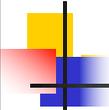


Programming Language Syntax

Read: Scott, Chapter 2.1 and 2.2

1



Announcements

- No class Friday this week
- HW1 will be out today
- Office hours schedule coming up
 - Check Submitty and course page:
<https://www.cs.rpi.edu/~milanova/csci4430/>

Programming Languages CSCI 4430, A Milanova

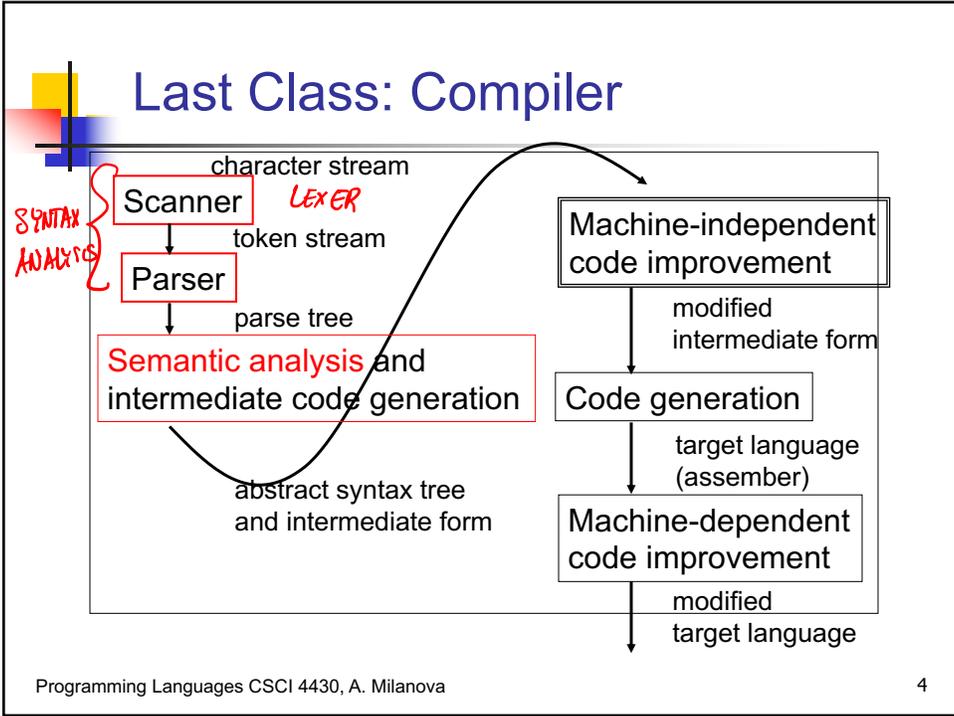
2

Lecture Outline

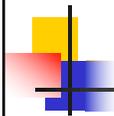
- *Formal languages*
- *Regular expressions*
- *Context-free grammars*
 - *Derivation*
 - *Parse*
 - *Parse trees*
 - *Ambiguity*
- *Expression grammars*

Programming Languages CSCI 4430, A. Milanova 3

3



4

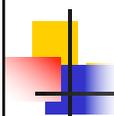


Syntax and Semantics

- **Syntax** is the form or structure of expressions, statements, and program units of a given language
 - Syntax of a Java **while** statement:
 - `while (boolean_expr) statement`
- **Semantics** is the meaning of expressions, statements and program units of a given language
 - Semantics of **while** (*boolean_expr*) *statement*
 - Execute *statement* repeatedly (0 or more times) as long as *boolean_expr* evaluates to `true`

5

5



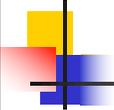
Formal Languages

- Theoretical foundations – Automata theory
- A **language** is a set of strings (also called sentences) over a finite alphabet
- A **generator** is a set of rules that generate the strings in the language
- A **recognizer** reads input strings and determines whether they belong to the language
- Languages are characterized by the complexity of generation/recognition rules

Programming Languages CSCI 4430, A. Milanova

6

6

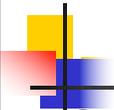


Question

- What are the classes of formal languages?
- The Chomsky hierarchy:
 - Regular languages
 - Context-free languages
 - Context-sensitive languages
 - Recursively enumerable languages

