| | Share | 0 | More | Next Blog» | | Create Blog | Sign In |

# Bart's blog

A blog about software development.

**SUNDAY, MARCH 06, 2011**

## 5. Building an AST

*This is part 5/9 of* **Creating your own programming language with ANTLR.**

## Index

- **Introduction**
- **AST construction in ANTLR**
- **Tree operators and rewrite rules**
- **Imaginary tokens**
- **A proper AST**

## Introduction

In this part we're going to continue with the combined grammar file that already contains all lexer- and parser rules. If you haven't already got the necessary files, you can download the project **here**, which contains all the source code, ANTLR jar, build script etc.

## AST construction in ANTLR

By default, ANTLR's generated parser emits tokens of type **CommonToken** as a 1 dimensional list. To tell our parser it should create **CommonTree** tokens that can have parent- and child tokens, we need to set the option `output=AST` in the `options` section of our grammar. This section must come directly *after* the grammar declaration:

```
view plain  print  ?
01.  grammar TL;
02.
03.  options {
04.    output=AST;
05.  }
06.
07.  @parser::header {
08.    package tl.parser;
09.  }
10.
11.  @lexer::header {
12.    package tl.parser;
13.  }
14.
15.  parse
16.    : block EOF
17.    ;
18.
19.  // ...
```

Now edit the **test.tl** file and paste the following in it:

```
view plain  print  ?
01.  a = 2 + 5;
02.  b = a * 2;
```

Also edit the **Main.java** file to look like this:

```
        view plain  print  ?
01.   package tl;
02.
03.   import tl.parser.*;
04.   import org.antlr.runtime.*;
05.   import org.antlr.runtime.tree.*;
06.   import org.antlr.stringtemplate.*;
07.
08.   public class Main {
09.     public static void main(String[] args) throws Exception {
10.        // create an instance of the lexer
11.        TLLexer lexer = new TLLexer(new ANTLRFileStream("test.tl"));
12.
13.        // wrap a token-stream around the lexer
14.        CommonTokenStream tokens = new CommonTokenStream(lexer);
15.
16.        // create the parser
17.        TLParser parser = new TLParser(tokens);
18.
19.        // invoke the entry point of our parser and generate a DOT image of the tree
20.        CommonTree tree = (CommonTree)parser.parse().getTree();
21.        DOTTreeGenerator gen = new DOTTreeGenerator();
22.        StringTemplate st = gen.toDOT(tree);
23.        System.out.println(st);
24.     }
25.   }
```

The code above will cause a **DOT file** to be printed to the console of the generated AST. It will produce the following output:
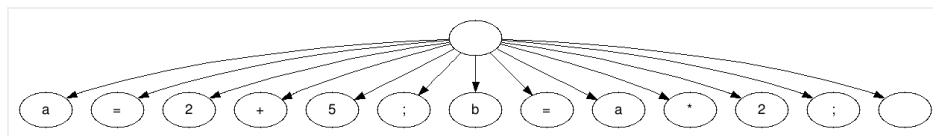
```
digraph {

  n0 [label="A"];
  n1 [label="B"];
  n2 [label="C"];
  n3 [label="2"];
  n4 [label="+"];
  n5 [label="5"];
  n6 [label=";"];
  n7 [label="b"];
  n8 [label="="];
  n9 [label="a"];
  n10 [label="*"];
  n11 [label="2"];
  n12 [label=";"];
  n13 [label=""];

  n0 -> n1 // "" -> "a"
  n0 -> n2 // "" -> "="
  n0 -> n3 // "" -> "2"
  n0 -> n4 // "" -> "+"
  n0 -> n5 // "" -> "5"
  n0 -> n6 // "" -> ";"
  n0 -> n7 // "" -> "b"
  n0 -> n8 // "" -> "="
  n0 -> n9 // "" -> "a"
  n0 -> n10 // "" -> "*"
  n0 -> n11 // "" -> "2"
  n0 -> n12 // "" -> ";"
  n0 -> n13 // "" -> ""

}
```

There are many viewers to generate images from such DOT files. I like **graphviz-dev.appspot.com** as a quick online tool. After you have pasted the output from Main in the **graphviz-dev.appspot.com** tool, you'll see an image like this appear:



This is still just a 1 dimensional list of tokens (but now of type `CommonTree`). We'll have to do a bit more: ANTLR does not magically know what rules need to be the root, and what tokens will need to be omitted from our tree.

