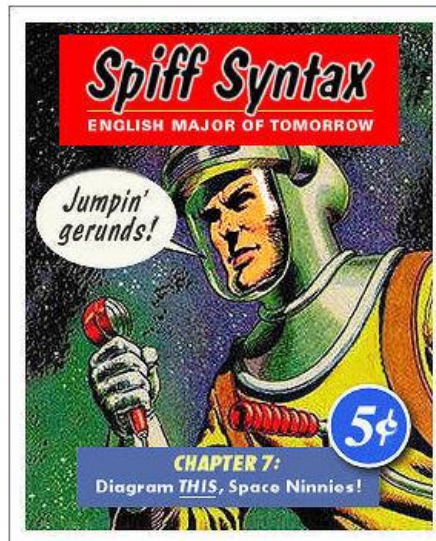


Chapter 3b: Syntax

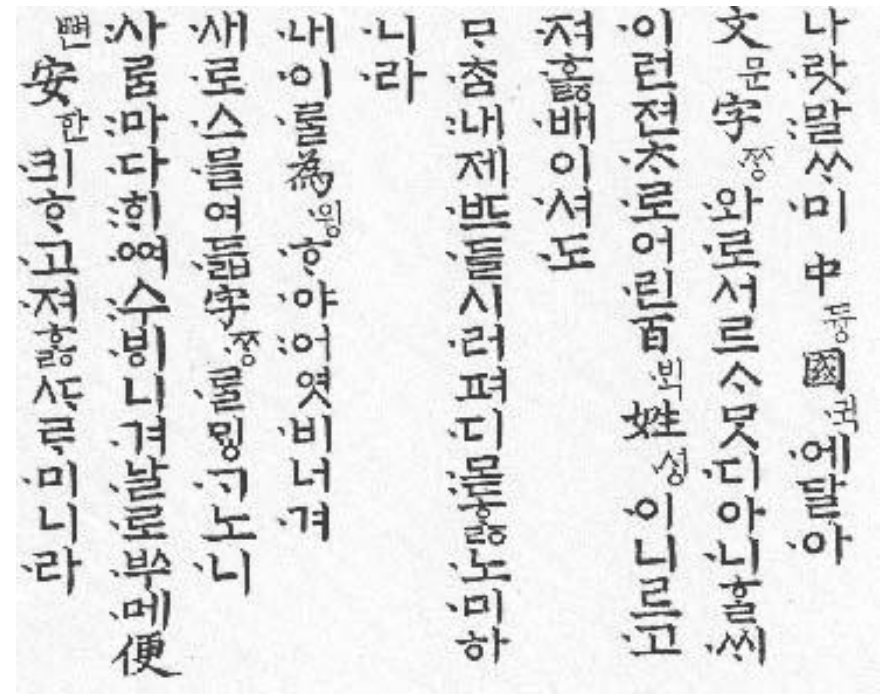
CMSC 124, 1st Semester, AY 2009-10



Chapter 3b: Syntax

Describing Syntax

- **Language** consists of a set of strings (syntactically correct programs) of characters from some alphabet of symbols.
- Strings of a language are called **sentences** or **statements**.
- **Syntax rules** specify which strings of characters from the language's alphabet are in the language.

