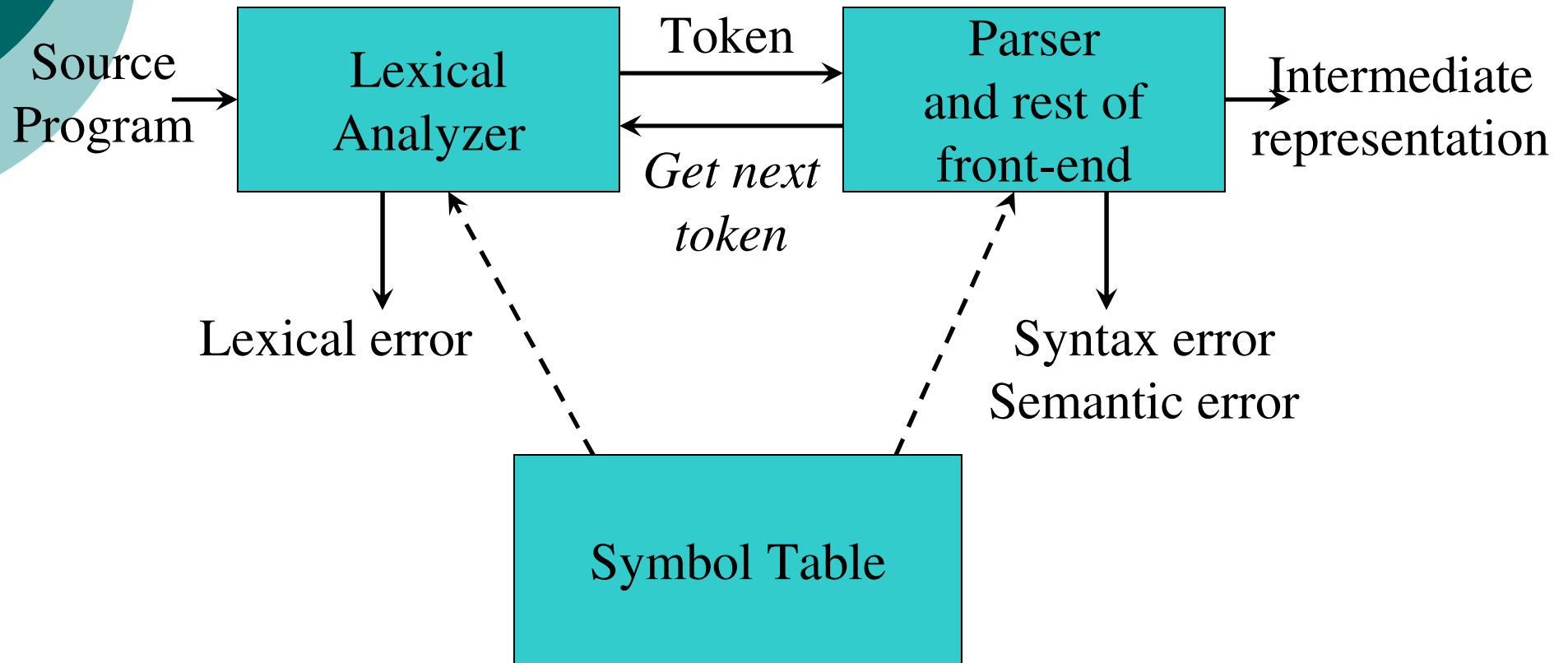




Syntax Analyzer

Introduction





Syntax Analyzer

- Creates the syntactic structure of the given source program.
- This syntactic structure is mostly a *parse tree*.
- Syntax Analyzer is also known as *parser*.
- The syntax of a programming is described by a *context-free grammar (CFG)*. We will use **BNF (Backus-Naur Form)** notation in the description of CFGs.



Syntax Analyzer Cont...

- The syntax analyzer (parser) checks whether a given source program satisfies the rules implied by a context-free grammar or not.
 - If it satisfies, the parser creates the parse tree of that program.
 - Otherwise the parser gives the error messages.
- A context-free grammar
 - gives a precise syntactic specification of a programming language.
 - the design of the grammar is an initial phase of the design of a compiler.
 - a grammar can be directly converted into a parser by some tools.



