

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".

## PLCC

PLCC is the name of a tool set that takes a language specification file and generates source files for an interpreter for the language.

The specification file is in three sections:

1. lexical specification
2. syntax
3. semantics

A line containing a single ‘%’ is used to separate the lexical specification section from the syntax section and to separate the syntax section from the semantics section.

A specification file can consist of just the lexical specification, in which case the tool set creates only a token scanner (`Scan`) for the language.

A specification file can consist of just the lexical specification and syntax, in which case the tool set creates only the scanner and parser (`Parse`) for the language.

The format of a language specification file is shown here:

```
# lexical specification
...
%
# syntax
...
%
# semantics
...
```

PLCC – Tokens

Comments may appear in a language specification file beginning with ‘#’ and cover the rest of the line, except inside a token specification regular expression or in Java code. This section.

In the lexical specification section of a language specification file, a line is either a comment or a token specification. These specifications have the following form:

```
{skip|token} NAME _' re'
```

where *NAME* is a string of uppercase letters, decimal digits, and underscores, starting with an uppercase letter, and *re* is a *regular expression* as defined by the Java `Pattern` class.

Here are some regular expression examples:

re	matches
d	the letter d
\d	a single decimal digit
d+	one or more ds
\d+	one or more decimal digits
.	any character
\.	the dot character (.)
.*	zero or more characters
%.*	the character % followed zero or more characters
[a-z]	any lowercase letter in the range a to z
\w	any letter (lowercase or uppercase), digit, or underscore
[A-Za-z0-9_]	the same as \w

PLCC regular expressions do not match anything that crosses a line boundary, so the specification `' .*'` matches everything up to the end of the current line, including the end

PLCC – Tokens

A skip specification is used to identify things in a program that are otherwise not part of the syntax structure of the program. We generally skip whitespace (spaces, tabs, and newlines) and language-defined comments. The format of a comment is language-dependent, but in PLCC, a comment starts with a ‘%’ character and continues to the end of the line. Here are the skip specifications for whitespace and comments:

```
skip WHITESPACE '\s+'
skip COMMENT '%.*'
```

A token specification is used to identify things in a program that are meaningful to the syntax structure of the program. Examples are a language reserved word (such as `if`), a left bracket (such as `[`), a numeric literal (such as `346`), or a variable name (such as `xyz`). Here are specifications that match these tokens.

```
token IF 'if'
token LBRACK '['
token LIT '\d+'
token VAR '[A-Za-z]\w*'
```

For token specifications, the initial ‘token’ can be omitted. Notice that regular expressions for characters such as `[` must be quoted with a backslash ‘\’ if they are to be treated as literal characters.

## PLCC – Tokens

The PLCC tool set uses the skip and token specifications to create `Token.java` and `Scan.java`.

The `Token.java` file defines the `Token.Match` enum class – each skip and token specification name. It also defines the structure object consisting of a `match` field of type `Token.Match`, a `str` field of type `String`, a `lno` field of type `int`, and a `line` field of type `String`. The `str` field consists of the characters in the input stream that match the token – the token’s *lexeme*, the `lno` field contains the line number on which the token appears, and the `line` field contains the input stream starting at the line number where the token appears. The `str` field always contains at least one character.

The `Scan.java` file defines an object that is constructed from an `BufferedReader`). The `cur()` method delivers the current `Token` object (skipping over strings that match skip specifications), and the `adv()` method advances to the next token. Multiple calls to `cur` without any intervening `adv` returns the same token repeatedly.

Upon encountering the end of the input stream, `cur` returns a specification object whose `match` field is `$EOF` and whose `toString()` represents the end of the input. The `cur` method is *lazy*, meaning that the `Scan` class does not read from the input stream until necessary to satisfy an explicit `cur` method call.

## PLCC – Tokens

If there are characters remaining in the input stream, the `cur()` method returns the current token and skip regular expressions one-by-one, in the order in which they appear in the lexical specification. Each regular expression is matched against the remaining unmatched input, up to (and including) the end of the current line. If a match is found, the regular expression, the next specification is tried.

