

# CMSC 330: Organization of Programming Languages

## Context-Free Grammars

### Motivation (cont'd)

- We want to describe program structure precisely
- Regular expressions are not enough
  - No regular expression for balanced pairs of ( )'s
    - {"()", "(())", "()()", ...} is not a regular language
- Instead, we'll use **context-free grammars**
  - These are *almost* enough for C, C++, Java

CMSC 330

3

### Context-Free Grammars (CFGs)

- But CFGs can do a lot more!
  - $S \rightarrow (S) | \epsilon$  // generates balanced pairs of ( )'s
- In fact, CFGs subsume REs, DFAs, NFAs
  - There is a CFG that generates any regular language
  - But REs are a better notation for regular languages
- CFGs can specify programming language syntax
  - CFGs (mostly) describe the parsing process

CMSC 330

5

### Motivation

- Programs are just strings of text
  - But they're strings that have a certain structure
- Informal description of syntax of a C program
  - A C program is a list of **declarations** and **definitions**
  - A **function definition** contains **parameters** and a **body**
  - A **function body** is a sequence of **statements**
  - A **statement** is an **expression**, **if**, **goto**, etc.
  - An **expression** may be **assignment**, **addition**, **subtraction**, etc

CMSC 330

2

### Context Free Grammar (CFG)

- A way of generating sets of strings or languages
- Grammar:  $S \rightarrow 0S | 1S | \epsilon$ 
  - Means every **S** may be replaced by **0S**, **1S**, or  **$\epsilon$**
  - Example
    - $S \Rightarrow 0S$  // using  $S \rightarrow 0S$
    - $\Rightarrow 01S$  // using  $S \rightarrow 1S$
    - $\Rightarrow 011S$  // using  $S \rightarrow 1S$
    - $\Rightarrow 011$  // using  $S \rightarrow \epsilon$
- Grammar is same as regular expression  $(0|1)^*$ 
  - Generates / accepts the same set of strings

CMSC 330

4

### Formal Definition

- A context-free grammar **G** is a 4-tuple:
  - $\Sigma$  – a finite set of *terminal* or *alphabet* symbols
    - Often written in lowercase
  - $N$  – a finite, nonempty set of *nonterminal* symbols
    - Often written in uppercase
    - It must be that  $N \cap \Sigma = \emptyset$
  - $P$  – a set of *productions* of the form  $N \rightarrow (\Sigma|N)^*$ 
    - Informally this means that the nonterminal can be replaced by the string of zero or more terminals or nonterminals to the right of the  $\rightarrow$
    - Can think of productions as rewriting rules
  - $S \in N$  – the *start symbol*

CMSC 330

6

## Backus-Naur Form

- Context-free grammar production rules are also called Backus-Naur Form or **BNF**
  - A production like  $A \rightarrow B c D$  is written in BNF as  $\langle A \rangle ::= \langle B \rangle c \langle D \rangle$  (Non-terminals written with angle brackets and  $::=$  instead of  $\rightarrow$ )
  - Often used to describe language syntax
- BNF was designed by
  - John Backus
    - Chair of the Algol committee in the early 1960s
  - Peter Naur
    - Secretary of the committee, who used this notation to describe Algol in 1962

CM/SC 330

7

## Informal Definition of Acceptance

- A string is **accepted** by a CFG if there is
  - Some sequence of applying productions (**rewrites**) starting at the start symbol that generates the string
- Example
  - Grammar:  $S \rightarrow 0S \mid 1S \mid \epsilon$
  - Sequence generating the string **010**
    - $S \Rightarrow 0S \Rightarrow 01S \Rightarrow 010S \Rightarrow 010$
- Terminology
  - Such a sequence of rewrites is a **derivation** or **parse**
  - Discovering the derivation is called **parsing**

CM/SC 330

8

## Derivations

- Notation
  - $\Rightarrow$  indicates a derivation of one step
  - $\Rightarrow^+$  indicates a derivation of one or more steps
  - $\Rightarrow^*$  indicates a derivation of zero or more steps
- Example
  - $S \rightarrow 0S \mid 1S \mid \epsilon$
- For the string **010**
  - $S \Rightarrow 0S \Rightarrow 01S \Rightarrow 010S \Rightarrow 010$
  - $S \Rightarrow^+ 010$
  - $S \Rightarrow^* S$

CM/SC 330

9

## Practice

- Try to make a grammar which accepts
    - $0^*1^*$      $0^n1^n$  where  $n \geq 0$      $0^n1^m$  where  $m \leq n$
  - $S \rightarrow A \mid B$   
 $A \rightarrow 0A \mid \epsilon$      $S \rightarrow 0S1 \mid \epsilon$      $S \rightarrow 0S1 \mid 0S \epsilon$   
 $B \rightarrow 1B \mid \epsilon$
- Give some example strings from this language
    - $S \rightarrow 0 \mid 1S$ 
      - 0, 10, 110, 1110, 11110, ...
    - What language is it?
      - $1^*0$

CM/SC 330

10

## Example

- $S \rightarrow aS \mid T$   
 $T \rightarrow bT \mid U$   
 $U \rightarrow cU \mid \epsilon$
- A derivation:
  - $S \Rightarrow aS \Rightarrow aT \Rightarrow aU \Rightarrow acU \Rightarrow ac$ 
    - Abbreviated as  $S \Rightarrow^+ ac$
  - $S \Rightarrow T \Rightarrow U \Rightarrow \epsilon$
- Is there any derivation
  - $S \Rightarrow^+ ccc$  ?     $S \Rightarrow^+ Sa$  ?
  - $S \Rightarrow^+ bab$  ?     $S \Rightarrow^+ bU$  ?

