

SYNTAX ANALYSIS – Part I

Based on Chapter 4 of Aho, Lam, Sethi, Ullman:

Compilers: Principles, Techniques, & Tools

2nd Ed, Addison Wesley, 2007

Table of Contents

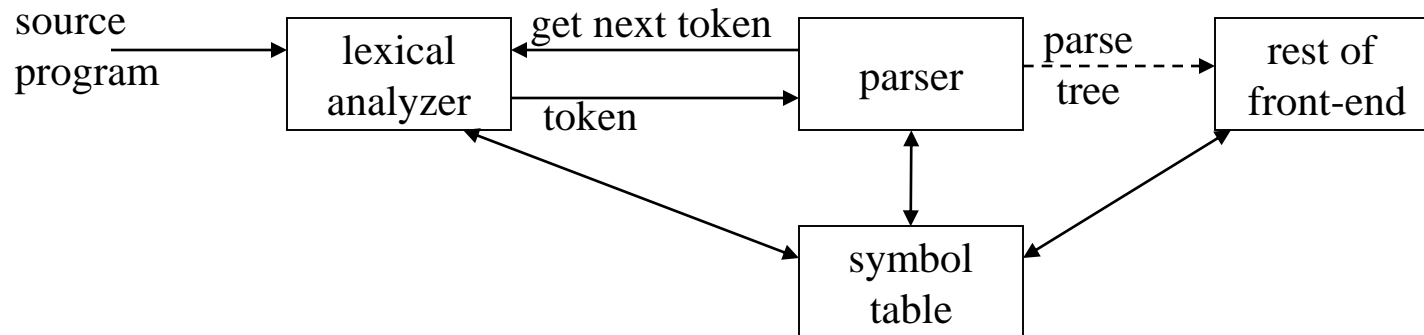
- Introduction
- The Role of Parser
- Context-Free Grammars
- Derivations
- Parse Trees
- Writing a Grammar
- Ambiguity
- Left Recursion Elimination
- Left Factoring

Introduction

- Every programming language has rules that prescribe the syntactic structure of well-formed programs
- The syntax of programming language constructs can be described by **context-free grammars**
- Formal grammars
 - a precise and easy-to-understand syntactic specification
 - automatic construction of an efficient parser
 - syntax-directed translation
 - evolution of programming language
- Parsing
 - generating a parse tree
 - various LL and LR parsing techniques
 - parser generator: yacc(bison), JavaCC, etc.

The Role of Parser

- Interaction of lexical analyzer with parser



- Secondary tasks
 - reporting errors
 - recovery of errors

Context-Free Grammars

- Recursive Structure of High-level Language

$\langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{exp} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

- this form of statement cannot be specified using regular expressions

- Definition

- a CFG G contains $\langle T, N, S, P \rangle$

- T : a set of terminal symbols (tokens)
- N : a set of non-terminal symbols
- S : a distinguished non-terminal symbol that denotes the language defined by the grammar
- P : a set of production rules

$A \rightarrow \alpha \quad (A \text{ in } N, \alpha \text{ in } (T \cup N)^*)$

Context-Free Grammars

- Example: Arithmetic Expressions

$expr \rightarrow expr + term$

$expr \rightarrow expr - term$

$expr \rightarrow term$

$term \rightarrow term * factor$

$term \rightarrow term / factor$

$term \rightarrow factor$

$factor \rightarrow (expr)$

$factor \rightarrow id$

$expr, term, factor$: non-terminal symbols
() - **id** + - * / : terminal symbols
 $expr$: start symbol

- BNF(Backus Naur Form) and Algol60

Context-Free Grammars

- Notational Conventions
 - Terminal symbols
 - lower-case letters early in the alphabet: a, b, c
 - operator symbols: +, -, etc.
 - punctuation symbols: () , etc.
 - digits: 0, 1, ..., 9
 - boldface strings such as **id** or **if**
 - Non-terminal symbols
 - upper-case letters early in the alphabet: A, B, C
 - the letter S (usually start symbol)
 - lower-case italic names such as *expr* or *stmt*
 - Upper-case letter late in the alphabet, X, Y, Z: grammar symbols
 - Lower-case letter late in the alphabet, u, v, ..., z: strings of terminals
 - Lower-case Greek letters, α , β , γ : strings of grammar symbols

