

UNIVERSITÉ DE SHERBROOKE  
Faculté de génie  
Département de génie électrique et génie informatique

## **RAPPORT DE L'APP 5**

Éléments de compilation  
GIF340

Présenté à  
Ahmed Khoumsi

Présenté par  
Équipe numéro 2  
Matthieu Daoust – DAOM2504  
William Pépin – PEPW3101  
Gabriel Vachon – VACG3501

Sherbrooke – 5 avril 2023

# TABLE DES MATIÈRES

|           |  |           |
|-----------|--|-----------|
| <b>1.</b> | <b>Analyse Lexicale</b>  | <b>1</b>  |
| 1.1       | Opérateurs   | 1         |
| 1.2       | Parenthèses  | 2         |
| 1.3       | Nombre   | 2         |
| 1.4       | Variable   | 3         |
| 1.5       | Construction de l'automate complet                                 | 4         |
| <b>2.</b> | <b>Analyse syntaxique</b>  | <b>5</b>  |
| 2.1       | Réalisation de la grammaire de base                                | 6         |
| 2.2       | Transformation de la grammaire pour un algorithme LL(1)            | 7         |
| 2.2.1     | Définition des contraintes d'une grammaire dite LL(1)              | 7         |
| 2.2.2     | Modification de la grammaire pour respecter la première contrainte | 7         |
| 2.2.3     | Modification de la grammaire pour respecter la deuxième contrainte | 9         |
| 2.3       | Grammaire transformé   | 10        |
| <b>3.</b> | <b>Résultats des tests de l'implémentation</b>                     | <b>11</b> |
| <b>4.</b> | <b>Références</b>  | <b>13</b> |

## LISTE DES FIGURES

|   |   |
|---|---|
| Figure 1 : Automate à état fini de la détection des opérateurs                            | 1 |
| Figure 2 : Automate a état fini de la détection des parenthèses                           | 2 |
| Figure 3 : Automate à état fini pour un opérande à chiffres                               | 3 |
| Figure 4 : Automate à état fini pour un opérande contenant des caractères maj, min, « _ » | 4 |
| Figure 5 : Automate complet de l'analyseur lexicale                                       | 5 |
| Figure 6 : Preuve mathématique de la première règle                                       | 8 |
| Figure 7 : Démarches de la modification de la grammaire pour la première règle            | 8 |
| Figure 8 : Démarches pour le respect de la deuxième règle de la grammaire                 | 9 |

## LISTE DES TABLEAUX

|   |    |
|---|----|
| Tableau 1 : Définition d'expression arithmétique par une grammaire [2]                  | 6  |
| Tableau 2 : Définition de la grammaire LL(1)  | 9  |
| Tableau 3 : Définition de la grammaire LL(1)  | 10 |
| Tableau 4 : Résultats des tests de l'implémentation de l'analyseur lexical              | 11 |
| Tableau 5 : Résultats des tests de l'implémentation de la méthode de descente récursive | 12 |

# 1. ANALYSE LEXICALE

L'analyse lexicale permet de détecter et d'identifier les unités lexicales présentes dans un code source. En effectuant une analyse caractère par caractère, le programme permet de fournir une série d'unité lexicale à l'analyseur syntaxique. Elle permet aussi de trouver des erreurs lexicales, soit des erreurs dans l'écriture des identificateurs, des mots clés, des opérateurs, des nombres, etc. Dans le cadre de la problématique, l'analyseur lexical doit identifier les opérateurs suivants :  $+$   $-$   $*$   $/$ , les unités de contrôle suivant :  $($   $,$   $)$ , les nombres naturels et les variables définies dans le guide [1]. Les sections suivantes présentent les expressions régulières ainsi que les automates requis afin de détecter les unités lexicales du contexte défini, en finissant par la construction d'un automate complet qui définit la logique globale de l'analyseur lexical.

