

# Specifying Syntax with BNFs (Backus Naur Forms)

Programming Languages

CS 214



# Sentence Structure

Consider the following “sentences”

- *there is hair in my soup*
- *there is soup in my hair*
- *is there soup in my hair*

What does each “sentence” mean?

What about this “sentence”?

- *hair soup there my is in*

→ The \_\_\_\_\_ determines the meaning...

What determines \_\_\_\_\_ ?



# Grammar & Syntax

Every language has a set of rules — its \_\_\_\_\_ or \_\_\_\_\_ — that specifies the word-sequences that form valid sentences.

The “sentence”:

- *hair soup there my is in*

is not valid, because it violates English’s grammar rules (i.e., it contains \_\_\_\_\_).

Grammar and syntax help us decode a sentence’s meaning:

*syntax were read to easier sentence would unimportant this if be  
this sentence would be easier to read if syntax were unimportant*



# Semantics

When a sentence has correct syntax, our brains can determine what the words in the sentence mean: their \_\_\_\_\_.

Consider: *time flies like an arrow, fruit flies like a banana*

→ \_\_\_\_\_ and \_\_\_\_\_ have two different meanings

We decode a word's meaning using its \_\_\_\_\_.

Since a word's semantics (meaning) depends on its context, English is a \_\_\_\_\_

If a sentence contains syntax errors, we can't understand it, because syntax specifies what words can be adjacent, and a word's semantics depends on the words surrounding it.



# Program Sentences

A program is a “sentence” in a programming language.

A program’s “meaning” depends on the order of its symbols.

To be valid, the order must obey the syntax rules of the language; for example:

$$\mathbf{x = y + 1 ;}$$

is a valid statement in some languages (e.g., C++, Java)  
but not in others (e.g., Lisp, Ada).

The meanings of the “words” in a program are determined by the \_\_\_\_\_ of the language in which it’s written.





# BNF

The \_\_\_\_\_ is a tool for specifying the syntax of a high level language (HLL).

Example: A BNF giving the structure of C++ identifiers is:

```
<identifier> ::= <first_symbol> <valid_sequence>
<first_symbol> ::= <letter> | _
<valid_sequence> ::= <valid_symbol> <valid_sequence> | ε
<valid_symbol> ::= <letter> | <digit> | _
<letter> ::= A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q |
           R | S | T | U | V | W | X | Y | Z | a | b | c | d | e | f | g | h | i |
           j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z
<digit> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
```

A correct BNF specifies all valid “sentences”, and prohibits all invalid “sentences”.



# Examples:

Is *R2D2* a valid sentence in our <identifier> language?

<identifier>

<first\_symbol> <valid\_sequence>

<letter> <valid\_sequence>

R <valid\_sequence>

R <valid\_symbol> <valid\_sequence>

R <digit> <valid\_sequence>

R 2 <valid\_sequence>

R 2 <valid\_symbol> <valid\_sequence>

R 2 <letter> <valid\_sequence>

R 2 D <valid\_sequence>

R 2 D <valid\_symbol> <valid\_sequence>

R 2 D <digit> <valid\_sequence>

R 2 D 2 <valid\_sequence>

R 2 D 2  $\epsilon$

This sequence of steps  
is called a \_\_\_\_\_.

“Sentences” not conforming to the BNF are \_\_\_\_\_.



# Formal Definitions

Let  $\Sigma$  be a set of symbols.

- A \_\_\_\_\_  $\Sigma$  is a finite sequence of zero or more symbols from the set  $\Sigma$ .
- Symbols whose meaning is predefined are called \_\_\_\_\_.
  - Symbols like  $A$ ,  $_$ ,  $\delta$ , etc. are terminals in our  $\langle \text{identifier} \rangle$  BNF.
- Symbols whose meanings must be defined are called \_\_\_\_\_, and are enclosed in angle-brackets ( $\langle$  and  $\rangle$ ).
  - Symbols like  $\langle \text{identifier} \rangle$ ,  $\langle \text{letter} \rangle$ , etc. are non-terminals.
  - Like variables, non-terminals usually describe what they represent.
- One symbol is designated as the \_\_\_\_\_.
  - The symbol  $\langle \text{identifier} \rangle$  is the starting non-terminal in our BNF.



# Formal Definitions (ii)

- Each non-terminal must be defined by a \_\_\_\_\_ (rule):

$$\langle \text{LHS} \rangle ::= \langle \text{RHS} \rangle$$

where:  $\langle \text{LHS} \rangle$  is the nonterminal being defined,  
 $\langle \text{RHS} \rangle$  is a string of terminals and/or nonterminals defining  $\langle \text{LHS} \rangle$ .