CM/SC 330

11

## Example (cont'd)

- $S \rightarrow aS \mid T$   
 $T \rightarrow bT \mid U$   
 $U \rightarrow cU \mid \epsilon$
- Generates what language?
- Do other grammars generate this language?
  - $S \rightarrow ABC$   
 $A \rightarrow aA \mid \epsilon$   
 $B \rightarrow bB \mid \epsilon$   
 $C \rightarrow cC \mid \epsilon$
  - So grammars are not unique

CM/SC 330

12

## Example: Arithmetic Expressions (Limited)

- $E \rightarrow a \mid b \mid c \mid E+E \mid E-E \mid E^*E \mid (E)$ 
  - An expression  $E$  is either a letter  $a$ ,  $b$ , or  $c$
  - Or an  $E$  followed by  $+$  followed by an  $E$
  - etc.
- This **describes** or **generates** a set of strings
  - $\{a, b, c, a+b, a+a, a^*c, a-(b^*a), c^*(b+d)\}$
- Example strings not in the language
  - $d, c(a), a+, b^{**}c$ , etc.

CM/SC 330

13

## Example, Formally

- Formally, the grammar we just showed is
 

$\Sigma = \{ +, -, *, (, ), a, b, c \}$	// terminals
$N = \{ E \}$	// nonterminals
$P = \{ E \rightarrow a, E \rightarrow b, E \rightarrow c,$	// productions
$E \rightarrow E-E, E \rightarrow E+E,$	
$E \rightarrow E^*E, E \rightarrow (E) \}$	
$S = E$	// start symbol

CM/SC 330

14

## Uniqueness of Grammars

- Grammars are not unique
  - Different grammars can generate same set of strings
- Following grammar generates the same set of strings as the previous grammar:
 
$$\begin{aligned} E &\rightarrow E+T \mid E-T \mid T \\ T &\rightarrow T^*P \mid P \\ P &\rightarrow (E) \mid a \mid b \mid c \end{aligned}$$

CM/SC 330

15

## Notational Shortcuts

- A production is of the form
  - left-hand side (LHS)  $\rightarrow$  right hand side (RHS)
- If not specified
  - Assume LHS of first listed production is the start symbol
- Productions with the same LHS
  - Are usually combined with  $|$
- If a production has an empty RHS
  - It means the RHS is  $\epsilon$

$S \rightarrow ABC$  //  $S$  is start symb  
 $A \rightarrow aA$   
 $| b$  //  $A \rightarrow b$   
 $|$  //  $A \rightarrow \epsilon$

CM/SC 330

16

## Sentential Forms and Derivations

- A **sentential form** is a string of terminals and nonterminals produced from that start symbol
- Inductively
  - The start symbol is a sentential form for a grammar
  - If  $\alpha A \delta$  is a sentential form for a grammar, where  $\alpha$  and  $\delta \in (N \cup \Sigma)^*$ , and  $A \rightarrow \gamma$  is a production, then  $\alpha \gamma \delta$  is a sentential form for the grammar
    - In this case, we say that  $\alpha A \delta$  *derives*  $\alpha \gamma \delta$  in one step, which is written as  $\alpha A \delta \Rightarrow \alpha \gamma \delta$

CM/SC 330

17

## Sentential Forms Example

- Given grammar
  - $S \rightarrow 0S \mid 1S \mid \epsilon$
- Possible derivations
  - $S \Rightarrow 0S \Rightarrow 01S \Rightarrow 010S \Rightarrow 010$ 
    - $S, 0S, 01S, 010S, 010$  are sentential forms
  - $S \Rightarrow 1S \Rightarrow 11S \Rightarrow 111S \Rightarrow 111$ 
    - $S, 1S, 11S, 111S, 111$  are sentential forms
  - $S \Rightarrow \epsilon$ 
    - $S, \epsilon$  are sentential forms
- In other words
  - If  $S \Rightarrow^* \alpha$ , then  $\alpha$  is a sentential form

CM/SC 330

18

## The Language Generated by a CFG

- The language generated by a grammar  $G$  is

$$L(G) = \{ \omega \mid \omega \in \Sigma^* \text{ and } S \Rightarrow^+ \omega \}$$

- $S$  is the start symbol of the grammar
- $\Sigma$  is the alphabet for that grammar
- In other words
  - All sentential forms with only terminals
  - All strings over  $\Sigma$  that can be derived from the start symbol via one or more productions

CM/SC 330

19

## Parse Trees

- A parse tree shows how a string is produced by a grammar
  - Root node is the start symbol
  - Each interior node is a nonterminal
  - Children of node are symbols on r.h.s of production applied to that nonterminal
  - Leaves are all terminal symbols
- Reading the leaves left-to-right shows the string corresponding to the tree

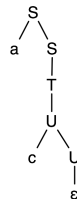
CM/SC 330

20

## Example

$$S \Rightarrow aS \Rightarrow aT \Rightarrow aU \Rightarrow acU \Rightarrow ac$$

$$\begin{aligned} S &\rightarrow aS \mid T \\ T &\rightarrow bT \mid U \\ U &\rightarrow cU \mid \epsilon \end{aligned}$$



