



*COSC252: Programming Languages:*

*Language Basics and Overview*

Jeremy Bolton, PhD

Asst Teaching Professor

GEORGETOWN  
UNIVERSITY

# *Notes*

Notes: Read ALSU 1 – 2

Topics covered in CH1 will be revisited in more detail throughout

Sign up for Gradiance and begin Assignment #2

# *Outline*

- I. Language Basics
  - I. Definitions
- II. Grammars and Syntactic Analysis
  - I. BNF
  - II. Derivations
  - III. Parse trees
- III. Tokens and Lexical Analysis

# *The Basics of Programming Languages*



- Before we can learn about each of these steps, we will formalize the concepts and nomenclature

# *Basic Elements of a Language*

- Language: A set of sentences.
  - Exhaustive list (roster notation) is not practical. Usually a syntax is used.
- Sentence: a string of lexemes in a language (as specified by the syntax of the language)
- Syntax: A specification that can be used to determine whether a sentence is in a set (or not).
- Lexeme: Lowest level, integral, syntactic unit that has meaning
  - e.g. word, number, operator, ...
- Token: A lexeme group or lexeme category

# *The Basics of Programming Languages*



- Lexical Analysis: identifying and categorizing lexemes
- Syntactic Analysis: determining if sentence (string of lexemes) is in language.

# *Example: Identify the lexemes and tokens*

- Example: English Language
- Sample Sentence:
  - John has a cat, and Mary has a dog.

Lexeme	Token (lexeme category)
John	noun
has	Verb
a	article
cat	noun
,	comma
and	conjunction
...	

## *Example: Identify the lexemes and tokens*

- Example: Programming Language
- Sample Sentence:
  - `sum = sum + 1.0;`

Lexeme	Token (lexeme category)
sum	variable
=	assignment operator
sum	variable
+	plus operator
1.0	double
;	semicolon



# *Discussion: Syntactic Analysis*

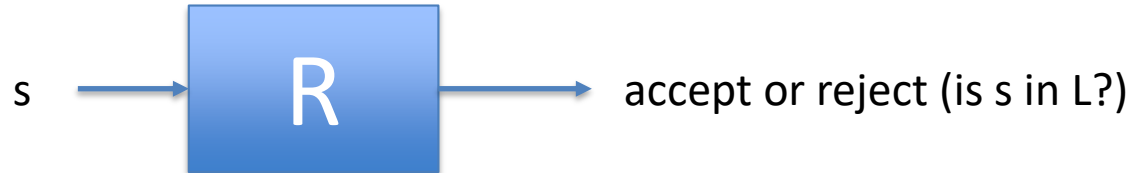
- Identifying and categorizing lexemes is often an easy task, but how do we determine whether a sentence is in a language?
- Example: Language L is the set of all English sentences.
  - $L = \{s \mid s \text{ is an English sentence}\}$
  - When given a sentence s, how do we determine whether it is in the English Language?
    - Let's say s = "The cat has a hat."
    - Option: Compare s to all sentences in L
      - Brute Force, Exhaustive search
      - No – not practical

# *Rules for a Language*

- Syntax Rules are used to characterize a language
- Lets try to develop some simple rules for English Sentences:

# Syntax

- How can we use syntax to define or characterize a language?
  - There is no “perfect” solution or recipe
- There are 2 categories for syntax implementation and therefore 2 categories for language characterization.
  - Recognizers



- Generators



# *Generators vs Recognizers*

- Both generators and recognizers are useful in different applications, but each is limited
- Recognizers
  - EG compiler, Clearly useful
  - Compilers are generally very complex: Not good for *communicating* a language or describing a language
- Generators
  - Useful for learning or describing a language
  - (Examples upcoming)

# *Recognizers*

- Recognizers take as input a sentence and determine whether the sentence is in the language.
  - Example: compiler
  - A compiler is an algorithmic representation of syntax
  - A compiler is a syntactic analyzer
    - The parsing portion of the compiler

# *Using Generators to Characterize a Language*

- Noam Chomsky
  - MIT Linguist
  - Published work on categories of Languages / Grammars, 1950s
    - Two categories are commonly used in Programming Languages
      - Context-Free Languages
        - » Often used to characterize programming language sentence structure
      - Regular Languages
        - » Often used to characterize the structure of lexemes / tokens
- John Backus and Peter Naur developed a formal notation for generating a Context Free Language, (similar to the notation used by Chomsky)
  - BNF (Bachus Naur Form)

# *Grammars*

- Generators are often implemented as a set of rules, called grammars.
- BNF grammars are grammars consisting of
  - terminals: lexemes (integral syntactic units)
  - non-terminals: abstract compositions of terminals
    - Have at least two possible forms
  - productions: composition rules for non-terminals.

# *BNF Example*

- BNF production for English grammar example

$\langle \textit{sentence} \rangle \rightarrow \langle \textit{noun} \rangle \langle \textit{verb} \rangle \langle \textit{adjective} \rangle \langle \textit{period} \rangle$

Each line in BNF is a rule or *production*. Each production consists of a left hand side (LHS), followed by an arrow (“is produced by”), followed by a right hand side (RHS).

- The production defines the LHS abstraction, or *non-terminal*
  - Here sentence is defined to be produced by a noun followed by a verb followed by an adjective followed by a period.
- Note this production requires that we define noun, verb, ... as well.



# *BNF Example*

- BNF snippet for English grammar example

$\langle \textit{sentence} \rangle \rightarrow \langle \textit{noun} \rangle \langle \textit{verb} \rangle \langle \textit{adjective} \rangle \langle \textit{period} \rangle$

$\langle \textit{noun} \rangle \rightarrow \textit{Bob}$

$\langle \textit{verb} \rangle \rightarrow \textit{is}$

$\langle \textit{adjective} \rangle \rightarrow \textit{tall}$

$\langle \textit{period} \rangle \rightarrow .$

Each line in BNF is a rule or *production*. Each production consists of a left hand side (LHS), followed by an arrow (“derives”), followed by a right hand side (RHS).

Each non-terminal is enclosed by a “< >”. Each non-terminal should have a defining production.

All other character strings, *lexemes*, are called *terminals* in this context.

# *BNF Example*

- BNF for English grammar example

$\langle \textit{sentence} \rangle \rightarrow \langle \textit{noun} \rangle \langle \textit{verb} \rangle \langle \textit{adjective} \rangle \langle \textit{period} \rangle$

$\langle \textit{noun} \rangle \rightarrow \textit{Bob}$

$\langle \textit{verb} \rangle \rightarrow \textit{is}$

$\langle \textit{adjective} \rangle \rightarrow \textit{tall}$

$\langle \textit{period} \rangle \rightarrow .$

What sentence(s) can be generated using this BNF grammar?  
What sentence(s) are in this Language?  
How many sentence(s) are in this Language?

# *BNF Example*

- BNF snippet for Programming grammar example

$\langle \textit{assignment} \rangle \rightarrow \langle \textit{variable} \rangle = \langle \textit{expression} \rangle$

$\langle \textit{expression} \rangle \rightarrow \langle \textit{variable} \rangle + \langle \textit{variable} \rangle$

$\langle \textit{expression} \rangle \rightarrow \langle \textit{variable} \rangle - \langle \textit{variable} \rangle$

Use the “|” symbol to represent the logical OR. Multiple definitions of each non-terminal can be written in the same production.

