

CS 314 Principles of Programming Languages

Guang Wang

Department of Computer Science

Rutgers University

1. What do you expect to be covered in recitations?
2. I will answer your questions ASAP, normally, in 24 hours.
3. Two parts:
 - a. go over(understand the conceptions better)
 - b. problem sets/non-graded quizzes(prepare for the midterms and the final)

Why this course?

- **Why Learn More than One Programming Language?**

Each language encourages thinking about a problem in a particular way.

- Each language provides (slightly) different expressiveness & efficiency.
- \Rightarrow The language should match the problem

“No free lunch!”

Imperative Paradigm

A program is: A sequence of state-changing actions

Manipulate an abstract machine with:

- – Variables that name **memory locations**
- – Arithmetic and logical **operations**
- – Evaluate and assign operations
- – Explicit **control flow** statements

Key operations: *Assignment, Call*

- – also *Go To, Go To if 0*
- – or *If, While*

Functional Paradigm

A program is: Composition of functions on data

- **Characteristics**
 - Name values, not memory locations
 - Value binding through **parameter passing**
 - Recursion rather than iteration
- **Key operations: *Function Application, Function Abstraction***
 - Based on the Lambda Calculus

Eg: **Scheme**

Logic Paradigm

A program is: Formal logic specification of the problem

- **Characteristics**

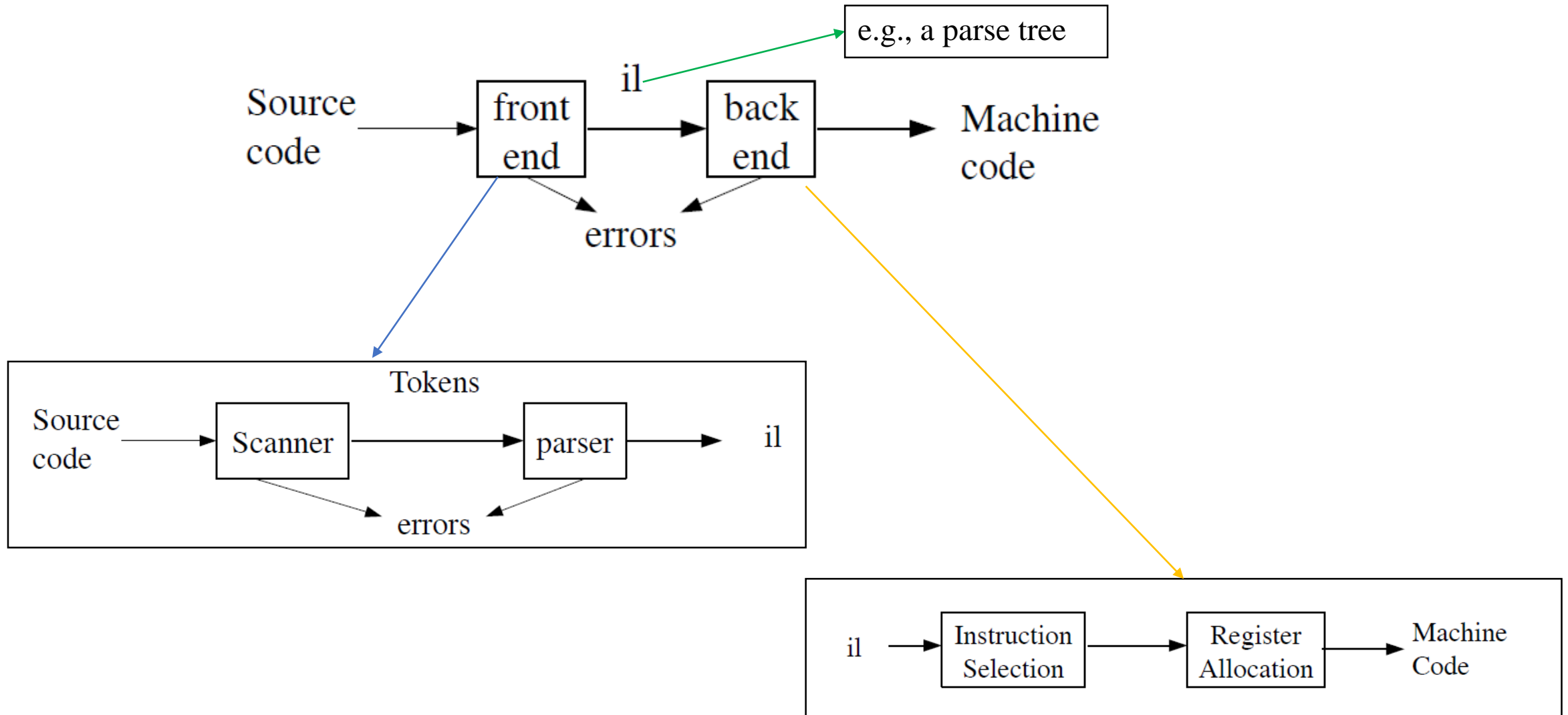
- Programs say *what* properties the solution must have, not *how* to find it
- Solutions are obtained through a specialized form of *theorem-proving*

