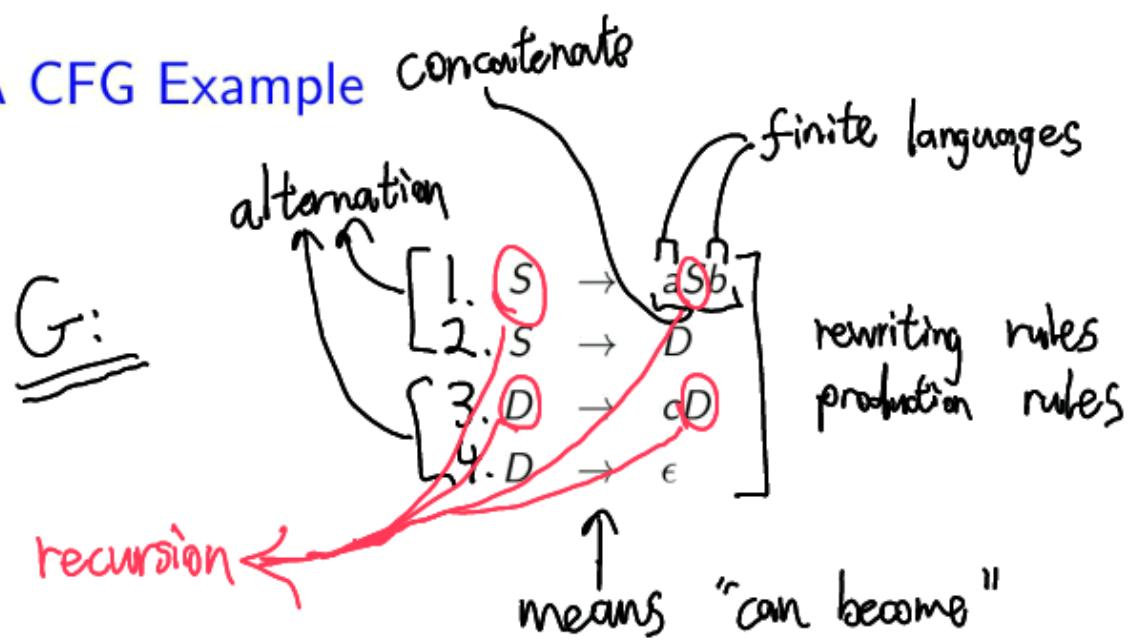
Context-Free Grammars



A way to specify a context-free language.

He jumps high. ✓
↑ ↑ ↑
Noun verb adverb
Subject

A CFG Example



word $acb \in L(G)$

Derivation:

$$\overline{S} \xrightarrow{\textcircled{1}} aSb \xrightarrow{\textcircled{2}} aDb \xrightarrow{\textcircled{3}} acDb \xrightarrow{\textcircled{3}} accDB \\ \xrightarrow{\textcircled{4}} acc\beta$$