CM/SC 330

21

## Leftmost and Rightmost Derivation

- Leftmost derivation
  - Leftmost nonterminal is replaced in each step
- Rightmost derivation
  - Rightmost nonterminal is replaced in each step
- Example
  - Grammar
    - $S \rightarrow AB, A \rightarrow a, B \rightarrow b$
  - Leftmost derivation for "ab"
    - $S \rightarrow AB \Rightarrow aB \Rightarrow ab$
  - Rightmost derivation for "ab"
    - $S \rightarrow AB \Rightarrow Ab \Rightarrow ab$

CM/SC 330

22

## Parse Tree For Derivations

- Parse tree may be same for both leftmost & rightmost derivations
  - Example Grammar:  $S \rightarrow a \mid SbS$  String:  $aba$ 
    - Leftmost Derivation
      - $S \Rightarrow SbS \Rightarrow abS \Rightarrow aba$
    - Rightmost Derivation
      - $S \Rightarrow SbS \Rightarrow Sba \Rightarrow aba$
  - Parse trees don't show order productions are applied
- Every parse tree has a unique leftmost and a unique rightmost derivation



CM/SC 330

23

## Parse Tree For Derivations (cont.)

- Not every string has a unique parse tree
  - Example Grammar:  $S \rightarrow a \mid SbS$  String:  $ababa$ 
    - Leftmost derivation
      - $S \Rightarrow SbS \Rightarrow abS \Rightarrow abSbS \Rightarrow ababS \Rightarrow ababa$
    - Another leftmost derivation
      - $S \Rightarrow SbS \Rightarrow SbSbS \Rightarrow abSbS \Rightarrow ababS \Rightarrow ababa$



CM/SC 330

24

## Ambiguity

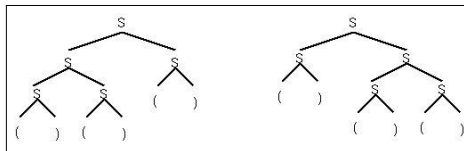
- A grammar is **ambiguous** if a string may have multiple leftmost (or rightmost) derivations
    - Equivalent to multiple parse trees
    - Can be hard to determine
1.  $S \rightarrow aS \mid T$   
 $T \rightarrow bT \mid U$                       **No**  
 $U \rightarrow cU \mid \epsilon$
2.  $S \rightarrow SS \mid () \mid (S)$                       **?**

CM/SC 330

25

## Ambiguity (cont'd)

- 2 **different** parse trees for the same string:  $()()()$
- 2 distinct leftmost derivations :
  - $S \Rightarrow SS \Rightarrow SSS \Rightarrow ()SS \Rightarrow ()()S \Rightarrow ()()()$
  - $S \Rightarrow SS \Rightarrow ()S \Rightarrow ()SS \Rightarrow ()()S \Rightarrow ()()()$



CM/SC 330

26

## More on Leftmost/Rightmost Derivations

- Is the following derivation leftmost or rightmost?
  - $S \Rightarrow aS \Rightarrow aT \Rightarrow aU \Rightarrow acU \Rightarrow ac$
  - Both! At most one non-terminal in each sentential form, so there's no choice which non-terminals to expand
- How about the following derivation?
  - $S \Rightarrow SbS \Rightarrow SbSbS \Rightarrow SbabS \Rightarrow ababS \Rightarrow ababa$
  - Neither! Selects left, center, left, and rightmost nonterminals

CM/SC 330

27

## Tips for Designing Grammars

- Use recursive productions to generate an arbitrary number of symbols
  - $A \rightarrow xA \mid \epsilon$                       Zero or more  $x$ 's
  - $A \rightarrow yA \mid y$                       One or more  $y$ 's
- Use separate nonterminals to generate disjoint parts of a language, and then combine in a production
  - $G = S \rightarrow AB$
  - $A \rightarrow aA \mid \epsilon$
  - $B \rightarrow bB \mid \epsilon$
  - $L(G) = a^*b^*$

CM/SC 330

28

## Tips for Designing Grammars (cont'd)

- To generate languages with matching, balanced, or related numbers of symbols, write productions which generate strings from the middle
  - $\{a^n b^n \mid n \geq 0\}$  (not a regular language!)
  - $S \rightarrow aSb \mid \epsilon$
  - Example:  $S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb$
  - $\{a^n b^{2n} \mid n \geq 0\}$
  - $S \rightarrow aSbb \mid \epsilon$

CM/SC 330

29

## Tips for Designing Grammars (cont'd)

$\{a^n b^m \mid m \geq 2n, n \geq 0\}$   
 $S \rightarrow aSbb \mid B \mid \epsilon$   
 $B \rightarrow bB \mid b$

The following grammar also works:

$S \rightarrow aSbb \mid B$   
 $B \rightarrow bB \mid \epsilon$

How about the following?

$S \rightarrow aSbb \mid bS \mid \epsilon$

CM/SC 330

30

## Tips for Designing Grammars (cont'd)

$\{a^n b^m a^{n+m} \mid n \geq 0, m \geq 0\}$

Rewrite as  $a^n b^m a^n$ , which now has matching superscripts (two pairs)

Would this grammar work?

