

---

IN2009

# **Language Processors**

---

Week 4

# **Parsing with JavaCC**

Igor Siveroni

# Session Plan

---

## Session 4: Parsing with JavaCC

- EBNF
- Parsing with JavaCC
- Lookahead
- JavaCC grammars, semantic actions and values
- Simple expression evaluator

# Extended BNF (EBNF)

---

A few additional operators to shorten definitions:

- $e_1 \mid e_2 \mid e_3 \mid \dots$  : choice of  $e_1$ ,  $e_2$ ,  $e_3$ , etc
- $(\dots)$  bracketing allowed
- $[ \dots ]$  : the expression in  $[ \dots ]$  may be omitted  
 $[M] = (M)?$
- $( e )^+$  : One or more occurrences of  $e$
- $( e )^*$  : Zero or more occurrences of  $e$

# Extended BNF (EBNF)

---

Note that these may be nested within each other, so we can have:

$$(( e_1 \mid e_2 )^* [ e_3 ] ) \mid e_4$$

Examples:

IfStatement  $\rightarrow$  **if** “(“ Expression “)” StatementBlock  
[ **else** StatementBlock ]

StatementBlock  $\rightarrow$  “{“ (Statement)+ “}”

# Expression grammar in EBNF

---

$E \rightarrow E + T$	$T \rightarrow T * F$	$F \rightarrow \text{id}$	<i>Original</i>
$E \rightarrow E - T$	$T \rightarrow T / F$	$F \rightarrow \text{num}$	
$E \rightarrow T$	$T \rightarrow F$	$F \rightarrow "( E )"$	

---

$E \rightarrow T E'$	$T \rightarrow F T'$	$F \rightarrow \text{id}$	<i>Left-recursion eliminated</i>
		$F \rightarrow \text{num}$	
$E' \rightarrow + T E'$	$T' \rightarrow * F T'$	$F \rightarrow "( E )"$	
$E' \rightarrow - T E'$	$T' \rightarrow / F T'$		
$E' \rightarrow$	$T' \rightarrow$		

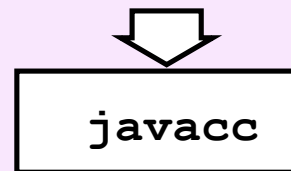
---

$E \rightarrow T ( + T \mid - T )^*$	$T \rightarrow F ( * F \mid / F )^*$	$F \rightarrow \text{id}$	<i>EBNF</i>
		$F \rightarrow \text{num}$	
		$F \rightarrow "( E )"$	

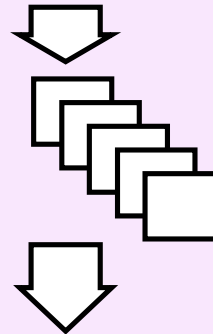
# JavaCC

---

JavaCC specification  
(includes both token specifications *and*  
grammar, possibly with actions)



Parsing and lexical  
analysis Java class files  
(with actions if specified)



combined token matcher and parser (executing actions if specified)

# JavaCC

---

- **JavaCC is a *parser generator*.** Given as input a set of token definitions, a programming language syntax grammar, and a set of actions written in Java, it produces a Java program which will perform lexical analysis to find tokens and then parse the tokens according to the grammar and execute the actions as appropriate.

# JavaCC

---

- it works on LL(1) grammars (no need to understand this definition), which are similar to those that recursive descent works for.
- it requires a non-ambiguous grammar with left-recursion removed, so we use the techniques from earlier this session.

For the record: LL(1) grammar

*Left-to-right* parse, *leftmost* derivation, 1 symbol lookahead



# JavaCC input file format (.jj)

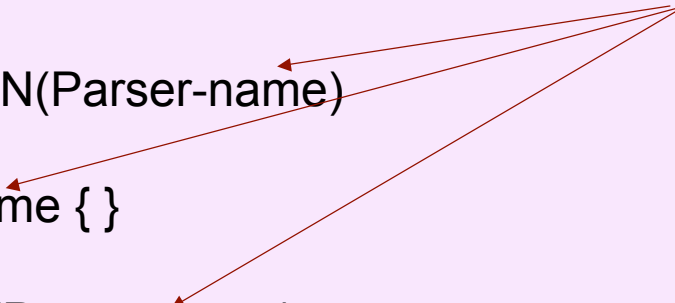
---

PARSER\_BEGIN(Parser-name)

class Parser-name { }

PARSER\_END(Parser-name)

