

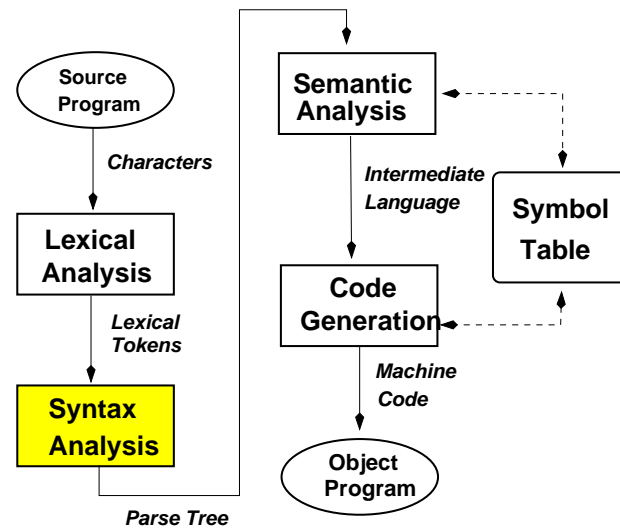
# CSC 488S/CSC 2107S Lecture Notes

These lecture notes are provided for the personal use of students taking CSC488H1S or CSC2107S in the Winter 2012/2013 term at the University of Toronto

Copying for purposes other than this use and all forms of distribution are expressly prohibited.

©David B. Wortman, 2002,2003,2004,2008,2009,2010,2012,2013

©Marsha Chechik, 2005,2006,2007



0

63

## Reading Assignment

*Fischer, Cytron, LeBlanc*

## Chapter 4

## Syntax Analysis

- The *syntax* of a language defines the sequences of language tokens that form legal sentences in the language.
- Examples:  

$X := Y * Z + W / X$	Legal
$X =: * Y Z + / W X$	Illegal
<code>for ( i = 0 ; i &lt; N ; i++ )</code>	Legal
<code>for ( i 0 ; i = &lt; N ; ++ )</code>	Illegal
- Syntax analysis is the process of analyzing a sequence of lexical tokens to determine if they are a legal sentence in some given language.
- Syntax analysis is based on a precise formal definition of the source language.

64

65

## 66

- 68

## 68

- 68

## 67



## Notation

## Notation

## Definitions for Parsing

- Let  $N$  Set of nonterminal symbols  
 $\Sigma$  Set of terminal symbols  
 $S$  The distinguished nonterminal start (goal) symbol  
 $A$  A *non-terminal* symbol in  $N$   
 $\alpha, \beta, \gamma, \omega$  Strings in  $(N \cup \Sigma)^*$
- The most general (*context sensitive*) form of a production rule is  
 $\gamma A \omega \rightarrow \gamma \alpha \omega$   
 In the left context  $\gamma$  and the right context  $\omega$ ,  $A$  derives  $\alpha$ .
- For context free grammars,  $\gamma, \omega$  are null and all rules are of the form:  
 $A \rightarrow \alpha$

70

## Context Free Grammar Example

- For the language definition:  

$$\begin{array}{lcl} S & \rightarrow & A B \\ A & \rightarrow & a A \\ & & | \quad a \\ B & \rightarrow & B b \\ & & | \quad b \end{array}$$
- Non terminals:  $N = \{ A, B, S \}$
- Terminals:  $\Sigma = \{ a, b \}$
- Goal symbol:  $S$
- Productions:  
 $P = \{ S \rightarrow A B, A \rightarrow a A, A \rightarrow a, B \rightarrow B b, B \rightarrow b \}$

72

## Programming Language Grammars

- A *context free grammar*  $G$  used to formally define programming languages is a 4-tuple  $G = (N, \Sigma, S, P)$   
 $\Sigma$  A finite vocabulary of *terminal* symbols.  
 This is the set of lexical tokens produced by lexical analysis.  
 $N$  A finite set of *nonterminal* symbols, the nonterminal vocabulary  
 $S$  The start symbol a distinguished nonterminal ( $S \in N$ )  
 $S$  is also called the *goal symbol*  
 $P$  A finite set of *productions*. The productions are also called *rewriting rules*.  
 $V$  The *vocabulary* of the grammar is  $\Sigma \cup N$
- Each of the productions in  $P$  has the form:  

$$A \rightarrow X_1 \dots X_m \quad m \geq 0$$
 where  
 $A \in N$   
 $X_i \in V, 1 \leq i \leq m, m \geq 0$   
 Note  $A \rightarrow$  is a valid production ( $m = 0$ )