- **Key operations: *Unification, NonDeterministic Search***

- Based on First Order Predicate Logic

Eg: **Prolog rules**

Compiler



Scanner

- **Maps** characters \rightarrow tokens
 - Tokens: basic unit of syntax
 - E.g., $x = x + y$ becomes a list of tokens:
 $\langle \text{id}, x \rangle \langle \text{operator}, \text{assign} \rangle \langle \text{id}, x \rangle \langle \text{operator}, + \rangle \langle \text{id}, y \rangle$
 - Typical *token* types:
number, **id**, **operator** (e.g., +), **keyword** (e.g., do, else)

Parser

- **Parse**: determine the grammatical structure of a sequence of tokens
- **Grammar**: set of rule like
 - assignment \rightarrow variable '=' expression ';'
- Analogy with grammars of human languages.
 - e.g., English: Sentence \rightarrow Subject Verb Object

Defining a language

- 1. Syntax
 - tokens defined by a *Regular Expression* or *Finite State Automaton*
 - Larger scale structure defined by a *Context Free Grammar*
- 2. Semantics
 - several formal approaches but in practice we just use English to explain the meaning

Formal languages

- Alphabet: a set of symbols
 - {a, b} or {if, while, for}
- String: a sequence of symbols from the alphabet
 - abbab
- **Language**: a set of strings
 - {ab, abab, ababab, ...}
 - Each string: **finite length**
 - Set: possibly **infinite number** of strings

Grammar

- A set of **Terminal Symbols**

a,b

- A set of **Non-terminal Symbols**

MultiAB,AB

- A set of **Rules of the form**

Non-terminal \Rightarrow sequence of Symbols

MultiAB \Rightarrow MultiAB AB

MultiAB \Rightarrow AB

AB \Rightarrow a b

Grammar Rule(Production)

A rule is of the form

Left-side \Rightarrow Right-side

Left-side is a **single Nonterminal**,

Right-side is a **sequence of Terminals and/or Non-terminals**

The Language of a Grammar

Any string you can produce by the following process:

Initialize a list of symbols to be just the **Start Symbol**

Repeatedly:

- Find the **first Non-terminal** in the list
- Find **a rule** whose Left-hand side is this non-Terminal
- Replace the Non-terminal with the **Right-hand side** of that rule

Until the list contains only Terminals

The Language of a Grammar

Grammar GAB:

MultiAB \Rightarrow MultiAB AB

MultiAB \Rightarrow AB

AB \Rightarrow a b



Start symbol

Example:

Proof that **a b a b a b a b is in L(GAB):**

MultiAB

MultiAB AB

MultiAB AB AB

AB AB AB

a b AB AB

a b a b AB

a b a b a b

Derivation in a Grammar

- Example: Write a leftmost derivation of **aabb** in G2?

G2:

$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \mid \langle A \rangle \langle B \rangle$

$\langle A \rangle \Rightarrow \mathbf{a} \mid \mathbf{a} \langle A \rangle$

$\langle B \rangle \Rightarrow \mathbf{b} \mid \langle B \rangle \mathbf{b}$

Answer:

$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \langle B \rangle$

$\Rightarrow \mathbf{a} \langle A \rangle \langle B \rangle$

$\Rightarrow \mathbf{aa} \langle B \rangle$

$\Rightarrow \mathbf{aa} \langle B \rangle \mathbf{b}$

$\Rightarrow \mathbf{aabb}$

Q: How about derivation of **aaab** in G2? **abbb?? baa??**

Answer:

$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \langle B \rangle$
 $\Rightarrow \mathbf{a} \langle A \rangle \langle B \rangle$
 $\Rightarrow \mathbf{aa} \langle A \rangle \langle B \rangle$
 $\Rightarrow \mathbf{aa} \mathbf{a} \langle B \rangle$
 $\Rightarrow \mathbf{aaab}$