Discussion of example

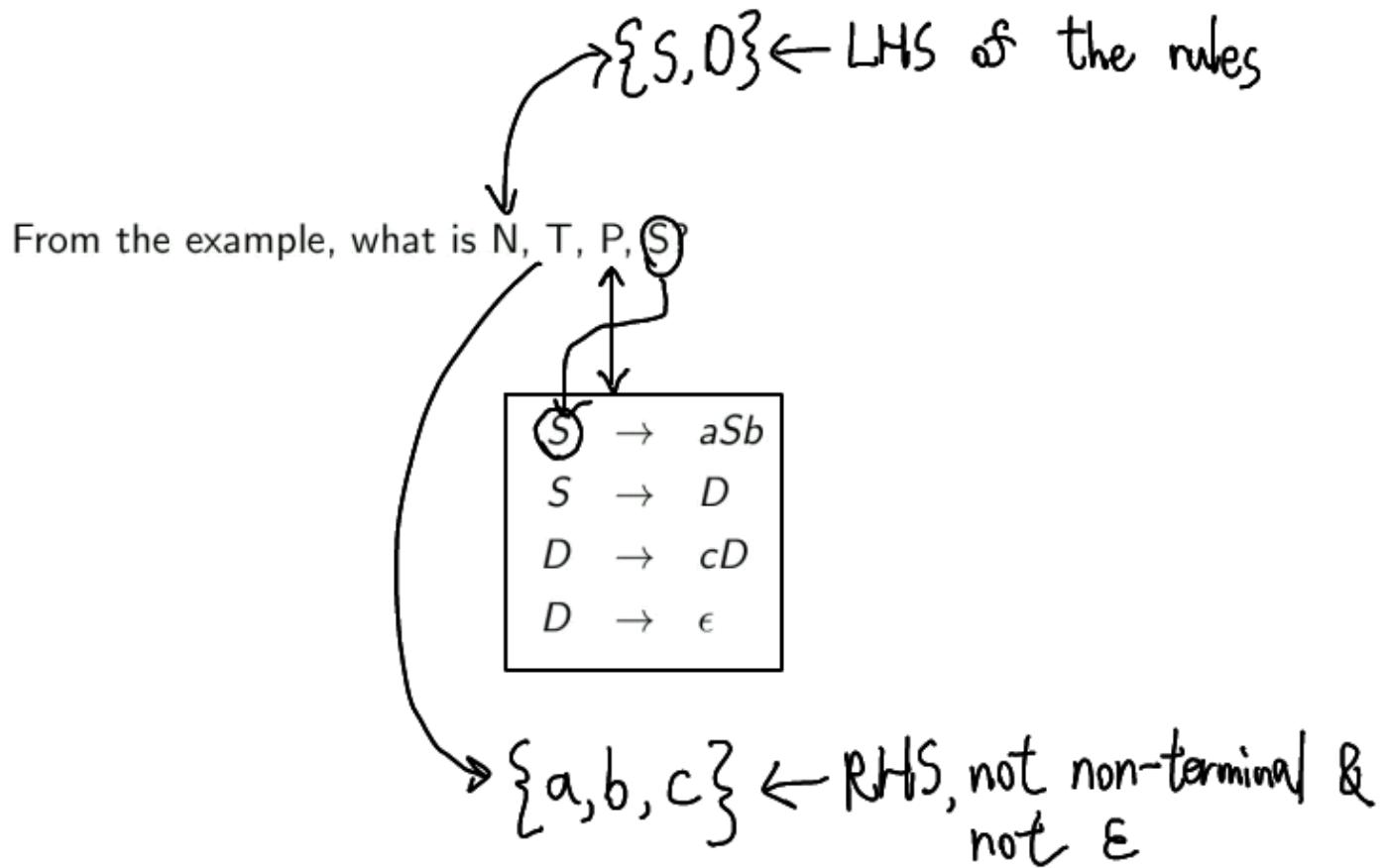
- ▶ G : a context free grammar
- ▶ $L(G)$: the language specified by CFG G
 ↳ set of words
- ▶ a word:
 a sequence of tokens that can be derived by CFG G
- ▶ a derivation:
 (informally) a sequence of rewriting steps from G
- ▶ alternation and concatenation
- ▶ recursion vs. repetition
 ↳ is strictly more powerful than

Formal Definition

A context-free grammar (CFG) consists of

- ▶ N - a finite set of non-terminals
↳ "not the ending"
- ▶ T - a finite set of terminals
↳ "ending"
- ▶ P - a finite set of production rules
 $A \rightarrow B, A \in N, B \in (N \cup T)^*$
- ▶ S - start symbol, $S \in N$
by convention: S is the LHS of the first rule

Example, more formally



Example: balanced parentheses

Example words:

ϵ , $()$, $(())$, $()()$,

- CFG:
1. $S \rightarrow (S)$
 2. $S \rightarrow \epsilon$
 3. $S \rightarrow SS$

A sample derivation: $(()())$

$$S \xrightarrow{①} (S) \xrightarrow{③} (SS) \xrightarrow{①} ((S)S) \xrightarrow{①} ((S)(S)) \xrightarrow{②} ((S)(S)) \\ \xrightarrow{②} (()())$$

10

Binary expressions

Words are composed of binary numbers (no leading zeros, other than 0) with + or - signs in infix notation.