71

## Some Informal Definitions

- The **productions** in a grammar are a set of rules (in some notation) that define the sequences of terminal symbols that make up legal sentences in the language defined by the grammar.
- A grammar is **unambiguous** if there is exactly **one sequence of rules** for forming **each** legal sentence in the language. A **language** is ambiguous if **all** grammars for the language are ambiguous.  
**Being unambiguous is good.**
- A grammar is **context free** (CF) if every rule can be applied completely independent of the surrounding context.  
**Being context free is good.**
- A language is **deterministic** if it is always possible to determine which rule to apply next when parsing every sentence in the language.  
**Being deterministic is good.**

73

## Formal Definitions for Parsing

- Define: **Sentential form**:  
The set of strings produced by the start symbol.  
 $\{ w \mid S \Rightarrow^* w \}$   
 $w$  contains all strings that could occur during parse.
- Define: **Language**:  
The set of terminal strings produced by the start symbol  
 $\{ w \mid S \Rightarrow^* w \} \cap \Sigma^*$   
This is all legal sentences (programs) in the language.
- Define: **Rewriting**: The replacement of a nonterminal symbol  $A$  with the right side of one of the production rules for  $A$  e.g.  $A \rightarrow \alpha \beta \gamma$

74

## Nullability – Deriving the Empty String

- Define: **nullable**  
A string  $\alpha$  in  $N^+$  is called nullable iff  
 $\alpha \Rightarrow^* \lambda$
- Nullability is a significant issue for the construction of parsers.  
When the string  $\alpha$  actually produces  $\lambda$  then it provides **no** information that the parser can use to make a parsing decision.
- A grammar can be processed to determine which non terminal symbols can potentially produce  $\lambda$ , See Fischer, Cytron, LeBlanc Figure 4.7
- Example:  

$$\begin{aligned} A &\rightarrow B C D \\ B &\rightarrow \lambda \\ C &\rightarrow B \mid c \\ D &\rightarrow B C \mid d \\ A &\Rightarrow BCD \Rightarrow CD \Rightarrow BD \Rightarrow D \Rightarrow BC \Rightarrow C \Rightarrow B \Rightarrow \lambda \end{aligned}$$

76

- Define: **Derivation**

Given  $A \rightarrow \gamma$  and sentential form  $\alpha A \beta$  then

$\alpha A \beta \Rightarrow \alpha \gamma \beta$  is one derivation step

and

$\alpha A \beta \Rightarrow^* \alpha \gamma \beta$  derives in zero or more steps

$\alpha A \beta \Rightarrow^+ \alpha \gamma \beta$  derives in one or more steps

- A derivation can be represented as a *parse tree*.  
Applying a rule  $A \rightarrow \alpha \beta \gamma$   
adds a node  $A$  with children  $\alpha$ ,  $\beta$  and  $\gamma$  to the parse tree
  - If there is more than one nonterminal that can be rewritten in the sentential form then there are two conventions for choosing the nonterminal
    - **Leftmost derivation** expand the nonterminals left to right
    - **Rightmost derivation** expand the nonterminals right to left
- See slides 93 and 125 for examples.

75

## Rewriting and Left/Right Recursion

- A production of the form  $A \rightarrow \alpha \beta \gamma$  means  
 $\alpha \beta \gamma$  can be derived from  $A$  in a top down parse  
 $\alpha \beta \gamma$  can be reduced to  $A$  in a bottom up parse.
  - Define: **recursive nonterminal**  
A nonterminal symbol  $A$  is recursive if  $A \Rightarrow^* \alpha A \beta$
  - Define: **left recursion**  
A nonterminal symbol  $A$  is left recursive if  $A \Rightarrow^* A \beta$
  - Define: **right recursive**  
A nonterminal symbol  $A$  is right recursive if  $A \Rightarrow^* \alpha A$
  - A grammar is ambiguous if any nonterminal is both left and right recursive.  
Example:  

$$expression \Rightarrow expression + expression$$
- In general the ambiguity of a context free grammar is not computable.