$\langle \textit{expression} \rangle \rightarrow \langle \textit{variable} \rangle + \langle \textit{variable} \rangle \mid$   
 $\langle \textit{variable} \rangle - \langle \textit{variable} \rangle$

## *Listing out non-terminals*

- Non-terminals can generally take on many forms
- How can we practically list out all possible forms of a non-terminal?
- Example  $\langle \text{expr} \rangle$

$\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle + \langle \text{var} \rangle \mid$   
 $\quad \langle \text{var} \rangle + \langle \text{var} \rangle + \langle \text{var} \rangle$   
 $\quad \langle \text{var} \rangle + \langle \text{var} \rangle + \langle \text{var} \rangle + \langle \text{var} \rangle$   
 $\quad \dots$

# *Listing Out Non-Terminals*

$\langle expr \rangle \rightarrow \langle var \rangle + \langle var \rangle \mid$   
 $\langle var \rangle + \langle var \rangle + \langle var \rangle$   
 $\langle var \rangle + \langle var \rangle + \langle var \rangle + \langle var \rangle$   
...

The above is tedious and impractical. We can avoid such definitions by recursively defining non-terminals.

$\langle expr \rangle \rightarrow \langle expr \rangle + \langle var \rangle \mid \langle var \rangle$

# Generating Sentences using BNF

- **Start Symbol:** One unique abstraction used to start a derivation
- Derivation: a repeated application of BNF rules

A simple but complete grammar

$\langle \textit{program} \rangle \rightarrow \textit{begin} \langle \textit{stmt} \rangle \textit{end}$

$\langle \textit{stmt} \rangle \rightarrow \langle \textit{var} \rangle = \langle \textit{expr} \rangle$

$\langle \textit{expr} \rangle \rightarrow \langle \textit{var} \rangle \mid \langle \textit{var} \rangle + \langle \textit{var} \rangle \mid \langle \textit{var} \rangle - \langle \textit{var} \rangle$

$\langle \textit{var} \rangle \rightarrow x \mid y \mid z$

# Generation / Derivation Example

- Derivation: a repeated application of BNF rules
  - Starting with start symbol, repeatedly replace a non-terminal using one of the production rules until no non-terminals remain
  - Only reduce one non-terminal per line of derivation – it is best not to perform multiple steps per line
- Try to generate: `begin x = y + z end`

*< program > → begin < stmt > end*

*< stmt > → < var > = < expr >*

*< expr > → < var > | < var > + < var > | < var > - < var >*

*< var > → x | y | z*

• *< program >*

# Derivation Practice

- Try to derive sentence: `begin y = x - y end`

*< program >* → *begin < stmt > end*

*< stmt >* → *< var > = < expr >*

*< expr >* → *< var > | < var > + < var > | < var > - < var >*

*< var >* → *x | y | z*

- *< program >*

- *begin < stmt > end*

- *begin < var > = < expr > end*

- *begin y = < expr > end*

- *begin y = < var > - < var > end*

- *begin y = x - < var > end*

- *begin y = x - y end*



# Derivation Practice

- Try to derive sentence: `begin y = x - y + z end`

*< program > → begin < stmt > end*

*< stmt > → < var > = < expr >*

*< expr > → < var > | < var > + < var > | < var > - < var >*

*< var > → x | y | z*

- *< program >*
- *begin < stmt > end*
- *begin < var > = < expr > end*
- *begin y = < expr > end*
- *begin y = < var > - < var > end*
- ...
- Cannot! Need a more general production rule
- What rule can we add?

# Derivation example with recursive production

- Try to derive sentence: `begin y = x - y + z end`

$\langle \text{program} \rangle \rightarrow \text{begin } \langle \text{stmt} \rangle \text{ end}$

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$

$\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{expr} \rangle + \langle \text{var} \rangle$   
 $\mid \langle \text{expr} \rangle - \langle \text{var} \rangle$

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

- $\langle \text{program} \rangle$
- $\text{begin } \langle \text{stmt} \rangle \text{ end}$
- $\text{begin } \langle \text{var} \rangle = \langle \text{expr} \rangle \text{ end}$
- $\text{begin } y = \langle \text{expr} \rangle \text{ end}$
- $\text{begin } y = \langle \text{expr} \rangle + \langle \text{var} \rangle \text{ end}$
- $\text{begin } y = \langle \text{expr} \rangle - \langle \text{var} \rangle + \langle \text{var} \rangle \text{ end}$
- $\text{begin } y = \langle \text{var} \rangle - \langle \text{var} \rangle + \langle \text{var} \rangle \text{ end}$
- $\text{begin } y = x - \langle \text{var} \rangle + \langle \text{var} \rangle \text{ end}$
- $\text{begin } y = x - y + \langle \text{var} \rangle \text{ end}$
- $\text{begin } y = x - y + z \text{ end}$

## *LMD vs RMD*

- Note: Multiple derivations *may* exist per sentence.
  - Many sentences have multiple derivations
- Two common derivation strategies
  - Left Most Derivation: During derivation, always reduce left-most non-terminal at each step.
  - Right Most Derivation: During derivation, always reduce right-most non-terminal at each step.
- Note: Derivation order has no effect on the language generated

# *LMD vs RMD: begin $y = x - y$ end*

## LMD

- *< program >*
- *begin < stmt > end*
- *begin < var > = < expr > end*
- *begin y = < expr > end*
- *begin y = < var > - < var > end*
- *begin y = x - < var > end*
- *begin y = x - y end*

## RMD

- *< program >*
- *begin < stmt > end*
- *begin < var > = < expr > end*
- *begin < var > = < var > - < var > end*
- *begin < var > = < var > - y end*
- *begin < var > = x - y end*
- *begin y = x - y end*

Note: Any intermediate representation derived from the start symbol called a *sentinel form*. If derived via a {LMD, RMD} it is a {left, right} *sentinel form*.