Programming Languages CSCI 4430, A. Milanova 7

7



Formal Languages

- Generators and recognizers become more complex as languages become more complex
 - Regular languages
 - Describe PL **tokens** (e.g., keywords, identifiers, numeric literals)
 - Generated by a **regular expression**
 - Recognized by a **finite automaton** (scanner)
 - Context-free languages
 - Describe more complex PL constructs (e.g., expressions and statements)
 - Generated by a **context-free grammar**
 - Recognized by a **push-down automaton** (parser)
 - Even more complex constructs

Programming Languages CSCI 4430, A. Milanova 8

8

Formal Languages

- Key application of formal languages: enable proof of relative difficulty of computational problems
- Our focus: formal languages provide the formalism for describing PL constructs
 - A compelling application of formal languages!
 - Building a scanner
 - Building a parser
 - Central issue: build efficient, linear-time parsers

Programming Languages CSCI 4430, A. Milanova 9

9

A Single Pass

*position = initial + rate * 60;*

Scanner

id = ...

Parser

- Scanner emits next token
- Parser consumes token and adds to parse tree
- Example shows bottom-up parse tree construction

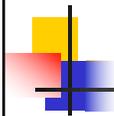
assign

```

graph TD
    assign --> id
    assign --> eq["="]
    assign --> expr
    expr --> dots["..."]
  
```

Programming Languages CSCI 4430, A. Milanova 10

10

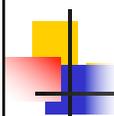


Lecture Outline

- *Formal languages*
- *Regular expressions*
- *Context-free grammars*
 - *Derivation*
 - *Parse*
 - *Parse trees*
 - *Ambiguity*
- *Expression grammars*

Programming Languages CSCI 4430, A. Milanova 11

11

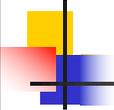


Regular Expressions

- Simplest structure
- Formalism to describe the simplest programming language constructs, the **tokens**
 - Each symbol (e.g., “+”, “-”) is a token
 - An identifier (e.g., position, rate, initial) is a token
 - A numeric constant (e.g., 60) is a token
- Recognized by a finite automaton

Programming Languages CSCI 4430, A. Milanova 12

12

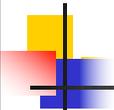


Regular Expressions

- A regular expression is one of the following:
 - Base case:
 - A character, e.g., a $\{a\}$
 - The empty string, denoted by ϵ $\{\epsilon\}$

Programming Languages CSCI 4430, A. Milanova 13

13



Regular Expressions

- A regular expression is one of the following:
 - Inductive case:
 - Two regular expressions next to each other: $R_1 R_2$
 - Meaning: $R_1 R_2$ generates the language of strings that are made up of any string generated by R_1 , followed by any string generated by R_2
 - Two regular expressions separated by |: $R_1 | R_2$
 - Meaning: $R_1 | R_2$ generates the language that is the union of the strings generated by R_1 with the strings generated by R_2
 - R^*

Programming Languages CSCI 4430, A. Milanova 14

14

Question

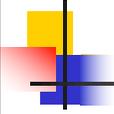
- What is the language defined by reg. exp.

$(a|b)(aa|bb)$? $(a|b)(aa|bb)$
 $\{a, b\}$ $\{aa, bb\}$
 $\{aaa, abb, baa, bbb\}$

- We saw concatenation and alternation. What operation is still missing?

Regular Expressions

- A Regular Expression is one of the following:
 - A character, e.g., a
 - The empty string ϵ
 - $R_1 R_2$
 - $R_1 | R_2$
 - Regular expression followed by a Kleene star, R^*
 - Meaning: the concatenation of zero or more strings generated by R
 - E.g., a^* generates $\{\epsilon, a, aa, aaa, \dots\}$
 - E.g., $(a|b)^*$ generates all strings of a 's and b 's

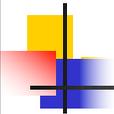


Regular Expressions

- Precedence
 - Kleene * has highest precedence $(ab|c)$ (ab^*)
 - Followed by concatenation $ab|c$ is $(ab)|c$
 - Followed by alternation |
- E.g., $ab|c$ is $(ab)|c$ not $a(b|c)$
 - Generates language $\{ab, c\}$ not $\{ab, ac\}$
- E.g., ab^* generates $\{a, ab, abb, \dots\}$ not $\{\epsilon, ab, abab, ababab, \dots\}$

Programming Languages CSCI 4430, A. Milanova 17

17



Question

- What is the language generated by regular expression $(0|1)^*1$?

