

# Project INFO F 403 : Compilateur Perl

RODRIGUEZ Paul, VACCARI Eric

2 mars 2013

## Table des matières

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Unités lexicales</b>                             | <b>3</b> |
| 1.1      | Tableau . . . . .                                   | 3        |
| 1.2      | Remarques . . . . .                                 | 4        |
| <b>2</b> | <b>DFA</b>  | <b>5</b> |
| 2.1      | Variables, comparateurs, blocs, littéraux . . . . . | 5        |
| 2.2      | Else, elsif et identifier . . . . .                 | 6        |
| 2.3      | Opérateurs et divers . . . . .                      | 7        |
| 2.4      | Remarques . . . . .                                 | 7        |
| <b>3</b> | <b>Grammaire LL(1)</b>                              | <b>7</b> |

# 1 Unités lexicales

## 1.1 Tableau

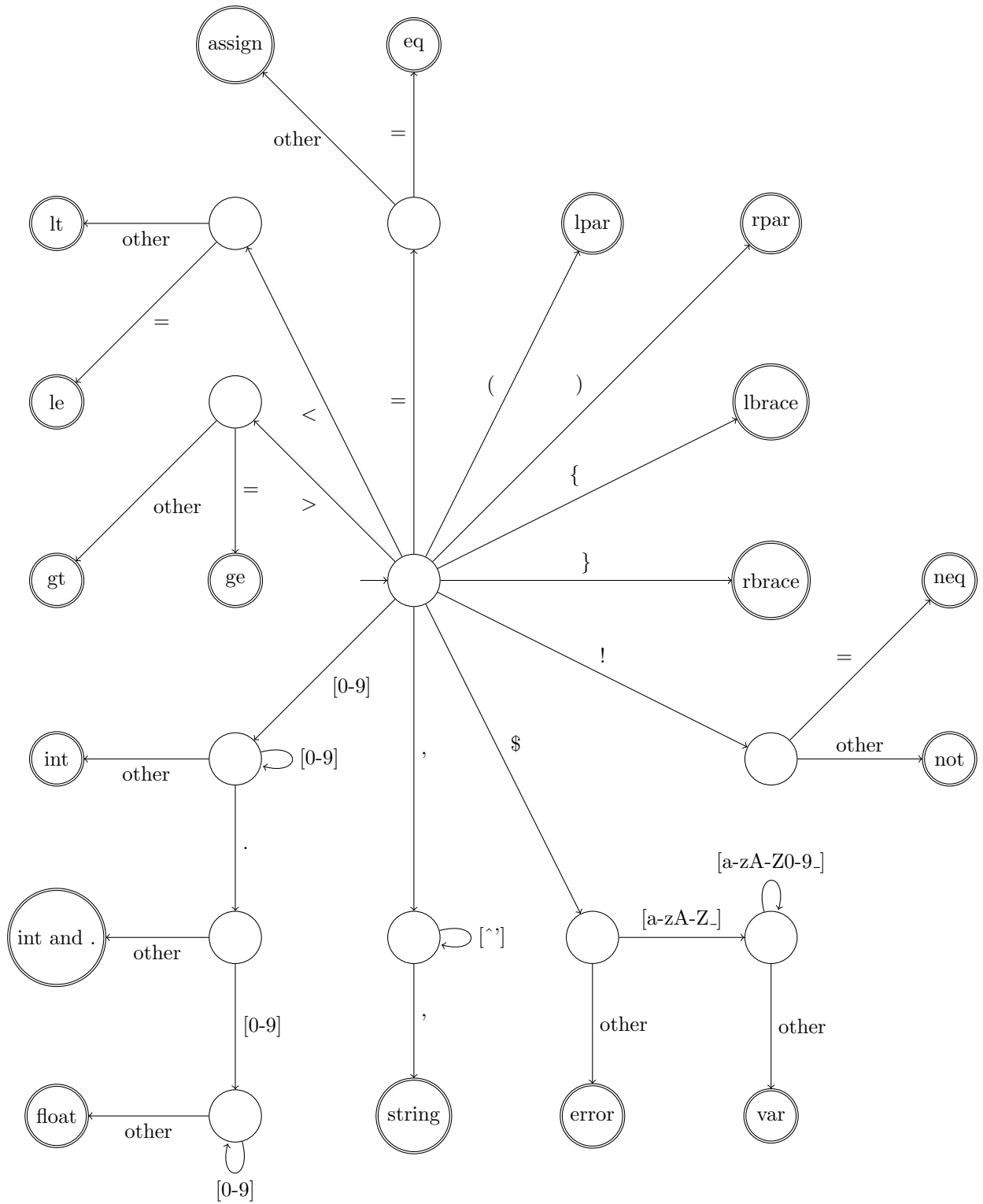
| Nom            | Regex                                  |
|----------------|--|
| var            | <code>[\$[a-zA-Z_][a-zA-Z0-9_]*</code> |
| identifier     | <code>[a-zA-Z_][a-zA-Z0-9_]*</code>    |
| integer        | <code>[0-9]+</code>                    |
| float          | <code>{integer}\.{integer}</code>      |
| string         | <code>'[^']*'</code>                   |
| space          | <code>[\t\n ]</code>                   |
| comment        | <code>#.*\n</code>                     |
| lbrace         | <code>\{</code>                        |
| rbrace         | <code>\}</code>                        |
| lpar           | <code>\(</code>                        |
| rpar           | <code>\)</code>                        |
| semicolon      | <code>;</code>                         |
| call_mark      | <code>&amp;</code>                     |
| plus           | <code>\+</code>                        |
| minus          | <code>\-</code>                        |
| times          | <code>\*</code>                        |