*< program >* → *begin < stmt > end*

*< stmt >* → *< var > = < expr >*

*< expr >* → *< var > | < var > + < var > | < var > - < var >*

*< var >* → *x | y | z*

# *In-class Exercise*

*< program >* → *begin < stmt > end*

*< stmt >* → *< var > = < expr >*

*< expr >* → *< var > | < expr > + < var > | < expr > - < var >*

*< var >* → *x | y | z*

- Derive (LMD or RMD): *begin x = y + z - x end*

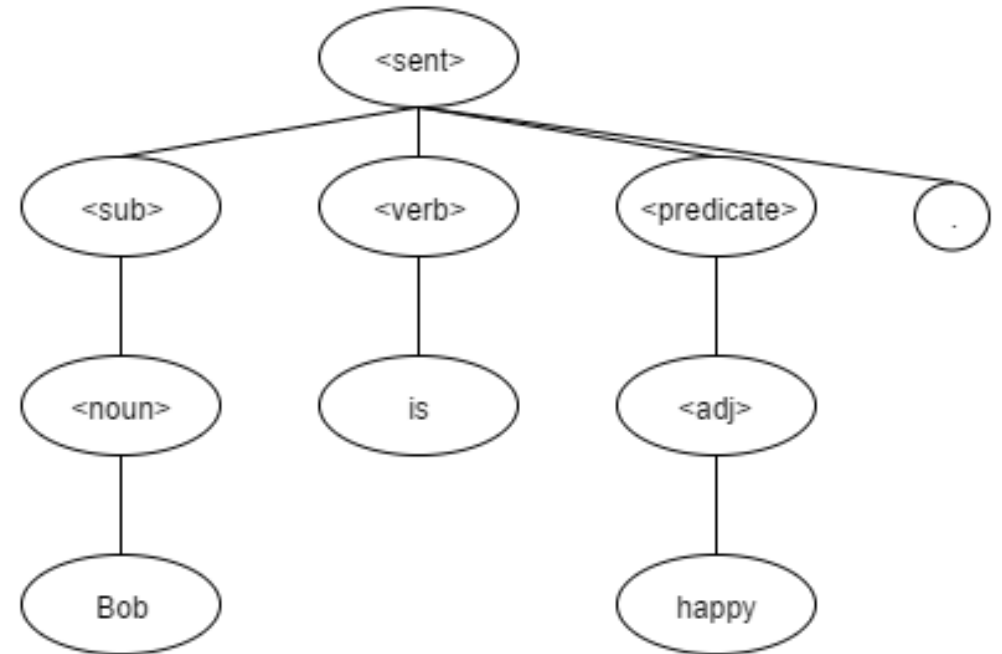
## *In Class Exercise*

- Construct a grammar that allows programs with operators +, -, /, and \*.

# Parse Trees

- Hierarchical representation of the syntactic structure of a sentence
  - Based on syntax rules / derivation
  - *LHS of production is parent node and RHS are children nodes*
  - Example: Bob is happy.

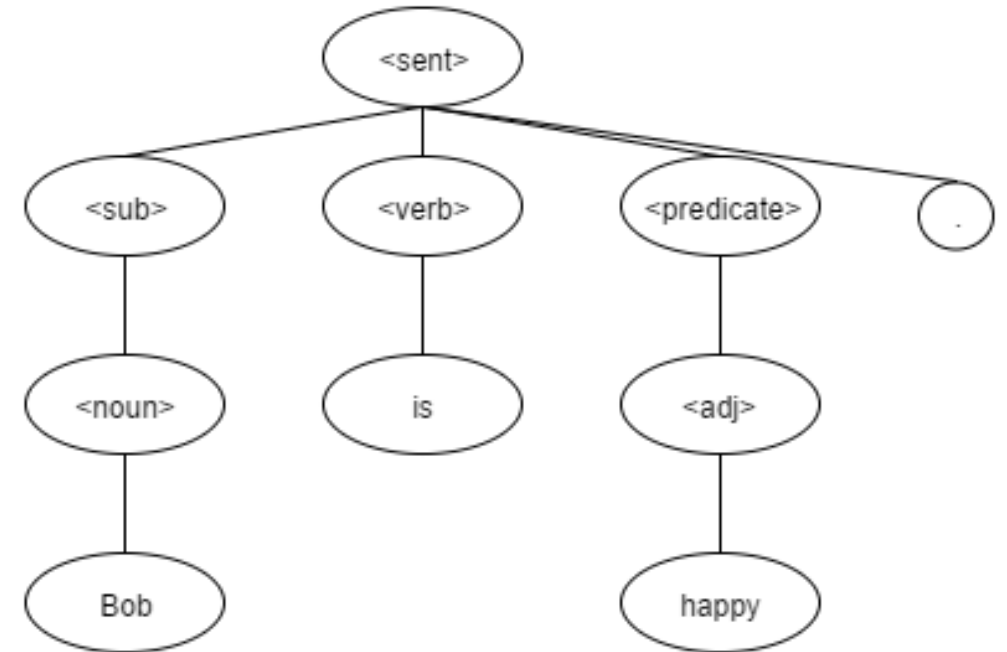
$\langle sent \rangle \rightarrow \langle sub \rangle \langle verb \rangle \langle predicate \rangle .$   
 $\langle sub \rangle \rightarrow \langle noun \rangle$   
 $\langle predicate \rangle \rightarrow \langle art \rangle \langle noun \rangle \mid \langle adj \rangle$   
 $\langle noun \rangle \rightarrow Bob \mid Mary$   
 $\langle art \rangle \rightarrow a \mid the$   
 $\langle verb \rangle \rightarrow is$   
 $\langle adj \rangle \rightarrow tall \mid happy \mid sad$



# Parse Trees

- Observations
  - Root: start symbol
  - Internal nodes are non-terminals
  - Leaf nodes are terminals

$\langle sent \rangle \rightarrow \langle sub \rangle \langle verb \rangle \langle predicate \rangle .$   
 $\langle sub \rangle \rightarrow \langle noun \rangle$   
 $\langle predicate \rangle \rightarrow \langle art \rangle \langle noun \rangle \mid \langle adj \rangle$   
 $\langle noun \rangle \rightarrow Bob \mid Mary$   
 $\langle art \rangle \rightarrow a \mid the$   
 $\langle verb \rangle \rightarrow is$   
 $\langle adj \rangle \rightarrow tall \mid happy \mid sad$





# Parse Tree (programming) Example

- Example
  - Based on syntax rules / derivation
  - Example: begin x = x + y end

<program>

$\langle program \rangle \rightarrow begin \langle stmt \rangle end$

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$

$\langle expr \rangle \rightarrow \langle var \rangle \mid \langle expr \rangle + \langle var \rangle$

