

Parsing/Syntactic Analysis

CIS*4650 (Winter 2020)

Review

- Compiler: translates a source program into target machine code
 - Front-End Analysis: language-specific, focused on analysis, systematic solutions with tools
 - Back-End Synthesis: machine-specific, focused on synthesis, and ad hoc solutions through coding
- Front-End Analysis:
 - Scanning/Lexical Analysis: break an input stream into a token sequence
 - Parsing/Syntactic Analysis: parse the phrase structure of a program
 - Semantic Analysis: calculate the meaning and generate the intermediate code of a program

What is a Parser?

- Validate the ordering of tokens and determine the phrase structure (syntax tree):



Limits of Regular Expressions

- Languages with the pattern “ $a^n b^n$ ” are not regular languages, since they can’t be expressed by a finite number of states.

digits = $[0-9]^+$

sum = (digits “+”)* digits

e.g., $28 + 301 + 9$

digits = $[0-9]^+$

sum = expr “+” expr

expr = “(” sum “)” | digits

e.g., $(109 + 23)$, 61 , $(1 + (250 + 3))$

expr = “(” expr “+” expr “)” | digits

expr = “(” (“(” expr “+” expr “)” | digits) “+” expr “)” | digits

...

Towards a Context-Free Grammar

- Allow recursive rules:

$\text{sum} = \text{expr} \text{ "+" } \text{expr}$
 $\text{expr} = \text{"(" sum "}" \mid \text{digits}$

- Remove alternations:

$\text{expr} = a b (c \mid d) e$
 $\text{aux} = c \mid d$
 $\text{expr} = a b \text{ aux } e$

$\text{aux} = c$
 $\text{aux} = d$
 $\text{expr} = a b \text{ aux } e$

- Remove Kleene closures:

$\text{expr} = (a b c)^*$

$\text{expr} = a b c \text{ expr}$
 $\text{expr} = \varepsilon$

Context-Free Grammars

- Terminals: atomic symbols/tokens for a language (shown in lower cases)
- Non-terminals: symbols/variables that specify the phrases of the language (shown in UPPER CASES)
 - Non-terminals can be described recursively.
- Productions: rules relating variables to their definitions
 - Syntax: symbol \rightarrow symbol symbol ... symbol
 - Terminals can't appear on the left-hand-side (LHS)
- Start symbol (S): a special symbol that starts all derivations

Example CFG

- Arithmetic expressions with $+$ and $-$ operators, id tokens, and balanced parentheses:

$S \rightarrow \text{EXPR}$

$\text{EXPR} \rightarrow \text{EXPR} + \text{EXPR}$

$\text{EXPR} \rightarrow \text{EXPR} - \text{EXPR}$

$\text{EXPR} \rightarrow (\text{EXPR})$

$\text{EXPR} \rightarrow \text{id}$

- Why $+$, $-$, $($, and $)$ need not be quoted?

Derivations

- Derivation: begin with the start symbol and repeatedly replace a non-terminal with one of its RHS's.
- Sentence: a derivation becomes a sentence if every symbol in it is a terminal.

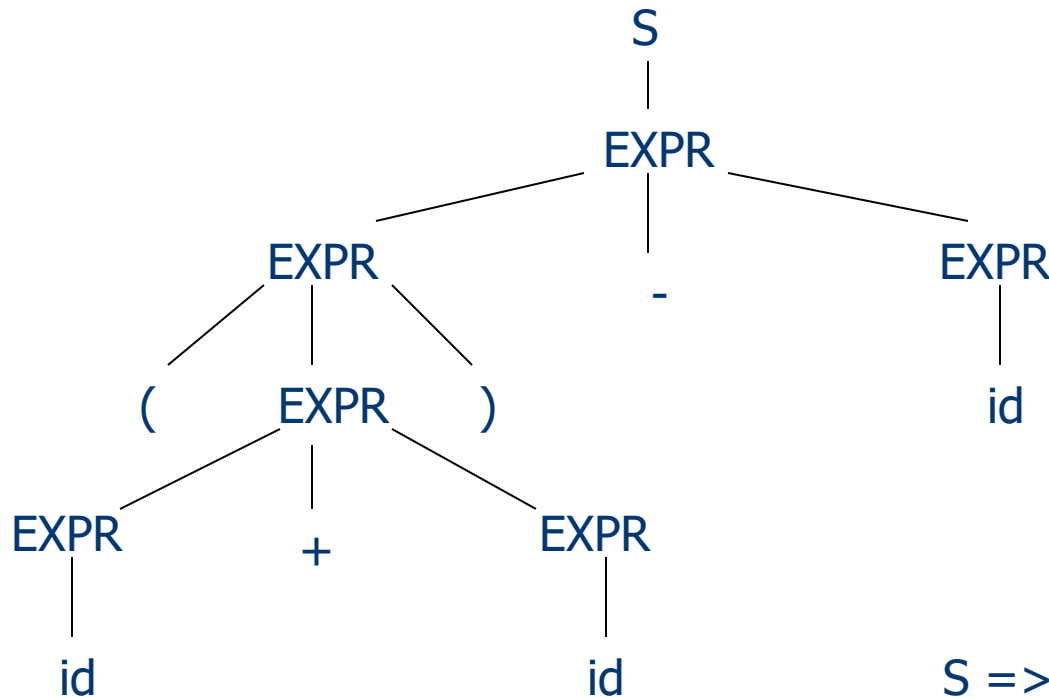
$S \Rightarrow \text{EXPR}$
 $\Rightarrow \text{EXPR} + \text{EXPR}$
 $\Rightarrow \text{EXPR} + \text{id}$
 $\Rightarrow \text{id} + \text{id}$

