

Principles of Programming Languages

COMP3031: Syntax and Grammars

Prof. Dekai Wu

Department of Computer Science and Engineering
The Hong Kong University of Science and Technology
Hong Kong, China



Fall 2012

Part I

Language Description

“Able was I ere I saw Elba.” — about Napoléon

How do you know that this is English, and not French or Chinese?

A language has 2 parts:

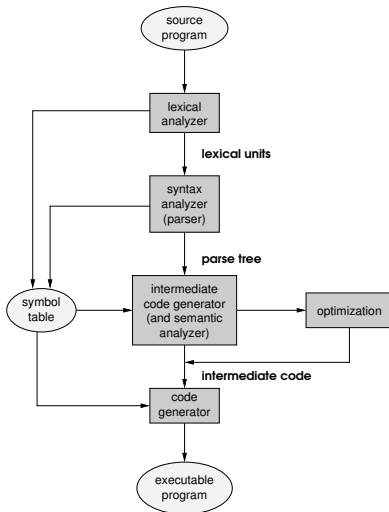
① Syntax

- **lexical syntax**
 - describes how a sequence of *symbols* makes up *tokens* (*lexicon*) of the language
 - checked by a *lexical analyzer*
- **grammar**
 - describes how a sequence of *tokens* makes up a valid *program*.
 - checked by a *parser*

② Semantics

specifies the *meaning* of a program

Compilation



Example 1: English Language

A word = some combination of the 26 letters, a,b,c, ...,z.

One form of a sentence = Subject + Verb + Object.

e.g. The student wrote a great program.

Example 2: Date Format

A date like 06/04/2010 may be written in the general format:

$$D D / D D / D D D D$$

where $D = 0,1,2,3,4,5,6,7,8,9$

But, does 03/09/1998 mean Sept 3rd, or March 9th?

Example 3: Real Numbers (Simplified)

Examples of reals: 0.45 12.3 .98

Examples of non-reals: $2+4i$ 1a2b $8 <$

Informal rules:

- In general, a real number has three parts:
 - an integer part (I)
 - a dot “.” symbol ($.$)
 - a fraction part (F)
- valid forms: $I.F$, $.F$
- I and F are strings of digits
- I may be empty but F cannot
- a digit is one of $\{ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 \}$

Expression: Examples

$$a + b$$

$$3 * a + b/c$$

$$\frac{-b + \sqrt{b^2 - 4 * a * c}}{2 * a}$$

$$\frac{a * (1 - R^n)}{1 - R}$$

```
if (x > 10) then
    x /= 10
else
    x *= 2
```

c.f. “While I was coming to school, I saw a car accident.”
The sentence is in the form of: “While E_1, E_2 .”

Expression Notation: Example 4

Goal: Add a to b .

Infix : $a + b$

Prefix : $+ab$

Postfix : $ab+$

Abstract Syntax Tree



Abstract syntax tree is *independent* of notation.

- A **constant** or **variable** is an expression.
- In general, an expression has the form of a **function**:

$$E \triangleq \mathbf{Op} (E_1, E_2, \dots, E_k)$$

where **Op** is the operator, and E_1, E_2, \dots, E_k are the operands.

- An operator with k operands is said to have an **arity** of k ; and **Op** is an k -ary operator.

unary operator : $-x$

binary operator : $x + y$

ternary operator : $(x > y) ? x : y$

Infix, Prefix, Postfix, Mixfix

- **Infix** : $E_1 \text{ Op } E_2$ (must be binary operator!)