$\mid \langle expr \rangle - \langle var \rangle$

$\langle var \rangle \rightarrow x \mid y \mid z$

# *In Class -- Parse Tree (programming) Example*

- Example: create parse tree

$x = x * y + z$

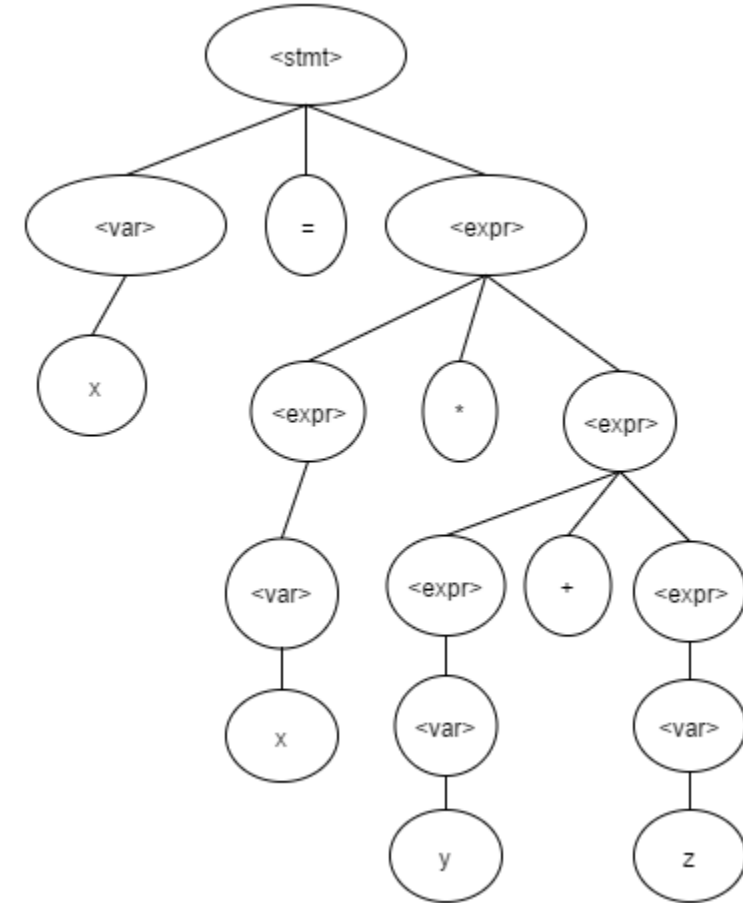
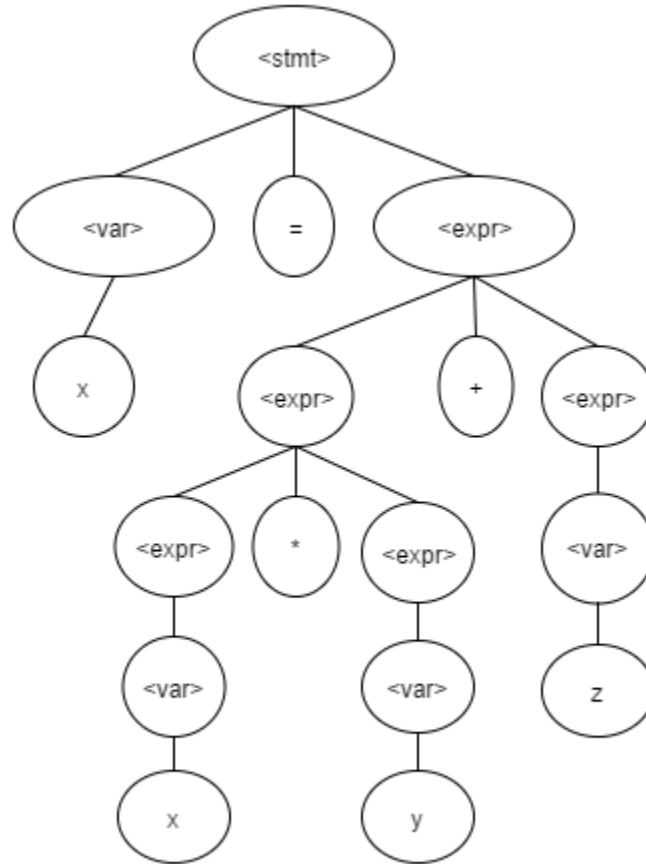
$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$

$\langle expr \rangle \rightarrow \langle var \rangle \mid \langle expr \rangle + \langle expr \rangle \mid \langle expr \rangle * \langle expr \rangle$

$\langle var \rangle \rightarrow x \mid y \mid z$

# Ambiguous Grammars

- A grammar that generates a sentence with two or more possible parse trees is called an **ambiguous grammar**
- A grammar is ambiguous if it produces more than one left most derivation or more than one right most derivation
- Example:  
 $x = x * y + z$



$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle * \langle \text{expr} \rangle$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z$

# Ambiguous Grammars

- Inspection: Why is this grammar ambiguous ...

$$\begin{aligned} \langle stmt \rangle &\rightarrow \langle var \rangle = \langle expr \rangle \\ \langle expr \rangle &\rightarrow \langle var \rangle \mid \boxed{\langle expr \rangle} + \boxed{\langle expr \rangle} \mid \boxed{\langle expr \rangle} * \boxed{\langle expr \rangle} \\ \langle var \rangle &\rightarrow x \mid y \mid z \end{aligned}$$

- The  $\langle expr \rangle$  non terminal can be expanded on either side of the  $+$  or  $*$  operator. We can derive the either side first **possibly** giving rise to two distinct parse trees
- Note: Syntactically this may be OK, but in practice, this may cause semantic issues

# *Ambiguous Grammar : English\**

- I see Bob on the roof with binoculars.

# Disambiguating Grammars

- This grammar is ambiguous as there is the same non-terminal on either side of the '+' and '\*' operators, which may permit multiple parse trees

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$

$\langle expr \rangle \rightarrow \langle var \rangle \mid \langle expr \rangle + \langle expr \rangle \mid \langle expr \rangle * \langle expr \rangle$

$\langle var \rangle \rightarrow x \mid y \mid z$

- Fix:
  - Change the **associativity** of the operator within the productions.
  - Make each operator either *left associative* or *right associative*, but not both.
  - Associativity can dictate derivation and disambiguate a grammar.

# Operator Associativity

- A production is left recursive if its LHS appears at the beginning of the RHS. This implies left associativity.
  - The following grammar is left associative and unambiguous.

$$\begin{aligned} \langle stmt \rangle &\rightarrow \langle var \rangle = \langle expr \rangle \\ \langle expr \rangle &\rightarrow \langle var \rangle \mid \langle expr \rangle + \langle var \rangle \mid \langle expr \rangle * \langle var \rangle \\ \langle var \rangle &\rightarrow x \mid y \mid z \end{aligned}$$

- A production is right recursive if its LHS appears at the beginning of the RHS. This implies right associativity.
  - The following grammar is right associative and unambiguous.

$$\begin{aligned} \langle stmt \rangle &\rightarrow \langle var \rangle = \langle expr \rangle \\ \langle expr \rangle &\rightarrow \langle var \rangle \mid \langle var \rangle + \langle expr \rangle \mid \langle var \rangle * \langle expr \rangle \\ \langle var \rangle &\rightarrow x \mid y \mid z \end{aligned}$$

# Parse Tree Example

- Lets attempt our previous example using a non-ambiguous grammar.
  - Example:  $x = x * y + z$

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$

$\langle expr \rangle \rightarrow \langle var \rangle \mid \langle expr \rangle + \langle var \rangle \mid \langle expr \rangle * \langle var \rangle$

$\langle var \rangle \rightarrow x \mid y \mid z$



# *Does associativity solve all language issues?*

- NO – Associativity may help with ambiguity issues, but *semantic issues* may still remain.
  - Example 1: operator precedence.
  - Example 2: the dangling “else” problem.
- Operator Precedence Concern: apply semantics to “  $x = x + y * z$  ”
  - Semantic Goal: evaluate the multiplication first!
  - Note: Semantics is applied to the parsed structure – parse tree.
    - The lower in the parse tree, the higher the precedence. That is, items lower on the parse tree (along the same branch) are evaluated first

