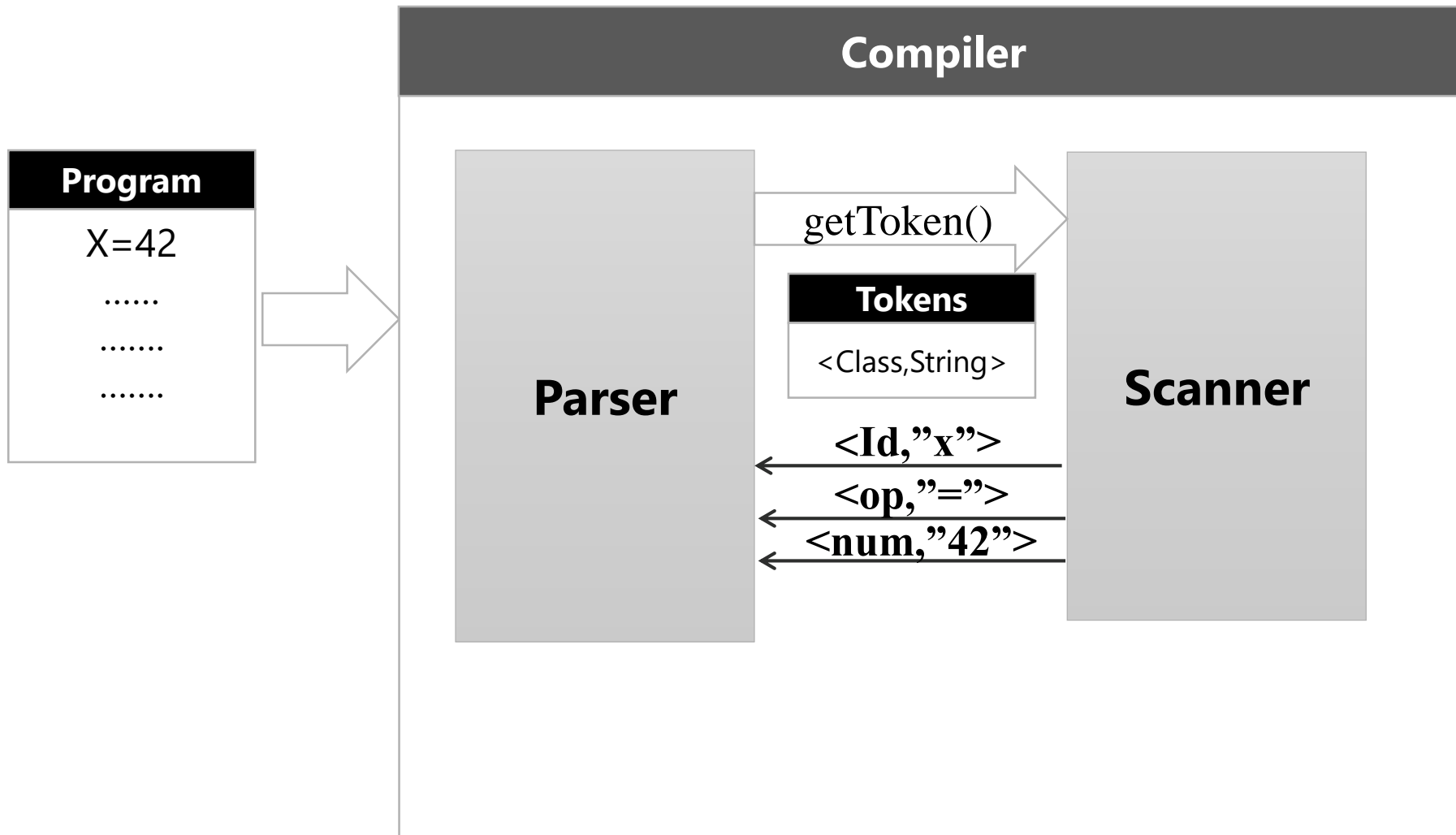


Parser

Lecture Five

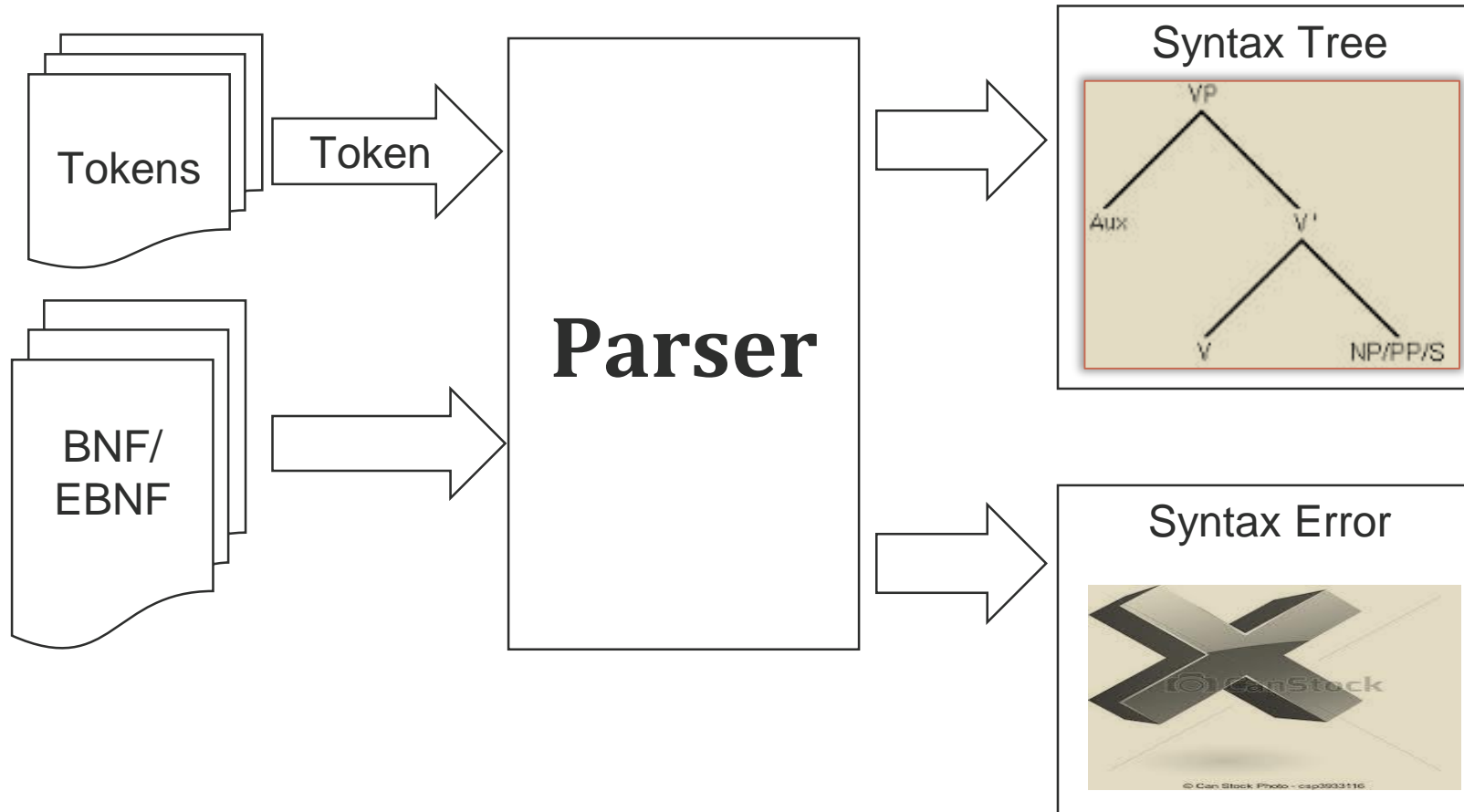


Parser (syntax Analyzer)

Lecture 5 Topics

- The Parsing Process
- Context-Free Grammars
- Derivations
- Parse Trees
- Abstract Syntax Trees
- Ambiguous Grammars
- Dealing with ambiguity
- EBNF (Extended BNF)

The Parsing Process



The Parsing Process

- Input: *tokens*
- Output: A *syntax tree*
- lexical structure of tokens is specified by regular expressions, the syntax of programs is specified by the *grammar rules* of a *context-free grammar*.
- Parser must distinguish between valid and invalid programs (strings of tokens).
- $(1 + + 2) + 3$
- The parser verifies that the string can be generated by the grammar for the source language

Regular Languages

- The languages used till now is Regular Language
- A language is *regular* if it is accepted by some DFA.