Parsers

1. **Top-Down Parser**

- the parse tree is created top to bottom, starting from the root.
- LL parsing

2. **Bottom-Up Parser**

- the parse is created bottom to top; starting from the leaves
 - LR parsing
- Both top-down and bottom-up parsers scan the input from left to right (one symbol at a time).



Error Handling

- A good compiler should assist in identifying and locating errors.
- Errors can be
 - *Lexical errors*: such as misspelling an identifier, keyword or operator.
 - *Syntax errors*: such as an arithmetic expression with unbalanced parentheses.
 - *Semantic errors*: such as an operator applied to an incompatible operand.
 - *Logical errors*: such as an infinitely recursive call.



Error Recovery Strategies

- *Panic mode*
 - Discard input until a token in a set of designated synchronizing tokens is found
- *Phrase-level recovery*
 - Perform local correction on the input to repair the error
- *Error productions*
 - Augment grammar with productions for erroneous constructs
- *Global correction*
 - Choose a minimal sequence of changes to obtain a global least-cost correction



Context Free Grammar



Grammars

- Inherently recursive structures of a programming language are defined by a context-free grammar.
- Context-free grammar is a 4-tuple $G = (N, T, P, S)$ where
 - T is a finite set of tokens (*terminal* symbols)
 - N is a finite set of *nonterminals*
 - P is a finite set of *productions* of the form
$$\alpha \rightarrow \beta$$
where $\alpha \in N$ and $\beta \in (N \cup T)^*$
 - $S \in N$ is a designated *start symbol*



Example

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

Notational Conventions Used

- Terminals

$a, b, c, \dots \in T$

specific terminals: **0**, **1**, **id**, **+**

- Nonterminals

$A, B, C, \dots \in N$

specific nonterminals: *expr*, *term*, *stmt*

- Grammar symbols

$X, Y, Z \in (N \cup T)$

- Strings of terminals

$u, v, w, x, y, z \in T^*$

- Strings of grammar symbols

$\alpha, \beta, \gamma \in (N \cup T)^*$



Derivations

$E \Rightarrow E+E$

- $E+E$ derives from E
 - we can replace E by $E+E$
 - to be able to do this, we have to have a production rule $E \rightarrow E+E$ in our grammar.

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

- A sequence of replacements of non-terminal symbols is called a **derivation** of $\text{id}+\text{id}$ from E .

Derivations Cont...

In general a derivation step is

$\alpha A \beta \Rightarrow \alpha \gamma \beta$ if there is a production rule $A \rightarrow \gamma$ in our grammar

where α and β are arbitrary strings of terminal and non-terminal symbols

$\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ (α_n derives from α_1 or α_1 derives α_n)

\Rightarrow : derives in one step

\Rightarrow^* : derives in zero or more steps

\Rightarrow^+ : derives in one or more steps

+



Derivation Example

- If we always choose the left-most non-terminal in each derivation step, this derivation is called as **left-most derivation**.
- If we always choose the right-most non-terminal in each derivation step, this derivation is called as **right-most derivation**.