ALL STRINGS OF 0'S AND 1'S THAT END IN 1.

0101
- What about $0^*(10^*10^*)^*$?

ALL STRINGS OF 0'S AND 1'S WITH EVEN NUMBER OF 1'S.

E.g. 11, 10001 are in the language. 1 is not.

Programming Languages CSCI 4430, A. Milanova 18

18

Regular Expressions in Programming Languages

- Describe tokens
- Let

$letter \rightarrow a|b|c| \dots |z$

$digit \rightarrow 1|2|3|4|5|6|7|8|9|0$

- Which token is this?

- $letter (letter | digit)^*$? IDENTIFIER
- $digit digit^*$? INTEGER CONSTANT (NON NEG)
- $digit^* . digit digit^*$? REAL CONSTANT (NON NEG)

Regular Expressions in Programming Languages

- Which token is this:

$6.02E+23 \approx 6.02 \times 10^{23}$
 decimal exponent

6.02

$number \rightarrow integer | real$ / 6.02 is recognized by "decimal" 's automaton

$real \rightarrow integer exponent | decimal (exponent | \epsilon)$

$decimal \rightarrow digit^* (. digit | digit .) digit^*$

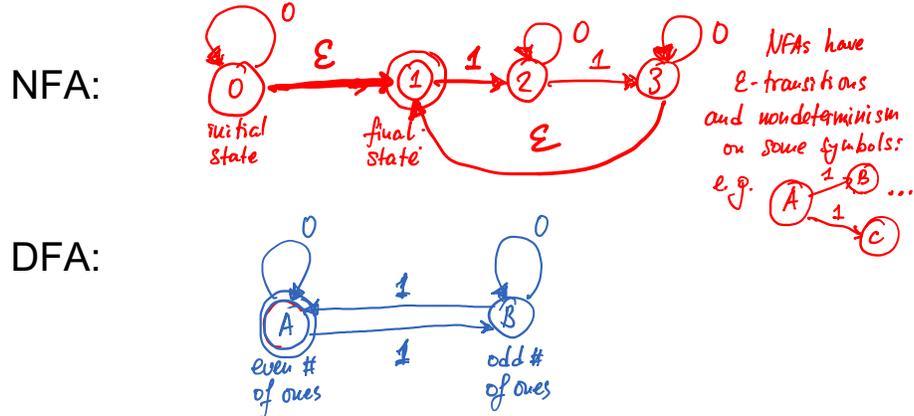
$exponent \rightarrow (e | E) (+ | - | \epsilon) integer$

$integer \rightarrow digit digit^*$

$digit \rightarrow 1|2|3|4|5|6|7|8|9|0$

Exercises

Regular expression: $0^*(10^*10^*)^*$

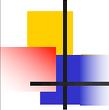


21

Lecture Outline

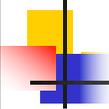
- Formal languages
- Regular expressions
- **Context-free grammars**
 - Derivation
 - Parse
 - Parse trees
 - Ambiguity
- Expression grammars

22



Context-Free Grammars

- Unfortunately, regular languages cannot specify all constructs in programming
- E.g., can we write a regular expression that specifies valid arithmetic expressions?
 - `id * (id + id * (number - id))`
(Handwritten annotations: blue underlines under 'id', 'id', and 'number'; red 'open' and 'close' labels with arrows pointing to the parentheses.)
 - Among other things, we need to ensure that parentheses are matched!
 - Answer is no. We need **context-free languages** and context-free grammars!