- Different productions defining the same non-terminal:

$$\langle \text{NT}_i \rangle ::= \text{Def}_1$$

$$\langle \text{NT}_i \rangle ::= \text{Def}_2$$

...

$$\langle \text{NT}_i \rangle ::= \text{Def}_n$$

can be written in shorthand using the \_\_\_\_\_ operator:

$$\langle \text{NT}_i \rangle ::= \text{Def}_1 \mid \text{Def}_2 \mid \dots \mid \text{Def}_n$$



# Formal Definitions (iii)

- A BNF is a quadruple:  $(\Sigma, N, P, S)$ , where:
  - $\Sigma$  is the set of symbols in the BNF;
  - $N$  is the subset of  $\Sigma$  that are nonterminals in the language;
  - $P$  is the set of productions defining the symbols in  $N$ ; and
  - $S$  is the element of  $N$  that is the starting nonterminal.
- A \_\_\_\_\_ is a sequence of strings, beginning with the starting nonterminal  $S$ , in which each successive string replaces a nonterminal with one of its productions, and in which the final string consists solely of terminals.
  - A derivation is sometimes called a \_\_\_\_\_.



# Formal Definitions (iv)

- A BNF \_\_\_\_\_ is a tree, such that
  - the root of the tree is the starting nonterminal (S) in the BNF;
  - the children of the root are the symbols (L to R) in a production whose <LHS> is S, the starting nonterminal;
    - each terminal child is a leaf; and
    - each nonterminal child is the root of a derivation tree for that nonterminal.
- The act of building a derivation tree for a sentence (to check its correctness) is called \_\_\_\_\_ that sentence.
  - The set of valid sentences in a language is the set of all sentences for which a parse tree exists!
- A \_\_\_\_\_ is a derivation built by *always expanding the left-most non-terminal* in a production.

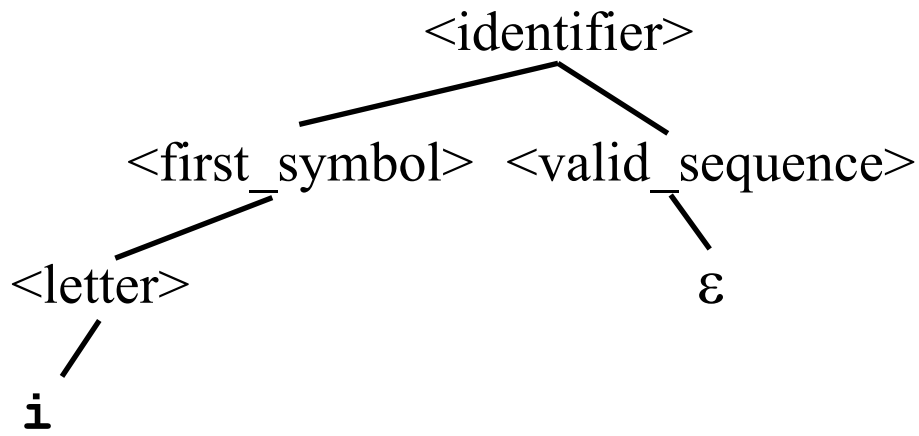


# Examples

Do parse trees (using our <identifier> BNF) exist for these?

*i*





# Examples

Do parse trees (using our <identifier> BNF) exist for these?

