

# Compiler Construction

## Week 03- Context-Free Grammars

### A few necessary definitions

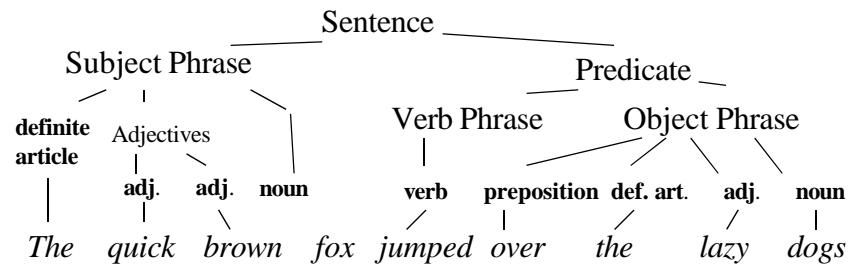
Parse - *vt*, to resolve (as a sentence) into component parts of speech and describe them grammatically

Grammar - *n*, the study of the classes of words, their inflections, and their functions and relations in the sentence

Syntax - *n*, the way in which words are put together to form phrases, clauses or sentences

# The Parsing Process

Syntactic Analysis (or ***Parsing***) involves breaking a program into its ***syntactic*** components



## The Parsing Process (continued)

Nb: In the previous example,  
***subject phrase, predicate, adjectives***, etc. were *nonterminals*.

***definite article, adjective, noun, verb***, etc. were *terminals*

A language is a set of sentences formed the set of basic symbols.

A grammar is the set of rules that govern how we determine if these sentences are part of the language or not.

## The Parsing Process (continued)

The analysis is based purely on **syntax**. A syntactically correct sentence can be nonsensical:

### **Example:**

A group of trout were flying east, where they hunted down camels for their dinner.

## Parsing as a procedure

The parser takes tokens from scanner as necessary and produces a tree structure (or analyzes the program as if it were producing one). It is called as a procedure of the main program:

```
struct parsenoderec      *parsetree;  
parsetree = parse( );
```

In most real cases, the parser actually returns a pointer to an abstract syntax tree or some other intermediate representation.

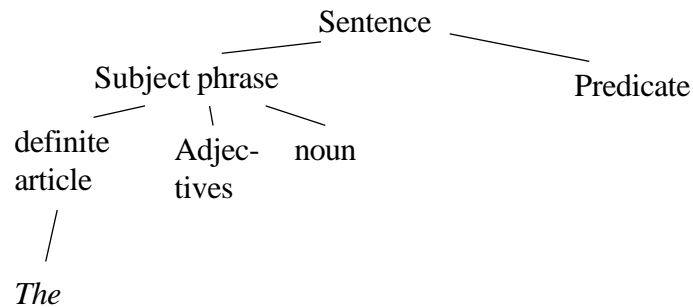
## Error recovery during parsing

- The parser will (or certainly *should*) spot any and all syntactic errors in the program.
- This requires us to consider how we will handle recovery from any errors encountered:
  - We can consider any error fatal and point it out to the user immediately and terminate execution.
  - We can attempt to find a logical place within the program where we can resume parsing so that we can spot other potential errors as well.

## Types of Parsers

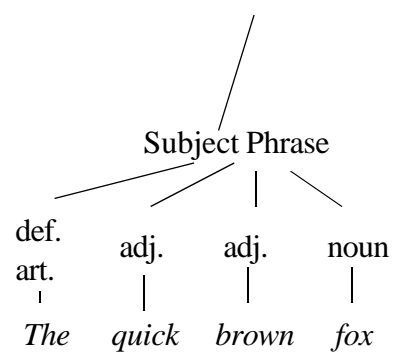
- Parsers can be either top-down or bottom-up:
  - Top-down parsers build the parse-tree starting from the root building until all the tokens are associated with a leaf on the parse tree.
  - Bottom-up parsers build the parse-tree starting from the leaves, assembling the tree fragments until the parse tree is complete.

## Top-down Parsers



Top-down parsing assumes a certain minimum structure as we start building the parse tree

## Bottom-up parsers



Bottom-up parsers *shift* by each token, *reducing* them into a non-terminal as the grammar requires.

