

# **Principles of Programming Languages**

**Topic: Formal Languages, part A**

# Review

- **See Sakai for**
  - Syllabus
  - Rules and Procedures for Exams
  - Exam Schedule
- **This class will teach you**
  - some new ways of thinking about programs (new *paradigms*)
  - some common principles underlying most languages

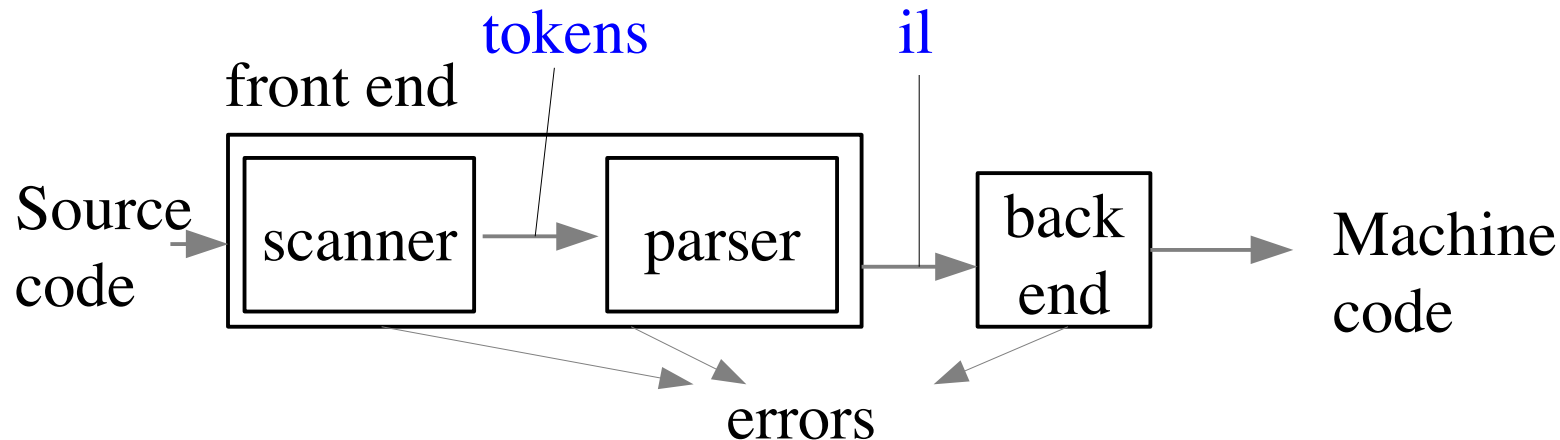
# Review

- A programming language should:
  - be as easy as possible for people to read, write, and learn
  - be as easy as possible for a compiler to compile into efficient machine code or an interpreter to execute
  - be as powerful as possible for the task at hand:  
scale, novelty, paradigm, ...

**“The best chainsaw is no good at cutting paper.  
The best scissors are no good at cutting logs.”**

# Traditional Two-Pass Compiler

- Pass: read and write entire program



**il = intermediate language**

# Scanner

- **Maps characters → tokens**
  - **Tokens: basic unit of syntax**
  - **E.g., `x = x + y;` becomes**  
**`<id, x> <operator, assign> <id, x> <operator + > <id, y>`**
  - **Typical token types:**  
**number, id, operator (e.g., +), keyword (e.g., do, else)**

# Parser

- **Parse:** determine the grammatical structure of a sequence of tokens
- **Grammar:** set of rule like  
assignment  $\rightarrow$  variable '=' expression ';'
- **Analogy with grammars of human languages.**
  - e.g., English: Sentence  $\rightarrow$  Subject Verb Object

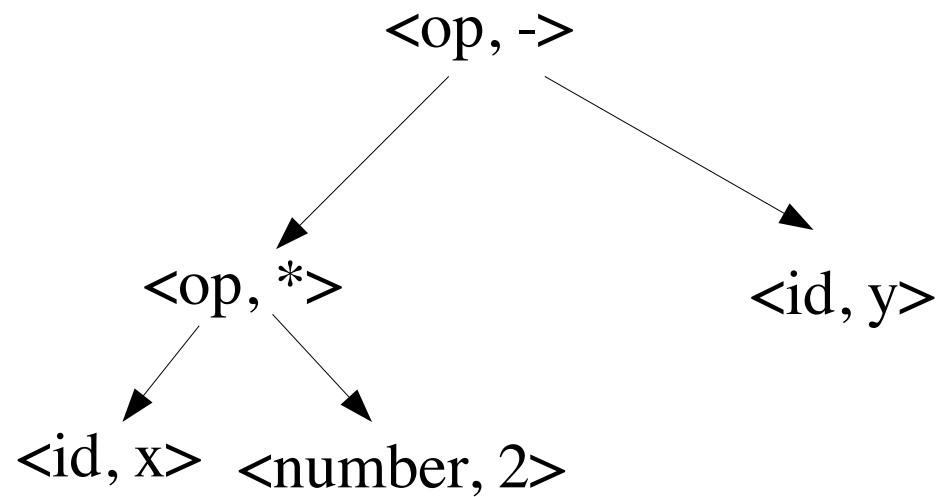
The dog bit Bob

Bob bit the dog     $\leftarrow$  different meaning

Dog Bob the bit     $\leftarrow$  meaningless

# Abstract Syntax Tree

For:  $x * 2 - y$



# Defining a Language

- To define a computer language we need to say
  - What are the parts of the **text** of the program and how are they related: *syntax*
  - What are the parts of the **behavior** specified by the program and how are they related: *semantics*
- Syntactic structure should mirror semantic structure



# Text and Behavior

- Text of the program

**a = b + c ;**

- Behavior of the program

**calculate:**

**add the value of variable b  
to the value of variable c**

**and store the result in variable a**

# Parts

- **Text of the program**

**a** **=** **b** **+** **c** **;**

- **Behavior of the program**

**calculate:**

**add**

**the value of variable b**

**to the value of variable c**

**and store the result in variable a**

# Defining a language

- **Semantics:**
  - several formal approaches but in practice we just use English to explain the meaning
- **Syntax:**
  - tokens defined by a Regular Expression or Finite State Automaton
  - Larger scale structure defined by a Context Free Grammar