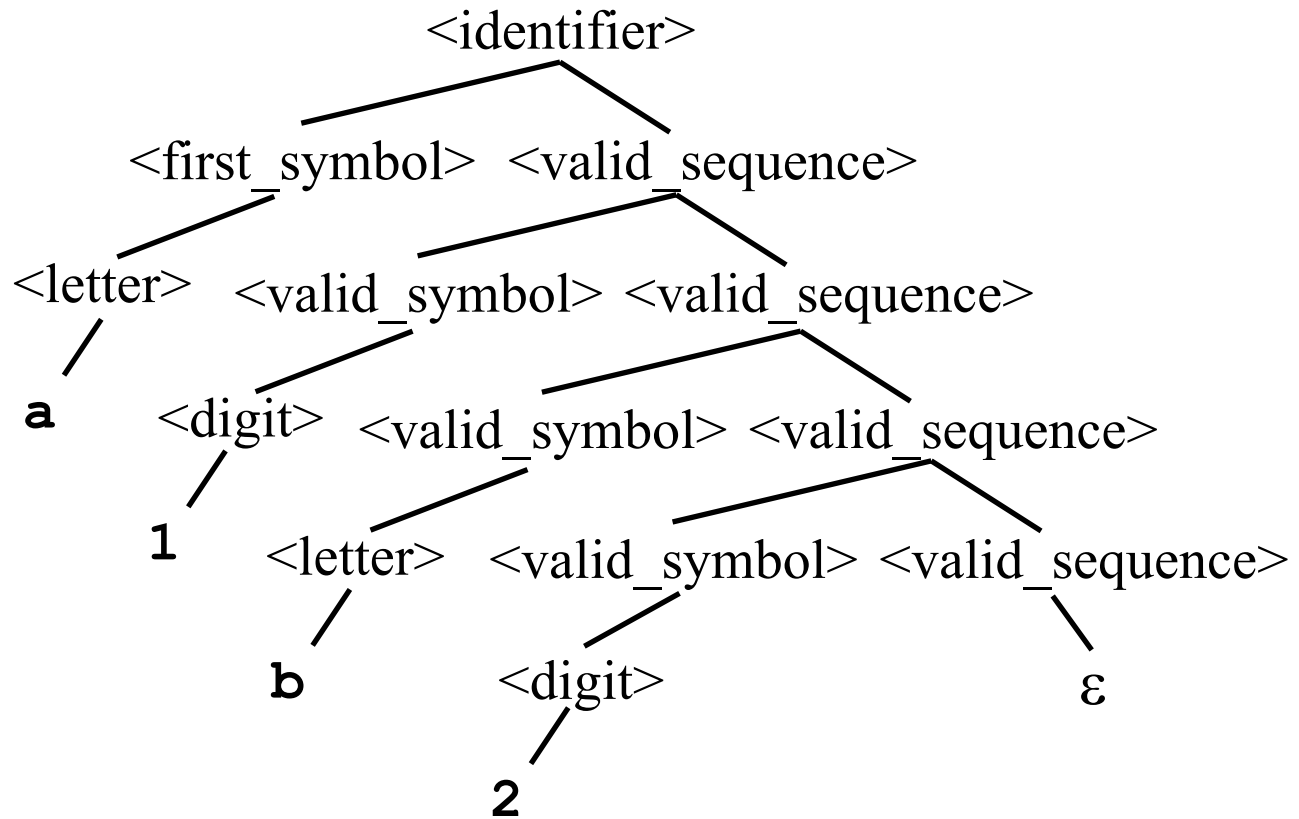
*i*

*a1b2*



~~a~~ ~~1~~ ~~b~~ ~~2~~

---



# Examples

Do parse trees (using our <identifier> BNF) exist for these?

*i*

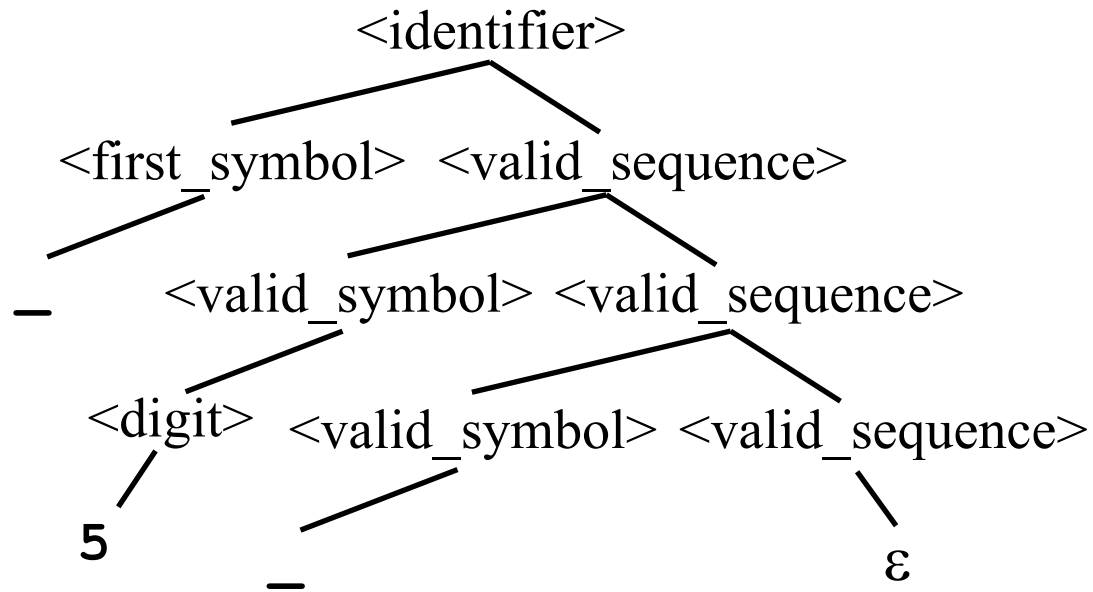
*a1b2*

5



# LLP

---



# Examples

Do parse trees (using our <identifier> BNF) exist for these?

*i*

*a1b2*

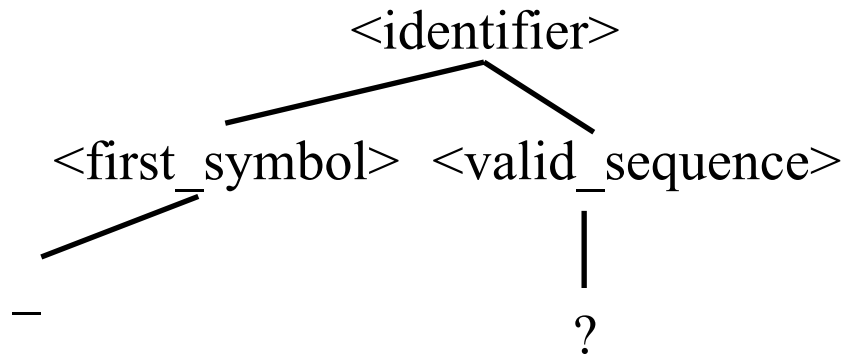
5

\$



**L** \$ \_

---



Parse Error: Expecting letter, digit, or \_ but saw \$

# Recursive Productions

Our identifier-BNF permits identifiers to have arbitrary lengths through the use of \_\_\_\_\_:

$\langle \text{valid\_sequence} \rangle ::= \langle \text{valid\_symbol} \rangle \langle \text{valid\_sequence} \rangle \mid \varepsilon$

This production is recursive because the non-terminal on the  $\langle \text{LHS} \rangle$  appears in the production (on the  $\langle \text{RHS} \rangle$ ).