$S \rightarrow aSa \mid B$   
 $B \rightarrow bBa \mid ba$

Corrected:

$S \rightarrow aSa \mid B$   
 $B \rightarrow bBa \mid \epsilon$

The outer  $a^n a^n$  are generated first, then the inner  $b^m a^m$

CM/SC 330

31

## Tips for Designing Grammars (cont'd)

4. For a language that's the union of other languages, use separate nonterminals for each part of the union and then combine

$\{a^n(b^n|c^m) \mid m > n \geq 0\}$

Can be rewritten as

$\{a^n b^m \mid m > n \geq 0\} \cup$   
 $\{a^n c^m \mid m > n \geq 0\}$

CM/SC 330

32

## Tips for Designing Grammars (cont'd)

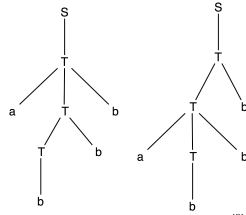
$\{a^n b^m \mid m > n \geq 0\} \cup \{a^n c^m \mid m > n \geq 0\}$

$S \rightarrow T \mid U$   
 $T \rightarrow aTb \mid Tb \mid b$   
 $U \rightarrow aUc \mid Uc \mid c$

T generates the first set

U generates the second set

- What about the string **abbb**?
  - Ambiguous!



CM/SC 330

33

## Tips for Designing Grammars (cont'd)

$\{a^n b^m \mid m > n \geq 0\} \cup \{a^n c^m \mid m > n \geq 0\}$

Will this fix the ambiguity?

$S \rightarrow T \mid U$   
 $T \rightarrow aTb \mid bT \mid b$   
 $U \rightarrow aUc \mid cU \mid c$

- It's not ambiguous, but it can generate invalid strings such as **babb**

CM/SC 330

34

## Tips for Designing Grammars (cont'd)

$\{a^n b^m \mid m > n \geq 0\} \cup \{a^n c^m \mid m > n \geq 0\}$

Unambiguous version

$S \rightarrow T \mid V$   
 $T \rightarrow aTb \mid U$   
 $U \rightarrow Ub \mid b$   
 $V \rightarrow aVc \mid W$   
 $W \rightarrow Wc \mid c$

CM/SC 330

35

## CFGs for Languages

- Recall that our goal is to describe programming languages with CFGs

- We had the following example which describes limited arithmetic expressions

$E \rightarrow a \mid b \mid c \mid E+E \mid E-E \mid E^*E \mid (E)$

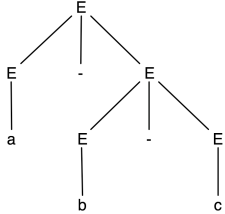
- What's wrong with using this grammar?
  - It's ambiguous!

CM/SC 330

36

## Example: a-b-c

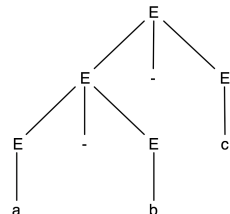
$E \Rightarrow E-E \Rightarrow a-E \Rightarrow a-E-E \Rightarrow a-b-E \Rightarrow a-b-c$



Corresponds to  $a-(b-c)$

CMSC 330

$E \Rightarrow E-E \Rightarrow E-E-E \Rightarrow a-E-E \Rightarrow a-b-E \Rightarrow a-b-c$

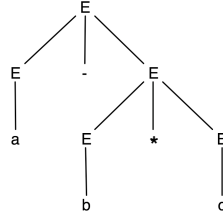


Corresponds to  $(a-b)-c$

37

## Example: a-b\*c

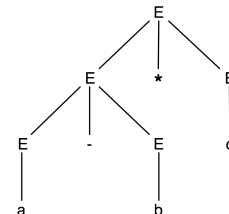
$E \Rightarrow E-E \Rightarrow a-E \Rightarrow a-E^*E \Rightarrow a-b^*E \Rightarrow a-b^*c$



Corresponds to  $a-(b^*c)$

CMSC 330

$E \Rightarrow E-E \Rightarrow E-E^*E \Rightarrow a-E^*E \Rightarrow a-b^*E \Rightarrow a-b^*c$



Corresponds to  $(a-b)^*c$

38

## Another Example: If-Then-Else

$\langle \text{stmt} \rangle ::= \langle \text{assignment} \rangle \mid \langle \text{if-stmt} \rangle \mid \dots$

$\langle \text{if-stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \langle \text{stmt} \rangle \mid \text{if } \langle \text{expr} \rangle \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

– (Here  $\langle \rangle$ 's are used to denote nonterminals and  $::=$  for productions)

- Consider the following program fragment:

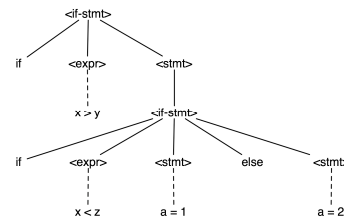
```
if (x > y)
  if (x < z)
    a = 1;
  else a = 2;
```

– Note: Ignore newlines

CMSC 330

39

## Parse Tree #1

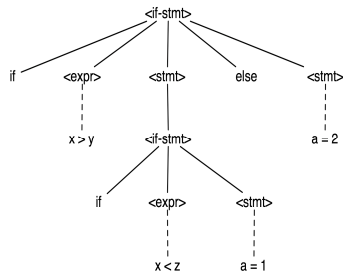


- Else belongs to inner if

CMSC 330

40

## Parse Tree #2



- Else belongs to outer if

CMSC 330

41

## Dealing With Ambiguous Grammars

- Ambiguity is bad
  - Syntax is correct
  - But semantics differ depending on choice
    - Different associativity  $(a-b)-c$  vs.  $a-(b-c)$
    - Different precedence  $(a-b)^*c$  vs.  $a-(b^*c)$
    - Different control flow  $\text{if (if else) vs. if (if) else}$

- Two approaches
  - Rewrite grammar
  - Use special parsing rules
    - Depending on parsing method (learn in CMSC 430)

