

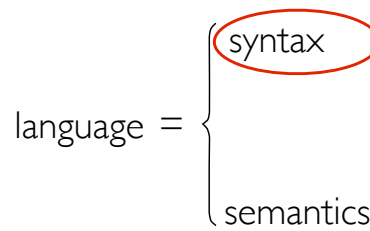
# CS 320: Principles of Programming Languages

Mark P Jones, Portland State University

Spring 2019

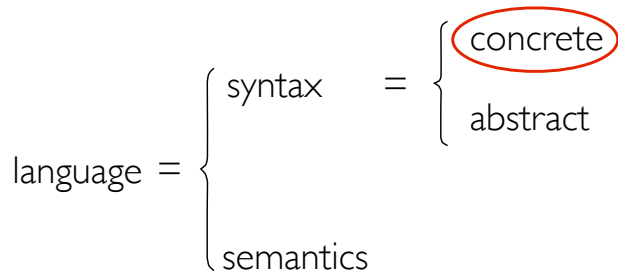
Week 2: Describing Syntax

1



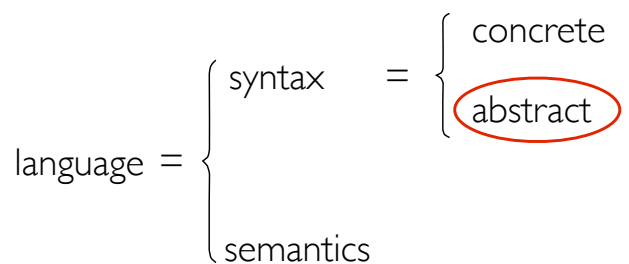
**syntax:** the written/spoken/symbolic/physical form; how things are communicated

2



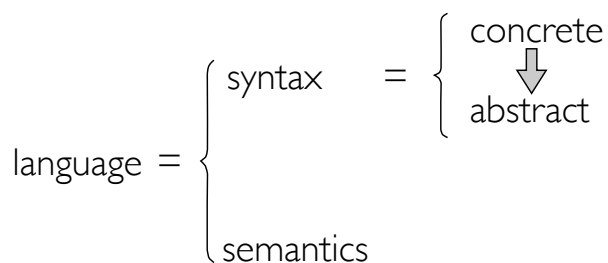
**concrete syntax:** the representation of a program text in its source form as a sequence of bits/bytes/characters/lines

3



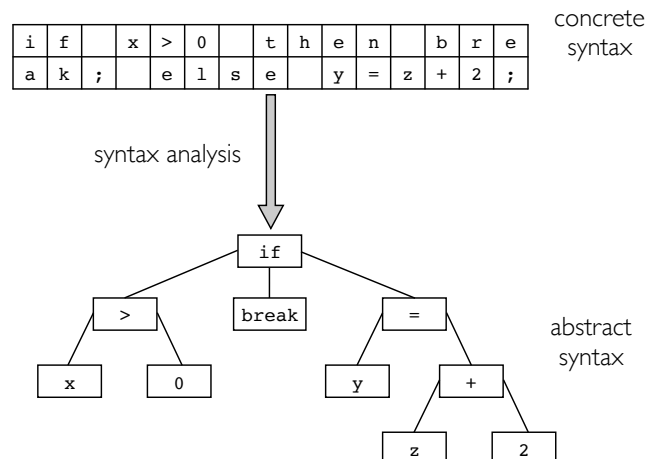
**abstract syntax:** the representation of a program structure, independent of written form

4

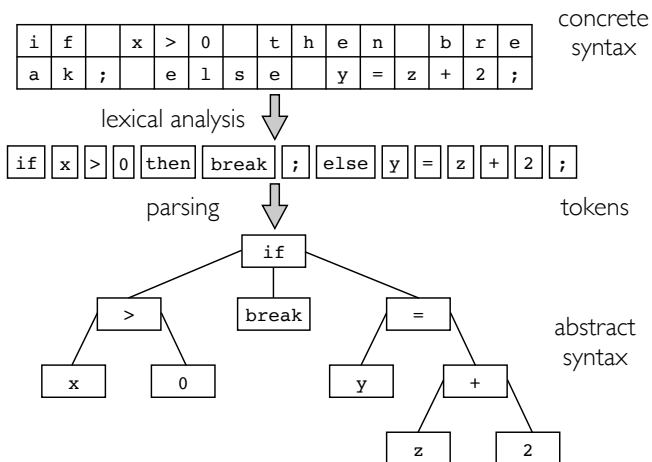


**syntax analysis:** This is one of the areas where theoretical computer science has had major impact on the practice of software development

5



6



# Wanted!

- We want methods for describing language syntax that are:
  - clear; precise, and unambiguous
  - expressive (e.g., finite descriptions of infinite languages)
  - suitable for use in the implementation of practical syntax analysis tools (lexical analyzers, parsers, ...)
- Formal language theory provides a rigorous foundation for specifying syntax. In particular:
  - Regular languages are well-suited to describing lexical syntax (what sequences of characters are valid tokens?)
  - Context-free languages are well-suited to describing grammatical structure (what sequences of tokens are valid sentences?)

# Formal languages

- Pick a set,  $A$ , of **symbols**, which we refer to as the **alphabet**
  - For lexical analysis, “symbols” are typically characters
  - For parsing, “symbols” are typically tokens
- The set of all finite strings of symbols taken from  $A$  is written as  $A^*$
- A **language** (over  $A$ ) is a subset of  $A^*$

## Examples

- If  $A = \{0, 1\}$ , then:
 
$$A^* = \{ "", "0", "1", "00", "01", "10", "11", "000", \dots \}$$
- Bytes form a finite language over  $A$ :
 
$$\text{Bytes} = \{ b_0 b_1 b_2 b_3 b_4 b_5 b_6 b_7 \mid b_i \in \{0, 1\} \}$$
- Even length bitstreams form an infinite language over  $A$ 

$$\text{Evens} = \{ "", "00", "01", "10", "11" \} \cup \{ xy \mid x, y \in \text{Evens} \}$$

