

Copyright (C) 2021 Timothy Fossum

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.1 or any later version published by the Free Software Foundation. A copy of the license is included in the file "COPYING", entitled "GNU Free Documentation License".

Specifying Syntax

A *language*, in computer science theory, is a set of strings, where a string is a finite sequence of *symbols* chosen from a given *alphabet*. Computer scientists are interested in part with how to specify languages using things such as Nondeterministic Finite Automata (NFAs), Context-Free Grammars (CFGs), and Turing Machines.

A *programming language* also defines a language in the theory sense. The strings in a programming language are called *programs*, and the symbols are called *tokens* (which we discussed in Slide Set 0). The *syntax* of a programming language is a set of rules used to specify the programs in the language. Most programming languages use a context-free grammar (or something close to it) to specify what sequences of tokens belong to the language.

A programming language also defines the run-time behavior of a program, which is its *semantics*, which we discuss at length later.

Our first step is to describe a formal way in which we can define a programming language. We start out with two simple examples of how to describe familiar data structures. (These are not “programs” in the theory sense of the word, but they will at least get us started with how to specify syntax.)

Backus-Naur Form (BNF)

BNF is a meta-language used to specify a context-free grammar. Every modern programming language uses some sort of BNF notation to define its syntax. We work first with examples of languages that define two simple constructs: lists and trees. Remember that the “alphabet” of a programming language is a set of tokens, so any definition of the syntax of a language must first specify the tokens. In Slide Set 0, we showed how to specify tokens using PLCC.

Our first example is to define a language whose “programs” are lists. Some sample “programs” in this language are:

```
(3 4 5)
(    7 11 )
()
```

Backus-Naur Form (BNF) (continued)

Here is a BNF definition of this language, using three BNF formulas:

```
<lon>    ::= LPAREN <nums> RPAREN
<nums>    ::= NUM <nums>
<nums>    ::=
```

In these formulas, the token names LPAREN and RPAREN stand for left parenthesis ‘(’ and right parenthesis ‘)’, respectively. The token name NUM stands for an (unsigned) decimal number. [Think about what regular expression you would use to match these.] Conforming to PLCC specifications, we use all-UPPERCASE letters for our token names.

BNF: (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

Every BNF formula has the form

$$LHS ::= RHS$$

The *LHS* (Left-Hand Side) of a BNF formula always has the form `<nonterm-symbol>` where `nonterm-symbol` is an identifier (letters, digits, and underscores). PLCC requires that the first character of `<nonterm-symbol>` be a lowercase letter. A `<nonterm-symbol>` expression is called a *nonterminal*. In the example above, the nonterminals are `<lon>` and `<nums>`.

The *RHS* (Right-Hand Side) of a BNF formula is a (possibly empty) sequence of terminal token names and nonterminals.

Notes: The term *syntactic category* is sometimes used instead of *nonterminal*, and the term *terminal* is sometimes used instead of *token name*. When using a token name such as `LPAREN`, some languages use BNF for the corresponding actual character string such as `'('`. In the example above, you can see that we have also skipped whitespace.

BNF (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

BNF has some shortcuts **that we do not use** but that you may encounter in other languages. These shortcuts are usually called Extended BNF, or simply EBNF. For example, instead of writing two different formulas with `<nums>` or `<lon>`, you can use *alternation* notation “`|`”:

$$\langle \text{nums} \rangle ::= \text{NUM} \langle \text{nums} \rangle \mid \epsilon$$

Note: ‘ ϵ ’ means the empty string.

One could also use the *Kleene star* notation to define `<nums>`:

$$\langle \text{nums} \rangle ::= \{ \text{NUM} \}^*$$

Note: We use a variant of the Kleene star notation later.

Parsing (syntactic derivation):

```
<lon> ::= LPAREN <nums> RPAREN  
<nums> ::= NUM <nums>  
<nums> ::=
```

A BNF grammar defines the set of all “legal” token sequences that grammar rules. This set is the *language* of the grammar, where “language” in the sense of computational theory. In the context of languages, these legal token sequences are called “programs”. We also use *syntax rules* to refer to the rules given by a BNF grammar, and we use *syntactically correct* to refer to token sequences that conform to the grammar. Finally, when there is little chance for confusion, we use the term *string* (in the sense of computational theory) to refer to a finite length sequence of tokens.

If we had a set of grammar rules for the Java programming language, then the set of all sentences that conform to this grammar would be the set of syntactically correct Java programs.

In these notes, we casually use the term *token* to refer both to the name of the abstraction (like ‘NUM’) used in a BNF rule and to the concrete lexeme (like ‘23’) used in a program. In most instances, you should have no difficulty understanding which meaning is intended.

Parsing (syntactic derivation):

```
<lon> ::= LPAREN <nums> RPAREN  
<nums> ::= NUM <nums>  
<nums> ::=
```

A grammar can be used:

- to construct syntactically correct sentences, or
- to check to see if a particular target sentence belongs to the language (i.e., is syntactically correct).

A programmer’s job is to construct syntactically correct programs in a programming language, illustrating the first use of a grammar. This activity is called *programming*.

A compiler’s job is, in part, to check a particular sentence for syntactic correctness, illustrating the second use of a grammar. This activity is called *checking or parsing*.

Given a grammar, we want to generate an algorithm to carry out parsing. The program [written in Java, or C, or whatever] that carries out this algorithm is called, unsurprisingly, a *parser*. In this course, all of our parsers are Java programs.

Parsing (continued):

Given a BNF grammar, we describe here an algorithm (written in English) that parses a sentence in the language using what's called a *leftmost* algorithm. The algorithm returns “success” if the parse is successful, “failure” otherwise.

1. Find the *start symbol* of the grammar. **For all of the grammar rules, the start symbol is the first LHS nonterminal in the rule.** This nonterminal becomes the initial *sentential form*. (A “sentential form” is a sentence-like thing that is an ordered list of nonterminals, token names, and lexemes. A “sentence” only has token names. When all of the nonterminal symbols have been removed and the tokens have been replaced by lexemes using the steps given below, the result is a sentence.) Set the *unmatched* sentence to the target sentence.
2. Repeat Step 3 (see the next page) until the sentential form is a sentence (consists only of lexemes (no nonterminals, no token names)).