The recursive production provides for \_\_\_\_\_ of the non-terminal being defined, which is useful what is being defined can be appear 0 or more times.

The  $\varepsilon$ -production provides both:

- a \_\_\_\_\_ for trivial instances of the non-terminal, and
- an \_\_\_\_\_ to terminate the recursion.



# Writing BNFs: A 3-Step Process

- A. Start with the non-terminal you're defining.
- B. Build the productions to define the non-terminal:
  1. Start with the question: \_\_\_\_\_
  2. If a construct is \_\_\_\_\_:
    - a. create a new nonterminal for that construct;
    - b. add a production for the non-optional case;
    - c. add an  $\epsilon$ -production for the optional case.
  3. If a construct can be \_\_\_\_\_:
    - a. create a new nonterminal for that construct;
    - b. add an  $\epsilon$ -production for the zero-reps case;
    - c. add a recursive production for the other cases.
- C. For each nonterminal in the  $\langle \text{RHS} \rangle$  of every production:  
\_\_\_\_\_.



# Example: C++ *if* Statement

A. Create a non-terminal for what we're defining:

`<if_stmt>`

B. Build a production to define it: what comes first?

1. keyword `if`
2. open parentheses
3. an expression
4. close parentheses
5. a statement
6. (optional) an `else` and a statement

`<if_stmt>` ::= `if ( <expr> ) <stmt> <else_part>`

C. Repeat B for each undefined non-terminal:

`<else_part>` ::= `else <stmt> | ε`

(We'll see how to define `<expr>` and `<stmt>` a bit later...)



Our dog is named Minton.

He ate the shuttlecock from one of our games...

Bad Minton!

# A Small Problem

The C++ if statement presents a small problem: Consider...

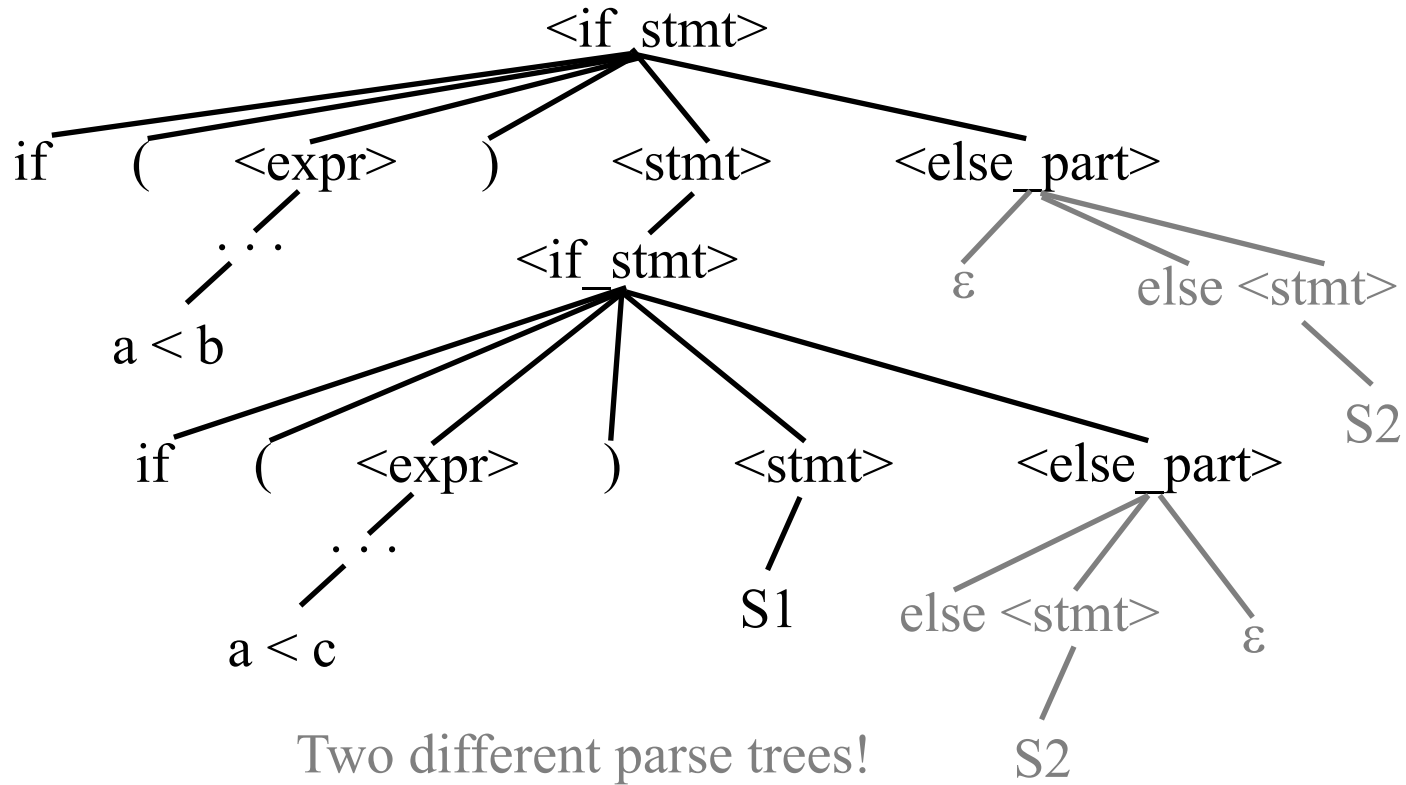
```
if ( a < b) if (a < c) S1 else S2
```

