

## Computer Language Engineering

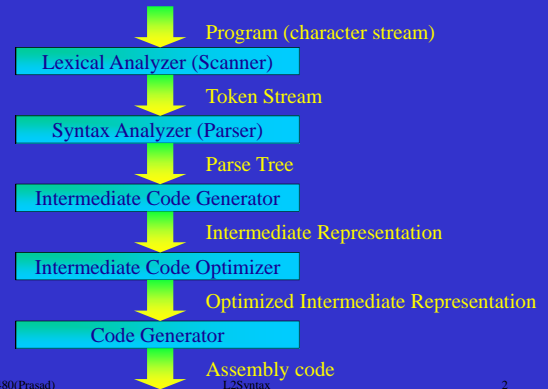
- How to give instructions to a computer?
  - **Programming Languages.**
- How to make the computer carryout the instructions efficiently?
  - **Compilers.**

cs480(Prasad)

L2Syntax

1

## Anatomy of a Compiler



cs480(Prasad)

L2Syntax

2

토큰단위로 인식

## What is a Lexical Analyzer?

Source program text → Tokens

- Examples of Token
  - Operators      = + - > ( { := == <>
  - Keywords      if while for int double
  - Numeric literals    43 6.035 -3.6e10 0x13F3A
  - Character literals   'a' '~' '\'
  - String literals    "3.142" "Fall" "\\\" = empty"
- Examples of non-token
  - White space      space(' ') tab('\t') end-of-line('\n')
  - Comments        /\*this is not a token\*/

cs480(Prasad)

L2Syntax

3

## Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0		v	a	r	l		<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	---	---	---	---	--	---	---	--

**for**      ID("var1")   eq\_op   Num(10)   ID("var1")   leq\_op

- Partition input program text into sequence of tokens, attaching corresponding attributes.
  - E.g., C-lexeme "015" token NUM attribute 13
- Eliminate white space and comments

cs480(Prasad)

L2Syntax

4

## Syntax and Semantics of a programming language

- Syntax
  - What is the structure of the program?
  - Textual representation.
    - Formally defined using context-free grammars (Backus-Naur Formalism)
- Semantics
  - What is the meaning of a program?
    - Harder to give a mathematical definition.

cs480(Prasad)

L2Syntax

5

## Input to and output of a parser

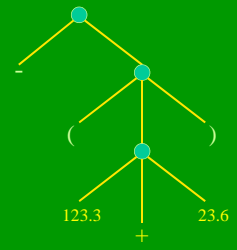
Input: - (123.3 + 23.6)

### Token Stream

minus\_op  
left\_paren\_op  
num(123.3)  
plus\_op  
num(23.6)  
right\_paren\_op

Syntax Analyzer (Parser)

### Parse Tree



cs480(Prasad)

L2Syntax

6

## Example: A CFG for expressions

- Simple arithmetic expressions with + and \*
  - 8.2 + 35.6
  - 8.32 + 86 \* 45.3
  - (6.001 + 6.004) \* (6.035 \* -(6.042 + 6.046))
- Terminals (or tokens)
  - **num** for all the numbers
  - plus\_op, minus\_op, times\_op, left\_paren\_op, right\_paren\_op
- What is the grammar for all possible expressions?

cs480(Prasad)

L2Syntax

7

## Example: A CFG for expressions

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow ( \langle \text{expr} \rangle )$   
 $\langle \text{expr} \rangle \rightarrow - \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \text{num}$   
 $\langle \text{op} \rangle \rightarrow +$   
 $\langle \text{op} \rangle \rightarrow *$

cs480(Prasad)

L2Syntax

8

## Context-Free Grammars (CFGs)

- *Terminals*
  - Symbols for strings or tokens { **num**, (, ), +, \*, - }
- *Nonterminals*
  - Syntactic variables { **<expr>**, **<op>** }
- *Start symbol*
  - A special nonterminal **<expr>**
- *Productions*
  - The manner in which terminals and nonterminals are combined to form strings.
  - A nonterminal in LHS and a string of terminals and nonterminals in RHS. **<expr>** → - **<expr>**

cs480(Prasad)