Parsing (continued):

3. (a) If the leftmost unmatched term in the sentential form is a **token name**, match it with the leftmost string in the unmatched sentence. [Here, the token name in the sentential form describes exactly the string to be matched. For example, the token name LPAREN matches the string ‘(’ and NUM matches the string ‘42’.] If there is no match, return failure. If there is a match, replace the leftmost token name in the sentential form with its matching string (its lexeme) from the unmatched sentence, and remove the matched string from the unmatched sentence.
- (b) If the leftmost unmatched term in the sentential form is a **nonterminal**, choose a rule from the grammar with this nonterminal as its LHS. Replace the nonterminal with the (possibly empty) RHS of the chosen rule. (The choice of rule to choose depends on finding a rule that is most likely to lead to a successful derivation. For some grammars, there is only one choice: the grammar is said to be *predictive*. All of the grammars we use in this course are predictive.) If no rule can apply, return failure.
4. If the unmatched sentence is empty, return success: the target sentence has been successfully parsed. Otherwise, return failure.

It is possible, for some grammars, that Step 3 loops indefinitely. However, for all of the grammars that we encounter in this course, this step exits in a finite number of iterations.

Parsing (continued):

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

We perform a *leftmost derivation* of the target string `(_ 1 4 _ 6 _)`. Always start with the first nonterminal in the grammar (the *start symbol* – in this case it’s `<lon>`) as the sentential form.

sentential form	unmatched sentence	algorithms
<code><lon></code>	<code>(1 4 6)</code>	1
\Rightarrow <u><code>LPAREN</code></u> <code><nums></code> <code>RPAREN</code>	<code>(1 4 6)</code>	3b
\Rightarrow <u><code>(</code></u> <code><nums></code> <code>RPAREN</code>	<code>1 4 6)</code>	3a
\Rightarrow <code>(</code> <u><code>NUM</code></u> <code><nums></code> <code>RPAREN</code>	<code>1 4 6)</code>	3b
\Rightarrow <code>(</code> <u><code>1 4</code></u> <code><nums></code> <code>RPAREN</code>	<code>6)</code>	3a
\Rightarrow <code>(</code> <code>1 4</code> <u><code>NUM</code></u> <code><nums></code> <code>RPAREN</code>	<code>6)</code>	3b
\Rightarrow <code>(</code> <code>1 4</code> <u><code>6</code></u> <code><nums></code> <code>RPAREN</code>	<code>)</code>	3a
\Rightarrow <code>(</code> <code>1 4</code> <code>6</code> <u><code>RPAREN</code></u>	<code>)</code>	3b (<code><nums></code>)
\Rightarrow <code>(</code> <code>1 4</code> <code>6</code> <u><code>)</code></u>		3a
\Rightarrow <i>success!</i>		4

In the above derivation, the leftmost unmatched token name or nonterminal is shown. The underlined parts are the matched lexemes for the unmatched token name as described in rule 3a, or the chosen substitutions for the LHS nonterminals as described in rule 3b.

A derivation ends successfully when there are no token names or nonterminals in the sentential form and no unmatched lexemes.

Parsing (continued):

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

Next we attempt a parse of ‘`(_ 1 4 _ (_ 6 _)`’ using this grammar. We start with the first nonterminal `<lon>` as the sentential form:

sentential form	unmatched sentence
<code><lon></code>	<code>(1 4 (6)</code>
\Rightarrow <u><code>LPAREN</code></u> <code><nums></code> <code>RPAREN</code>	<code>(1 4 (6)</code>
\Rightarrow <code>(</code> <u><code><nums></code></u> <code>RPAREN</code>	<code>1 4 (6)</code>
\Rightarrow <code>(</code> <u><code>NUM</code></u> <code><nums></code> <code>RPAREN</code>	<code>1 4 (6)</code>
\Rightarrow <code>(</code> <code>1 4</code> <u><code><nums></code></u> <code>RPAREN</code>	<code>(6) – try first <num></code>
\Rightarrow <code>(</code> <code>1 4</code> <u><code>NUM</code></u> <code><nums></code> <code>RPAREN</code>	?? NUM doesn’t match
\Rightarrow <code>(</code> <code>1 4</code> <u><code><nums></code></u> <code>RPAREN</code>	<code>(6) – try second <num></code>
\Rightarrow <code>(</code> <code>1 4</code> <u><code>RPAREN</code></u>	?? RPAREN doesn’t match
\Rightarrow <i>failure!</i>	

We can conclude that the string ‘`(_ 1 4 _ (_ 6 _)`’ does not conform to the grammar specifications, so it is *not* a “list of numbers”.

BNF (continued)

A *tree* is another data type that can be specified in BNF form:

```
<tree> ::= NUM
<tree> ::= LPAREN SYMBOL <tree> <tree> RPAREN
```

Here, SYMBOL is a token name that represents a string of alphanumeric characters starting with a letter. [Think of how to specify this using a regular expression.]
NUM, LPAREN and RPAREN token names are as before.

Here are some examples of trees that you can check for syntactic correctness using the parsing algorithm given above.

```
3
( bar 1 ( foo 1 2 ) )
( bar ( biz 3 4 ) ( foo 1 2 ) )
```

BNF (continued)

```
<tree> ::= NUM
<tree> ::= LPAREN SYMBOL <tree> <tree> RPAREN
```

In a given sentential form, the nonterminal `<tree>` can be replaced by the right-hand side of *either* rule having `<tree>` as its left-hand side. For example, the string `'(bar 1 (foo 1 2))'` gives this:

```
<tree>
⇒ LPAREN SYMBOL <tree> <tree> RPAREN
⇒ ( bar <tree> <tree> RPAREN      [*see below...]
⇒ ( bar NUM <tree> RPAREN
⇒ ( bar 1 <tree> RPAREN
⇒ ( bar 1 LPAREN SYMBOL <tree> <tree> RPAREN
⇒ ...
```

Note: this line collapses two **match steps into one:*

```
LPAREN ⇒ ' ('
SYMBOL ⇒ bar
```

Parsers

Recall our grammar for a list of numbers (directory **LON**):

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

We have seen how to take a sentence in this language and use a left-to-right algorithm to parse it. What can we do to automate this algorithm?

A *parser* is a program that takes as input a sentence in a language and outputs a parse of the sentence, producing either success or failure. But deriving a parser from the BNF specification of a language is conceptually simple: it is to write a program to carry out the steps of the parsing algorithm. The difficulty is the algorithm's essential *nondeterminism*: given a nonterminal and a grammar rule having this nonterminal as its left-hand side, how do we replace this nonterminal with the right-hand side of the rule? (For some grammars there may even be more than one “correct” rule that leads to a successful parse.) Such grammars are said to be *ambiguous*.)