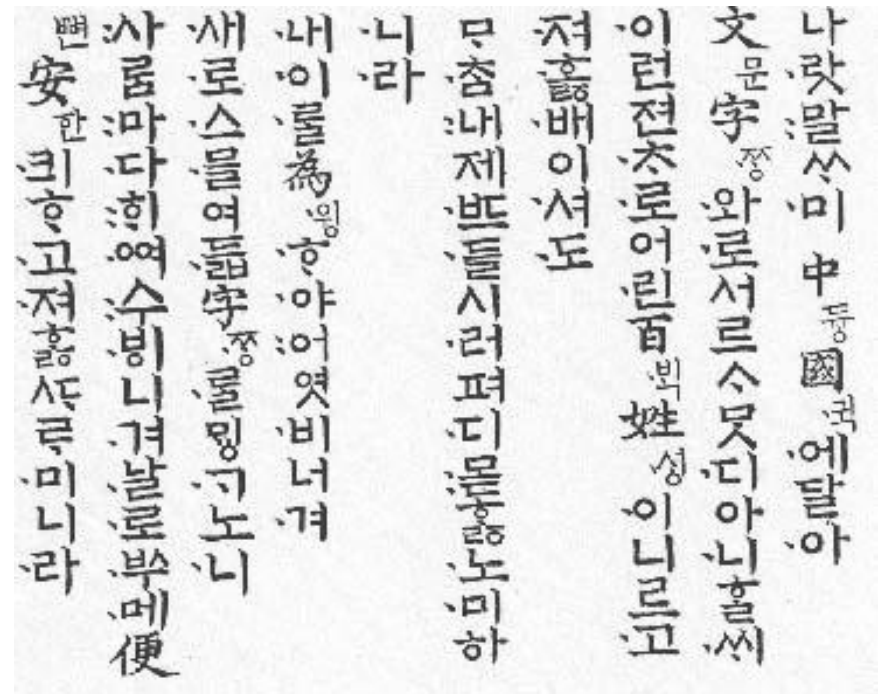


Chapter 3b: Syntax

Describing Syntax

- **Grammar**

- Formal definition of the syntax of the language.
- It naturally defines the hierarchical structure of many PL's.



나랏말싸미 중^중국^국에 달아
文^문字^자와로서 르스 못 다 아 니 할 씨
이런 전^전초^초로 어 린 백^백姓^성이 니 르 고
저 흠 배 이 셔 도
무 촌 내 제 브 들 시 러 퍼 디 문^문호^호노 미 하
니 라
내 이 를 爲^爲 하 야 어 엇 비 너 겨
새 로 스 를 여 름 字^자를 寫^寫하 노 니
사 례 마 다 회 여 수 빙 니 겨 날 로 부 메 便^便
편^편安^안한 크 히 고 저 흠 사 례 미 니 라

Chapter 3b: Syntax

Grammar

Definition: A **grammar** $\langle \Sigma, N, P, S \rangle$ consists of four parts:

1. A finite set Σ of **terminal symbols** of tokens.
2. A finite set N of **non-terminal symbols** or syntactic categories.
3. A finite set P of **productions** or **rules** that describe how each non-terminal is defined in terms of terminal symbols and non-terminals.
4. A distinguished non-terminal S , the **start symbol**, that specifies the category being defined.

Chapter 3b: Syntax

Sample English Grammar

- **Terminals**

- the, boy, ran, ate, cake, a, an

- **Non-Terminals**

- <sentence>, <subject>, <predicate>, <verb>, <article>, <noun>

- **Start Symbol**

- One of the non-terminals.

- **Rules/Productions** (A finite set of replacement rules)

<sentence> ::= <subject> <predicate>

<subject> ::= <article> <noun>

<predicate> ::= <verb> <article> <noun>

<verb> ::= ran | ate

<article> ::= the | a | an

<noun> ::= boy | girl | cake

Try Later!

1. The boy ate the cake.
2. A cake ate the boy.

Chapter 3b: Syntax

Backus-Naur Form

- A grammar used to express the rules/production.
- Originally developed for the syntactic definition of Algol-60
- The BNF grammar is a set of rules or productions of the form:

`left-side ::= right-side`

where

- ✓ `left-side` is a non-terminal.
- ✓ `right-side` is a string of non-terminals and terminals.



John Backus

Chapter 3b: Syntax

Backus-Naur Form

- The BNF grammar is a set of rules or productions of the form:
$$\text{left-side} ::= \text{right-side}$$
- A **terminal** represents the atomic symbols in the language.
- A **non-terminal** represents for other symbols as defined to the right of the symbol “ $::=$ ”.
- The operator “ $::=$ ” is read as “**produces**”.
- The form above is read “**the non-terminal left-side produces right-side**”.

- “ $|$ ” is interpreted as alternative.
- “ $\{ \}$ ” denotes possible repetition of the enclosed symbols 0 or more times.

Eg:

$$A ::= B \mid \{C\} \mid 9$$

“A produces B”

“A produces a string of 0 or more C’s”

“A produces 9”

Something to Ponder

Given the grammar, how do determine if a particular string is a member of the lang.?