$$a + b, a * b, a - b, a / b, a == b, a < b.$$

- **Prefix** : $\text{Op } E_1 E_2 \dots E_k$

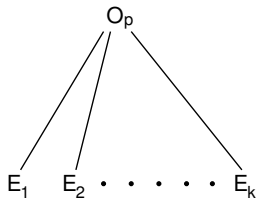
$$+ab, *ab, -ab, /ab, == ab, < ab.$$

- **Postfix** : $E_1 E_2 \dots E_k \text{ Op}$

$$ab+, ab*, ab-, ab/, ab==, ab < .$$

- **Mixfix** : e.g. **if** E_1 **then** E_2 **else** E_3

Abstract Syntax Tree



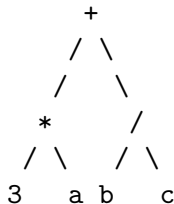
Expression Notation: Example 5

abstract syntax tree

infix : $3 * a + b / c$

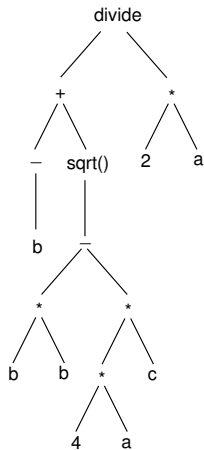
prefix : $+ * 3a / bc$

postfix : $3a * bc / +$



Note: Prefix and postfix notation does not require parentheses.

Expression Notation: Example 6



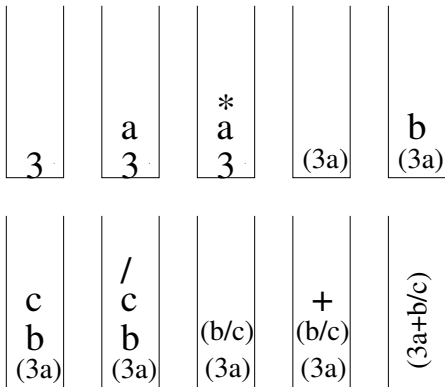
infix : $(-b + \sqrt{b^2 - 4 * a * c}) / (2 * a)$

prefix : $/ + - b \sqrt{ - * b b * * 4 a c * 2 a }$

postfix : $b - b b * 4 a * c * - \sqrt{ + 2 a * / }$

Postfix Evaluation: By a Stack

- **infix** expression: $3 * a + b/c$.
- **postfix** expression: $3a * bc/+$.



Precedence and Associativity in C++

Operator	Description	Associativity
[] . →	array element structure member pointer	LEFT
- ++ -- *	minus increment decrement indirection	RIGHT
* / %	multiply divide mod	LEFT
+ -	add subtract	LEFT
==	logical equal	LEFT
=	assignment	RIGHT

Example: $1/2 + 3 * 4 = (1/2) + (3 * 4)$

because $*$, $/$ has a *higher precedence* over $+$, $-$.

Precedence rules decide which operators run first. In general,

$$x P y Q z = x P (y Q z)$$

if operator Q is at a higher precedence level than operator P .

Associativity: Binary Operators

Example: $1 - 2 + 3 - 4 = ((1 - 2) + 3) - 4$
because $+$, $-$ are *left associative*.

Associativity decides the grouping of operands with operators of the *same* level of precedence.

In general, if **binary** operator P , Q are of the **same** precedence level:

$$x P y Q z = x P (y Q z)$$

if operator P , Q are both **right associative**;

$$x P y Q z = (x P y) Q z$$

if operator P , Q are both **left associative**.

Question : What if $+$ is left while $-$ is right associative?

Associativity: Unary Operators

- Example in C++: $*a++ = *(a++)$
because all unary operators in C++ are right-associative.
- In Pascal, all operators including unary operators are left-associative.
- In general, unary operators in many languages may be considered as non-associative as it is not important to assign an associativity for them, and their usage and semantics will decide their order of computation.

Question : Which of infix/prefix/postfix notation needs precedence or associative rules?

Summary on Syntax

- ✓ Will describe a language by a formal syntax and an informal semantics
- ✓ Syntax = lexical syntax + grammar
- ✓ Expression notation: infix, prefix, postfix, mixfix
- ✓ Abstract syntax tree: independent of notation
- ✓ Precedence and associativity of operators decide the order of applying the operators

Part II

Grammar

Grammar: Motivation

What do the following sentences really mean?

- 路不通行不得在此小便
- “I saw a small kid on the beach with a binocular.”
- What is the final value of x?

```
x = 15
if (x > 20) then
if (x > 30) then
x = 8
else
x = 9
```

- 楊乃武與小白菜

Ambiguity in semantics is often caused by **ambiguous grammar** of the language.