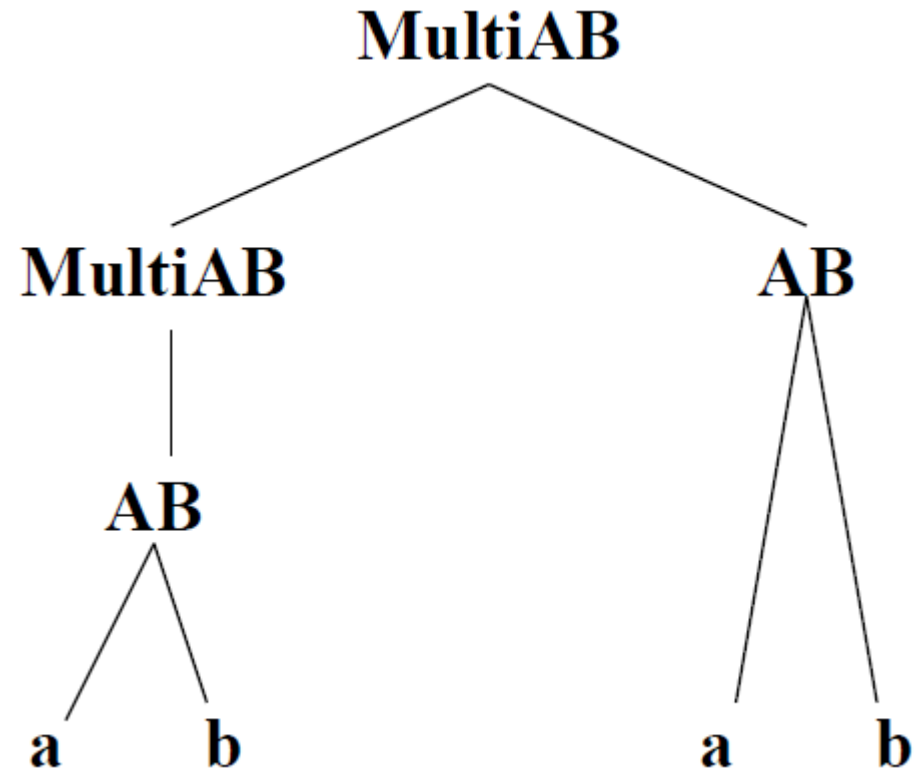
Answer:

$\langle \text{Stmt} \rangle \Rightarrow \langle A \rangle \langle B \rangle$
 $\Rightarrow \mathbf{a} \langle B \rangle \mathbf{b}$
 $\Rightarrow \mathbf{a} \langle B \rangle \mathbf{b} \mathbf{b}$
 $\Rightarrow \mathbf{abbb}$

Answer: Can not derivate **baa** or
any strings start with b.

Parse Trees

- Each internal node is a **Non-terminal**; its children make up the right-hand side of one of the productions for that **Non-terminal**.



Parse Tree

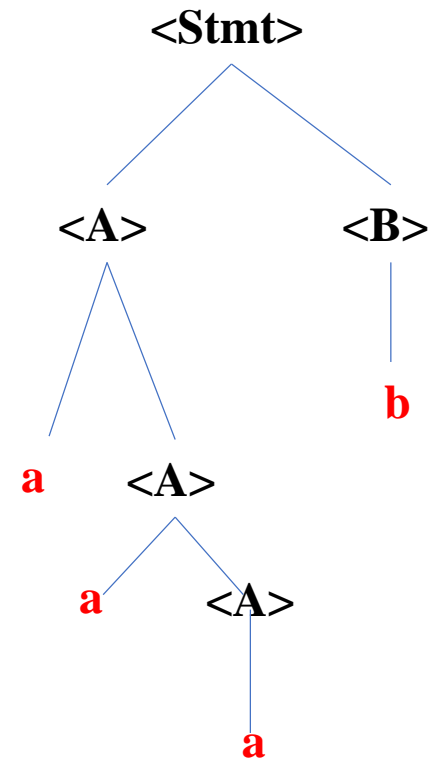
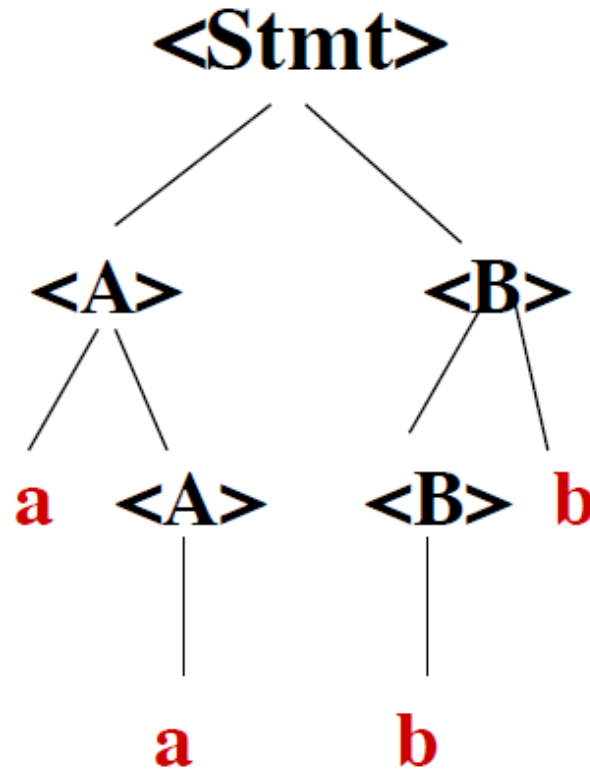
- Example: Draw a parse tree for **aabb** and **aaab** in G2?