Something to Ponder

Chapter 3b: Syntax

Derivations

- Use derivations to determine.
(a valid syntax in PL's)
- Derivation is a sequence of sentential forms starting from the start symbol.
- Either leftmost or rightmost derivation.
- Replace any non-terminal by a right hand side value using any rule.
- Symbol used: " \Rightarrow "

Eg:

$\langle \text{sentence} \rangle ::= \langle \text{subject} \rangle \langle \text{predicate} \rangle$
 $\langle \text{subject} \rangle ::= \langle \text{article} \rangle \langle \text{noun} \rangle$
 $\langle \text{predicate} \rangle ::= \langle \text{verb} \rangle \langle \text{article} \rangle \langle \text{noun} \rangle$
 $\langle \text{verb} \rangle ::= \text{ran} \mid \text{ate}$
 $\langle \text{article} \rangle ::= \text{the} \mid \text{a} \mid \text{an}$
 $\langle \text{noun} \rangle ::= \text{boy} \mid \text{girl} \mid \text{cake}$

Try Now!

1. The boy ate the cake.
2. A cake ate the boy.

$$\begin{array}{l} \frac{\Delta \vec{R} = \Delta \vec{V}}{R \quad \vec{V}} \\ \frac{\Delta \vec{R} = \frac{\Delta \vec{V}}{R \cdot \Delta t} \quad \Delta \vec{R} = \frac{\vec{V}}{\Delta t} \quad \frac{\Delta \vec{V} = \vec{a}}{\Delta t}}{\vec{V} = \frac{\vec{a}}{R}} \\ \frac{\vec{a} = \frac{\vec{V}^2}{R}}{R} \end{array}$$

Chapter 3b: Syntax

Derivations

- **Derivation for #1:**

<sentence> => <subject> <predicate>
=> <article> <noun> <predicate>
=> the <noun> <predicate>
=> the boy <predicate>
=> the boy <verb> <article> <noun>
=> the boy ate <article> <noun>
=> the boy ate the <noun>
=> the boy ate the cake



- Also from <sentence>, the statement **“a cake ate the boy”** can also be derived. What does that mean?
- Syntax **does not imply** correct semantics.

Chapter 3b: Syntax

Derivations: Another Example

- Show that 010 is a member of the following grammar:

$$\checkmark \langle B \rangle ::= 0\langle B \rangle \mid 1\langle B \rangle \mid 0 \mid 1$$

- Rightmost or Leftmost Derivation

$$\begin{aligned} \checkmark \langle B \rangle &\Rightarrow 0\langle B \rangle \\ &\Rightarrow 01\langle B \rangle \\ &\Rightarrow 010 \end{aligned}$$

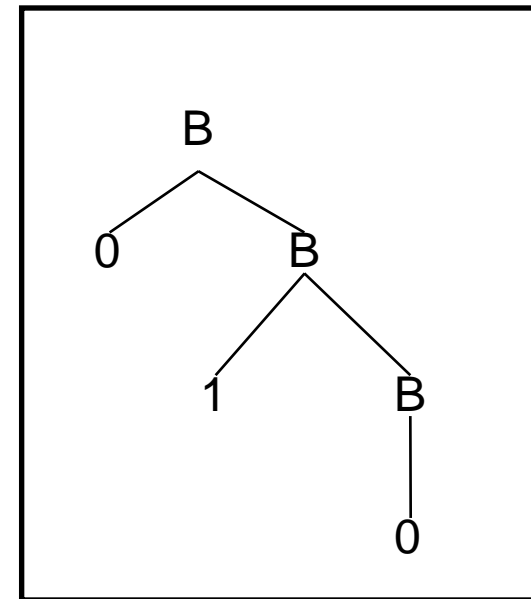
Therefore, it is
a valid string!



Chapter 3b: Syntax

Derivation/Parse Tree

- Graphically shows how the start symbol of a grammar derives a string in the language.
- Using a parse tree, show that 010 is a member of the grammar:
 - **$B \rightarrow 0B \mid 1B \mid 0 \mid 1$**