## 1.1 OPÉRATEURS

Définition de l'alphabet, du langage et des expressions régulières :

$$\begin{aligned}\Sigma &= \{+, -, *, /\} & (1) \\ L_{\Sigma} &= \{(+, -, *, /)\} \\ ER_+ &= + & ER_- &= - & ER_* &= * & ER_/_ &= /\end{aligned}$$

Automate :

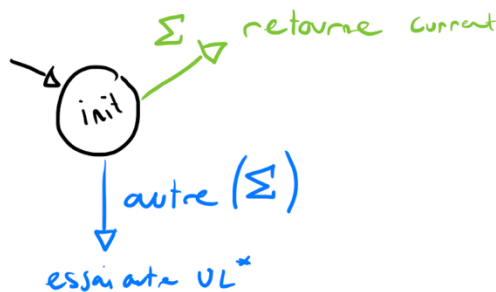


Figure 1 : Automate à état fini de la détection des opérateurs

\*La fonction « essai autre UL » n'est pas une vraie fonction de l'automate, cette flèche sera définie à la fin de la section, lors de la création de l'automate complet.

En somme, l'alphabet des opérateurs constitue des opérateurs directement et l'automate détecte l'opérateur et le retourne. Dans un autre cas, elle vérifie si ce n'est pas une autre unité lexicale.

## 1.2 PARENTHÈSES

Définition de l'alphabet, du langage et des expressions régulières :

$$\begin{aligned}\Sigma &= \{ (, ) \} \\ L_{\Sigma} &= \{ (, ) \} \\ ER_{(} &= ( \quad ER_{)} = )\end{aligned}\tag{2}$$

Automate :

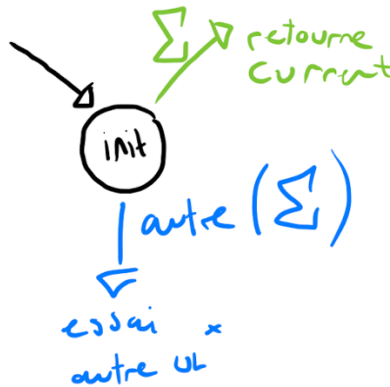


Figure 2 : Automate a état fini de la détection des parenthèses

\*La fonction « essai autre UL » n'est pas une vraie fonction de l'automate, cette flèche sera définie à la fin de la section, lors de la création de l'automate complet.

En somme, l'alphabet des parenthèses constitue des parenthèses ouvrantes et fermantes. L'automate détecte alors les parenthèses individuellement. Dans un autre cas, elle vérifie si ce n'est pas une autre unité lexicale.

## 1.3 NOMBRE

Définition de l'alphabet, du langage et des expressions régulières :

$$\begin{aligned}\Sigma &= \{0,1,2,3,4,5,6,7,8,9\} \\ L_{\Sigma} &= \mathbb{N} \\ ER &= (0|1|2|3|4|5|6|7|8|9)^+\end{aligned}\tag{3}$$

Automate :

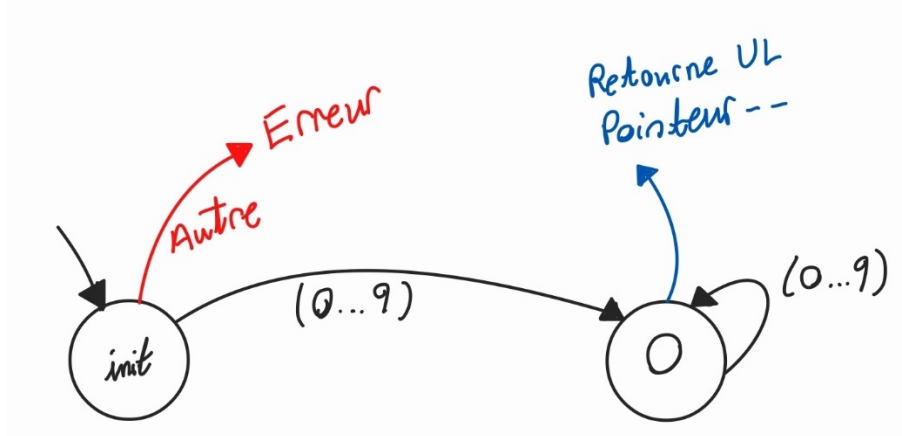


Figure 3 : Automate à état fini pour un opérande à chiffres

L'automate à état fini ci-dessus permet d'identifier un opérande composé de chiffres. À l'état initial, l'automate s'attend à recevoir un chiffre de l'ensemble mentionné plus haut. S'il n'en reçoit pas ou il reçoit un caractère non représenté dans l'ensemble, une erreur est retournée. Il arrive ensuite à l'état final, ici la chaîne de l'opérande a au moins un caractère. Si un caractère qui n'est pas représenté dans l'ensemble est lu, l'unité lexicale est retournée et on recule le pointeur.