$S \Rightarrow \text{EXPR}$
 $\Rightarrow \text{EXPR} - \text{EXPR}$
 $\Rightarrow (\text{EXPR}) - \text{EXPR}$
 $\Rightarrow (\text{EXPR}) - \text{id}$
 $\Rightarrow (\text{EXPR} + \text{EXPR}) - \text{id}$
 $\Rightarrow (\text{EXPR} + \text{id}) - \text{id}$
 $\Rightarrow (\text{id} + \text{id}) - \text{id}$

Parse Trees

- Graphically show a derivation with LHS connected to its RHS components:

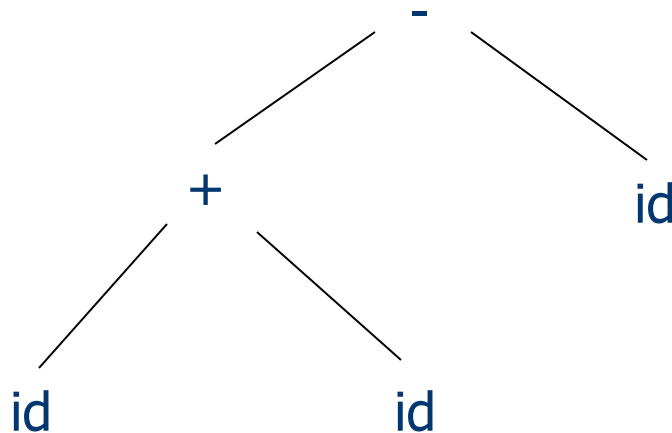
➤ Several derivations can have the same parse tree.



$S \Rightarrow^* (id + id) - id$

Abstract Syntax Trees

- Only capture the information needed for semantic analysis and code generation:
 - Also called syntax trees for short



Derivation Sequence

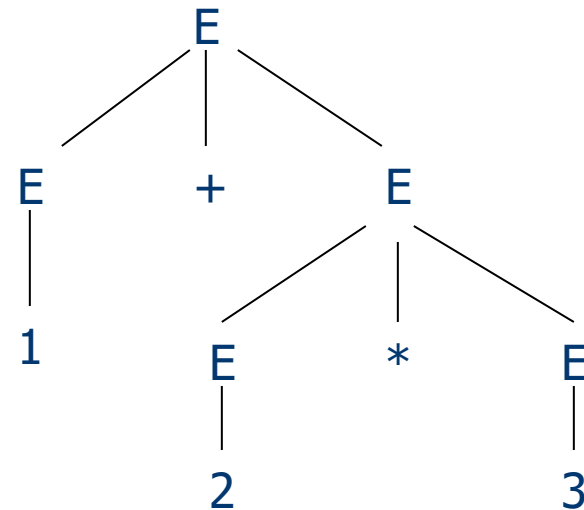
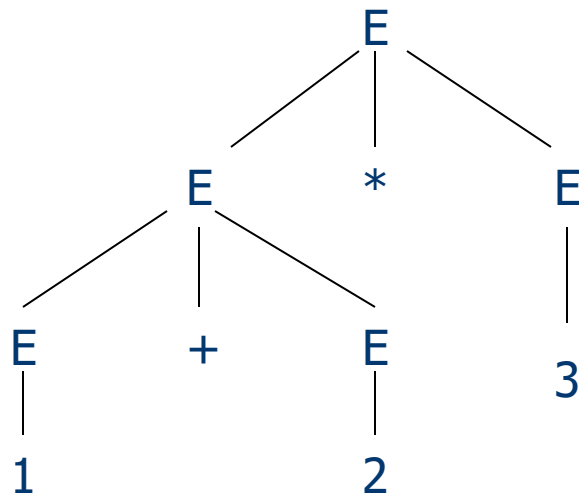
- Many different possible derivations of the same sentence
 - If more than one non-terminal appears on the LHS of a production, we can choose which to expand next.
- Two obvious conventions:
 - Leftmost derivation:
 - Choose leftmost non-terminal to expand
 - Top-down parsing process
 - Rightmost derivation:
 - Choose rightmost non-terminal to expand
 - Bottom-up parsing process