Derivations

- A derivation is viewed as the process by which a grammar defines a language
 - top-down construction of a parse tree

- General form (one-step derivation)

$\alpha A \beta \Rightarrow \alpha \gamma \beta$ if $A \rightarrow \gamma$ is a production

$\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ “ α_1 derives α_n ”

- Example

$E \rightarrow E + E \mid E * E \mid (E) \mid - E \mid \mathbf{id}$

$E \Rightarrow - E$

“E derives -E”

$E \Rightarrow - E \Rightarrow - (E) \Rightarrow - (\mathbf{id})$

“a derivation of - (**id**) from E

Derivations

- More notations

$\xRightarrow{*}$: “derives in zero or more steps”

$\xRightarrow{+}$: “derives in one or more steps”

$\alpha \xRightarrow{*} \alpha$

$\alpha \xRightarrow{*} \beta$ & $\beta \Rightarrow \gamma$, then $\alpha \xRightarrow{+} \gamma$

- $L(G)$, language generated by G

$\{w \mid S \xRightarrow{+} w, \text{ } w \text{ is a string of terminals}\}$

- Sentential form

$S \xRightarrow{*} \alpha$: α is a sentential form of G

$E \Rightarrow - E \Rightarrow - (E) \Rightarrow - (E+E \Rightarrow - (id + E) \Rightarrow -(id + id)$

- Rightmost and leftmost derivations

- rightmost(leftmost) non-terminal is replaced at each step

$E \Rightarrow - E \Rightarrow - (E) \Rightarrow - (E+E) \Rightarrow - (E + id) \Rightarrow - (id + id)$

- Leftmost(rightmost) sentential form $S \xRightarrow{*}_{lm} \alpha$

Examples

- Example 1

$$E \rightarrow (E) \mid a$$

$$L(G) = \{a, (a), ((a)), \dots\} = \{ ({}^n a) {}^n \mid n \geq 0 \}$$

$$\text{e.g. } E \Rightarrow (E) \Rightarrow ((E)) \Rightarrow ((a))$$

$$\text{cf. } E \rightarrow (E) \rightarrow \text{infinite recursion}$$

- Example 2

$$E \rightarrow E + a \mid a$$

$$L(G) = \{ a, a + a, a + a + a, \dots \}$$

$$\text{e.g. } E \Rightarrow E + a \Rightarrow E + a + a \Rightarrow E + a + a + a \Rightarrow a + a + a + a$$

Examples

- Example 3: statements

statement \rightarrow *if-stmt* | **other**

if-stmt \rightarrow **if** (*exp*) *statement*
| **if** (*exp*) *statement* **else** *statement*

exp \rightarrow 0 | 1

E.g. **if** (1) **other** **else** **if** (0) **other** **else** **other**

statement \Rightarrow *if-stmt* \Rightarrow **if** (*exp*) *statement* **else** *statement*

\Rightarrow **if** (1) *statement* **else** *statement* \Rightarrow **if** (1) **other** **else** *statement*

\Rightarrow **if** (1) **other** **else** *if-stmt*

\Rightarrow **if** (1) **other** **else** **if** (*exp*) *statement* **else** *statement*

\Rightarrow **if** (1) **other** **else** **if** (0) *statement* **else** *statement*

\Rightarrow **if** (1) **other** **else** **if** (0) **other** **else** *statement*

\Rightarrow **if** (1) **other** **else** **if** (0) **other** **else** **other**

Examples

- Example 4: Balanced parentheses

$A \rightarrow (A)A \mid \varepsilon$

e.g. $A \Rightarrow (A)A \Rightarrow ()A \Rightarrow ()(A)A \Rightarrow ()((A)A)A$
 $\Rightarrow ()((()A)A \Rightarrow ()((())A \Rightarrow ()((()))$

or

$A \rightarrow AB \mid B$

$B \rightarrow (A) \mid ()$

$A \Rightarrow AB \Rightarrow BB \Rightarrow ()B \Rightarrow ()(A) \Rightarrow ()(B) \Rightarrow ()(())$

