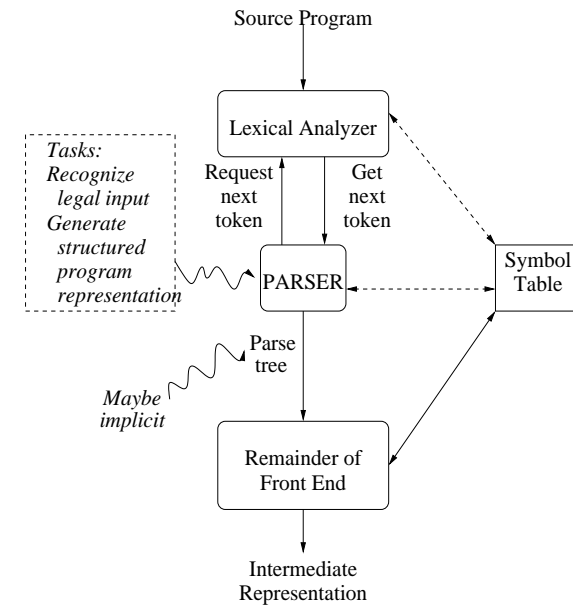


# CS321 Languages and Compiler Design I

## Winter 2012

### Lecture 6

## SYNTAX ANALYSIS (PARSING)



1

PSU CS321 W'12 LECTURE 6 © 1992–2012 ANDREW TOLMACH

2

## SYNTAX ANALYSIS

Specify legal program formats using **context-free grammar (CFG)**

- Use **Backus-Naur Form (BNF)** as notation.
- Gives precise, readable specification.
- Many CFG's have efficient **parsers**.
- Parser **recognizes** syntactically legal programs and **rejects** illegal ones.

Successful parse also captures **hierarchical** structure of programs (expressions, blocks, etc.).

- Convenient representation for further semantic checking (e.g., types) and for code generation.

We can use another program generator (e.g., yacc or CUP) to generate parser automatically from grammar.

## GRAMMARS

A **Context-free Grammar** is described by a set of **productions** (also called **rewrite rules**), e.g.:

$$stmt \rightarrow \text{if } expr \text{ then } stmt \text{ else } stmt$$

$$expr \rightarrow expr + expr \mid expr * expr \mid (expr) \mid -expr \mid id$$

Grammars contain **terminals** ( $\equiv$  **tokens**) (e.g. if,+,id) and **non-terminals** (e.g., *expr,stmt*).

Grammars have a distinguished **start symbol** (ordinarily listed first, e.g., *stmt*).

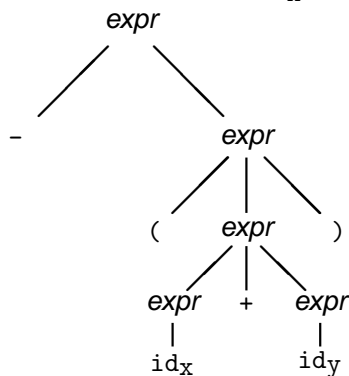
The language **generated** by a grammar is the set of **sentences** (strings) of terminals that can be **derived** by repeated application of productions, beginning with **start symbol**.

- We write  $L(G)$  for the language generated by grammar  $G$ .

## PARSE TREES

A **parse tree** is the graphical representation of a derivation.

Example tree for derivation of sentence  $-(id_x + id_y)$ :



Each application of a production corresponds to an **internal node**, labeled with a **non-terminal**.

Leaves are labeled with **terminals**, possibly with attributes.

The derived sentence is found by reading leaves left-to-right.

## DERIVATIONS

Can “linearize” a parse tree into a sequence of one-step derivations.

Example:

$$\begin{aligned}
 expr &\Rightarrow -\ expr \\
 &\Rightarrow -(expr) \\
 &\Rightarrow -(expr + expr) \\
 &\Rightarrow -(id_x + expr) \\
 &\Rightarrow -(id_x + id_y)
 \end{aligned}$$

or

$$expr \xRightarrow{*} -(id_x + id_y)$$

Here  $\Rightarrow$  is pronounced “derives” and  $\xRightarrow{*}$  is pronounced “derives in zero or more steps.”