# Example: Operator Precedence

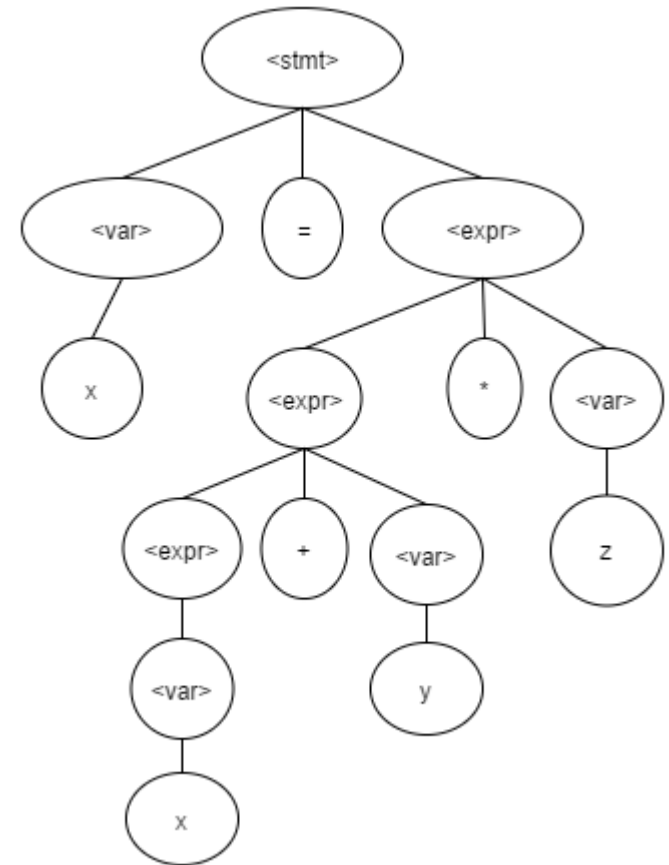
- Parse “  $x = x + y * z$  ” , using the following grammar

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$

$\langle expr \rangle \rightarrow \langle var \rangle \mid \langle expr \rangle + \langle var \rangle \mid \langle expr \rangle * \langle var \rangle$

$\langle var \rangle \rightarrow x \mid y \mid z$

- Using Left Associative Productions leads to improper parse tree
  - Addition is lower in the parse tree and is evaluated first
- Question
  - Does using a Right Associative production fix the issue?
  - For this one example, yes, but for the general case **NO!**
  - Associativity will change the resulting parse tree based on the order each of the operators or symbols, but cannot be used alone to enforce operator precedence.
- Perspective: Parsing and evaluation as a depth first traversal
  - Post Order!
  - Parsing can be viewed as a downward operation – traversing the tree downward. Applying semantics can be viewed as an upward operation, performed when traversing back up the tree.
- So ... How can we enforce operator precedence?



# Operator Precedence

- How can we enforce operator precedence?
  - With the addition of non-terminals and rules
- Example: Require operands of the '+' operator to consist of variables or products of variables (but not the reverse – do not allow operands of products to consist of additions of variables). This will ensure that multiplication is lower on the parse tree.

No enforced precedence

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle var \rangle \mid \langle expr \rangle + \langle var \rangle \mid \langle expr \rangle * \langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$

'\*' has precedence over '+'

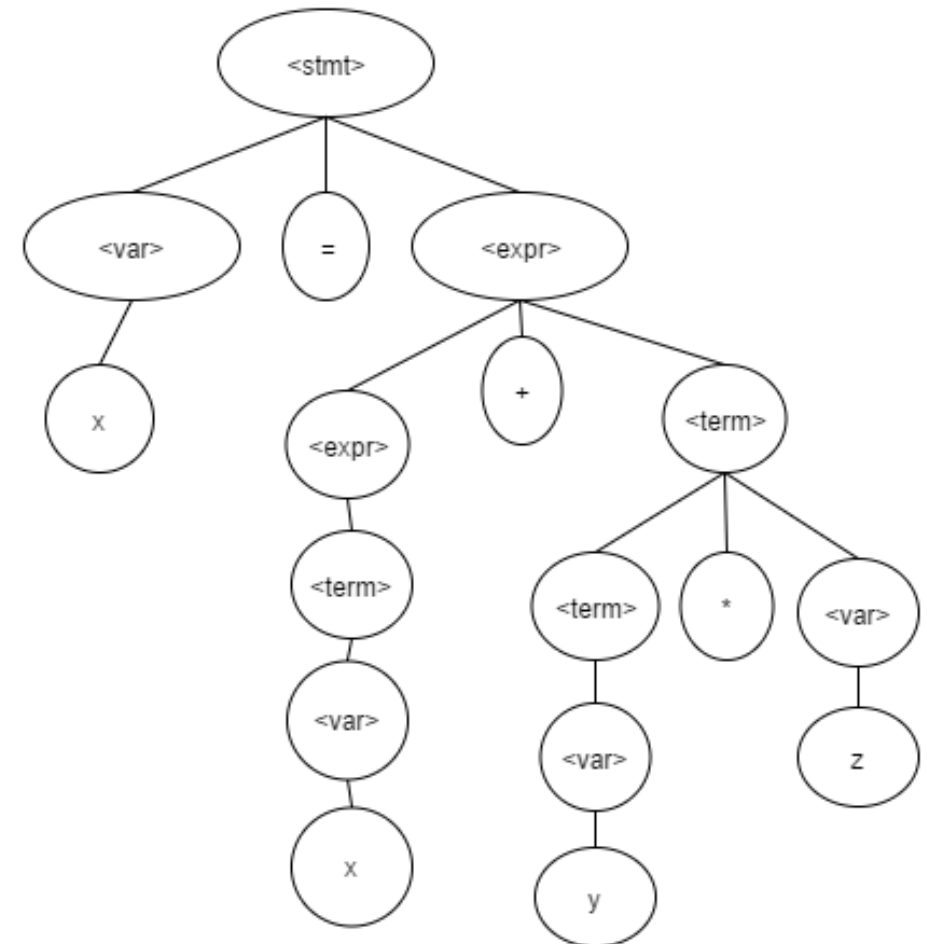
$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle$   
 $\langle term \rangle \rightarrow \langle var \rangle \mid \langle term \rangle * \langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$

# Example: Operator Precedence

- Parse “  $x = x + y * z$  ”, using the following grammar

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{term} \rangle * \langle \text{var} \rangle$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z$

- Since the operands of expressions are composed of terms, (and terms are not composed of expressions), terms will appear lower in the parse tree



# *In Class Example*

- Given the following grammar, parse the following sentences (if possible).

$$\begin{aligned} \langle \text{stmt} \rangle &\rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle \\ \langle \text{expr} \rangle &\rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{var} \rangle \mid \langle \text{expr} \rangle - \langle \text{var} \rangle \\ \langle \text{term} \rangle &\rightarrow \langle \text{var} \rangle \mid \langle \text{term} \rangle * \langle \text{var} \rangle \mid \langle \text{term} \rangle / \langle \text{var} \rangle \\ \langle \text{var} \rangle &\rightarrow x \mid y \mid z \end{aligned}$$

- Sentences:

$$x = x * y / z$$

$$y = z + x / y$$

$$z = -x / y + z$$

## *Example: non-ideal grammar*

- Using the following grammar, try to parse sentences

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{term} \rangle + \langle \text{var} \rangle \mid \langle \text{expr} \rangle - \langle \text{var} \rangle$

$\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle * \langle \text{var} \rangle \mid \langle \text{term} \rangle / \langle \text{var} \rangle$

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

- Which of these sentences are in the language? (ie which can we parse)

$x = x * y / z$

$y = z + x / y$

$z = -x / y + z$

- What is “wrong” with this grammar? How can we “fix” it?