Example words: 1001, 10+1, 11-11110+0

CFG:

1. $E \rightarrow E + E$
2. $E \rightarrow E - E$
3. $E \rightarrow B$
4. $B \rightarrow \emptyset$
5. $B \rightarrow D$
6. $D \rightarrow 1$
7. $D \rightarrow D\emptyset$
8. $D \rightarrow D1$

$E = \text{expression}$
 $B = \text{binary } \#$
 $D = \text{digit (non-zero)}$



Derivations

Two derivations of 10+1 using the previous grammar:

Leftmost derivation

$$\begin{aligned} E &\xrightarrow{\textcircled{1}} E+E \\ &\xrightarrow{\textcircled{2}} B+E \\ &\Rightarrow D+E \\ &\Rightarrow D\emptyset+E \\ &\Rightarrow \emptyset+E \\ &\Rightarrow \emptyset+B \\ &\Rightarrow \emptyset+D \\ &\Rightarrow \emptyset+1 \end{aligned}$$

Rightmost derivation

$$\begin{aligned} E &\xrightarrow{\textcircled{1}} E+E \\ &\xrightarrow{\textcircled{2}} E+B \\ &\Rightarrow E+D \\ &\Rightarrow E+1 \\ &\Rightarrow B+1 \\ &\Rightarrow D+1 \\ &\Rightarrow D\emptyset+1 \\ &\Rightarrow \emptyset+1 \end{aligned}$$

Formal Definitions

We say that $\alpha A \beta$ directly derives $\alpha \gamma \beta$ if there exists a production $A \rightarrow \gamma$.
Also called a derivation step.

$$A \in N$$

$$\alpha, \beta, \gamma \in (N \cup T)^*$$

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$

We say that $\alpha A \beta$ derives $\alpha \gamma \beta$ if

$$\alpha A \beta \Rightarrow \alpha \theta_1 \beta \Rightarrow \alpha \theta_2 \beta \Rightarrow \dots \Rightarrow \underline{\alpha \gamma \beta}$$

Write this as

$$\alpha A \beta \xrightarrow{*} \alpha \gamma \beta$$

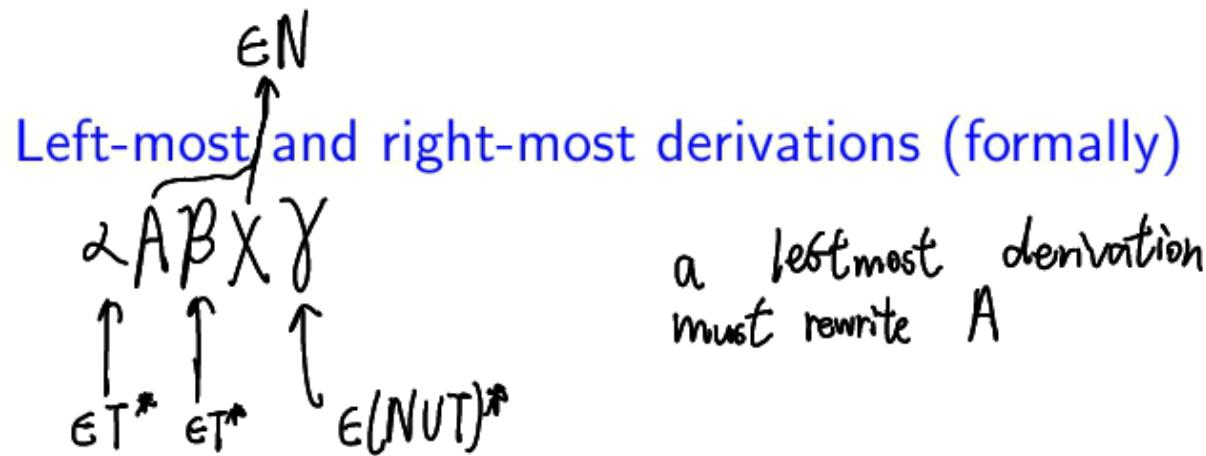
Formal Definitions (continued)

