

Formal Languages - 2

- Context-free PLs
- Grammars
 - Derivation
 - Parsing and parse trees
 - Ambiguity
 - Precedence and Associativity
- Deterministic parsing techniques
 - TD parsing - LL(1)
 - First and Follow sets
 - Parse table construction

Context-free PLs

- Describe most of the constructs in real PLs
- Form of rules
 - Each lhs contains one nonterminal
 - Each rhs contains a sequence of terminals and/or nonterminals
- PLs describable by context-free grammars are recognized by push-down automata (analogous to an FSA with a stack)

Context-free Grammars

- Can be used to generate correct sentences in the PL (*derivation*)
- Can be used to recognize syntactically correct sentences in the PL (*parse*)
 - Can be automated efficiently in a compiler (*LL or LR parsing*)
 - Is insufficient to describe all constructs of a real PL
 - E.g., type checking with declaration

Derivation

- 1 $\langle \text{letter} \rangle ::= a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z$
- 2 $\langle \text{digit} \rangle ::= 0|1|2|3|4|5|6|7|8|9$
- 3 $\langle \text{identifier} \rangle ::= \langle \text{letter} \rangle | \langle \text{identifier} \rangle \langle \text{letter} \rangle | \langle \text{identifier} \rangle \langle \text{digit} \rangle$
- 4 $\langle 0\text{-assign-stmt} \rangle ::= \langle \text{identifier} \rangle = 0$

Can we generate $x2 = 0$ from these rules?

- $$\begin{aligned} \langle 0\text{-assign-stmt} \rangle &\rightarrow 4 \quad \langle \text{identifier} \rangle = 0 && \text{sentential form} \\ &\rightarrow 3c \quad \langle \text{identifier} \rangle \langle \text{digit} \rangle = 0 \\ &\rightarrow 3a \quad \langle \text{letter} \rangle \langle \text{digit} \rangle = 0 \\ &\rightarrow 1 \ x \quad \langle \text{digit} \rangle = 0 \\ &\rightarrow 2 \quad x \ 2 = 0 && \text{sentence} \end{aligned}$$

YES! *leftmost* or *canonical derivation*.

Parse

- 1 $\langle \text{letter} \rangle ::= a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z$
 2 $\langle \text{digit} \rangle ::= 0|1|2|3|4|5|6|7|8|9$
 3 $\langle \text{identifier} \rangle ::= \langle \text{letter} \rangle | \langle \text{identifier} \rangle \langle \text{letter} \rangle | \langle \text{identifier} \rangle \langle \text{digit} \rangle$
 4 $\langle \text{O-assign-stmt} \rangle ::= \langle \text{identifier} \rangle = 0$

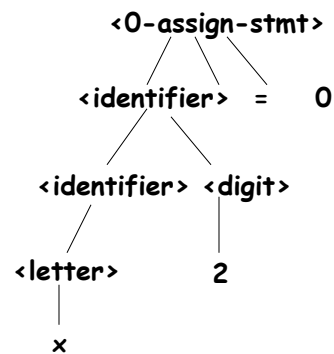
Can we recognize $x2 = 0$ as belonging to this PL?

- $x2 = 0 \rightarrow \langle \text{letter} \rangle 2 = 0$ rule 1
 $\rightarrow \langle \text{identifier} \rangle 2 = 0$ rule 3a
 $\rightarrow \langle \text{identifier} \rangle \langle \text{digit} \rangle = 0$ rule 2
 $\rightarrow \langle \text{identifier} \rangle = 0$ rule 3c
 $\rightarrow \langle \text{O-assign-stmt} \rangle$ rule 4

A *parse* of the sentence $x2 = 0$.

Parse Tree

- $x2 = 0 \rightarrow \langle \text{letter} \rangle 2 = 0$ rule 1
 $\rightarrow \langle \text{identifier} \rangle 2 = 0$ rule 3a
 $\rightarrow \langle \text{identifier} \rangle \langle \text{digit} \rangle = 0$ rule 2
 $\rightarrow \langle \text{identifier} \rangle = 0$ rule 3c
 $\rightarrow \langle \text{O-assign-stmt} \rangle$ rule 4



In parse tree, each internal node is a nonterminal; its children are the rhs of a rule for that nonterminal. Frontier of the tree is a sentence or valid PL construct.