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

There are several industrial-strength tool sets in common use that convert a grammar-like grammar specification into a parser. Many of these tool sets are designed for C as the target language: input to such a tool set is a grammar specification, and the output is a set of target language programs that, when compiled and executed, act as a parser. Learning how to use some of these tool sets can be a daunting task, and the generated parsers can be obscure. Furthermore, most of these tool sets are divided into two separate parts: a lexical analysis (scanning) part and a parsing (parsing) part.

Parsers (continued)

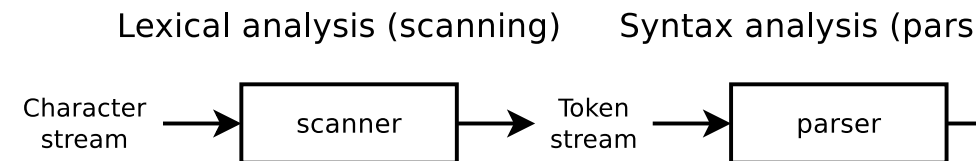
The PLCC tool set is designed to make it easy to build a scanner and a parser for a large collection of programming languages. While it is not “industrial strength” like many of the standard tool sets, it is easy to learn. The target language is Java, so all of the source code that PLCC generates consists of Java programs. To use PLCC, you only need to be comfortable reading and writing Java programs.

To say that a parser returns only success or failure is cold comfort if you want to “run” a program once it parses successfully. So a parser does one of two things: it “runs” the program *as it carries out the parse* (which is called *interpretation on the fly*), or it produces some form of output that can later be used to run the program *once the parse is complete*. The latter approach, since generating a parser is simpler if it is divorced from the task of carrying out run-time behavior during the parse.

The input to a PLCC parser is a token stream – specifically, the token stream generated by an instance of the `Scanner` class (see Slide Set 0). The output of a PLCC parser is a *parse tree* of a program: more specifically, it is a Java object that represents the parse tree.

Parsers (continued)

The following diagram shows how the output of a scanner (a program that implements lexical analysis) becomes input to a parser (a program that implements syntax analysis). The output of a parser is a *parse tree* – which we describe later.



For a particular language whose PLCC grammar file specifies the tokens and BNF grammar, the `plccmk` program generates and compiles source files that implement a scanner and parser for the language specified in the grammar file.

We start with the list-of-numbers (LON) example to show how this works.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

From each BNF grammar rule for a language, PLCC generates a Java class rule such as

```
<lon> ::= LPAREN <nums> RPAREN
```

PLCC generates the Java class `Lon`. PLCC gets the name `Lon` from the LHS nonterminal `<lon>` of this grammar rule by converting the first letter of this nonterminal to uppercase.

PLCC determines the *fields* of a Java class generated from a BNF grammar rule from the RHS terms of the rule – specifically, the terms that are enclosed in angle brackets `<...>`. The RHS of the `lon` grammar rule given above consists of two terms in angle brackets: `<nums>`. The corresponding field name in the `Lon` class is `nums` (Java field names always begin with a lowercase letter). The `nums` field is `Nums`, because `Nums` is the name of the Java class corresponding to the nonterminal `<nums>`.

Since `<nums>` appears as the LHS nonterminal on the second and third grammar rules for this language, PLCC generates a corresponding `Nums` class

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums> ::= NUM <nums>
<nums> ::=
```

However, there are two grammar rules with `<nums>` as the LHS nonterminal. PLCC generates a Java class name `Nums` automatically (by converting the first character of the LHS nonterminal to uppercase), but **PLCC must generate a *unique* Java class name for each grammar rule**. PLCC accomplishes this by annotating the LHS nonterminal on each of these lines with a unique identifier that is different from `Nums` and that distinguishes one from the other. We modify the grammar rules with annotations as follows:

```
<nums>:NumsNode ::= NUM <nums>
<nums>:NumsNull ::=
```

(Any Java class names are OK for these annotations, but good naming conventions suggest using `Node` and `Null` and the names must be unique among the Java class names that PLCC generates). The RHS entries of these rules are unchanged.

With these modifications, PLCC generates two new classes, `NumsNode` and `NumsNull`. Both of these classes are declared to extend the `Nums` class, so an instance of a `NumsNode` or `NumsNull` example, is also automatically an instance of its parent `Nums` class.

The `NumsNode` class has a field named `nums`, which is an instance of the `Nums` class (the first letter of `<nums>`). The `NumsNull` has no fields.

Parsers (continued)

Here is a complete text file representation of **Language LON**, including syntax specification sections, suitable for processing by `plccr`. We have included the annotations for the two `<nums>` rules as described in the previous slide. Also notice that a line containing a single `'%'` separates the lexical specification from the syntax specification (BNF rules).

```
# Language specification for a list of numbers
# Lexical spec
skip WHITESPACE '\s+'
NUM '\d+'
LPAREN '\('
RPAREN '\)'
%
# Syntax spec (BNF rules)
<lon> ::= LPAREN <nums> RPAREN
<nums>:NumsNode ::= NUM <nums>
<nums>:NumsNull ::=
%
```

Parsers (continued)

```
# Language specification for a list of numbers
# Lexical spec
skip WHITESPACE '\s+'
NUM '\d+'
LPAREN '\('
RPAREN '\)'
%
# Syntax spec (BNF rules)
<lon> ::= LPAREN <nums> RPAREN
<nums>:NumsNode ::= NUM <nums>
<nums>:NumsNull ::=
%
```

To summarize, the top lines of the file, up to the first percent (`%`), constitute the *lexical specification* of the language. For this language, these lines specify that whitespace (including spaces, tabs, and newlines) should be skipped, that `NUM` is a string of one or more decimal digits, and that `LPAREN` and `RPAREN` are the characters `'('` and `')'`, respectively.

The remaining lines of the file constitute the *syntax specification* of the language, given in BNF form. `PLCC` takes this file as input and generates a collection of source files in a Java subdirectory that implement a scanner and a parser for the language.

Parsers (continued)

```
# Language specification for a list of numbers
# Lexical spec
skip WHITESPACE '\s+'
NUM '\d+'
LPAREN '\('
RPAREN '\)'
%
# Syntax spec (BNF rules)
<lon>          ::= LPAREN <nums> RPAREN
<nums>:NumsNode ::= NUM <nums>
<nums>:NumsNull ::=
%
```

Observe that the LPAREN and RPAREN token specifications use `\\` for backslashes. This is because parentheses have a special meaning in regular expressions, and the backslash turns off this special meaning. For example, the regular expression `'\\('` really matches the left parenthesis. (You should take this opportunity to review how to write and interpret regular expressions in the `JavaPattern` class.)

To complete this example, assume that you have created a directory `LON` and that in this directory you have created a file named `grammar.l` with the lines appearing above. In your `LON` directory, run the `plccmk` script

Parsers (continued)

`plccmk`

Your script output should appear as follows:

Nonterminals (* indicates start symbol):

*<lon>
<nums>

Abstract classes:

Nums

Source files created:

NumsNode.java

...

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN  
<nums>:NumsNode ::= NUM <nums>  
<nums>:NumsNull ::=
```