If the regular expression is a skip specification that matches at least one character and if no prior token specification has found a matching token, the matched part is skipped and processing continues on the remaining input, *starting over from the first lexical specification*. If a prior token specification has found a matching token candidate, the skip specification is ignored.

If a token specification match of length at least one occurs, and if it is longer than the match length of the previous token candidate (if any), the token specification is chosen as the current token candidate. If not, the current candidate is retained.

## PLCC – Tokens

If – after iterating through all of the lexical specifications – a token has been identified (first longest match), the token candidate is used to instantiate the `Token` class. The `cur()` method advances the input stream by the number of characters matched, saves the `Token` object for possible future use, and returns the `Token` object. If no token candidate has been identified, the `cur()` method returns a special `$ERROR` `Token` object. The `str` field of the `Token` object is a string of the form `‘!ERROR(. . .)’`, where the `‘. . .’` part is (a representation of) the current input stream character where the match failed.

If the `cur()` method finds that there is a non-null saved `Token` object from a prior call to `cur()`, it simply returns that `Token` object without advancing the input stream. Otherwise, processing occurs as described above.

The `adv()` method replaces the saved `Token` object with `null`, so that the next call to `cur()` will be forced to get the next token from the input stream. If it first detects that the saved token object is already `null`, it calls `cur()` to get the next input token.)

As described earlier, the `cur` method returns a special `$EOF` token when it reaches the end of the `BufferedReader` input.

## PLCC – Syntax

The second section of a PLCC language specification file is the syntax section, consisting of grammar rules in the style of Backus-Naur Form (BNF), as shown in Slide Set 1. Recall that a BNF grammar rule has the following form:

$$\text{LHS} ::= \text{RHS}$$

where LHS (the *Left Hand Side*) is a nonterminal and the RHS (the *Right Hand Side*) is a sequence of nonterminals and terminals. The individual parts of a grammar rule, including the `‘ ::= ’` part, are separated by whitespace.

Nonterminals in PLCC are identifiers enclosed between angle brackets. The identifier must begin with a lowercase letter and can consist of any number of additional letters, digits, or underscores. Identifiers that match Java keywords should be avoided.

Terminals in PLCC must begin with an uppercase letter and can consist of any number of more additional uppercase letters, digits, or underscores. A terminal name is the name of a token in the lexical specification section.

**PLCC associates every grammar rule with a unique Java class name derived from the LHS of the grammar rule by converting nonterminals to uppercase.** Class names that match standard Java class names should be avoided.

### PLCC – Syntax (continued)

If the LHS is a simple nonterminal, the Java class name associated with the BNF rule is the nonterminal name with its first letter converted to uppercase. For example

```
<proc> ::= PROC LPAREN <formals> RPAREN <exp>
```

the Java class name is `Proc`.

If the LHS is a nonterminal annotated by adding a colon and a Java class name, then the class name associated with the BNF rule is the annotated Java class name. In this case, an abstract base class is also created whose Java class name is the nonterminal name with its first letter converted to uppercase (as described in Slide 1.1) with the annotated Java class name as a subclass. In this example

```
<exp>:AppExp ::= DOT <exp> LPAREN <rands> RPAREN
```

the Java class name associated with this BNF rule is `AppExp`, and the base class `Exp`.

### PLCC – Syntax (continued)

PLCC requires that there are no duplicates of Java class names associated with the given grammar rules. In particular, if two grammar rules have the same nonterminal name, then their left-hand sides must have annotations that create distinct Java class names. For example, the grammar rules on Slide 1.3

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

```
<nums> ::=
```

would not be acceptable to PLCC. The following annotations fix this problem.

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

```
<nums>:NumsNull ::=
```

When LHS rules are annotated in this way, the nonterminal class associated with the rule (in this example) becomes an abstract Java class, and the annotated Java class names are used to generate classes that extend the abstract class. In this example, the `NumsNode` and `NumsNull` classes are declared to extend the abstract class `Nums`.

## PLCC – Syntax (continued)

The RHS entries in a grammar rule are used to declare public field (in the abstract) class associated with the grammar rule. Only those RHS entries in angle brackets ‘< . . . >’ correspond to fields in the class.