*Parser-name must be the same in all three places*



*/\* Lexical items (i.e. token definitions) – see previous examples \*/*

Token-definitions

*/\* Grammar rules – in a stylised form of EBNF (see next slide). \*/*

Syntax-definitions

# JavaCC Syntax-definitions

---

A BNF production: *Non-Terminal-Name* -> *Right-Hand-Side* is written:

```
java_return_type Non-Terminal-Name ( java_parameter_list ) : (1)
{ java_block }                                           (2)
{ expansion_choices }                                    (3)
```

(1) gives the name of the non-terminal being defined

The rest of (1) looks like a Java method declaration. Using this feature we can cause values to be passed up and down the parse tree while the parse takes place (up via return values and down via parameters).

(2) `java_block`: introduces some Java code which is usually used to declare variables for use in the production

(3) is the EBNF definition and actions...see next slide

# JavaCC EBNF

## expansion\_choices

---

expansion   expansion   ...	where the ` ' separates alternatives.
expansion expansion ...	matches first expansion then second and so on
( expansion_choices )*	matches zero or more expansion_choices
( expansion_choices )+	matches one or more expansion_choices
( expansion_choices )?	matches expansion_choices or empty string
[ expansion_choices ]	ditto (ie same as ?)
regexp	matches the token matched by the regexp
java_id = regexp	ditto, assigning token to java_id
non-terminal-name (...)	matches the non-terminal
java_id = non-terminal-name (...)	ditto, assigning returned value to java_id

The java\_id will usually be declared in the java\_block.

Any of these expansions may be followed by some Java code written in {...} and this code (often called an action) will be **executed** when the generated parser matches the expansion.

# JavaCC EBNF example

---

$$E \rightarrow T ( + T \mid - T )^*$$
$$T \rightarrow F ( * T \mid / T )^*$$
$$F \rightarrow \text{num}$$
$$F \rightarrow ( E )$$

```
void E() :  
{  
  {  
    T() ( "+" T() | "-" T() ) *  
  }  
}
```

```
void T() :  
{  
  {  
    F() ( "*" F() | "/" F() ) *  
  }  
}
```

```
void F() :  
{  
  {  
    <NUM>  
    | "(" E() ")"  
  }  
}
```

```
TOKEN :  
{  
  < NUM: (["0"-"9"])+ >  
}
```

# JavaCC example: Exp.jj file

---

```
PARSER_BEGIN(Exp)

public class Exp {
}

PARSER_END(Exp)

SKIP :
{
    " " | "\t" | "\r"
}

TOKEN :
{
    < NUM: (["0"-"9"])+ > | < EOL: "\n" >
}

void S() :
{
{
    E() <EOL>
    | <EOL>
    | <EOF>
}
}
```

```
void E() :
{
{
    T() ( "+" T() | "-" T() ) *
}
}

void T() :
{
{
    F() ( "*" F() | "/" F() ) *
}
}

void F() :
{
{
    <NUM>
    | "(" E() ")"
}
}
```

# JavaCC example: Main.java file

---

Another class can instantiate the Exp object and call the S() method:

```
public class Main {
    public static void main(String args[]) throws ParseException {
        Exp parser = new Exp(System.in);
        try {
            System.out.println("Type in an expression on a single line.");
            parser.S();
            System.out.println("Expression parser - parse successful");
        } catch (ParseException e) {
            System.out.println("Expression parser - error in parse");
        }
    }
}
```

# JavaCC BNF example

$E \rightarrow T E'$

$T \rightarrow F T'$

$F \rightarrow \text{num}$

$F \rightarrow (E)$

$E' \rightarrow + T E'$

$T' \rightarrow * F T'$

$T' \rightarrow / F T'$

$E' \rightarrow - T E'$

$T' \rightarrow$

$E' \rightarrow$

```
void E() :
{
{
    T() Eprime()
}
}

void Eprime() :
{
{
    ( "+" T() Eprime() )
  | ( "-" T() Eprime() )
  | {} /* empty */
}
}
```

```
void T() :
{
{
    F() Tprime()
}
}

void Tprime() :
{
{
    ( "*" F() Tprime() )
  | ( "/" F() Tprime() )
  | {} /* empty */
}
}
```

```
void F() :
{
{
    <NUM>
  | "(" E() ")"
}
}
```

```
TOKEN :
{
    < NUM: (["0"-"9"])+ >
}
```