Grammars are not Unique

1 $\langle \text{letter} \rangle ::= a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z$

2 $\langle \text{digit} \rangle ::= 0|1|2|3|4|5|6|7|8|9$

3' $\langle \text{id} \rangle ::= \langle \text{letter} \rangle | \langle \text{id} \rangle \langle \text{letterordigit} \rangle$

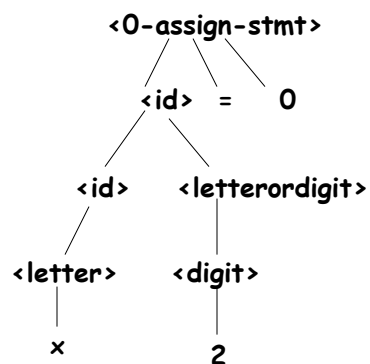
4' $\langle 0\text{-assign-stmt} \rangle ::= \langle \text{id} \rangle = 0$

5' $\langle \text{letterordigit} \rangle ::= \langle \text{letter} \rangle | \langle \text{digit} \rangle$

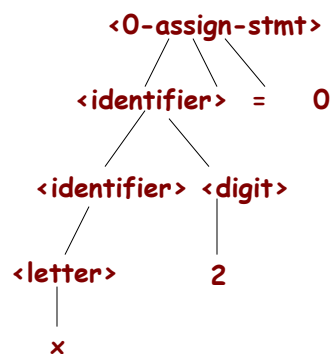
This grammar generates the same language (i.e, set of trees whose frontiers are the same), but has different parse trees than the previous grammar.

Example

2nd grammar tree



1st grammar tree



Many grammars can correspond to 1 PL, but only 1 PL should correspond to any *useful* grammar!

Definitions - Review

- *Grammar*
 - <finite set of terminals, non-terminals, production rules, special symbol>
- *Context-free grammar*
 - corresponds to PLs whose rules have only 1 nonterminal on the lhs
- *Sentence*
 - a finite sequence of terminals, constructed according to the rules of the grammar for that PL
- *Sentential form*
 - a finite sequence of terminals and non-terminals, constructed according to the rules of the grammar for that PL
- *Derivation*
 - A step by step procedure that substitutes righthandsides of productions for the nonterminal on their left, eventually leading to a sequence of terminals that is a sentence in a PL.
- *Parse* (basically a reverse derivation)

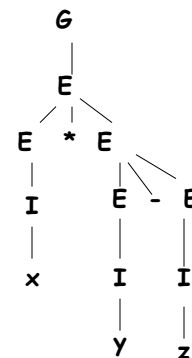
Formal-2, CS5314, © BGRyder

9

Ambiguity in PL Definition

$1 \ G ::= E$ $2 \ E ::= E - E$ $3 \ E ::= E * E$ $4 \ E ::= I$
 $5 \ I ::= a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z$

G	$\rightarrow 1$	E
	$\rightarrow 3$	$E * E$
	$\rightarrow 4$	$I * E$
	$\rightarrow 5$	$x * E$
	$\rightarrow 2$	$x * E - E$
	$\rightarrow 4$	$x * I - E$
	$\rightarrow 5$	$x * y - E$
	$\rightarrow 4$	$x * y - I$
	$\rightarrow 5$	$x * y - z$



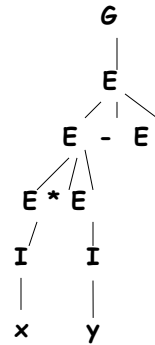
Formal-2, CS5314, © BGRyder

10

Ambiguity in PL Definition

$1 G ::= E$ 2 3 4
 $E ::= E - E \mid E * E \mid I$
 $5 I ::= a \mid b \mid c \mid d \mid e \mid f \mid g \mid h \mid i \mid j \mid k \mid l \mid m \mid n \mid o \mid p \mid q \mid r \mid s \mid t \mid u \mid v \mid w \mid x \mid y \mid z$

$G \rightarrow 1$ E
 $\rightarrow 2$ $E - E$
 $\rightarrow 3$ $E * E - E$
 $\rightarrow 4$ $I * E - E$
 $\rightarrow 5$ $x * E - E$
 $\rightarrow 4$ $x * I - E$
 $\rightarrow 5$ $x * y - E$
 $\rightarrow 4$ $x * y - I$
 $\rightarrow 5$ $x * y - z$