# Designing a Grammar

- Each operator with different precedence should have its own non-terminal with production(s)
- The productions should be designed to enforce the correct precedence
- Example: Add the unary negation operator to the existing grammar

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\langle term \rangle \rightarrow \langle var \rangle \mid \langle term \rangle * \langle var \rangle \mid \langle term \rangle / \langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\langle term \rangle \rightarrow \langle neg \rangle \mid \langle term \rangle * \langle neg \rangle \mid \langle term \rangle / \langle neg \rangle$   
 $\langle neg \rangle \rightarrow \langle var \rangle \mid -\langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$

# *In-Class Example*

- Everyone try: Update the following grammar to allow for parenthesis within a mathematical expression
  - Questions to ask yourself
    - What precedence level is a parenthetical expression?
    - Given that, where should this new rule go within the BNF?
      - Should we add a new non-terminal and production, or a new production for an existing non-terminal?
      - What is the form of all possible parenthetical statements, ie what is the form of the production?

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\langle term \rangle \rightarrow \langle neg \rangle \mid \langle term \rangle * \langle neg \rangle \mid \langle term \rangle / \langle neg \rangle$   
 $\langle neg \rangle \rightarrow \langle var \rangle \mid -\langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$



# Example (cont)

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\langle term \rangle \rightarrow \langle neg \rangle \mid \langle term \rangle * \langle neg \rangle \mid \langle term \rangle / \langle neg \rangle$   
 $\langle neg \rangle \rightarrow \langle var \rangle \mid -\langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$

- Placing the parenthetical near the top of the BNF hierarchy is intuitive in many ways; HOWEVER, given its general use and *semantics*, this construct should be lower in the hierarchy.
- Why? BNF is a hierarchical construct, building a abstractions out of lower level abstractions. The lower the abstraction the more atomic (indivisible) ... in the context of operators: the higher the precedence.

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\mid ( \langle expr \rangle )$   
 $\langle term \rangle \rightarrow \langle neg \rangle \mid \langle term \rangle * \langle neg \rangle \mid \langle term \rangle / \langle neg \rangle$   
 $\langle neg \rangle \rightarrow \langle var \rangle \mid -\langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$

- Will this generate?:

$$x = x + ( y + z )$$

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\mid ( \langle expr \rangle ) \mid \langle expr \rangle + ( \langle expr \rangle ) \mid ( \langle expr \rangle ) + \langle expr \rangle \mid \dots$   
 $\langle term \rangle \rightarrow \langle neg \rangle \mid \langle term \rangle * \langle neg \rangle \mid \langle term \rangle / \langle neg \rangle$   
 $\langle neg \rangle \rightarrow \langle var \rangle \mid -\langle var \rangle$   
 $\langle var \rangle \rightarrow x \mid y \mid z$

- Will this generate?:

$$x = x + ( y + z ) \text{ yes}$$

$$x = x * ( x + y ) \text{ no}$$

# Example (cont)

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{neg} \rangle \mid \langle \text{term} \rangle * \langle \text{neg} \rangle \mid \langle \text{term} \rangle / \langle \text{neg} \rangle$   
 $\langle \text{neg} \rangle \rightarrow \langle \text{var} \rangle \mid -\langle \text{neg} \rangle$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z$

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{neg} \rangle \mid \langle \text{term} \rangle * \langle \text{neg} \rangle \mid \langle \text{term} \rangle / \langle \text{neg} \rangle$   
 $\langle \text{neg} \rangle \rightarrow \langle p \rangle \mid -\langle \text{neg} \rangle$   
 $\langle p \rangle \rightarrow (\langle \text{expr} \rangle) \mid \langle \text{var} \rangle$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z$

- When adding a construct to a grammar (with its own level of precedence ... indivisibility), add a new production at the appropriate level of the hierarchy.

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{neg} \rangle \mid \langle \text{term} \rangle * \langle \text{neg} \rangle \mid \langle \text{term} \rangle / \langle \text{neg} \rangle$   
 $\langle \text{neg} \rangle \rightarrow \langle \text{var} \rangle \mid -\langle \text{neg} \rangle$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z \mid (\langle \text{expr} \rangle)$

Can we generate?:  $x = -(x + y)$   
YES

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{neg} \rangle \mid \langle \text{term} \rangle * \langle \text{neg} \rangle \mid \langle \text{term} \rangle / \langle \text{neg} \rangle$   
 $\langle \text{neg} \rangle \rightarrow \langle \text{var} \rangle \mid -\langle \text{neg} \rangle \mid (\langle \text{expr} \rangle)$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z$

Can we generate?:  $x = -(x + y)$   
YES

# *Else-Matching: the dangling else problem*

- Consider the following BNF for <if stmts>

$$\langle \text{if\_stmt} \rangle \rightarrow \text{if}(\langle \text{expr} \rangle) \langle \text{stmt} \rangle \mid \\ \text{if}(\langle \text{expr} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$$
$$\langle \text{stmt} \rangle \rightarrow \langle \text{if\_stmt} \rangle \mid \langle \text{assign} \rangle$$

...

- Consider the following code snippet and parse using grammar above:

```
if( x )
if( y )
    x = x + y
else
    y = x + y
```

# *Else-Matching: the dangling else problem*

- Consider the following BNF for <if stmts>

$$\begin{aligned} \langle \text{if\_stmt} \rangle \rightarrow & \text{if}(\langle \text{expr} \rangle) \langle \text{stmt} \rangle \mid \\ & \text{if}(\langle \text{expr} \rangle) \langle \text{stmt} \rangle \\ \text{else } \langle \text{stmt} \rangle \end{aligned}$$
$$\langle \text{stmt} \rangle \rightarrow \langle \text{if\_stmt} \rangle \mid \langle \text{assign} \rangle$$

...

- Consider the following code snippet and parse using grammar above:

```
if( x )
if( y )
    x = x + y
else
    y = x + y
```

# *Else-Matching: the dangling else problem (solution match else with nearest if)*

$\langle if\_stmt \rangle \rightarrow \langle matched \rangle \mid \langle unmatched \rangle$

$\langle matched \rangle \text{ if}(\langle expr \rangle) \langle matched \rangle \text{ else } \langle matched \rangle \mid$   
 $\langle assign \rangle$

$\langle unmatched \rangle \rightarrow \text{if}(\langle expr \rangle) \langle if\_stmt \rangle \mid$   
 $\text{if}(\langle expr \rangle) \langle matched \rangle \text{ else } \langle unmatched \rangle$

- If time, try to parse the following

if( x )

if( y )

    x = x + y

else

    y = x + y

## *If-else discussion*

- Observation: the dangling else problem arises from the combination of a matched and unmatched `if_stmt` . Allowing “flexible” matchings would produce an ambiguous grammar
- One can add a disambiguating rule to a grammar to enforce a matching scheme / rule. Note: This is not the only solution.
- What else might one do to disambiguate such a grammar?