(The above text shows just the syntax section of the `grammar` file examining.) In your `LON` directory, change to the subdirectory named `plccmk` command has created this subdirectory and populated it with code generated by `plcc.py`. In this subdirectory you will find (things) the following Java source files:

```
Lon.java  
Nums.java  
NumsNode.java  
NumsNull.java
```

Each of these corresponds to one or more of the grammar rule lines. The line beginning with `<lon>` results in the file `Lon.java` being created in the `Lon` Java subdirectory. As you can see from looking at the Java code in the `Lon` directory for the `NumsNode` and `NumsNull` classes, both of these correspond to the `Nums` abstract class. This is because the `<nums>` nonterminal is the LHS of two grammar rules.

Parsers (continued)

For every BNF grammar rule, PLCC creates a Java class uniquely named after the rule. For BNF rules that have the same nonterminal appearing on the left-hand side of multiple rules, PLCC creates an abstract class based on the nonterminal. The annotated class names become derived classes of this abstract class.

The first BNF nonterminal in a language's syntax specification section is the *start symbol* for the language. Given a program in the language, the parser creates an instance of the class of the start symbol, which is the root of the parse tree for the program.

Consider Language `LON`, whose grammar file syntax specification is shown below. The start symbol is `<lon>`, and the corresponding Java class is `Lon`. Given a program in Language `LON`, the output of the parser is a *parse tree* rooted at the `Lon` class. We will see shortly how to interpret this instance as a *parse tree*.

```
<lon> ::= LPAREN <nums> RPAREN  
<nums>:NumsNode ::= NUM <nums>  
<nums>:NumsNull ::=
```

Parsers (continued)

From a particular language defined by a language specification file, PLCC generates a parser with class name `Parse` that uses the Java classes created from processing the syntax specification sections. To run this parser in language directory `LON`, for the `Parse` class file using `-cp Java` on the `java` command line so the Java interpreter can find the `.class` files. This program displays a prompt `-->_`, after which you enter programs from standard input, “programs” in Language `LON` to parse. (Standard input typically comes from the keyboard, sometimes called the *console*.) If the parse for a particular program succeeds, the program displays the string `"OK"`. If the parse fails, it displays an error message indicating where the parse failed.

For example, in the `LON` directory (after having run `plccmk`), enter the following:

```
java -cp Java Parse
```

With a standard input of `(14 6)`, your input and output looks like this:

```
--> ( 14 6 )
OK
```

With a standard input of `(14 (6)`, your input and output looks like this:

```
--> ( 14 ( 6 )
%%% Parse error: Nums cannot begin with LPAREN
```

You can also use the `parse` shell script, which is equivalent to `java -cp Java Parse`.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums>:NumsNode ::= NUM <nums>
<nums>:NumsNull ::=
```

When invoked with command-line arguments, the `Parse` program processes the arguments given as filename arguments and processes them as if they were entered from standard input.

Three command-line arguments have special meaning when running the `Parse` program:

- The `-n` command-line argument disables displaying the `-->_` prompt, reading programs from standard input.
- The `-t` command-line argument toggles displaying a *parse trace*. When `-t` is present, the programs are parsed. A parse trace is a text representation of the parse tree for the programs: an example of such a parse trace appears on Slide 1.3. When `-t` is not present, the parser does not display the parse trace.
- The `-v` command-line argument toggles displaying the command-line files when processing them in left-to-right order. When `-v` is present, the program does not display the name.

The same command-line arguments pertain to the `Rep` program which begins with the next slide.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN  
<nums>:NumsNode ::= NUM <nums>  
<nums>:NumsNull ::=
```

PLCC also generates an interactive parser/evaluator called `Rep` in the `Java` subdirectory along with the `Parse` program. `Rep` executes an “Eval-Print” loop that displays a prompt, *Reads* program input from the user, and parses it – producing a parse tree, *Evaluates* the parse tree (computes its value), and *Prints* the result obtained from “running” the program. This result is a representation of the “value” of the program’s parse tree, computed by the Java class defined by the start symbol of the language’s BNF.

How to determine the “value” of a parse tree is the essence of *semantics*, which we describe in detail for each of the languages we consider in this course. In some cases, this value may be just a re-cast version of the program text; in other cases, it may be the numeric value of a function application.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN  
<nums>:NumsNode ::= NUM <nums>  
<nums>:NumsNull ::=
```

By default, the value of a parse tree is simply the `toString()` value of the `tree` Java instance. For Language LON, this value is a Java `String` object, `Lon@xxxx`, where the `xxxx` part is the hexadecimal address of where the object resides in memory. Observe that `Lon` is the Java class that defines the BNF rules of the LON BNF rules.

The ‘-t’ command-line argument turns on a parse trace when used with `Rep` or `Parse` programs. The next slide shows the output using this feature.

Parsers (continued)

Here's a sample interaction using this trace feature for Language output edited for the sake of readability:

```
$ java -cp Java Rep -t
--> (14 6)
<lon>
| LPAREN "("
| <nums>:NumsNode
| | NUM "42"
| | <nums>:NumsNode
| | | NUM "6"
| | | <nums>:NumsNull
| RPAREN ")" "
Lon@372f7a8d
```

In this example, the root of the parse tree is a `Lon` object whose `nums` field is an instance of `NumsNode`. This instance in turn has a `nums` field that is an instance of `NumsNode`. And finally, this instance has a `nums` field that is an instance of `NumsNull`. A `NumsNull` object has no fields. The trace above shows each token is matched by the parser, along with the lexeme from the input that matched the token. The line '`Lon@372f7a8d`' shows the result of the parse, which is a `Lon` object representing the root of the parse tree. The '@' part represents the location in memory where the `Lon` object resides.

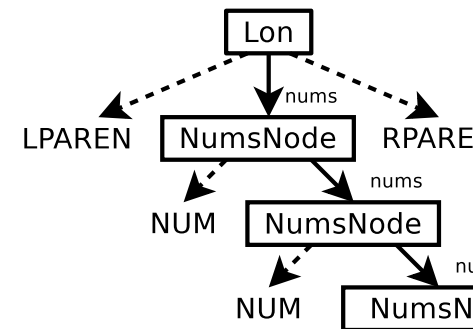
Parsers (continued)

Using this trace feature, we can construct a visual representation of the parse tree whose root is an object of type `Lon`. The diagrams below show the raw output on the left and the resulting parse tree on the right. Note that the parse tree is represented by boxes. The dashed lines point to tokens that are not part of the conceptual parse tree but that do *not* correspond to fields in the non-terminal (this is the case when the BNF entries on its RHS are tokens that are not non-terminals, like `LPAREN`.)

Rep -t
output