# Local Lookahead (using JavaCC)

<i>Statement</i>	→ { <i>Statement</i> * }	BlockStm
<i>Statement</i>	→ <i>id</i> = <i>Exp</i> ;	AssignmentStm
<i>Statement</i>	→ <i>id</i> [ <i>Exp</i> ] = <i>Exp</i> ;	ArrayAssignmentStm

Given token <ID>, the parser needs to know which rule to apply. It needs to lookahead the next token:

```
void Statement() :
{
{
    BlockStm()
|
    LOOKAHEAD(2)
    AssignmentStm()
|
    LOOKAHEAD(2)
    ArrayAssignmentStm()
}
```

```
void AssignmentStm() :
{
{ <ID> "=" Expression() ";" }

void ArrayAssignmentStm() :
{
{
    <ID> "[" Expression() "]" "="
    Expression() ";"
}
```



# Syntactic lookahead

<i>Exp</i> → <i>PrimaryExp</i> <i>op</i> <i>PrimaryExp</i>	OpExpression
<i>Exp</i> → <i>PrimaryExp</i> [ <i>Exp</i> ]	ArrayLookup
<i>Exp</i> → <i>PrimaryExp</i> . <b>length</b>	ArrayLength
<i>Exp</i> → <i>PrimaryExp</i> . <i>Id</i> ( <i>ExpList</i> )	MethodCall
<i>Exp</i> → <i>PrimaryExp</i>	PrimaryExpression

```
void Expression() :
{
{
    LOOKAHEAD( PrimaryExpression() <OP> ) OpExpression()
|
...
    LOOKAHEAD( PrimaryExpression() "[" ) ArrayLookup()
|
    LOOKAHEAD( PrimaryExpression() "." "length" ) ArrayLength()
|
    LOOKAHEAD( PrimaryExpression() "." Identifier() "(" )
    MethodCall()
|
    PrimaryExpression()
}
}
```

# Syntactic lookahead

---

```
void OpExpression() :  
{  
  {  
    PrimaryExpression() "&&" PrimaryExpression()  
  }  
}
```

```
void ArrayLength() :  
{  
  {  
    PrimaryExpression() "." "length"  
  }  
  
void MethodCall() :  
{  
  {  
    PrimaryExpression() "."  
  }  
  Identifier()  
  "(" ( ExpressionList() )? ")"  
}
```

```
void PrimaryExpression() :  
{  
  {  
    IntegerLiteral()  
    |  
    TrueLiteral()  
    |  
    FalseLiteral()  
    |  
    Identifier()  
    | ...  
  }  
}
```

# Semantic actions

---

- Each terminal and non-terminal has their own type of semantic value.
- Terminal (token) semantic values are the tokens returned by the lexical analyser (type Token in JavaCC).

E.g.

$t = \text{<RATIONAL>}$

where  $t$  is of type Token.

# Semantic actions

---

- Non-terminals semantic values are given depending on what you want the rules to do.
- Semantic action for rule  $A \rightarrow B C D$ 
  - returns type associated with  $A$
  - can build this from values associated with  $B, C, D$
- JavaCC allows us to intersperse actions within rules (written in {...})

# Example: simple expression evaluator

```
TOKEN :  
{  
  < NUM: (["0"-"9"])+ > | < EOL: "\n" >  
}
```

```
int S() :  
{ int s; }  
{  
  s=E() <EOL> { return s; }  
  | <EOL>  
  | <EOF>  
}
```

```
int E() :  
{ int e; int t; }  
{  
  e=T() ( "+" t=T() { e=e+t; }  
          | "-" t=T() { e=e-t; } ) *  
  { return e; }  
}
```

```
int T() :  
{ int t; int f; }  
{  
  t=F() ( "*" f=F() { t=t*f; }  
          | "/" f=F() { t=t/f; } ) *  
  { return t; }  
}
```

```
int F() :  
{ Token t; int result; }  
{  
  t=<NUM>  
  { return Integer.parseInt(t.image); }  
  | "(" result=E() ")"  
  { return result; }  
}
```

# JavaCC actions

---

non-terminals can deliver values

we can declare some variables to use in actions

we can assign to variables from terminals and non-terminals

```
int E() :  
{ int e; int t; }  
{  
  e=T() ( "+" t=T() { e=e+t; }  
    | "-" t=T() { e=e-t; } )*  
  { return e; }  
}
```

we can write any Java code in actions

this is where the non-terminal value is delivered