Grammar

- A grammar is a formalism to describe the strings of a (formal) language
- A grammar consists of a set of **terminals**, set of **nonterminals**, a set of **productions**, and a **start symbol**
 - **Terminals** are the characters in the alphabet
 - **Nonterminals** represent language constructs
 - **Productions** are rules for forming syntactically correct constructs
 - **Start symbol** tells where to start applying the rules

Notation

$A \rightarrow A0|B0$
 $B \rightarrow A1|B1$
 $A \rightarrow \epsilon$
 $B \rightarrow \epsilon$

Specification of identifier:

Regular expression: *letter* (*letter* | *digit*)*

DESCRIPTIVE NAMES FOR NONTERMINALS

BNF: $\langle digit \rangle ::= 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0$
 $\langle letter \rangle ::= a | b | c | \dots | x | y | z$
 $\langle id \rangle ::= \langle letter \rangle | \langle id \rangle \langle letter \rangle | \langle id \rangle \langle digit \rangle$

Textbook and slides: *digit* $\rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0$
letter $\rightarrow a | b | c | d | \dots | z$
id $\rightarrow letter | id letter | id digit$

Nonterminals shown in *italic*

Terminals shown in **typewriter**

Programming Languages CSCI 4430, A. Milanova 25

25

Regular Grammars

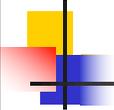
- Regular grammars generate regular languages
- The rules in regular grammars are of the form:
 - Each left-hand-side (lhs) has exactly one nonterminal
 - Each right-hand-side (rhs) is one of the following
 - A single terminal symbol or
 - A single nonterminal symbol or
 - A nonterminal followed by a terminal

$\leadsto 1222$

e.g., $1 2^+ | 0^+$ $S \rightarrow A | B$
 $A \rightarrow 1 | A 2$ "A" GENERATES 12^*
 $B \rightarrow 0 | B 0$ "B" GENERATES 00^*

Programming Languages CSCI 4430, A. Milanova 26

26



Question

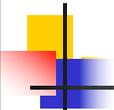
- Is this a regular grammar:
 - $S \rightarrow 0 A$
 - $A \rightarrow S 1$
 - $S \rightarrow \epsilon$

S ⇒ 0A ⇒ 0S1 ⇒ 01
(GENERATED STRING 01)

- No, this is a context-free grammar
 - It generates $0^n 1^n$, the canonical example of a context-free language
 - rhs should be nonterminal followed by a terminal, thus, $S \rightarrow 0 A$ is not a valid production

Programming Languages CSCI 4430, A. Milanova 27

27

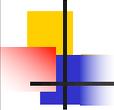


Lecture Outline

- *Formal languages*
- *Regular expressions*
- *Context-free grammars*
 - *Derivation*
 - *Parse*
 - *Parse trees*
 - *Ambiguity*
- *Expression grammars*

Programming Languages CSCI 4430, A. Milanova 28

28



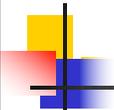
Context-free Grammars (CFGs)

- **Context-free grammars** generate context-free languages
 - Most of what we need in programming languages can be specified with CFGs
- Context-free grammars have rules of the form:
 - Each left-hand-side has exactly one nonterminal
 - Each right-hand-side contains an arbitrary sequence of terminals and nonterminals
- A context-free grammar
 - e.g., $0^n 1^n, n \geq 1$

$S \rightarrow 0 S 1$ *A SINGLE NONTERMINAL AT LHS.*
 $S \rightarrow 0 1$ *A SEQUENCE OF TERMINALS AND NONTERMINALS AT RHS.*

29

29



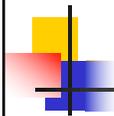
Question

- Examples of non-context-free languages?
 - E.g., $a^n b^m c^n d^m$ $n \geq 1, m \geq 1$
 - E.g., $w c w$ where w is in $(0 | 1)^*$
 - E.g., $a^n b^n c^n$ $n \geq 1$ (canonical example)