A Formal Description: Example 7

1. $\langle \textit{real-number} \rangle ::= \langle \textit{integer-part} \rangle . \langle \textit{fraction} \rangle$
2. $\langle \textit{integer-part} \rangle ::= \langle \textit{empty} \rangle \mid \langle \textit{digit-sequence} \rangle$
3. $\langle \textit{fraction} \rangle ::= \langle \textit{digit-sequence} \rangle$
4. $\langle \textit{digit-sequence} \rangle ::= \langle \textit{digit} \rangle \mid \langle \textit{digit} \rangle \langle \textit{digit-sequence} \rangle$
5. $\langle \textit{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

This is the **context-free grammar** of real numbers written in the **Backus-Naur Form**.

Context Free Grammar (CFG)

A **context-free grammar** has 4 components:

- 1 **A set of tokens or terminals:**
atomic symbols of the language.

English : a, b, c,, z

Reals : 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, .

- 2 **A set of nonterminals:**
variables denoting language constructs.

English : $\langle \textit{Noun} \rangle$, $\langle \textit{Verb} \rangle$, $\langle \textit{Adjective} \rangle$, . . .

Reals : $\langle \textit{real-number} \rangle$, $\langle \textit{integer-part} \rangle$, $\langle \textit{fraction} \rangle$,
 $\langle \textit{digit-sequence} \rangle$, $\langle \textit{digit} \rangle$

- ③ A set of rules called **productions**:
for generating expressions of the language.

nonterminal ::= a string of terminals and nonterminals

English : $\langle \textit{Sentence} \rangle ::= \langle \textit{Noun} \rangle \langle \textit{Verb} \rangle \langle \textit{Noun} \rangle$

Reals : $\langle \textit{integer-part} \rangle ::= \langle \textit{empty} \rangle | \langle \textit{digit-sequence} \rangle$

Notice that CFGs allow only a **single** non-terminal on the left-hand side of any production rules.

- ④ A nonterminal chosen as the **start symbol**:
represents the main construct of the language.

English : $\langle \textit{Sentence} \rangle$

Reals : $\langle \textit{real-number} \rangle$

The set of strings that can be generated by a CFG makes up a **context-free language**.

Backus-Naur Form (BNF)

One way to write context-free grammar.

- **Terminals** appear as they are.
- **Nonterminals** are enclosed by \langle and \rangle .
e.g.: $\langle \textit{real-number} \rangle$, $\langle \textit{digit} \rangle$.
- The special **empty string** is written as $\langle \textit{empty} \rangle$.
- **Productions** with a common nonterminal may be abbreviated using the special “or” symbol “|”.

e.g. $X ::= W_1, X ::= W_2, \dots, X ::= W_n$

may be abbreviated as $X ::= W_1 | W_2 | \dots | W_n$

Top-Down Parsing: Example 8

- A **parser** checks to see if a given expression or program can be derived from a given grammar.

Check if “.5” is a valid real number by finding from the CFG of Example 6 a **leftmost derivation** of “.5”:

< real-number >

\Rightarrow *< integer-part > . < fraction >* [Production 1]

\Rightarrow *< empty > . < fraction >* [Production 2]

\Rightarrow *. < fraction >* [By definition]

\Rightarrow *. < digit-sequence >* [Production 3]

\Rightarrow *. < digit >* [Production 4]

\Rightarrow *.5* [Production 5]

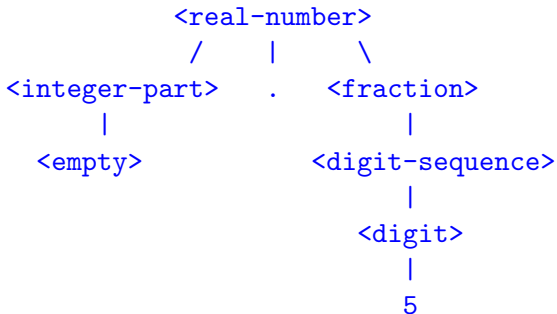
Bottom-Up Parsing: Example 9

Check if “.5” is a valid real number by finding from the CFG of Example 6 a **rightmost derivation** of “.5” in reverse:

.5 = $\langle \text{empty} \rangle . 5$ [By definition]
 $\Rightarrow \langle \textit{integer-part} \rangle . 5$ [Production 2]
 $\Rightarrow \langle \textit{integer-part} \rangle . \langle \textit{digit} \rangle$ [Production 5]
 $\Rightarrow \langle \textit{integer-part} \rangle . \langle \textit{digit-sequence} \rangle$ [Production 4]
 $\Rightarrow \langle \textit{integer-part} \rangle . \langle \textit{fraction} \rangle$ [Production 3]
 $\Rightarrow \langle \textit{real-number} \rangle$ [Production 1]

Parse Tree: Example 10 [Real Numbers]

A parse tree of “.5” generated by the CFG of Example 6.



A **parse tree** shows how a string is generated by a CFG — the **concrete syntax** in a tree representation.

- Root = **start symbol**.
- Leaf nodes = **terminals** or **<empty>**.
- Non-leaf nodes = **nonterminals**
- For any subtree, the **root** is the left-side nonterminal of some production, while its **children**, if read from left to right, make up the right side of the production.
- The **leaf** nodes, read from left to right, make up a string of the language defined by the CFG.

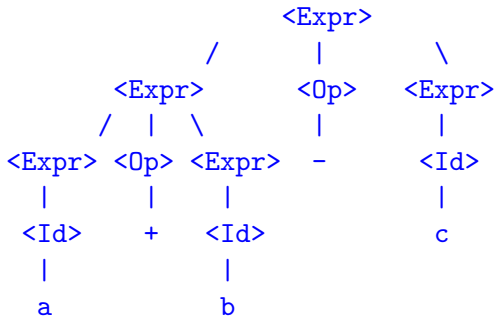
Example 11: CFG/BNF [Expression]

$\langle \text{Expr} \rangle ::= \langle \text{Expr} \rangle \langle \text{Op} \rangle \langle \text{Expr} \rangle$
 $\langle \text{Expr} \rangle ::= (\langle \text{Expr} \rangle)$
 $\langle \text{Expr} \rangle ::= \langle \text{Id} \rangle$
 $\langle \text{Op} \rangle ::= + \mid - \mid * \mid / \mid =$
 $\langle \text{Id} \rangle ::= a \mid b \mid c$

1. Terminals: $a, b, c, +, -, *, /, =, (,)$
2. Nonterminals: $\text{Expr}, \text{Op}, \text{Id}$
3. Start symbol: Expr

Parse Tree : Example 12 [Expression]

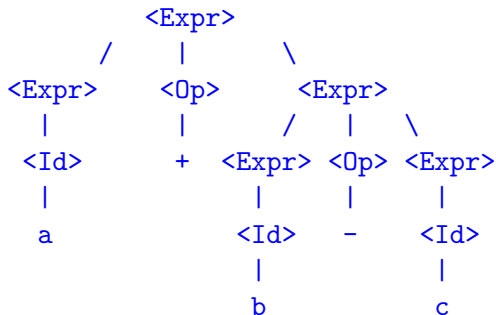
A parse tree of “ $a + b - c$ ” generated by the CFG of Example 10:



Question: What is the difference between a parse tree and an abstract syntax tree?

Ambiguous Grammar: Example 13

A grammar is (syntactically) **ambiguous** if some string in its language is generated by more than one parse tree.



Solution: Rewrite the grammar to make it unambiguous.

Handle Left Associativity: Example 14

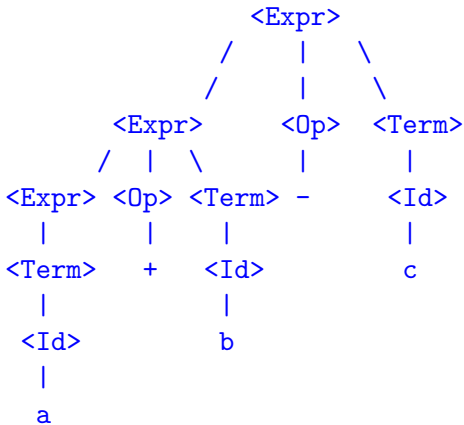
CFG of Example 10 cannot handle “ $a + b - c$ ” correctly.