Nb: Until we finish building the predicate, we have no reason to reduce anything into the nonterminal *Sentence*

## Types of Parsers (continued)

- Parsers can be either **table-driven** or **handwritten**:
  - Table-driven parsers perform the parsing using a driver procedure and a table containing pertinent information about the grammar. The table is usually generated by automated software tools called **parser generators**.
  - Handwritten parsers are hand-coded using the grammar as a guide for the various parsing procedures.

## Types of Parsers (continued)

- LL(1) and LR(1) parsers are table-driven parsers which are top-down and bottom-up respectively.
- Recursive-descent parsers are top-down hand-written parsers.
- Operator-precedence parsers are bottom-up parsers which are largely handwritten for parsing expressions.

## Context-Free Grammars

A context-free grammar is defined by the 4-tuple:

$$G = (T, N, S, P)$$

where

**T** = The set of *terminals* (e.g., the tokens returned by the scanner)

**N** = The set of *nonterminals* (denoting structures within the language such as *DeclarationSection*, *Function*).

**S** = The *start symbol* (in most instances, our program).

**P** = The set of *productions* (rules governing how tokens are arranged into syntactic units).

## Context-Free Grammars

- Context-free grammars are well-suited to programming languages because they restrict the manner in which programming construct can be used and thus simplify the process of analyzing its use in a program.
- They are called context-free because the manner in which we parse any nonterminal is independent of the other symbols surrounding it (i.e., parsing is done without respect to *context*)
- The grammars of most programming languages are explicitly context-free (although a few have one or two context-sensitive elements).

## Distinction between syntax and semantics

- Syntax refers to features of sentence structure as it appears in languages.
- Semantics refers to the meaning of such structures.
- The parser will analyze the syntax of a program, not its semantics.
  - E. g., the parser does not do type-checking.
  - Semantic actions will frequently be associated with specific productions, but are not actually part of the parser.

## Backus-Naur Form

BNF (**B**ackus-**N**aur **F**orm) is a metalanguage for describing a context-free grammar.

- The symbol  $::=$  (or  $\rightarrow$ ) is used for *may derive*.
- The symbol  $|$  separates alternative strings on the right-hand side.

Example      $E ::= E + T \mid T$   
                   $T ::= T * F \mid F$   
                   $F ::= \text{id} \mid \text{constant} \mid (E)$

where E is *Expression*, T is *Term*, and F is *Factor*



## A simple grammar

Start Symbol  $\Rightarrow$   $S ::= A B c$

$$A ::= a A \mid b$$
$$B ::= A b \mid a$$

The strings *abbbc*, *aaabac*, *aaaababbc* are all generated by this grammar. Can you determine how?

## Another simple grammar

$S ::= a \mid (b S S)$

Sample strings generated by this grammar include :

$(b a a)$      $(b (b a a) a)$      $a$

## The Empty String

- Productions within a grammar can contain  $\epsilon$ , the empty string.
- $A \rightarrow B$  is equivalent to  $A \rightarrow B\epsilon$
- It is also possible to write the production  $A \rightarrow \epsilon$ ; such productions become particularly useful in top-down parsing.

## Derivations

- A derivation is a series of replacements where the nonterminal on the left of a production is replacement by a string of symbols from the right-hand side of a production.
- This may be done in one step or in many steps.

### Example

For the grammar

$$S ::= Aa$$
$$A ::= Ab \mid c$$
$$S \Rightarrow Aa \Rightarrow Aba \Rightarrow Abba \Rightarrow cbba$$

*cbba* is ultimately derived from *S*

## Derivations (continued)

- There are several different notations used to indicate occurs:

$A \Rightarrow \alpha$      $A$  derives  $\alpha$  in one step

$A \Rightarrow^* \alpha$      $A$  derives  $\alpha$  in zero or more steps

$A \Rightarrow^+ \alpha$      $A$  derives  $\alpha$  in one or more steps

- Example

$$S \Rightarrow Aa \Rightarrow Aba \Rightarrow Abba \Rightarrow cbba$$

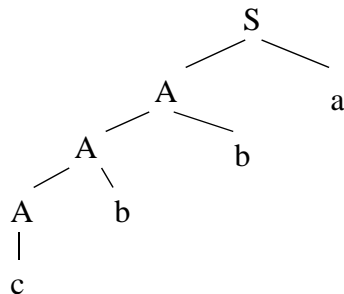
We can say that  $S \Rightarrow^* cbba$

## Derivations (continued)

- If the start symbol  $S$  derives a string  $\beta$  which contains nonterminals,  $\beta$  is a sentential form.
- If  $S$  derives a string  $\beta$  which contains only terminals,  $\beta$  is a sentence.