Programming Languages CSCI 4430, A. Milanova

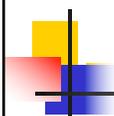
30

30



Context-free Grammars

- Can be used to generate strings in the context-free language (**derivation**)
- Can be used to recognize well-formed strings in the context-free language (**parse**)
- In programming languages and compilers, we are concerned with two special CFGs, called LL and LR grammars



Derivation

Simple context-free grammar for expressions:

$expr \rightarrow id \mid (expr) \mid expr \ op \ expr$

$op \rightarrow + \mid *$

A derivation generates (derives) strings. E.g.,:

$expr \Rightarrow expr \ op \ expr$

$\Rightarrow expr \ op \ id$

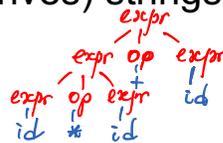
$\Rightarrow expr \ + \ id$

$\Rightarrow expr \ op \ expr \ + \ id$

$\Rightarrow expr \ op \ id \ + \ id$

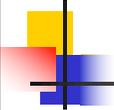
$\Rightarrow expr \ * \ id \ + \ id$

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



← **sentential form**

← **sentence, string or yield**

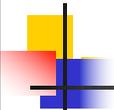


Derivation

- A **derivation** starts from the start symbol and at each step replaces a nonterminal with the right-hand side of a production for that nonterminal
 - E.g., $expr\ op\ \underline{expr}$ derives $expr\ op\ id$
We replaced the right (underlined) $expr$ with id due to production $expr \rightarrow id$
- An intermediate sentence is called a **sentential form**
 - E.g., $expr\ op\ id$ is a sentential form

Programming Languages CSCI 4430, A. Milanova 33

33

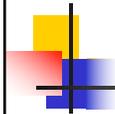


Derivation

- Resulting sentence is called **yield**
 - E.g., $id*id+id$ is the yield of our derivation
- What is a **leftmost derivation**?
 - Replaces the **leftmost** nonterminal in the sentential form at each step
- What is a **rightmost derivation**?
 - Replaces the **rightmost** nonterminal in the sentential form at each step
- Some derivations are neither left- nor rightmost

Programming Languages CSCI 4430, A. Milanova 34

34

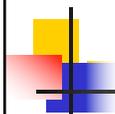


Question

- What kind of derivation is this:

$expr \Rightarrow expr\ op\ \underline{expr}$
 $\Rightarrow expr\ \underline{op}\ id$
 $\Rightarrow \underline{expr} + id$
 $\Rightarrow expr\ op\ \underline{expr} + id$
 $\Rightarrow expr\ \underline{op}\ id + id$
 $\Rightarrow \underline{expr} * id + id$
 $\Rightarrow id * id + id$

- A rightmost derivation. At each step we replace the rightmost nonterminal

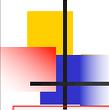


Question

- What kind of derivation is this:

$expr \Rightarrow expr\ op\ \underline{expr}$
 $\Rightarrow expr\ \underline{op}\ id$
 $\Rightarrow \underline{expr} + id$
 $\Rightarrow \underline{expr}\ op\ expr + id$
 $\Rightarrow id\ op\ \underline{expr} + id$
 $\Rightarrow id\ \underline{op}\ id + id$
 $\Rightarrow id * id + id$

- Neither leftmost nor rightmost



Parse

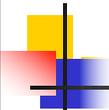
PROCESS OF RECOGNIZING A STRING
AS BELONGING TO THE LANGUAGE

$expr \rightarrow id \mid (expr) \mid expr \ op \ expr$
 $op \rightarrow + \mid *$

- A parse takes a string and constructs a derivation (if one exists). E.g.:
 - $\underline{id} \ * \ id \ + \ id \Rightarrow \underline{expr} \ * \ id \ + \ id$
 - $\Rightarrow \underline{expr} \ op \ \underline{id} \ + \ id$
 - $\Rightarrow \underline{expr} \ op \ \underline{expr} \ + \ id$
 - $\Rightarrow \underline{expr} \ + \ id$
 - $\Rightarrow \underline{expr} \ op \ \underline{id}$
 - $\Rightarrow \underline{expr} \ op \ \underline{expr}$
 - $\Rightarrow \underline{expr}$