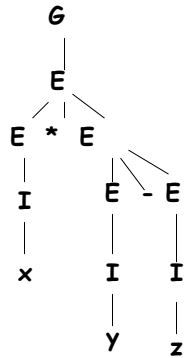


Formal-2, CS5314, © BGRyder

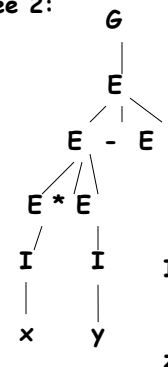
11

Comparison

Tree 1:



Tree 2:



*Which tree is correct?
 Can we rewrite the grammar to only
 generate one of them?*

Formal-2, CS5314, © BGRyder

12

Ambiguity

- If there are 2 different canonical derivations (or alternatively, 2 parse trees) for the same sentence then the grammar is *ambiguous*
 - Solution
 - Change grammar to reflect operator precedence
i.e., $X*Y-Z$ means $((X*Y) - Z)$
- There is no algorithm which can tell if an arbitrary context-free grammar is ambiguous
- Also no algorithm to tell if 2 arbitrary context-free grammars generate the same language
 - But can tell if 2 regular languages are equivalent!

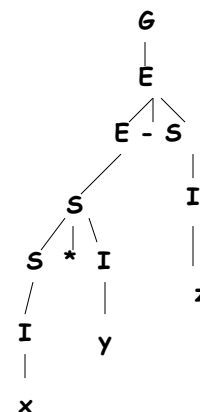
Formal-2, CS5314, © BGRyder

13

A Better Grammar

$G ::= E$
 $E ::= S \mid E - S$
 $S ::= I \mid S * I$
 $I ::= a|b|c|d|e|f|g|h|i|j|k|l|m|n$
 $\quad |o|p|q|r|s|t|u|v|w|x|y|z$

Note: since S is operand of $-$ operation, this forces $$ to have **higher precedence** than $-$.*



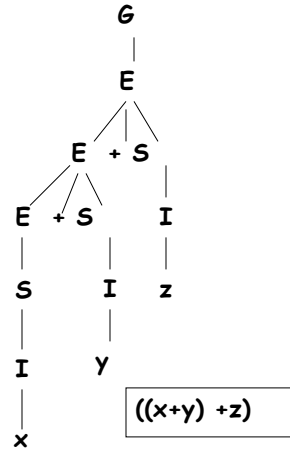
Formal-2, CS5314, © BGRyder

14

Associativity in the Grammar

$G ::= E$
 $E ::= E + S \mid S$
 $S ::= I \mid S * I$
 $I ::= a|b|c|d|e|f|g|h|i|j|k|l|m|n$
 $\quad |o|p|q|r|s|t|u|v|w|x|y|z$

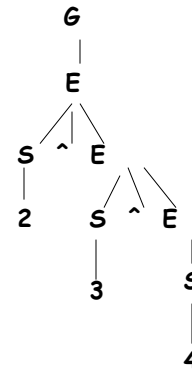
How to parse $x+y+z$?
Tree shows that $+$ is left associative because E 's rule is left recursive.



Right Associativity

$G ::= E$
 $E ::= S \wedge E \mid S$
 $S ::= 0|1|2|3|4|5|6|7|8|9$

What is 2^3^4 ?
 8^4 or 2^{81} ?



TD Parsing

Elimination of left recursion to prevent infinite loops in the parse .

$$E \rightarrow E \alpha | \beta \quad \Longrightarrow \quad \begin{array}{l} E \rightarrow \beta A \\ A \rightarrow \alpha A | \varepsilon \end{array}$$

Example:

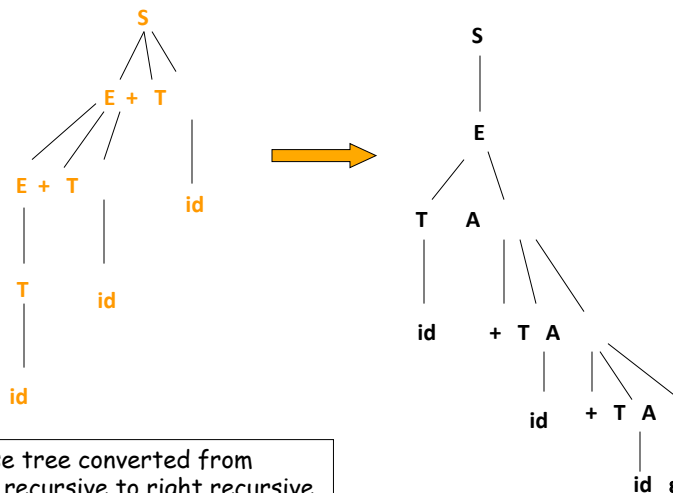
$$\begin{array}{l} S \rightarrow E \\ E \rightarrow E + T \\ E \rightarrow T \\ T \rightarrow id \end{array} \quad \Longrightarrow \quad \begin{array}{l} S \rightarrow E \\ E \rightarrow T A \\ A \rightarrow + T A | \varepsilon \\ T \rightarrow id \end{array}$$