⇒ Add a left recursive production.

$$\begin{aligned} \langle \text{Expr} \rangle & ::= \langle \text{Expr} \rangle \langle \text{Op} \rangle \langle \text{Term} \rangle \\ \langle \text{Expr} \rangle & ::= \langle \text{Term} \rangle \\ \langle \text{Term} \rangle & ::= (\langle \text{Expr} \rangle) | \langle \text{Id} \rangle \\ \langle \text{Op} \rangle & ::= + | - | * | / | = \\ \langle \text{Id} \rangle & ::= a | b | c \end{aligned}$$

Handle Left Associativity ..

Now there is only one parse tree for “ $a + b - c$ ”:



Handling Right Associativity: Example 15

CFG of Example 10 cannot handle “ $a = b = c$ ” correctly.

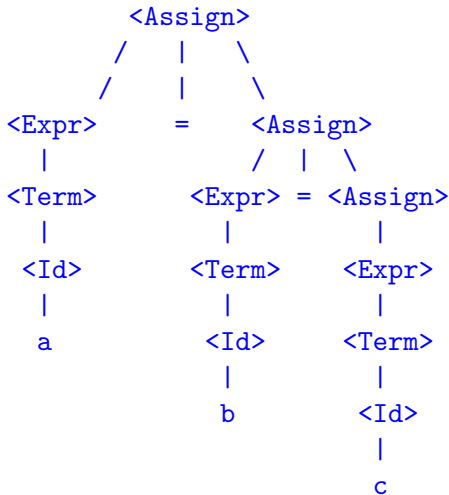
⇒ Add a right recursive production.

$$\begin{aligned} \langle \textit{Assign} \rangle &::= \langle \textit{Expr} \rangle = \langle \textit{Assign} \rangle \\ \langle \textit{Assign} \rangle &::= \langle \textit{Expr} \rangle \\ \langle \textit{Expr} \rangle &::= \langle \textit{Expr} \rangle \langle \textit{Op} \rangle \langle \textit{Term} \rangle \mid \langle \textit{Term} \rangle \\ \langle \textit{Term} \rangle &::= (\langle \textit{Expr} \rangle) \mid \langle \textit{Id} \rangle \\ \langle \textit{Op} \rangle &::= + \mid - \mid * \mid / \\ \langle \textit{Id} \rangle &::= a \mid b \mid c \end{aligned}$$

Question: this grammar will accept strings like “ $a + b = c - d$ ”.
Try to correct it.

Handling Right Associativity ..

Now there is only one parse tree for “ $a = b = c$ ”:



Handling Precedence: Example 16

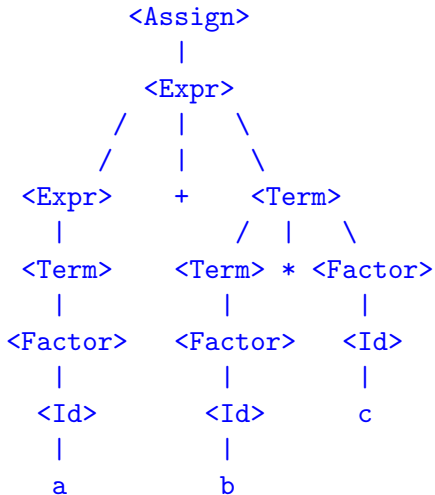
CFG of Example 10 cannot handle “ $a + b * c$ ” correctly.