CMSC 330

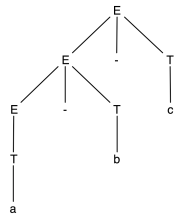
42

## Fixing the Expression Grammar

- Idea: Require that the right operand of all of the operators not have an operator in it, unless it's parenthesized

$E \rightarrow E+T \mid E-T \mid E^*T \mid T$   
 $T \rightarrow a \mid b \mid c \mid (E)$

- Now only one parse tree for  $a-b-c$ 
  - Left associative
  - Exercise: Give a derivation for the string  $a-(b-c)$



CM/SC 330

43

## What if We Wanted Right-Associativity?

- Left-recursive productions are used for left-associative operators
- Right-recursive productions are used for right-associative operators

- Left:

$E \rightarrow E+T \mid E-T \mid E^*T \mid T$   
 $T \rightarrow a \mid b \mid c \mid (E)$

- Right:

$E \rightarrow T+E \mid T-E \mid T^*E \mid T$   
 $T \rightarrow a \mid b \mid c \mid (E)$

CM/SC 330

44

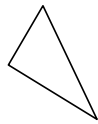
## Parse Tree Shape

- The kind of recursion/associativity determines the shape of the parse tree

left recursion



right recursion



- Exercise: draw a parse tree for  $a-b-c$  in the prior grammar in which subtraction is right-associative

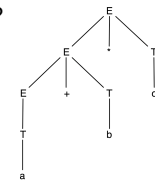
CM/SC 330

45

## A Different Problem

- How about the string  $a+b^*c$  ?

$E \rightarrow E+T \mid E-T \mid E^*T \mid T$   
 $T \rightarrow a \mid b \mid c \mid (E)$



- Doesn't have correct precedence for  $*$ 
  - When a nonterminal has productions for several operators, they effectively have the same precedence
- How can we fix this?

CM/SC 330

46

## Final Expression Grammar

$E \rightarrow E+T \mid E-T \mid T$  lowest precedence operators  
 $T \rightarrow T^*P \mid P$  higher precedence  
 $P \rightarrow a \mid b \mid c \mid (E)$  highest precedence (parentheses)

- Exercises:
  - Construct tree and left and right derivations for
    - $a+b^*c$   $a^*(b+c)$   $a^*b+c$   $a-b-c$
  - See what happens if you change the last set of productions to  $P \rightarrow a \mid b \mid c \mid E \mid (E)$
  - See what happens if you change the first set of productions to  $E \rightarrow E+T \mid E-T \mid T \mid P$

CM/SC 330

47

## Regular expressions and CFGs

	Description	Machine
regular languages	regular expressions	DFAs, NFAs
context-free languages	context-free grammars	pushdown automata (PDAs)

- Programming languages are neither regular nor context-free
  - Usually almost context-free, with some hacks

CM/SC 330

48

## Pushdown Automaton (PDA)

- A **pushdown automaton** (PDA) is an abstract machine similar to the DFA
  - Has a finite set of states
  - Also has a *pushdown stack*
- Moves of the PDA are as follows:
  - An input symbol is read and the top symbol on the stack is read
  - Based on both inputs, the machine
    - Enters a new state, and
    - Writes zero or more symbols onto the pushdown stack
  - String accepted if the stack is empty at end of string

CM/SC 330

49

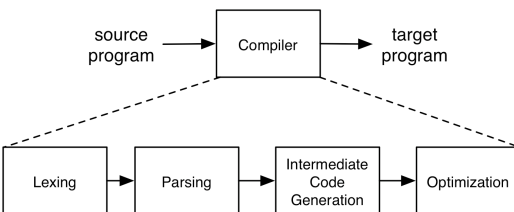
## Power of PDAs

- PDAs are more powerful than DFAs
  - $a^n b^n$ , which cannot be recognized by a DFA, can easily be recognized by the PDA
    - Stack all  $a$  symbols and, for each  $b$ , pop an  $a$  off the stack.
    - If the end of input is reached at the same time that the stack becomes empty, the string is accepted
- As with NFA, we can also have a NDPDA
  - NDPDA are more powerful than DPDA
  - NDPDA can recognize even length palindromes over  $\{0,1\}^*$ , but a DPDA cannot. Why? (Hint: Consider palindromes over  $\{0,1\}^2\{0,1\}^*$ )
- It is true, but less clear, that the languages accepted by NDPDAs are equivalent to the context-free languages

CM/SC 330

50

## Steps of Compilation



CM/SC 330

51

## Parsing

- There are many efficient techniques for turning strings into parse trees or ASTs
  - They all have strange names, like LL(k), SLR(k), LR(k)...
  - Take CMSC 430 for more details
- We will look at one very simple technique: *recursive descent parsing*
  - This is a “top-down” parsing algorithm because we’re going to begin at the start symbol and try to produce the string

CM/SC 330

52

## Recursive Descent Parsing

- Goal
  - Determine if we can produce the string to be parsed from the grammar’s start symbol
- Approach
  - Recursively replace nonterminal with RHS of production
- At each step, we’ll keep track of two facts
  - What tree node are we trying to match?
  - What is the **lookahead** (next token of the input string)?
    - Helps guide selection of production used to replace nonterminal

CM/SC 330

53

## Recursive Descent Parsing (cont.)

- At each step, 3 possible cases
  - If we’re trying to match a terminal
    - If the lookahead is that token, then succeed, advance the lookahead, and continue
  - If we’re trying to match a nonterminal
    - Pick which production to apply based on the lookahead
  - Otherwise fail with a parsing error