# Example: A more complex grammar with PL like syntax

- Note operator precedence and composition rules of the productions.
- Exercise:
  - How can we permit sequences of statements?
  - How can we permit while loops?
  - How can we permit Boolean expressions?

$\langle \text{program} \rangle \rightarrow \langle \text{stmt} \rangle$   
 $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{neg} \rangle \mid \langle \text{term} \rangle * \langle \text{neg} \rangle \mid \langle \text{term} \rangle / \langle \text{neg} \rangle$   
 $\langle \text{neg} \rangle \rightarrow \langle \text{var} \rangle \mid -\langle \text{var} \rangle$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z$

$\langle \text{program} \rangle \rightarrow ?$   
 $\langle \text{while} \rangle \rightarrow ?$   
 $\langle \text{stmtList} \rangle \rightarrow ?$   
 $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{bool} \rangle \rightarrow ?$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{neg} \rangle \mid \langle \text{term} \rangle * \langle \text{neg} \rangle \mid \langle \text{term} \rangle / \langle \text{neg} \rangle$   
 $\langle \text{neg} \rangle \rightarrow \langle \text{var} \rangle \mid -\langle \text{var} \rangle$   
 $\langle \text{var} \rangle \rightarrow x \mid y \mid z$

# *Parsing*

- Goals
  - Find syntax errors
  - Produce or trace parse tree
- Types of parsers
  - Top-down: build tree from top down
  - Bottom-up: build tree from bottom up
    - Later!



# *Top-Down Parsers*

- Produces a hierarchical representation of a left-most derivation
  - Traces or builds tree in pre-order
- Recursive Descent Parsing
  - A Top – Down Parser implemented such that a set of recursive methods is used to parse the input.
- Predictive Parser – A recursive descent parser, that may require a lookahead symbol(s) to correctly predict a production at various stages of a derivation.

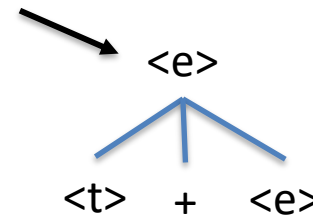
# Predictive Parsers

- Crux: Must choose / **predict** correct RHS when expanding a node in the parse tree.

– Eg: 3+4



- For efficient design best to use next token or **lookahead** only one token



$\langle e \rangle \Rightarrow \langle t \rangle + \langle e \rangle \mid \langle f \rangle - \langle e \rangle \mid \langle t \rangle$

# *Top Down Parsing*

- Formally, a top down parser must do the following:
  - Given a sentence in the form  $xA\alpha$  , the parser must choose the correct  $A \Rightarrow$  rule .
    - Notation
      - Lowercase: terminal sequence
      - Uppercase: non-terminal
      - Greek: sequence of terminals and non-terminals
    - Its best if this decision can be made based on the first tokens of A's RHSs.
      - This is true if all first terminals of A productions are different.
- Implementation of Top Down Parser
  - Recursive decent
  - Table driven
    - Later!

# *Left Recursion and Top Down Parsing:* *( also - Using BNF as a design template)*

```
/* term
Parses strings in the language generated by the rule:
<term> -> <term> * <factor> | ...
*/
void term() {

/* Parse the first factor */
term();

/* As long as the next token is * or /,
next token and parse the next factor */
if(nextToken == MULT_OP || nextToken == DIV_OP) {
    MultDiv();
    factor();
}
} /* End of function term */
```

# *Recursive Decent Parsing*

```
/* term
Parses strings in the language generated by the rule:
<term> -> <factor> { (* | /) <factor> }
*/
void term() {

    /* Parse the first factor */
    factor();

    /* As long as the next token is * or /,
       next token and parse the next factor */
    while (nextToken == MULT_OP || nextToken == DIV_OP) {
        MultDiv();
        factor();
    }
} /* End of function term */
```

# *Limitations of Recursive Descent Parsers*

- **Left Recursion**

- Modify Grammar rules to remove direct left recursion (or similarly use appropriate EBNF)

For each nonterminal, A,

1. Group the A-rules as  $A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$   
where none of the  $\beta$ 's begins with A

2. Replace the original A-rules with

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \varepsilon$$

NOTE: operator association not explicitly specified in this format (can still be left associative). We will see during implementation

- **Prediction**

- Pairwise disjointedness: assure first terminals are different for each A Production
  - The inability to determine the correct RHS on the basis of one token of lookahead
  - FIRST()
- Left factoring

# *Eliminate Left Recursion Example*

$E \Rightarrow E + T \mid E - T \mid T$

$T \Rightarrow 0 \mid 1 \mid 3 \mid \dots \mid 9$

$E \Rightarrow TE'$

$E' \Rightarrow +TE' \mid -TE' \mid \varepsilon$

$T \Rightarrow 0 \mid 1 \mid 3 \mid \dots \mid 9$

For each nonterminal,  $A$ ,

1. Group the  $A$ -rules as  $A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

where none of the  $\beta$ 's begins with  $A$

2. Replace the original  $A$ -rules with

$A \rightarrow \beta_1A' \mid \beta_2A' \mid \dots \mid \beta_nA'$

$A' \rightarrow \alpha_1A' \mid \alpha_2A' \mid \dots \mid \alpha_mA' \mid \varepsilon$

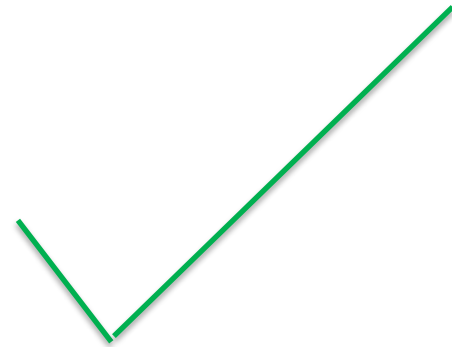
# Prediction: FIRST to determine pairwise disjointedness

- Def:  $FIRST(\alpha) = \{a \mid \alpha \Rightarrow^* a\beta\}$   
(If  $\alpha \Rightarrow^* \epsilon$ ,  $\epsilon$  is in  $FIRST(\alpha)$ )
  - For each nonterminal, A, in the grammar that has more than one RHS, for each pair of rules,  $A \rightarrow \alpha_i$  and  $A \rightarrow \alpha_j$ , it must be true that  $FIRST(\alpha_i) \cap FIRST(\alpha_j) = \phi$

- EG

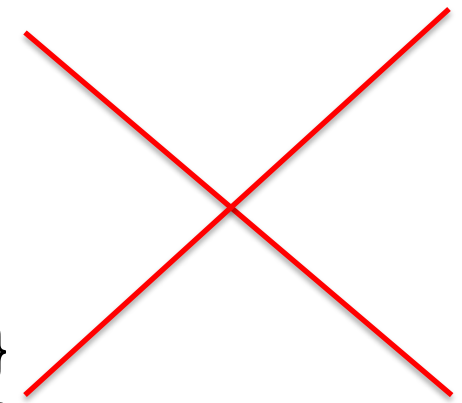
$A \Rightarrow aB \mid bB \mid BA$   
 $B \Rightarrow c \mid d$