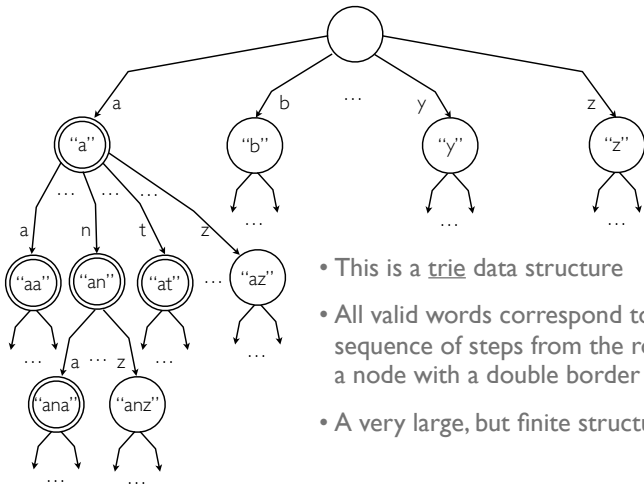
## English as a formal language

- With alphabet  $A = \{“a”, “b”, “c”, \dots, “z”\}$ :
  - $A^*$  = all strings with zero or more letters from  $A$
  - Words = the subset of  $A^*$  containing only those strings that are valid words in English  
e.g., “language”  $\in$  Words, “jubmod”  $\notin$  Words
  - English = the subset of Words\* containing only those strings that are valid sentences in English  
[“languages”, “are”, “interesting”]  $\in$  English
- But how do we specify which subsets to use?

## A list of English words (/usr/share/dict/words)

[illegible]

## Let's give this another trie ...



- This is a trie data structure
- All valid words correspond to a sequence of steps from the root to a node with a double border
- A very large, but finite structure

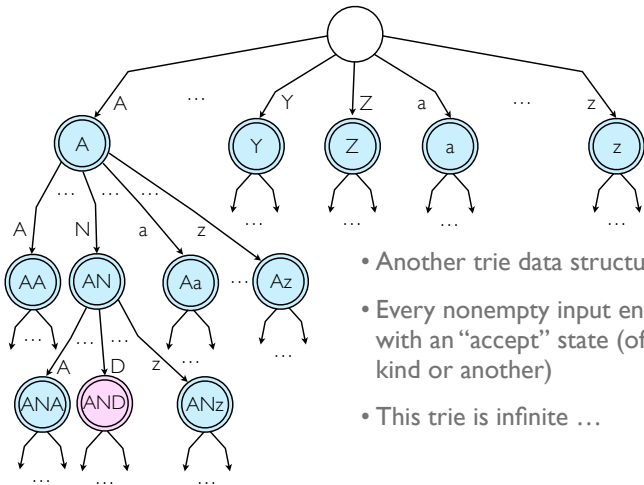
13

## Prop as a formal language

- With alphabet  $A = \{ "A", "B", "C", \dots, "Z", "(", ")", \dots \}$ :  
 $A^*$  = all strings with zero or more letters from A
- Tokens = Keywords: "AND", "OR", "NOT";  
 Literals: "TRUE", "FALSE";  
 Punctuation: "(" and ")"  
 Variable names: all nonempty subsets of  $A^*$  containing only letters
- Prop = the subset of Tokens\* corresponding to valid circuits
- But how do we specify these details?

14

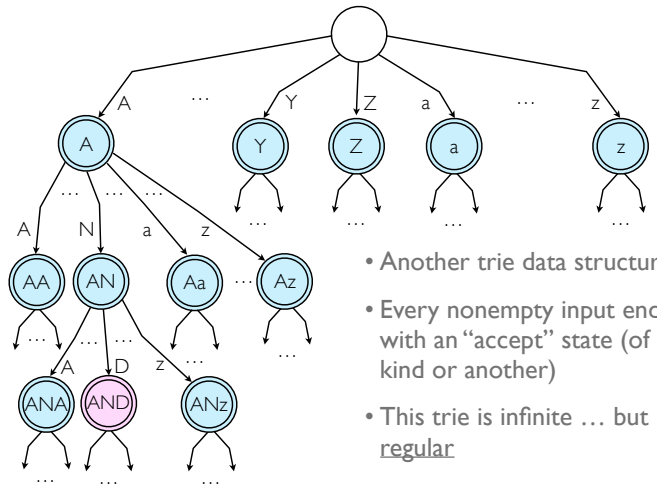
## Tokens in Prop (a subset)



- Another trie data structure
- Every nonempty input ends with an "accept" state (of one kind or another)
- This trie is infinite ...

15

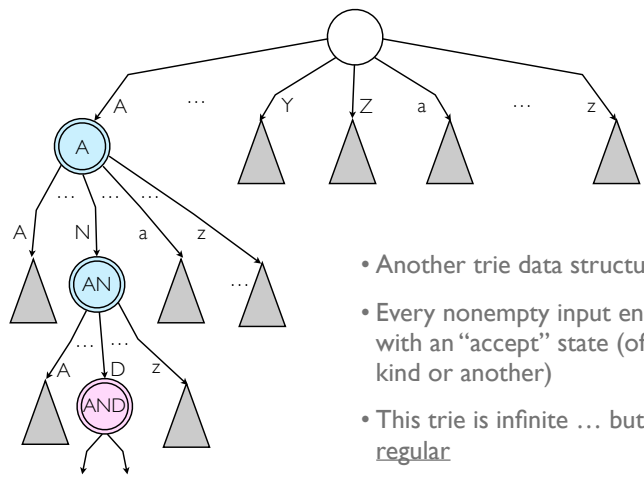
## Tokens in Prop (a subset)



- Another trie data structure
- Every nonempty input ends with an "accept" state (of one kind or another)
- This trie is infinite ... but regular

16

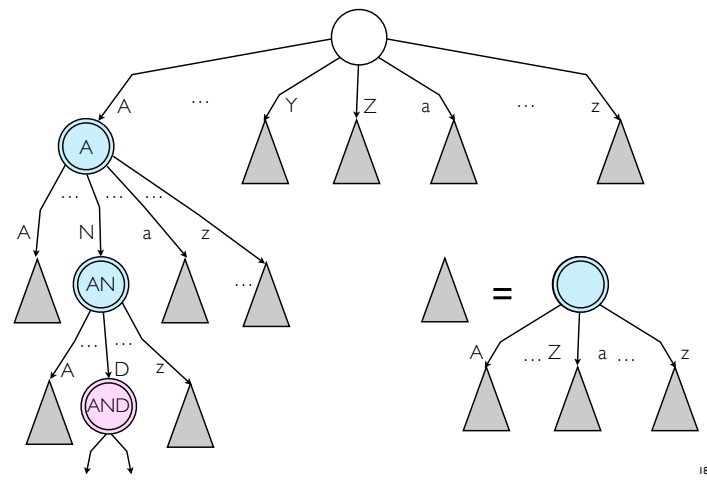
## Tokens in Prop (a subset)