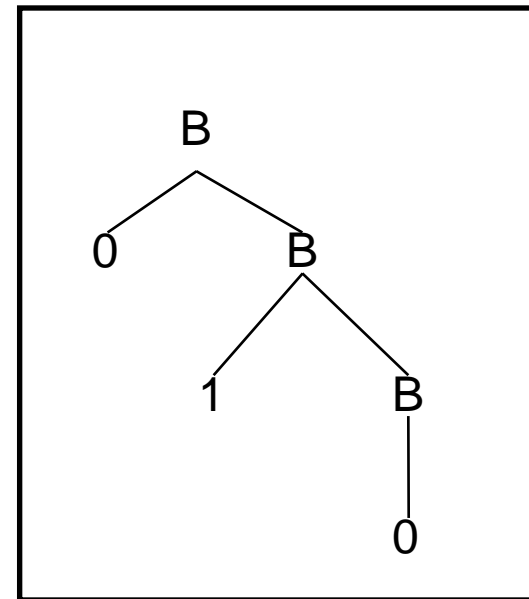


Chapter 3b: Syntax

Derivation/Parse Tree

Properties of a Parse Tree

- The **root** is labeled by the start symbol.
- Each **leaf** is labeled by a token or by ϵ .
- Each **interior node** is labeled by a non-terminal.
- If A is the non-terminal labeling some interior node and x_1, x_2, \dots, x_n are the children of that node from left to right then **A \rightarrow x_1, x_2, \dots, x_n** is a **production**.



Chapter 3b: Syntax

Derivation Tree: Sample PL Grammar

$\langle \text{expression} \rangle ::= \langle \text{term} \rangle \mid \langle \text{expression} \rangle \langle \text{addoperator} \rangle \langle \text{term} \rangle$

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

$\langle \text{factor} \rangle ::= \langle \text{identifier} \rangle \mid \langle \text{literal} \rangle \mid (\langle \text{expression} \rangle)$

$\langle \text{identifier} \rangle ::= a \mid b \mid c \mid \dots \mid z$

$\langle \text{literal} \rangle ::= 0 \mid 1 \mid 2 \mid \dots \mid 9$

$\langle \text{addoperator} \rangle ::= + \mid - \mid \text{or}$

$\langle \text{multoperator} \rangle ::= * \mid / \mid \text{div} \mid \text{mod} \mid \text{and}$

Generate a parse tree for the string $a + b * c$.

Chapter 3b: Syntax

Leftmost and Rightmost Derivations

Leftmost Derivation

- In each step, the leftmost non-terminal is replaced.

Eg.

Consider the grammar:

$G = (\{S, A\}, \{a, b\}, P, S)$

$P = \{ S \rightarrow aAS \mid a, A \rightarrow SbA \mid SS \mid ba \}$

Derive the string **aabbbaa**.

Rightmost Derivation

- In each step, the rightmost non-terminal is replaced.

Something to Ponder

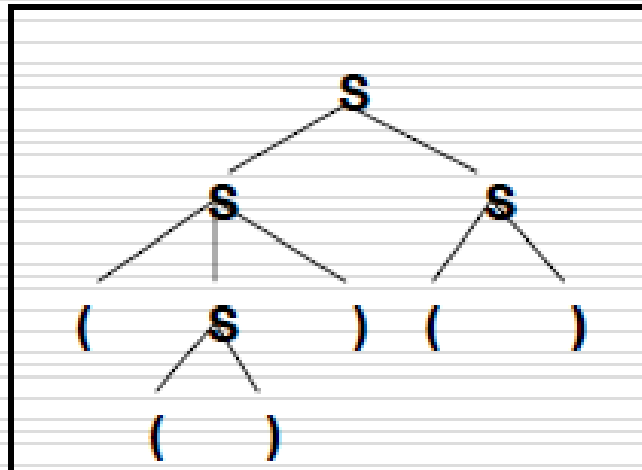
A certain statement can be generated by
2 or more distinct LM/RM derivations.
What's the implication?

Something to Ponder

Chapter 3b: Syntax

Different Derivations, Same Parse Tree

- Derivations may not be **unique**.
 - **Grammar:** $SS \rightarrow SS \mid (S) \mid ()$
 - **Sentence:** $((()))()$
 - **Derivations:**
 - $S \Rightarrow SS \Rightarrow (S)S \Rightarrow (())S \Rightarrow (())()$
 - $S \Rightarrow SS \Rightarrow S() \Rightarrow (S)() \Rightarrow (())()$
 - Different derivations but same parse tree.



Something to Ponder

Another statement can be generated
by 2 or more distinct parse trees.
What's the implication?

Something to Ponder

Chapter 3b: Syntax

Ambiguity

- A grammar that generates a sentence for which there are 2 or more distinct parse trees is said to be ambiguous.