G2:

$\langle \text{Stmt} \rangle \Rightarrow \langle \text{A} \rangle \mid \langle \text{A} \rangle \langle \text{B} \rangle$

$\langle \text{A} \rangle \Rightarrow \mathbf{a} \mid \mathbf{a} \langle \text{A} \rangle$

$\langle \text{B} \rangle \Rightarrow \mathbf{b} \mid \langle \text{B} \rangle \mathbf{b}$



Grammars are not Unique

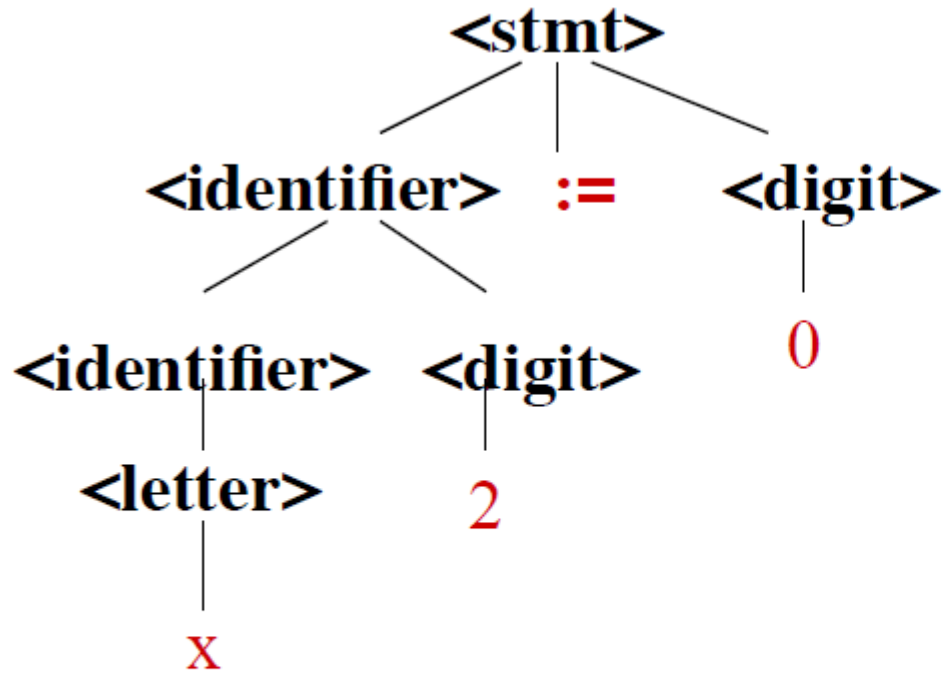
Consider a grammar G:

$\langle \text{stmt} \rangle ::= \langle \text{ident} \rangle := \langle \text{digit} \rangle$
 $\langle \text{ident} \rangle ::= \langle \text{letter} \rangle$
 $\langle \text{ident} \rangle ::= \langle \text{ident} \rangle \langle \text{letter} \rangle$
 $\langle \text{ident} \rangle ::= \langle \text{ident} \rangle \langle \text{digit} \rangle$
 $\langle \text{letter} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$
 $\langle \text{digit} \rangle ::= 0 \mid 1 \mid \dots \mid 8 \mid 9$

