

An Example Grammar

$\langle \text{program} \rangle \rightarrow \langle \text{stmts} \rangle$
 $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle$
 $\quad \quad \quad | \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle$
 $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
 $\langle \text{var} \rangle \rightarrow a \mid b \mid c \mid d$
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle + \langle \text{term} \rangle$
 $\quad \quad \quad | \langle \text{term} \rangle - \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle$
 $\quad \quad \quad | \text{const}$

1

An Exemplar Derivation

$\langle \text{program} \rangle \Rightarrow \langle \text{stmts} \rangle$
 $\quad \quad \quad \Rightarrow \langle \text{stmt} \rangle$
 $\quad \quad \quad \Rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
 $\quad \quad \quad \Rightarrow a = \langle \text{expr} \rangle$
 $\quad \quad \quad \Rightarrow a = \langle \text{term} \rangle + \langle \text{term} \rangle$
 $\quad \quad \quad \Rightarrow a = \langle \text{var} \rangle + \langle \text{term} \rangle$
 $\quad \quad \quad \Rightarrow a = b + \langle \text{term} \rangle$
 $\quad \quad \quad \Rightarrow a = b + \text{const} \leftarrow \text{sentence}$

2

Sentential Forms

- Every string of symbols in the derivation is a **sentential form**
- A **sentence** is a sentential form that has only terminal symbols
- A **leftmost derivation** is one in which the leftmost non-terminal in each sentential form is the one that is expanded next in the derivation

3

Sentential Forms

- A **left-sentential form** is a sentential form that occurs in the leftmost derivation
- A **rightmost derivation** works right to left instead
- A **right-sentential form** is a sentential form that occurs in the rightmost derivation
- Some derivations are neither leftmost nor rightmost

4

Why BNF?

- Provides a clear and concise syntax description
- The parse tree can be generated from BNF
- Parsers can be based on BNF and are easy to maintain

5

Context-Free Grammars

- The syntax of simple arithmetic expression

```
expr -> id | number | -expr | (expr)
      | expr op expr
op   -> + | - | * | /
```

- What are the terminal symbols and nonterminal symbols?
- What is the start symbol?

6

One Possible Derivation

```

expr => expr op expr
     => ...
     => id + number
  
```

7

Another Example

<pre> <program> -> <stmts> <stmts> -> <stmt> <stmt> ; <stmts> <stmt> -> <var> = <expr> <var> -> a b c d <expr> -> <term> + <term> <term> - <term> <term> -> <var> const </pre>	<ul style="list-style-type: none"> • $G = \{T, N, S, P\}$ • What are the terminals? • What are the nonterminals? • What is the start symbol? • Possible strings?
--	--

8

Parse Tree

- A **parse tree** is
 - a hierarchical representation of a derivation
 - to represent the structure of the derivation of a terminal string from some non-terminal
 - to describe the hierarchical syntactic structure of programs for any language

9

An Example

- Given the simple assignment statement syntax
 - $\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$
 - $\langle \text{id} \rangle \rightarrow A \mid B \mid C$
 - $\langle \text{expr} \rangle \rightarrow \langle \text{id} \rangle + \langle \text{expr} \rangle$
 - $\mid \langle \text{id} \rangle * \langle \text{expr} \rangle$
 - $\mid (\langle \text{expr} \rangle)$
 - $\mid \langle \text{id} \rangle$
- With leftmost derivation, how is $A = B * (A + C)$ generated?

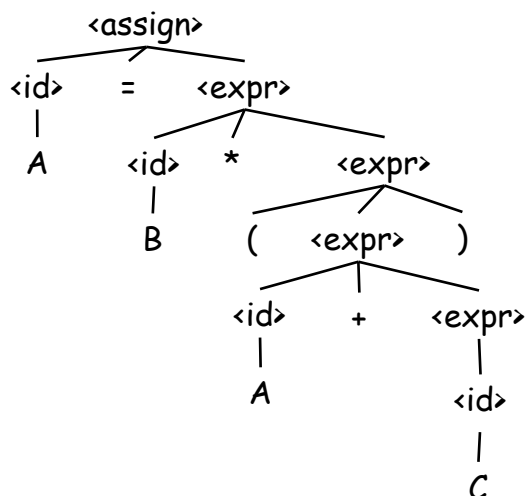
10

Derivation for $A = B * (A + C)$

$\langle \text{assign} \rangle \Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{id} \rangle * \langle \text{expr} \rangle$
 $\Rightarrow A = B * \langle \text{expr} \rangle$
 $\Rightarrow A = B * (\langle \text{expr} \rangle)$
 $\Rightarrow A = B * (\langle \text{id} \rangle + \langle \text{expr} \rangle)$
 $\Rightarrow A = B * (A + \langle \text{expr} \rangle)$
 $\Rightarrow A = B * (A + \langle \text{id} \rangle)$
 $\Rightarrow A = B * (A + C)$

11

The Parse Tree for $A = B * (A + C)$



12

Parse Tree

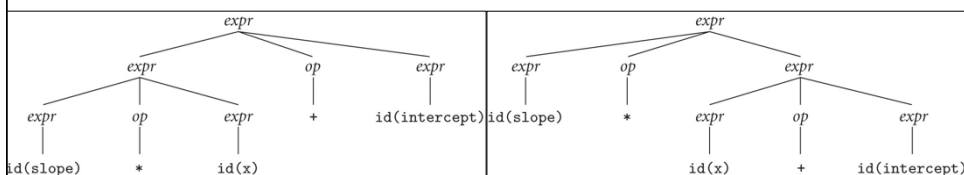
- A grammar is **ambiguous** if it generates a sentential form that has two or more distinct parse trees

13

An Ambiguous Grammar

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

- Parse trees for "slope * x + intercept":



14

What goes wrong?