G derives $w \in T^*$ if $S \Rightarrow^* w$

$L(G) = \{w \in T^* : S \Rightarrow^* w\}$
↑ language specified by G

L is context-free if

\exists a CFG $\cdot L(G) = L$



a leftmost derivation
must rewrite A

Derivations as Proofs

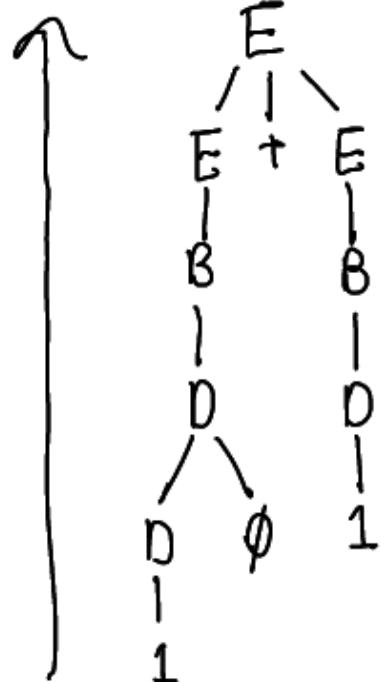
Given a CFG G & word w

$$S \Rightarrow \begin{matrix} : \\ : \\ \Rightarrow w \end{matrix} \left. \begin{matrix} \\ \\ \end{matrix} \right\} \text{constructive}$$

Parse Trees

Example:

$10+1$
using
 G_B



Internal nodes:

Non-terminal

Leaf nodes:

Terminals

Discussion:

All derivatives of $10+1$ are in this parse tree

Meaning of a parse tree

"Recursive Descent Parsing"

"Syntax Directed Translation"

$$E \rightarrow B \Leftrightarrow E.\text{val} = B.\text{val}$$

$$D_1 \rightarrow D_2 \emptyset \Leftrightarrow D_1.\text{val} = D_2.\text{val} * 2$$

$$D \rightarrow 1 \Leftrightarrow D.\text{val} = 1$$

$|0 + 1|$ means 3

Problems that grammars can encounter: Ambiguity

A real-world example of ambiguity:

Sally was given a book by Joyce.

Ambiguity in CFGs: Formal Definition

- ▶ A string x is ambiguous if...

