

CSCI-C311 Programming Languages

Context-Free Grammars and Parsing

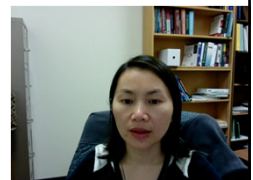
Dr. Hang Dinh



1

Outline and Reading

- After this lecture, you will learn
 - Context-free grammar
 - Parse tree
 - Classes of parsing algorithms
 - Difference between top-down parsing and bottom-up parsing
- Reading
 - Scott 4e - Section 2.1.2, 2.1.3
 - Scott 4e –Section 2.3 (but not 2.3.1 and 2.3.2)



2

Definition of Context-Free Grammars (CFG)

- Recall: CFG is used in syntax analysis, the 2nd phase of compilation
- Formally, a CFG consists of
 - A set *T* of *terminals*
 - A set *N* of *non-terminals* (or *variables*)
 - A *start symbol* *S* (a non-terminal)
 - A set of *productions* (or *rules*), each production is of the form:
$$x \rightarrow w$$
where *x* is a non-terminal and *w* is a string of terminals and non-terminals.
- The first rule of a CGF must have the start symbol on the left-hand side



3

CFG Notation

- Several rules of the same left-hand variables can be abbreviated into a single rule using the | meta-symbol :
 - Original rules: $X \rightarrow w_1$
 $X \rightarrow w_2$
...
 $X \rightarrow w_k$
 - Abbreviated rule: $X \rightarrow w_1 \mid w_2 \mid \dots \mid w_k$

The order of these rules doesn't matter.

$op \rightarrow + \mid * \mid := \mid -$

is shorthand for

$op \rightarrow +$
 $op \rightarrow -$
 $op \rightarrow *$
 $op \rightarrow :=$

or

$op \rightarrow +$
 $\rightarrow -$
 $\rightarrow *$
 $\rightarrow :=$



4

Example of CFG

- The expression grammar

- Terminals: `id number + - * / ()`
- Non-terminals: *expr op*
- Start symbol: *expr*
- The set of rules:

$$\begin{aligned} \textit{expr} &\longrightarrow \textit{id} \mid \textit{number} \mid - \textit{expr} \mid (\textit{expr}) \\ &\quad \mid \textit{expr} \textit{op} \textit{expr} \\ \textit{op} &\longrightarrow + \mid - \mid * \mid / \end{aligned}$$

- Convention: In this lecture, all non-terminals are written *italic*.



5

Backus-Naur-Form (BNF)

- The notation for CFG is sometimes called Backus-Naur Form (BNF)
 - Kleene star and meta-level parentheses of regular languages are not allowed
- Extended BNF allows
 - Kleene star, and meta-level parentheses
 - Kleene plus (+), which indicates one or more instances of symbol or group of symbols in front of it.

- Example:

- *idlist* \Rightarrow *id (, id)** is shorthand for

$$\begin{aligned} \textit{idlist} &\rightarrow \textit{id} \\ \textit{idlist} &\rightarrow \textit{idlist} , \textit{id} \end{aligned}$$

- *idlist* \Rightarrow *id (, id)+* is shorthand for

$$\begin{aligned} \textit{idlist} &\rightarrow \textit{id} , \textit{id} \\ \textit{idlist} &\rightarrow \textit{idlist} , \textit{id} \end{aligned}$$



6

Avoid Confusion in Extended BNF

- The () in the following grammar are characters that are part of the language defined by the grammar

$$\begin{aligned} \text{expr} &\longrightarrow \text{id} \mid \text{number} \mid - \text{expr} \mid (\text{expr}) \\ &\quad \mid \text{expr op expr} \\ \text{op} &\longrightarrow + \mid - \mid * \mid / \end{aligned}$$

- To avoid confusing with the meta-symbols () in extended BNF, some authors use quote marks around non-meta-symbol characters
 - $\text{idlist} \rightarrow \text{id} (' , ' \text{id})^*$
 - $\text{expr} \rightarrow ' (' \text{expr} ')'$