```
<lon>
| LPAREN "("
| <nums>:NumsNode
| | NUM "42"
| | <nums>:NumsNode
| | | NUM "6"
| | | <nums>:NumsNull
| RPAREN ")"
```

Parse tree



Parsers (continued)

To make the connection clear between a grammar rule and its PLCC-generated J class notes we display the grammar rules in the following way:

```
<lon>                ::= _LPAREN_<nums>_RPAREN
                        Lon (Nums_nums)
<nums>:NumsNode      ::= _NUM_<nums>
                        NumsNode (Nums_nums)
<nums>:NumsNull      ::= _
                        NumsNull ()
```

The item in the box following a grammar rule is the *Java signature of the Java* corresponding to the class PLCC generates from the grammar rule. So the box

```
Lon (Nums_nums)
```

means that the constructor for the PLCC-generated class `Lon` has a single param `Nums`.

There is a one-to-one correspondence between the types and formal param structor and the types and field names in the class. More specifically, when invoked, the constructor body simply copies the values of its parameters into the c names of the instance being constructed.

As you can see from above, the `NumsNull` constructor takes no parameters, an class has no corresponding fields.

Parsers (continued)

```
<lon>                ::= LPAREN <nums> RPAREN
<nums>:NumsNode      ::= NUM <nums>
<nums>:NumsNull      ::=
```

The parse tree for a program in this language accurately represents the list of numbers given as input – count the number of instances of `NUM` in the trace given on Slide 1.30 above – but the actual `NUM` tokens are not in the parse tree. The problem is that the parse tree instances of `Nums` do not have fields corresponding to their `NUM` tokens.

Parsers (continued)

To remedy this situation, we want to define a field in the `NumsNode` class that captures the `NUM` token obtained by the parser, which in turn allows us to access the token's lexeme. We already use angle brackets for all *nonterminal* symbols of grammar rules, and these nonterminals automatically become fields of the corresponding class. So to capture *tokens* on the RHS, we use angle brackets for *terminals* as well. This means that the `NumsNode` line now looks like this:

```
<nums>:NumsNode ::= <NUM> <nums>
```

The '`<NUM>`' entry creates a field named `num` (which is the token converted to lowercase) of type `Token` (since `NUM` is the name of the token in the grammar file).

Observe that it's unnecessary to capture tokens that always have the same lexeme, like `LPAREN`: every instance of the `LPAREN` token looks like every other instance, so there's no need to distinguish among them. This is not so with tokens that take on multiple lexeme values, such as `NUM`. Indeed, for a list of numbers, knowing exactly *what* numbers are in the list can be essential – for example, to calculate the sum of the numbers in the list. If these items do not appear in the PLCC-generated classes, their values do not appear in the parse tree, and their values are not retrievable after the parse.

Parsers (continued)

The revised grammar now looks like this:

```
<lon> ::= LPAREN <nums> RPAREN
<nums>:NumsNode ::= <NUM> <nums>
<nums>:NumsNull ::=
```

Since we now have two fields in the `NumsNode` class, we need to update the signature of the `NumsNode` constructor, as shown here (compare with the previous version):

```
<lon> ::= _LPAREN_<nums>_RPAREN
               Lon (Nums_nums)
<nums>:NumsNode ::= _<NUM>_<nums>
               NumsNode (Token_num, _Nums_nums)
<nums>:NumsNull ::= _
               NumsNull ()
```

Parsers (continued)

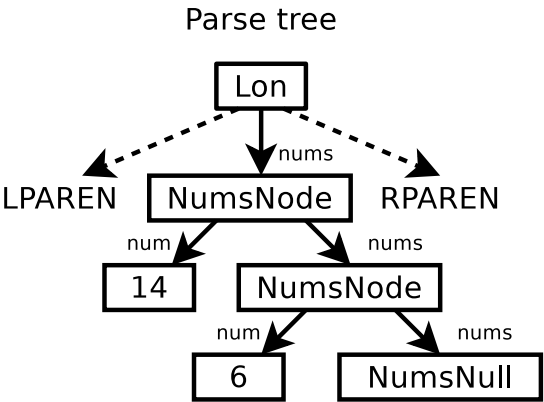
```
<lon> ::= LPAREN <nums> RPAREN
<nums>:NumsNode ::= <NUM> <nums>
<nums>:NumsNull ::=
```

We now have two fields in the NumsNode class: num (of type Token) and nums. When the parser generates a NumsNode object, it fills in its num field with the value object. The toString() method applied to this object retrieves the lexeme of the token which is a string of decimal digits. The rest of the NumsNode object is the same

The following diagram shows the parse tree for the string

(14 6)

using this revised grammar. Notice that each NumsNode object now has a num field containing decimal digits.



Parsers (continued)

Let's return to our tree example (directory TREE). Here is a modified parse tree as originally given on slide 1.12 that takes into account the PLCC for unique class names and the introduction of fields for the NUM tokens. Notice the Java signatures for the class constructors shown in

```
<tree>:Leaf ::= _<NUM>
                Leaf(Token_num)
<tree>:Interior ::= _LPAREN_<SYMBOL>_<tree>_<tree>_RPAREN
                Interior(Token_symbol, _Tree_tree, _Tree_tree)
```

But there is a problem with the Interior constructor signature. It does not allow multiple field names or constructor formal parameters with the same name. Specifically, in this case, there cannot be two field names with the same name. Here is what PLCC has to say about this when given the above grammar (the code has been folded for clarity):

```
duplicate field name tree in rule RHS
LPAREN <SYMBOL> <tree> <tree> RPAREN
```

Parsers (continued)

The solution is to use different identifiers for these field names at responding constructor formal parameters. PLCC allows duplication (in angle brackets) to be annotated – much like we have seen for nonterminal names – with alternate names that avoid this conflict.

In the case we are considering, we can resolve this issue in the `<tree>:Interior` grammar rule by using the identifier `left` for the first `<tree>` field and the identifier `right` for the second `<tree>` field. Here is the updated rule:

```
<tree>:Interior ::= _LPAREN _<SYMBOL> _<tree>left _<tree>right  
Interior(Token_symbol, _Tree_left, _Tree_right)
```

In summary, this grammar rule creates:

- a class `Interior`
- that extends the abstract class `Tree` and
- that has a field `symbol` of type `Token`,
- a field `left` of type `Tree`, and
- a field `right` of type `Tree`.

Parsers (continued)

As a result, the following grammar is acceptable to PLCC:

```
<tree>:Leaf ::= _<NUM>  
Leaf(Token_num)  
<tree>:Interior ::= _LPAREN _<SYMBOL> _<tree>left _<tree>right  
Interior(Token_symbol, _Tree_left, _Tree_right)
```

When processed by PLCC, we get the following interaction using the `loop`:

```
$ java -cp Java Rep  
--> 3  
Leaf@15db9742  
--> (foo 5 8)  
Interior@6d06d69c  
--> (foo (bar 13 23) 8)  
Interior@7852e922  
--> (goo blah 99) % blah is the culprit here  
%% Parse error: Tree cannot begin with SYMBOL  
-->
```

Examine the Java code for the `Interior` class. You will see the following:

```
public Token symbol;  
public Tree left;  
public Tree right;
```

Parsers (continued)

Let's return to the list-of-numbers example.

```
<lon>          ::= LPAREN <nums> RPAREN
<nums>:NumsNode ::= <NUM> <nums>
<nums>:NumsNull ::=
```

The `Nums` abstract class is extended by both `NumsNode` and `NumsNull`.

```
<nums>:NumsNode ::= <NUM> <nums>
<nums>:NumsNull ::=
```

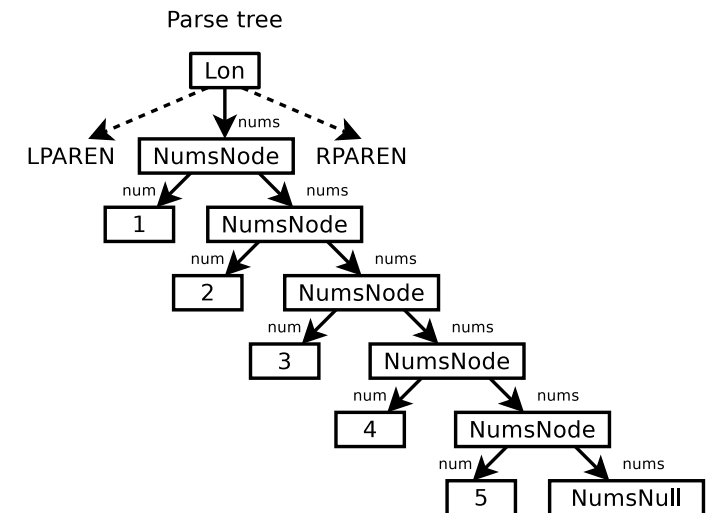
Clearly, when parsing the `<nums>` grammar rules, you get `NumsNode` instances but just one final `NumsNull` instance. RHS of the first `<nums>` rule has `<nums>` as its last entry (which rule is *right recursive*), this rule results in a recursive parsing loop when there are no more `NUM` tokens in the program.

Parsers (continued)

The parse tree for the following list-of-numbers

(1 2 3 4 5)

is shown here:



Because of this right-recursive looping, the nodes in the parse tree double as the number of elements in the list grows. Clearly a list of `num` entries would have 100 `NumsNode` nodes in the tree, weighted significantly right. How can we structure our tree nodes to avoid this?

Parsers (continued)

This sort of looping occurs frequently in programming language specifications. PLCC has a way to encode this. Instead of having two `<nums>` rules, the first being right-recursive and the second having an empty RHS, we can write two rules using a special ‘`**=`’ notation:

```
<nums> ::= _**=_<NUM>
         Nums(List<Token>_numList)
```

The parser accumulates all of the NUM tokens into a single `numList`. The `numList` field name is obtained from the NUM token name by converting characters to lowercase and appending the string `List`.

Note: The use of ‘`**`’ in the notation we have just introduced shows the *Kleene star* repetition notation used in EBNF as well as in regular expressions.

The modified list-of-numbers grammar is in the directory LON2 as follows:

```
<lon> ::= LPAREN <nums> RPAREN
<nums> **= <NUM>
```

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums> **= <NUM>
```

Here is an (edited) example of a parse trace for the list of numbers `(3 5 8 13)`:

```
--> (3 5 8 13)
<lon>
| LPAREN "("
| <nums>
| | NUM "3"
| | NUM "5"
| | NUM "8"
| | NUM "13"
| RPAREN ")"
```

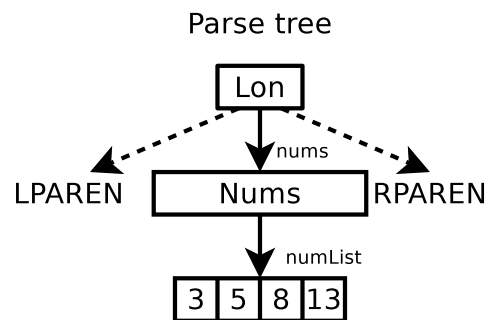
Compare this with the parse trace using the previous grammar that did not use the `**=` construct. The previous parse trace drifts to the right as additional tokens are encountered. Using the `**=` construct, the parse trace becomes much more compact.

PLCC grammar rules that use this construct are called *repeating* grammar rules. Repeating rules are useful in specifying most of our languages.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN  
<nums> **= <NUM>
```

The parse tree for the list of numbers (3_5_8_13) is shown here with compartments to represent a Java List structure:



This parse tree has fewer nodes, and the individual NUM token values are aggregated together into a single numList structure.

Observe that an RHS entry of the form <UPPERCASE> in a BNF corresponds to a leaf node in the parse tree, and an RHS entry <lowercase> corresponds to a subtree of the parse tree. Keep this in mind as you examine the BNF grammar rules for a particular language and the corresponding Java source files.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN  
<nums> **= <NUM>
```

How exactly can we have the Rep program access and display the NUM fields from the parse tree? The Rep program first parses the program to obtain an instance of the start symbol class – an instance of Lon, in this case. Then it evaluates (evaluates) this instance, which as we have seen defaults to displaying the instance like 'Lon@...'.