| divide         | <code>\/</code>                        |
| not            | <code>!</code>                         |
| notletters     | <code>not</code>                       |
| lazy_and       | <code>&amp;&amp;</code>                |
| lazy_or        | <code>  </code>                        |
| equals         | <code>==</code>                        |
| eq             | <code>eq</code>                        |
| different      | <code>!=</code>                        |
| ne             | <code>ne</code>                        |
| lower          | <code>&lt;</code>                      |
| lt             | <code>lt</code>                        |
| greater        | <code>&gt;</code>                      |
| gt             | <code>gt</code>                        |
| lower_equals   | <code>&lt;=</code>                     |
| le             | <code>le</code>                        |
| greater_equals | <code>&gt;=</code>                     |
| ge             | <code>ge</code>                        |
| comma          | <code>,</code>                         |
| concat_mark    | <code>\.</code>                        |
| assign_mark    | <code>=</code>                         |
| sub            | <code>sub</code>                       |
| if             | <code>if</code>                        |
| else           | <code>else</code>                      |
| elsif          | <code>elsif</code>                     |
| unless         | <code>unless</code>                    |
| return         | <code>return</code>                    |
| defined        | <code>defined</code>                   |
| int            | <code>int</code>                       |
| length         | <code>length</code>                    |
| print          | <code>print</code>                     |
| scalar         | <code>scalar</code>                    |
| substr         | <code>substr</code>                    |