Examples

- Example 5: statement sequence

$stmt_sequence \rightarrow stmt ; stmt_sequence \mid stmt$

$stmt \rightarrow s$

$L(G) = \{ s, s ; s, s ; s ; s, \dots \}$; as a separator

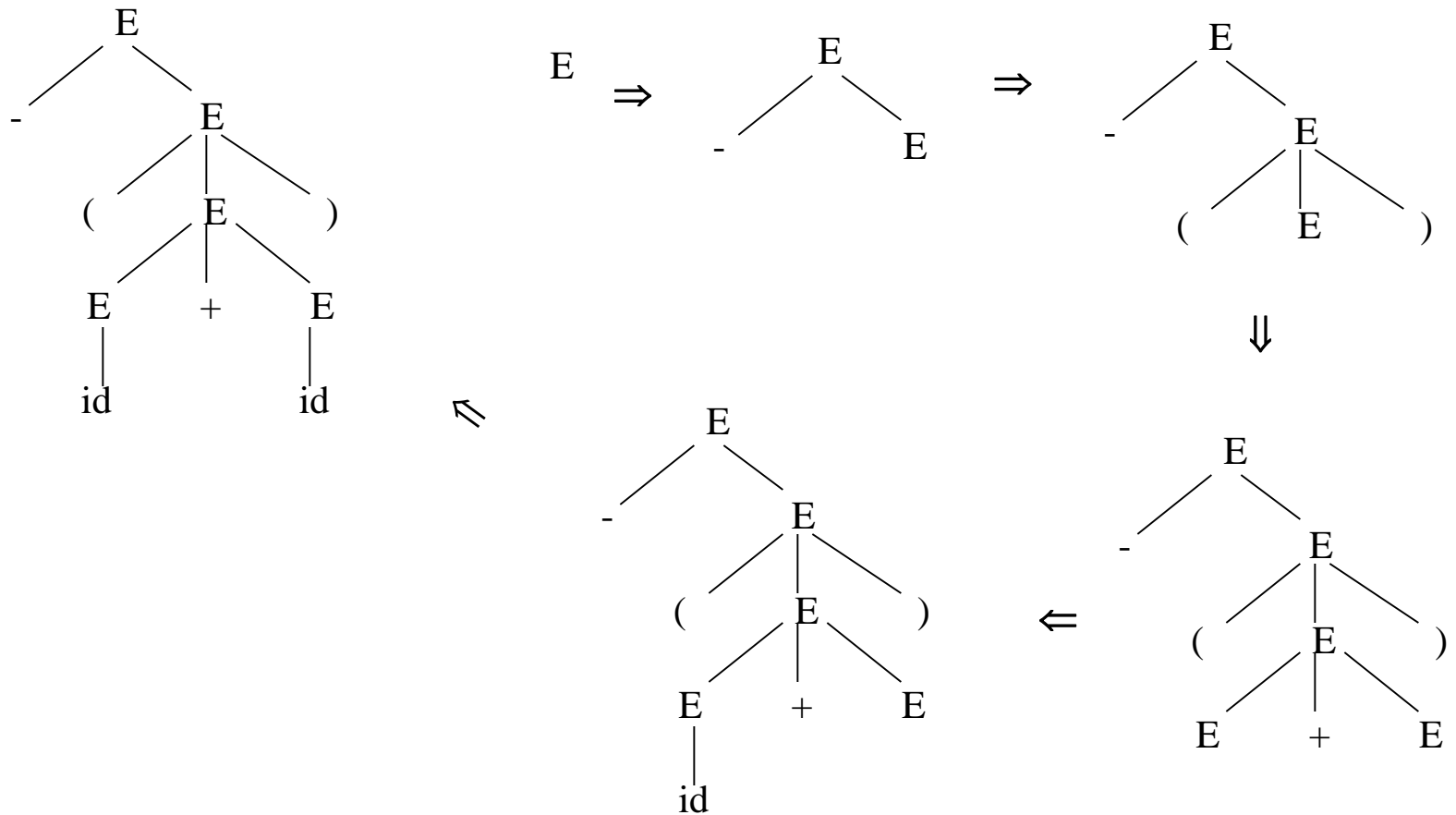
$stmt_sequence \rightarrow stmt ; stmt_sequence \mid \epsilon$