Take a moment to build a parse tree for this “sentence”...



~~if~~ ~~(~~ ~~a < b~~ ~~)~~ ~~if~~ ~~(~~ ~~a < c~~ ~~)~~ ~~S1~~ else S2

---



# Ambiguities

When a “sentence” has multiple parse trees, it is \_\_\_\_\_  
(i.e., it has multiple interpretations and/or meanings).

The two parses reflect different ways to associate the **else**:

```
if (a < b)
    if (a < c)
        ↑
        else
            S2
            S1
```

```
if (a < b)
    ↑
    if (a < c)
        else
            S2
            S1
```

The grammar cannot resolve which is meant,  
so C++ uses a *semantic rule* to resolve the ambiguity:



Did you hear about the guy who wanted a brain transplant?

He changed his mind!



What is a British physicist's favorite food?

Fission chips!

# Writing BNFs: A 3-Step Process

- A. Start with the non-terminal you're defining.
- B. Build the productions to define the non-terminal:
  1. Start with the question: \_\_\_\_\_
  2. If a construct is \_\_\_\_\_:
    - a. create a new nonterminal for that construct;
    - b. add a production for the non-optional case;
    - c. add an  $\epsilon$ -production for the optional case.
  3. If a construct can be \_\_\_\_\_:
    - a. create a new nonterminal for that construct;
    - b. add an  $\epsilon$ -production for the zero-reps case;
    - c. add a recursive production for the other cases.
- C. For each nonterminal in the  $\langle \text{RHS} \rangle$  of every production:  
\_\_\_\_\_.



# Example: C++ *do* Statement

A. Create a non-terminal for what we're defining:

`<do_stmt>`

B. Build a production to define it: what comes first?

1. keyword `do`
2. a statement
3. keyword `while`
4. open parentheses
5. an expression
6. close parentheses
7. a semicolon

`<do_stmt> ::= do <stmt> while ( <expr> );`

C. Repeat B for each undefined non-terminal...

`<stmt>` isn't too bad, so let's tackle it next.



# Example: C++ $\langle \text{stmt} \rangle$

A. Create a non-terminal for what we're defining:

$\langle \text{stmt} \rangle$

B. Build a production to define it: what comes first?

- It depends on the kind of statement being described...
- There are seven different kinds of C++ statements, so let's introduce a new non-terminal for each one:

$$\begin{aligned} \langle \text{stmt} \rangle & ::= \langle \text{compound\_stmt} \rangle \mid \langle \text{selection\_stmt} \rangle \mid \\ & \quad \langle \text{iteration\_stmt} \rangle \mid \langle \text{expression\_stmt} \rangle \mid \\ & \quad \langle \text{jump\_stmt} \rangle \mid \langle \text{labeled\_stmt} \rangle \mid \langle \text{declaration\_stmt} \rangle \end{aligned}$$

C. Repeat B for each undefined non-terminal...

$$\begin{aligned} \langle \text{compound\_stmt} \rangle & ::= \{ \langle \text{stmt\_list} \rangle \} \\ \langle \text{stmt\_list} \rangle & ::= \langle \text{stmt} \rangle \langle \text{stmt\_list} \rangle \mid \varepsilon \end{aligned}$$


# Example: C++ <stmt> (ii)

<selection_stmt>	::= <if_stmt>   <switch_stmt>
<iteration_stmt>	::= <while_stmt>   <do_stmt>   <for_stmt>
<expression_stmt>	::= <opt_expr> ;
<opt_expr>	::= <expr>   $\epsilon$
<jump_stmt>	::= break ;   continue ;   return <opt_expr> ;   goto <identifier> ;
<labeled_stmt>	::= <identifier> : <stmt>   case <literal> : <stmt>   default : <stmt>
<declaration_stmt>	::= <object_dec>   <function_dec>   <class_dec>
<object_dec>	::= <modifier> <type> <identifier> <initializer> <more_ids> ;
<modifier>	::= register   const   static   auto   extern   mutable   $\epsilon$
<type>	::= char   wchar_t   bool   short   int   long   signed   unsigned   float   double   <identifier>   <identifier> :: <type>
<initializer>	::= = <expr>   $\epsilon$
<more_ids>	::= , <identifier> <initializer> <more_ids>   $\epsilon$
...	