## 1.2 Remarques

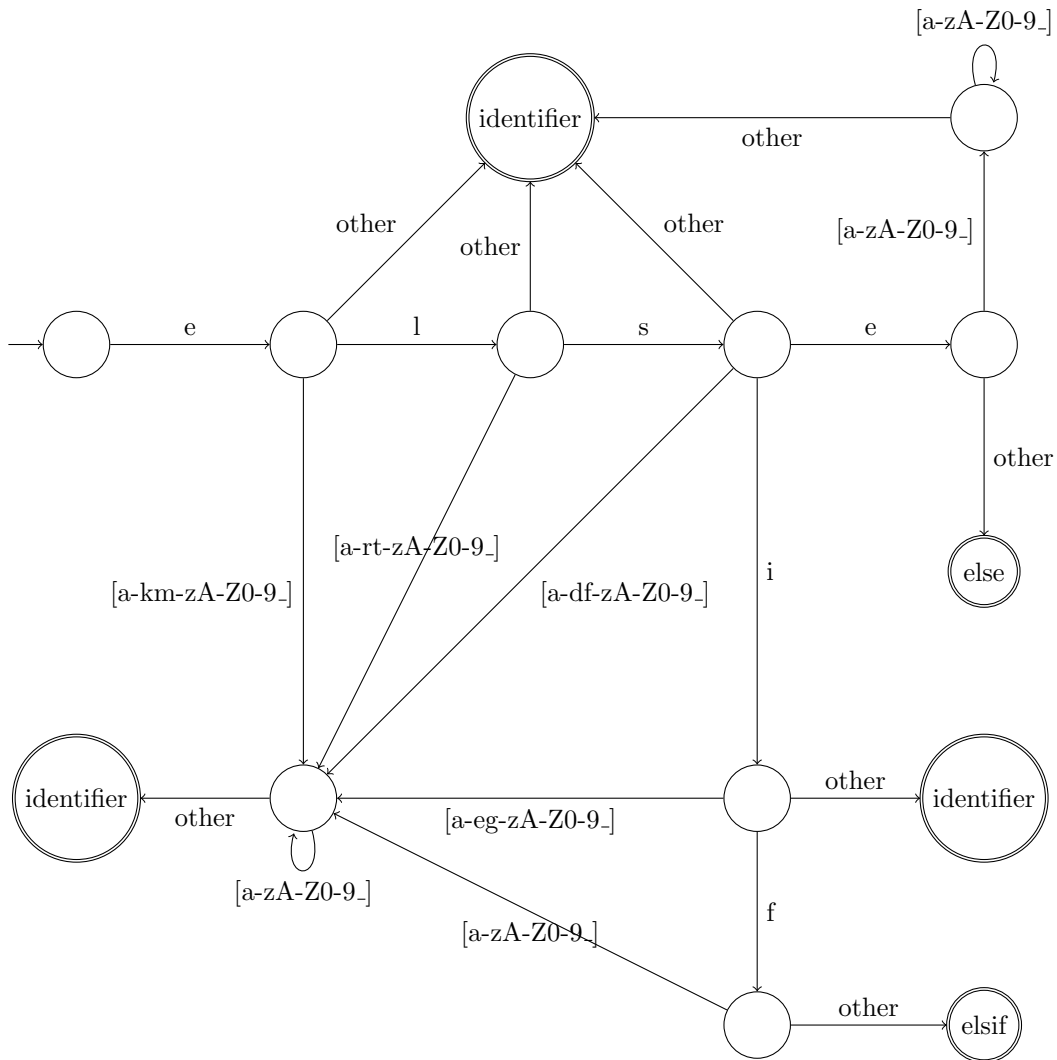
La syntaxe complète de Perl concernant les noms de variables est beaucoup plus compliquée mais concerne des fonctionnalités (packages) hors du cadre de ce projet, ce pourquoi nous nous sommes limités aux règles les plus simples.

## 2 DFA

### 2.1 Variables, comparateurs, blocs, littéraux

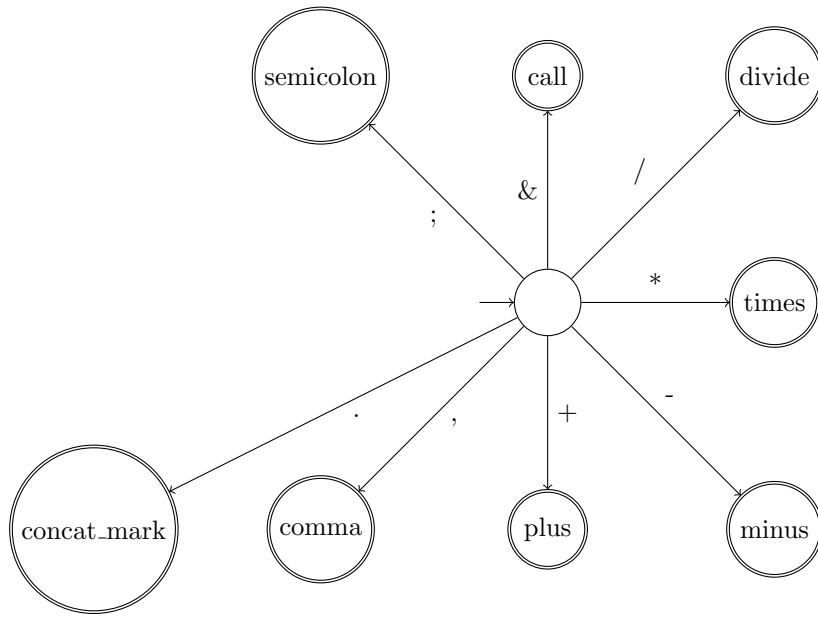


## 2.2 Else, elsif et identifier



Nous avons décidé de ne représenter que ces deux exemples, tous les mots clés fonctionnent sur le même principe.