This example gives a **leftmost derivation**, i.e., at each step, the leftmost non-terminal is replaced.

## DERIVATIONS (2)

We can define **rightmost derivation** analogously:

$$\begin{aligned}
 expr &\Rightarrow -\ expr \\
 &\Rightarrow -(expr) \\
 &\Rightarrow -(expr + expr) \\
 &\Rightarrow -(expr + id_y) \\
 &\Rightarrow -(id_x + id_y)
 \end{aligned}$$

Every parse tree has a **unique** leftmost derivation and a **unique** rightmost derivation.

(These are usually, but not necessarily, different.)

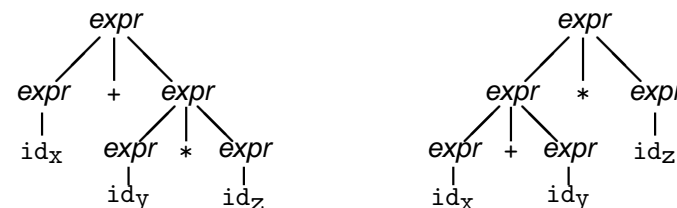
## AMBIGUITY

**BUT** a given **sentence** in  $L(G)$  can have more than one parse tree. Grammars  $G$  for which this is true are called **ambiguous**.

Example: with our grammar, the sentence

$id_x + id_y * id_z$

has two parse trees:



We might think that the left tree is the “correct” one, but nothing in the grammar says this.

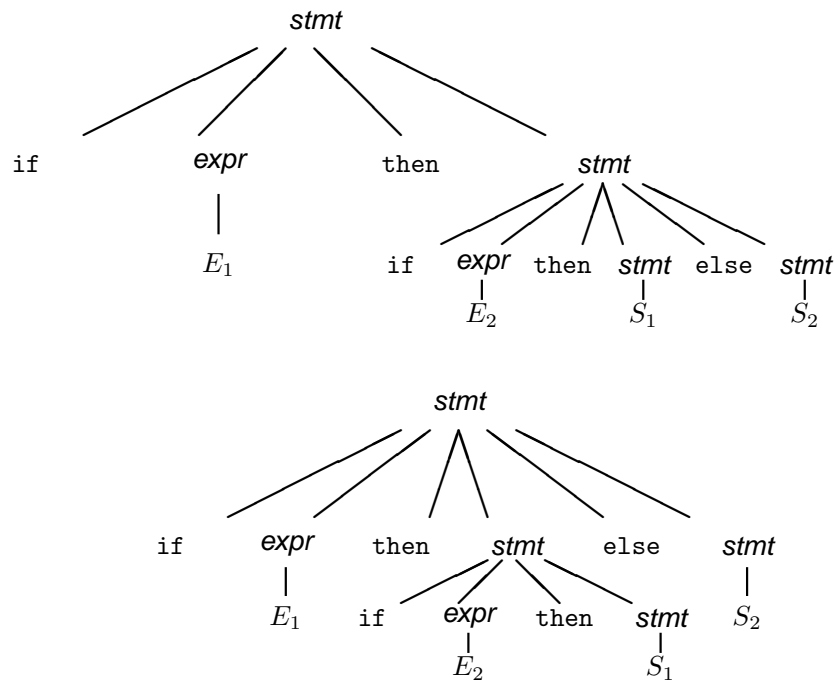
## RESOLVING AMBIGUITIES

Ambiguous grammars can be a significant problem in practice, because we rely on the parse tree to capture the basic structure of a program.

To avoid the problems of ambiguity, we can try to:

- Rewrite grammar
- Use “disambiguating rules” when we implement parser for grammar.

When using these techniques, we must be careful not to change the language generated by our grammar. And often, we must do extra work to make sure that the (unique) parse trees we end up with are the ones we wanted.