77

## Top Down vs. Bottom Up Parsing

- Top Down parsers start with the start symbol for the language (  $S$  ) and attempt to find a sequence of production rules that transforms the start symbol into a given sequence of input tokens.  
Conceptually it builds the parse tree from the root to the leaves.  
A top down parse is called a **derivation**.  
Recursive Descent and LL(k) are top down parsing techniques.
- Bottom Up parsers start with a sequence of input tokens and attempt to find a sequence of production rules that will reduce the entire sequence of terminal symbols to the start symbol.  
Conceptually it builds the parse tree from the leaves to the root.  
A bottom up parse is called a **reduction** or a **reverse derivation**.  
LR(k), SLR(k) and LALR(k) are bottom up parsing techniques.
- For unambiguous, deterministic, context free languages, top down and bottom up parsers should produce the same parse tree for a given sequence of input tokens.

78

- For LALR(1) and SLR(1) parsers semantic and/or code generation actions can only be invoked when a reduction is applied.  
This is equivalent to requiring that such actions occur only at the right end of a rule.
- It is usually necessary to add extra productions to a grammar that provide a hook for attaching actions. YACC/Bison do this automatically.

### Example:

Reference Grammar:

ifStatement  $\rightarrow$  if expression **then** statement **else** statement **end**

Compiler Grammar

ifStatement  $\rightarrow$  ifHead statement elsePart **end** fixFwdBranch

ifHead  $\rightarrow$  if expression **then** checkBoolean emitBranchFalse

elsePart  $\rightarrow$

$\rightarrow$  elseHead statement

elseHead  $\rightarrow$  **else** emitFwdBranch fixFwdBranch

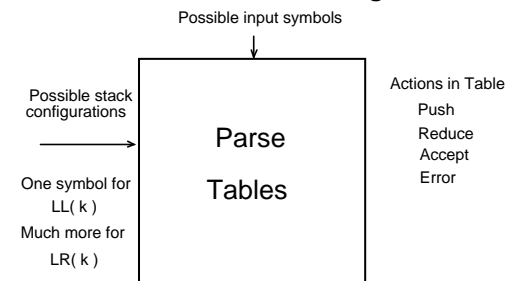
80

## Grammar Design for Compiling

- A typical programming language will have two grammars.
  - A *reference grammar* that is made available to *users* of the language.
  - A *compiler grammar* that is used by the *implementation* of a language.
  - These two grammars **should** describe exactly the same language.
- The design of the compiler grammar
  - Affects the effort required to implement semantic analysis and code generation.
  - Determines the order in which constructs in the language will be processed (for a given parsing algorithm).
  - Determines the class of parsers that must be used to perform syntax analysis.
- Compiler implementors should use the degrees of freedom that they have in designing the compiler grammar wisely to make it easier to implement semantic analysis and code generation.

79

## Table-Driven Parsing



- Conceptually a table for each parser state. Allow multiple actions per state so tables can be merged to one table for all parser states.
- Assume input stream ends with  $k$  end markers  $\$^k$ .  
Final state is accepting or rejecting.
- Top Down - stack represents what we expect to see.  
Bottom Up - stack represents what we've seen so far.

81

## Practical Parsing Techniques

- **Recursive Descent** is a top down parsing technique with the potential for backtracking. Recursive descent parsers are usually written as a collection of interacting recursive functions. Recursive descent is useful for dealing with complicated or poorly designed languages (e.g. C (gcc), C++, Java, Fortran).
- **LL(1)** is a table-driven top down parsing technique.  
An LL(1) parser is a *deterministic push down automaton (DPDA)*
- **LR(1)** is a table-driven bottom up parsing technique.  
The parser uses a Shift/Reduce algorithm to record the state of its reduction.  
An LR(1) parser is a *deterministic push down automaton (DPDA)*  
**SLR(1)** and **LALR(1)** are practical versions of LR(1)

There are many other parsing techniques, but the ones listed here are used in most compilers.

82

## Follow Sets

- Follow sets are a data structure that are often used in making parsing decisions. For example if a nonterminal symbol  $A$  is *nullable* then the parser may look for *terminal symbols* in  $Follow(A)$  to make a parsing decision.
- Define:  $follow(A)$   

$$Follow(A) = \{ b \in \Sigma \mid S \Rightarrow^+ \alpha A b \beta \}$$

For a nonterminal symbol  $A$ , the set  $follow(A)$  is the set of **terminal symbols** that can **immediately** follow  $A$  in **any sentential form** derived from the goal symbol  $S$
- For *small* grammars it is possible to compute Follow sets by inspection, **but** this is much harder than computing First sets, since all possible instances of nullability must be accounted for.  
The calculation of Follow sets can be automated, See *Fischer, Cytron, LeBlanc* Figure 4.11
- More generally  $Follow_k(\alpha)$  uses k-symbol lookahead.