## 2.3 Opérateurs et divers



## 2.4 Remarques

Certains tokens sont identifiables dès que leur dernier caractère a été lu (par exemple les accolades), d'autres nécessitent la lecture du caractère suivant le dernier (par exemple, pour terminer un entier il faut lire autre chose qu'un chiffre). Dans ce deuxième cas, après avoir identifié le token la lecture du dernier caractère est annulée, il servira comme premier caractère du token suivant.

## 3 Grammaire LL(1)

|      |   |  |
|------|---|--|
| [1]  | $\langle \text{PROGRAM} \rangle$            | $\rightarrow \langle \text{PROGRAM\_F} \rangle \langle \text{PROGRAM\_V} \rangle$  |
| [2]  | $\langle \text{PROGRAM\_V} \rangle$         | $\rightarrow \langle \text{PROGRAM\_F} \rangle \langle \text{PROGRAM\_V} \rangle$  |
| [3]  |   | $\rightarrow \epsilon$   |
| [4]  | $\langle \text{PROGRAM\_F} \rangle$         | $\rightarrow \langle \text{FUNCTION} \rangle$  |
| [5]  |   | $\rightarrow \langle \text{INSTRUCTION} \rangle$   |
| [6]  | $\langle \text{FUNCTION} \rangle$           | $\rightarrow \text{SUB IDENTIFIER } \langle \text{FUNCTION\_ARGUMENT} \rangle$<br>$\text{LBRACE } \langle \text{INSTRUCTION\_LIST} \rangle \text{ RBRACE}$ |
| [7]  | $\langle \text{FUNCTION\_ARGUMENT} \rangle$ | $\rightarrow \text{LPAR } \langle \text{ARGUMENT\_LIST} \rangle \text{ RPAR}$  |
| [8]  |   | $\rightarrow \epsilon$   |
| [9]  | $\langle \text{ARGUMENT\_LIST} \rangle$     | $\rightarrow \text{VAR } \langle \text{ARGUMENT\_LIST\_V} \rangle$   |
| [10] |   | $\rightarrow \epsilon$   |
| [11] | $\langle \text{ARGUMENT\_LIST\_V} \rangle$  | $\rightarrow \text{COMMA VAR } \langle \text{ARGUMENT\_LIST\_V} \rangle$   |
| [12] |   | $\rightarrow \epsilon$   |
| [13] | $\langle \text{INSTRUCTION\_LIST} \rangle$  | $\rightarrow \langle \text{INSTRUCTION} \rangle \langle \text{INSTRUCTION\_LIST} \rangle$  |
| [14] |   | $\rightarrow \epsilon$   |
| [15] | $\langle \text{INSTRUCTION} \rangle$        | $\rightarrow \langle \text{EXPRESSION} \rangle \langle \text{INSTRUCTION\_F} \rangle \text{ SEMICOLON}$  |

