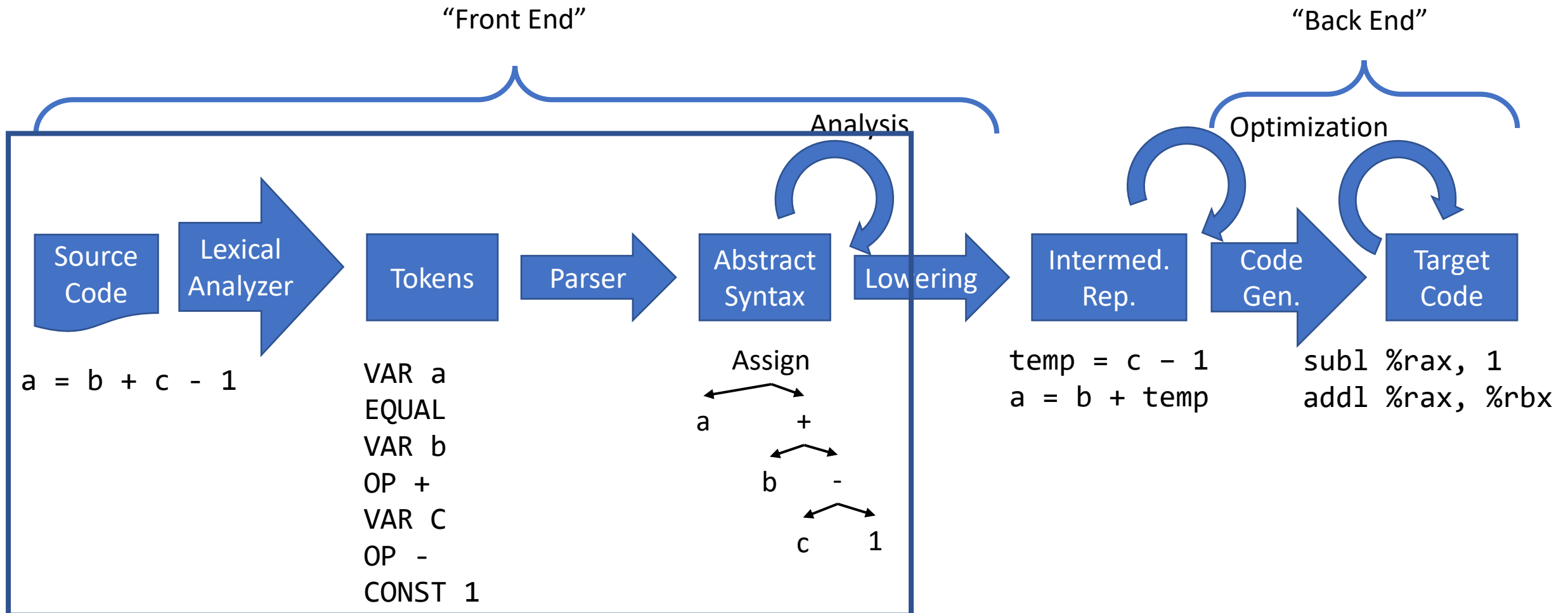


# CS443: Compiler Construction

Lecture 1: Lexing and Parsing

# Compilers translate code in phases



# Terminology

- *Lexical analysis* “lexing”
- Performed by *lexical analyzer* “lexer”
- Produces stream of *tokens*

# Tokens are specified using a *regular* grammar

- Regular expressions  $R$ :

- $\epsilon$             Empty string
- $abc$             Exactly the string  $abc$             Literal
- $R_1R_2$          $R_1$  followed by  $R_2$             Concatenation
- $R_1 \mid R_2$      $R_1$  or  $R_2$                     Alternation
- $R^*$             Zero or more  $R$                 Kleene Star
  
- $R^+$             One or more  $R$
- $R?$             Optional  $R$
- $[a-z]$          $a, b, c, d, \dots, z$

Tokens are specified using a *regular* grammar

digit ::= [0-9]

alpha ::= [a-z]

ident ::= alpha (alpha | digit)\*

num ::= digit<sup>+</sup>

ident → IDENT s

num → NUM s

“while” → WHILE

“+” → PLUS...