7

Derivations in CFG

- A CFG shows how to generate a syntactically valid string of terminals
 - Begin with start symbol S
 - Choose a rule $S \rightarrow w_1$
 - Replace S with w_1
 - Choose a nonterminal A in resulting string
 - Choose a rule $A \rightarrow w_2$
 - Replace A with w_2 , so on....

$$\begin{aligned} \text{expr} &\longrightarrow \text{id} \mid \text{number} \mid - \text{expr} \mid (\text{expr}) \\ &\quad \mid \text{expr op expr} \\ \text{op} &\longrightarrow + \mid - \mid * \mid / \end{aligned}$$

$$\begin{aligned} \text{expr} &\Rightarrow \text{expr op } \underline{\text{expr}} \\ &\Rightarrow \text{expr } \underline{\text{op}} \text{ id} \\ &\Rightarrow \underline{\text{expr}} + \text{id} \\ &\Rightarrow \text{expr op } \underline{\text{expr}} + \text{id} \\ &\Rightarrow \text{expr } \underline{\text{op}} \text{ id} + \text{id} \\ &\Rightarrow \underline{\text{expr}} * \text{id} + \text{id} \\ &\Rightarrow \text{id} * \text{id} + \text{id} \\ &\quad (\text{slope}) \quad (\text{x}) \quad (\text{intercept}) \end{aligned}$$

- Grammar on the left generates the string "**slope * x + intercept**"



8

Derivations in CFG

- The \Rightarrow metasymbol means “derives”
 - Indicates a replacement operation
- *Derivation*
 - is a series of replacement operations that shows how to derive a string of terminals from the start symbol
 - A string along the way is a *sentential form*.
 - Final sentential form is called the *yield*
- The \Rightarrow^* metasymbol means
 - “derives after zero or more replacements”

$$\begin{aligned} \text{expr} &\Rightarrow \text{expr op expr} \\ &\Rightarrow \text{expr op id} \\ &\Rightarrow \text{expr} + \text{id} \\ &\Rightarrow \text{expr op expr} + \text{id} \\ &\Rightarrow \text{expr op id} + \text{id} \\ &\Rightarrow \text{expr} * \text{id} + \text{id} \\ &\Rightarrow \text{id} * \text{id} + \text{id} \\ &\quad (\text{slope}) \quad (\text{x}) \quad (\text{intercept}) \end{aligned}$$

$$\text{expr} \Rightarrow^* \text{slope} * \text{x} + \text{intercept}$$



9

Leftmost vs. Rightmost Derivations

- *Right-most* derivation
 - At each step, choose to replace the *right-most* nonterminal with the right-hand side of some production
- *Left-most* derivation
 - At each step, choose to replace the *left-most* nonterminal with the right-hand side of some production

$$\begin{aligned} \text{expr} &\Rightarrow \text{expr op expr} \\ &\Rightarrow \text{expr op id} \\ &\Rightarrow \text{expr} + \text{id} \\ &\Rightarrow \text{expr op expr} + \text{id} \\ &\Rightarrow \text{expr op id} + \text{id} \\ &\Rightarrow \text{expr} * \text{id} + \text{id} \\ &\Rightarrow \text{id} * \text{id} + \text{id} \\ &\quad (\text{slope}) \quad (\text{x}) \quad (\text{intercept}) \end{aligned}$$