## Parse Trees

A parse tree is a graphical representation of such a derivation:



# Abstract Syntax Tree

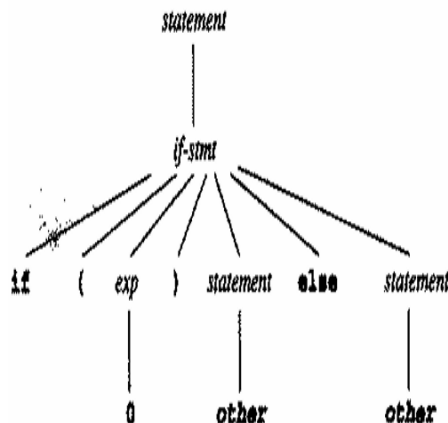
- An abstract syntax tree represents abstractions of actual source code token sequences, and token sequences cannot be recovered from them (unlike parse trees).
- Abstract syntax trees is a tree representation of a shorthand notation called abstract syntax.

## Parse Tree vs Abstract Syntax Tree

$statement \rightarrow if-stmt \mid other$   
 $if-stmt \rightarrow if ( exp ) statement$   
 $\quad \quad \quad \mid if ( exp ) statement else statement$   
 $exp \rightarrow 0 \mid 1$

String: **if {0} other else other**

Parse Tree:



Syntax Tree:

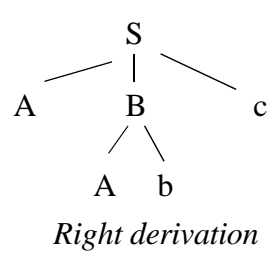
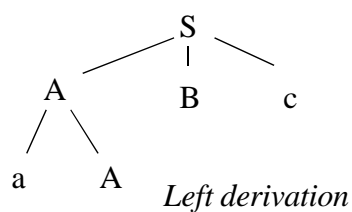


## Left and right derivations

Remember our grammar:

$$S ::= A B c$$
$$A ::= a A \mid b$$
$$B ::= A b \mid a$$

How do we parse the string *abbbc*?



## Languages and Grammars

- A grammar is just a way of describing a language.
- There are actually an infinite number of grammars for a particular language.
- 2 grammars are equivalent if they describe the same language.
  - This becomes extremely important when parsing top-down.
  - Most programming language manuals contain a grammar in BNF, which we may modify to fit our parsing method better.

## Ambiguous grammars

- While there may be an infinite number of grammars that describe a given language, their parse trees may be very different.
- A grammar capable of producing two different parse trees for the same sentence is called **ambiguous**. Ambiguous grammars are highly undesirable.

## Is it IF-THEN or IF-THEN-ELSE?

The IF-THEN=ELSE ambiguity is a classical example of an ambiguous grammar.

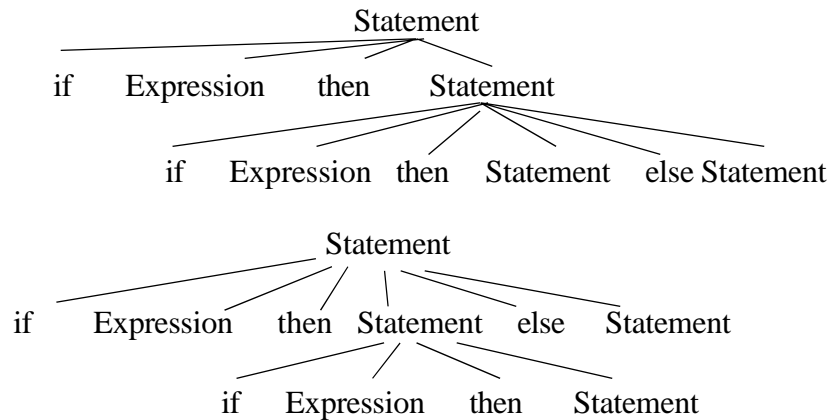
*Statement ::=       if Expression then Statement else Statement  
                      | if Expression then Statement*