# Formal Languages

**For now, we will use the following terms in a special technical sense:**

- **Alphabet:** a set of symbols  
 $\{a, b\}$  or  $\{\text{if, while, for, ...}\}$
- **String:** a sequence of symbols from the alphabet  
abbab
- **Language:** a set of strings
  - $\{ab, abab, ababab, ...\}$
  - Each string: finite length
  - Set: possibly infinite number of strings

# Specifying a Language

- Languages can be infinite sets so we can't just list all strings in the language, must describe the set
  - “any string of alternating ‘a’ and ‘b’ that starts with ‘a’ and ends with ‘b’ ”
- Formalisms for specifying a language
  - – Grammars
  - Regular Expressions
  - Automata

# Grammar

- A grammar is a formalism for specifying a language
  - A set of *Terminal Symbols*, the alphabet of the language  
**a, b**
  - A set of *Non-terminal Symbols*, the grammatical categories of the language, one of which is the “start” symbol  
**MultiAB, AB**
  - A set of *Rules* of the form  
*Non-terminal  $\Rightarrow$  sequence of Symbols*  
**MultiAB  $\Rightarrow$  MultiAB AB**  
**MultiAB  $\Rightarrow$  AB**  
**AB  $\Rightarrow$  a b**

# Grammar Rules

- Another term for Rule is *Production*

- A rule is of the form

*Left-side*  $\Rightarrow$  *Right-side*

where *Left-side* is a single Non-terminal,

and *Right-side* is a sequence of Terminals and/or Non-terminals

**Assignment  $\Rightarrow$  Variable '=' Expression ';'**

# Grammar

- Usually we just give the Rules, since you can see from them what the Terminals and Non-terminals are.

Grammar GAB:

**MultiAB  $\Rightarrow$  MultiAB AB**

**MultiAB  $\Rightarrow$  AB**

**AB  $\Rightarrow$  a b**

start symbol



# Grammar for (Part of) English

$\langle S \rangle \Rightarrow \langle NP \rangle \langle VP \rangle$

$\langle NP \rangle \Rightarrow \langle Name \rangle \mid \langle Det \rangle \langle Noun \rangle$

$\langle VP \rangle \Rightarrow \langle Verb \rangle \mid \langle Verb \rangle \langle NP \rangle$

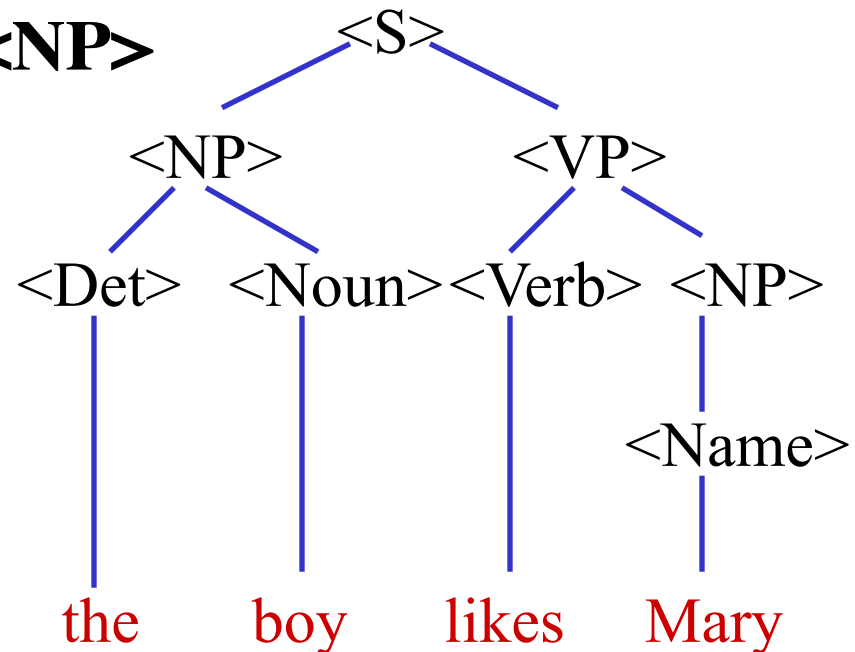
$\langle Name \rangle \Rightarrow \text{John} \mid \text{Mary}$

$\langle Det \rangle \Rightarrow \text{a} \mid \text{the}$

$\langle Det \rangle \Rightarrow \text{some} \mid \text{every}$

$\langle Noun \rangle \Rightarrow \text{boy} \mid \text{girl}$

$\langle Verb \rangle \Rightarrow \text{runs} \mid \text{likes}$



# **The Language of a Grammar**

**Initialize a list of symbols to be just the Start Symbol**

**Repeatedly:**

- Find the first Non-terminal in the list**
- Find a rule whose Left-hand side is this non-Terminal**
- Replace the Non-terminal with the Right-hand side of that rule**

**Until the list contains only Terminals**

# The Language of a Grammar

**Grammar GAB:**

**MultiAB  $\Rightarrow$  MultiAB AB**

**MultiAB  $\Rightarrow$  AB**

**AB  $\Rightarrow$  a b**

**Derivation of a b a b from**

**GAB**

**MultiAB**

**MultiAB AB**

**AB AB**

**a b AB**

**a b a b**

# The Language of a Grammar

- Proof that  $a b a b a b$  is in  $L(G_{AB})$ :

