



Programming Language Syntax

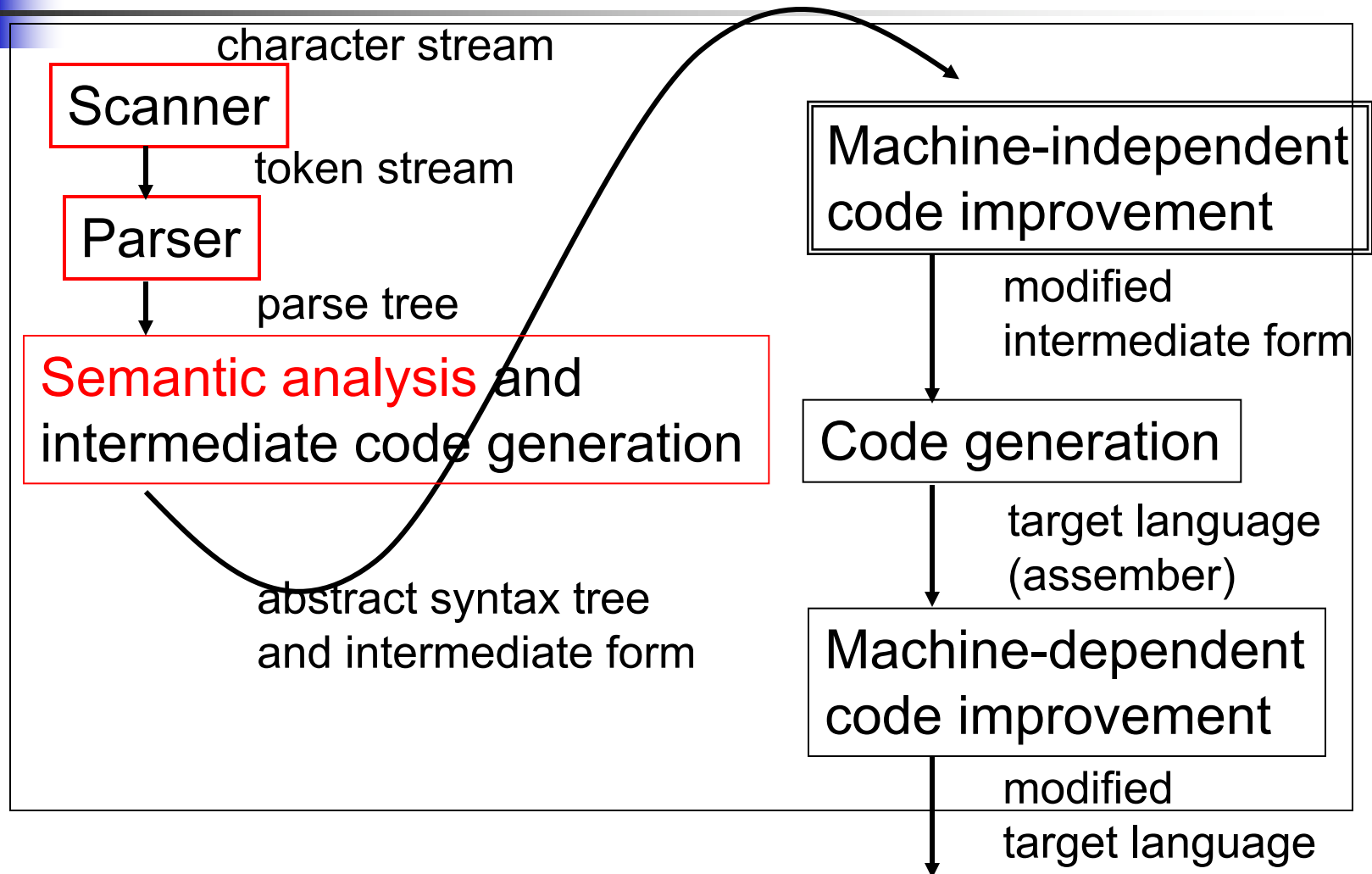
Read: Scott, Chapter 2.1



Lecture Outline

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

Last Class: Compiler





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`



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
 - E.g., regular languages
 - E.g., context-free languages



Question

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



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 **Regular Expressions**
 - 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