Ambiguous Grammars

- A grammar is ambiguous if we can derive a sentence with two or more different parse trees.

$E \rightarrow \text{id}$
 $E \rightarrow \text{num}$
 $E \rightarrow E * E$
 $E \rightarrow E / E$

$E \rightarrow E + E$
 $E \rightarrow E - E$
 $E \rightarrow (E)$



Ambiguous Grammars

- Ambiguity can also occur for the same operations:

$E \rightarrow \text{id}$

$E \rightarrow \text{num}$

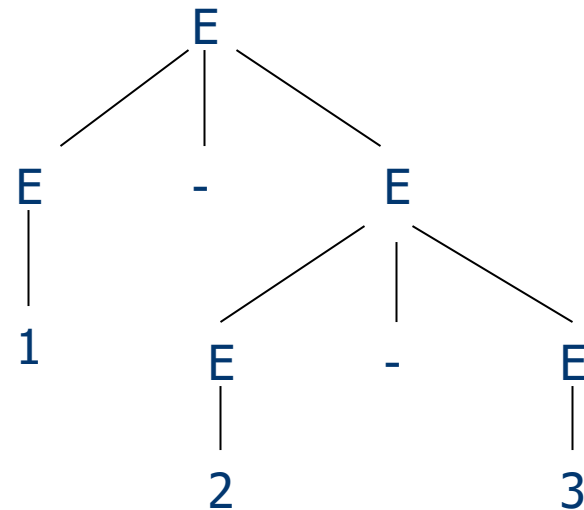
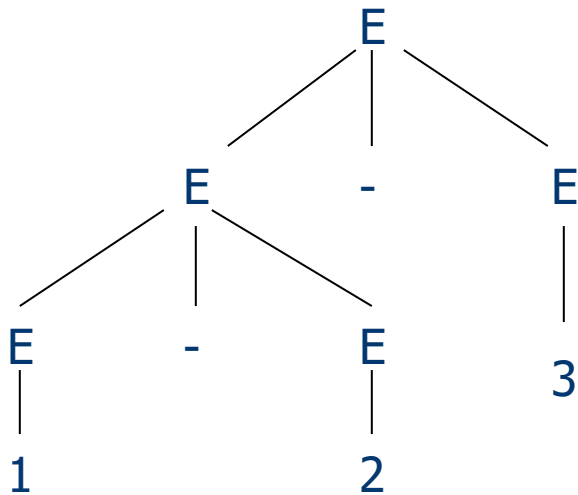
$E \rightarrow E * E$

$E \rightarrow E / E$

$E \rightarrow E + E$

$E \rightarrow E - E$

$E \rightarrow (E)$



Resolving Ambiguities

- Explicitly state which parse tree is correct

- No change required to the grammar

- Precedence:

- Stated order of derivation based on operator
- Sub-trees are evaluated before the root expression
 - The order of derivations is opposite to the order of evaluations.

- Associativity:

- Stated order of derivation based on location
- Left associative: derivation from the first choice
- Right associative: derivation from the last choice

Resolving Ambiguities

- An unambiguous grammar for expressions:

$E \rightarrow E + T$

$E \rightarrow E - T$

$E \rightarrow T$

$T \rightarrow T * F$

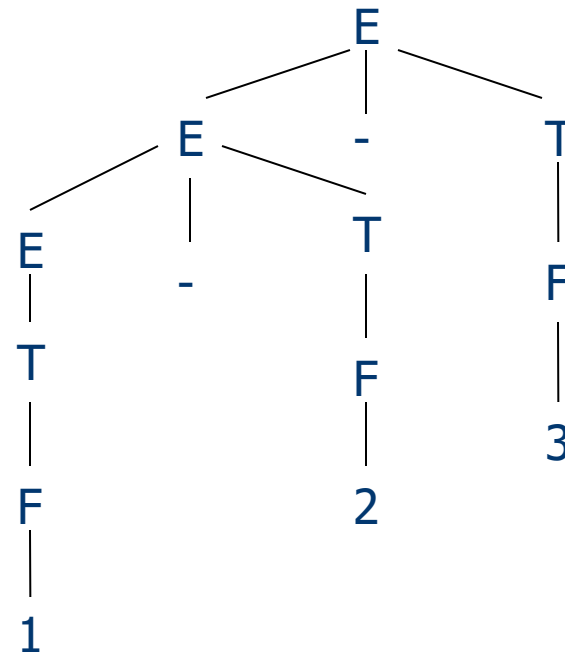
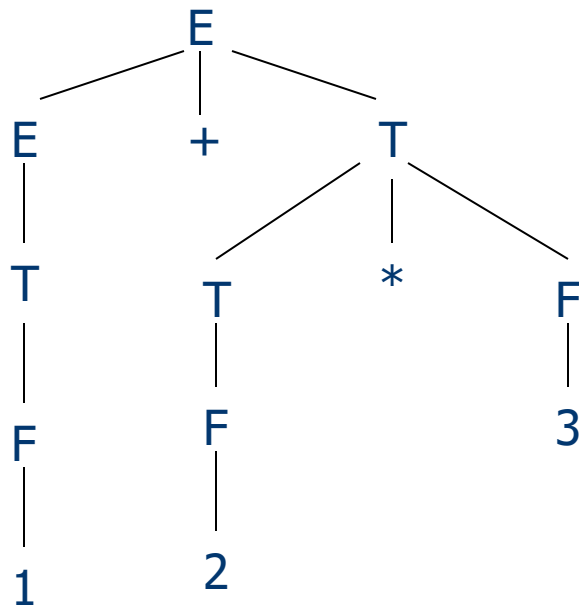
$T \rightarrow T / F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Resolving Ambiguities

● Associativity of operators:

- Left recursion specifies left associativity:

$$E \rightarrow E - T \quad T \rightarrow T / F$$

- Right recursion specifies right associativity:

$$\begin{array}{ll} F \rightarrow B ** F & B \rightarrow \text{id} \\ F \rightarrow B & B \rightarrow \text{num} \\ & B \rightarrow (E) \end{array}$$

Backus Normal Form

- BNF is equivalent to context-free grammars.
- In BNF, non-terminals are represented by pointed brackets:

`<program> -> begin <stmt-list> end`

`<stmt-list> -> <stmt> | <stmt> ; <stmt-list>`

`<stmt> -> id := <expr>`

`<expr> -> <expr> + <term> | <expr> - <term> | <term>`

`<term> -> <term> * <factor> | <term> / <factor> | <factor>`

`<factor> -> (<expr>) | id | num`

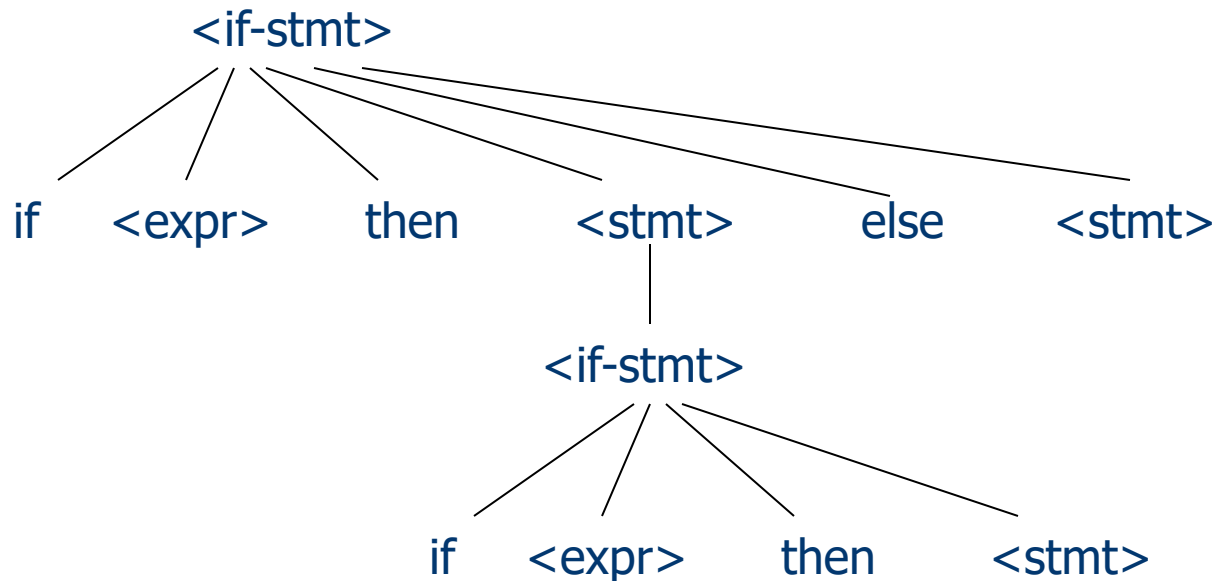
- Extended BNF includes shorthand notations for repetitive and optional constructs.