**MultiAB**

**MultiAB AB**

**MultiAB AB AB**

**AB AB AB**

**a b AB AB**

**a b a b AB**

**a b a b a b**

# Parse Trees

- The order in which Non-terminals are replaced does not matter

**MultiAB**

**MultiAB AB**

**AB AB**

**a b AB**

**a b a b**

**MultiAB**

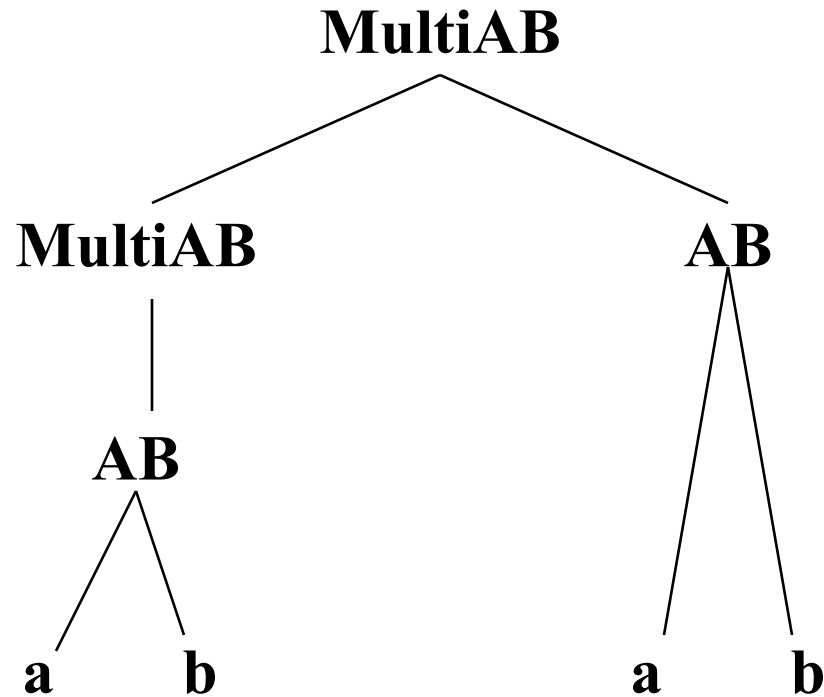
**MultiAB AB**

**MultiAB a b**

**AB a b**

**a b a b**

# Parse Trees



**Each internal node is a Non-terminal; its children make up the right-hand side of one of the productions for that Non-terminal.**

# Derivation in a Grammar

Question: Write a leftmost derivation of **aabb** in G2?

G2:

$$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \mid \langle A \rangle \langle B \rangle$$
$$\langle A \rangle \Rightarrow \mathbf{a} \mid \mathbf{a} \langle A \rangle$$
$$\langle B \rangle \Rightarrow \mathbf{b} \mid \langle B \rangle \mathbf{b}$$

Answer:

$$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \langle B \rangle$$
$$\Rightarrow \mathbf{a} \langle A \rangle \langle B \rangle$$
$$\Rightarrow \mathbf{aa} \langle B \rangle$$
$$\Rightarrow \mathbf{aa} \langle B \rangle \mathbf{b}$$
$$\Rightarrow \mathbf{aa} \mathbf{bb}$$

# Derivation in a Grammar

Question: Write a  
leftmost derivation of  
**baa** in G2?

G2:

$$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \mid \langle A \rangle \langle B \rangle$$
$$\langle A \rangle \Rightarrow \mathbf{a} \mid \mathbf{a} \langle A \rangle$$
$$\langle B \rangle \Rightarrow \mathbf{b} \mid \langle B \rangle \mathbf{b}$$



# Derivation in a Grammar

**Question:** Write a  
leftmost derivation of  
**baa** in G2?

**Answer:**  
**Impossible!**  
No way to get b left of a

**G2:**

$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \mid \langle A \rangle \langle B \rangle$

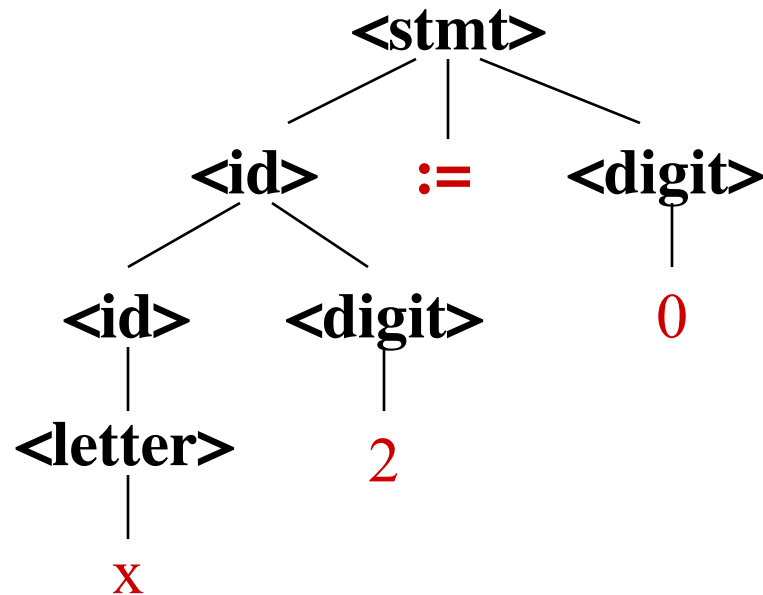
$\langle A \rangle \Rightarrow \mathbf{a} \mid \mathbf{a} \langle A \rangle$

$\langle B \rangle \Rightarrow \mathbf{b} \mid \langle B \rangle \mathbf{b}$

# Derivation vs Parse

- **Derivation:**  $\langle \text{start} \rangle \Rightarrow \dots \Rightarrow \text{list of Terminals}$
- **Parse:**  $\text{list of Terminals} \Rightarrow \dots \Rightarrow \langle \text{start} \rangle$

# Parse Trees



**Each internal node is a nonterminal; its children make up the right-hand side of one of the productions for that nonterminal.**

# Parse Trees

Draw a parse tree for  
**aabb** in G2?