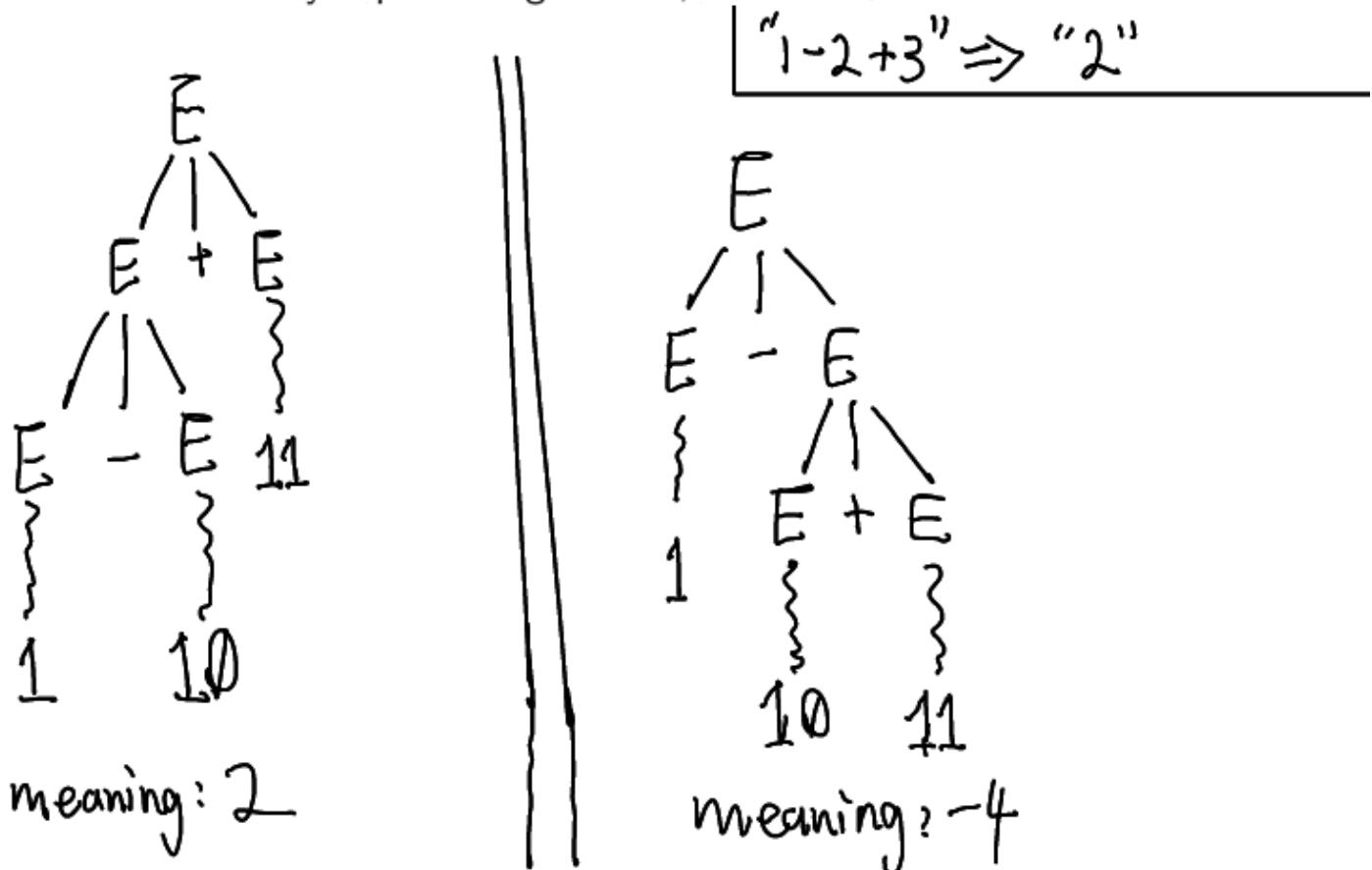
there is more than one
parse tree for x .

- ▶ A CFG G is ambiguous if...

there is at least one string
 $x \in L(G)$ which is ambiguous.

Problems that grammars can encounter: Ambiguity

Consider the binary expression grammar, and 1-10+11



Problems that grammars can encounter: Ambiguity

Terminology:

also may have:
 ≥ 2 leftmost derivations
 ≥ 2 rightmost derivations
for the same word

Left-recursion
+
Right-recursion
↓
ambiguous

Rightmost symbol of RHS
is the LHS
↑
Right recursion
↑

Fixing ambiguity

Let's rewrite our binary expression grammar differently.

$$E \rightarrow E + E$$

$$\left. \begin{array}{l} 1. E \rightarrow B + E \\ 2. E \rightarrow B - E \end{array} \right]$$

$$3. E \rightarrow B$$

$$4. E \rightarrow O$$

$$5. B \rightarrow D$$

$$6. D \rightarrow I$$

$$7. D \rightarrow DD$$

$$8. D \rightarrow DI$$

leftmost symbol of RHS
is the LHS
↓

Same as before

unambiguous

Ambiguity in programming languages

- ▶ Some programming languages have ambiguous grammars

- ▶ Pascal has a “dangling else”

The diagram shows a hand-drawn state transition graph. It starts with a rounded rectangle containing the word "if". From this node, two arrows branch out: one labeled "A" leading to another "if" node, and another labeled "B" leading to a "then" node. From the "if" node, an arrow labeled "else" points back to the original "if" node, creating a loop.

- ▶ “Fixed” by way of a footnote

RTFM

Associativity ← right associative

1-10+11

the
parse
tree

only one
parse tree :)

... it's the wrong
tree :(



means: -4

Fixing associativity problems

1. $E \rightarrow E + B \quad \leftarrow \text{left-recursion} \Leftrightarrow \text{left-associative}$

2. $E \rightarrow E - B$

3.

4.