## 1.4 VARIABLE

Définition de l'alphabet, du langage et des expressions régulières :

$$\Sigma = \{A, B, C, D, \dots, Z, a, b, c, d, \dots, z\} \quad (4)$$

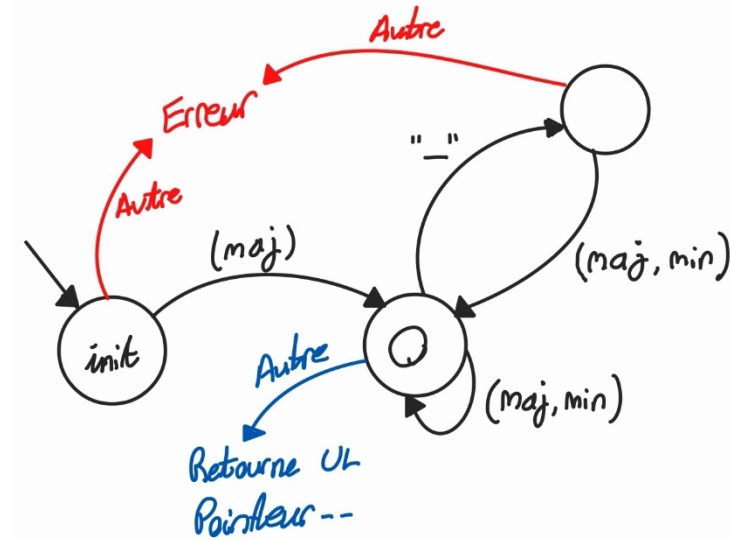
$$\Sigma_{maj,min} = \{A, B, C, D, \dots, Z, a, b, c, d, \dots, z\}$$

$$\Sigma_{maj} = \{A, B, C, D, \dots, Z\}$$

$$\Sigma_{min} = \{a, b, c, d, \dots, z\}$$

$$ER = \Sigma_{maj}[-?(\Sigma_{maj,min})]^*$$

**Automate :**



**Figure 4 : Automate à état fini pour un opérande contenant des caractères maj, min, « \_ »**

L'automate ci-dessus permet d'identifier une chaîne de caractères contenant des majuscules, minuscules ou tirets bas. À l'état initial, il est impératif d'avoir une lettre majuscule, si un caractère autre est identifié, une erreur est lancée [le premier caractère ne peut commencer par un tiret bas]. Après le premier caractère, une majuscule, une minuscule ou un tiret bas peut être ajouté. Il n'y a pas de limite de caractères majuscules ou minuscules dans la chaîne, cependant, il ne peut y avoir plusieurs tirets bas consécutifs et la chaîne ne peut finir par ce caractère. Pour cette raison, si l'état est « tiret\_bas », une erreur est lancée si le caractère n'est pas une lettre. À l'état final, si un autre caractère autre qu'une lettre ou un tiret bas est entré, on retourne l'unité lexicale et on recule le pointeur.

## 1.5 CONSTRUCTION DE L'AUTOMATE COMPLET

La réalisation suivante constitue l'automate complet de l'analyseur lexicale, elle incorpore les automates de chacune des sections dans un seul automate, sans ambiguïté. Les ensembles sont définies avec les noms suivants :

- Op = ensemble des opérateurs ( $\Sigma_{op} = \{+, -, *, /\}$ );
- Co = ensembles des parenthèses ( $\Sigma_{co} = \{ (, ) \}$ );
- No = ensemble des nombres ( $\Sigma_{no} = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ );
- Maj = ensemble des majuscules ( $\Sigma_{maj} = \{A, B, C, D, \dots, Z\}$ );
- Min = ensemble des minuscules ( $\Sigma_{min} = \{a, b, c, d, \dots, z\}$ );
- X ensemble des variables ( $\Sigma_x = \Sigma_{maj} \cup \Sigma_{min} \cup \{ \_ \}$ ).