A field can correspond to an RHS token (e.g. <NUM>) or an RHS nonterminal (<nums>).

If the field corresponds to an RHS *token*, its field *name* defaults to the token name with all of its letters in lowercase. For example, the field name corresponding to the RHS token <NUM> is `num`, and its field *type* is `Token`. If the field corresponds to an RHS *nonterminal*, its field *name* defaults to the nonterminal name (with a change). For example, the field name corresponding to the RHS nonterminal <nums> would be `nums`, and its field *type* is the underlying type of the nonterminal – in this case, `Nums` – obtained by converting the first character of the nonterminal name to uppercase.

## PLCC – Syntax (continued)

A BNF grammar rule may have the same nonterminal name appear more than once on its RHS, as in

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

(see Slide 1.2). PLCC requires that there are no duplicates of the same nonterminal associated with the RHS of a given grammar rule, so the above grammar rule is not acceptable to PLCC. To solve the problem of duplicate field names, we can annotate the duplicate field entries by appending an alternate field name. This identifier will do, but the convention is to have it start with a lowercase letter. In the above example, `left` becomes the name of the corresponding field. The following annotated grammar rule is acceptable:

```
<tree>:Interior ::= LPAREN <SYMBOL> <tree>left <tree>right
```

The same requirement that there be no duplicate field entries applies to rules associated with tokens instead of nonterminals. So for a BNF rule such as

```
<hhmmss> ::= <TWOD> COLON <TWOD> COLON <TWOD> NL
```

we would annotate the three <TWOD> token fields with different names:

```
<hhmmss> ::= <TWOD>hh COLON <TWOD>mm COLON <TWOD>ss NL
```

## PLCC – Syntax (continued)

Repeating grammar rules – ones that have ‘`**=`’ instead of ‘`: :=`’ – have fields similar to non-repeating grammar rules, except that their LHS is a `List` of the appropriate type, and their field names have the suffix `list` appended.

The following grammar rule (see Slide 1.36)

```
<nums> **= <NUM>
```

corresponds to the Java class `Nums`, having a single field `numList` of type `List<Token>`.

Similarly, the following grammar rule (we will encounter this later)

```
<letDecls> **= LET <VAR> EQUALS <exp>
```

corresponds to the Java class `LetDecls` having one field `varList` of type `List<Token>` and another field `expList` of type `List<Exp>`.

**A repeating rule cannot be the first rule in a grammar file. A nonterminal class cannot define a repeating rule. A repeating rule cannot have an empty RHS.**

## PLCC – Parsing

PLCC generates a unique Java class for every grammar rule. Each class has a static `parse` method that is called with a `Scan` object parameter and returns an instance of the class with the class fields (defined by the grammar rule described above) populated with appropriate values. For an RHS field corresponding to a token, the field value – a `Token` – comes directly from the `Scan` object of the input file being parsed. For an RHS field corresponding to a nonterminal, the field value comes from calling the `parse` method on the nonterminal class.

Similar remarks apply to repeating rules, where the `parse` method populates the `List` fields in the class. The members of the `List` fields appear in the same order that their corresponding syntax entities appear in the grammar rule.

An abstract Java class generated by a nonterminal that appears on the LHS of more than one grammar rule also defines a static `parse` method. This method takes the current token (delivered by the `Scan` object) and determines which grammar rule corresponds to that token. It then returns the value obtained by calling the `parse` method on the derived class corresponding to the selected rule. The result is an instance of the derived class, which is also an instance of the given abstract class.

## PLCC – Parsing (continued)

Here's a simple example grammar:

```
<tree>:Leaf      ::= <NUM>
<tree>:Interior ::= LPAREN <SYMBOL> <tree>left <tree>right
```

The abstract Tree class looks similar to this:

```
public abstract class Tree {

    public static Tree parse(Scan scn$) {
        Token t$ = scn$.cur();
        Token.Match match$ = t$.match;
        switch(match$) {
            case NUM:
                return Leaf.parse(scn$);
            case LPAREN:
                return Interior.parse(scn$);
            default:
                throw new PLCCException(
                    "Parse error",
                    "Tree cannot begin with " + t$.errSt
                );
        }
    }
}
```