Dangling Else Problem

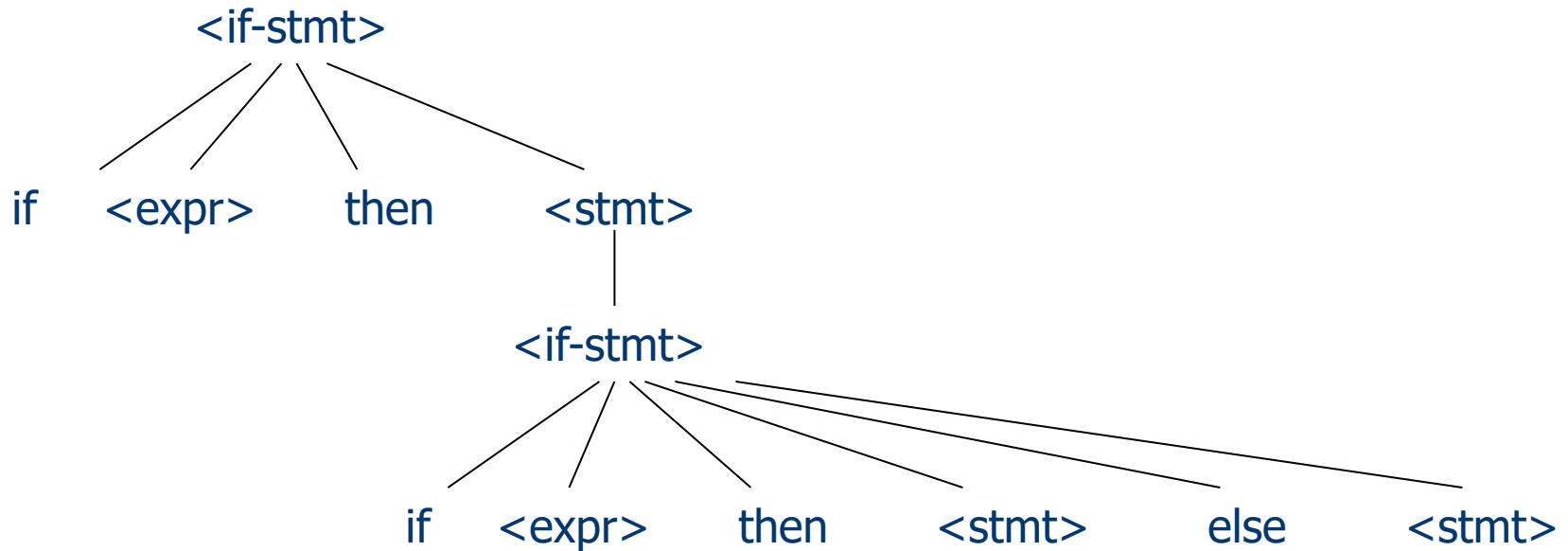
● Ambiguous grammar for if-statements:

$\langle \text{stmt} \rangle \rightarrow \langle \text{if-stmt} \rangle \mid \langle \text{other-stmt} \rangle$

$\langle \text{if-stmt} \rangle \rightarrow \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \mid$
 $\text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$



Dangling Else Problem



○ Unambiguous grammar for if-statements:

`<stmt>` \rightarrow `<matched>` | `<unmatched>`

`<matched>` \rightarrow `if <expr> then <matched> else <matched> |`
`<other-stmt>`

`<unmatched>` \rightarrow `if <expr> then <stmt> |`
`if <expr> then <matched> else <unmatched>`

Chomsky's Language Hierarchy

● Classes of grammars:

➤ Regular grammar

$A \rightarrow zB \mid z$ or $A \rightarrow Bz \mid z$ but not both

➤ Context-free grammar

$A \rightarrow B$ where A is a non-terminal and B is any string

➤ Context-sensitive grammar

$xAz \rightarrow xBz$ where A is a non-terminal and B is any string

➤ Unrestricted grammar

Also called recursively enumerable

● Context issues in programming languages:

➤ Variables are declared before being used

➤ Disambiguating rules

➤ Deferred to a parser generator or semantic analyzer