- Another trie data structure
- Every nonempty input ends with an "accept" state (of one kind or another)
- This trie is infinite ... but regular

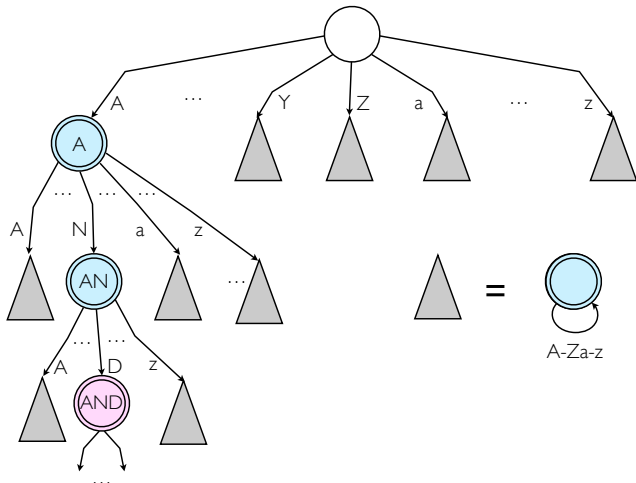
17

## Tokens in Prop (a subset)



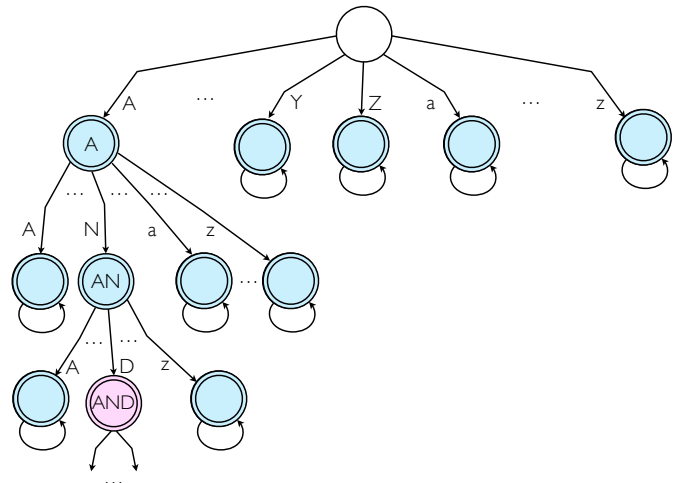
18

## Tokens in Prop (a subset)



19

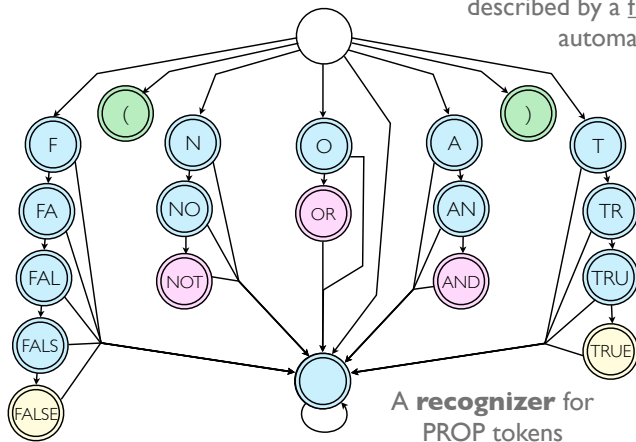
## Tokens in Prop (a subset)



20

## Tokens in Prop (all of them!)

An infinite language, described by a finite automaton!



A **recognizer** for PROP tokens

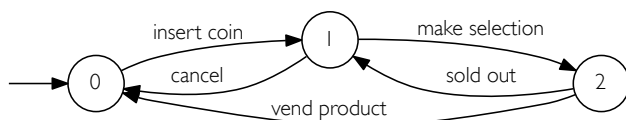
21

## Finite Automata and Regular Languages

22

## Terminology

A **finite automaton** (or state machine) describes a system that can **transition** (or move) between a finite set of **states** in response to particular **inputs** (or actions).



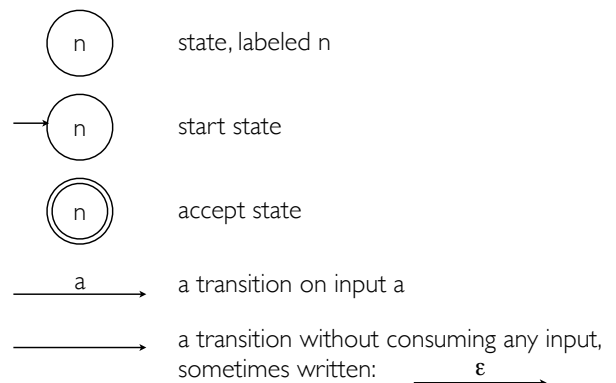
A **deterministic** finite automaton (DFA) has *at most one* way to transition out of any state in response to a given input

A **nondeterministic** finite automaton (NFA) may allow *multiple* transitions from some state in response to a given input

23

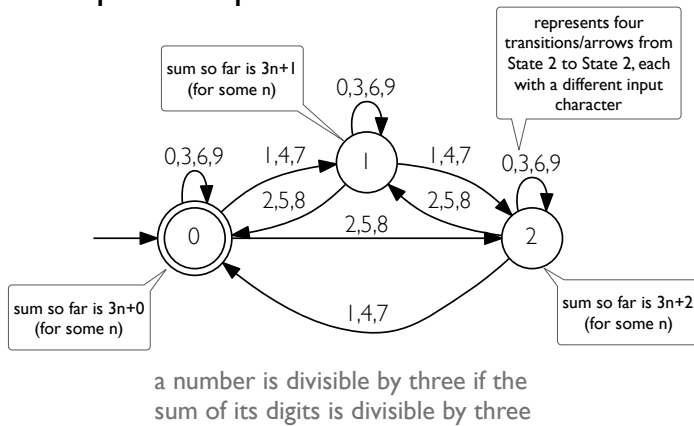
## Finite automata building blocks

We use the following symbols to describe automata:



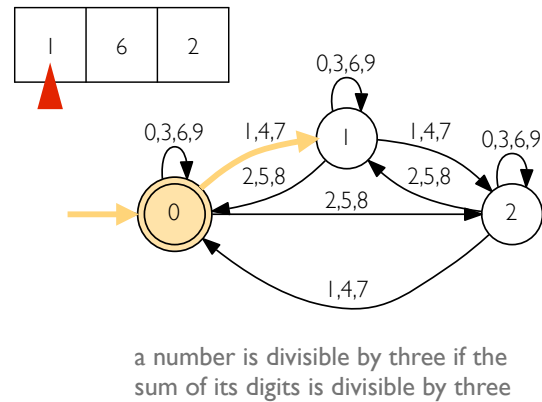
24

## Example: multiples of 3



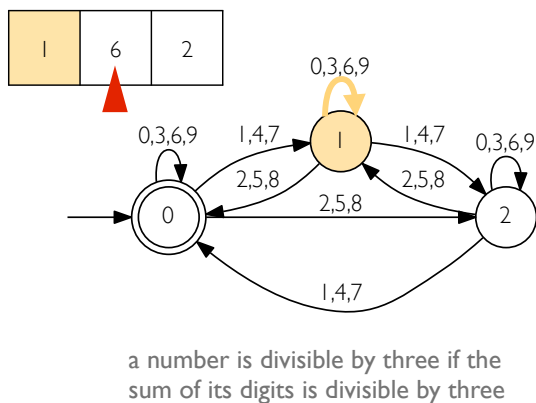
25

## Example: multiples of 3



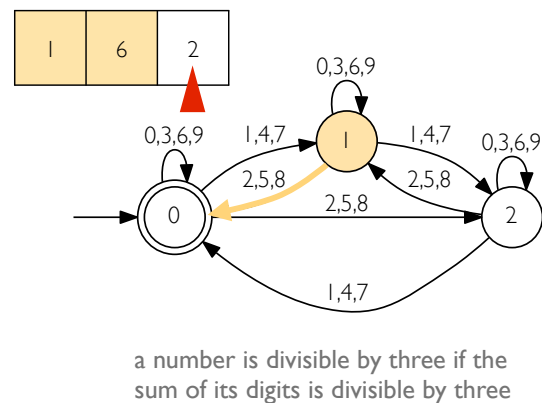
26

## Example: multiples of 3



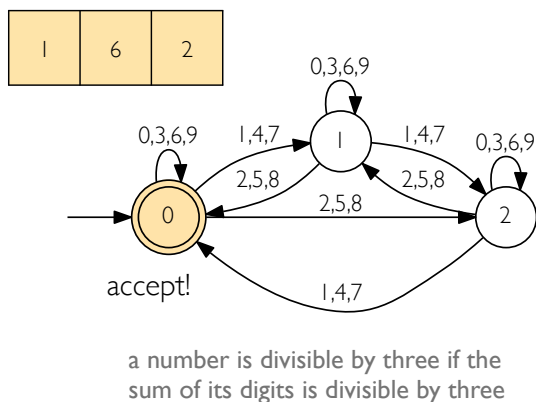
27

## Example: multiples of 3



28

## Example: multiples of 3



29

## Example: recognizers for lexical analysis

The "if" keyword:

Identifiers:

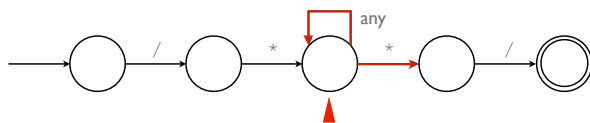
Integer literals: (without leading zeros)

30

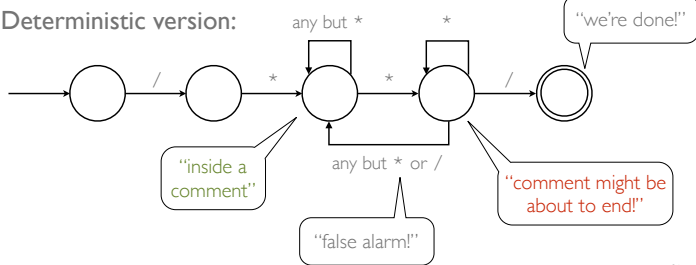
## Handling nondeterminism

Comments in C:

`/* y = x * z */`

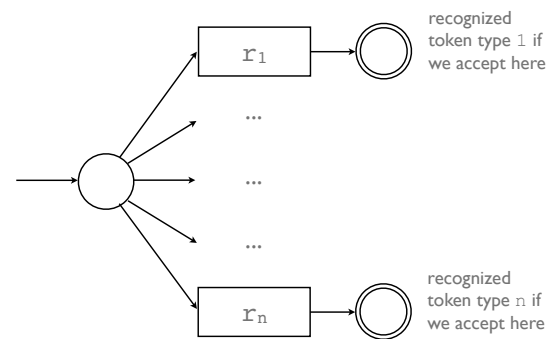


Deterministic version:



31

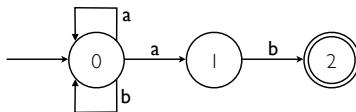
## Combining multiple recognizers



32

## Finite automata as language recognizers

- Every finite automaton defines a formal language:
  - The set of strings corresponding to paths from the start state to an accept state
- Example: what language does this automaton describe?

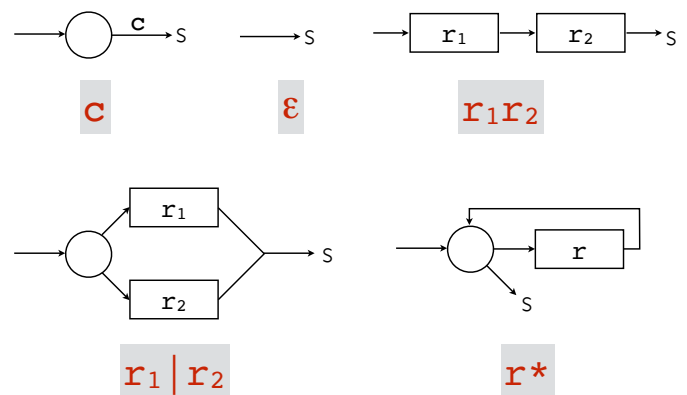


Answer: {"ab", "aab", "bab", "aaab", "abab", "baab", "bbab", ...}

- However:
  - "a" is not included (does not reach an accept state)
  - "cab" is not included (gets stuck at state 0)