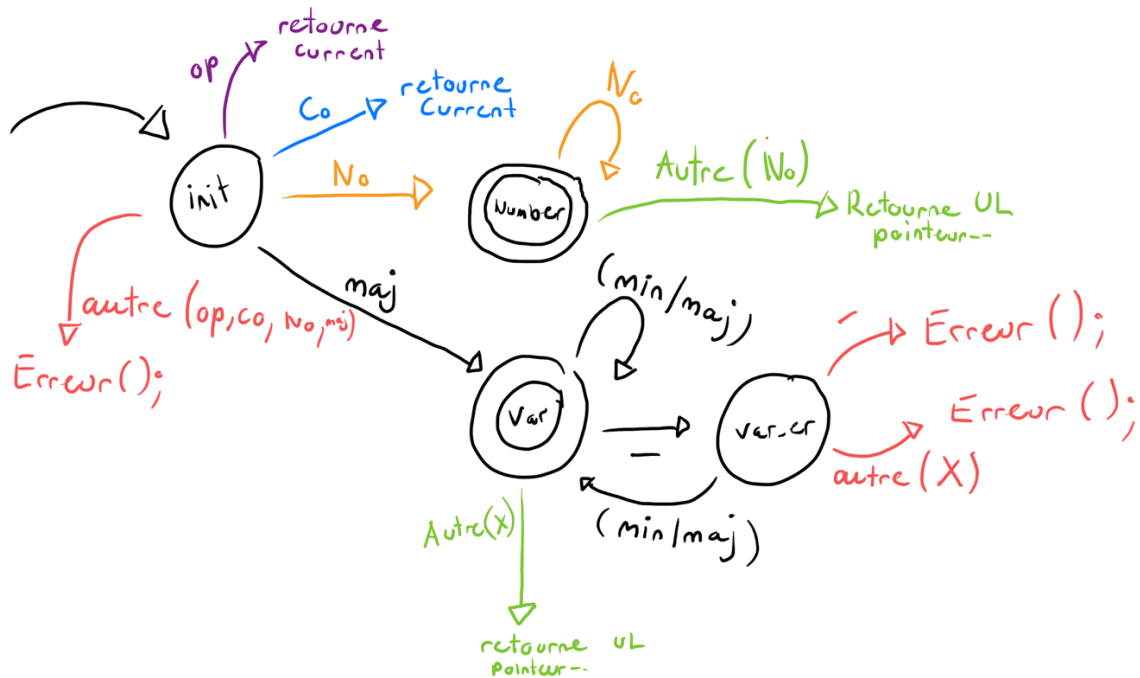


Figure 5 : Automate complet de l'analyseur lexicale

Les états sont les suivants : *init*, *number*, *var* et *var\_er*. L'état « *init* » est le point d'entrée de l'automate, elle permet de filtrer l'alphabet de l'analyseur lexical en incorporant une vérification sur les opérateurs, les parenthèses, les nombres et les variables, si le caractère lu n'est pas dans l'alphabet de ses 4 types d'unités lexicales, une erreur est retournée. Si c'est un opérateur ou une parenthèse, l'état retourne directement l'unité lexicale. Si c'est un nombre, l'automate passe à l'état « *number* » et effectue la logique de vérification des nombres. Si c'est une Majuscule, l'automate passe à l'état « *var* » et effectue la logique de vérification des variables. Cette logique retourne une erreur si la variable termine par « *\_* » ou contient deux « *\_* » à la suite de l'autre.

## 2. ANALYSE SYNTAXIQUE

L'analyseur syntaxique permet d'effectuer l'analyse des unités lexicales afin de trouver des erreurs de syntaxes. Elle permet de définir le comportement du compilateur de sorte à accepter [ou refuser] des expressions arithmétiques en fonction d'une grammaire. Une fois le regroupement des unités lexicales en structures grammaticales, l'analyseur syntaxique peut générer un arbre syntaxique abstrait qui représente la structure hiérarchique du flot d'unités lexicales. [2] Ainsi, pour réaliser un analyseur syntaxique, il est nécessaire de définir une



grammaire rigoureuse qui permet la vérification efficace des flots d'unités lexicales. L'algorithme utilisé pour la réalisation de cet analyseur syntaxique est l'algorithme LL [1].