CM/SC 330

54

## Example

$E \rightarrow id = n \mid \{ L \}$

$L \rightarrow E ; L \mid \epsilon$

– Here  $n$  is an integer and  $id$  is an identifier

- One input might be
  - $\{ x = 3 ; \{ y = 4 ; \} ; \}$
  - This would get turned into a list of tokens  
 $\{ x = 3 ; \{ y = 4 ; \} ; \}$
  - And we want to turn it into a parse tree

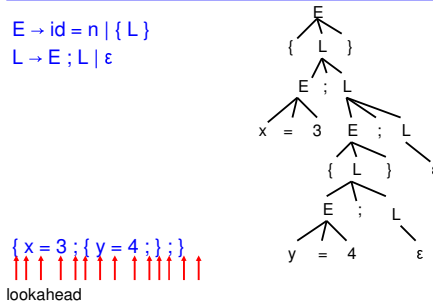
CMSC 330

55

## Example (cont'd)

$E \rightarrow id = n \mid \{ L \}$

$L \rightarrow E ; L \mid \epsilon$



CMSC 330

56

## Recursive Descent Parsing (cont.)

- Key step
  - Choosing which production should be selected
- Two approaches
  - Backtracking
    - Choose some production
    - If fails, try different production
    - Parse fails if all choices fail
  - Predictive parsing
    - Analyze grammar to find FIRST sets for productions
    - Compare with lookahead to decide which production to select
    - Parse fails if lookahead does not match FIRST

CMSC 330

57

## First Sets

- Motivating example
  - The lookahead is  $x$
  - Given grammar  $S \rightarrow xyz \mid abc$ 
    - Select  $S \rightarrow xyz$  since 1st terminal in RHS matches  $x$
  - Given grammar  $S \rightarrow A \mid B$   $A \rightarrow x \mid y$   $B \rightarrow z$ 
    - Select  $S \rightarrow A$ , since  $A$  can derive string beginning with  $x$
- In general
  - Choose a production that can derive a sentential form beginning with the lookahead
  - Need to know what terminal may be **first** in any sentential form derived from a nonterminal / production

CMSC 330

58

## First Sets

- Definition
  - **First( $\gamma$ )**, for any terminal or nonterminal  $\gamma$ , is the set of initial terminals of all strings that  $\gamma$  may expand to
  - We'll use this to decide what production to apply
- Examples
  - Given grammar  $S \rightarrow xyz \mid abc$ 
    - $\text{First}(xyz) = \{ x \}$ ,  $\text{First}(abc) = \{ a \}$
    - $\text{First}(S) = \text{First}(xyz) \cup \text{First}(abc) = \{ x, a \}$
  - Given grammar  $S \rightarrow A \mid B$   $A \rightarrow x \mid y$   $B \rightarrow z$ 
    - $\text{First}(x) = \{ x \}$ ,  $\text{First}(y) = \{ y \}$ ,  $\text{First}(A) = \{ x, y \}$
    - $\text{First}(z) = \{ z \}$ ,  $\text{First}(B) = \{ z \}$
    - $\text{First}(S) = \{ x, y, z \}$

CMSC 330

59

## Calculating First( $\gamma$ )

- For terminal  $a$ ,  $\text{First}(a) = \{ a \}$
- For a nonterminal  $N$ :
  - If  $N \rightarrow \epsilon$ , then add  $\epsilon$  to  $\text{First}(N)$
  - If  $N \rightarrow \alpha_1 \alpha_2 \dots \alpha_n$ , then (note the  $\alpha_i$  are all the symbols on the right side of one single production):
    - Add  $\text{First}(\alpha_1 \alpha_2 \dots \alpha_n)$  to  $\text{First}(N)$ , where  $\text{First}(\alpha_1 \alpha_2 \dots \alpha_n)$  is defined as
      - $\text{First}(\alpha_1)$  if  $\epsilon \notin \text{First}(\alpha_1)$
      - Otherwise  $(\text{First}(\alpha_1) - \epsilon) \cup \text{First}(\alpha_2 \dots \alpha_n)$
    - If  $\epsilon \in \text{First}(\alpha_i)$  for all  $i$ ,  $1 \leq i \leq n$ , then add  $\epsilon$  to  $\text{First}(N)$

CMSC 330

60

## Examples

$E \rightarrow id = n \mid \{ L \}$	$E \rightarrow id = n \mid \{ L \} \mid \epsilon$
$L \rightarrow E ; L \mid \epsilon$	$L \rightarrow E ; L \mid \epsilon$
$First(id) = \{ id \}$	$First(id) = \{ id \}$
$First("=") = \{ "=" \}$	$First("=") = \{ "=" \}$
$First(n) = \{ n \}$	$First(n) = \{ n \}$
$First("{") = \{ "{" \}$	$First("{") = \{ "{" \}$
$First("}") = \{ "}" \}$	$First("}") = \{ "}" \}$
$First(";") = \{ ";" \}$	$First(";") = \{ ";" \}$
$First(E) = \{ id, "{" \}$	$First(E) = \{ id, "{", \epsilon \}$
$First(L) = \{ id, "{", \epsilon \}$	$First(L) = \{ id, "{", ";", \epsilon \}$

CM/SC 330

61

## Recursive Descent Parser Implementation

- For terminals, create function `match(a)`
  - If lookahead is `a` it consumes the lookahead by advancing the lookahead to the next token, and returns
  - Otherwise fails with a parse error if lookahead is not `a`
  - In algorithm descriptions, consider `parse_a`, `parse_term(a)` to be aliases for `match(a)`
- For each nonterminal `N`, create a function `parse_N`
  - Called when we're trying to parse a part of the input which corresponds to (or can be derived from) `N`
  - `parse_S` for the start symbol `S` begins the parse