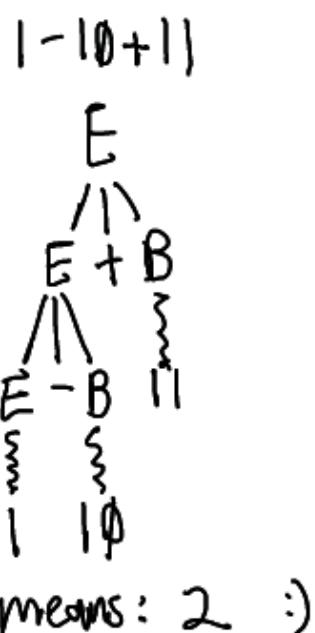
5.

6.

7.

8.

same
as
before



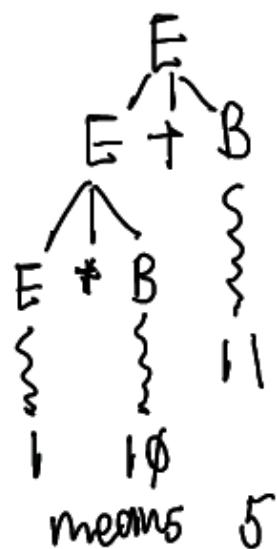
Precedence Problems

Adding multiplication to the grammar.

$$0. E \rightarrow E * B$$

Consider $\underline{\underline{1}} * 10 + 11$.