G2:

$$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \mid \langle A \rangle \langle B \rangle$$
$$\langle A \rangle \Rightarrow \mathbf{a} \mid \mathbf{a} \langle A \rangle$$
$$\langle B \rangle \Rightarrow \mathbf{b} \mid \langle B \rangle \mathbf{b}$$

# Parse Trees

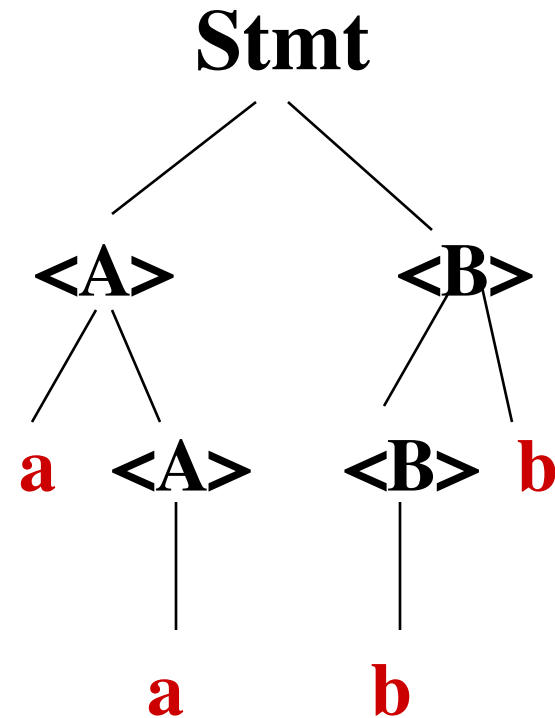
Draw a parse tree for  
**aabb** in G2?

G2:

$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \mid \langle A \rangle \langle B \rangle$

$\langle A \rangle \Rightarrow a \mid a \langle A \rangle$

$\langle B \rangle \Rightarrow b \mid \langle B \rangle b$



# Grammars are not Unique

- Consider a grammar G:

$\langle \text{stmt} \rangle ::= \langle \text{ident} \rangle := \langle \text{digit} \rangle$   
 $\langle \text{ident} \rangle ::= \langle \text{letter} \rangle$   
 $\langle \text{ident} \rangle ::= \langle \text{ident} \rangle \langle \text{letter} \rangle$   
 $\langle \text{ident} \rangle ::= \langle \text{ident} \rangle \langle \text{digit} \rangle$   
 $\langle \text{letter} \rangle ::= a | b | c | \dots | x | y | z$   
 $\langle \text{digit} \rangle ::= 0 | 1 | \dots | 8 | 9$

- Consider a grammar G':

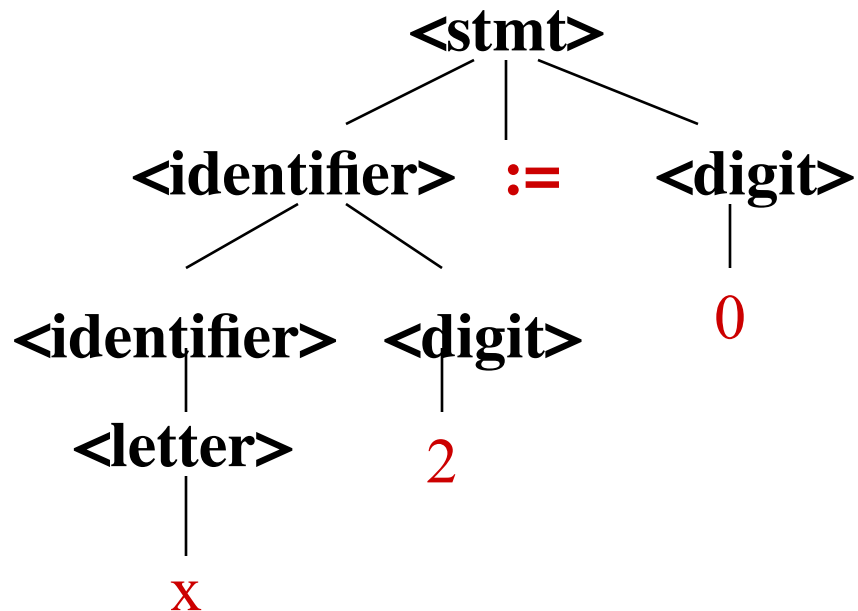
$\langle \text{stmt} \rangle ::= \langle \text{ident} \rangle := \langle \text{digit} \rangle$   
 $\langle \text{ident} \rangle ::= \langle \text{letter} \rangle$   
 $\langle \text{ident} \rangle ::= \langle \text{ident} \rangle \langle \text{char} \rangle$   
 $\langle \text{char} \rangle ::= \langle \text{letter} \rangle | \langle \text{digit} \rangle$   
 $\langle \text{letter} \rangle ::= a | b | c | \dots | x | y | z$   
 $\langle \text{digit} \rangle ::= 0 | 1 | \dots | 8 | 9$

The grammar G' generates the same language as G, but it has different parse trees.

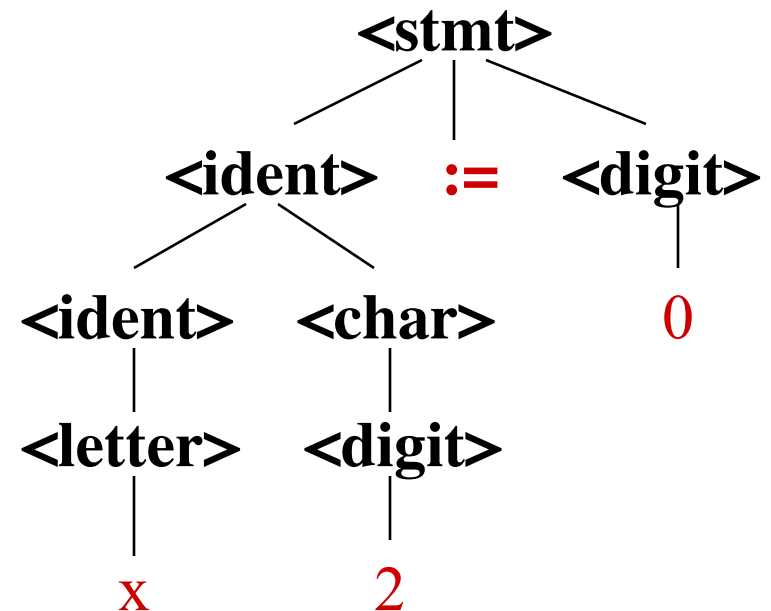
# G vs G'

- **G**  
...  
**<identifier> ::= <letter> | <identifier> <letter> | <identifier><digit>**  
...
- **G': ...**  
**<ident> ::= <letter> | <ident> <char>**  
**<char> ::= <letter> | <digit>**  
...