$stmt \rightarrow s$

$L(G) = \{ \epsilon, s ; , s ; s ; , s ; s ; s ; , \dots \}$; as a terminator

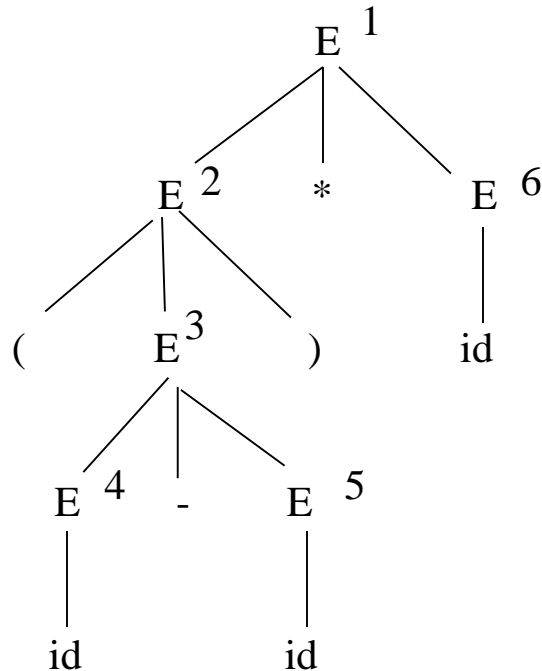
Parse Trees and Derivations

- A parse tree can be viewed as a graphical representation for a derivation without considering replacement order

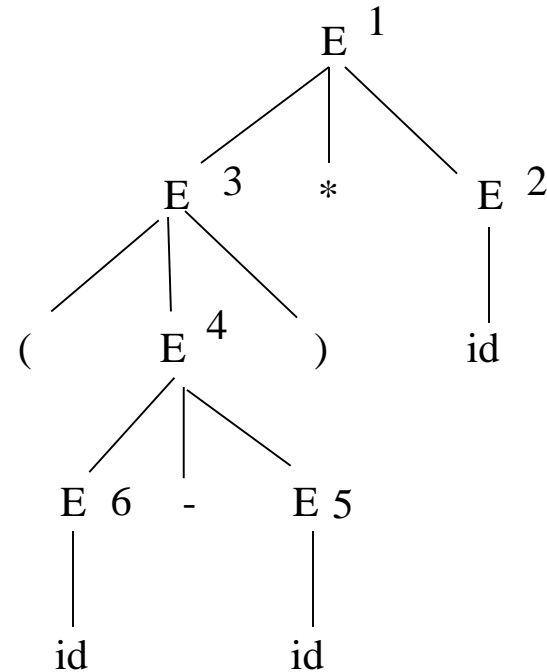


Parse Trees and Derivations

- Derivation and parse tree : $(a - b) * c$



leftmost derivation-
preorder numbering



rightmost derivation –
reverse in postordering

Ambiguity

- A grammar that produces more than one parse tree for some sentence is said to be *ambiguous*
- Example : $\text{id} + \text{id} * \text{id}$