CM/SC 330

62

## Parser Implementation (cont.)

- The body of `parse_N` for a nonterminal `N` does the following
  - Let  $N \rightarrow \beta_1 \mid \dots \mid \beta_k$  be the productions of `N`
    - Here  $\beta_i$  is the entire right side of a production- a sequence of terminals and nonterminals
  - Pick the production  $N \rightarrow \beta_i$  such that the lookahead is in `First( $\beta_i$ )`
    - It must be that  $First(\beta_i) \cap First(\beta_j) = \emptyset$  for  $i \neq j$
    - If there is no such production, but  $N \rightarrow \epsilon$  then return
    - Otherwise fail with a parse error
  - Suppose  $\beta_i = \alpha_1 \alpha_2 \dots \alpha_n$ . Then call `parse_α1()`; ... ; `parse_αn()` to match the expected right-hand side, and return

CM/SC 330

63

## Parser Implementation (cont.)

- Parse is built on procedure calls
- Procedures may be (mutually) recursive

CM/SC 330

64

## Recursive Descent Parser

- Given grammar  $S \rightarrow xyz \mid abc$ 
  - $First(xyz) = \{ x \}$ ,  $First(abc) = \{ a \}$

### Parser

```

parse_S() {
  if (lookahead == "x") {
    match("x"); match("y"); match("z"); // S → xyz
  }
  else if (lookahead == "a") {
    match("a"); match("b"); match("c"); // S → abc
  }
  else error();
}

```

CM/SC 330

65

## Recursive Descent Parser

- Given grammar  $S \rightarrow A \mid B$   $A \rightarrow x \mid y$   $B \rightarrow z$ 
  - $First(A) = \{ x, y \}$ ,  $First(B) = \{ z \}$

### Parser

```

parse_S() {
  if ((lookahead == "x" ||
      lookahead == "y"))
    parse_A(); // S → A
  else if (lookahead == "z")
    parse_B(); // S → B
  else error();
}

parse_A() {
  if (lookahead == "x")
    match("x"); // A → x
  else if (lookahead == "y")
    match("y"); // A → y
  else error();
}

parse_B() {
  if (lookahead == "z")
    match("z"); // B → z
  else error();
}

```

CM/SC 330

66

## Example

$E \rightarrow id = n \mid \{ L \}$        $\text{First}(E) = \{ id, "{" \}$   
 $L \rightarrow E ; L \mid \epsilon$

```

parse_E() {
    if (lookahead == "id") {
        match("id");
        match("="); // E → id = n
        match("n");
    }
    else if (lookahead == "{" {
        match("{");
        parse_L(); // E → { L }
        match("}");
    }
    else error();
}

parse_L() {
    if ((lookahead == "id" ||
        lookahead == "{")) {
        parse_E();
        match(";"); // L → E ; L
        parse_L();
    }
    else ; // L → ε
}
    
```

CM/SC 330

67

## Things to Notice

- If you draw the execution trace of the parser
  - You get the parse tree

### Examples

– Grammar

$S \rightarrow xyz$

$S \rightarrow abc$

– String "xyz"

parse\_S()

match("x")

match("y")

match("z")

**S**

**/!\**

**x y z**

### Grammar

$S \rightarrow A \mid B$

$A \rightarrow x \mid y$

$B \rightarrow z$

### String "x"

parse\_S()

parse\_A()

match("x")

**S**

**|**

**A**

**|**

**x**

CM/SC 330

68

## Things to Notice (cont.)

- This is a **predictive** parser
  - Because the lookahead determines exactly which production to use
- This parsing strategy may fail on some grammars
  - Possible infinite recursion
  - Production First sets overlap
  - Production First sets contain  $\epsilon$
- Does not mean grammar is not usable
  - Just means this parsing method not powerful enough
  - May be able to change grammar

CM/SC 330

69

## Left Factoring

- Consider parsing the grammar  $E \rightarrow ab \mid ac$ 
  - $\text{First}(ab) = a$
  - $\text{First}(ac) = a$
  - Parser cannot choose between RHS based on lookahead!
- Parser fails whenever  $A \rightarrow \alpha_1 \mid \alpha_2$  and
  - $\text{First}(\alpha_1) \cap \text{First}(\alpha_2) \neq \epsilon \text{ or } \emptyset$
- Solution
  - Rewrite grammar using **left factoring**

CM/SC 330

70

## Left Factoring Algorithm

- Given grammar
  - $A \rightarrow x\alpha_1 \mid x\alpha_2 \mid \dots \mid x\alpha_n \mid \beta$
- Rewrite grammar as
  - $A \rightarrow xL \mid \beta$
  - $L \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$
- Repeat as necessary
- Examples
  - $S \rightarrow ab \mid ac \quad \Leftrightarrow S \rightarrow aL \quad L \rightarrow b \mid c$
  - $S \rightarrow abcA \mid abB \mid a \quad \Leftrightarrow S \rightarrow aL \quad L \rightarrow bcA \mid bB \mid \epsilon$
  - $L \rightarrow bcA \mid bB \mid \epsilon \quad \Leftrightarrow L \rightarrow bL' \mid \epsilon \quad L' \rightarrow cA \mid B$