# Grammars are not Unique



Parse Tree for G



Parse Tree for G'

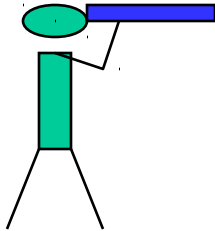


# Ambiguity within a grammar

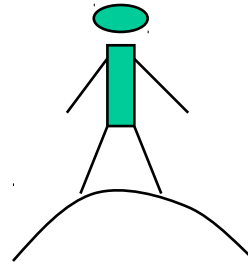
- John saw the man on the hill with the telescope

- Who had the telescope?

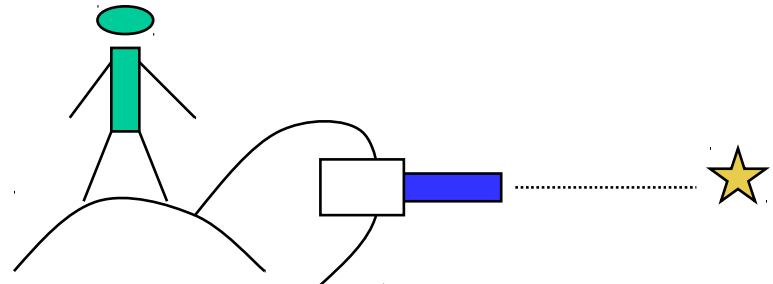
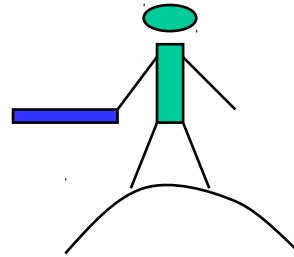
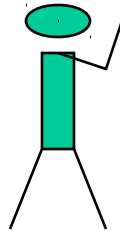
- John



- The man



- The hill



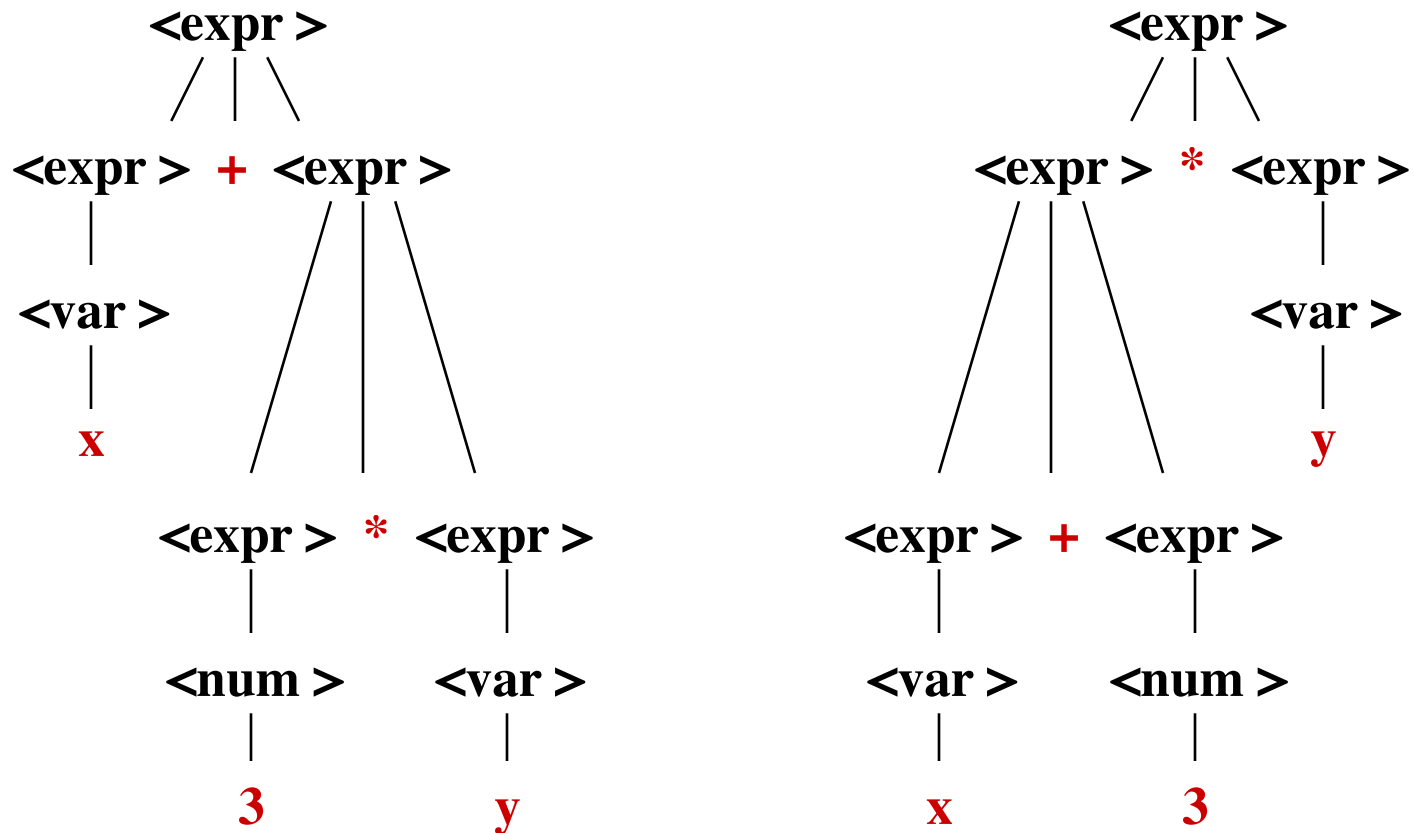
# Arithmetic Expressions

Here is a grammar for arithmetic expressions:

$$\begin{aligned} \langle \text{expr} \rangle ::= & \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid \\ & \langle \text{expr} \rangle * \langle \text{expr} \rangle \mid \langle \text{expr} \rangle / \langle \text{expr} \rangle \mid \\ & \langle \text{var} \rangle \mid \langle \text{num} \rangle \end{aligned}$$
$$\langle \text{var} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$$
$$\langle \text{num} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

Using this grammar, how would we parse:  $x + 3 * y$  ?