Formal Languages

- Main 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

A Single Pass

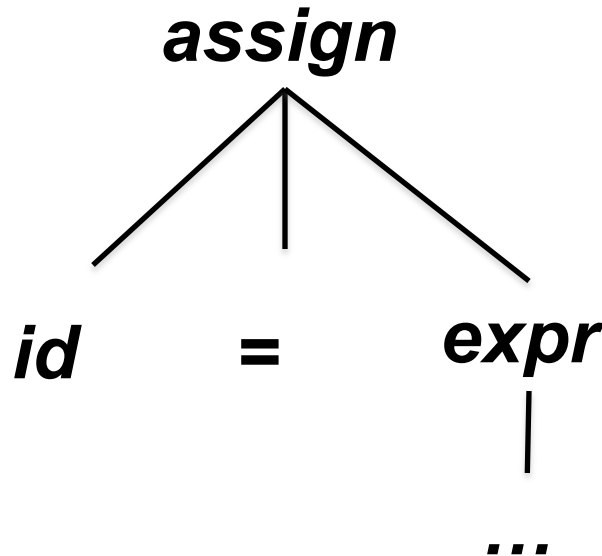
*position = initial + rate * 60;*

Scanner

- Scanner emits next token
- Parser consumes the token and continues building the parse tree (typically bottom up)

id = ...

Parser





Lecture Outline

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



Regular Expressions

- Simplest structure
- Formalism to describe the simplest programming language constructs, the **tokens**
 - each symbols (e.g., “+”, “-”) is a token
 - an identifier (e.g., position, rate, initial) is a token
 - a numeric constant (e.g., 59) is a token
 - etc.
- Recognized by a finite automaton



Regular Expressions

- A Regular Expression is one of the following:
 - A character, e.g., **a**
 - The empty string, denoted by **ϵ**
 - 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**



Question

- What is the language defined by reg. exp.
 $(a \mid b) (a a \mid b b) ?$
- 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, denoted by ϵ
 - $R_1 R_2$
 - $R_1 \mid R_2$
 - 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
- Followed by concatenation
- Followed by alternation |
- E.g., **a b | c** is **(a b) | c** not **a (b | c)**
 - Generates {ab, c} not {ab, ac}
- E.g., **a b*** generates {a, ab, abb, ...} not { ϵ , ab, abab, ababab, ...}



Question

- What is the language defined by regular expression $(0 \mid 1)^* 1$?
- What about $0^* (1 0^* 1 0^*)^* ?$

Regular Expressions in Programming Languages

- Describe tokens

- Let

letter \rightarrow **a|b|c| ... |z**

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

- Which token is this?

1. *letter (letter | digit)** ?

2. *digit digit ** ?

3. *digit *. digit digit ** ?

Regular Expressions in Programming Languages

- Which token is this:

number \rightarrow *integer* | *real*

real \rightarrow *integer exponent* | *decimal (exponent | ϵ)*

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

exponent \rightarrow (*e* | *E*) (*+* | *-* | ϵ) *integer*

integer \rightarrow *digit digit**

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



Lecture Outline

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



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))`
 - 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

Specification of identifier:

Regular expression: $\textit{letter} (\textit{letter} | \textit{digit})^*$

BNF: $\langle \textit{digit} \rangle ::= 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0$

$\langle \textit{letter} \rangle ::= a | b | c | \dots | x | y | z$

$\langle \textit{id} \rangle ::= \langle \textit{letter} \rangle | \langle \textit{id} \rangle \langle \textit{letter} \rangle | \langle \textit{id} \rangle \langle \textit{digit} \rangle$

Textbook and slides:
(also BNF)

Nonterminals shown in *italic*

digit → 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0

letter → a | b | c | d | ... | z

id → *letter* | *id letter* | *id digit*

Terminals shown in **typewriter**



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

e.g., $1\ 2^* \mid 0^+$

$$S \rightarrow A \mid B$$
$$A \rightarrow 1 \mid A\ 2$$
$$B \rightarrow 0 \mid B\ 0$$



Question

- Is this a regular grammar:

$$S \rightarrow 0 A$$

$$A \rightarrow S 1$$

$$S \rightarrow \varepsilon$$

- 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



Lecture Outline

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



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. $0^n 1^n, n \geq 1$
$$\begin{array}{l} S \rightarrow 0 S 1 \\ S \rightarrow 0 1 \end{array}$$



Question

- Examples of a 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)



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:

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

We can generate (derive) expressions:

$$\text{expr} \Rightarrow \text{expr op } \underline{\text{expr}}$$
$$\Rightarrow \text{expr } \underline{\text{op}} \text{ id}$$
$$\Rightarrow \underline{\text{expr}} + \text{id}$$
$$\Rightarrow \text{expr op } \underline{\text{expr}} + \text{id} \quad \longleftarrow \text{ sentential form}$$
$$\Rightarrow \text{expr } \underline{\text{op}} \text{ id} + \text{id}$$
$$\Rightarrow \underline{\text{expr}} * \text{id} + \text{id}$$
$$\Rightarrow \text{id} * \text{id} + \text{id} \quad \longleftarrow \text{ sentence, string or yield}$$



Derivation

- A **derivation** is the process that starts from the start symbol, and at each step, replaces a nonterminal with the right-hand-side of a production
 - E.g., *expr op 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



Derivation

- The resulting sentence is called **yield**
 - E.g., **id*id+id** is the yield of our derivation
- What is a **left-most derivation**?
 - Replaces the **left-most** nonterminal in the sentential form at each step
- What is a **right-most derivation**?
 - Replaces the **right-most** nonterminal in the sentential form at each step
- There are derivations that are neither left- nor right-most



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 right-most derivation. At each step we replace the right-most 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 left-most nor right-most



Parse

Recall our context-free grammar for expressions:

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

- A parse is the reverse of a derivation

$$\begin{aligned} \text{id} * \text{id} + \text{id} &\Rightarrow \text{expr} \underline{*} \text{id} + \text{id} \\ &\Rightarrow \text{expr op} \underline{\text{id}} + \text{id} \\ &\Rightarrow \underline{\text{expr op expr}} + \text{id} \\ &\Rightarrow \text{expr} \underline{+} \text{id} \\ &\Rightarrow \text{expr op} \underline{\text{id}} \\ &\Rightarrow \underline{\text{expr op expr}} \\ &\Rightarrow \text{expr} \end{aligned}$$



Parse

- A parse starts with the string of terminals, and at each step, replaces the right-hand-side (rhs) of a production with the left-hand-side (lhs) of that production. E.g.,

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

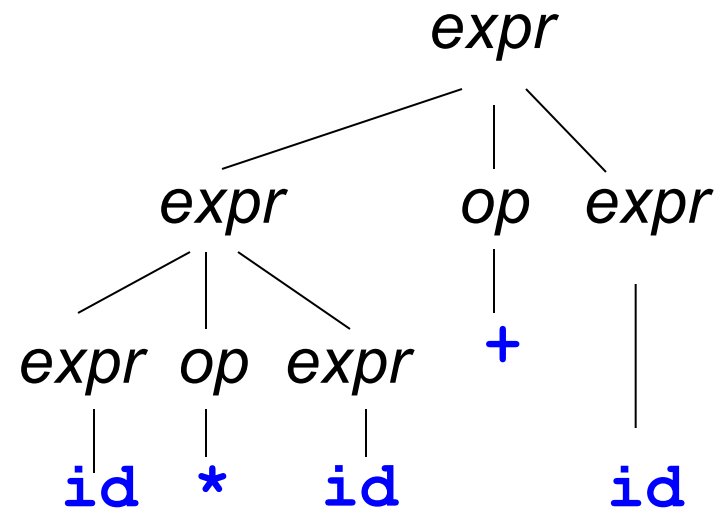
Here we replaced $expr\ op\ expr$ (the rhs of production $expr \rightarrow expr\ op\ expr$) with $expr$ (the lhs of the production)