# Example: C++ Assignment Exprs

A. Create a non-terminal for what we're defining:

`<assign_expr>`

B. Build a production to define it: what comes first?

- Many things might come first (variable, pointer, array, ...) so let's introduce a non-terminal to hide the details.
- Next comes an assignment operator; then an expression:

`<assign_expr> ::= <lvalue> <assign_op> <assign_expr> | ...`

C. Repeat B for each undefined non-terminal...

`<lvalue> ::= <unary_op> <identifier> <id_suffix>`

`<unary_op> ::= * | ε`

`<id_suffix> ::= [ <expr> ] <id_suffix> | . <lvalue> <id_suffix> |  
-> <lvalue> <id_suffix> | ε`

`<assign_op> ::= = | += | -= | *= | /= | %= | <<= | >>= | &= | |= | ^=`



# Expressions

C++ expressions are complicated:

- Each of its 52 operators must be included
- Each of its 17 precedence levels must be included
- Associativity rules must be enforced
- The grammar/BNF must be unambiguous

Example: For the expression:

$$2 + 3 * 4$$

our grammar must ensure that:

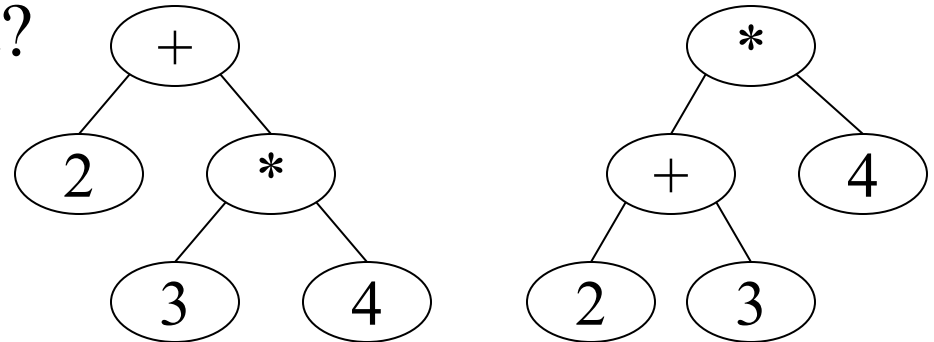
$$2 + (3 * 4) = 14$$

is evaluated, not:

$$(2 + 3) * 4 = 20$$

Which parse tree is correct?

Our grammar must *only* generate the correct one.



# Grammar and Precedence

To ensure that higher precedence operators appear lower in the parse tree, we must build \_\_\_\_\_ into our BNF:

```

<expr>          ::= ... | <add_expr> | ...
<add_expr>     ::= <mul_expr> | <add_expr> + <mul_expr> |
                  <add_expr> - <mul_expr>
<mul_expr>     ::= <value> | <mul_expr> * <expr> |
                  <mul_expr> / <expr> | <mul_expr> % <expr>
    
```

Rules like these ensure that “multiply-level” operators will be applied before “add-level” operators...

```

          <expr>
          <add_expr>
          <add_expr> + <mul_expr>
<mul_expr> <mul_expr> * <expr>
<value>    <value>    <add_expr>
           2          3          <mul_expr>
                               <value>
                               4
    
```

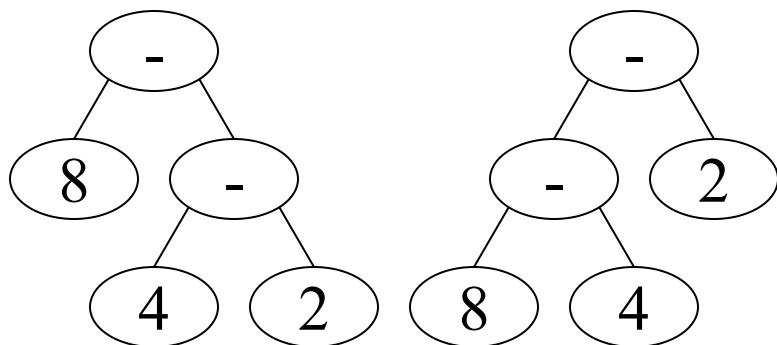


# Grammar and Associativity

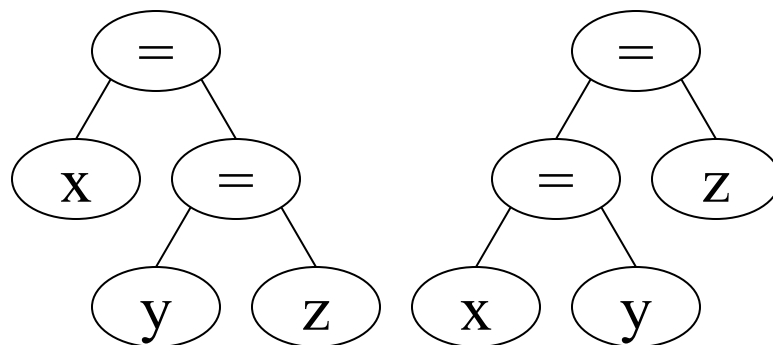
Associativity gives ordering of equal-precedence operators.