Eg:

Consider the following grammar:

<assign> -> <id> = <expr>

<id> -> A | B | C

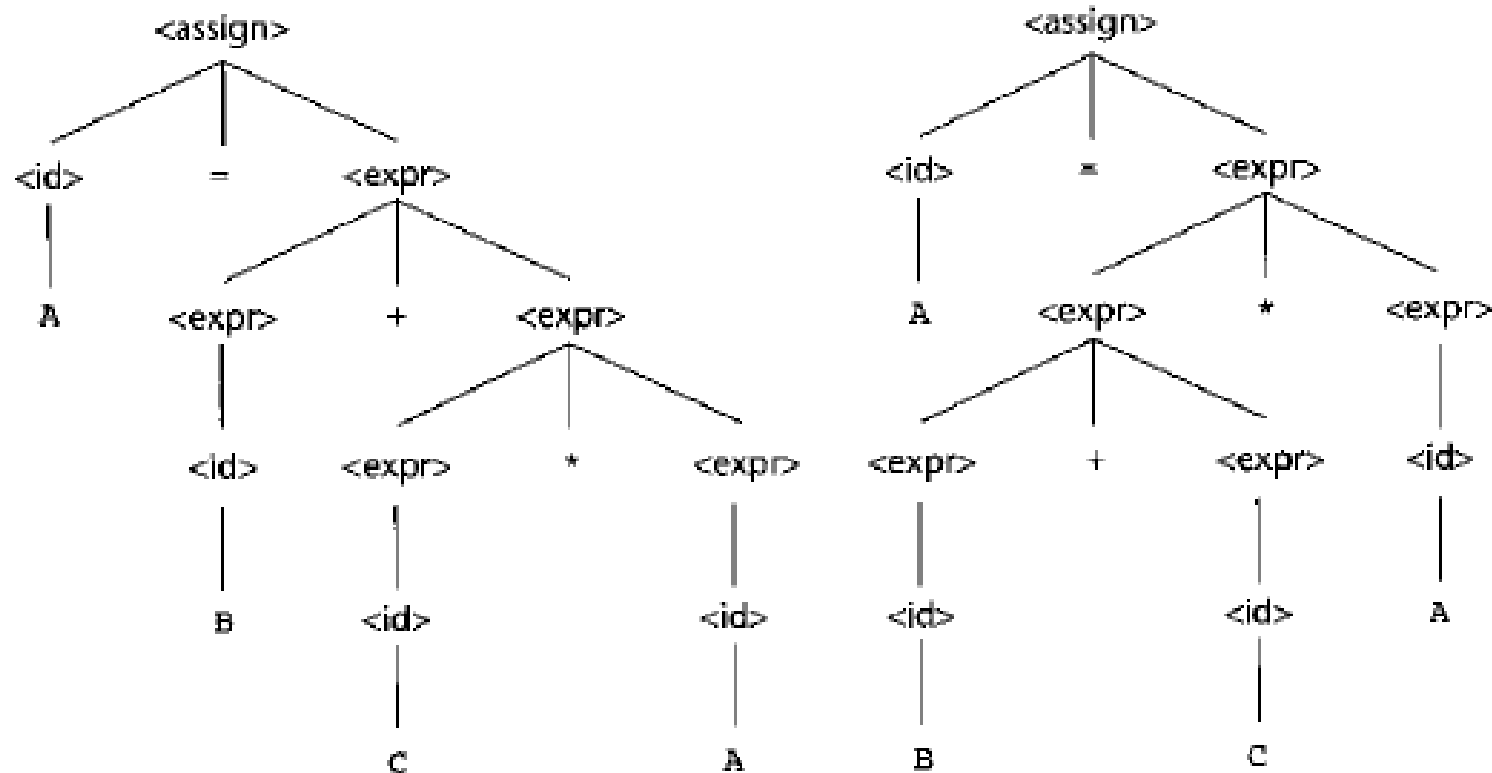
<expr> -> <expr> + <expr> | <expr> * <expr> | (<expr>) | <id>

Show the derivation for the sentence:

A = B + C * A

Chapter 3b: Syntax

Ambiguity



Two distinct parse trees for **$A = B + C * A$**

Chapter 3b: Syntax

The Problem of Ambiguity

- Syntactic ambiguity of language structures is a problem for compilers often base the semantics of those structures on their **syntactic form**.
- The compiler decides what code to generate for a stmt by examining its parse tree.
- Meaning of the structure cannot be determined **uniquely** if there are >1 parse tree.