How would you parse the following string?

```
IF x > 0
  THEN IF y > 0
        THEN z := x + y
        ELSE z := x;
```

## Is it IF-THEN or IF-THEN-ELSE? (continued)

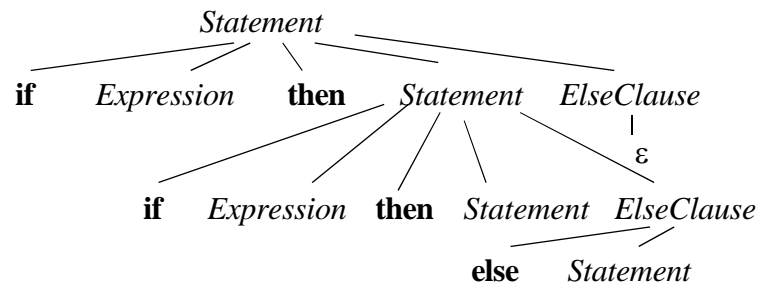
There are two possible parse trees:



## Is it IF-THEN or IF-THEN-ELSE? (continued)

*Statement* ::= **if** *Expression* **then** *Statement* *ElseClause*

*ElseClause* ::= **else** *Statement* /  $\epsilon$





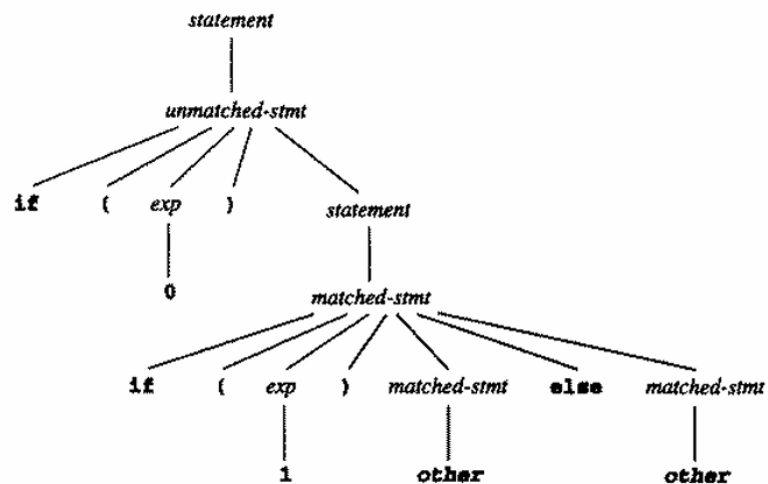
## The Dangling Else Problem

Using disambiguating rule called as **most closely nested rule**.

A solution to the dangling else ambiguity in the BNF itself is more difficult than the previous ambiguities we have seen. A solution is as follows:

```
statement → matched-stmt | unmatched-stmt
matched-stmt → if ( exp ) matched-stmt else matched-stmt | other
unmatched-stmt → if ( exp ) statement
                | if ( exp ) matched-stmt else unmatched-stmt
exp → 0 | 1
```

This works by permitting only a *matched-stmt* to come before an **else** in an if-statement, thus forcing all else-parts to be matched as soon as possible. For instance, the associated parse tree for our sample string now becomes



which indeed associates the else-part with the second if-statement.

## Operator Precedence

Most programming languages have an order of precedence for operators. It would be helpful if this could be encoded into the language's grammar

E. g., let's take a look at the order of precedence in Pascal:

Highest	<i>Unary +, Unary - , NOT</i> <i>*, /, DIV, MOD, AND</i> <i>+, -, OR</i>
Lowest	<i>=, &lt; &gt;, &gt; =, &lt; =, &gt;, &lt;</i>

## Operator Precedence (continued)

This can be encoded in our grammar by considering first a production for our highest level of precedence:

$$\textbf{Factor} ::= \textbf{Unary-operator Unary-Factor} \\ \textbf{/ Unary-Factor}$$

Let's now consider the next-highest level:

$$\textbf{Term} ::= \textbf{Term Multiplicative-operator Factor} \\ \textbf{/ Factor}$$

### Operator Precedence (continued)

Now let's consider the next-level:

***Expr. ::= Expr. Add.-op Term / Term***

And finally,

***Rel.-Expr. ::= Rel.-Expr Rel.-op Expr. / Expr.***

Once we add the production

***Unary-Factor ::= Identifier / Constant / (Rel.Expr.)***

we have a complete expression for Pascal.