84

## First Sets

- First sets are a data structure used as a decision making tool in many parsers.  
For example if a top down parser is trying to find a *non-terminal* symbol  $B$  then it will look for *terminal symbols* in  $first(B)$
- Define:  $First(\alpha)$   

For a string  $\alpha$  ( $\alpha \in (N \cup \Sigma)^*$ )

$$First(\alpha) = \{ b \in \Sigma \mid \alpha \Rightarrow^+ b \beta \}$$

The set  $First(\alpha)$  is the set of **terminal symbols** that can occur at the beginning of strings derived from  $\alpha$ .
- First Sets can also be defined recursively:
 

$First(c\beta) = \{c\}$	$c$ is a terminal symbol
$First(\lambda) = \{ \}$	
$First(B\gamma) = First(B)$	$B$ is not nullable
$= First(B) \cup First(\gamma)$	$B$ nullable
- For small grammars First sets can often be determined by inspection as *long as nullability is taken into account*. For larger grammars see the algorithm in *Fischer, Cytron, LeBlanc* Figure 4.8 .
- More generally  $First_k(\alpha)$  uses k-symbol lookahead.

83

## Informal definition of Follow Sets

Follow Sets can also be defined intuitively:

- |                                       |  |
|---------------------------------------|--|
| If $S$ is the goal symbol             | Add $\{ \$ \}$ to $Follow(S)$  |
| If $S \Rightarrow^+ \alpha A$         | Add $\{ \$ \}$ to $Follow(A)$  |
| If $S \Rightarrow^+ \alpha A c \beta$ | Add $\{ c \}$ to $Follow(A)$   |
| If $S \Rightarrow^+ \alpha A B \beta$ | Add $First(B)$ to $Follow(A)$ <span style="float: right;"><math>B</math> not nullable</span>               |
| If $S \Rightarrow^+ \alpha A B \beta$ | Add $First(B) \cup First(\beta)$ to $Follow(A)$ <span style="float: right;"><math>B</math> nullable</span> |
| If $A \rightarrow \alpha B$           | Add $Follow(A)$ to $Follow(B)$   |

85

## First and Follow Sets Example

Grammar:

$A \rightarrow B C c$   
 $\rightarrow e D B$   
 $B \rightarrow \lambda$   
 $\rightarrow b C D E$   
 $C \rightarrow D a B$   
 $\rightarrow c a$   
 $D \rightarrow \lambda$   
 $\rightarrow d D$   
 $E \rightarrow e A f$   
 $\rightarrow c$

B and D are *nullable*

86

## First Sets (by inspection)

$A \rightarrow B C c$   
 $A \rightarrow e D B$   
 $First(A) = First(B) \cup First(Cc) \cup \{e\} = \{a, b, c, d, e\}$   
 $B \rightarrow \lambda$   
 $B \rightarrow b C D E$   
 $First(B) = First(\lambda) \cup \{b\} = \{b\}$   
 $C \rightarrow D a B$   
 $C \rightarrow c a$   
 $First(C) = First(D) \cup First(aB) \cup \{c\} = \{a, c, d\}$   
 $D \rightarrow \lambda$   
 $D \rightarrow d D$   
 $First(D) = First(\lambda) \cup \{d\} = \{d\}$   
 $E \rightarrow e A f$   
 $E \rightarrow c$   
 $First(E) = \{e\} \cup \{c\} = \{c, e\}$

87

## Follow Sets (by inspection)

$A \rightarrow B C c$   
 $E \rightarrow e A f$   
 $Follow(A) = \{ \$ \} \cup \{ f \} = \{ f, \$ \}$   
 $A \rightarrow B C c$   
 $A \rightarrow e D B$   
 $C \rightarrow D a B$   
 $Follow(B) = First(C) \cup Follow(A) \cup Follow(C) = \{ a, c, d, e, f, \$ \}$   
 $A \rightarrow B C c$   
 $B \rightarrow b C D E$   
 $Follow(C) = \{ c \} \cup First(D) \cup First(E) = \{ c, d, e \}$   
 $A \rightarrow e D B$   
 $B \rightarrow b C D E$   
 $C \rightarrow D a B$   
 $D \rightarrow d D$   
 $Follow(D) = First(B) \cup Follow(A) \cup First(E) \cup \{ a \} = \{ a, b, c, e, f, \$ \}$   
 $B \rightarrow b C D E$   
 $Follow(E) = Follow(B) = \{ a, c, d, e, f, \$ \}$

88