## 2.1 RÉALISATION DE LA GRAMMAIRE DE BASE

Afin de réaliser une grammaire de base pour un analyseur syntaxique, celle-ci doit être de type « non-contextuelle ». C.-à-d. qu'elle doit respecter avoir deux caractéristiques suivantes : un seul symbole doit être présent sur la partie gauche de la règle et la partie droite de la règle ne doit pas nécessairement débiter par un symbole terminal. Le chapitre 3 des notes de cours offre une grammaire non contextuelle afin de modéliser des expressions arithmétiques [2] :

- $V_t$  = ensemble des symboles terminaux;
- $V_n$  = ensemble des symboles non-terminaux;
- $S$  = symbole initiale;
- $P$  = ensemble fini des règles de production;
- $Id$  = opérande.

**Tableau 1 : Définition d'expression arithmétique par une grammaire [2]**

**$V_t = \{id, nombre, +, -, *, /, (, )\}$**

**$V_n = \{E, T, F\}$**

**$S = E$**

**$P \{$**

| Symbole |               | Dérivation | Numéro de la règle |
|---------|---------------|------------|--------------------|
| E       | $\Rightarrow$ | T          | 1                  |
| E       | $\Rightarrow$ | T+T        | 2                  |
| E       | $\Rightarrow$ | T-T        | 3                  |
| T       | $\Rightarrow$ | F          | 4                  |
| T       | $\Rightarrow$ | F*F        | 5                  |
| T       | $\Rightarrow$ | F/F        | 6                  |
| F       | $\Rightarrow$ | (E)        | 7                  |
| F       | $\Rightarrow$ | id         | 8                  |
| F       | $\Rightarrow$ | nombre     | 9                  |

**$\}$**

Toutefois, cette grammaire ne respecte pas les contraintes d'une grammaire LL (1). Les sections suivantes présentent la démarche afin de transformer la grammaire suivante en grammaire qui respecte les contraintes définies par une grammaire LL (1).

## 2.2 TRANSFORMATION DE LA GRAMMAIRE POUR UN ALGORITHME LL(1)

### 2.2.1 DÉFINITION DES CONTRAINTES D'UNE GRAMMAIRE DITE LL(1)

Afin d'obtenir une grammaire LL (1), il est nécessaire de respecter les trois contraintes suivantes (assumant une grammaire non contextuelle) :

- 1- Pour chaque règle ayant la même partie gauche, l'ensemble premier de chaque règle ne doit pas avoir d'élément(s) en commun entre les ensembles [2] ;

$$A \rightarrow \alpha 1 \quad (5)$$

$$A \rightarrow \alpha 2$$

$$PREMIER(\alpha 1) \cap PREMIER(\alpha 2) = \emptyset$$

- 2- Pour chaque règle capable de devenir un ensemble vide, l'ensemble premier de cette règle ne doit pas avoir d'élément(s) en commun avec l'ensemble suivant de cette même règle [2].

$$A \rightarrow \alpha 1 \mid \epsilon \quad (6)$$

$$PREMIER(\alpha 1) \cap SUIVANT(\alpha 1) = \emptyset$$

Autrement dit, la grammaire LL (1) ne doit pas être récursive à gauche, et ne doit pas être ambiguë. Les sections suivantes vérifient les contraintes ci-dessus afin de modifier la grammaire pour qu'elle respecte les contraintes d'une grammaire LL (1).

### 2.2.2 MODIFICATION DE LA GRAMMAIRE POUR RESPECTER LA PREMIÈRE CONTRAINTES

La première contrainte n'est pas respectée ici. La raison est la suivante : l'ensemble des premiers des différentes dérivations de E partagent tous les premiers de T, c'est la même chose pour les dérivations de T où les premiers de T qui contient l'ensemble des premiers de F. cependant les ensembles des premiers des règles de dérivation de F sont tous des ensembles différents, ce qui respecte la première règle.