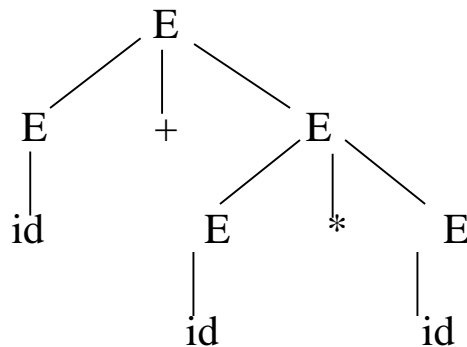
$E \Rightarrow E + E$

$\Rightarrow \text{id} + E$

$\Rightarrow \text{id} + E * E$

$\Rightarrow \text{id} + \text{id} * E$

$\Rightarrow \text{id} + \text{id} * \text{id}$



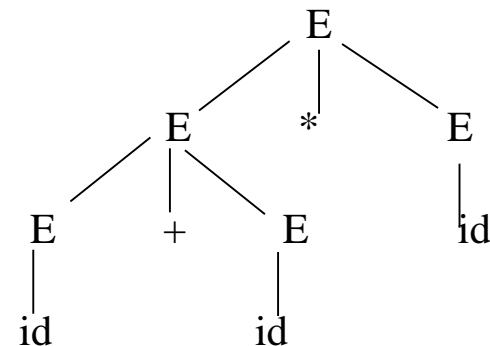
$E \Rightarrow E * E$

$\Rightarrow E + E * E$

$\Rightarrow \text{id} + E * E$

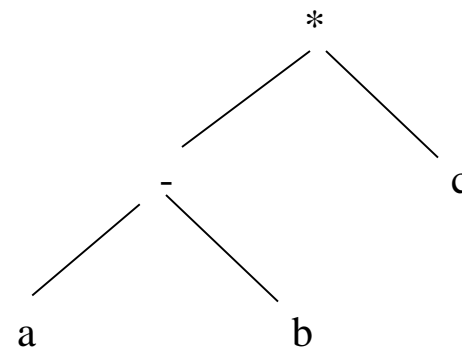
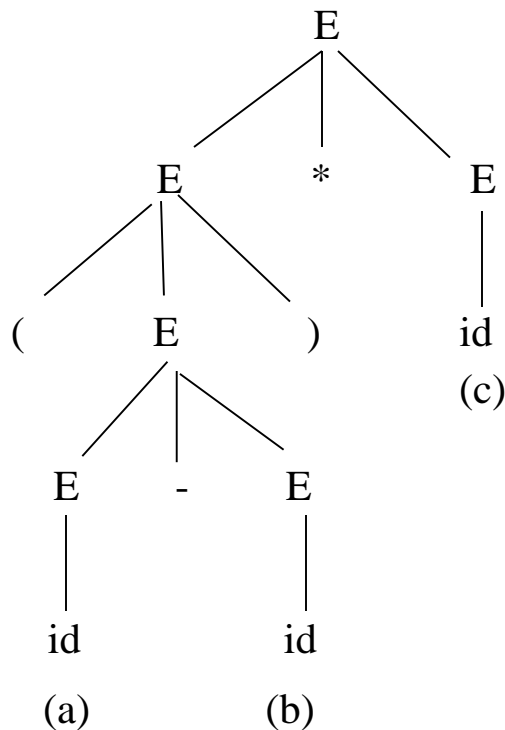
$\Rightarrow \text{id} + \text{id} * E$

$\Rightarrow \text{id} + \text{id} * \text{id}$



Syntax Trees

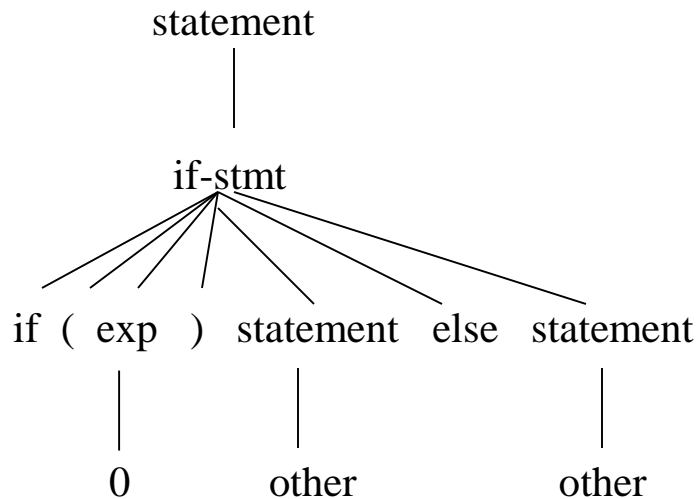
- Parse tree contains much more information than is necessary for a compiler to produce executable code
- Syntax tree(or abstract syntax tree) is an abstraction of the actual source code token sequences