# Two Parse Trees



# Two Parse Trees

- in a PL we want to base meaning on parse so
- ambiguous parse -> ambiguous meaning -> 

# Precedence

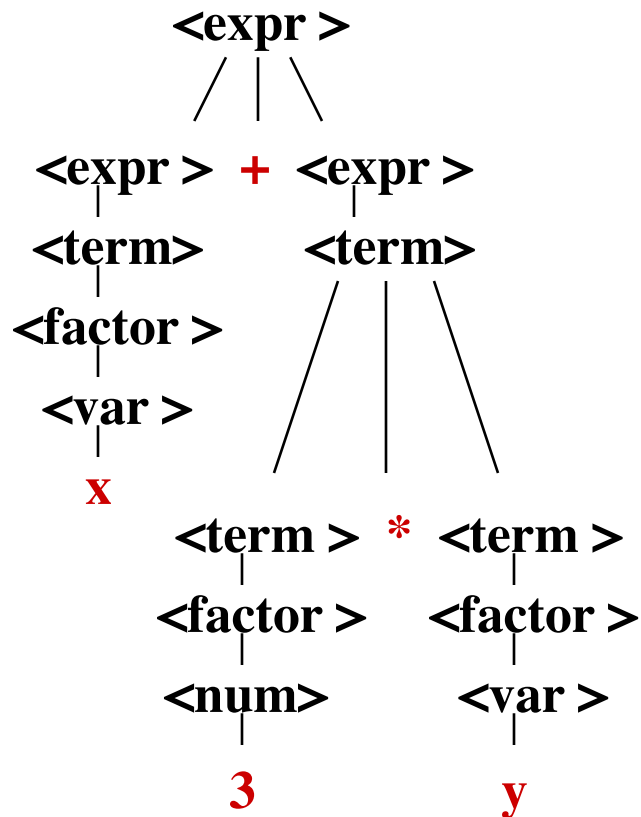
Modify the grammar to add precedence:

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid \langle \text{term} \rangle$   
 $\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{term} \rangle \mid \langle \text{term} \rangle / \langle \text{term} \rangle \mid \langle \text{factor} \rangle$   
 $\langle \text{factor} \rangle ::= \langle \text{var} \rangle \mid \langle \text{num} \rangle \mid ( \langle \text{expr} \rangle )$   
 $\langle \text{var} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$   
 $\langle \text{num} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

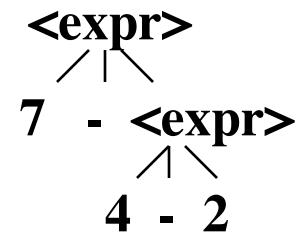
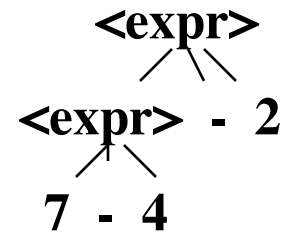
Using this grammar, how would we parse:  $x + 3 * y$  ?

Using this grammar, how would we parse:  $7 - 4 - 2$  ?

# Only One Parse Tree



But there are two parse trees for the second example:



# Associativity

Modify the grammar to add associativity:

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle \mid$   
 $\langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{term} \rangle / \langle \text{factor} \rangle \mid$   
 $\langle \text{factor} \rangle$

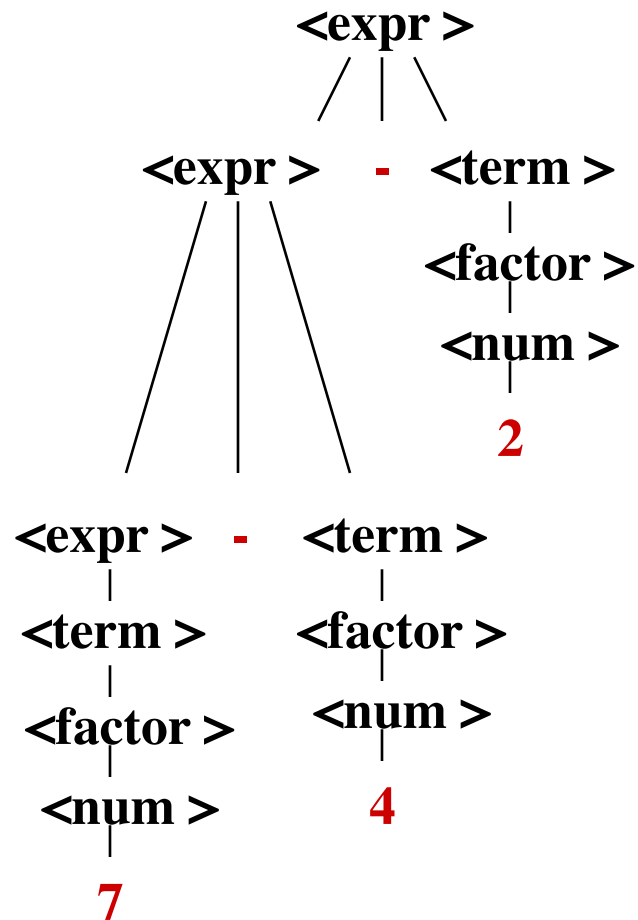
$\langle \text{factor} \rangle ::= \langle \text{var} \rangle \mid \langle \text{num} \rangle \mid ( \langle \text{expr} \rangle )$

$\langle \text{var} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$

$\langle \text{num} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Using this grammar, how would we parse: 7 - 4 - 2 ?