33

## Capturing some recurring patterns



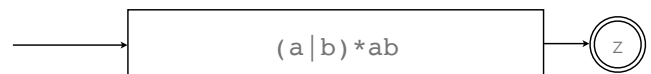
34

## Regular expressions

- A widely used notation for describing patterns in text strings: emacs, vi, grep, awk, perl, ruby, python, javascript, ...
- Also good for describing the lexical structure of programming languages ...
- Every regular expression corresponds to a DFA
- Every regular expression defines a language (the language that is recognized by the corresponding DFA)

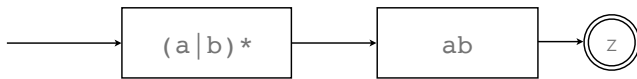
35

## A recognizer for $(a \mid b)^* ab$



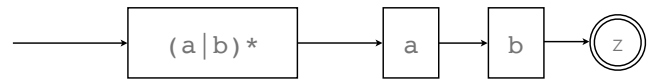
36

# A recognizer for $(a|b)^*ab$



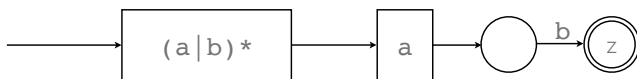
37

# A recognizer for $(a|b)^*ab$



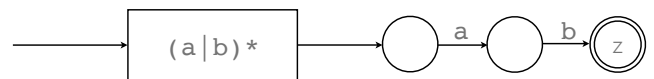
38

# A recognizer for $(a|b)^*ab$



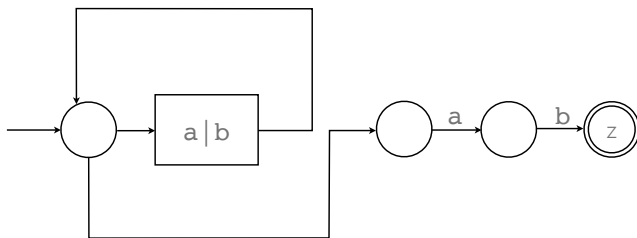
39

# A recognizer for $(a|b)^*ab$



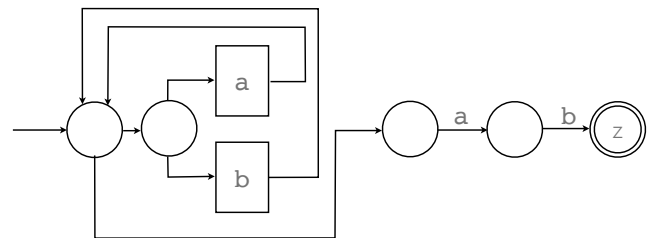
40

# A recognizer for $(a|b)^*ab$



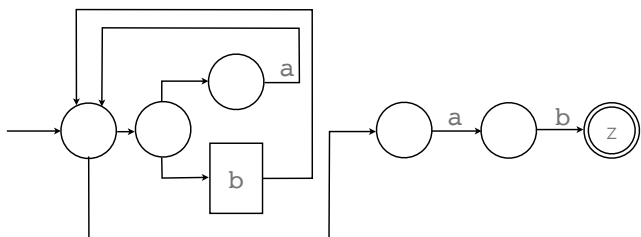
41

# A recognizer for $(a|b)^*ab$



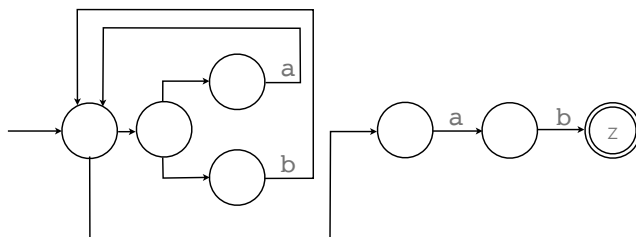
42

## A recognizer for $(a|b)^*ab$



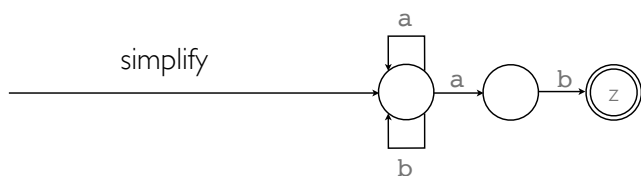
43

## A recognizer for $(a|b)^*ab$



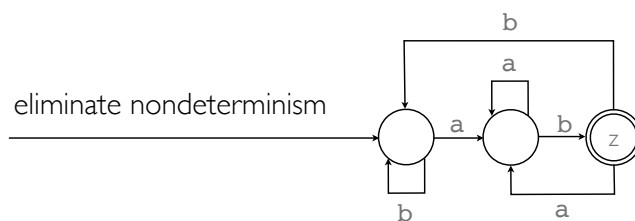
44

## A recognizer for $(a|b)^*ab$



45

## A recognizer for $(a|b)^*ab$



46

## Regular expressions

Expression	Meaning
$\epsilon$	<b>Empty:</b> matches the empty string
$c$	<b>Constant:</b> matches the single character $c$
$r_1   r_2$	<b>Alternatives:</b> matches text matching $r_1$ or text matching $r_2$
$r_1 r_2$	<b>Sequencing:</b> matches text matching $r_1$ followed by text matching $r_2$
$r^*$	<b>Repetition:</b> matches a sequence of zero or more items, each of which matches $r$
$(r)$	<b>Grouping:</b> matches text matching $r$

47

## Derived forms

Expression	Meaning
$r^+$	<b>Repetition:</b> a sequence of one or more items, each of which matches $r$ . $r^+ = rr^*$ .
$r?$	<b>Optional:</b> optional text matching $r$ . $r? = (r   \epsilon)$ .
$[abc]$	<b>Character classes:</b> short for $(a b c)$ . Also allows ranges of characters. e.g., $[a-zA-Z]$ .
$.$	<b>Wildcard:</b> matches any character (except newline)

48

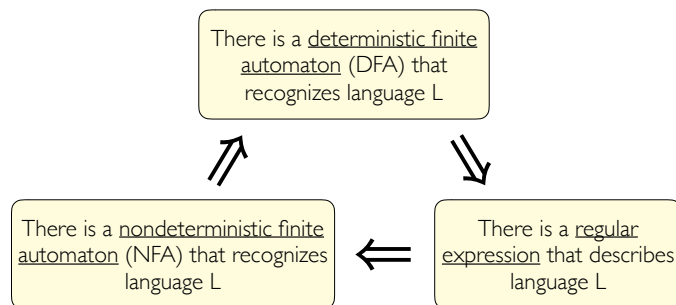


## Programming language examples

Regular expression	Describes
<code>if</code>	the keyword <code>if</code>
<code>[A-Za-z][A-Za-z0-9_]*</code>	identifiers
<code>0 [1-9][0-9]*</code>	integer literals (without leading zeros)
<code>[\t\n]*</code>	whitespace
<code>{int}{frac}?{exp}?</code> where <code>{int}</code> = <code>0 [1-9][0-9]*</code> <code>{frac}</code> = <code>\.{int}</code> <code>{exp}</code> = <code>[Ee](\+ -){int}</code>	floating point literals (using auxiliary definitions)