This default behavior resides in the Java class _Start. This class contains a method named \$run() that simply displays the toString() value of the instance (see the Java code for _Start.java in the Java subdirectory of the LON). Since the Java class Lon extends the _Start class, evaluating an instance of a Lon object defaults to evaluating the \$run() method of the _Start class, which gives us the behavior we have already seen. Rep simply calls the \$run() method on the Lon instance obtained by parsing the program.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN  
<nums> **= <NUM>
```

So we only need to redefine the `$run()` method in the `Lon` class to have the behavior we want: to display the values of the `NUM` fields from the `p` parameter.

Fortunately, PLCC allows us to add methods to PLCC-generated Java files including the added methods in the grammar file. Every time we run PLCC, these added methods are incorporated into the Java source files automatically.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN  
<nums> **= <NUM>
```

For example, if we want to define the `$run()` method in the `Lon` class, we put the following lines into the grammar file *in the semantics section following the syntax section*. The semantics section is separated from the syntax section by a line containing a single ‘%’.

```
Lon  
%%  
    public void $run() {  
        ...  
    }  
%%
```

The `Lon` line tells PLCC that we are adding a method to the `Lon` class. The two lines containing `%%` bracket the Java code to be added. This code can be used to add Java code to *any* PLCC-generated Java file arising from the grammar lines. PLCC inserts this added code at the end of the class, after the Java code generated automatically by PLCC for the class. In the `Lon` class, the above Java code appears at the end of the automatically generated `Lon.java` class code.

Parsers (continued)

```
<lon> ::= LPAREN <nums> RPAREN
<nums> **= <NUM>
```

Recall that we want to define the `$run()` method in the `Lon` class that displays the values of the tokens in the `numList` field. This field, an `List<Token>`, is accessible in the `Lon` class as follows:

```
nums.numList
```

So we can display the `NUM` entries by iterating over this list, displaying (as Java Strings) as we do so. Here is the completed version of the specification that adds the `$run()` method to the `Lon.java` file to achieve this result. Since the parse tree does not include the tokens for parentheses, we add back the parentheses in the output to make the output look “pretty”.

```
Lon
%%%
    public void $run() {
        System.out.print("(");
        for (Token tok: nums.numList)
            System.out.print(tok.toString() + " ");
        System.out.println(")");
    }
%%%
```

Parsers (continued)

Here is the complete grammar file in code directory `LON2` with the following content:

```
# Lexical specification
skip WHITESPACE '\s+'
NUM '\d+'
LPAREN '\('
RPAREN '\)'
%
# Grammar
<lon> ::= LPAREN <nums> RPAREN
<nums> **= <NUM>
%

Lon
%%%
    public void $run() {
        System.out.print("(");
        for (Token tok: nums.numList)
            System.out.print(tok.toString() + " ");
        System.out.println(")");
    }
%%%
```

Parsers (continued)

Continuing with our list-of-numbers example, suppose we want to require that the numbers in our list are separated by commas, like this

(5, 8, 13, 21)

How can we devise a grammar that accommodates these separators?

PLCC provides a way to specify a token that serves as a *separator in a repeating grammar rule*. The separator must be a bare token (no brackets) at the end of the rule, preceded by a '+' character. So to separate the items in our lists by a comma (with token name COMMA) the specification looks like:

```
skip WHITESPACE '\s+'
NUM '\d+'
LPAREN '\('
RPAREN '\)'
COMMA ','
%
<lon> ::= LPAREN <nums> RPAREN
<nums> **= <NUM> +COMMA
%
```

While this grammar uses a COMMA separator between NUM items, it generates *exactly the same parse trees* for this language as for the previous

Parsers (continued)

Slide set 1a provides a summary of PLCC features. You should prepare yourself with this section in preparation for material in subsequent slides of this course.

Static Properties of Variables

A *variable* in a program is a symbol that has an associated value at run-time. One of the main issues in determining the behavior of a program is determining *how* to find the value of a variable at run-time. At any instance in time, the value associated with a variable is called the *value* of the variable.

An *expression* is a language construct that has a value at run-time. A variable, by definition, is an expression, but other language constructs can also have values: for example `x + y` if `x` and `y` are numeric-valued variables.

A programming language that is designed solely for the purpose of evaluating expressions is called an *expression-based language*. Many of the languages we construct in these notes are expression-based, especially early on. Scheme, ML, and Haskell are examples of expression-based languages used in practice. Expression-based languages do their looping principally by using *recursion*.

A programming language whose language constructs are designed to “do something” (e.g., assign the value of an expression to a variable or displaying the value of an expression to the output) is called an *imperative language*. In an imperative language, a language construct that “does something” is called a *statement*; Imperative languages are therefore often called *statement-based languages*. Java, and Python are examples of imperative languages used in practice. Imperative languages do their looping principally using a form of “goto”.

Expression-based languages get their power from defining and applying *function* constructs. A language describing such languages is *functional*.

Static Properties of Variables (continued)

Determining the value of an expression at run-time is at the heart of many issues in programming, particularly so in expression-based languages. Since many expressions involve variables, evaluating an expression requires determining the values of its constituent variables – in other words, finding the values bound to the variables.

At run-time, how can you find the value bound to a variable? There are two main approaches:

- if the location of the value bound to a variable can be determined statically (i.e., if the variable appears *in the text of a program*), we call it *static binding*.
- if the location of the value bound to a variable can only be determined dynamically (i.e., if the variable is accessed *during program execution*), we call it *dynamic binding*.

Almost all programming languages commonly in use today use *static binding* principally because it is easier to reason (or prove things) about programs with static bindings. You will have the opportunity to explore dynamic binding in your homework.

Static Properties of Variables (continued)

We consider only static bindings for now, so the program text tells us where variable bindings occur.

For a given variable, the *scope* of the variable is the region of code in which the variable's binding can be determined. Consider the following Java program:

```
public class Foo {
    public static int y;
    public int z;
    public static void main(String [] args) {
        // args is local to main
        Foo f = new Foo(); // f is local to main
        int x = 1; // x is local in main
        Foo.y = 2; // y is static throughout in Foo
        f.z = 3; // z is known only within instances of Foo
    }
}
```

In the above code, the scope of `y` is *global*, from its declaration as a `public static` variable to the end of the class. The scope of `z` is *global*, known only within (and throughout) instances of the class `Foo`. The scope of `f` and `x` are *local*, from their declarations to the end of the `main` method. The scope of `args` is also local, from the beginning of the method body to the end of the method.

Static Properties of Variables (continued)

It is possible for one symbol to have multiple bindings depending on where it occurs. Consider:

```
public class Bar {
    public static int x;
    public static void main(String [] args) {
        x = 3;
        System.out.println(x);
        { // beginning of block
            int x = 4;
            System.out.println(x);
        } // end of block
        System.out.println(x);
    }
}
```

When this program is run, the output appears as follows:

```
3
4
3
```

This is because the “`int x = 4;`” line defines a new variable `x` bound to the value 4. This new binding is from its point of declaration to the end of the *block* in which it is defined, as shown by the comments. In this case, we say that the definition of `x` in the block *shadows* the global definition. We say that the block definition *punches a hole in the scope* of the global definition.