⇒ Add one nonterminal (plus appropriate productions) for each precedence level.

$$\begin{aligned} \langle \textit{Assign} \rangle & ::= \langle \textit{Expr} \rangle = \langle \textit{Assign} \rangle \mid \langle \textit{Expr} \rangle \\ \langle \textit{Expr} \rangle & ::= \langle \textit{Expr} \rangle + \langle \textit{Term} \rangle \\ \langle \textit{Expr} \rangle & ::= \langle \textit{Expr} \rangle - \langle \textit{Term} \rangle \mid \langle \textit{Term} \rangle \\ \langle \textit{Term} \rangle & ::= \langle \textit{Term} \rangle * \langle \textit{Factor} \rangle \\ \langle \textit{Term} \rangle & ::= \langle \textit{Term} \rangle / \langle \textit{Factor} \rangle \mid \langle \textit{Factor} \rangle \\ \langle \textit{Factor} \rangle & ::= (\langle \textit{Expr} \rangle) \mid \langle \textit{Id} \rangle \\ \langle \textit{Id} \rangle & ::= a \mid b \mid c \end{aligned}$$

Handling Precedence ..

Now there is only one parse tree for “ $a + b * c$ ”:



Tips on Handling Precedence/Associativity

- **left** associativity \Rightarrow **left-recursive** production
- **right** associativity \Rightarrow **right-recursive** production
- n levels of precedence
 - **divide** the operators into n groups
 - write productions for each group of operators
 - start with operators with the **lowest** precedence
- In all cases, introduce **new** non-terminals whenever necessary.
- In general, one needs a new non-terminal for each new group of operators of different associativity and different precedence.

Dangling-Else: Example 17

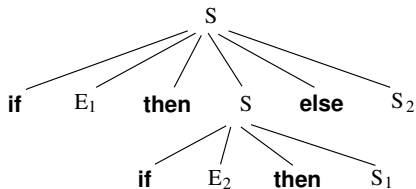
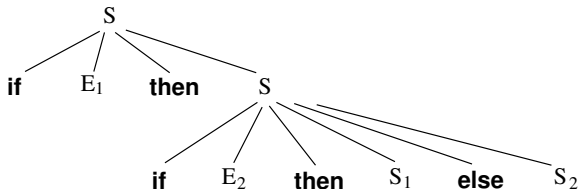
Consider the following grammar:

$$\langle S \rangle ::= \text{if } \langle E \rangle \text{ then } \langle S \rangle$$
$$\langle S \rangle ::= \text{if } \langle E \rangle \text{ then } \langle S \rangle \text{ else } \langle S \rangle$$

- How many parse trees can you find for the statement:

$$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$$

Dangling-Else ..



- Ambiguity is often a property of a grammar, not of a language.

Solution: matching an “else” with the nearest unmatched “if” .
i.e. the first case.

More CFG Examples

1

$$\begin{aligned}\langle S \rangle &::= \langle A \rangle \langle B \rangle \langle C \rangle \\ \langle A \rangle &::= a \langle A \rangle \mid a \\ \langle B \rangle &::= b \langle B \rangle \mid b \\ \langle C \rangle &::= c \langle C \rangle \mid c\end{aligned}$$

2

$$\begin{aligned}\langle S \rangle &::= \langle A \rangle a \langle B \rangle b \\ \langle A \rangle &::= \langle A \rangle b \mid b \\ \langle B \rangle &::= a \langle B \rangle \mid a\end{aligned}$$

3

$$\begin{aligned}\langle \text{stmts} \rangle &::= \langle \text{empty} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle \\ \langle \text{stmt} \rangle &::= \langle \text{id} \rangle := \langle \text{expr} \rangle \\ &\mid \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \\ &\mid \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ &\mid \text{while } \langle \text{expr} \rangle \text{ do } \langle \text{stmt} \rangle \\ &\mid \text{begin } \langle \text{stmts} \rangle \text{ end}\end{aligned}$$

Non-Context Free Grammars: Examples

$$\begin{aligned} \langle S \rangle &::= \langle B \rangle \langle A \rangle \langle C \rangle \mid \langle C \rangle \langle A \rangle \langle B \rangle \\ b \langle A \rangle &::= c \langle A \rangle \langle B \rangle \mid \langle B \rangle \\ c \langle A \rangle &::= b \langle A \rangle \langle C \rangle \mid \langle C \rangle \\ \langle B \rangle &::= b \\ \langle C \rangle &::= c \end{aligned}$$

①

$$\Rightarrow L = \{ (cb)^n, b(cb)^n, (bc)^n, c(bc)^n \}.$$

②

$$L = \{ w cw \mid w \text{ is a string of } a\text{'s or } b\text{'s} \}.$$

This language abstracts the problem of checking that an identifier is declared before its use in a program.