CM/SC 330

71

## Left Recursion

- Consider grammar  $S \rightarrow Sa \mid \epsilon$ 
  - $\text{First}(Sa) = a$ , so we're ok as far as which production
  - Try writing parser
 

```

parse_S() {
    if (lookahead == "a") {
        parse_S();
        match("a"); // S → Sa
    }
    else { }
}
                    
```
  - Body of `parse_S()` has an infinite loop
    - if (lookahead = "a") then `parse_S()`
  - Infinite loop occurs in grammar with **left recursion**

CM/SC 330

72

## Right Recursion

- Consider grammar  $S \rightarrow aS \mid \epsilon$ 
  - Again,  $\text{First}(aS) = a$
  - Try writing parser

```
parse_S() {
    if (lookahead == "a") {
        match("a");
        parse_S(); // S → aS
    }
    else {}
}
```
  - Will `parse_S()` infinite loop?
    - Invoking `match()` will advance lookahead, eventually stop
  - Top down parsers handles grammar w/ **right recursion**

CMSC 330

73

## Algorithm To Eliminate Left Recursion

- Given grammar
  - $A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \beta$ 
    - Why must  $\beta$  exist?
- Rewrite grammar as
  - $A \rightarrow \beta L$
  - $L \rightarrow \alpha_1 L \mid \alpha_2 L \mid \dots \mid \alpha_n L \mid \epsilon$
- Replaces left recursion with right recursion
- Repeat as necessary

CMSC 330

74

## Eliminating Left Recursion (cont.)

- Examples
  - $S \rightarrow Sa \mid \epsilon$        $\Leftrightarrow S \rightarrow L$      $L \rightarrow aL \mid \epsilon$
  - $S \rightarrow Sa \mid Sb \mid c$      $\Leftrightarrow S \rightarrow cL$      $L \rightarrow aL \mid bL \mid \epsilon$
- May need more powerful algorithms to eliminate **mutual recursion** leading to left recursion
  - $S \rightarrow Aa \mid b$
  - $A \rightarrow Sb$

CMSC 330

75

## Expr Grammar for Top-Down Parsing

$E \rightarrow T E'$   
 $E' \rightarrow \epsilon \mid + E$   
 $T \rightarrow P T'$   
 $T' \rightarrow \epsilon \mid * T$   
 $P \rightarrow n \mid ( E )$

- Notice we can always decide what production to choose with only one symbol of lookahead

CMSC 330

76

## Tradeoffs with Other Approaches

- Recursive descent parsers are easy to write
  - The formal definition is a little clunky, but if you follow the code then it's almost what you might have done if you weren't told about grammars formally
  - They're unable to handle certain kinds of grammars
- Recursive descent is good for a simple parser
  - Though tools can be fast if you're familiar with them
- Can implement top-down predictive parsing as a table-driven parser
  - By maintaining an explicit stack to track progress

CMSC 330

77

## Tradeoffs with Other Approaches

- More powerful techniques need tool support
  - Can take time to learn tools (lex/flex, yacc/bison)
- Main alternative is bottom-up, shift-reduce parser
  - Replaces RHS of production with LHS (nonterminal)
  - Example grammar
    - $S \rightarrow aA, A \rightarrow Bc, B \rightarrow b$
  - Example parse
    - $abc \Rightarrow aBc \Rightarrow aA \Rightarrow S$
    - Derivation happens in reverse
  - Something to look forward to in CMSC 430

CMSC 330

78

## What's Wrong With Parse Trees?

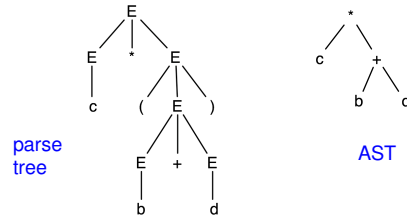
- Parse trees contain too much information
  - Example
    - Parentheses
    - Extra nonterminals for precedence
  - This extra stuff is needed for parsing
- But when we want to **reason** about languages
  - Extra information gets in the way (too much detail)

CM/SC 330

79

## Abstract Syntax Trees (ASTs)

- An **abstract syntax tree** is a more compact, abstract representation of a parse tree, with only the essential parts

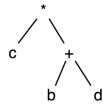


CM/SC 330

80

## Abstract Syntax Trees (cont.)

- Intuitively, ASTs correspond to the data structure you'd use to represent strings in the language
  - Note that grammars describe trees
  - So do OCaml datatypes (which we'll see later)
  - $E \rightarrow a \mid b \mid c \mid E+E \mid E-E \mid E^*E \mid (E)$



CM/SC 330

81

## Producing an AST

- To produce an AST, we can modify the `parse()` functions to construct the AST along the way
  - `match(a)` returns an AST node (leaf) for `a`
  - `Parse_A` returns an AST node for `A`
    - AST nodes for RHS of production become children of LHS node
- Example
  - $S \rightarrow aA$

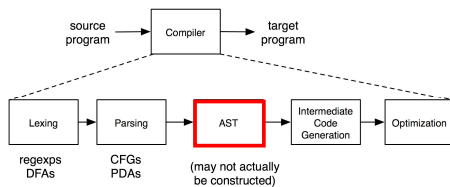
```
Node parse_S() {
  Node n1, n2;
  if (lookahead == "a") {
    n1 = match("a");
    n2 = parse_A();
    return new Node(n1, n2);
  }
}
```



CM/SC 330

82

## The Compilation Process



CM/SC 330

83

## Summary

- Learned a little about parsing
  - Recursive descent parser
  - Predictive parsing using FIRST sets
- Rewriting grammars for predicative parsing
  - Left factoring
  - Eliminating left recursion
- Abstract syntax trees (ASTs)

CM/SC 330

84