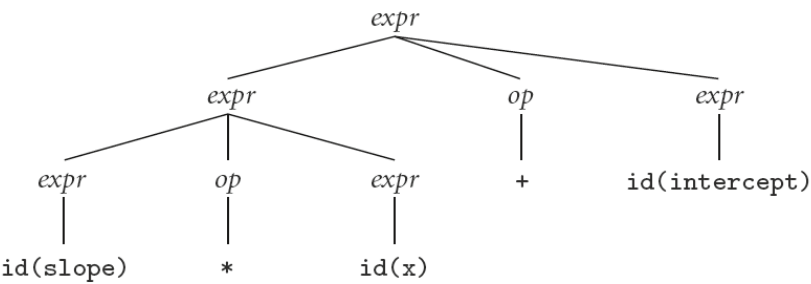


10

Parse Trees

- Represent a derivation graphically as a parse tree
 - Parse tree for "**slope * x + intercept**"

$expr \Rightarrow expr\ op\ expr$
 $\Rightarrow expr\ op\ id$
 $\Rightarrow expr\ +\ id$
 $\Rightarrow expr\ op\ expr\ +\ id$
 $\Rightarrow expr\ op\ id\ +\ id$
 $\Rightarrow expr\ * id\ +\ id$
 $\Rightarrow id\ * id\ +\ id$
 $(slope)\ (x)\ (intercept)$



- Many derivations can be represented by the same parse tree.

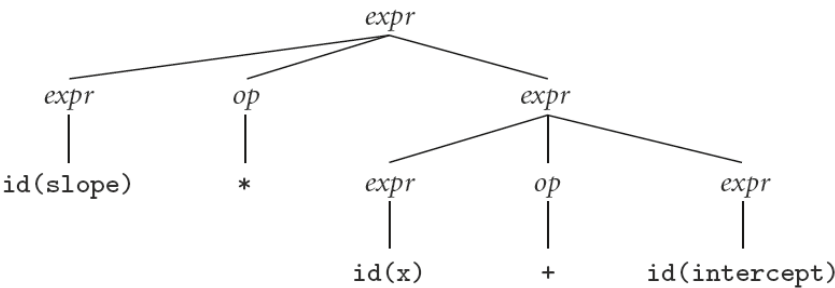


11

Ambiguous Grammars

$expr \rightarrow id \mid number \mid -\ expr \mid (\ expr)$
 $\mid expr\ op\ expr$
 $op \rightarrow + \mid - \mid * \mid /$

- Alternative parse tree for "**slope * x + intercept**"



- A grammar that generates different parse trees for the same string of terminals is said to be *ambiguous*.



12

Example of Unambiguous Grammars

- Unambiguous version of expression grammar, which captures
 - *Precedence*: multiplication, division before addition, subtraction
 - *Associativity*: operators group left to right,
so that **3+4+5** means **(3+4) +5** rather than **3+ (4+5)**

- 1. $expr \rightarrow term \mid expr\ add_op\ term$
- 2. $term \rightarrow factor \mid term\ mult_op\ factor$
- 3. $factor \rightarrow id \mid number \mid -\ factor \mid (\ expr\)$
- 4. $add_op \rightarrow + \mid -$
- 5. $mult_op \rightarrow * \mid /$

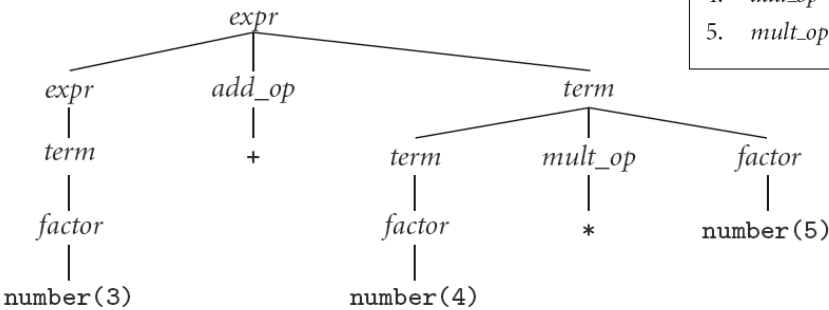


13

Example of Unambiguous Grammars

- Parse tree in unambiguous expression grammar for **3 + 4 * 5**

- 1. $expr \rightarrow term \mid expr\ add_op\ term$
- 2. $term \rightarrow factor \mid term\ mult_op\ factor$
- 3. $factor \rightarrow id \mid number \mid -\ factor \mid (\ expr\)$
- 4. $add_op \rightarrow + \mid -$
- 5. $mult_op \rightarrow * \mid /$