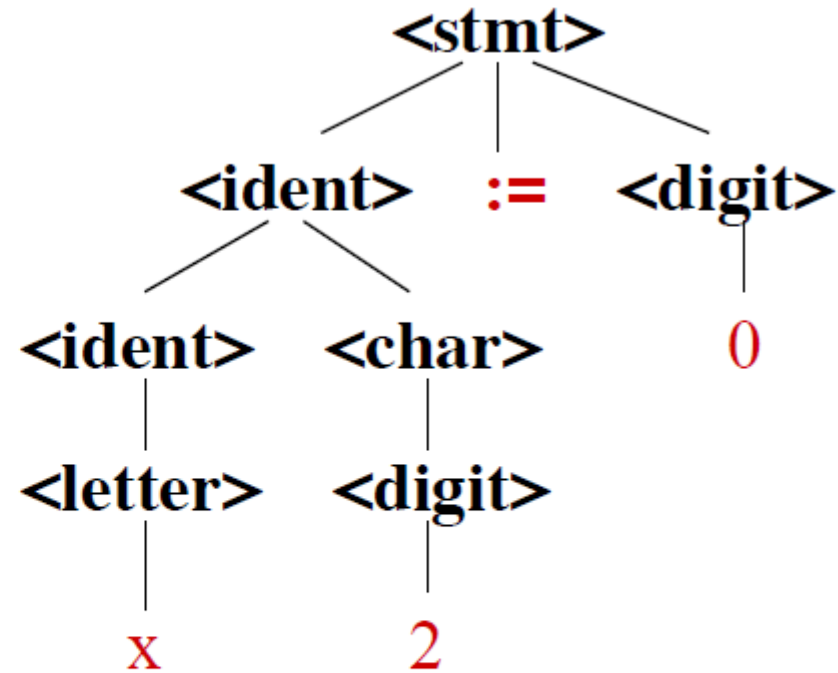
Consider a grammar G':

$\langle \text{stmt} \rangle ::= \langle \text{ident} \rangle := \langle \text{digit} \rangle$
 $\langle \text{ident} \rangle ::= \langle \text{letter} \rangle$
 $\langle \text{ident} \rangle ::= \langle \text{ident} \rangle \langle \text{char} \rangle$
 $\langle \text{char} \rangle ::= \langle \text{letter} \rangle \mid \langle \text{digit} \rangle$
 $\langle \text{letter} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$
 $\langle \text{digit} \rangle ::= 0 \mid 1 \mid \dots \mid 8 \mid 9$