$FIRST(aB) = \{a\}$   
 $FIRST(bB) = \{b\}$   
 $FIRST(B) = \{c,d\}$



$A \Rightarrow aB \mid bB \mid BA$   
 $B \Rightarrow c \mid a$

$FIRST(aB) = \{a\}$   
 $FIRST(bB) = \{b\}$   
 $FIRST(B) = \{c,a\}$





## *Left Factoring for Prediction*

- Sometimes we can resolve a nondisjointed FIRST through left factoring by factoring out similar FIRSTs

$\langle \text{variable} \rangle \rightarrow \text{identifier} \mid \text{identifier} [\langle \text{expression} \rangle]$

$\langle \text{variable} \rangle \rightarrow \text{identifier} \langle \text{new} \rangle$

$\langle \text{new} \rangle \rightarrow \varepsilon \mid [\langle \text{expression} \rangle]$

# *EBNF notation*

- Extended BNF
  - Uses curly braces to indicate repetition (in place of explicit recursion)
  - Uses ( | ) to indicate logical OR applied to individual symbols rather than entire RHSs

- Example:

BNF

$$\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$$

EBNF

$$\langle expr \rangle \rightarrow \langle term \rangle \{ (+ \mid -) \langle term \rangle \}$$

# Regular Languages

- Both context free languages and regular languages are useful in programming languages (Chomsky)
- We have seen that BNF is a great way to formalize context free grammars (to define a context free language)

$$\begin{aligned} \langle \text{stmt} \rangle &\rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle \\ \langle \text{expr} \rangle &\rightarrow \langle \text{term} \rangle \mid \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle \\ \langle \text{term} \rangle &\rightarrow \langle \text{neg} \rangle \mid \langle \text{term} \rangle * \langle \text{neg} \rangle \mid \langle \text{term} \rangle / \langle \text{neg} \rangle \\ \langle \text{neg} \rangle &\rightarrow \langle \text{var} \rangle \mid -\langle \text{var} \rangle \\ \langle \text{var} \rangle &\rightarrow x \mid y \mid z \end{aligned}$$

- BNF defines sequences of tokens (sentences) for a languages.
- Regular languages are often used to model the tokens of a language

# *Use of Language Models for PLs*

- In general: RLs are used to model the tokens
  - Each token class is a regular language. *The input alphabet is the set of all characters.*
  - The job of the tokenizer is to “recognize” each token class
    - In a left to right scan of the input, the tokenizer can recognize the beginning and end of each token (given the rules: fsm of each token class). Thus creating a token list.
- In general: CFLs are used to model the programming language (int terms of the tokens / lexemes)
  - Each programming language is a CFL. *The input alphabet is the set of all tokens.*
  - The job of the parser is to “recognize” and produce a parse tree.
    - In a left to right scan of the input (token sequence), the parser will determine if the token sequence conforms with the rules of the grammar, thus creating a parse tree.

*Usage of RLs and CFLs in PLs (revisited)*  
*Use RLs to define rules for tokens. Use CFL for PL rules.*

- Could you use BNF to define the set of all possible variables or ints? YES. But instead we will use regular languages and regular expressions. (More on this later)

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\langle term \rangle \rightarrow \langle neg \rangle \mid \langle term \rangle * \langle neg \rangle \mid \langle term \rangle / \langle neg \rangle$   
 $\langle neg \rangle \rightarrow \langle var \rangle \mid -\langle var \rangle$   
 $\langle var \rangle \rightarrow$  see regular expression for variables.  
 $\langle int \rangle \rightarrow$  see regular expression for integers.

Groups of abstractions: compositions of terminals and nonterminals. We use CFGs to model these constructs.

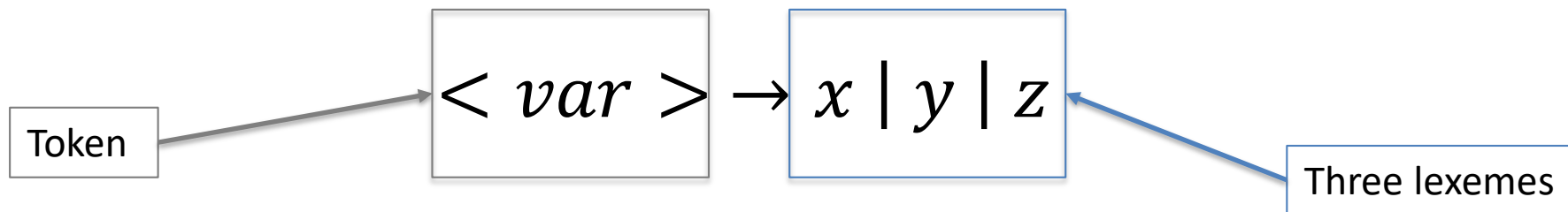
It is intuitive to define our language in terms of its smallest syntactic unit (the most integral units that have meaning). However – these “lexemes” consist of parts (chars) themselves. Solution: Use 1 model for the (smallest syntactic units) lexemes and 1 model for the language.

Groups of tokens (sets of terminals)

We will use regular grammars to model these constructs

# Regular Languages

- Remember: Tokens are categories of lexemes and lexemes are terminals (the smallest syntactic unit)
  - If the number of lexemes per token category are small, then we need only list them out. That is, there is no need to use a Regular Language to model all of the lexemes in the Token Class (Regular Language)
- In our previous example, there were only 3 possible variable tokens.



- If a token class has many possible lexemes, then a formal rule system can be used to help define all possible lexemes (rather than list them all out!). For example, try listing out all possible variable names in C++.

# *Use Regular Languages to define rules for tokens.*

- Could you use BNF to define the set of all possible variables or ints? YES. But instead we will use regular languages and regular expressions. (More on this later)

$\langle stmt \rangle \rightarrow \langle var \rangle = \langle expr \rangle$   
 $\langle expr \rangle \rightarrow \langle term \rangle \mid \langle expr \rangle + \langle term \rangle \mid \langle expr \rangle - \langle term \rangle$   
 $\langle term \rangle \rightarrow \langle neg \rangle \mid \langle term \rangle * \langle neg \rangle \mid \langle term \rangle / \langle neg \rangle$   
 $\langle neg \rangle \rightarrow \langle var \rangle \mid -\langle var \rangle$   
 $\langle var \rangle \rightarrow$  see regular expression for variables.  
 $\langle int \rangle \rightarrow$  see regular expression for integers.

Groups of abstractions: compositions of terminals and nonterminals. We use CFGs to model these constructs.

Groups of tokens  
(sets of terminals)

We will use regular grammars to model these constructs

# *Regular Languages*

- Regular Languages used to model Tokens
- Regular Expression is a generator for a regular language (it defines a regular language)
  - Consists of characters and regular operations
  - Just as BNF is a generator and define a context free language.
- Regular Operations:
  - Concatenation
  - Repetition
    - Repeat 0 or more times: \*
    - Repeat 1 or more times: +
  - Selection: |
  - Other common symbols used
    - Optional: ?
    - Any character: .
    - Short hand for 1 element in a set: [ ]



# *Appendix*



*GEORGETOWN UNIVERSITY*