49

## Key things to know (see CS 311 for proofs, etc.)



In any/all of these situations, we say that L is a **regular language**

50

## Wanted! (reprise)

- We want methods for describing language syntax that are:
  - clear, precise, and unambiguous
  - expressive (e.g., finite descriptions of infinite languages)
  - suitable for use in the implementation of practical syntax analysis tools (lexical analyzers, parsers, ...)
- Regular expressions work really well for describing lexical syntax ...
- But the set of regular languages that they describe is quite limited ...

51

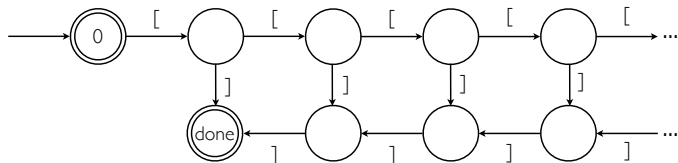
## A simple language of Brackets

- Brackets =  $\{ "" \} \cup \{ [ b ] \mid b \in \text{Brackets} \}$
- So the "words" in Brackets are:
  - `"", "[ ]", "[ [ ]", "[ [ [ ]", "[ [ [ [ ]", "[ [ [ [ [ ]", ...`
- In other words, nested pairs of bracket characters:
  - a sequence of n open brackets ...
  - ... followed by exactly n close brackets
- A subset of any language that uses parentheses/brackets
- But is it a regular language?

52

## Brackets is not regular

- If brackets is regular, then we can recognize it using a **finite** automaton that would look something like this:



- If  $s_n$  is the state that we can reach after n open brackets and  $n \neq m$ , then  $s_n \neq s_m$
- So this machine must have infinitely many states (at least one distinct state  $s_n$  for each n)
- So Brackets cannot be regular

53

## Iteration vs recursion

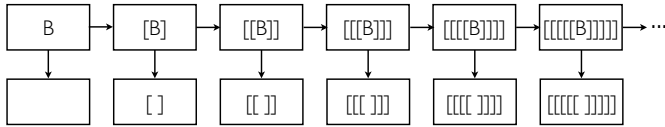
- Regular expressions don't allow recursion, just iteration
- But it is easy enough to give a simple recursive characterization for  $B \in \text{Brackets}$ :

$B \rightarrow \epsilon$  meaning: B is empty

$B \rightarrow [ B ]$  meaning: B consists of an initial [ symbol, another element from Brackets; and a closing ] symbol

54

## Generating brackets



- We have two rewrite rules:
  - $B \rightarrow \epsilon$  replace a B with the empty string
  - $B \rightarrow [B]$  replace a B with the string  $[B]$
- Either rule can be used to rewrite an occurrence of B
- We say that “B **derives** s” if the string s can be obtained from B by repeated rewriting/replacement

55

## Context-Free Grammars and Languages

56

## Context-free grammars (CFGs)

A context-free grammar  $G = (T, N, P, S)$  consists of

a set T of **terminal** symbols (“tokens”)

a set N of **nonterminal** symbols

a set P of **productions**, each of which is a rule of the form  $n \rightarrow w$  where  $n \in N$ , and  $w \in (T \cup N)^*$

a **start symbol**  $S \in N$

57

## Example

- A CFG for Brackets:  $(\{[, ]\}, \{B\}, \{B \rightarrow \epsilon, B \rightarrow [B]\}, B)$
- In practice, it is often sufficient just to write down the productions for a CFG:
  - the sets of terminals and nonterminals can usually be inferred from the productions
  - the start symbol can either be identified explicitly, or assumed as the first nonterminal
- For example, we can describe Brackets by:
  - $B \rightarrow \epsilon$
  - $B \rightarrow [B]$

58

## More examples

Prop:	$P \rightarrow \text{TRUE}$	Regexps:	$R \rightarrow c$
	$P \rightarrow \text{FALSE}$		$R \rightarrow \epsilon$
	$P \rightarrow \text{VAR}$		$R \rightarrow R R$
	$P \rightarrow (P)$		$R \rightarrow R \mid R$
	$P \rightarrow \text{AND } P P$		$R \rightarrow R *$
	$P \rightarrow \text{OR } P P$	CFGs:	$G \rightarrow$
	$P \rightarrow \text{NOT } P$		$G \rightarrow P ; G$
Arithmetic:	$E \rightarrow n$		$P \rightarrow n \rightarrow W$
	$E \rightarrow (E)$		$W \rightarrow$
	$E \rightarrow E + E$		$W \rightarrow n W$
	$E \rightarrow E * E$		$W \rightarrow t W$

59

## Context-free grammars and languages

- What is the relationship between context-free grammars and languages (i.e., sets of strings)?
- A **derivation** is a sequence of strings:
 
$$S_1 \rightarrow S_2 \rightarrow S_3 \rightarrow S_4 \rightarrow S_5 \rightarrow S_6 \rightarrow S_7 \rightarrow \dots$$
 in which each string  $s_{i+1}$  is obtained from the previous string  $s_i$  by choosing a production  $n \rightarrow w$  and replacing an occurrence of n in  $s_i$  with w
- In this case, we say that  $s_1$  **derives**  $s_i$  for each  $i = 1, 2, \dots$
- We say that a context-free grammar  $G = (T, N, P, S)$  **generates** the language that contains all strings in  $T^*$  that can be derived from S

60

## Brackets is a “context-free language”

- Any language that is generated from a context-free grammar is said to be a **context-free language**
- Sample derivations for Brackets:

$B \rightarrow$   
 $B \rightarrow [B] \rightarrow [ ]$   
 $B \rightarrow [B] \rightarrow [[B]] \rightarrow [[ ]]$   
 $\dots$   
 $B \rightarrow [B] \rightarrow [[B]] \rightarrow [[[B]]] \rightarrow [[[[B]]]] \rightarrow [[[[[ ]]]]]$   
 $\dots$