## A CLASSIC AMBIGUITY: THE “DANGLING ELSE”

Suppose we want *else* clauses to be optional in *if* statements. Here's a possible grammar:

$$\begin{aligned} \textit{stmt} &\rightarrow \textit{if expr then stmt} \\ &| \textit{if expr then stmt else stmt} | \dots \end{aligned}$$

Given this grammar, a statement of the form

*if E*<sub>1</sub> *then if E*<sub>2</sub> *then S*<sub>1</sub> *else S*<sub>2</sub>

has two possible parse trees...

## RESOLVING AMBIGUITY BY REWRITING GRAMMAR

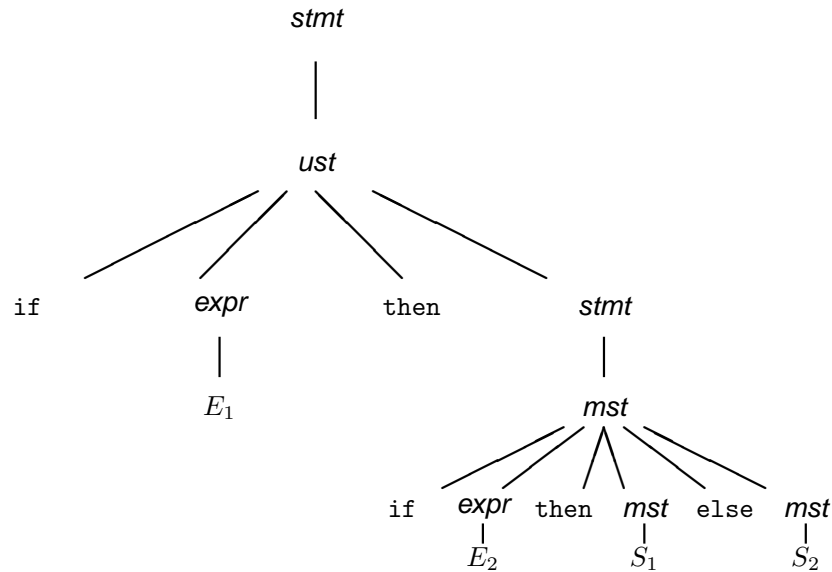
For most languages, we want the first tree (*else* goes with most recent *then*), but grammar is ambiguous.

Solution: rewrite grammar using new non-terminals *mst* (“matched statement”) and *ust* (“unmatched statement”).

$$\begin{aligned} \textit{stmt} &\rightarrow \textit{mst} | \textit{ust} \\ \textit{mst} &\rightarrow \textit{if expr then mst else mst} \\ &| \dots \\ \textit{ust} &\rightarrow \textit{if expr then stmt} \\ &| \textit{if expr then mst else ust} \end{aligned}$$

Now only one parse is possible.

Assuming *S*<sub>1</sub>, *S*<sub>2</sub> are not unmatched *if* statements, we get...



A grammar such as

$$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \\ E \uparrow E \mid (E) \mid - E \mid \text{id}$$

is ambiguous about order of operations.

Want to define

- **Precedence** - which operation is done first (“binds more tightly”) ?

- **Associativity** - is

$$X \text{ op}_1 Y \text{ op}_2 Z$$

equivalent to

$$(X \text{ op}_1 Y) \text{ op}_2 Z \quad (\text{left-associativity})$$

or to

$$X \text{ op}_1 (Y \text{ op}_2 Z) \quad (\text{right-associativity})$$

assuming  $\text{op}_1$  and  $\text{op}_2$  have same precedence?

## STANDARD PRECEDENCE AND ASSOCIATIVITY

The “usual” rules (based on common usage in written math) give the following precedences, highest first:

- (unary minus)
- $\uparrow$  (exponentiation)
- \* /
- + -

All the binary operators are left-associative except exponentiation ( $\uparrow$ ).

Note that these rules are just a matter of **convention**; a programming language designer might choose different ones.