---

# Tree operators and rewrite rules

There are two ways to create a proper AST in ANTLR:

1. by using tree operators: `^` and `!`
2. or by using rewrite rules: `... -> ^(...)`

## 1. Tree operators

You can tell ANTLR to make a certain node the root by placing a `^` after it in your parser rules, and to exclude certain nodes from the tree, place a `!` after it.
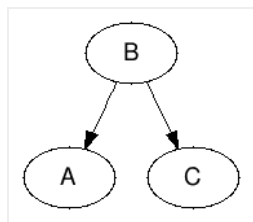
Take the following parser rule for example:

```
view plain  print  ?
01.  bar
02.   : A B C D
03.   ;
```

If we wanted to let `B` become the root, and wanted to exclude token `D` from the tree, we'd do that like this:

```
view plain  print  ?
01.  bar
02.   : A B^ C D!
03.   ;
```

This will produce the tree:



## 2. Rewrite rules

The second option to create an AST, is to use rewrite rules. A rewrite rule is placed after a parser rule (but before the `';'`!). It consists of an arrow, `->`, followed by `^( ... )` where the first token (or rule) inside the parenthesis will become the root of the tree. Tokens that need to be removed from the AST will simply be omitted from the rewrite rule.

To create the same tree as the example above using a rewrite rule, we'd write:

```
view plain  print  ?
01.  bar
02.   : A B C D -> ^(B A C)
03.   ;
```

I.e.: `B` will be the root, and `A` and `C` its children (`D` is removed).

### What to use?

The first option using tree operators may look easier to use. However, when your parser rules become larger, the (small) operators will become less conspicuous. Also, the rewrite rules are more flexible. If we wanted to let token `C` become the first child of parent token `B` in the previous example, tree operators wouldn't be able to do this, but a rewrite rule would simply look like this:

```
view plain  print  ?
01.  bar
02.   : A B C D -> ^(B C A)
03.   ;
```

Another advantage is that the translation of a parser grammar to a tree grammar (the next part of this tutorial) is more straight forward. But you'll see that for yourself later on.

But sometimes it *is* more convenient to use tree operators than rewrite rules. For example, with an add expression:

```
view plain  print  ?
01.  addExpr
02.   : multiplyExpr (('+' | '-') multiplyExpr)*
03.   ;
```
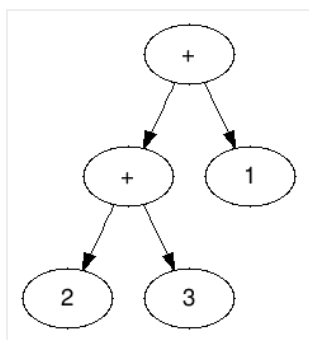
we would like to make the + or − the root of the expression. With rewrite rules, that would look like this:

```
view plain   print   ?
01.   addExpr
02.     :   (multiplyExpr -> multiplyExpr)
03.         ( '+' e=multiplyExpr -> ^('+' $addExpr $e)
04.         | '-' e=multiplyExpr -> ^('-' $addExpr $e)
05.         )*
06.     ;
```

which is no trivial rewrite syntax! And I'll leave it unexplained because in our grammar, we'll go for the far less complex:

```
view plain   print   ?
01.   addExpr
02.     :  multiplyExpr (('+' | '-')^ multiplyExpr)*
03.     ;
```

Yes, that is correct: the single ^ produces the exact tree as the complex looking rewrite rule! For the expression 1+2+3, both produce the tree:



---

# Imaginary tokens

Sometimes a parser rule doesn't have an obvious root candidate. In that case we can insert a so called *imaginary token*. Take the block rule for example:

```
view plain   print   ?
01.   block
02.     :  (statement | functionDecl)* (Return expression ';')?
03.     ;
```

which matches zero or more statements or function declarations, optionally ending with a return statement. To mark the first statement as the root would make little sense. What we're going to do is create a couple of *imaginary tokens* to make a *block tree* always look the same. The first we need to do is put a tokens { ... } section directly below the options { ... } section (and above the header sections!):

```
view plain   print   ?
01.   grammar TL;
02.
03.   options {
04.     output=AST;
05.   }
06.
07.   tokens {
08.     BLOCK;
09.     RETURN;
10.     STATEMENTS;
11.   }
12.
13.   @parser::header {
14.     package tl.parser;
15.   }
16.
17.   // ...
```

Imaginary tokens should start with a capital, just like lexer rules. I find it convenient to let them be distinguishable from lexer rules by using *only* capitals instead of just capitalizing the first letter like I did with lexer rules.

Our `block` rule including a rewrite rule would look like:

```
view plain   print   ?
01.   block
02.     :   (statement | functionDecl)* (Return expression ';')?
03.         -> ^(BLOCK ^(STATEMENTS statement*) ^(RETURN expression?))
04.     ;
```

This makes sure every *block-tree* will have exactly two children: `STATEMENTS` and `RETURN`.