Left-Most and Right-Most Derivations

Left-Most Derivation

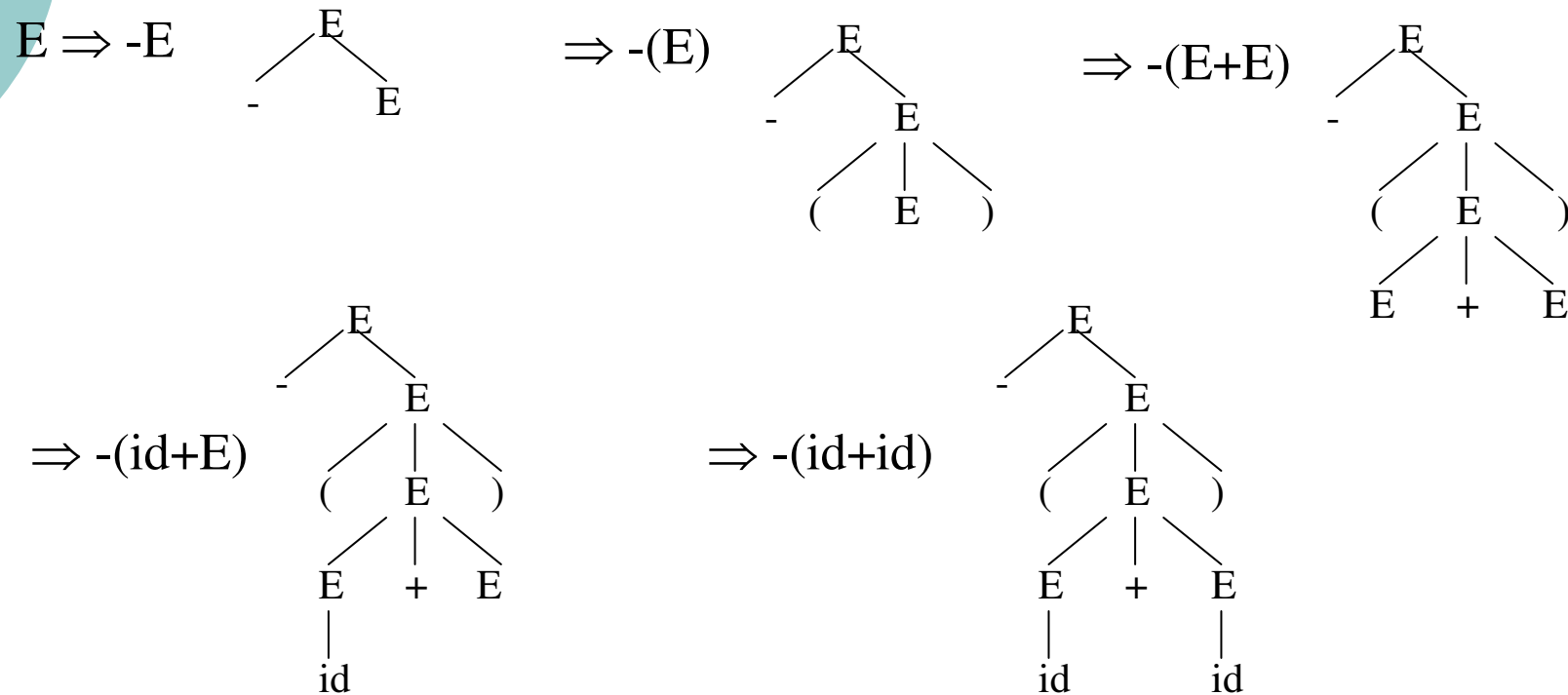
$$\begin{array}{ccccccc} \text{lm} & & \text{lm} & & \text{lm} & & \text{lm} \\ E & \Rightarrow & -E & \Rightarrow & -(E) & \Rightarrow & -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id) \end{array}$$

Right-Most Derivation

$$\begin{array}{ccccccc} \text{rm} & E & \xRightarrow{\text{rm}} & -E & \Rightarrow & -(E) & \Rightarrow -(E+E) \xRightarrow{\text{rm}} -(E+id) \xRightarrow{\text{rm}} -(id+id) \end{array}$$

Parse Tree

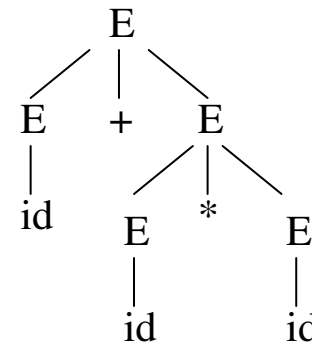
- Inner nodes of a parse tree are non-terminal symbols.
 - The leaves of a parse tree are terminal symbols.
-
- A parse tree can be seen as a graphical representation of a derivation.