L2Syntax

9

## English Language

- Letters: **a,b,c, ...**
- Alphabet: { **a,b, ..., z** } U { **A, B, ..., Z** }
- Tokens (Terminals) : **English words**
- Nonterminals: **<Sentence>**, **<Subject>**, **<Clause>**, ...
- Context-free aspect: **Replacing <Noun> by <Proper Noun> or <Common Noun>**
- Context-sensitive aspects: **Agreement among words w.r.t. number, gender, tense, etc.**

cs480(Prasad)

L2Syntax

10

## Example Derivation

**<expr>** ⇒ **<expr> <op> <expr>**  
⇒ **num <op> <expr>**  
⇒ **num \* <expr>**  
⇒ **num \* ( <expr> )**  
⇒ **num \* ( <expr> <op> <expr> )**  
⇒ **num \* ( num <op> <expr> )**  
⇒ **num \* ( num + <expr> )**  
⇒ **num \* ( num + num )**  
⇒ **num \* ( num + num )**

cs480(Prasad)

L2Syntax

11

## Another Example Derivation

**<expr>** ⇒ **<expr> <op> <expr>**  
⇒ **<expr> \* <expr>**  
⇒ **<expr> \* ( <expr> )**  
⇒ **<expr> \* ( <expr> <op> <expr> )**  
⇒ **<expr> \* ( num <op> <expr> )**  
⇒ **<expr> \* ( num + <expr> )**  
⇒ **<expr> \* ( num + num )**  
⇒ **num \* ( num + num )**

cs480(Prasad)

L2Syntax

12

## Parse/Derivation Tree

- Graphical Representation of the parsed structure
- Shows the sequence of derivations performed
  - Internal nodes are non-terminals.
  - Leaves are terminals.
  - Each parent node is LHS and the children are RHS of a production.
- *Abstracts* the details of sequencing of the rule applications, but *preserves* decomposition of a non-terminal.

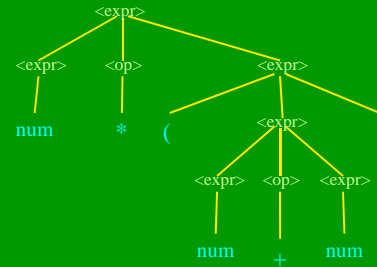
cs480(Prasad)

L2Syntax

13

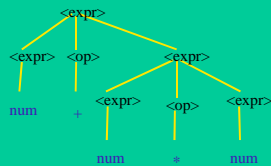
## Parse/Derivation Tree Example

num \* ( num + num )



## Expression – One possible parse tree

num + num \* num



124 + (23.5 \* 86) = 2145

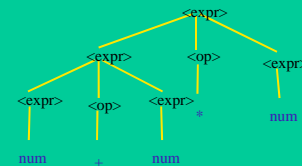
cs480(Prasad)

L2Syntax

15

## Expression – Another possible parse tree

num + num \* num



(124 + 23.5) \* 86 = 12685

cs480(Prasad)

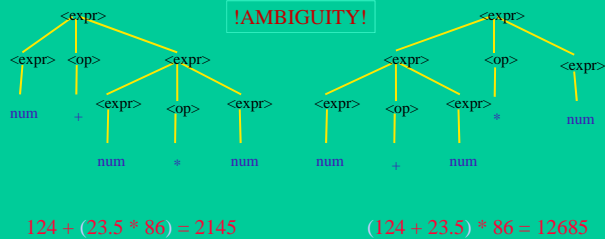
L2Syntax

16

## Same string – Two distinct parse (derivation) trees

num + num \* num

!AMBIGUITY!



cs480(Prasad)

L2Syntax

17

## Ambiguous Grammar

- Applying different derivation orders produces different parse trees.
  - This can lead to ambiguous/unexpected results.
  - Note that multiple derivations leading to the same parse tree is not a problem.
- A CFG is *ambiguous* if the same string can be associated with two distinct parse trees.
  - E.g., The expression grammar can be shown to be *ambiguous* using expressions containing *binary infix operators* and *non-fully parenthesized expressions*.

cs480(Prasad)