### Operator Precedence (continued)

In general, we can start from the lowest order of precedence and work our way to highest in this fashion

$\text{ExprA} ::= \text{ExprA opA ExprB} \mid \text{Exprb}$

$\text{ExprB} ::= \text{ExprB opB ExprC} \mid \text{ExprC}$

.....

$\text{ExprZ} ::= \text{Identifier} \mid \text{Const} \mid \text{.....}$

### Operator Precedence and Associativity

$exp \rightarrow exp \text{ addop } term \mid term$

$addop \rightarrow + \mid -$

$term \rightarrow term \text{ mulop } factor \mid factor$

$mulop \rightarrow *$

$factor \rightarrow ( exp ) \mid \text{number}$

In the above grammar the operators are set to be left associative and higher precedence operators are set on the lower level.

To do: Add an operator like ^ (power of) in the existing expression grammar which is right associative and has highest precedence among +, - and \*.

### DEFINING CFG FOR RELATIONAL EXPRESSIONS

$A-E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow id \mid ( A-E )$

$R-E \rightarrow A-E \quad R-OPR \quad A-E$

$R-OPR \rightarrow < \mid > \mid <= \mid >= \mid == \mid !=$

$2 < 8 \rightarrow \text{valid}$

$a >= h \rightarrow \text{valid}$

$(a+b) <= (c+d) \rightarrow \text{valid}$

CFG for LOGICAL EXPRESSION:

Logical Operators → and, or, not

Binary → and, or

Unary → not

$L-E \rightarrow A-E \text{ LOPR-A } A-E \mid L-OPRB \ A-E$

$L-OPRA \rightarrow \text{and} \mid \text{or}$

$L-OPRB \rightarrow \text{not}$

$(A\&\&B) \mid \mid (! B) \rightarrow \text{valid expression}$

$(a < b) \&\& (a < c) \rightarrow \text{valid}$

$L-E \rightarrow E \&\& E \mid E \text{ " " " " } E \mid ! E$

$E \rightarrow A-E \mid R-E \mid L-E$

$((a > b \mid \mid a < c) \&\& a < d)$

$! ( a > b \&\& c < d)$

$!!! a \quad !! (a < b)$

### Expression grammar with compound statement

IF  $\rightarrow$  if (EXPR) STMT ELSE-PART

ELSE-PART  $\rightarrow \epsilon \mid$  else STMT

EXPR  $\rightarrow 0 \mid 1$

STMT  $\rightarrow$  SIMPLE  $\mid$  COMPUND

SIMPLE  $\rightarrow$  ASSIGN-S  $\mid$  IF  $\mid$  FOR  $\mid$  WHILE  $\mid$  Do-WHILE  
 $\mid$  SWITCH  $\mid$  RETURN  $\mid$  GOTO ...

COMPUND  $\rightarrow \{ \text{STMTS} \}$

STMTS  $\rightarrow \epsilon \mid$  STMT STMTS

CFG of FOR:

FOR  $\rightarrow$  for ( EXPR;EXPR ;EXPR ) STMT

CFG of WHILE:

WHILE  $\rightarrow$  while (EXPR) STMT

CFG of DO-WHILE:

DO-WHILE  $\rightarrow$  do COMPOUND while (EXPR) ;

DEFINING FUNCTION CALL:

name( parameters)

sum(a,b);

a = sum(a,b)+h;

cout<<sum(a,b)<<endl;

abc(a+B,)

display()

FN-CALL  $\rightarrow$  identifier ( PARAM-LIST)