## PLCC – Parsing (continued)

PLCC generates the Interior class as follows:

```
// <tree>:Interior ::= LPAREN <SYMBOL> <tree>left <tree>right
public class Interior extends Tree {

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

    public Interior (Token symbol, Tree left, Tree right) {
        this.symbol = symbol;
        this.left = left;
        this.right = right;
    }

    public static Interior parse(Scan scn$) {
        scn$.match(Token.Match.LPAREN);
        Token symbol = scn$.match(Token.Match.SYMBOL);
        Tree left = Tree.parse(scn$);
        Tree right = Tree.parse(scn$);
        scn$.match(Token.Match.RPAREN);
        return new Interior(symbol, left, right);
    }
}
```



## PLCC – Parsing (continued)

The parse method for a repeating grammar rule uses a loop. For the `<nums> → * <NUM>`, PLCC generates the following Java class `Nums`:

```
public class Nums {

    public List<Token> numList;

    public Nums (List<Token> numList) {
        this.numList = numList;
    }

    public static Nums parse(Scan scn$) {
        List<Token> numList = new ArrayList<Token>();
        while (true) {
            Token t$ = scn$.cur();
            Token.Match match$ = t$.match;
            switch(match$) {
                case NUM:
                    numList.add(scn$.match(Token.Match.NUM));
                    continue;
                default:
                    return new Nums(numList);
            }
        }
    }
}
```

Similar code is generated by repeating rules with a token separator.

## PLCC – Parsing (continued)

While not shown in the examples above, PLCC adds a `Trace` parameter name to the signature of all parse methods. If this parameter is not null – when the ‘–trace’ option is given to the `Parse` or `Rep` programs – parsing the program will display a “parse trace” to standard error (by default). In Language LON, an example parse trace for the expression `(1_3_5)` looks like this:

```
1: <lön>
1: | LPAREN "("
1: | <nums>:NumsNode
1: | | NUM "1"
1: | | <nums>:NumsNode
1: | | | NUM "3"
1: | | | <nums>:NumsNode
1: | | | | NUM "5"
1: | | | | <nums>:NumsNull
1: | RPAREN ")"
```

Each line in the parse trace displays either a nonterminal or a token. If it displays a nonterminal, it shows that the parser calls the `parse` method for that nonterminal (possibly in the `Parse` or `Rep` programs). If it displays a token, then the scanner also displays the token’s lexeme. The decimal number at the beginning of each line is the line number in the source file where the parse found the token; the number of vertical bars indicates the recursive depth of the parse.

## PLCC – Semantics

For a given “program” in the language defined by its specifications, the `$run()` method in the *start symbol class* – the class determined by the first non-terminal in the language BNF grammar rules – returns an instance of this class which is the root of the parse tree for the program. For example, given the

```
<tree>:Leaf      ::= <NUMBER>
<tree>:Interior ::= LPAREN <SYMBOL> <tree>left <tree>right
```

the `parse` method defined in the `Tree` class returns an instance of `Tree` which is the root of the parse tree. (In what follows, we use the term *root* to refer to the root of the parse tree).

**The runtime semantics of any PLCC program is the behavior of calling the `void_$run()` method on the parse tree.** For any PLCC program, the start symbol class extends a special `_Start` class generated automatically by the PLCC compiler. Because the `_Start` class defines a `$run()` method whose behavior is to display the `toString()` representation of this object to standard output, and because the parse tree is an instance of the `_Start` class, the `$run()` behavior defined on the parse tree defaults to displaying this `toString` representation. Here is the code for `$run()` in the `_Start` class:

```
public void $run() {
    System.out.println(this.toString());
}
```

## PLCC – Semantics (continued)

To get a behavior different from this default behavior, the `$run()` method may be redefined in the start symbol class or any of its subclasses. For example, the `$run()` method may be redefined in the `Tree` class so that it displays a human-readable `toString` representation of the object rather than the standard formative default representation. Here is code for the `$run()` method in the `Tree` class that redefines this behavior:

```
public void $run() {
    System.out.println("Tree: " + this.toString());
}
```

Since the `Tree` class appears more than once on the LHS of the grammar rules, each derived class `Leaf` and `Interior` can define its respective `toString()` method.