# Only One Parse Tree



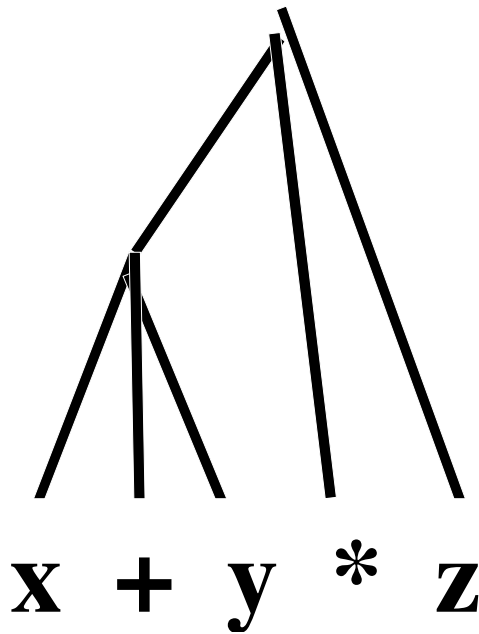


# Solution

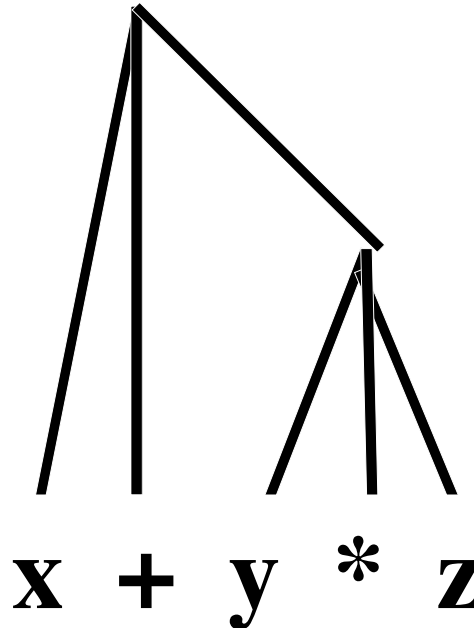
- **Solution: encode precedence & associativity in grammar**
  - non-terminal for each level of precedence
    - + - <term>
    - \* / <factor>
  - for each level of precedence:
    - $\langle nt1 \rangle ::= \langle nt2 \rangle$  (nt2 higher precedence)
    - $\langle nt1 \rangle ::= \langle nt1 \rangle + \langle nt2 \rangle$  (left associative)

# How it works

- Higher in the tree means lower precedence



$(x+y) * z$



$x + (y * z)$

# How it works: Precedence

$\text{NT1} \Rightarrow \text{NT2} \mid \text{NT1} + \text{NT2}$  (NT1 is Start Symbol)

$\text{NT2} \Rightarrow \text{NT3} \mid \text{NT2} * \text{NT3}$

$\text{NT3} \Rightarrow \text{Number}$

To get  $1+2*3$ :

Can only get  $+$  from an NT1 so start

NT1

NT1 + NT2

# How it works: Precedence

$\text{NT1} \Rightarrow \text{NT2} \mid \text{NT1} + \text{NT2}$  (NT1 is Start Symbol)

$\text{NT2} \Rightarrow \text{NT3} \mid \text{NT2} * \text{NT3}$

$\text{NT3} \Rightarrow \text{Number}$   $1+2*3$

$\text{NT1} + \text{NT2}$  \* is right of +, so must come from here



$\text{NT1} + \text{NT2} * \text{NT3}$

# How it works: associativity

$NT1 \Rightarrow NT2 \mid NT1 + NT2$  (NT1 is Start Symbol)

$NT2 \Rightarrow NT3 \mid NT2 * NT3$

$NT3 \Rightarrow \text{Number}$   $1+2+3$

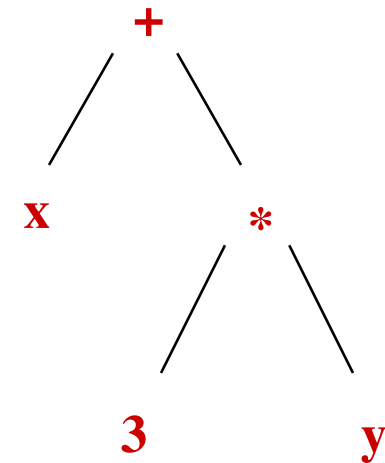
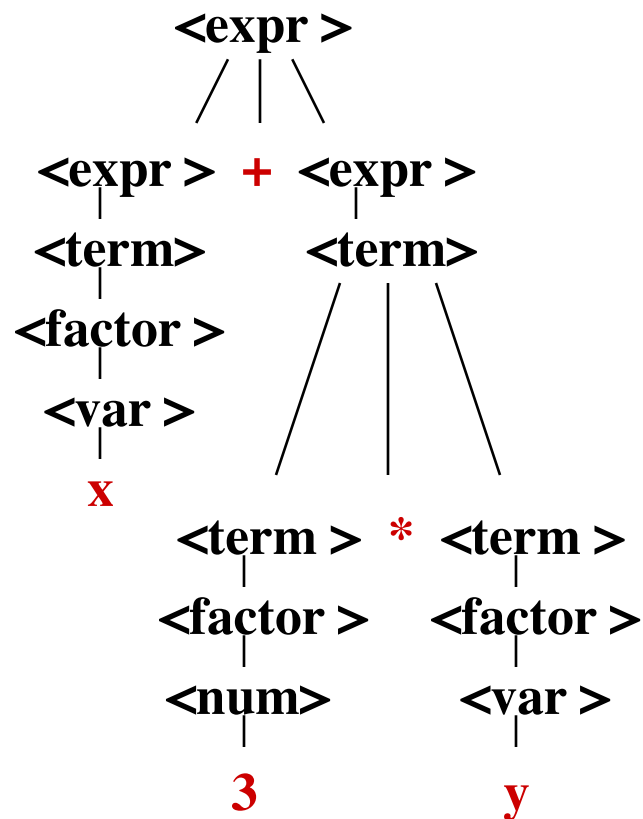
Can only get + from an NT1