$$\begin{array}{lcl}
 E \rightarrow T & \longrightarrow & Pr(E) \subseteq Pr(T) \\
 E \rightarrow T+T & \longrightarrow & Pr(E) \subseteq Pr(T) \\
 E \rightarrow T-T & \longrightarrow & Pr(E) \subseteq Pr(T)
 \end{array} \left. \vphantom{\begin{array}{l} E \rightarrow T \\ E \rightarrow T+T \\ E \rightarrow T-T \end{array}} \right\} Pr(T) \cap Pr(T) \cap Pr(T) \neq \emptyset$$

donc ne respecte pas la première règle

$$\begin{array}{lcl}
 T \rightarrow F & \longrightarrow & Pr(T) \subseteq Pr(F) \\
 T \rightarrow F * F & \longrightarrow & Pr(T) \subseteq Pr(F) \\
 T \rightarrow F / F & \longrightarrow & Pr(T) \subseteq Pr(F)
 \end{array} \left. \vphantom{\begin{array}{l} T \rightarrow F \\ T \rightarrow F * F \\ T \rightarrow F / F \end{array}} \right\} Pr(F) \cap Pr(F) \cap Pr(F) \neq \emptyset$$

donc ne respecte pas la première règle

$$\begin{array}{lcl}
 F \rightarrow (E) & \longrightarrow & Pr(F) \subseteq Pr(E) \\
 F \rightarrow id & \longrightarrow & id \in Pr(F) \\
 F \rightarrow nombre & \longrightarrow & nombre \in Pr(F)
 \end{array} \left. \vphantom{\begin{array}{l} F \rightarrow (E) \\ F \rightarrow id \\ F \rightarrow nombre \end{array}} \right\} Pr(E) \cap id \cap nombre \cap "(" = \emptyset$$

donc respecte la première règle

Figure 6 : Preuve mathématique de la première règle

La grammaire peut être modifiée ainsi pour convenir à la première règle :

$$\begin{array}{lcl}
 E \rightarrow T \mid T+E \mid T-E & \longrightarrow & Pr(E) \subseteq Pr(T) \\
 T \rightarrow F \mid F * T \mid F / T & \longrightarrow & Pr(T) \subseteq Pr(F) \\
 F \rightarrow (E) \mid id \mid nombre & \longrightarrow & \{ "(", id, nombre \} \in Pr(F) \\
 & \longrightarrow & Pr(F) \subseteq Pr(E)
 \end{array}$$

Figure 7 : Démarches de la modification de la grammaire pour la première règle

La grammaire respectant la première règle est indiquée au tableau suivant.

Tableau 2 : Définition de la grammaire LL(1)

$V_t = \{\text{id, nombre, +, -, *, /, (, )}\}$

$V_n = \{E, T, F\}$

$S = E$

$P \{$

| Symbole | Dérivation  | Numéro de la règle |
|---------|---|--------------------|
| E       | $\Rightarrow T \mid T+E \mid T-E$                   | 1                  |
| T       | $\Rightarrow F \mid F*T \mid F/T$                   | 2                  |
| F       | $\Rightarrow (E) \mid \text{id} \mid \text{nombre}$ | 3                  |

$\}$

### 2.2.3 MODIFICATION DE LA GRAMMAIRE POUR RESPECTER LA DEUXIÈME CONTRAINTE

La deuxième contrainte est respectée dans ce cas-ci, car il n'y a aucune règle qui contient d'ensemble vide. Puisqu'il n'est pas possible d'obtenir d'ensemble vide dans aucune des règles, la deuxième règle est respectée par défaut.