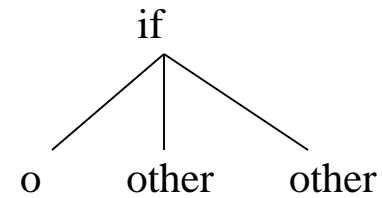
Syntax Tree

Syntax Trees

- Example: if-statement



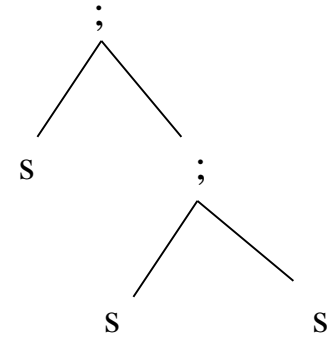
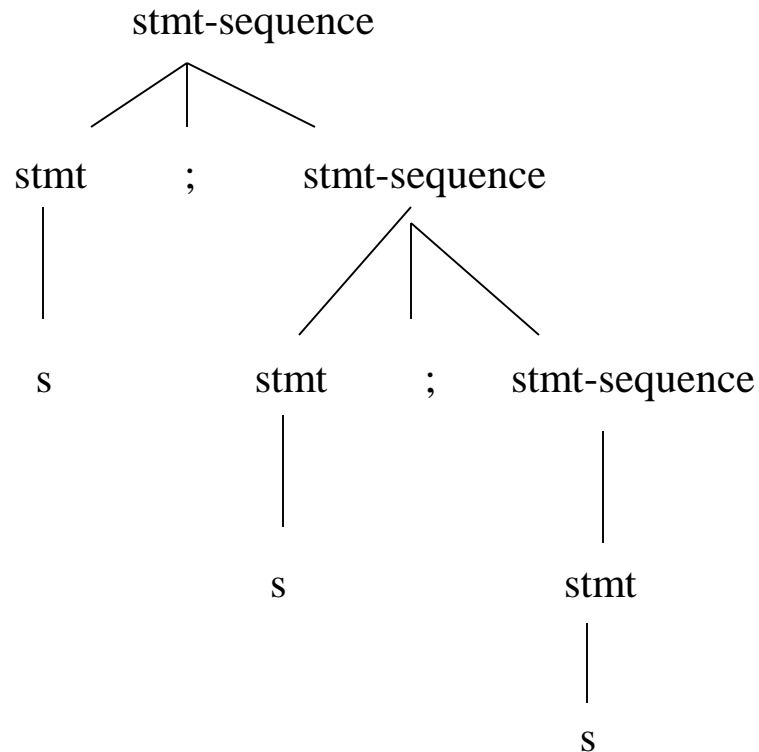
Parse Tree



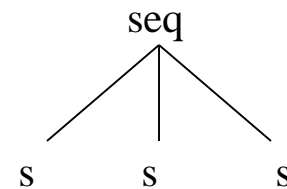
Syntax Tree

Syntax Trees

- Example: Statement Sequence



or



Writing a Grammar

- Grammatical rules of programming languages
 - lexical rule : regular expressions - lexical analyzer
 - syntactic rule: CFG - parser
 - semantic rule(context-sensitive syntax): semantic analyzer
 - e.g. requirement that identifiers should be declared before they are used
- Each parsing method can handle grammars only of a certain form
 - rewriting initial grammars
 - using operator precedence and associativity
 - left-recursion elimination
 - left factoring

Regular Expressions vs CFG

- Regular Grammar : $\langle T, N, S, P \rangle$
 - right regular grammar (right linear grammar)
 - $A \rightarrow aB$
 - $A \rightarrow a$
 - $A \rightarrow \epsilon$
 - left regular grammar (left linear grammar)
 - $A \rightarrow Ba$
 - $A \rightarrow a$
 - $A \rightarrow \epsilon$
- A regular grammar is a CFG
 - but, simple and easy to implement efficiently

Regular Expressions vs CFG

- Every regular expression can be converted into a CFG
 - example : $(a|b)^*abb$
 $A_0 \rightarrow aA_0 \mid bA_0 \mid aA_1$
 $A_1 \rightarrow bA_2$
 $A_2 \rightarrow bA_3$
 $A_3 \rightarrow \varepsilon$
- Converting a NFA to a regular grammar
 - state $i \Rightarrow$ non-terminal A_i
 - transition $j = T(i, a) \Rightarrow$ production $A_i \rightarrow aA_j$
 - transition $j = T(i, \varepsilon) \Rightarrow$ production $A_i \rightarrow A_j$
 - $i \in F \Rightarrow$ production $A_i \rightarrow \varepsilon$
 - $i = \text{start state} \Rightarrow A_i : \text{start symbol}$

Eliminating Ambiguity

- Using Operator Precedence and Associativity

- Example:

- $E \rightarrow E + E$

- $E \rightarrow E - E$

- $E \rightarrow E * E$

- $E \rightarrow (E)$

- $E \rightarrow -E$

- $E \rightarrow \text{id}$

- Grouping operators of equal precedences

- Lower precedence \Rightarrow higher in parse tree

- Left Associativity \Rightarrow left recursion

- $E \rightarrow E \text{ **addop** } T \mid T$

- $T \rightarrow T \text{ **multop** } F \mid F$

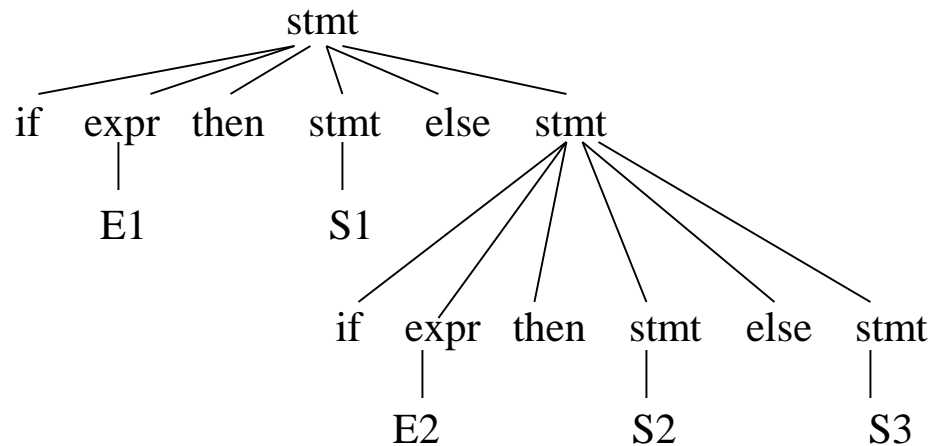
- $F \rightarrow \text{**id**} \mid (E) \mid - E$

Eliminating Ambiguity

- Dangling-else grammar

stmt \rightarrow **if** expr **then** stmt
| **if** expr **then** stmt **else** stmt
| S

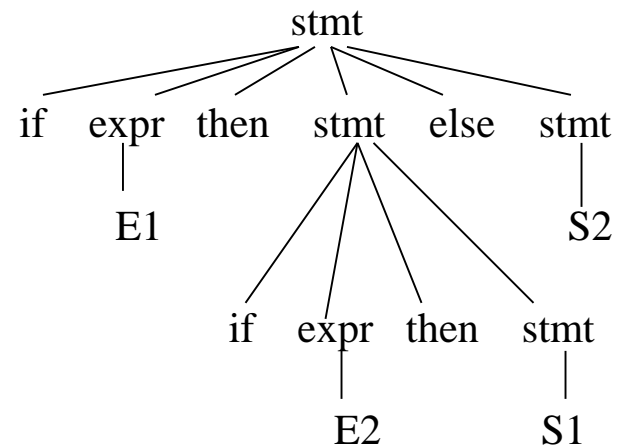
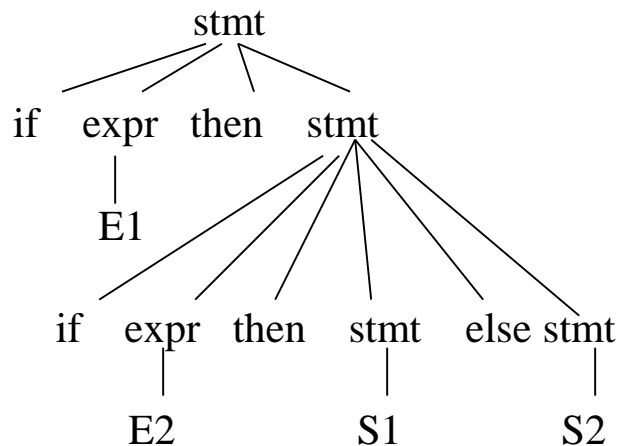
if E1 **then** S1 **else if** E2 **then** S2 **else** S3