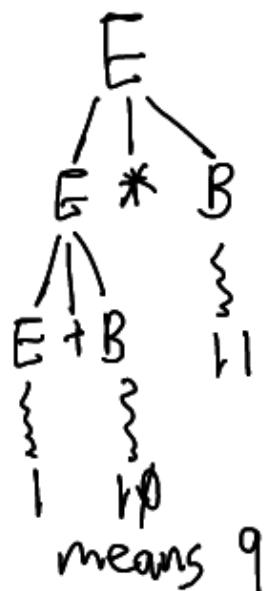
means: 5



Consider $1 + 10 * 11$

means 7

27



Fixing precedence problems: grammar

weaken the root \Rightarrow later the evaluation
 \Rightarrow lower precedence

$$S \rightarrow S + P$$

$$S \rightarrow S - P$$

$$S \rightarrow P$$

$$P \rightarrow P * B$$

$$P \rightarrow P / B$$

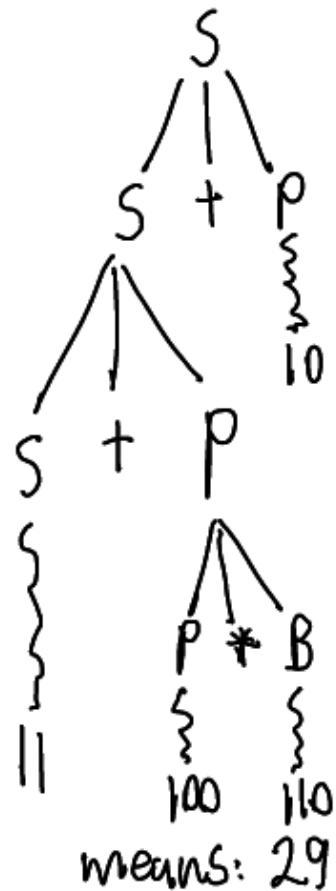
$$P \rightarrow B$$

$$B \rightarrow \dots \text{ (as before)}$$

$S = \text{sum}$
 $P = \text{product}$

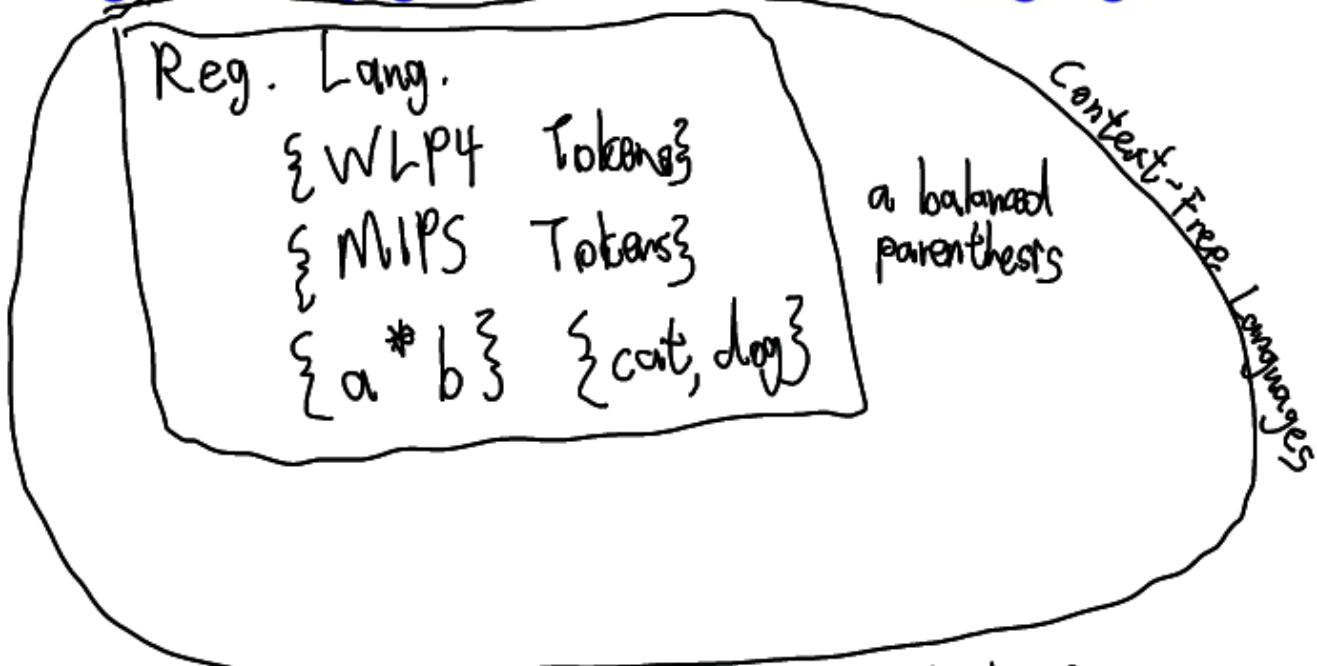
Fixing precedence problems: parse tree

$$11 + 100 * 110 + 10 \Rightarrow \text{means: } 29$$



Ex. add ()
to the grammar
 $(10+11)*10$

Regular Languages vs. Context-Free Languages



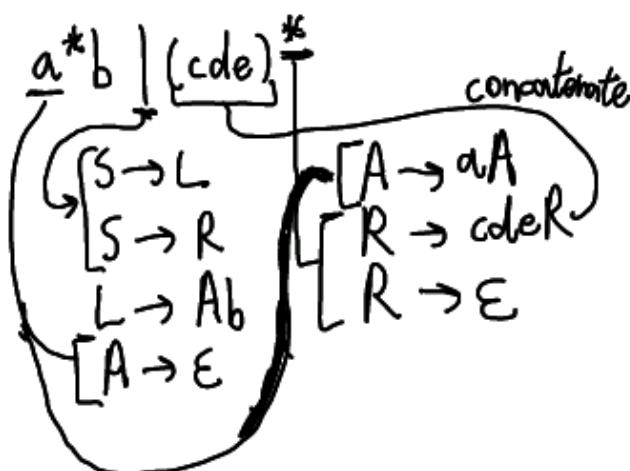
Every regular language is context-free.

"Proof": \exists regular expression for the reg. lang.

e.g. $a^*b | (cde)^*$ \leftarrow write a CFG

30

for this



\therefore All reg. lang. are context free.

Assignment Tips

- ▶ WLP4 grammar ← read it
- ▶ .cfg file format ← read the spec
- ▶ recursion is your friend
 - CFG has recursive rule

Summary

- CFG are a way to specify a language
 - ↓
 - describe the words
in the language
 $L(G)$
- To recognize \Rightarrow would like a derivation
 - [Is $w \in L(G)$? to answer, I need $S \Rightarrow^* w$
 - like to do this automatically