If we didn't use `STATEMENTS`, all zero or more statements would be direct children of the root `BLOCK` node, making it harder to evaluate since we'd always have to inspect if we're at the end, which could possibly be the return statement. Now we know that all zero or more statements are in the left tree, and the optional return statement is in the right tree.

You may have noticed that I removed the `functionDecl` from the tree. This is no accident: at the parsing stage, we're going to create a `java.uti.Map` containing all methods so that when evaluating our language, we can invoke functions before they've been defined in the source. We'll handle that at a later stage of this tutorial.

Note that you can't use a rewrite rule *and* tree operators in a single parser rule: it's always just one of the two.

---

## A proper AST

Okay, having said all that, here's how our grammar looks like with tree operators and rewrite rules:

```
view plain   print   ?
01.   grammar TL;
02.
03.   options {
04.     output=AST;
05.   }
06.
07.   tokens {
08.     BLOCK;
09.     RETURN;
10.     STATEMENTS;
11.     ASSIGNMENT;
12.     FUNC CALL;
13.     EXP;
14.     EXP LIST;
15.     ID LIST;
16.     IF;
17.     TERNARY;
18.     UNARY MIN;
19.     NEGATE;
20.     FUNCTION;
21.     INDEXES;
22.     LIST;
23.     LOOKUP;
24.   }
25.
26.   @parser::header {
27.     package tl.parser;
28.   }
29.
30.   @lexer::header {
31.     package tl.parser;
32.   }
33.
34.   parse
35.     :   block EOF -> block
36.     ;
37.
38.   block
39.     :   (statement | functionDecl)* (Return expression ';')?
40.         -> ^(BLOCK ^(STATEMENTS statement*) ^(RETURN expression?))
41.     ;
42.
43.   statement
44.     :   assignment ';'    -> assignment
45.     |   functionCall ';' -> functionCall
46.     |   ifStatement
47.     |   forStatement
48.     |   whileStatement
49.     ;
50.
51.   assignment
52.     :   Identifier indexes? '=' expression
53.         -> ^(ASSIGNMENT Identifier indexes? expression)
54.     ;
55.
56.   functionCall
57.     :   Identifier '(' exprList? ')' -> ^(FUNC CALL Identifier exprList?)
58.     |   Println '(' expression? ')'  -> ^(FUNC CALL Println expression?)
59.     |   Print '(' expression ')'     -> ^(FUNC CALL Print expression)
60.     |   Assert '(' expression ')'    -> ^(FUNC CALL Assert expression)
61.     |   Size '(' expression ')'      -> ^(FUNC_CALL Size expression)
```

```
 62.      ;
 63.
 64.   ifStatement
 65.      : ifStat elseIfStat* elseStat? End -> ^(IF ifStat elseIfStat* elseStat?)
 66.      ;
 67.
 68.   ifStat
 69.      : If expression Do block -> ^(EXP expression block)
 70.      ;
 71.
 72.   elseIfStat
 73.      : Else If expression Do block -> ^(EXP expression block)
 74.      ;
 75.
 76.   elseStat
 77.      : Else Do block -> ^(EXP block)
 78.      ;
 79.
 80.   functionDecl
 81.      : Def Identifier '(' idList? ')' block End {/* implemented later */}
 82.      ;
 83.
 84.   forStatement
 85.      : For Identifier '=' expression To expression Do block End
 86.        -> ^(For Identifier expression expression block)
 87.      ;
 88.
 89.   whileStatement
 90.      : While expression Do block End -> ^(While expression block)
 91.      ;
 92.
 93.   idList
 94.      : Identifier (',' Identifier)* -> ^(ID_LIST Identifier+)
 95.      ;
 96.
 97.   exprList
 98.      : expression (',' expression)* -> ^(EXP_LIST expression+)
 99.      ;
100.
101.   expression
102.      : condExpr
103.      ;
104.
105.   condExpr
106.      : (orExpr -> orExpr)
107.        ( '?' a=expression ':' b=expression -> ^(TERNARY orExpr $a $b)
108.        | In expression                     -> ^(In orExpr expression)
109.        )?
110.      ;
111.
112.   orExpr
113.      : andExpr ('||'^ andExpr)*
114.      ;
115.
116.   andExpr
117.      : equExpr ('&&'^ equExpr)*
118.      ;
119.
120.   equExpr
121.      : relExpr (('==' | '!=')^ relExpr)*
122.      ;
123.
124.   relExpr
125.      : addExpr (('>=' | '<=' | '>' | '<')^ addExpr)*