Programming Languages CSCI 4430, A. Milanova 37

37



Parse

- A bottom-up parse builds parse tree bottom-up, i.e., from leaves towards root (start symbol)
- Example was a bottom-up parse:
 - $\dots \Rightarrow \underline{expr \ op \ expr} \ + \ id$
 - $\Rightarrow \underline{expr} \ + \ id$

Replacing $expr \ op \ expr$ with $expr$ corresponds to adding node $expr$ to tree
- A top-down parse builds parse tree top-down, i.e., from root towards leaves (a lot more on this later)

Programming Languages CSCI 4430, A. Milanova 38

38

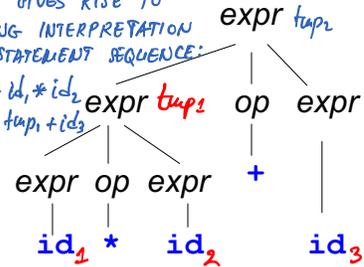
Parse Tree

$expr \rightarrow id \mid (expr) \mid expr \ op \ expr$
 $op \rightarrow + \mid *$

$expr \Rightarrow expr \ op \ expr$
 $\Rightarrow expr \ op \ id$
 $\Rightarrow expr \ + \ id$
 $\Rightarrow expr \ op \ expr \ + \ id$
 $\Rightarrow expr \ op \ id \ + \ id$
 $\Rightarrow expr \ * \ id \ + \ id$
 $\Rightarrow id \ * \ id \ + \ id$

PARSE TREE GIVES RISE TO
THE FOLLOWING INTERPRETATION
INTO A STATEMENT SEQUENCE:

$t_{op1} = id_1 * id_2$
 $t_{op2} = t_{op1} + id_3$



Internal nodes are nonterminals. Children are the rhs of a rule for that nonterminal. Leaf nodes are terminals.

39

39

Ambiguity

- A grammar is **ambiguous** if some string can be generated by two or more distinct parse trees
 - There is no algorithm that can tell if an arbitrary context-free grammar is ambiguous
- Ambiguity arises in programming language grammars
 - Arithmetic expressions
 - If-then-else: the dangling else problem
- Ambiguity is bad

Programming Languages CSCI 4430, A. Milanova

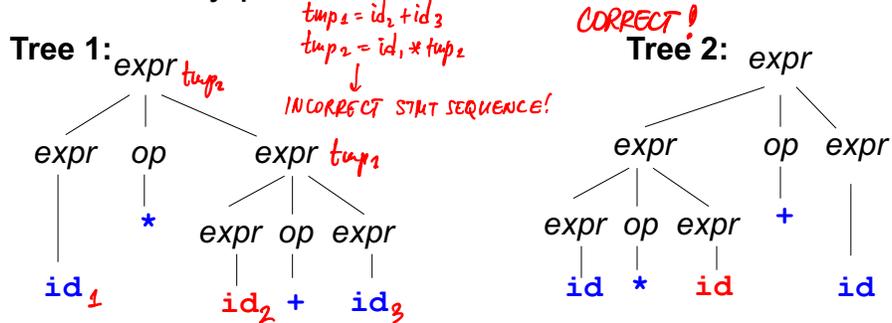
40

40

Ambiguity

$expr \rightarrow id \mid (expr) \mid expr \ op \ expr$
 $op \rightarrow + \mid *$

- How many parse trees for $id * id + id$?



- Which one is "correct"? *Tree 2.*

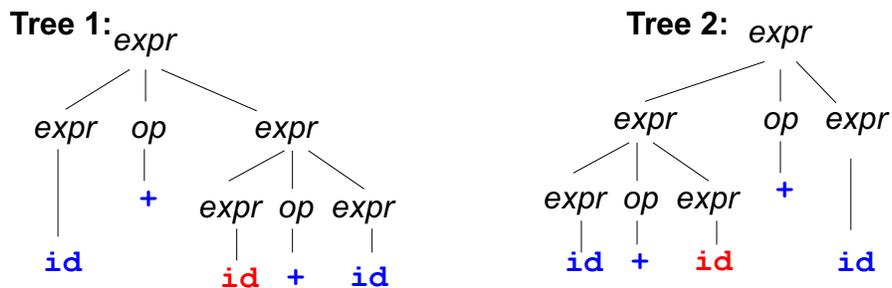
41

41

Ambiguity

$expr \rightarrow id \mid (expr) \mid expr \ op \ expr$
 $op \rightarrow + \mid *$

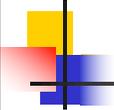
- How many parse trees for $id + id + id$?