Examples:

- The strings that represent floating point numbers in Java.
- Some Languages cannot be represented by DFA. They are not regular
- any language that requires **unbounded counting** cannot be regular

Examples:

- $L1 = \{0^n 1^n \mid n \geq 1\}$: The set of strings consisting of one or more 0's followed by the same number of 1's.
- $L2 = \{()^n \mid n > 0\}$: The language of balanced parentheses.
e. g. $()$, $()()$, $(())$, $(())()$, ...

Grammar

sentence \rightarrow <subject> <verb-phrase> <object>

subject \rightarrow This | Computers | I

verb-phrase \rightarrow <adverb> <verb> | <verb>

adverb \rightarrow never

verb \rightarrow is | run | am | tell

object \rightarrow the <noun> | a <noun> | <noun>

noun \rightarrow university | world | cheese | lies

leftmost derivation

- sentence → <subject> <verb-phrase> <object>
- sentence → This <verb-phrase> <object>
- sentence → This <verb> <object>
- sentence → This is <object>
- sentence → This is a <noun>
- sentence → This is a university

sentence → <subject> <verb-phrase> <object>

subject → This | Computers | I

verb-phrase → <adverb> <verb> | <verb>

adverb → never

verb → is | run | am | tell

object → the <noun> | a <noun> | <noun>

noun → university | world | cheese | lies

Context Free Grammar(CFG)

Context free grammar $G = (N, T, S, P)$

- N : set of Non-terminal symbols
- T : set of Terminal symbols
- S : Start symbol ($S \in N$)
- P : set of Production rules $\{\alpha \rightarrow \beta \mid \alpha \in N \wedge \beta \in (N \cup T)^*\}$
- $N \cap T = \emptyset$

Note :

- $(N \cup T)^* = \{N \cup T \cup \varepsilon\}$
- $::=$ may be used instead of \rightarrow

CFG Example

$$L1 = \{ 0^n 1^n \mid n \geq 1 \}$$

$$S \rightarrow 01$$

$$S \rightarrow 0S1$$

- Nonterminal = $\{S\}$.
- Terminals = $\{0, 1\}$.
- Start symbol = S .
- Productions =

$$S \rightarrow 01$$

$$S \rightarrow 0S1$$

CFG Example

$L2 = \{ (^n)^n \mid n > 0 \}$: The language of balanced parentheses.

$$S \rightarrow (S)$$

$$S \rightarrow \varepsilon$$

- $N = \{ S \}$.
- $T = \{ (,) \}$.
- Start symbol = S .
- Productions =

$$S \rightarrow (S)$$

$$S \rightarrow \varepsilon$$

CFG Example

$$\textit{exp} \rightarrow \textit{exp op exp} \mid (\textit{exp}) \mid \textit{number}$$
$$\textit{op} \rightarrow + \mid - \mid *$$

- $N = \{\textit{exp}, \textit{op}\}.$
- $T = \{ (,), \textit{number}, *, -, + \}.$
- Start symbol = $\textit{exp}.$
- Productions =

$$\textit{exp} \rightarrow \textit{exp op exp} \mid (\textit{exp}) \mid \textit{number}$$
$$\textit{op} \rightarrow + \mid - \mid *$$

Derivation

- Given a production:

$$A ::= X_1 X_2 \dots X_n$$

$aAb \Rightarrow aX_1 X_2 \dots X_nb$ is called a derivation

- A derivation is a sequence of replacements of structure names by choices on the right-hand sides of grammar rules.
 - It starts with a single structure name and ends with a string of token symbols
 - At each step in a derivation, a single replacement is made using one choice from a grammar rule.
- The set of all strings of token symbols obtained by derivations from the exp symbol is the *language defined by the grammar* of expressions. We can write this as:

$$L(G) = \{ s : \text{exp} \Rightarrow^* s \}$$

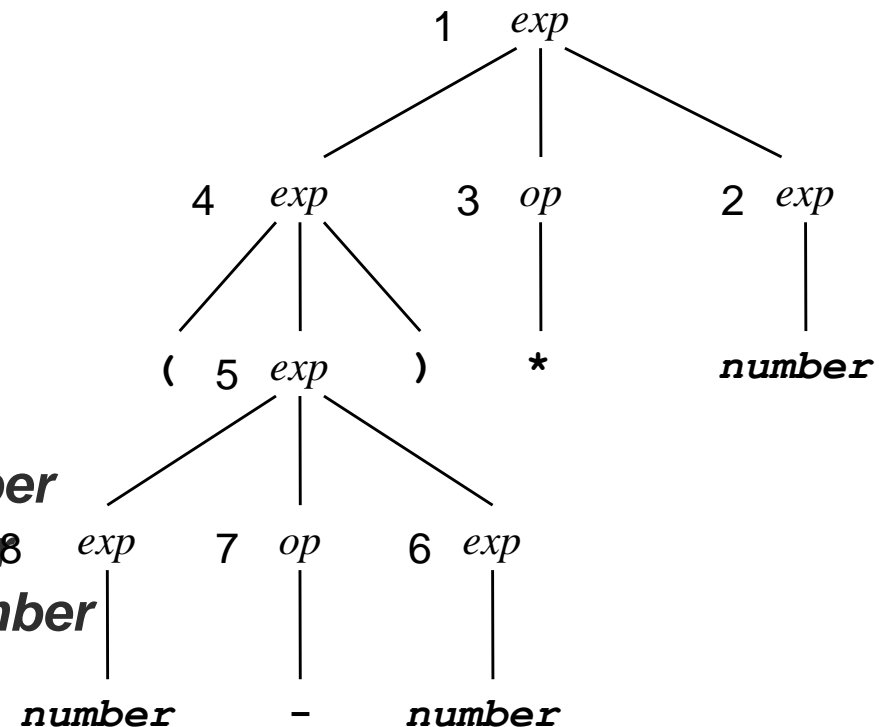
Derivation (Example)

$exp \rightarrow exp\ op\ exp \mid (exp) \mid \textit{number}$

$op \rightarrow + \mid - \mid *$

Derive $(34-3) * 42$

- (1) $exp \Rightarrow exp\ op\ exp$
- (2) $\Rightarrow exp\ op\ \textit{number}$
- (3) $\Rightarrow exp\ *\ \textit{number}$
- (4) $\Rightarrow (exp) *\ \textit{number}$
- (5) $\Rightarrow (exp\ op\ exp) *\ \textit{number}$
- (6) $\Rightarrow (exp\ op\ \textit{number}) *\ \textit{number}$
- (7) $\Rightarrow (exp - \textit{number}) *\ \textit{number}$
- (8) $\Rightarrow (\textit{number} - \textit{number}) *\ \textit{number}$



Right most derivation expand the Rightmost nonterminal in the production

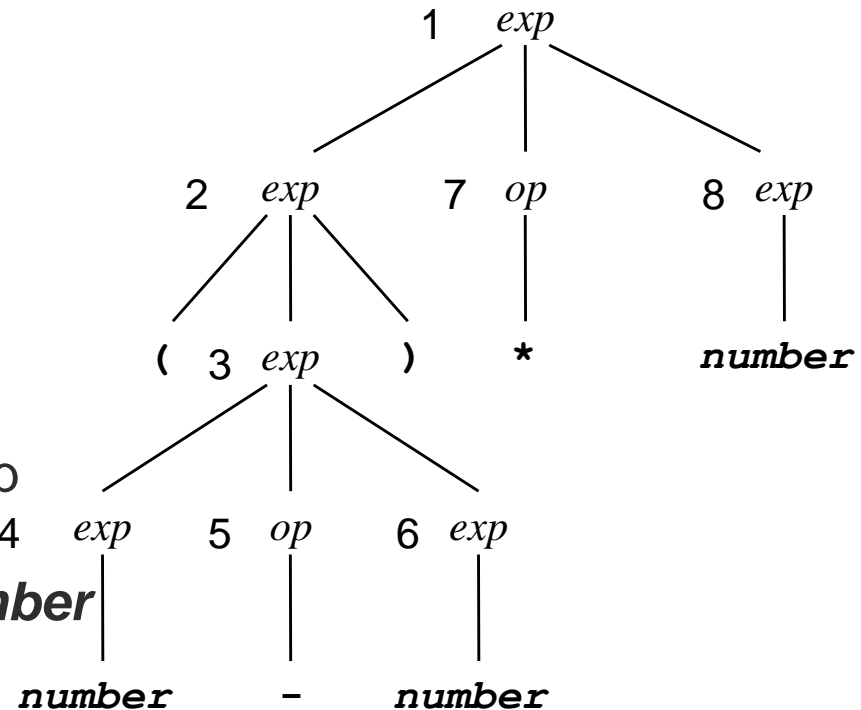
Derivation (Example)