Examples:  $8 - 4 - 2$  vs.  $x = y = z$

$-$  is \_\_\_\_\_;  
which is correct?



$=$  is \_\_\_\_\_;  
which is correct?



Associativity can be built into a grammar by using \_\_\_\_\_ for left-associative operators;  
and \_\_\_\_\_ for right-associative operators.



# Associativity Examples

Example: Since +, - are left-associative, we write:

$$\langle \text{add\_expr} \rangle ::= \langle \text{mul\_expr} \rangle | \langle \text{add\_expr} \rangle + \langle \text{mul\_expr} \rangle | \langle \text{add\_expr} \rangle - \langle \text{mul\_expr} \rangle$$

but since = is right-associative, we write what amounts to:

$$\langle \text{assign\_expr} \rangle ::= \langle \text{lvalue} \rangle \langle \text{assign\_op} \rangle \langle \text{assign\_expr} \rangle$$

These generate the correct parse trees for each expression:

$$\begin{array}{c} \langle \text{expr} \rangle \\ \langle \text{add\_expr} \rangle \\ \langle \text{add\_expr} \rangle - \langle \text{mul\_expr} \rangle \\ \langle \text{add\_expr} \rangle - \langle \text{mul\_expr} \rangle \quad \langle \text{value} \rangle \\ \langle \text{mul\_expr} \rangle \quad \langle \text{value} \rangle \quad 2 \\ \langle \text{value} \rangle \quad 4 \\ 8 \end{array}$$

$$\begin{array}{c} \langle \text{expr} \rangle \\ \langle \text{assign\_expr} \rangle \\ \langle \text{lvalue} \rangle \langle \text{assign\_op} \rangle \langle \text{assign\_expr} \rangle \\ \langle \text{identifier} \rangle = \langle \text{lvalue} \rangle \langle \text{assign\_op} \rangle \langle \text{assign\_expr} \rangle \\ x \quad \langle \text{identifier} \rangle = \langle \text{add\_expr} \rangle \\ y \quad \quad \quad \langle \text{mul\_expr} \rangle \\ \quad \quad \quad \langle \text{identifier} \rangle \\ \quad \quad \quad z \end{array}$$


Ann: Did you know that photons have mass?

Bob: I didn't even know they were Catholic!

# C++ Expressions

<expr>	::=	<assign_expr>   <expr> , <assign_expr>
<assign_expr>	::=	<lvalue> <assign_op> <assign_expr>   <cond_expr>
<cond_expr>	::=	<lor_expr>   <lor_expr> ? <expr> : <cond_expr>
<lor_expr>	::=	<land_expr>   <lor_expr>    <land_expr>
<land_expr>	::=	<bor_expr>   <land_expr> && <bor_expr>
<bor_expr>	::=	<xor_expr>   <bor_expr>   <xor_expr>
<xor_expr>	::=	<band_expr>   <xor_expr> ^ <band_expr>
<band_expr>	::=	<equ_expr>   <band_expr> & <equ_expr>
<equ_expr>	::=	<rel_expr>   <equ_expr> == <rel_expr>   <equ_expr> != <rel_expr>
<rel_expr>	::=	<shft_expr>   <rel_expr> < <shft_expr>   <rel_expr> > <shft_expr>   <rel_expr> <= <shft_expr>   <rel_expr> >= <shft_expr>
<shft_expr>	::=	<add_expr>   <shft_expr> << <add_expr>   <shft_expr> >> <add_expr>
<add_expr>	::=	<mul_expr>   <add_expr> + <mul_expr>   <add_expr> - <mul_expr>
<mul_expr>	::=	<ptr_expr>   <mul_expr> * <ptr_expr>   <mul_expr> / <ptr_expr>   <mul_expr> % <ptr_expr>
<ptr_expr>	::=	...



# Exercise

Build a parse tree for:  $a = x + y / z ;$





Why did the cross-eyed teacher lose his job?

He couldn't control his pupils!

# EBNF

The BNF is the most general tool for expressing syntax.

Another tool frequently used is the \_\_\_\_\_

The differences are:

- EBNF terminals are distinguished from non-terminals by
  - o Capitalizing the first-letter of non-terminals, AND
  - o Underlining, single-quoting, or bolding terminals (instead of surrounding non-terminals by angle-brackets).
- \_\_\_\_\_ may be used to denote grouping.
- { and } surround symbols that are repeated zero or more times.
  - \_\_\_\_\_!
- [ and ] surround symbols that are optional.
  - \_\_\_\_\_!