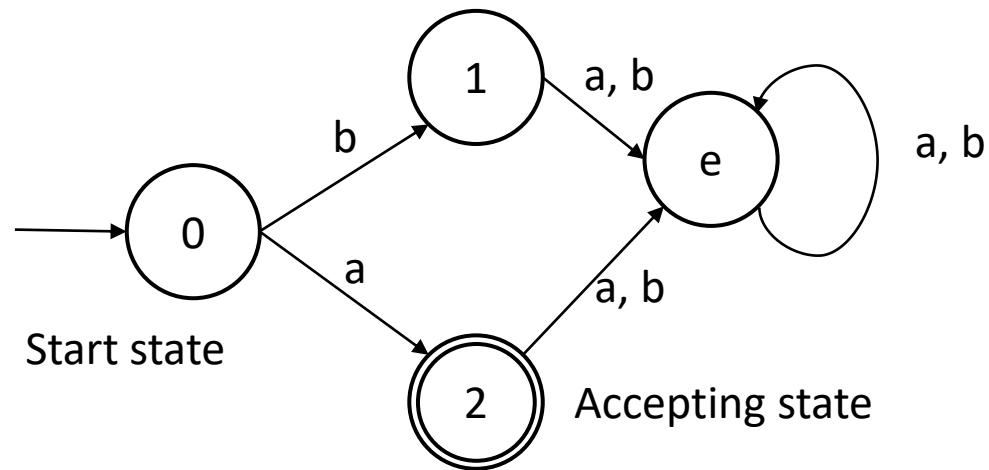
# Lexing examples

- while (i < 5)    WHILE; LPAREN; IDENT "i"; LT; NUM 5; RPAREN
- while i < 5)    WHILE; IDENT "i"; LT; NUM 5; RPAREN
- whole (i < 5)    IDENT "whole"; LPAREN; IDENT i; LT; NUM 5; RPAREN

Might be syntax errors  
during *parsing*. *Not* errors  
during lexing.

# Regex matching can be done by finite state machines (FSMs)

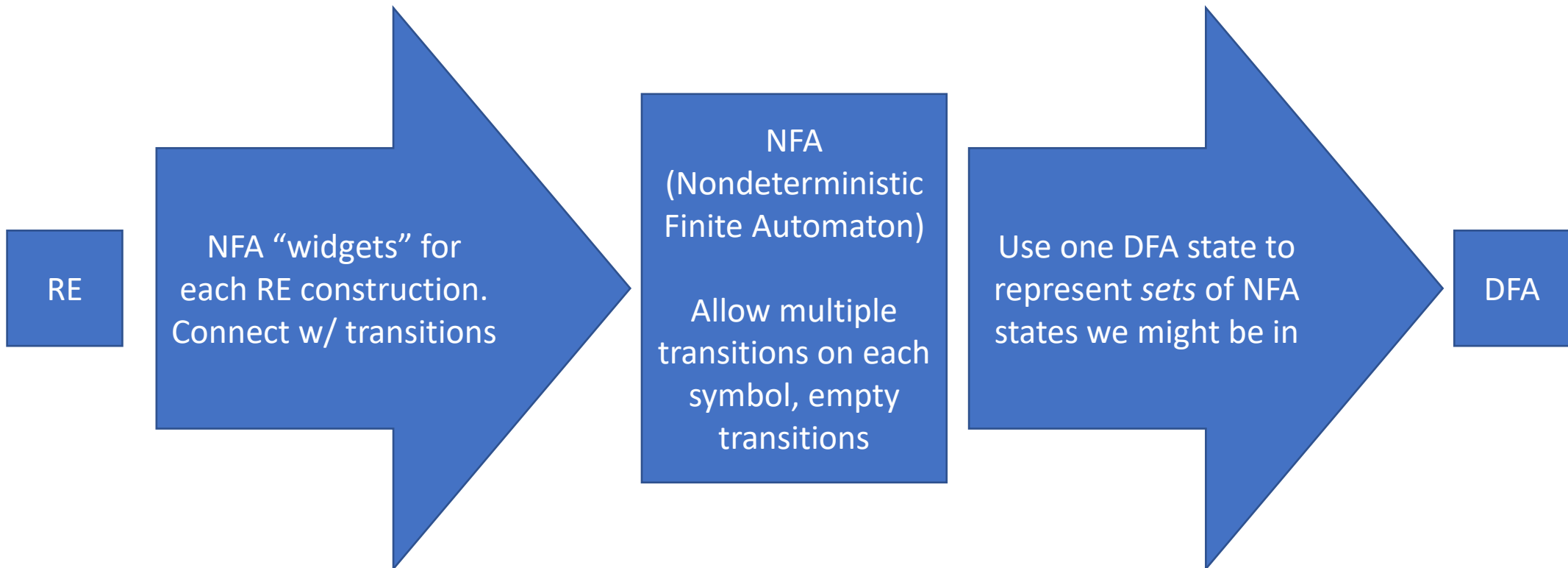
- Deterministic Finite Automaton (DFA)
- *States + Transition function + Start state + Set of Accepting states*