Parse Tree

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

$op \rightarrow + \mid *$

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



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



Ambiguity

■ 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

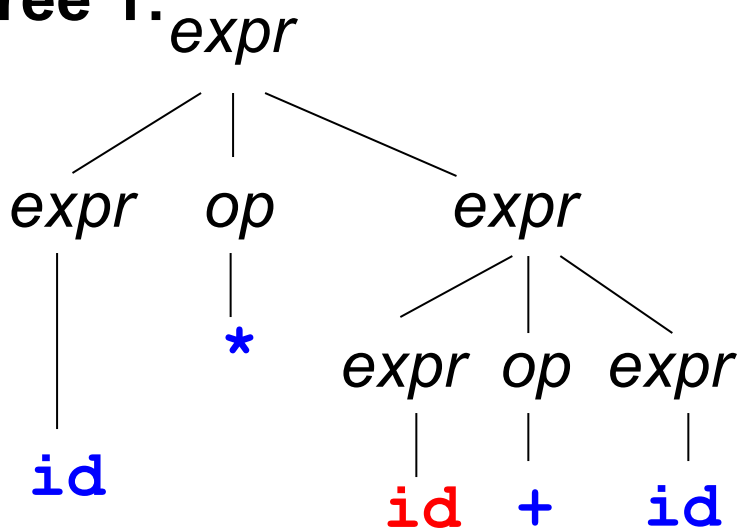
Ambiguity

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

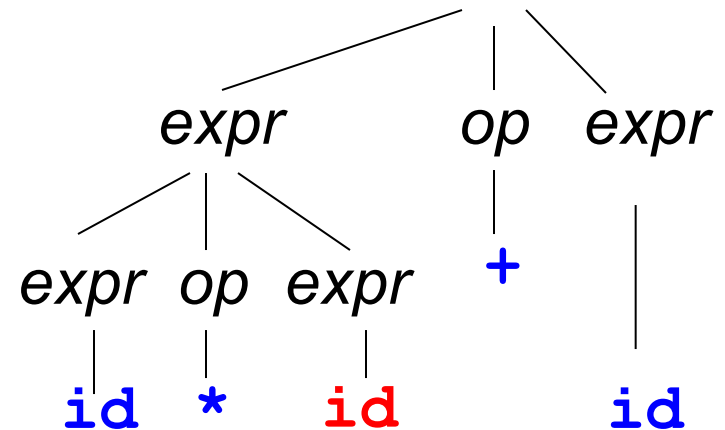
$op \rightarrow + \mid *$

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

Tree 1:



Tree 2:



- Which one is “correct”?

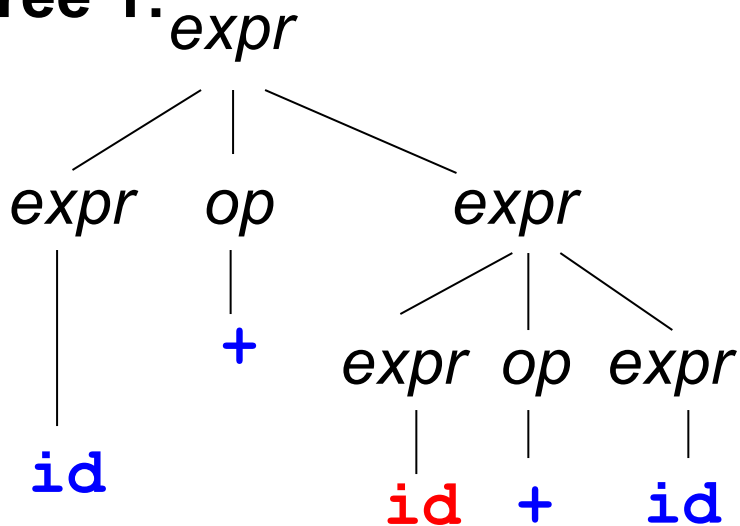
Ambiguity

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

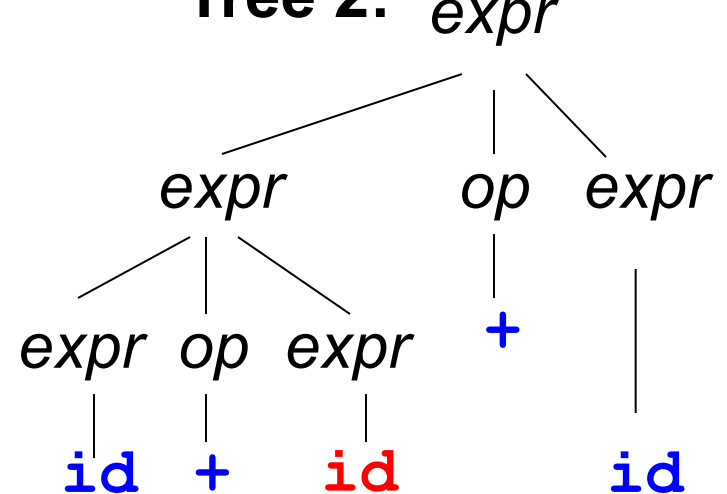
$op \rightarrow + \mid *$

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

Tree 1:



Tree 2:



- Which one is “correct”?