Chapter 3b: Syntax

Ambiguity Issue #1: Operator Precedence

- A grammar can describe a certain syntactic structure so that part of the structure's meaning can follow its parse tree.
- **A fact:**
"If an operator in an arithmetic expression is generated **lower** in the parse tree, it indicates that it has **higher** precedence over an operator produced higher up in the tree."
- A grammar can be rewritten to separate addition and multiplication operators.
- Rewriting would require additional non-terminals and some new rules.



Chapter 3b: Syntax

Ambiguity Issue #1: Operator Precedence

ORIGINAL GRAMMAR

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

MODIFIED GRAMMAR

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$
 $\langle \text{id} \rangle \rightarrow A \mid B \mid C$
 $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \mid$
 $\langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid$
 $\langle \text{factor} \rangle$
 $\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle$

NOW!

Generate a parse tree for **A = B + C * A**

Chapter 3b: Syntax

Ambiguity Issue #2: Operator Associativity

- Another interesting question is whether operator associativity is also **correctly described**.
 - Expressions with 2 or more adjacent occurrences of operators with equal precedence have those occurrences in proper hierarchical order?



Chapter 3b: Syntax

Ambiguity Issue #2: Operator Associativity

- Consider the following statements:
 - A = B + C + A**
 - How do we specify the associativity?

Left Associativity

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A \mid B \mid C$

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

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid$
 $\langle \text{factor} \rangle$

$\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle$

Left Recursive

A BNF rule has its left hand side (LHS) appear at the beginning of its RHS.

Chapter 3b: Syntax

Ambiguity Issue #2: Operator Associativity

- Consider the following statements:
 - A = B + C + A**
 - How do we specify the associativity?

Right Associativity

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A \mid B \mid C$

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

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid$
 $\langle \text{factor} \rangle$

$\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle$

Right Recursive

A BNF rule has its left hand side (LHS) appear at the right end of its RHS.

Chapter 3b: Syntax

Ambiguity Issue #2: Operator Associativity

- In addition and multiplication, it does not matter what kind of associativity to be used. **Why?**
- But in other operations, like **exponentiation**, the kind of associativity should be **defined**.
- In most PL's, the exponentiation operator is right associative.
 - A right recursive grammar rule is provided.

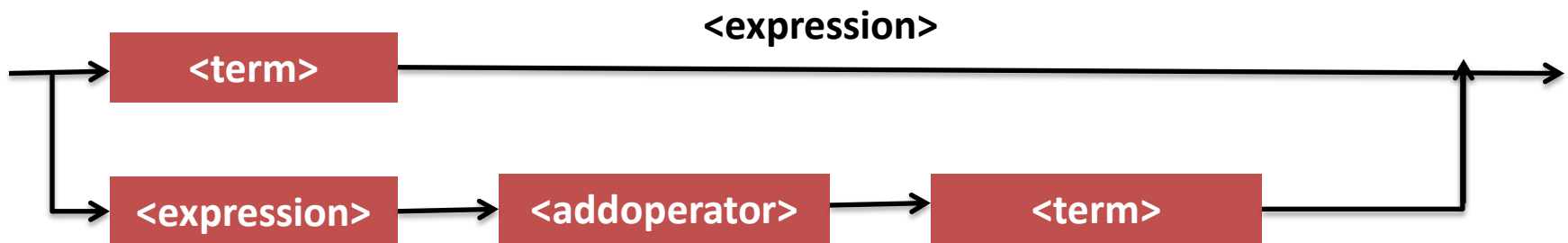
Eg:

```
<factor> -> <expr> ** <factor> |  
           <expr>  
<expr> -> (<expr>) | <id>
```

Chapter 3b: Syntax

Syntax Diagrams

- A **graphical way** used to represent BNF rules.
- For each grammar rule, an equivalent syntax diagram can be drawn.
- This was popularized in the design of Pascal.
- **Symbols:**
 - Rectangle nodes for non-terminals.
 - Circles for terminals.



Sample syntax diagram for
 $\langle \text{expression} \rangle ::= \langle \text{term} \rangle \mid \langle \text{expression} \rangle \langle \text{addoperator} \rangle \langle \text{term} \rangle$