For example, the `Leaf` class has a field `num` of type `Token`. The `toString()` method in this class can be easily written as follows:

```
public String toString() {
    return num.toString();
}
```

### PLCC – Semantics (continued)

For an instance of class `Interior`, an appropriate `toString` method is:

```
public String toString() {  
    return "("+symbol.toString()+" "+left+" "+right+  
}
```

This relies on the proper (recursive) `toString` behavior of the `left` and `right` fields, both of which are defined as instances of the `Tree` class.

Recall that PLCC generates a Java class for each of the BNF grammar symbols. In the syntax section of the language specification. We can define *semantics* for these classes by adding entries to the semantics section of the language specification file having the form

```
ClassName  
%%%  
...  
%%%
```

where `ClassName` stands for the name of a PLCC-generated class such as `Tree` in Language `TREE`. PLCC inserts the lines of code bracketed by `%%%` lines verbatim into the `ClassName.java` file. Most often, these lines are one or more Java methods that can be applied to instances of the class.

### PLCC – Semantics (continued)

The entire semantics section of the `TREE` language appears here:

```
Tree  
%%%  
    public void $run() {  
        System.out.println("Tree: " + this.toString());  
    }  
%%%  
  
Leaf # extends the Tree class  
%%%  
    public String toString() {  
        // 'num' is a Token field in the Leaf class  
        return num.toString();  
    }  
%%%  
  
Interior # extends the Tree class  
%%%  
    public String toString() {  
        // 'left' and 'right' are Tree fields in the Interior class  
        return "("+symbol.toString()+" "+left+" "+right+  
    }  
%%%
```

### PLCC – Semantics (continued)

As we observed above, PLCC automatically generates a Java source file of the classes – both abstract and non-abstract – that are derived from the grammar rules in the syntax section. In the semantics section of the specification file, if PLCC encounters an entry of the form

```
ClassName
%%%
...
%%%
```

where `ClassName` stands for a class that is *not* one of the automatically generated classes, then PLCC generates a new file `ClassName.java` containing the code bracketed by the ‘`%%%`’ lines. This makes it possible for PLCC to generate source files that can be used to augment the semantics of the language. We will use this to implement *environments*, which we describe later.

In this situation, there is no automatically generated source file, so the code bracketed by the ‘`%%%`’ lines must be a complete Java source file, not

### PLCC – Semantics (continued)

An `include` feature allows a PLCC language specification file to include the contents of other files, making them part of a single specification. Independently created files can be combined together to form a single language specification. The names of include files must be given in the semantics section of the specification file, and generally appear at the end of the specification. Here is an example from a grammar file:

```
...
include code
include env
include prim
include val
```

In this example, the entire contents of the `code` file will be read and appended to the `grammar` file, then the `env` file, and so forth. The `include` files will normally be representative of their purposes, but they do not otherwise play a role in the generated Java files.

### PLCC – Java predefined/reserved name conflicts

The Java class names generated by the LHS of rules in the grammar PLCC file must not conflict with standard Java class names. This means a grammar rule such as

```
<string> ::= ...
```

PLCC creates a Java class named `String`, which results in a Java compiler error since ‘`String`’ is a reserved class name.

Similarly, the names of fields in the RHS of rules in the grammar must not conflict with Java reserved words or predefined identifiers. This means a grammar rule such as

```
<foo> ::= <IF> <blah>null
```

PLCC creates a field named `if` in the `Foo` class, which results in a Java compiler error since `if` is a reserved word in Java. PLCC itself does not identify these errors, so these errors will only show up during compilation.

### PLCC – plcc/plccmk command line arguments

The `plcc` script runs the `plcc.py` program with input from the command line arguments given on the command line. The `plccmk` script does the same but the filename defaults to `grammar`, and also compiles all of the Java files in the `Java` directory. If `plccmk` is given the ‘`-c`’ command line argument, all files in the `Java` directory are removed before running `plcc.py`.