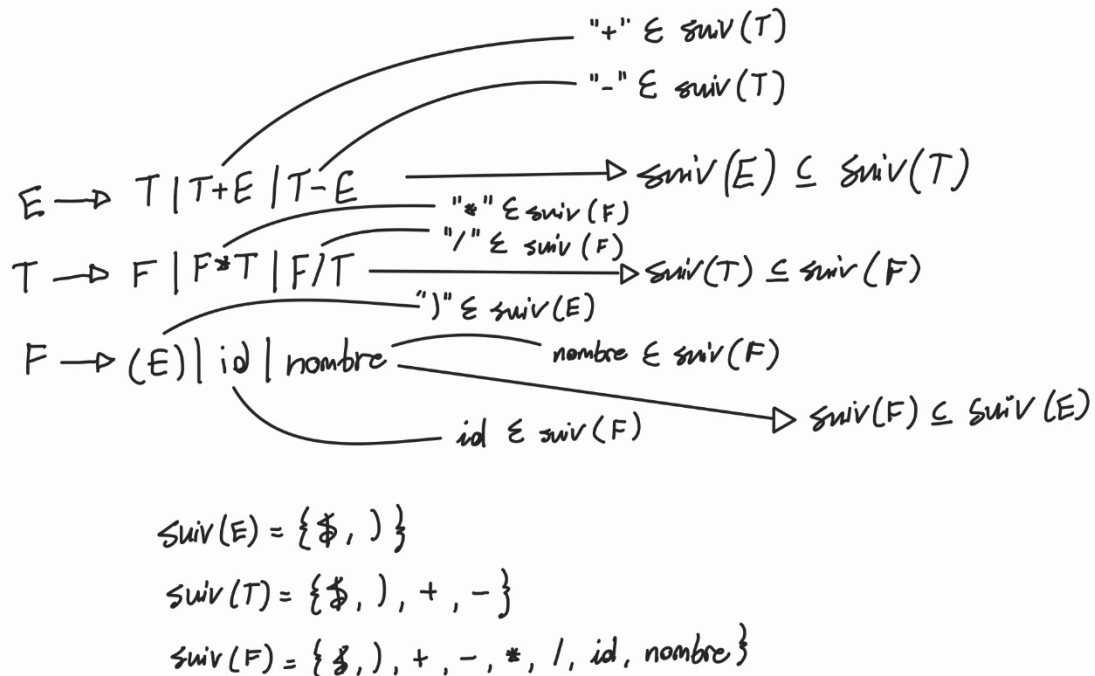


Figure 8 : Démarches pour le respect de la deuxième règle de la grammaire

La grammaire résultante est donc celle définie au **Tableau 2**.

## 2.3 GRAMMAIRE TRANSFORMÉ

En somme, en modifiant la grammaire de base, il est possible d'arriver à la grammaire suivante (même que le **Tableau 2**) :

**Tableau 3 : Définition de la grammaire LL(1)**

**Vt = {id, nombre, +, -, \*, /, (, )}**

**Vn = {E, T, F}**

**S = E**

**P {**

| Symbole |               | Dérivation        | Numéro de la règle |
|---------|---------------|-------------------|--------------------|
| E       | $\Rightarrow$ | T   T+E   T-E     | 1                  |
| T       | $\Rightarrow$ | F   F*T   F/T     | 2                  |
| F       | $\Rightarrow$ | (E)   id   nombre | 3                  |

**}**

Cette grammaire est valide pour un analyseur syntaxique utilisant une descente récursive comme algorithme.

### 3. RÉSULTATS DES TESTS DE L'IMPLÉMENTATION

Le présent tableau offre une synthèse des résultats des tests unitaires de l'implémentation du compilateur. Les méthodes sont décrites dans les classes de tests (*AnalLexTest* et *AnalSyntTest*) du projet.

**Tableau 4 : Résultats des tests de l'implémentation de l'analyseur lexical**

| ID | Test   | Description   | Résultats attendu                            | OK |
|----|--|---|--|----|
| 1  | <i>prochainTerminalW<br/>henParentheses()</i>                      | Vérifie la détection de parenthèses.  | Tableau d'unités lexicales contenant [(, )]. | OK |
| 2  | <i>prochainTerminalW<br/>henVariableStartsW<br/>ithLowercase()</i> | Vérifie l'erreur de variable commençant par une minuscule.                              | <i>AnalLexException</i>                      | OK |
| 3  | <i>prochainTerminalW<br/>henVariableEndsWit<br/>hUnderScore()</i>  | Vérifie l'erreur de variable terminant par un tiret bas.                                | <i>AnalLexException</i>                      | OK |
| 4  | <i>prochainTerminalW<br/>henNumbersOk()</i>                        | Vérifie un cas général d'un nombre.   | Tableau contenant les unités lexicales.      | OK |
| 5  | <i>prochainTerminalW<br/>henVariableHasTwo<br/>UnderScores()</i>   | Vérifie l'erreur de variable ayant deux tirets bas de suite.                            | <i>AnalLexException</i>                      | OK |
| 6  | <i>prochainTerminalW<br/>henVariableIsOk()</i>                     | Vérifie un cas général d'une variable.  | Tableau contenant les variables.             | OK |
| 7  | <i>prochainTerminalW<br/>henCharNotInAlpha<br/>bet()</i>           | Vérifie l'erreur d'un caractère qui n'est pas dans l'alphabet du compilateur.           | <i>AnalLexException</i>                      | OK |
| 8  | <i>prochainTerminalW<br/>henOperators()</i>                        | Vérifie la détection des 4 opérateurs (+, -, *, /).                                     | Tableau contenant les 4 opérateurs.          | OK |
| 9  | <i>prochainTerminalGe<br/>neralCase()</i>                          | Vérifie un cas général couvrant les unités lexicales. $((X\_a + Y\_b) * Z\_c / 59) - 4$ | Tableau contenant l'expression               | OK |