PARAM-LIST  $\rightarrow \epsilon \mid$  PARAMETERS

PARAMETERS  $\rightarrow$  PARAM  $\mid$  PARAM , PARAMETERS

PARAM  $\rightarrow$  EXPR

DEFINITION of FUNCTION DEFINITION:

2 parts:

Ret-type name(parameters)  $\rightarrow$  header

{

$\rightarrow$  body

}

int sum (int a, int b)

**CFG for Function definition:**

$\text{FN\_DEF} \rightarrow \text{FN-HDR FN-BODY}$

$\text{FN-HDR} \rightarrow \text{RET-TYPE identifier ( F-PARAM-LIST )}$

$\text{RET-TYPE} \rightarrow \varepsilon \mid \text{void} \mid \text{int} \mid \text{float} \mid \text{double} \mid \text{char} \mid \text{identifier....}$

$\text{F-PARAM-LIST} \rightarrow \varepsilon \mid \text{F-PARAM MORE-F-PARAM}$

$\text{F-PARAM} \rightarrow \text{PARAM-TYPE identifier}$

$\text{PARAM-TYPE} \rightarrow \text{int} \mid \text{float} \mid \text{double} \mid \text{char} \mid \text{bool....}$

$\text{MORE-F-PARAM} \rightarrow \varepsilon \mid \text{F-PARAM MORE-F-PARAM}$

**DEFINITION of FUNCTION DECLARATION:**

$\text{FN-DEC} \rightarrow \text{FN-HDR ;}$

Define C language SWITCH Statement using BNF Notation.  
You may assume that EXPR and STMT are already defined.



## Expression grammar in C

C has 14 levels of precedence, making its expression grammar more complex than that of most other languages:

$Expr ::= Expr, AssnExpr \mid AssnExpr$

$AssnExpr ::= CondExpr AssnOp AssnExpr \mid CondExpr$

$AssnOp ::= = \mid *= \mid /= \mid \% = \mid += \mid -= \mid < < = \mid > > = \mid \& =$   
 $\mid \wedge = \mid !=$

$CondExpr ::= LogORExpr \mid LogORExpr ? Expr : CondExpr$

$LogORExpr ::= LogORExpr \parallel LogANDExpr \mid LogANDExpr$

$LogANDExpr ::= LogANDExpr \&\& InclORExpr \mid InclORExpr$

$InclORExpr ::= InclORExpr \_ ExclORExpr \mid ExclORExpr$

Derivation  
Separator

C operator

## Expression grammar in C (continued)

$ExclORExpr ::= ExclORExpr \wedge ANDExpr \mid ANDExpr$

$ANDExpr ::= ANDExpr \& EQExpr \mid EQExpr$

$EQExpr ::= EQExpr == RelExpr \mid EQExpr != RelExpr \mid RelExpr$

$RelExpr ::= RelExpr >= ShiftExpr \mid RelExpr <= ShiftExpr$

$\mid RelExpr > ShiftExpr \mid RelExpr < ShiftExpr \mid ShiftExpr$

$ShiftExpr ::= ShiftExpr >> AddExpr \mid ShiftExpr << AddExpr$

$\mid ShiftExpr$

$AddExpr ::= AddExpr + MultExpr \mid AddExpr - MultExpr$

$\mid MultExpr$

$MultExpr ::= MultExpr * CastExpr \mid MultExpr / CastExpr$

$\mid MultExpr \% CastExpr \mid CastExpr$

## Expression grammar in C (continued)

*CastExpr* ::= **(typename)** *CastExpr* | *UnExpr*

*UnExpr* ::= *PostExpr* | ++*UnExpr* | --*UnExpr*

          | *UnOp* *CastExpr* | **sizeof** *UnExpr* | **sizeof**(**typename**)

*UnOp* ::= **&** | \* | + | - | ~ | !

*ExprList* ::= *ExprList*, *AssnExpr* | *AssnExpr*

*PostExpr* ::= *PrimExpr* | *PostExpr*[*Expr*] | *PosrExpr*(*ExprList*)

          | *PostExpr*.**id** | **Post Expr -> id** | *PostExpr* ++

          | *PostExpr* --

*PrimExpr* ::= *Literal* | (*Expr*) | **id**

*Literal* ::= **integer-constant** | **char-constant** | **float-constant**

          | **string-constant**