The first w = declaration of the identifier, and the second w = its use in the program.

Summary on Grammar

- ✓ Context-free grammar (CFG) is commonly used to specify most of the syntax of a programming language.
- ✓ However, most programming languages are not CFL!
- ✓ CFG is commonly written in Backus-Naur Form (BNF).
- ✓ $\text{CFG} = (\text{Terminals}, \text{Nonterminals}, \text{Productions}, \text{Start Symbol})$
- ✓ A program is valid if we may construct a parse tree, or a derivation from the grammar.
- ✓ Associativity and precedence of operations are part of the design of a CFG.
- ✓ Avoid ambiguous grammars by rewriting them or imposing parsing rules.

Part III

Regular Grammar, Regular Expression

Regular Grammars are a subset of CFGs in which all productions are in one of the following forms:

1 Right-Regular Grammar

$$\langle A \rangle ::= x$$
$$\langle A \rangle ::= x\langle B \rangle$$

2 Left-Regular Grammar

$$\langle A \rangle ::= x$$
$$\langle A \rangle ::= \langle B \rangle x$$

where A and B are **non-terminals** and x is a string of **terminals**.

RE Example 1: Right-Regular Grammar

$\langle S \rangle ::= a\langle A \rangle$

$\langle S \rangle ::= b\langle B \rangle$

$\langle S \rangle ::= \langle \text{empty} \rangle$

$\langle A \rangle ::= a\langle S \rangle$

$\langle B \rangle ::= bb\langle S \rangle$

What is the **regular language** this RG generates?

Regular Expressions

Regular expressions (RE) are succinct representations of RGs using the following notations.

Sub-Expression	Meaning
x	the single char 'x'
.	any single char except the newline
[abc]	char class consisting of 'a', 'b', or 'c'
[^abc]	any char except 'a', 'b', 'c'
r*	repeat "r" zero or more times
r+	repeat "r" 1 or more times
r?	zero or 1 occurrence of "r"
rs	concatenation of RE "r" and RE "s"
(r)s	"r" is evaluated and concatenated with "s"
r s	RE "r" or RE "s"
\x	escape sequences for white-spaces and special symbols: \b \n \r \t

Precedence of Regular Expression Operators

The following table gives the order of RE operator precedence from the highest precedence to the lowest precedence.

Function	Operator
parenthesis	()
counters	* + ? { }
concatenation	
disjunction	

RE Example 2: Regular Expression Notations

RE	Meaning
abc	the string "abc"
a+b+	$\{a^m b^n : m, n \geq 1\}$
a*b*c	$\{a^m b^n c : m, n \geq 0\}$
a*b*c?	$\{a^m b^n c \text{ or } a^m b^n : m, n \geq 0\}$
xy(abc)+	$\{xy(abc)^n : n \geq 1\}$
xy[abc]	$\{xya, xyb, xyc\}$
xy(a b)	$\{xya, xyb\}$

Questions: What are the following REs?

- foo|bar*
- foo|(bar)*
- (foo|bar)*

RE Example 3: Regular Expressions

- REs are commonly used for **pattern matching** in editors, word processors, commandline interpreters, etc.
- The REs used for **searching texts** in Unix (vi, emacs, perl, grep), Microsoft Word v.6+, and Word Perfect are almost identical.
- Examples:
 - identifiers in C++:
 - real numbers:
 - email addresses:
 - white spaces:
 - all C++ source or include files:

Summary on Regular Grammars

- ✓ There are algorithms to prove if a language is regular.
- ✓ There are algorithms to prove if a language is context-free too.
- ✓ English is not RL, nor CFL.
- ✓ REs are commonly used for text search.
- ✓ Different applications may extend the standard RE notations.