L2Syntax

18

## Removing Ambiguity

- Sometimes rewriting a grammar to reflect *operator precedence* with additional nonterminals will eliminate ambiguity.
  - \* more binding than +.
  - && has precedence over ||.
- One can view the rewrite as “breaking the symmetry” through “stratification”.

cs480(Prasad)

L2Syntax

19

## Eliminating Ambiguity

<expr> → <term> + <expr>  
 <expr> → <term>  
 <term> → <unit> \* <term>  
 <term> → <unit>  
 <unit> → num  
 <unit> → ( <expr> )

<expr> → <expr> <op> <expr>  
 <expr> → ( <expr> )  
 <expr> → num  
 <op> → +  
 <op> → \*

cs480(Prasad)

L2Syntax

20

### Extended BNF (Backus Naur Formalism)

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle ( + \langle \text{term} \rangle )^*$   
 $\langle \text{term} \rangle \rightarrow \langle \text{unit} \rangle ( * \langle \text{unit} \rangle )^*$   
 $\langle \text{unit} \rangle \rightarrow \text{num} \mid ( \langle \text{expr} \rangle )$

★ Kleene Star :  
Zero or more  
occurrences

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle ( \langle \text{op} \rangle \langle \text{expr} \rangle )^*$   
 $\langle \text{expr} \rangle \rightarrow ( \langle \text{expr} \rangle )$   
 $\langle \text{expr} \rangle \rightarrow \text{num}$   
 $\langle \text{op} \rangle \rightarrow + \mid *$

cs480(Prasad)

L2Syntax

21

### Expression Parser Fragment (String => Parse Tree)

```

Node expr() {
    // PRE: Expects lookahead token.
    // POST: Consume an Expression
    // and update Lookahead token.
    Node temp = term();
    while ( inTok.ttype == '+' ) {
        inTok.nextToken();
        Node templ = term();
        temp = new OpNode(temp, '+', templ);
    }
    return temp;
}

```

cs480(Prasad)

L2Syntax

22

### Arithmetic Expressions (ambiguous) (with variables)

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle$   
 $\quad \quad \quad \mid \langle \text{expr} \rangle * \langle \text{expr} \rangle$   
 $\quad \quad \quad \mid ( \langle \text{expr} \rangle )$   
 $\quad \quad \quad \mid \langle \text{variable} \rangle$   
 $\quad \quad \quad \mid \langle \text{constant} \rangle$

$\langle \text{variable} \rangle \rightarrow x \mid y \mid z$

$\langle \text{constant} \rangle \rightarrow 0 \mid 1 \mid 2$

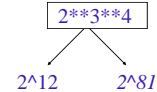
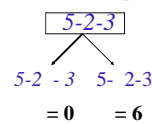
cs480(Prasad)

L2Syntax

23

### Resolving Ambiguity in Expressions

- Different operators : *precedence relation*
- Same operator : *associativity relation*



Left Associative ( (5-2)-3 ) Right Associative ( 2\*\*(3\*\*4) )

cs480(Prasad)

L2Syntax

24

## C++ Operator Precedence and Associativity

Level	Operator	Function	Level	Operator	Function
17R	::	global scope (unary)	12L	+, -	arithmetic operators
17L	::	class scope (binary)	11L	<<, >>	bitwise shift
16L	->, .	member selectors	10L	<, <=, >, >=	relational operators
16L	[]	array index	9L	==, !=	equality, inequality
16L	()	function call	8L	&	bitwise AND
16L	()	type construction	7L	^	bitwise XOR
15R	sizeof	size in bytes	6L		bitwise OR
15R	++, --	increment, decrement	5L	&&	logical AND
15R	~	bitwise NOT	4L		logical OR
15R	!	logical NOT	3L	?:	arithmetic if
15R	+, -	unary minus, plus	2R	=, *=, /=, %=, +=, -=	assignment operators
15R	*, &	dereference, address-of		<<=, >>=, &=,  =, ^=	
15R	()	type conversion (cast)	1L	,	comma operator
15R	new, delete	free store management			
14L	->*, *	member pointer select			
13L	*, /, %	multiplicative operators			

cs480(Prasad)

L2Syntax

25