Lecture Outline

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



Expression Grammars

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

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

$op \rightarrow + \mid *$

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:

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

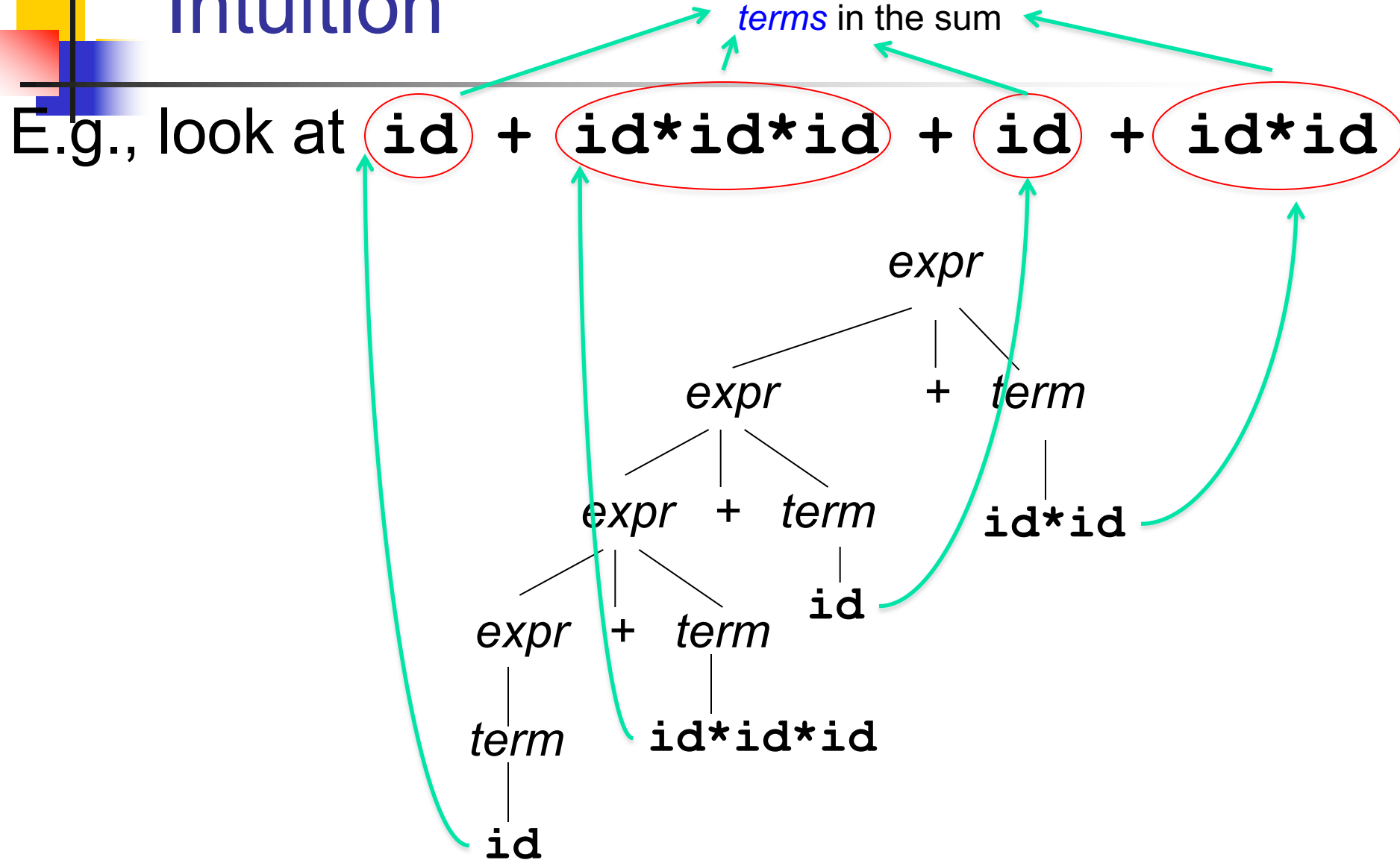
- Forces left associativity of $+$ and $*$
- Forces higher precedence of $*$ over $+$

Rewriting Expression Grammars: Intuition

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

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

Rewriting Expression Grammars: Intuition



Rewriting Expression Grammars: Intuition

- Another new nonterminal, *factor* and productions:
 - $term \rightarrow term * factor \mid factor$
 - $factor \rightarrow id \mid (expr)$



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 this grammar into an equivalent unambiguous grammar where
 - ^ has higher precedence than \times
 - ^ is right-associative
 - \times is left-associative



The End