14

CFG and Parsing

- CFG is a *generator* for a *context-free language* (CFL)
 - *The language of a CFG* is the set of strings of terminals generated by the CFG.
 - A CFL is simply the language of some context-free grammar.
- There is an infinite number of CFGs for every CFL
 - not all grammars are created equal, however
- A parser is a context-free language *recognizer*
- For any CFG, we can create a parser that runs in $O(n^3)$ time
 - *Early's algorithm*
 - *Cooke-Younger-Kasami (CYK) algorithm*



15

Parsing Algorithms

- $O(n^3)$ time is **too slow** for parsing sizable programs
- There are large classes of grammars for which we can build parsers that run in linear time $O(n)$
 - The two most important classes are called: **LL** and **LR**

Class	Direction of Scanning	Derivation discovered	Parse tree construction	Algorithm used
LL	Left-to-right	Left-most	Top-down	predictive
LR	Left-to-right	Right-most	Bottom-up	shift-reduce

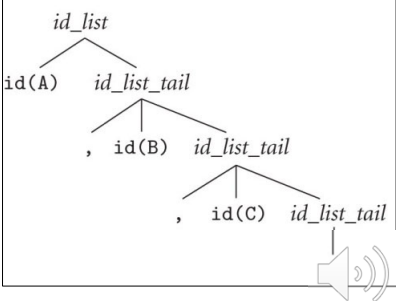
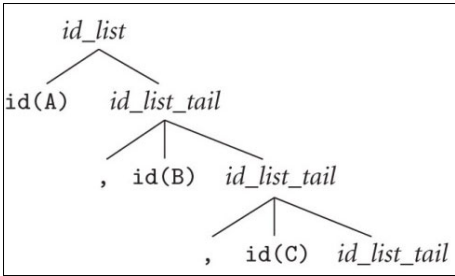
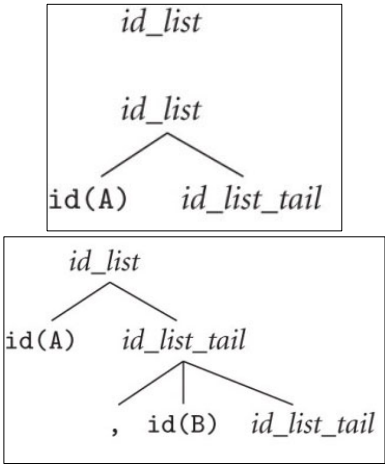
- **LR** parsers are more common than **LL** parsers



16

Example of Top-down Parsing

- Top-down parsing for input string **A, B, C;**



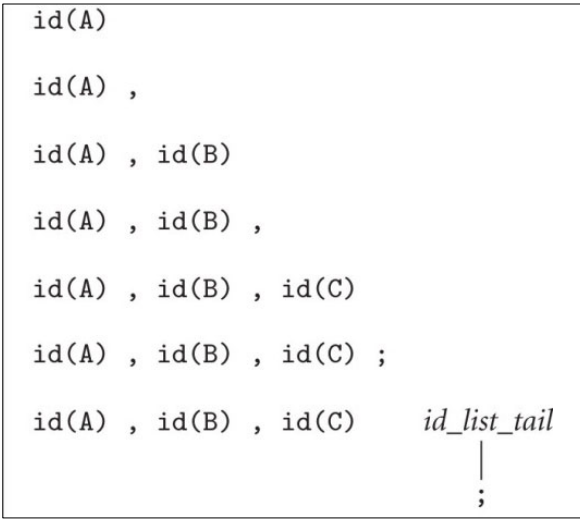
- Given the CFG for comma-separated list of identifiers ended by a ;