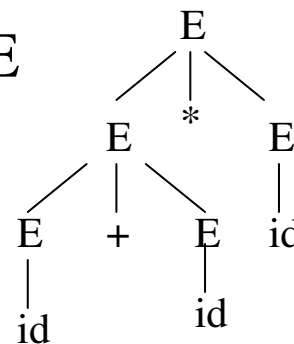
Ambiguity

- A grammar produces more than one parse tree for a sentence is called as an *ambiguous* grammar.

$E \Rightarrow E + E \Rightarrow id + E \Rightarrow id + E * E$
 $\Rightarrow id + id * E \Rightarrow id + id * id$



$E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow id + E * E$
 $\Rightarrow id + id * E \Rightarrow id + id * id$





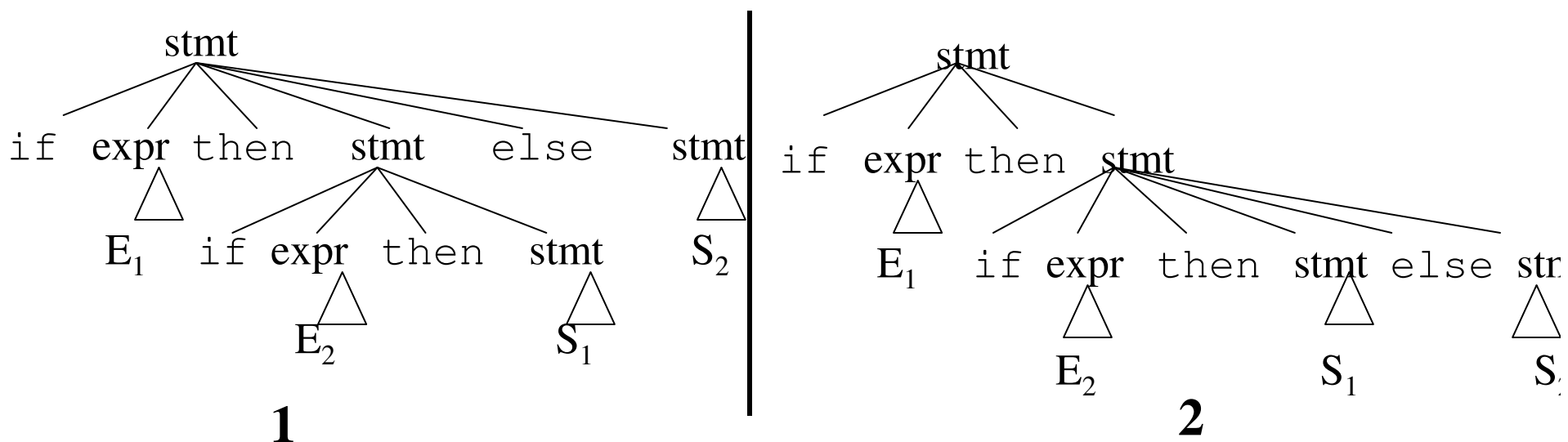
Ambiguity cont...

- For the most parsers, the grammar must be unambiguous.
- unambiguous grammar
 - ➔ unique selection of the parse tree for a sentence
- We should eliminate the ambiguity in the grammar during the design phase of the compiler.
- An unambiguous grammar should be written to eliminate the ambiguity.
- We have to prefer one of the parse trees of a sentence (generated by an ambiguous grammar) to disambiguate that grammar to restrict to this choice.

Ambiguity cont...

`stmt` \rightarrow `if` `expr` `then` `stmt` `|`
`if` `expr` `then` `stmt` `else` `stmt` `|` `otherstmts`

`if` E_1 `then` `if` E_2 `then` S_1 `else` S_2





Ambiguity cont...

- We prefer the second parse tree (else matches with closest if).
- So, we have to disambiguate our grammar to reflect this choice.
- The unambiguous grammar will be:

`stmt` \rightarrow `matchedstmt` | `unmatchedstmt`

`matchedstmt` \rightarrow `if` `expr` `then` `matchedstmt` `else` `matchedstmt` `lotherstmts`

`unmatchedstmt` \rightarrow `if` `expr` `then` `stmt` |
 `if` `expr` `then` `matchedstmt` `else` `unmatchedstmt`