|      |  |  |
|------|--|--|
| [16] |  | → RETURN $\langle$ EXPRESSION $\rangle$ $\langle$ INSTRUCTION_F $\rangle$ SEMICOLON  |
| [17] |  | → LBRACE $\langle$ INSTRUCTION_LIST $\rangle$ RBRACE   |
| [18] |  | → $\langle$ CONDITION $\rangle$ $\langle$ EXPRESSION $\rangle$ LBRACE $\langle$ INSTRUCTION_LIST $\rangle$<br>RBRACE $\langle$ CONDITION_END $\rangle$ |
| [19] | $\langle$ INSTRUCTION_F $\rangle$      | → $\langle$ CONDITION $\rangle$ $\langle$ EXPRESSION $\rangle$   |
| [20] |  | → $\epsilon$   |
| [21] | $\langle$ CONDITION $\rangle$          | → IF   |
| [22] |  | → UNLESS   |
| [23] | $\langle$ CONDITION_END $\rangle$      | → ELIF $\langle$ EXPRESSION $\rangle$ LBRACE $\langle$ INSTRUCTION_LIST $\rangle$<br>RBRACE $\langle$ CONDITION_END $\rangle$                          |
| [24] |  | → ELSE LBRACE $\langle$ INSTRUCTION_LIST $\rangle$ RBRACE  |
| [25] |  | → $\epsilon$   |
| [26] | $\langle$ EXPRESSION $\rangle$         | → NOTLETTERS $\langle$ EXPRESSION $\rangle$  |
| [27] |  | → $\langle$ EXPRESSION_TWO $\rangle$   |
| [28] | $\langle$ EXPRESSION_TWO $\rangle$     | → $\langle$ EXPRESSION_THREE $\rangle$ $\langle$ EXPRESSION_TWO_V $\rangle$  |
| [29] | $\langle$ EXPRESSION_TWO_V $\rangle$   | → ASSIGN_MARK $\langle$ EXPRESSION_THREE $\rangle$ $\langle$ EXPRESSION_TWO_V $\rangle$  |
| [30] |  | → $\epsilon$   |
| [31] | $\langle$ EXPRESSION_THREE $\rangle$   | → $\langle$ EXPRESSION_FOUR $\rangle$ $\langle$ EXPRESSION_THREE_V $\rangle$   |
| [32] | $\langle$ EXPRESSION_THREE_V $\rangle$ | → LAZY_OR $\langle$ EXPRESSION_FOUR $\rangle$ $\langle$ EXPRESSION_THREE_V $\rangle$   |
| [33] |  | → $\epsilon$   |
| [34] | $\langle$ EXPRESSION_FOUR $\rangle$    | → $\langle$ EXPRESSION_FIVE $\rangle$ $\langle$ EXPRESSION_FOUR_V $\rangle$  |
| [35] | $\langle$ EXPRESSION_FOUR_V $\rangle$  | → LAZY_AND $\langle$ EXPRESSION_FIVE $\rangle$ $\langle$ EXPRESSION_FOUR_V $\rangle$   |
| [36] |  | → $\epsilon$   |
| [37] | $\langle$ EXPRESSION_FIVE $\rangle$    | → $\langle$ EXPRESSION_SIX $\rangle$ $\langle$ EXPRESSION_FIVE_V $\rangle$   |
| [38] | $\langle$ EXPRESSION_FIVE_V $\rangle$  | → $\langle$ EXPRESSION_FIVE_F $\rangle$ $\langle$ EXPRESSION_SIX $\rangle$   |
| [39] |  | → $\epsilon$   |
| [40] | $\langle$ EXPRESSION_FIVE_F $\rangle$  | → DIFFERENT  |
| [41] |  | → EQ   |
| [42] |  | → EQUALS   |
| [43] |  | → NE   |
| [44] | $\langle$ EXPRESSION_SIX $\rangle$     | → $\langle$ EXPRESSION_SEVEN $\rangle$ $\langle$ EXPRESSION_SIX_V $\rangle$  |
| [45] | $\langle$ EXPRESSION_SIX_V $\rangle$   | → $\langle$ EXPRESSION_SIX_F $\rangle$ $\langle$ EXPRESSION_SEVEN $\rangle$  |
| [46] |  | → $\epsilon$   |
| [47] | $\langle$ EXPRESSION_SIX_F $\rangle$   | → GE   |
| [48] |  | → GREATER  |
| [49] |  | → GREATER_EQUALS   |
| [50] |  | → GT   |
| [51] |  | → LE   |
| [52] |  | → LOWER  |
| [53] |  | → LOWER_EQUALS   |
| [54] |  | → LT   |