## PLCC – Flags

PLCC has several internal name/value bindings that control its behavior, the value of `LL1` defaults to `True`; if it is changed to `False`, it will not check the grammar for being LL1. Here are some of the default

name	value	meaning
----	-----	-----
Token	True	(boolean) create a Token.java file
debug	0	(integer) larger values give more verbose output
destdir	'Java'	(string) Java source file destination
pattern	True	(boolean) create a scanner that uses regular expressions
LL1	True	(boolean) check for LL(1)
parser	True	(boolean) create a parser
semantics	True	(boolean) create semantics routines
nowrite	False	(boolean) when True, produce *.no* files

Names bound to integer or string values can be changed using comments:

```
... --debug=3 --destdir=JJJ ...
```

or at the top of the specification file on lines beginning with '!':

```
!debug=3
!destdir=JJJ
```

## PLCC – Flags (continued)

Names bound to a boolean value can be set to `True` on the command line, for example:

```
... --nowrite ...
```

or at the beginning of the language specification file:

```
!nowrite
```

Names bound to a boolean value can be set to `False` on the command line, for example:

```
... --LL1= ...
```

or at the top of the language specification file:

```
!LL1=
```

Bindings given in the language specification file always override bindings given on the command line.

## PLCC – Comments in the lexical specification section

In the lexical specification section, PLCC comments appearing in a specification line are defined as matching the Java regular expression `/*.**/`. Such comments are removed before attempting to process the remainder of the token or skip specification line. This means that a skip or expression containing a substring like `'_#'` will mistakenly be considered the start of a PLCC comment. In this case, use `'[_]#'` instead.

## PLCC – Hooks

PLCC creates Java source files from the BNF grammar rules in a location file. For a grammar rule whose LHS is a nonterminal such as `<formals>`, PLCC creates a Java file named `Formals.java`. The initial content appears as follows:

```
//Formals:top//
//Formals:import//
import java.util.*;

// <formals> **= <VAR> +COMMA // the original BNF rule
public class Formals /*Formals:class*/ {

    public static final String $className = "Formals";
    public static final String $ruleString = "<formals> **= <V";

    public List<Token> varList; // the RHS field

    public Formals(List<Token> varList) { // the Formals constructor
//Formals:init//
        this.varList = varList;
    }

    public static Formals parse(Scan scn$, Trace trace$) {
        ... code to parse the RHS of the <formals> grammar rule
    }

//Formals://
}
```

## PLCC – Hooks (continued)

The following lines in this file are called “hooks”:

```
//Formals:top//  
//Formals:import//  
/*Formals:class*/  
//Formals:init//  
//Formals//
```

In the semantics section of the language specification file, you can place these hooks with arbitrary text as needed.

Use the `:init` hook to add additional lines of Java code at the beginning of the `Formals` class constructor (called when parsing a program). For example, to check the `varList` field for duplicate identifiers, include the following code in your language specification file:

```
Formals:init  
%%  
    Env.checkDuplicates(varList, " in proc formals")  
%%
```

This code will then replace the `//Formals:init//` line appearing at the end of the `Formals` class constructor. (Since the “hook” appears as a comment in the original Java source code, it will not affect the compiled code if left

## PLCC – Hooks (continued)

You can use the `:import` hook to add Java `import` lines to source code that need to import Java packages other than the default `java.util`. You can see this used, for example, in Language ABC:

```
Program:import  
%%  
import abcdatalog.engine.bottomup.sequential.*;  
%%
```

Similar comments apply to the `:top` hook, which can be used to add a package comment or such.

The `:class` hook can be used to declare any interfaces that should be implemented in the class definition.



## PLCC – Hooks (continued)

You use the final hook in an automatically-generated Java source file definitions (and sometimes field declarations) that will be appended to the definition. We have already seen how this is used, for example, in Lon on Slide 1.48, which we repeat here:

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

For temporary testing purposes, you may disable adding the `$run()` definition into the `Lon.java` source file by replacing the `Lon` line with

```
Lon:ignore!
```

This is simpler (and easier to undo) than “commenting out” all of the code (including the ‘`%%%`’ lines) by turning them into PLCC ‘`#`’ comment lines. PLCC reads the lines in the grammar file bracketed by the ‘`%%%`’ lines and otherwise ignores them, making no changes to `Lon.java`.