$exp \rightarrow exp\ op\ exp \mid (exp) \mid \textit{number}$

$op \rightarrow + \mid - \mid *$

Derive $(34-3) * 42$

- (1) $exp \Rightarrow exp\ op\ exp$
- (2) $\Rightarrow (exp)\ op\ exp$
- (3) $\Rightarrow (exp\ op\ exp)\ op\ exp$
- (4) $\Rightarrow (\textit{number}\ op\ exp)\ op\ exp$
- (5) $\Rightarrow (\textit{number} - exp)\ op\ exp$
- (6) $\Rightarrow (\textit{number} - \textit{number})\ op\ exp$
- (7) $\Rightarrow (\textit{number} - \textit{number}) * exp\ 4$
- (8) $\Rightarrow (\textit{number} - \textit{number}) * \textit{number}$

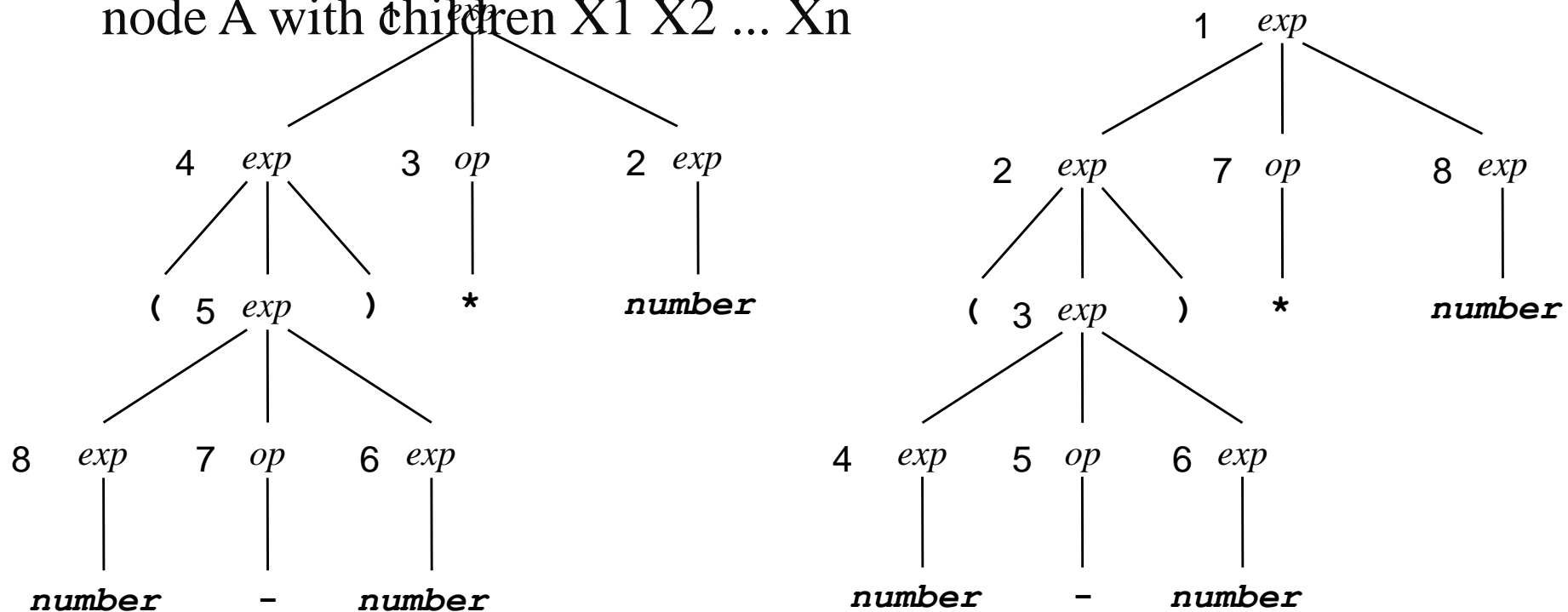


Left most derivation expand the leftmost nonterminal in the production

Parse Tree

Parse tree:

- the start symbol \rightarrow the root
- the terminals of the input sequence \rightarrow leafs
- for each production $A \rightarrow X_1 X_2 \dots X_n$ used in a derivation, a node A with children $X_1 X_2 \dots X_n$



Derivations

- Derivations allow us to replace any of the variables in a string. Leads to many different derivations of the same string.
- Leftmost Derivations: expand the leftmost nonterminal in the production
- Rightmost Derivations: expand the rightmost nonterminal in the production

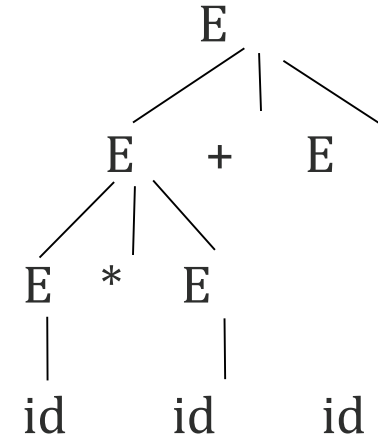
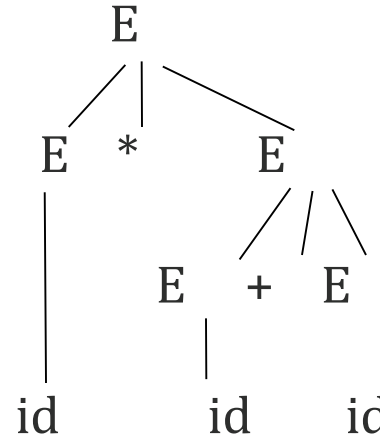
Ambiguous Grammars

$E ::= E + E$

| $E * E$

| id

String id*id+id



It is **ambiguous**: more than one parse tree for the same input depend on derivation used.

Ambiguous Grammars

Remove Ambiguity:

1. Make the grammar left-recursive

$$E \rightarrow E + E' \mid E'$$

2. Make the grammar right-recursive

$$E \rightarrow E' + E \mid E'$$

3. Make the grammar non-recursive

$$E ::= E' + E' \mid E'$$

Parsing

- A *leftmost* derivation corresponds to a (*top-down*) pre-order traversal of the parse tree.
- A *rightmost* derivation corresponds to a (*bottom-up*) post-order traversal, but in reverse.
- **Top-down parsers** construct leftmost derivations.
 - (LL = Left-to-right traversal of input, constructing a Leftmost derivation)
- **Bottom-up parsers** construct rightmost derivations in reverse order.
 - (LR = Left-to-right traversal of input, constructing a Rightmost derivation)

Ambiguous Grammar(priority)