# Examples

## 1. To specify a C++ block using EBNF:

- First comes a brace, then zero or more statements, then a brace:

Block ::= ‘{ { Stmt } ‘}

## 2. To specify a C++ int literal using EBNF:

- An optional sign, an optional base specifier, at least one digit

Int-Literal ::= [ Sign ] [ 0x ] Digit { Digit }

Sign ::= + | -

Digit ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

## 3. To specify a C++ do statement using EBNF:

- Keyword do, a statement, keyword while, an open-parentheses, an expression, a close-parentheses, a semicolon:

Do-Stmt ::= do Stmt while ‘( ‘ Expr ‘ ) ’ ;



# Problems

1. Specify a C++ identifier using EBNF:



# Specify a C++ Identifier Using EBNF

---

Identifier ::= FirstSymbol { ValidSymbol }

FirstSymbol ::= Letter | \_

Letter ::= A | B | C | D | E | F | G | H | I | J | K | L | M | N |  
P | Q | R | S | T | U | V | W | X | Y | Z | a | b | c |  
d | e | f | g | h | i | j | k | l | m | n | o | p | q | r |  
s | t | u | v | w | x | y | z

ValidSymbol ::= FirstSymbol | Digit

Digit ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

# Problems

1. Specify a C++ identifier using EBNF:
2. Specify a C++ while statement using EBNF:



# Specify a C++ while Stmt Using EBNF

---

WhileStmt ::= while '(' Expr ')' Stmt

Stmt ::= Expr ; | IfStmt | SwitchStmt | ForStmt | WhileStmt |  
DoStmt | BlockStmt | DecStmt | ...

# Problems

1. Specify a C++ identifier using EBNF:
2. Specify a C++ while statement using EBNF:
3. Specify a C++ if statement using EBNF:



# Specify a C++ if Stmt Using EBNF

---

IfStmt ::= if '(' Expr ')' Stmt [ else Stmt ]

# Why use BNFs instead of EBNFs?

The recursive productions in BNFs make it easier for a compiler to parse “sentences” in the language...

## Basic Parsing Algorithm:

0. Push  $S$  (the starting symbol) onto a stack.
  1. Get the first terminal symbol  $t$  from the input file.
  2. Repeat the following steps:
    - a. Pop the stack into  $topSymbol$ ;
    - b. If  $topSymbol$  is a nonterminal:
      - 1) Choose a production  $p$  of  $topSymbol$  based on  $t$
      - 2) If  $p \neq \epsilon$ :  
Push  $p$  right-to-left onto the stack.
    - c. Else if  $topSymbol$  is a terminal &&  $topSymbol == t$ :  
Get the next terminal symbol  $t$  from the input file.
    - d. Else  
Generate a ‘parse error’ message.
- while the stack is not empty.



# Example

Suppose our rules are:

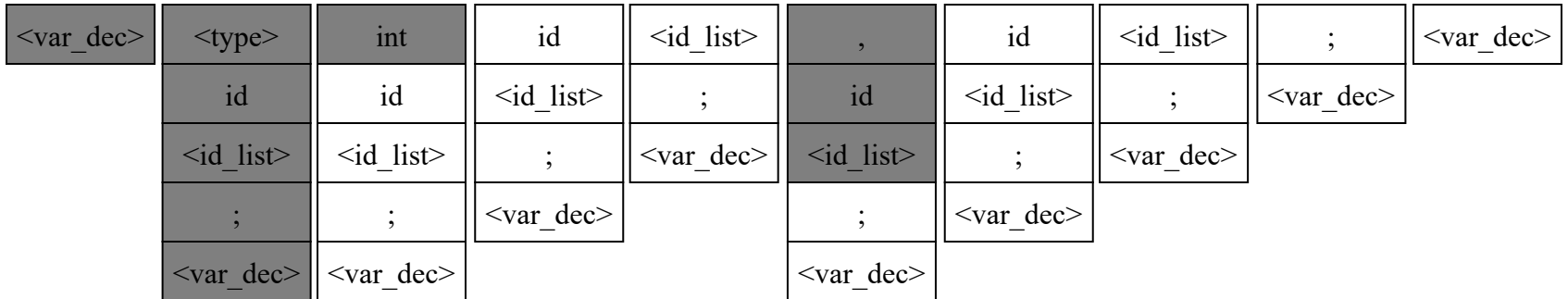
$\langle \text{var\_dec} \rangle ::= \langle \text{type} \rangle \text{ id } \langle \text{id\_list} \rangle ; \langle \text{var\_dec} \rangle \mid \epsilon$

$\langle \text{type} \rangle ::= \text{int} \mid \text{char} \mid \text{float} \mid \text{double} \mid \dots$

$\langle \text{id\_list} \rangle ::= , \text{ id } \langle \text{id\_list} \rangle \mid \epsilon$

Let's parse the declaration: **int x, y;**  
 assuming that  $\langle \text{var\_dec} \rangle$  is our starting symbol.

stack:



t: int

id (x)

,

id (y)

;



If pronouncing my b's as v's makes me sound Russian, then...  
... soviet!

# Summary

There are different ways to specify the syntax of a language.

Two of them are:

- \_\_\_\_\_
- \_\_\_\_\_

The \_\_\_\_\_ is simpler and easier to use, so it is frequently used in language guides and user manuals.

The \_\_\_\_\_ recursive- and  $\epsilon$ -productions simplify the task of parsing, so it is more useful to compiler-writers.