61

## Derivations and parse trees

62

## A language of arithmetic expressions

- Many computer languages are naturally described as context-free languages (i.e., using a context-free grammar)

- Example: a simple language of expressions:

$E \rightarrow n$  ( $n$  is an integer literal)  
 $E \rightarrow E + E$   
 $E \rightarrow E * E$   
 $E \rightarrow ( E )$

- Terminology:

$E$  is a nonterminal  
 $+$ ,  $*$ ,  $($ ,  $)$ , and  $n$  are terminals (i.e., tokens)

63

## Deriving expressions

For example,  $1 + (2 * 3)$  is an expression:

$E \rightarrow E + E$	$E \rightarrow E + E$
$\rightarrow E + (E)$	$\rightarrow 1 + E$
$\rightarrow E + (E * E)$	$\rightarrow 1 + (E)$
$\rightarrow E + (E * 3)$	$\rightarrow 1 + (E * E)$
$\rightarrow E + (2 * 3)$	$\rightarrow 1 + (2 * E)$
$\rightarrow 1 + (2 * 3)$	$\rightarrow 1 + (2 * 3)$

$E \rightarrow n$   
 $E \rightarrow E + E$   
 $E \rightarrow E * E$   
 $E \rightarrow ( E )$

one expression, multiple derivations

64

## Multiple derivations

- In a right-most derivation, we replace the right-most nonterminal at each step
- In a left-most derivation, we replace the left-most nonterminal at each step
- There are other choices between these extremes
- Does it matter which derivation we use?

65

## Deriving expressions

For example,  $1 + 2 * 3$  is an expression (no parentheses):

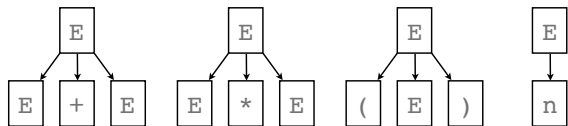
$E \rightarrow E + E$	$E \rightarrow E * E$
$\rightarrow E + E * E$	$\rightarrow E + E * E$
$\rightarrow E + E * 3$	$\rightarrow 1 + E * E$
$\rightarrow E + 2 * 3$	$\rightarrow 1 + 2 * E$
$\rightarrow 1 + 2 * 3$	$\rightarrow 1 + 2 * 3$

two valid derivations,  
but a fundamental difference this time

66

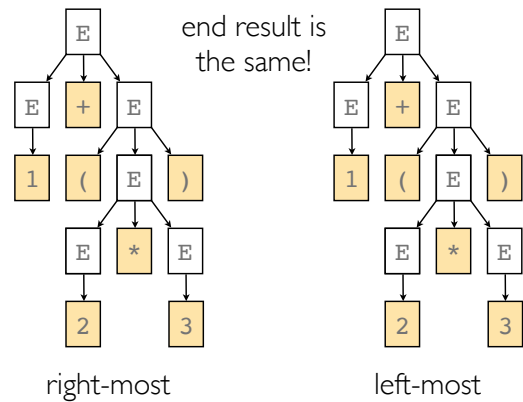
## Productions in graphical form

To understand the essential structure of a derivation, we will use a graphical notation to represent the productions in a grammar:



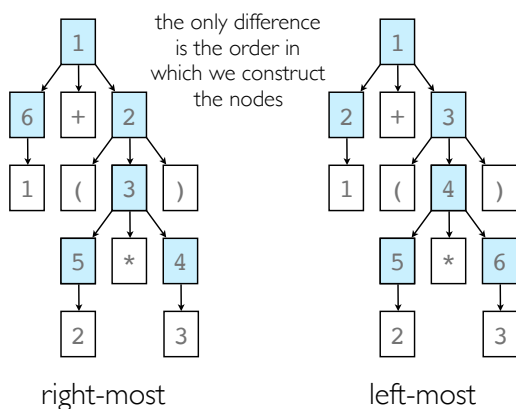
67

## $1+(2*3)$ revisited



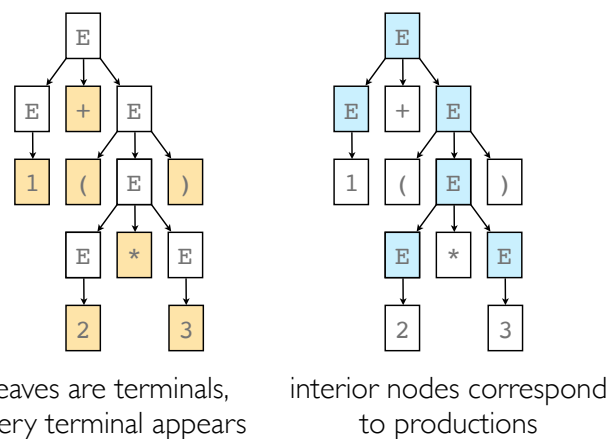
68

## Right-most vs left-most



69

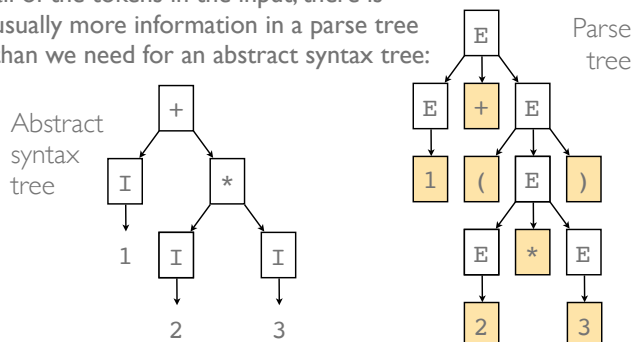
## Parse trees



70

## Parse trees vs abstract syntax trees

Because it describes a full derivation and all of the tokens in the input, there is usually more information in a parse tree than we need for an abstract syntax tree:



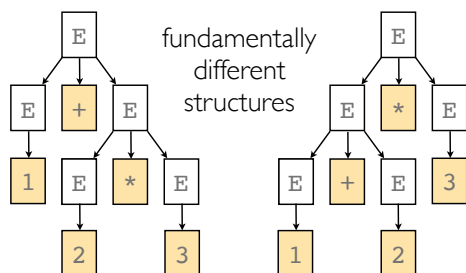
71

## CFGs and parse trees

- Context-free grammars don't just define languages (i.e., they don't just derive *sets of strings*) ...
- ... they actually define *sets of trees*!
- The strings in the corresponding context-free language can be recovered from the leaf nodes of each tree
- Parsing works in reverse: start with a string and try to construct the original tree
- What happens in there are distinct tree structures with the same sequence of symbols on their leaves?