Tableau 5 : Résultats des tests de l'implémentation de la méthode de descente réursive

| ID | Test  | Description  | Résultats attendu                | OK |
|----|---|--|----------------------------------|----|
| 1  | <i>AnalSyntTestSyntaxeOperateurFini()</i>           | Test d'erreur de syntaxe lorsqu'elle termine avec un opérateur               | <i>AnalSyntException</i>         | OK |
| 2  | <i>AnalSyntTestLectureAST1()</i>                    | Test de lecture du cas suivant : $(X_a + Y_b) * Z_c / 2$                     | $((X_a + Y_b) * (Z_c / 2))$      | OK |
| 3  | <i>AnalSyntTestLectureAST2()</i>                    | Test de lecture du cas suivant : $(X_a - Y_b / Z_c + 59) + 4$                | $((X_a + (Y_b / Z_c) + 59) - 4)$ | OK |
| 4  | <i>AnalSyntTestLectureAST3()</i>                    | Test de lecture du cas suivant : $4 * 4 * 4 + 12 / 12$                       | $((4 * (4 * 4)) + (12 / 12))$    | OK |
| 5  | <i>AnalSyntTestSyntaxeParentheseFermeanteTrop()</i> | Test d'erreur de syntaxe lorsqu'elle a une parenthèse fermante de trop       | <i>AnalSyntException</i>         | OK |
| 6  | <i>AnalSyntTestSyntaxeParentheseOuvranteTrop()</i>  | Test d'erreur de syntaxe lorsqu'elle a une parenthèse ouvrante de trop       | <i>AnalSyntException</i>         | OK |
| 7  | <i>AnalSyntTestPostfixCase1()</i>                   | Test de lecture <i>postfix</i> du cas suivant : $(X_a + Y_b) * Z_c / 2$      | $X_a Y_b Z_c * 2 /$              | OK |
| 8  | <i>AnalSyntTestPostfixCase2()</i>                   | Test de lecture <i>postfix</i> du cas suivant : $(X_a - Y_b / Z_c + 59) + 4$ | $X_a Y_b Z_c / - 59 + 4 +$       | OK |
| 9  | <i>AnalSyntTestPostfixCase3()</i>                   | Test de lecture <i>postfix</i> du cas suivant : $4 * 4 * 4 + 12 / 12$        | $44 * 4 * 12 12 / + "$           | OK |
| 10 | <i>AnalSyntTestEvaluationAST1()</i>                 | Test d'évaluation du cas suivant : $(X_a + Y_b) * Z_c / 2$                   | 19                               | OK |
| 11 | <i>AnalSyntTestEvaluationAST2()</i>                 | Test d'évaluation du cas suivant : $(X_a + Y_b / Z_c + 59) - 4$              | 69                               | OK |
| 12 | <i>AnalSyntTestEvaluationAST3()</i>                 | Test d'évaluation du cas suivant : $4 * 4 * 4 + 12 / 12$                     | 65                               | OK |
| 13 | <i>AnalSyntTestSyntaxeManqueVariable()</i>          | Tests d'évaluation du cas suivant : $(X_a + ) * Z_c / 2$                     | <i>AnalSyntException</i>         | OK |

## 4. RÉFÉRENCES

- [1] A. Khoumsi, Développement d'un analyseur syntaxique d'expressions arithmétiques, Sherbrooke: Faculté de génie, Université de Sherbrooke, 2023.
- [2] A. Khoumsi, GIF 340 : Éléments de compilation, Sherbrooke: Faculté de génie, Université de Sherbrooke, 2005.