- Which one is "correct"? *TREE 2. + and * are left associative.*

42

42

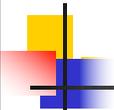


Lecture Outline

- *Formal languages*
- *Regular expressions*
- *Context-free grammars*
 - *Derivation*
 - *Parse*
 - *Parse trees*
 - *Ambiguity*
- *Expression grammars*

Programming Languages CSCI 4430, A. Milanova 43

43



Expression Grammars

- Generate expressions
 - Arithmetic expressions
 - Regular expressions
 - Other
- Terminals: operands, operators, and parentheses

$$\text{expr} \rightarrow \text{id} \mid (\text{expr}) \mid \text{expr op expr}$$
$$\text{op} \rightarrow + \mid *$$

Programming Languages CSCI 4430, A. Milanova 44

44

Handling Ambiguity

Our ambiguous grammar, slightly simplified:

$$\text{expr} \rightarrow \text{id} \mid (\text{expr}) \mid \text{expr} + \text{expr} \mid \text{expr} * \text{expr}$$

- Rewrite the grammar into unambiguous one:

$$\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term}$$
$$\text{term} \rightarrow \text{term} * \text{factor} \mid \text{factor}$$
$$\text{factor} \rightarrow \text{id} \mid (\text{expr})$$

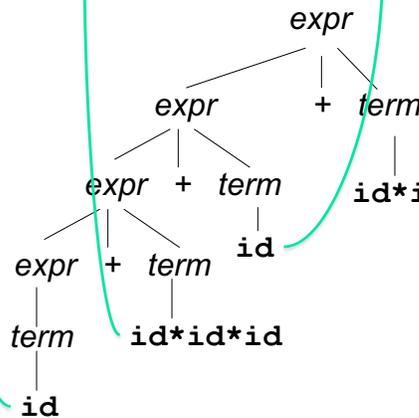
Rewriting Expression Grammars: Intuition

$$\text{expr} \rightarrow \text{id} \mid (\text{expr}) \mid \text{expr} + \text{expr} \mid \text{expr} * \text{expr}$$

- A new nonterminal, *term*
- $\text{expr} * \text{expr}$ becomes *term*. Thus, $*$ gets pushed down the tree, forcing higher precedence of $*$
- $\text{expr} + \text{expr}$ becomes $\text{expr} + \text{term}$. Pushes leftmost $+$ down the tree, forcing operand to associate with $+$ on its left
 - $\text{expr} \rightarrow \text{expr} + \text{expr}$ becomes $\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term}$

Rewriting Expression Grammars: Intuition

E.g., look at $\text{id} + \text{id} * \text{id} * \text{id} + \text{id} + \text{id} * \text{id}$



Programming Languages CSCI 4430, A. Milanova

47

47

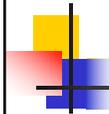
Rewriting Expression Grammars: Intuition

- To parse terms, we add a new nonterminal, *factor*, and productions:
 - $\text{term} \rightarrow \text{term} * \text{factor} \mid \text{factor}$
 - $\text{factor} \rightarrow \text{id} \mid (\text{expr})$

Programming Languages CSCI 4430, A. Milanova

48

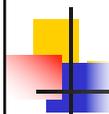
48



Exercise

$expr \rightarrow expr \times expr \mid expr \wedge expr \mid id$

- How many parse trees for $id \times id \wedge id \times id$?
 - No need to draw them all
- Rewrite above grammar into an equivalent unambiguous grammar such that
 - ▲ has higher precedence than \times
 - ▲ is **right-associative**
 - ▲ is left-associative



The End