72

## I+2\*3 revisited



73

## Ambiguity

- A grammar is **ambiguous** if there is a string in the corresponding language with more than one parse tree
- Example: Our grammar for expressions is ambiguous because the string "1+2\*3" has two distinct parse trees  
(We could find plenty of other examples of the same problem in this grammar, but finding just one is enough to demonstrate ambiguity)
- Ambiguity is a property of a grammar, NOT a language  
We can have multiple grammars describing the same language, some ambiguous and some unambiguous

74

## Dealing with ambiguity

- Does it matter?
- If any parse tree is as good as any other (i.e., they all have the same meaning), then just take whichever tree we get  
Example: for regular expressions,  $r_1(r_2r_3)$  matches exactly the same set of strings as  $(r_1r_2)r_3$ , so we can parse  $r_1r_2r_3$  either way
- If different trees have different meanings, then we need to choose between them:  
Disambiguating rules (e.g., operator precedence)  
Rewrite the grammar to avoid ambiguity

75

## Precedence and grouping (fixity)

- If  $\otimes$  has **higher precedence** than  $\oplus$ , then  $a \otimes b \oplus c$  should be parsed in the same way as  $(a \otimes b) \oplus c$
- If  $\otimes$  **groups/associates** to the left, then  $a \otimes b \otimes c$  should be parsed in the same way as  $(a \otimes b) \otimes c$
- If  $\otimes$  **groups/associates** to the right, then  $a \otimes b \otimes c$  should be parsed in the same way as  $a \otimes (b \otimes c)$
- If  $\otimes$  is **nonassociative**, then  $a \otimes b \otimes c$  should be treated as an error
- The combination of an operator's precedence and grouping is known as its **fixity**

76

## Order of operations

- There are widely used conventions for the precedence of standard arithmetic operations. (e.g., "PEMDAS")
- Parentheses first, then multiplications, then addition  
Example:  $(1+2)+3*4 == (1+2)+(3*4)$
- But the final decision about what fixity each operator should have rests with the language designer
- Let's suppose we want to resolve our ambiguities using:  
\* has higher precedence than +  
\* and + both group to the left

77

## An unambiguous grammar for expressions

Here is an unambiguous grammar for our language of expressions:

$E \rightarrow E + P$  expressions

$E \rightarrow P$

$P \rightarrow P * A$  products

$P \rightarrow A$

$A \rightarrow (E)$  atoms

$A \rightarrow n$  (n is an integer literal)

How were these properties met?  
\* has higher precedence than +  
\* and + both group to the left

78

## An unambiguous grammar for expressions

Here is an unambiguous grammar for our language of expressions:

$E \rightarrow E + P$       expressions are sums of products  
 $E \rightarrow P$   
  
 $P \rightarrow P * A$       products of atoms  
 $P \rightarrow A$   
  
 $A \rightarrow ( E )$       atoms  
 $A \rightarrow n$

79

## An unambiguous grammar for expressions

An expression is a sum of products, each of which is a product of atoms.

Example: if  $a_1, \dots, a_8$  are atoms, then:

$$\underbrace{a_1 * a_2 * a_3}_{\text{product}} + \underbrace{a_4 * a_5}_{\text{product}} + \underbrace{a_6}_{\text{product}} + \underbrace{a_7 * a_8}_{\text{product}}$$

$$\text{expr}$$

must be parsed as:

$(a_1 * a_2 * a_3) + (a_4 * a_5) + (a_6) + (a_7 * a_8)$

Multiply has been given a higher precedence than addition!

80

## An unambiguous grammar for expressions

The right argument of  $+$  must be a product

Example: if  $p_1, \dots, p_4$  are products, then:

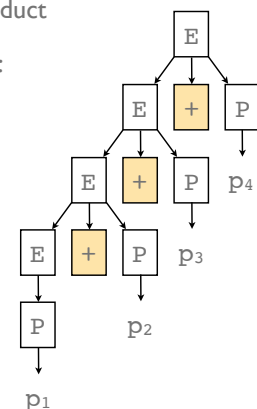
$$\underbrace{\underbrace{p_1}_{\text{prod.}} + \underbrace{p_2}_{\text{prod.}}}_{\text{expr}} + \underbrace{p_3}_{\text{prod.}} + \underbrace{p_4}_{\text{prod.}}$$

$$\text{expr}$$

must be parsed as:

$((p_1 + p_2) + p_3) + p_4$

In other words:  $+$  groups to the left



81

## Controlling grouping

$E \rightarrow E + P$   
 $E \rightarrow P$

recursion on the left  
results in left grouping

$P \rightarrow P * A$   
 $P \rightarrow A$

similarly for  $*$

for right grouping, change the recursion

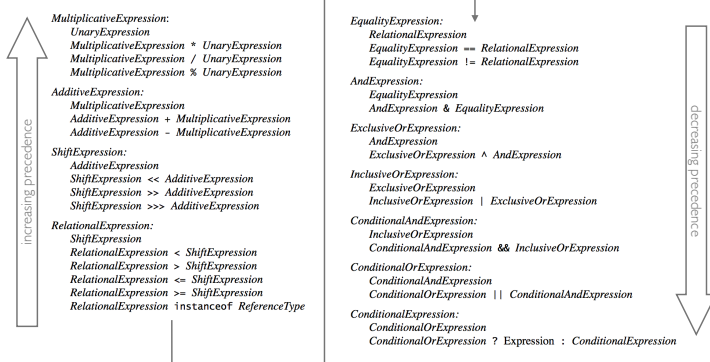
$E \rightarrow P + E$   
 $E \rightarrow P$

for nonassociative operators, no recursion on either side

$P \rightarrow A * A$   
 $P \rightarrow A$

82

## Fragments from the official Java grammar



83

## Context free languages summary

- Context free grammars can be used to describe a significantly larger family of languages than regular expressions
  - including many of the languages we encounter in practice
- Parse trees are graphical descriptions of derivations that
  - can reflect the grammatical structure of the input
  - can highlight ambiguities in the grammar
  - include more detail than is typically needed for an AST
- Grammars can be written so that operators are treated as having different precedence and grouping behaviors
  - This can help to avoid ambiguity

84