Eliminating Ambiguity

- Ambiguous if-stmt

if E1 then if E2 then S1 else S2



Eliminating Ambiguity

- Rewriting grammar

stmt \rightarrow matched_stmt

 | unmatched_stmt

matched_stmt \rightarrow **if** expr **then** matched_stmt **else** matched stmt

 | **other**

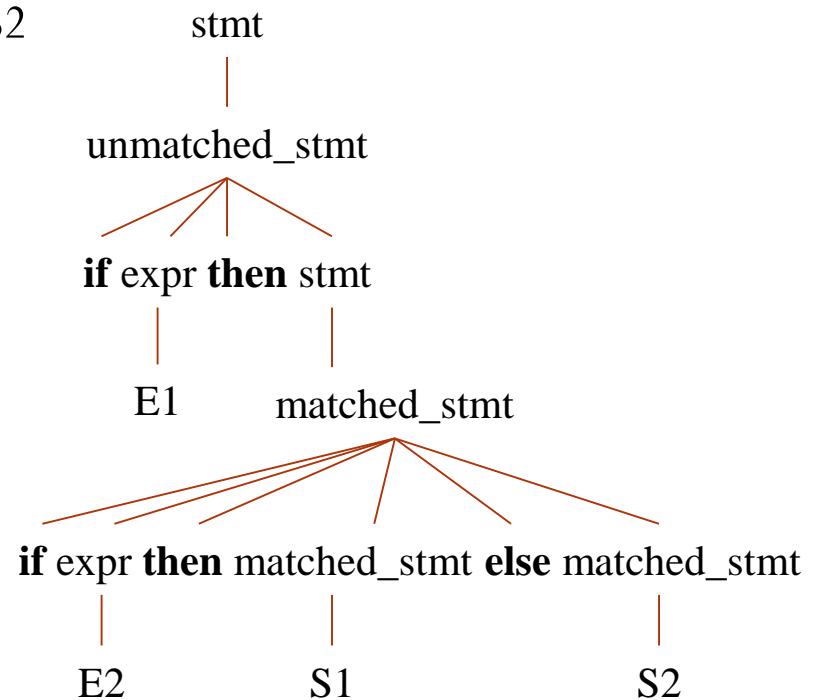
unmatched_stmt \rightarrow **if** expr **then** stmt

 | **if** expr **then** matched_stmt **else** unmatched_stmt

Eliminating Ambiguity

- Rewriting grammar

if E1 then if E2 then S1 else S2



Elimination of Left Recursion

- A grammar is left-recursive if there is $A \Rightarrow A \alpha$
 $\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term}$
- Top-down parsing cannot handle left-recursive grammars
- Eliminating left recursion

$$\begin{aligned} A \rightarrow A a \mid b &\Rightarrow A \rightarrow bA' \\ A' &\rightarrow aA' \mid \epsilon \end{aligned}$$

- Example

$$\begin{array}{ll} E \rightarrow E + T \mid T & \Rightarrow E \rightarrow TE' \\ T \rightarrow T * F \mid F & E' \rightarrow +TE' \mid \epsilon \\ F \rightarrow (E) \mid \text{id} & T \rightarrow FT' \\ & T' \rightarrow *FT' \mid \epsilon \\ & F \rightarrow (E) \mid \text{id} \end{array}$$

Elimination of Left Recursion

- General case (immediate left-recursion)

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

\Rightarrow

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

- It does not eliminate left recursion involving derivations of two or more steps
- Non-immediate left-recursion

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid \epsilon \quad \Rightarrow$$

$$S \Rightarrow Aa \Rightarrow Sda$$

$$S \rightarrow Aa \mid b$$

$$A \rightarrow SdA' \mid A' \quad // \text{ left-recursive}$$

$$A' \rightarrow cA' \mid \epsilon$$

Elimination of Left Recursion

- Eliminating left-recursion (non-immediate recursion)
 - Arrange the nonterminals in some order, A_1, A_2, \dots, A_n
 - for each A_i , ($i = 1$ to n)
 - for $j = 1$ to $i - 1$
 - replace $A_i \rightarrow A_j \gamma$
by $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$,
where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are A_j -productions
 - eliminate the immediate left recursion among the A_i -productions
- Example
 - $S \rightarrow Aa \mid b$
 - $A \rightarrow Ac \mid Sd \mid \epsilon$
 - 1. ordering: $S A$
 - 2. for S , no left-recursion
 - 3. for A , replace $A \rightarrow Sd$ with S -productions
 $A \rightarrow Ac \mid Aad \mid bd \mid \epsilon$
 - 4. elimination of left-recursion
 $A \rightarrow bdA' \mid A'$
 $A' \rightarrow cA' \mid adA' \mid \epsilon$

Left Factoring

- Two productions for a non-terminal A have a common prefix

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$$

\Rightarrow

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

- General case

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \alpha\beta_n \mid \gamma$$

\Rightarrow

$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \beta_n$$

- Example

$$S \rightarrow iEtS \mid iEtSeS \mid a$$

$$E \rightarrow b$$

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$