- The production rules do not capture the **associativity** and **precedence** of various operators
 - **Associativity** tells whether the operators group left to right or right to left
 - Is $10 - 4 - 3$ equal to $(10 - 4) - 3$ or $10 - (4 - 3)$?
 - **Precedence** tells some operators group more tightly than the others?
 - Is $\text{slope} * x + \text{intercept}$ equal to $(\text{slope} * x) + \text{intercept}$ or $\text{slope} * (x + \text{intercept})$?

15

Operator Associativity

- Single recursion in production rules
 - $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid \text{const}$
 - X Ambiguous**
 - $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle - \text{const} \mid \text{const}$
 - ✓ Unambiguous**
 - $\langle \text{expr} \rangle \rightarrow \text{const} - \langle \text{expr} \rangle \mid \text{const}$
 - ✓ Unambiguous (less desirable)**

16

Operator Precedence

- Use stratification in production rules
 - Intentionally put operators at different levels of parse trees

```

<expr> -> <expr> - <term> | <term>
<term> -> <term> / const | const

```

17

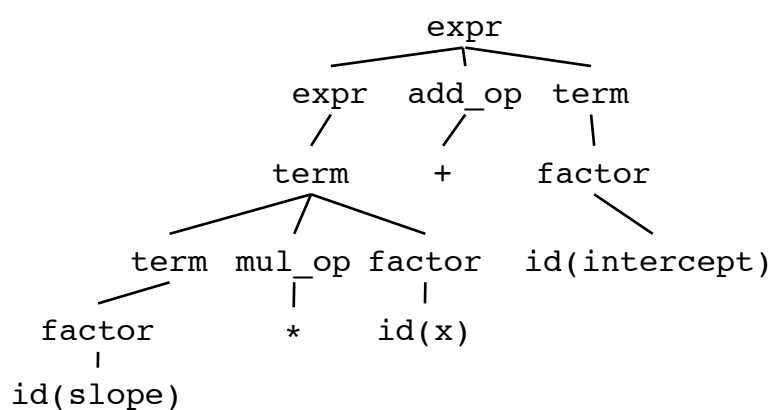
Improved Unambiguous Context-Free Grammar

1. $\text{expr} \rightarrow \text{expr add_op term} \mid \text{term}$
2. $\text{term} \rightarrow \text{term mul_op factor} \mid \text{factor}$
3. $\text{factor} \rightarrow \text{id} \mid \text{number} \mid \text{-factor} \mid (\text{expr})$
3. $\text{add_op} \rightarrow + \mid -$
4. $\text{mul_op} \rightarrow * \mid /$

18

Revisit "slope * x + intercept"

- Parse Tree



19

Extended BNF (EBNF)

- There are extensions of BNF to simplify representation
 - Kleene star $*$ or $\{ \}$ to represent repetition (0 or more)
 - $()$ to represent alternative parts
 - $[]$ to represent optional parts
 - $\text{id_list} \rightarrow \text{id} (, \text{id})^*$
 - $\text{proc_call} \rightarrow \text{id} ('[\text{expr_list}]')$

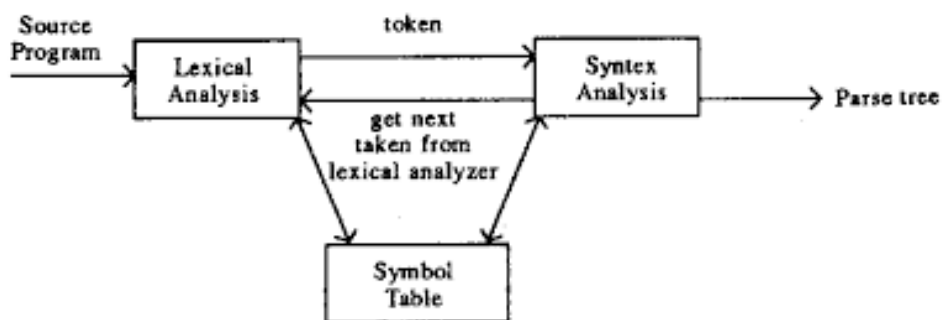
20

Lexical and Syntactic Analysis

- Two steps to discover the syntactic structure of a program
 - Lexical analysis (Scanner): to read the input characters and output a sequence of tokens
 - Syntactic analysis (Parser): to read the tokens and output a parse tree and report syntax errors if any

21

Interaction between lexical analysis and syntactic analysis



22

Scanner

- Pattern matcher for character strings
 - If a character sequence matches a pattern, it is identified as a token
- Responsibilities
 - Tokenize source, report lexical errors if any, remove comments and whitespace, save text of interesting tokens, save source locations, (optional) expand macros and implement preprocessor functions

23

Tokenizing Source

- Given a program, identify all lexemes and their categories (tokens)

24