Grammars are not Unique



Parse Tree for G



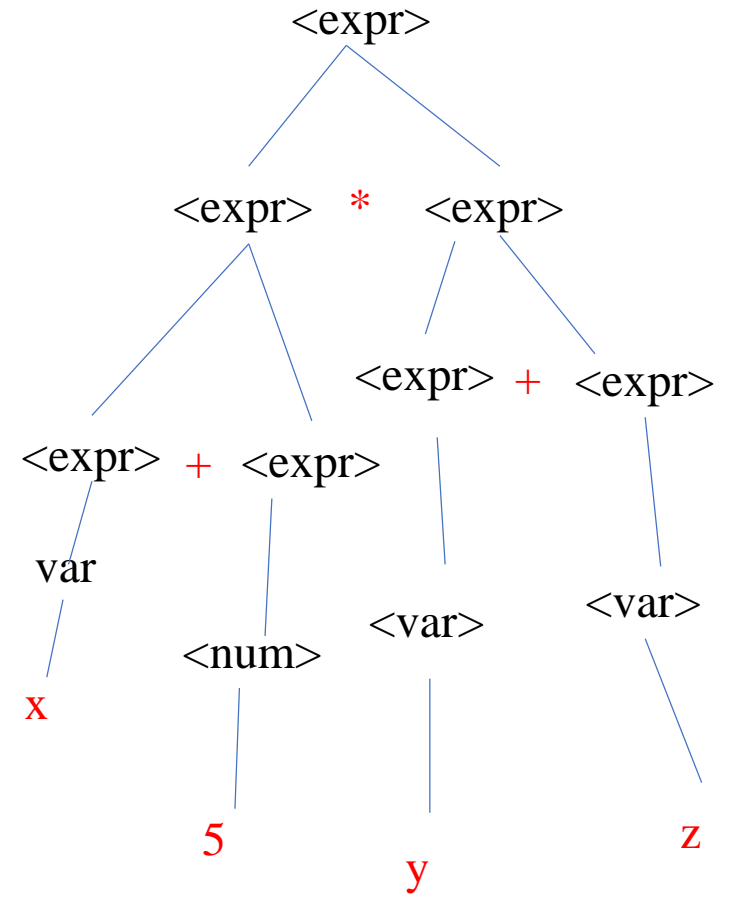
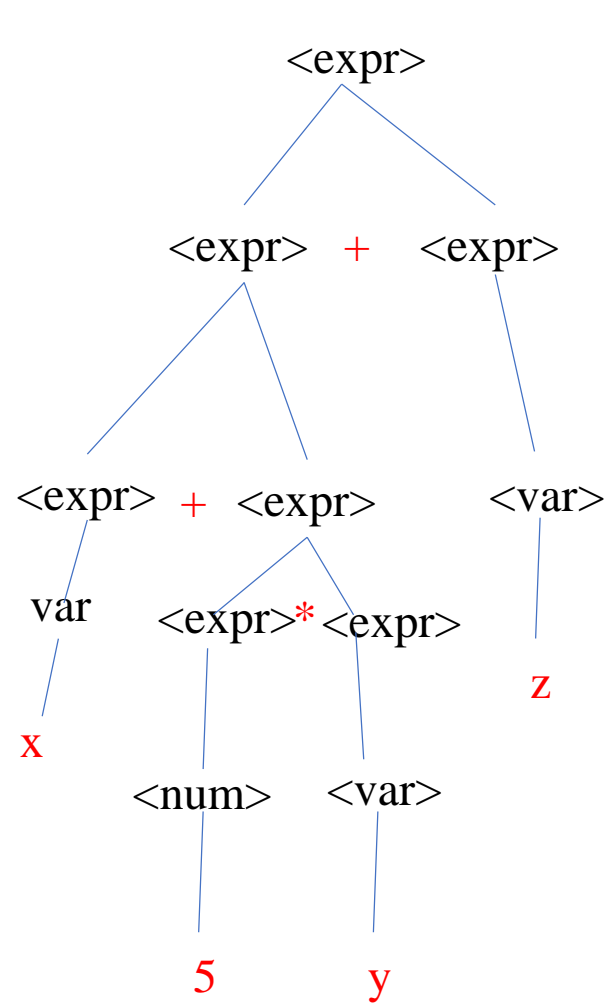
Parse Tree for G'

Arithmetic Expressions

Here is a grammar for arithmetic expressions:

- $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid$
- $\langle \text{expr} \rangle * \langle \text{expr} \rangle \mid \langle \text{expr} \rangle / \langle \text{expr} \rangle \mid \langle \text{var} \rangle \mid \langle \text{num} \rangle$
- $\langle \text{var} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$
- $\langle \text{num} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Example: Using this grammar, how would we parse: $x + 5 * y + z$?



Two different parse trees here

Precedence

Modify the grammar to add precedence:

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{term} \rangle \mid \langle \text{term} \rangle / \langle \text{term} \rangle \mid \langle \text{factor} \rangle$

$\langle \text{factor} \rangle ::= \langle \text{var} \rangle \mid \langle \text{num} \rangle \mid (\langle \text{expr} \rangle)$

$\langle \text{var} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$

$\langle \text{num} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

E.g.: $3+4*5$ means $3+(4*5)$ rather than $(3+4) * 5$

Associativity

Modify the grammar to add associativity:

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{term} \rangle / \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$

$\langle \text{factor} \rangle ::= \langle \text{var} \rangle \mid \langle \text{num} \rangle \mid (\langle \text{expr} \rangle)$

$\langle \text{var} \rangle ::= a \mid b \mid c \mid \dots \mid x \mid y \mid z$

$\langle \text{num} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

E.g.: 10-4-3 means (10-4)-3 rather than 10-(4-3)

Types of Grammars

Context Free Grammars:

- Every rule has a **single nonterminal** on the left-hand side:

$\langle A \rangle \Rightarrow \dots$

- Disallowed: $\langle X \rangle \langle A \rangle \Rightarrow \langle X \rangle a$

Regular Grammars:

- Rules all take the forms:

$\langle A \rangle \Rightarrow c$ or $\langle A \rangle \Rightarrow \langle B \rangle c$ (*left-linear*)

- Or rules all take the forms:

$\langle A \rangle \Rightarrow c$ or $\langle A \rangle \Rightarrow c \langle B \rangle$ (*right-linear*)

- Disallowed: $S \Rightarrow a S b$

- Cannot generate the language $\{ a^n b^n \mid n = 1, 2, 3, \dots \}$

Types of Grammars

- **Context Free Grammars** (CFGs) are used to specify the overall structure of a programming language:
 - if/then/else, ...
 - brackets: (), { }, begin/end, ...
- **Regular Grammars** (RGs) are used to specify the structure of tokens:
 - identifiers, numbers, keywords, ...
- RGs are a subset of CFGs

Ambiguous

- An ambiguous string is one that can be parsed into more than one different parse tree.
- Removing ambiguity:
 - 1: Precedence
 - 2: Associativity

Regular Expressions

	<u>RE Notation</u>	<u>Language</u>
	an empty RE	$\{ \}$
symbol a	a	$\{a\}$
null symbol	ϵ (means empty string)	$\{""\}$
R, S regular exprs	$R \mid S$	$L_R \cup L_S$
	$a \mid b$ (<i>alternation</i>)	$\{a, b\}$
R, S regular exprs	RS	$L_R L_S$
	ab (<i>concatenation</i>)	$\{ab\}$
	$(a \mid b)(c \mid d)$ (<i>concatenation of alternations</i>)	$\{ac, ad, bc, bd\}$

Regular Expressions

RE Notation

a

b

ab

a | b

ab | ac

(a | b)(a | c)

(abc | ε) d

Language

{a}

{b}

{ab}

{a, b}

{ab, ac}

{aa, ac, ba, bc}

{abcd, d}

Regular Expressions

RE Notation

Language

R regular expr R^*

a^*

$\{\epsilon\} \cup L_R \cup L_R L_R \cup L_R L_R L_R \cup \dots$

$\{\epsilon, a, aa, aaa, \dots\}$

R regular expr R^+

a^+

$L_R \cup L_R L_R \cup L_R L_R L_R \cup \dots$

$\{a, aa, aaa, \dots\}$

Note: $\epsilon a = a \epsilon = a$

Precedence is + *, concatenation, |

high to low

(all are left associative operators)

Regular Expressions

RE Notation

a^*

$ab^* = a(b^*)$

$(ab)^*$

$(a \mid b)^*$

a^+

$ab^+ = a(b^+)$

Language

$\{\epsilon, a, aa, aaa, \dots\}$

$\{a, ab, abb, abbb, \dots\}$

$\{\epsilon, ab, abab, ababab, \dots\}$

$\{\epsilon, a, b, aa, ab, ba, bb, \dots\}$

$\{a, aa, aaa, \dots\}$

$\{ab, abb, abbb, \dots\}$

Regular Expressions

- Find regex representing the language

Language	Regex
$\{0\}$	
$\{0,1\}$	
$\{0,01\}$	
$\{0, \epsilon\}\{001\}$	
$\{1\}^*\{10\}$	
$\{10,11,1100\}^*$	

Regular Expressions

Language	Regex
$\{0\}$	0
$\{0,1\}$	0 1
$\{0,01\}$	0 01
$\{0, \epsilon\}\{001\}$	(0 ϵ)001
$\{1\}^*\{10\}$	1*10
$\{10,11,1100\}^*$	(10 11 1100)*

RE's for PLs

- Let *letter* stand for $a|b|c|\dots|a|A|\dots|Z$ and *digit* stand for $0|1|2|3|4|5|6|7|8|9$
 - *letter (letter / digit)** is an identifier
 - *digit +* is an integer constant
 - *digit * . digit +* is real number. So *.5* is, but *5.* is not.

Regular Expressions

Q1: Construct a regular expression for binary numbers of length two

Q2: Construct a regular expression for binary numbers with even length

Q3: Construct a regular expression for floating point numbers that don't use scientific notation (e.g., 3.5, 0.15, -47.3).

- Answer1: $(1/0)(1/0)$
- Answer2: $((1/0)(1/0))^*$
- Answer3: $(-|\epsilon)(0 - 9)^+.(0 - 9)^+$

Key Points

1. Given a grammar, knowing how to prove a string in a $L(G)$ or not.
 - method 1: derivation
 - method 2: parse tree
2. Regular expression

Thanks!!!