|      |  |   |
|------|--|---|
| [55] | $\langle \text{EXPRESSION\_SEVEN} \rangle$       | $\rightarrow \langle \text{EXPRESSION\_EIGHT} \rangle \langle \text{EXPRESSION\_SEVEN\_V} \rangle$  |
| [56] | $\langle \text{EXPRESSION\_SEVEN\_V} \rangle$    | $\rightarrow \langle \text{EXPRESSION\_SEVEN\_F} \rangle \langle \text{EXPRESSION\_EIGHT} \rangle$<br>$\langle \text{EXPRESSION\_SEVEN\_V} \rangle$ |
| [57] |  | $\rightarrow \epsilon$  |
| [58] | $\langle \text{EXPRESSION\_SEVEN\_F} \rangle$    | $\rightarrow \text{PLUS}$   |
| [59] |  | $\rightarrow \text{MINUS}$  |
| [60] |  | $\rightarrow \text{CONCAT\_MARK}$   |
| [61] | $\langle \text{EXPRESSION\_EIGHT} \rangle$       | $\rightarrow \langle \text{EXPRESSION\_NINE} \rangle \langle \text{EXPRESSION\_EIGHT\_V} \rangle$   |
| [62] | $\langle \text{EXPRESSION\_EIGHT\_V} \rangle$    | $\rightarrow \langle \text{EXPRESSION\_EIGHT\_F} \rangle \langle \text{EXPRESSION\_NINE} \rangle$<br>$\langle \text{EXPRESSION\_EIGHT\_V} \rangle$  |
| [63] |  | $\rightarrow \epsilon$  |
| [64] | $\langle \text{EXPRESSION\_EIGHT\_F} \rangle$    | $\rightarrow \text{TIMES}$  |
| [65] |  | $\rightarrow \text{DIVIDE}$   |
| [66] | $\langle \text{EXPRESSION\_NINE} \rangle$        | $\rightarrow \langle \text{EXPRESSION\_NINE\_V} \rangle \langle \text{EXPRESSION\_TEN} \rangle$   |
| [67] | $\langle \text{EXPRESSION\_NINE\_V} \rangle$     | $\rightarrow \langle \text{EXPRESSION\_NINE\_F} \rangle \langle \text{EXPRESSION\_NINE\_V} \rangle$   |
| [68] |  | $\rightarrow \epsilon$  |
| [69] | $\langle \text{EXPRESSION\_NINE\_F} \rangle$     | $\rightarrow \text{NOT}$  |
| [70] |  | $\rightarrow \text{PLUS}$   |
| [71] |  | $\rightarrow \text{MINUS}$  |
| [72] | $\langle \text{EXPRESSION\_TEN} \rangle$         | $\rightarrow \text{LPAR} \langle \text{EXPRESSION} \rangle \text{RPAR}$   |
| [73] |  | $\rightarrow \langle \text{SIMPLE\_EXPRESSION} \rangle$   |
| [74] | $\langle \text{SIMPLE\_EXPRESSION} \rangle$      | $\rightarrow \langle \text{FUNCTION\_CALL} \rangle$   |
| [75] |  | $\rightarrow \text{INTEGER}$  |
| [76] |  | $\rightarrow \text{FLOAT}$  |
| [77] |  | $\rightarrow \text{STRING}$   |
| [78] |  | $\rightarrow \text{VAR}$  |
| [79] | $\langle \text{FUNCTION\_CALL} \rangle$          | $\rightarrow \text{CALL\_MARK IDENTIFIER LPAR} \langle \text{ARGUMENT\_CALL\_LIST} \rangle$<br>$\text{RPAR}$  |
| [80] | $\langle \text{ARGUMENT\_CALL\_LIST} \rangle$    | $\rightarrow \langle \text{EXPRESSION} \rangle \langle \text{ARGUMENT\_CALL\_LIST\_V} \rangle$  |
| [81] |  | $\rightarrow \epsilon$  |
| [82] | $\langle \text{ARGUMENT\_CALL\_LIST\_V} \rangle$ | $\rightarrow \text{COMMA} \langle \text{EXPRESSION} \rangle \langle \text{ARGUMENT\_CALL\_LIST\_V} \rangle$   |
| [83] |  | $\rightarrow \epsilon$  |

remarques : 1) on a enlevé la possibilité d'omettre les () autour des listes d'arguments a l'appel d'une fonction, c'est pas LL(1) car le follow de expression prend le follow de toutes les expressions

2) On a supprimé le "not", car sa priorité faible engendre un comportement non-LL(1), le follow de expression prend le follow de toutes les expressions