$id_list \rightarrow id\ id_list_tail$
 $id_list_tail \rightarrow ,\ id\ id_list_tail$
 $id_list_tail \rightarrow ;$

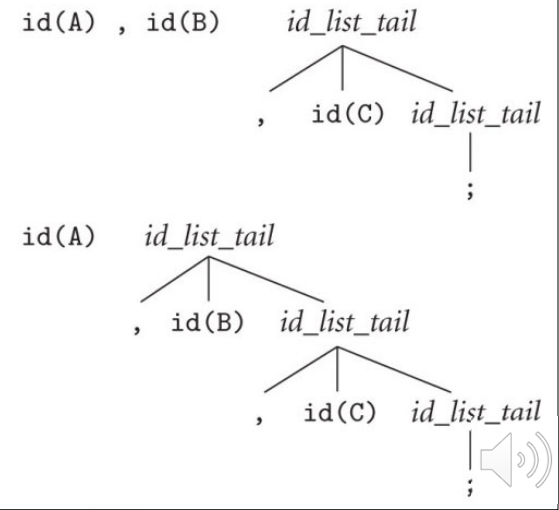
17

Example of Bottom-Up Parsing

- Top-down parsing for input string **A, B, C;**



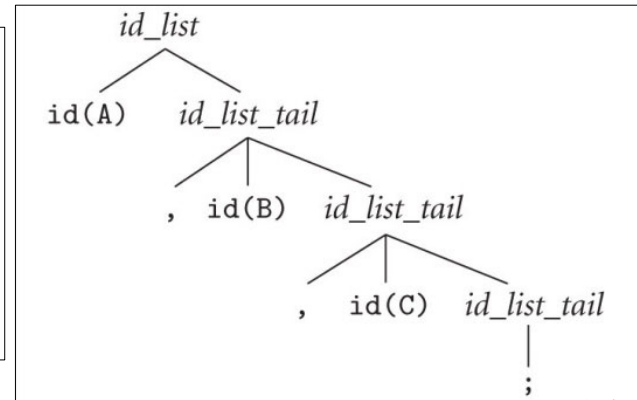
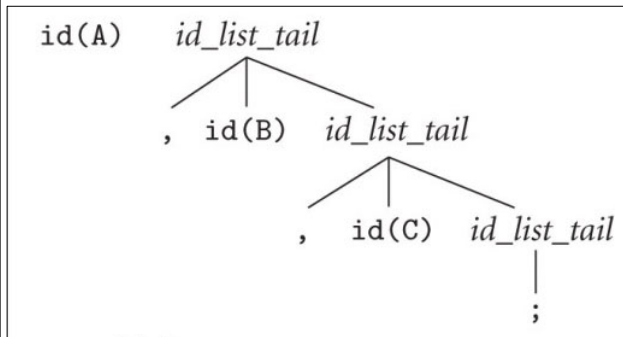
$id_list \rightarrow id\ id_list_tail$
 $id_list_tail \rightarrow ,\ id\ id_list_tail$
 $id_list_tail \rightarrow ;$



18

Example of Bottom-Up Parsing

- Top-down parsing for input string **A, B, C;**

$$id_list \longrightarrow \text{id } id_list_tail$$
$$id_list_tail \longrightarrow , id \ id_list_tail$$
$$id_list_tail \longrightarrow ;$$


19

Subclasses of Linear Time Parsers

- There are several important sub-classes of **LR** parsers
 - **SLR** (Simple LR)
 - **LALR** (Look-ahead LR)
 - “full **LR**”: Important for its generality
- Parser classes written with a number in () after main class name
 - Examples: **LL(2)** or **LALR(1)**
 - This number indicates how many tokens of look-ahead are required in order to parse.
 - Almost all real compilers use one token of look-ahead
 - We will cover **LL(1)** and **SLR(1)** in the last two lectures.



20