Accepts exactly “a”

# Can convert regexes to DFAs

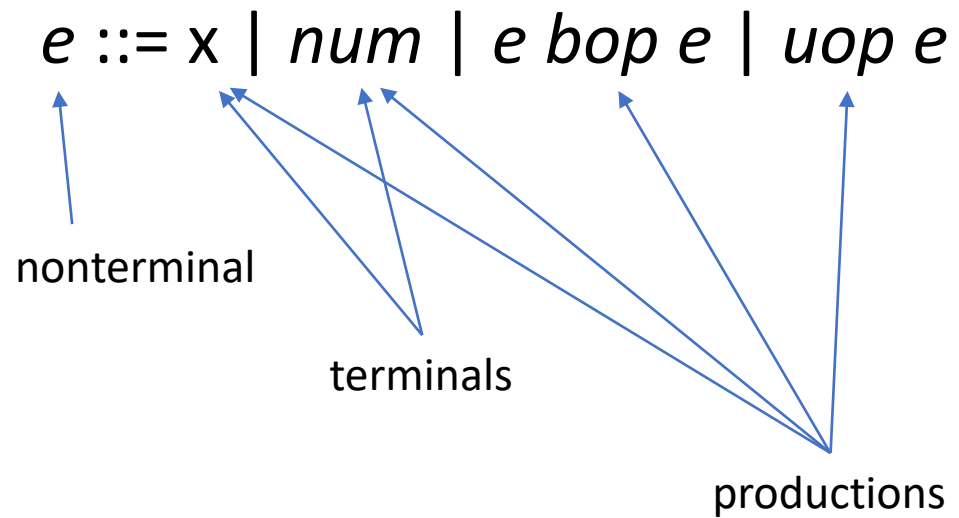
- Full algorithm in Appel, PDB. General idea:





- Draw DFA example

# BNF grammars are “context-free”



# Derivation: Expand one nonterminal at a time using a production

$e ::= n \mid e + e \mid (e)$

Input:  $1 + (2 + 3)$

$e$

$e + e$

$1 + e$

$1 + (e)$

$1 + (e + e)$

$1 + (2 + e)$

$1 + (2 + 3)$

Leftmost Derivation

$e$

$e + e$

$e + (e)$

$e + (e + e)$

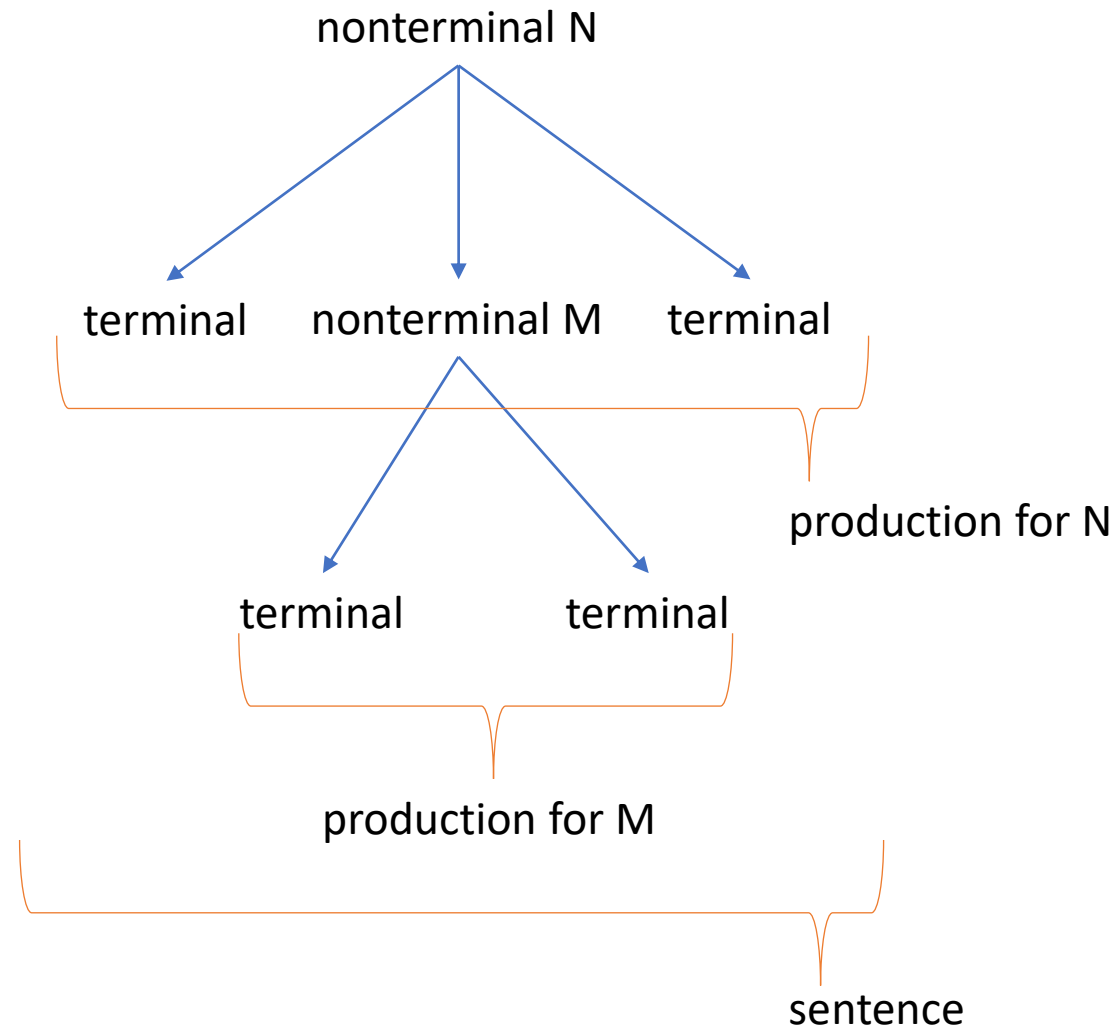
$e + (e + 3)$

$e + (2 + 3)$

$1 + (2 + 3)$

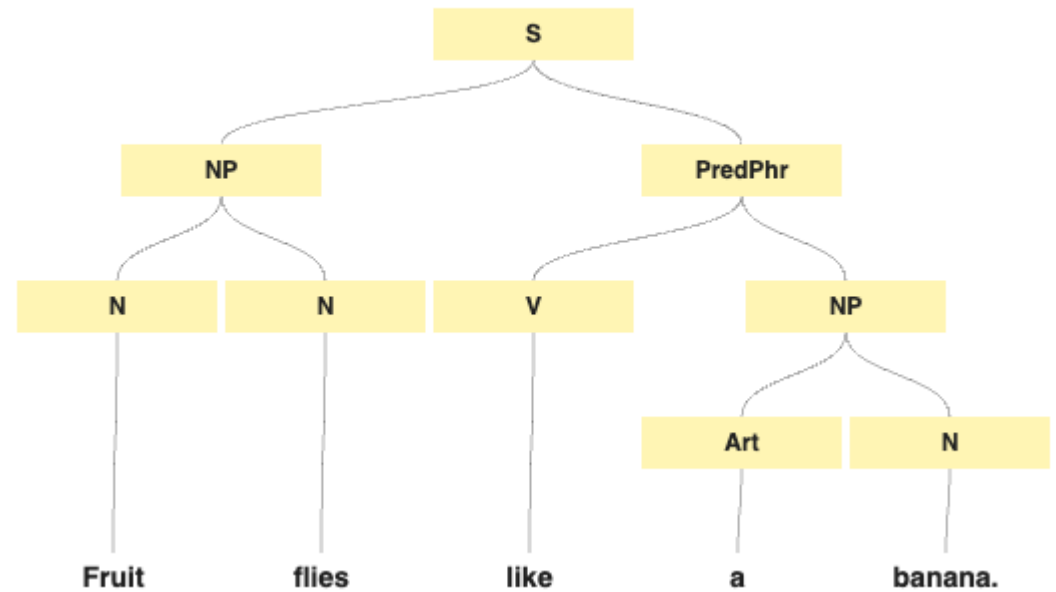
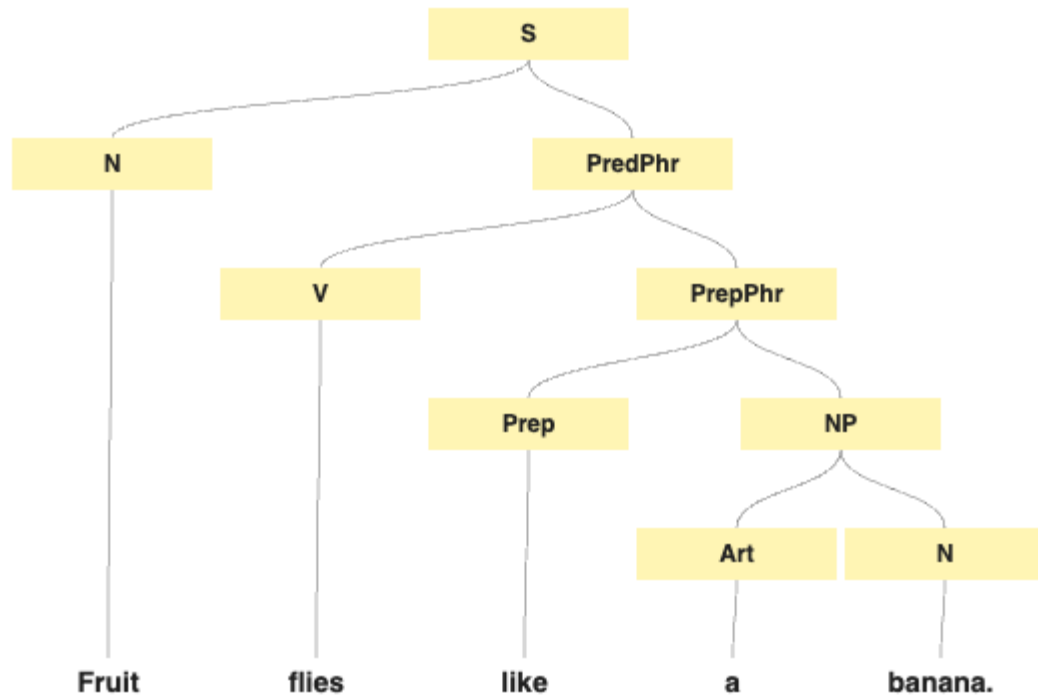
Rightmost Derivation

# Parsing: Produce a *parse tree* from a stream of tokens



# Ambiguous grammars allow multiple correct parse trees

- “Fruit flies like a banana”



# *Associativity* is one source of ambiguity

- $e ::= \text{num} \mid e + e \mid e - e$
- $1 + 2 + 3 + 4 \rightarrow ((1 + 2) + 3) + 4, (1 + 2) + (3 + 4), \dots$
- Solution:  
 $e ::= \text{num} \mid e + \text{num} \mid e - \text{num} \mid e + (e) \mid e - (e)$

# *Precedence* is one possible source of ambiguity

- Abstract syntax:  $e ::= e + e \mid e - e \mid e * e \mid e / e$

- $1 + 2 * 3 - 4$  ????

- Solution: Factoring out productions

- $f ::= \text{num} \mid (e)$

- $t ::= f \mid t * f \mid t / f$

- $e ::= t \mid e + t \mid e - t$

# Classic example: “dangling else”

- $s ::= \text{if } e \text{ } s \mid \text{if } e \text{ } s \text{ else } s$
- $\text{if } e_1 \text{ if } e_2 \text{ } s_1 \text{ else } s_2$ 
  - By convention:  $\text{if } e_1 \text{ (if } e_2 \text{ } s_1 \text{ else } s_2)$
- Solution:
  - $\text{closedstmt} ::= \text{if } e \text{ closedstmt else closedstmt}$ 
    - | ... (non-if stmts not ending with an *openstmt*)
  - $\text{openstmt} ::= \text{if } e \text{ closedstmt else openstmt} \mid \text{if } e \text{ stmt}$ 
    - | ... (non-if stmts ending with an *openstmt*)
  - $\text{stmt} ::= \text{openstmt} \mid \text{closedstmt}$