Can also left factor the grammar removing shared prefixes of right-hand-sides.

Formal-2, CS5314, © BGRyder

17

Parse Tree



Parse tree converted from left recursive to right recursive.

Formal-2, CS5314, © BGRyder

18

TD Parsing

- *Problem*: predicting which nonterminal to expand next, from a leading string of symbols
- *Idea*: generate parse tree top down so its frontier is always a sentential form
 - Use **First** and **Follow** sets to understand the shape of sentential forms possibly generated by the grammar

Formal-2, CS5314, © BGRyder

1
9

TD Stack Parser, EG

Stack	Input	Production
\$S	id+id+id\$	
\$E	id+id+id\$	$S \rightarrow E$
\$A T	id+id+id\$	$E \rightarrow T A$
\$A	+id+id\$	$T \rightarrow id$
\$A T	id+id\$	$A \rightarrow + T A$
\$A	+id\$	$T \rightarrow id$
\$A T	id\$	$A \rightarrow + T A$
\$A	\$	$T \rightarrow id$
\$	\$	$A \rightarrow \epsilon$

$S \rightarrow E$
$E \rightarrow T A$
$A \rightarrow + T A \mid \epsilon$
$T \rightarrow id$

See algm in ASU Fig 4.14, p 187

Formal-2, CS5314, © BGRyder

20

How to mechanize?

- Define α to be string of non-terminals and terminals
- $\text{First}(\alpha)$ is the set of terminals that begin strings derivable from α .
If $\alpha \xRightarrow{*} \epsilon$, then ϵ is in $\text{First}(\alpha)$.
- $\text{Follow}(A)$ is the set of terminals that can appear directly to the right of A in a sentential form
 $S \xRightarrow{*} \alpha A a \beta$ means a is in $\text{Follow}(A)$.
If A can be rightmost symbol in a sentential form, that is,
 $X \xRightarrow{*} \alpha A \delta$ where $\delta \xRightarrow{*} \epsilon$, then
 $\text{Follow}(X) \subseteq \text{Follow}(A)$ because whatever can follow an X can follow an A too.

Formal-2, CS5314, © BGRyder

2
1

Example

- $\text{First}(S) = \text{First}(E) = \text{First}(T) = \{id\}$
 - $\text{First}(A) = \{+, \epsilon\}$
 - $\text{Follow}(S) = \text{Follow}(E) = \text{Follow}(A) = \{\$\}$
 - $\text{Follow}(T) = \{+, \$\}$
- $$S \rightarrow E$$
- $$E \rightarrow T A$$
- $$A \rightarrow + T A \mid \epsilon$$
- $$T \rightarrow id$$

Formal-2, CS5314, © BGRyder

22

LL(k) Grammars

- Can choose next production to expand by during TD phase, by looking k symbols ahead into input
- Use **First sets** to choose production
- Use **Follow sets** to handle ϵ cases

Formal-2, CS5314, © BGRyder

23

Example: LL(1)

Nonterms\Inputs:	id	+	\$
S	$S \rightarrow E$		
E	$E \rightarrow T A$		
T	$T \rightarrow id$		
A		$A \rightarrow + T A$	$A \rightarrow \epsilon$

Ambiguous or left recursive grammars result in multiply defined entries in table - a problem!

$\text{First}(S) = \text{First}(E) = \text{First}(T) = \{id\}$
 $\text{First}(A) = \{+, \epsilon\}$
 $\text{Follow}(S) = \text{Follow}(E) = \text{Follow}(A) = \{\$, \}$
 $\text{Follow}(T) = \{+, \$\}$

$S \rightarrow E$
 $E \rightarrow T A$
 $A \rightarrow + T A \mid \epsilon$
 $T \rightarrow id$

Formal-2, CS5314, © BGRyder

24