NT1

NT1 + NT2

NT1 + NT2 + NT2

# Concrete vs. Abstract Syntax



Abstract Syntax

# Types of Grammars

- **Grammars can be classified into types**
  - **Type is based on form of rules**
  - **Different types -> parsing is a harder / easier computation**

# Types of Grammars

- **Context Free Grammars:**
  - Every rule has a single nonterminal on the left-hand side:  
 $\langle A \rangle \Rightarrow \dots$
  - Disallowed:  $\langle X \rangle \langle A \rangle \Rightarrow \langle X \rangle a$
- **Regular Grammars:**
  - Rules all take the forms:  
 $\langle A \rangle \Rightarrow c$  or  $\langle A \rangle \Rightarrow \langle B \rangle c$  (*left-linear*)
  - Or rules all take the forms:  
 $\langle A \rangle \Rightarrow c$  or  $\langle A \rangle \Rightarrow c \langle B \rangle$  (*right-linear*)
  - Disallowed:  $S \Rightarrow a S b$
  - Cannot generate the language  $\{ a^n b^n \mid n = 1, 2, 3, \dots \}$



# Types of Grammars

- **Context Free Grammars (CFGs)** are used to specify the overall structure of a programming language:
  - if/then/else, ...
  - brackets: ( ), { }, begin/end, ...
- **Regular Grammars (RGs)** are used to specify the structure of tokens:
  - identifiers, numbers, keywords, ...
- **RGs are a subset of CFGs**

# Extended BNF (EBNF)

**A language for defining the grammar of a language.**

**Write nonterminals as usual. (Variant: Write them with initial capital letters, or using a different font.)**

**Use additional *metasymbols*, as shortcuts:**

- **{...}** means repeat the enclosed text zero or more times
- **[...]** means the enclosed text is optional
- **(...)** is used for grouping, usually with the alternation symbol, e.g., **(... | ...)**.

**If { }, [ ], or ( ) are used as terminal symbols in the language being defined, then they must be quoted. (Variant: They must be underlined.)**

# Extended BNF (EBNF)

## Examples:

$\langle \text{expr} \rangle ::= \langle \text{term} \rangle \{ ( + | - ) \langle \text{term} \rangle \}$

$\langle \text{term} \rangle ::= \langle \text{factor} \rangle \{ ( * | / ) \langle \text{factor} \rangle \}$

$\langle \text{factor} \rangle ::= \langle \text{var} \rangle \mid \langle \text{num} \rangle \mid ' ( ' \langle \text{expr} \rangle ' ) '$

$\langle \text{if-stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle [ \text{else } \langle \text{stmt} \rangle ]$

$\langle \text{identifier} \rangle ::= \langle \text{letter} \rangle \{ ( \langle \text{letter} \rangle \mid \langle \text{digit} \rangle ) \}$

# Formal Language Theory

- **Offers a way to describe computation problems formulated as language recognition problems**
  - Enables proofs of relative difficulty of certain computational problems
- **Provides a mechanism to aid description of programming language constructs**
  - Regular expressions ~ PL tokens (e.g. real numbers, keywords)
  - Context-free grammars ~ PL statements

# Specifying a Language

- Formalisms for specifying a language
  - Grammars
  -  – Regular Expressions
  - Automata

# Regular Expressions

- **Formalism for describing simple PL constructs**
  - reserved words
  - identifiers
  - numbers
- **Simplest sort of grammatical structure**
- **Defined recursively**
- **Actually are expressions on languages, ie on sets of strings**

# Regular Expressions

	<u>RE Notation</u>	<u>Language</u>
	an empty RE	$\{ \}$
symbol $a$	$a$	$\{a\}$
null symbol	$\epsilon$	$\{\epsilon\}$
$R, S$ regular exprs	$R \mid S$	$L_R \cup L_S$
	<i><math>ab</math> (alternation)</i>	<i><math>\{a,b\}</math></i>
$R, S$ regular exprs	$RS$	$L_R L_S$
	<i><math>ab</math> (concatenation)</i>	<i><math>\{ab\}</math></i>

# Regular Expressions

## RE Notation

## Language

**a**

**{a}**

**b**

**{b}**

**ab**

**{ab}**

**a | b**

**{a, b}**

**ab | ac**

**{ab, ac}**

**(a | b)(c | d)**

**{ac, ad, bc, bd}**

**(abc | ε) d**

**{abcd, d}**



# Regular Expressions

## RE Notation   Language

<b>R regular expr</b>	<b><math>R^*</math></b>	$\{\epsilon\} \cup L_R \cup L_R L_R \cup L_R L_R L_R \dots$
	<i><math>a^*</math></i>	<i><math>\{\epsilon, a, aa, aaa, \dots\}</math></i>
<b>R regular expr</b>	<b><math>R^+</math></b>	$L_R \cup L_R L_R \cup L_R L_R L_R \dots$
	<i><math>a^+</math></i>	<i><math>\{a, aa, aaa, \dots\}</math></i>

**Note:  $\epsilon a = a \epsilon = a$**

**Precedence is + \* concatenation |**

**high                      to                      low**

**(all are left associative operators)**

# Regular Expressions

## RE Notation

**$a^*$**

**$ab^*$**

**$(ab)^*$**

**$(a \mid b)^*$**

**$a^+$**

**$ab^+$**

## Language

**$\{\epsilon, a, aa, aaa, \dots\}$**

**$\{a, ab, abb, abbb, \dots\}$**

**$\{\epsilon, ab, abab, ababab, \dots\}$**

**$\{\epsilon, a, b, aa, ab, ba, bb, \dots\}$**

**$\{a, aa, aaa, \dots\}$**

**$\{ab, abb, abbb, \dots\}$**