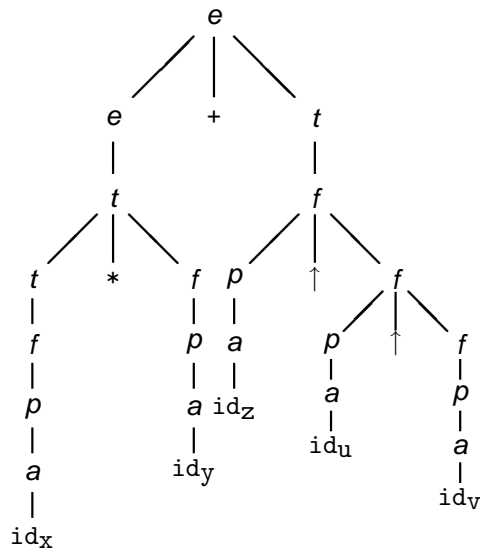
We can handle precedence/associativity information as “side-conditions” to ambiguous grammar when building a parser (by hand or via a parser generator).

## REWRITING ARITHMETIC GRAMMARS

Can build precedence/associativity into grammar using extra non-terminals, each corresponding to a separate level of precedence:

$$\begin{aligned} \text{atom} &\rightarrow (\text{expr}) \mid \text{id} \\ \text{primary} &\rightarrow \text{-primary} \mid \text{atom} \\ \text{factor} &\rightarrow \text{primary} \uparrow \text{factor} \mid \text{primary} \\ \text{term} &\rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor} \\ \text{expr} &\rightarrow \text{expr} + \text{term} \mid \text{expr} - \text{term} \mid \text{term} \end{aligned}$$

Example:  $id_x * id_y + id_z \uparrow id_u \uparrow id_v$



Various semi-standard extensions to BNF are often used in language manuals.

They allow grammar specifications to avoid explicit **recursion** and  $\epsilon$ -productions, by adding **optional** symbols, **repetition**, and **grouping**.

- $[a]$  means  $a \mid \epsilon$
- $\{a\}$  means  $\epsilon \mid a \mid aa \mid \dots$
- $(a \mid b)c$  means  $ac \mid bc$

To avoid confusion between new meta-symbols and terminals, we often enclose the latter in single quotes:

Example:

- $atom \rightarrow '( \textit{expr} )' \mid id$
- $primary \rightarrow -primary \mid atom$
- $factor \rightarrow primary \{ \uparrow primary \}$
- $term \rightarrow \{ factor ( * \mid / ) \} factor$
- $expr \rightarrow \{ term ( + \mid - ) \} term$

### ISSUES WITH EBNF

Must be careful not to **change** the generated language accidentally when going between EBNF and BNF.

Also, note that EBNF grammar given in example is ambiguous because we've lost the associativity information present in the BNF version (even though both generate same language).

Often, a language's grammar will be given in ambiguous EBNF together with separate informal specifications that resolve the ambiguities.

### PROPERTIES OF CFG'S

Any **regular** language can be described by a CFG.

Example:  $(a \mid b)^*abb$

- $A_0 \rightarrow aA_0 \mid bA_0 \mid aA_1$
- $A_1 \rightarrow bA_2$
- $A_2 \rightarrow bA_3$
- $A_3 \rightarrow \epsilon$

**But** we don't use CFG's for lexical analysis, because it's overkill.

Regular expressions are:

- easier to understand
- shorter
- always lead to efficient analyzers

**Any** CFL can be parsed by a computer program, but only **some** CFL's can be parsed **efficiently**.

We'll study both "bottom-up" and "top-down" parsing methods.

## CFG'S CAN'T DO EVERYTHING

Not every language is a CFL.

Example:  $L = \{wcv \mid w \in (a \mid b)^*\}$  is not CF.

- $L$  abstracts idea of variable declaration before use.
- So “semantic” analysis (type-checking) uses additional (mostly ad-hoc) techniques.