Static Properties of Variables (continued)

Here is the `Foo` class given on the second previous slide:

```
public class Foo {
    public static int y;
    public int z;
    public static void main(String [] args) {
        Foo f = new Foo(); // f is local to main
        int x = 1; // x is local in main
        y = 2; // y is static throughout in Foo
        f.z = 3; // z is known only within instances of Foo
    }
}
```

Consider just the main procedure in this class:

```
public static void main(String [] args) {
    Foo f = new Foo(); // f is local to main
    int x = 1; // x is local in main
    y = 2; // y is static throughout in Foo
    f.z = 3; // z is known only within instances of Foo
}
```

In this method, the identifiers `f` and `x` are explicitly defined. In these cases, identifiers *occur bound* in the main procedure.

However, the variable `y` is not defined anywhere in the procedure `main`. In this the identifier `y` *occurs free* in the main method, but it *occurs bound* in the enclos

The Lambda Calculus – OPTIONAL SECTION

The following grammar defines a formal language called “the lambda calculus” (which is similar to Turing Machines). The `PROC` token is the string `proc`, and `LPAREN`, `RPAREN`, `LBRACE`, `RBRACE`, and `DOT` tokens are straightforward. The examples below.

```
<exp> ::= <SYMBOL>
<exp> ::= PROC LPAREN <SYMBOL> RPAREN LBRACE <exp> RBRACE
<exp> ::= DOT <exp> LPAREN <exp> RPAREN
```

Consider the sentential form (remember what that means?) in this expression obtained from the second grammar rule, where `s` replaces `<SYMBOL>`:

```
proc(s) { <exp> }
```

The occurrence of the symbol `s` in this expression is called a *variable declaration* that *binds* all occurrences of `s` that appear in `<exp>` unless some other declaration of the same symbol `s` occurs in `<exp>`. We say that `<exp>` is the *scope* of the variable declaration for `s`.

The Lambda Calculus (continued)

Occurs Free, Occurs Bound (informal definitions):

A symbol x *occurs free* in an expression E if x appears somewhere that is not bound by any declaration of x in E . A symbol x *occurs bound* in E if x appears in E in such a way that is bound by a declaration of x in E . If the same symbol to occur both bound and free in different parts of E . (Note that the declaration itself is not considered free or bound.)

```
proc(x) {x}           ; x occurs bound
proc(x) {y}           ; y occurs free
.proc(x) {x} (x)      ; first x is bound, second is free
.proc(x) {x} (y)      ; y occurs free
proc(y) {.proc(x) {x} (y)} ; y occurs bound, x occurs free
proc(x) {.proc(y) {x} (y)} ; y occurs free, x occurs bound
.t (u)               ; t and u occur free
```

The Lambda Calculus (continued)

```
<exp> ::= <SYMBOL>
<exp> ::= PROC LPAREN <SYMBOL> RPAREN LBRACE <exp> RBRACE
<exp> ::= DOT <exp> LPAREN <exp> RPAREN
```

Formal definitions of *occurs free* and *occurs bound*:

For a Lambda Calculus expression E , a symbol x *occurs free* in E if

- **Rule 1:**
 E is a <SYMBOL> and E is the same as x .
 $x \text{ occurs free in } E$
- **Rule 2:**
 E is of the form $\text{proc}(y) \{E'\}$ where y is different from x and x occurs free in E' .
 $x \text{ occurs free in } E$
- **Rule 3:**
 E is of the form $\text{.}E_1(E_2)$ and x occurs free in E_1 or E_2 .
 $x \text{ occurs free in } E$

The Lambda Calculus (continued)

```
<exp> ::= <SYMBOL>
<exp> ::= PROC LPAREN <SYMBOL> RPAREN LBRACE <exp>
<exp> ::= DOT <exp> LPAREN <exp> RPAREN
```

For a Lambda Calculus expression E, a symbol *x* occurs bound in E

- **Rule 1:**

E is of the form `proc (y) {E' }` where *x* occurs bound in E' the same symbol and *y* occurs free in E'

```
proc (y) {proc (x) {x}}; x_is_bound
proc (y) {y}; y_is_bound
```

- **Rule 2:**

E is of the form `.E1 (E2)` and *x* occurs bound in E1 or E2

```
.proc (y) {proc (x) {x}} (y); x_is_bound
.proc (y) {x} (proc (y) {y}); y_is_bound
```

The Lambda Calculus (continued)

Lexical and Grammar specification for the Lambda Calculus:

```
# Lexical specification
skip WHITESPACE '\s+'
LPAREN '\('
RPAREN '\)'
LBRACE '\{'
RBRACE '\}'
DOT '\.'
PROC 'proc'
SYM '\w+'
%

# Grammar
<exp>:Var ::= <SYM>
<exp>:Proc ::= PROC LPAREN <SYM> RPAREN LBRACE <exp>
<exp>:App ::= DOT <exp>rator LPAREN <exp>rand RPAREN
%
```

In the Lambda Calculus, if a symbol is bound by a declaration, determine the precise declaration that binds the variable. The Lambda Calculus is of interest theoretically, but it has no practical value as a programming language.

Static Properties of Variables (continued)

Most imperative programming languages are *block structured* and using, another term for static scope rules. A *block* is a region of code containing one or more variable declarations and continuing to the end of the code block where the declarations are active. In C, C++, and Java, blocks are delimited by curly braces '{ ... }'.

In some languages, blocks may be *nested*, in which case variable bindings in inner blocks may be *shadowed* by bindings in inner blocks. Consider, for example, the following C++ code fragment:

```
{ int x = 3;
  { int x = 5;
    cout << x << endl;
  }
  cout << x << endl;
}
```

This code displays 5 and then 3.

In block structured languages, a variable in an expression is bound to the value of the variable with the same name in the *innermost* block that defines the variable. (This rule does not allow the same variable to be defined both in an outer block and an inner block.)

Static Properties of Variables (continued)

Let's return to our C++ example. The following picture shows the binding of the variable `x` in the C++ program fragment given on the previous slide:

```
{ int x = 3;
  { int x = 5;
    cout << x << endl;
  }
  cout << x << endl;
}
```

To determine the binding of a variable in an expression, cross the variable name outwards (up) until a variable declaration with the same variable name is found.

When defining procedures in block structured languages, the formal parameters are considered to be at the same lexical level as local variables. For example, in the following C++ example, the parameter `x` is at the same lexical level as the local variable `y`:

```
int foo(int x) {
    int y;
    ...
}
```