126.      ;
127.
128.   addExpr
129.      : mulExpr (('+' | '-')^ mulExpr)*
130.      ;
131.
132.   mulExpr
133.      : powExpr (('*' | '/' | '%')^ powExpr)*
134.      ;
135.
136.   powExpr
137.      : unaryExpr ('^'^ unaryExpr)*
138.      ;
139.
140.   unaryExpr
141.      : '-' atom -> ^(UNARY_MIN atom)
142.      | '!' atom -> ^(NEGATE atom)
143.      | atom
144.      ;
145.
146.   atom
147.      : Number
148.      | Bool
149.      | Null
150.      | lookup
151.      ;
152.
153.   list
154.      : '[' exprList? ']' -> ^(LIST exprList?)
155.      ;
156.
157.   lookup
158.      : functionCall indexes?      -> ^(LOOKUP functionCall indexes?)
159.      | list indexes?              -> ^(LOOKUP list indexes?)
160.      | Identifier indexes?        -> ^(LOOKUP Identifier indexes?)
```

```
161.     |   String indexes?                -> ^(LOOKUP String indexes?)
162.     |   '(' expression ')' indexes? -> ^(LOOKUP expression indexes?)
163.     ;
164.
165.    indexes
166.     :   ('[' expression ']')+ -> ^(INDEXES expression+)
167.     ;
168.
169.    Println  : 'println';
170.    Print    : 'print';
171.    Assert   : 'assert';
172.    Size     : 'size';
173.    Def      : 'def';
174.    If       : 'if';
175.    Else     : 'else';
176.    Return   : 'return';
177.    For      : 'for';
178.    While    : 'while';
179.    To       : 'to';
180.    Do       : 'do';
181.    End      : 'end';
182.    In       : 'in';
183.    Null     : 'null';
184.
185.    Or       : '||';
186.    And      : '&&';
187.    Equals   : '==';
188.    NEquals  : '!=';
189.    GTEquals : '>=';
190.    LTEquals : '<=';
191.    Pow      : '^';
192.    Excl     : '!';
193.    GT       : '>';
194.    LT       : '<';
195.    Add      : '+';
196.    Subtract : '-';
197.    Multiply : '*';
198.    Divide   : '/';
199.    Modulus  : '%';
200.    OBrace   : '{';
201.    CBrace   : '}';
202.    OBracket : '[';
203.    CBracket : ']';
204.    OParen   : '(';
205.    CParen   : ')';
206.    SColon   : ';';
207.    Assign   : '=';
208.    Comma    : ',';
209.    QMark    : '?';
210.    Colon    : ':';
211.
212.    Bool
213.     :   'true'
214.     |   'false'
215.     ;
216.
217.    Number
218.     :   Int ('.' Digit*)?
219.     ;
220.
221.    Identifier
222.     :   ('a'..'z' | 'A'..'Z' | '_') ('a'..'z' | 'A'..'Z' | '_' | Digit)*
223.     ;
224.
225.    String
226.    @after {
227.      setText(getText().substring(1, getText().length()-1).replaceAll("\\\\(.)", "$1"));
228.    }
229.     :   '"'  (~('"' | '\\')  | '\\' ('\\' | '"'))* '"'
230.     |   '\'' (~('\'' | '\\') | '\\' ('\\' | '\''))* '\''
231.     ;
232.
233.    Comment
234.     :   '//' ~('\r' | '\n')* {skip();}
235.     |   '/*' .* '*/'         {skip();}
236.     ;
237.
238.    Space
239.     :   (' ' | '\t' | '\r' | '\n' | '\u000C') {skip();}
240.     ;
241.
242.    fragment Int
243.     :   '1'..'9' Digit*
244.     |   '0'
245.     ;
246.
247.    fragment Digit
248.     :   '0'..'9'
249.     ;
```

And to test this, run the Ant target using the command line parameter `-emacs`, which casues Ant *not* to print all the `[java]` etc. stuff to the console, making it easier to copy and paste the DOT source.

Okay, the tree that was generated is a heck of a lot easier to interpret than the flat list of nodes! But before evaluating, we'll first have to create some sort of iterator that traverses our freshly constructed AST, which will be the subject of the next part **6. Creating a tree grammar**.

*P.S. You can download the project containing all source code created so far, build script, ANTLR jar etc. here.*

Posted by Bart Kiers at **7:58 AM**
Labels: **antlr**, **dsl**, **java**, **parsing**

## 1 comment:

**oleg** said...

Thanks for this wonderful post !!! It really makes things clear.

Keep up the good work :p

**7:28 PM**

**Post a Comment**

**Newer Post**                              **Home**                              **Older Post**

Subscribe to: **Post Comments (Atom)**

Powered by **Blogger**.