$E \rightarrow E + E \mid E * E \mid E \wedge E \mid \text{num}$

String: 4*3+2

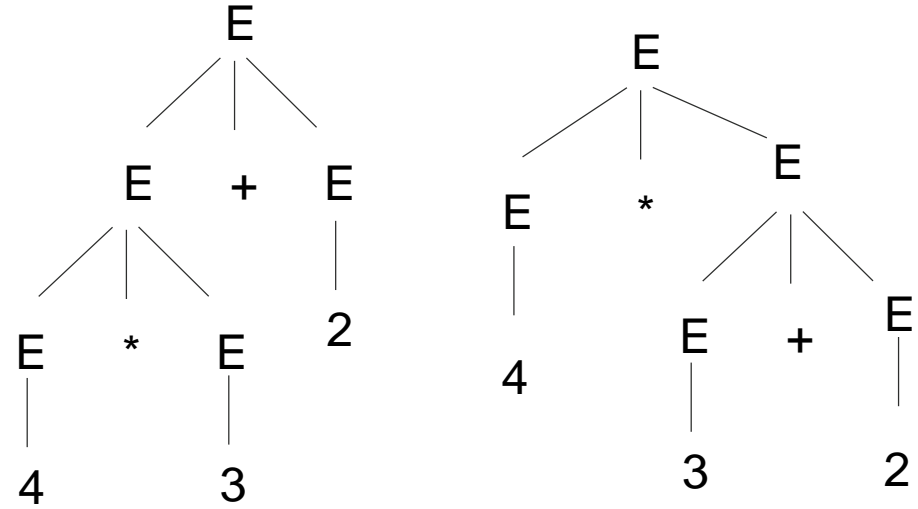
Remove Ambiguity

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow G \wedge F \mid G$

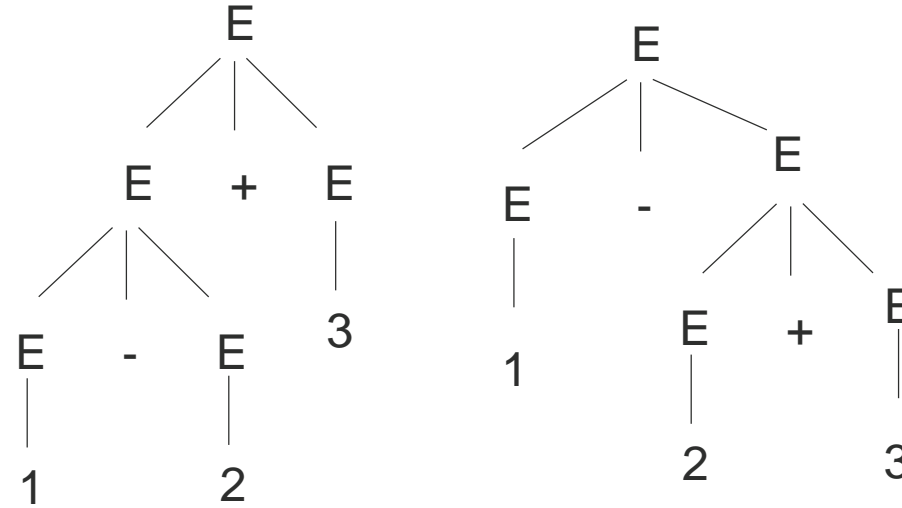
$G \rightarrow \text{num}$



Ambiguous Grammar (associativity)

$E \rightarrow E + E \mid E - E \mid \text{num}$

String: 1-2+3



Remove Ambiguity

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

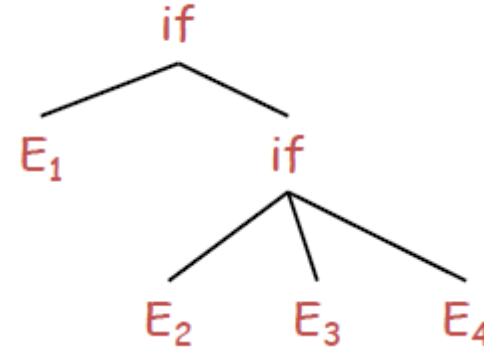
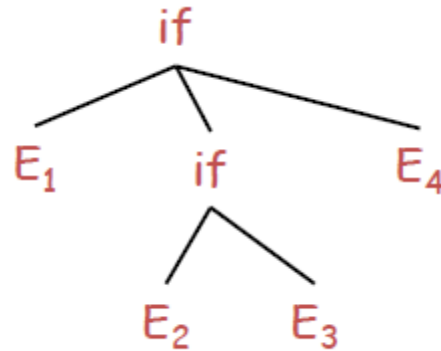
$T \rightarrow \text{num}$

Ambiguous Grammars

$E \rightarrow \text{if } E \text{ then } E$

$\mid \text{if } E \text{ then } E \text{ else } E$

String : if E1 then if E2 then E3 else E4



Unambiguous Grammar

$E \rightarrow \text{MIF} \mid \text{UIF}$

$\text{MIF} \rightarrow \text{if } E \text{ then MIF else MIF}$

$\text{UIF} \rightarrow \text{if } E \text{ then } E$

$\mid \text{if } E \text{ then MIF else UIF}$