

RÉPUBLIQUE ALGÉRIENNE DÉMOCRATIQUE ET POPULAIRE  
MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR ET DE LA RECHERCHE SCIENTIFIQUE  
UNIVERSITÉ DE GHARDAIA  
FACULTÉ DES SCIENCES ET DE LA TECHNOLOGIE  
DÉPARTEMENT DES MATHÉMATIQUES ET D'INFORMATIQUE

---

# ANALYSES FRONTALES DANS LES COMPILATEURS

COURS, EXERCICES ET PROJETS DE PROGRAMMATION

---

**Slimane OULAD-NAOUI**

*Maître de conférences*

*ouladnaoui@univ-ghardaia.dz*

---

*Année universitaire 2021-2022*

# Table des matières

<b>Liste des figures</b>	<b>ii</b>
<b>Liste des tables et algorithmes</b>	<b>iv</b>
<b>Introduction</b>	<b>3</b>
À l'attention des lecteurs . . . . .	3
Contexte d'apparition . . . . .	4
Structure générale d'un compilateur . . . . .	5
<b>1 Rappel sur les langages formels</b>	<b>10</b>
1.1 Fondements . . . . .	10
1.2 Alphabet, mots et langages . . . . .	10
1.3 Grammaires . . . . .	12
1.4 Automates finis . . . . .	14
1.4.1 Automates finis déterministes . . . . .	16
1.4.2 Automates finis non déterministes . . . . .	18
1.4.3 Automates finis non déterministes avec $\epsilon$ -transitions . . . . .	19
1.5 Expressions régulières . . . . .	21
1.6 Langages réguliers et propriétés de fermeture . . . . .	23
Exercices . . . . .	25
<b>2 Analyse lexicale</b>	<b>28</b>
2.1 Rôle de l'analyseur lexical . . . . .	28
2.2 Table des symboles . . . . .	30
2.3 Spécification lexicale . . . . .	30
2.3.1 Classes des unités lexicales . . . . .	30
2.3.2 Expressions régulières notation étendue . . . . .	31
2.4 Reconnaissance des unités lexicales . . . . .	34
2.4.1 Quelques problèmes dans la reconnaissance des tokens . . . . .	34
2.4.2 Approche de reconnaissance . . . . .	35
2.4.3 Conventions . . . . .	35
2.5 Implémentation d'analyseurs lexicaux . . . . .	36

2.5.1	Analyseur lexical en dur . . . . .	36
2.5.2	Implémentation par des automates finis . . . . .	37
2.5.3	Générateurs d'analyseurs lexicaux cas de Flex . . . . .	42
Exercices	. . . . .	45
<b>3</b>	<b>Analyse syntaxique</b>	<b>51</b>
3.1	Rôle de l'analyseur syntaxique . . . . .	51
3.2	Spécification syntaxique . . . . .	51
3.2.1	Grammaires hors contexte et dérivation vs réduction . . . . .	53
3.2.2	Qualités et formes particulières des grammaires . . . . .	55
3.3	Analyse syntaxique descendante . . . . .	59
3.3.1	Analyse par la descente récursive . . . . .	60
3.3.2	Analyse LL( $k$ ) . . . . .	63
3.4	Analyse syntaxique ascendante . . . . .	71
3.4.1	Analyse LR(0) . . . . .	72
3.4.2	Analyse SLR . . . . .	77
3.4.3	Analyse LR(1) canonique . . . . .	79
3.4.4	Analyse LALR(1) . . . . .	84
3.5	Générateurs d'analyseurs syntaxiques cas de Bison . . . . .	84
3.5.1	Partie déclaration C . . . . .	85
3.5.2	Partie déclaration Bison . . . . .	85
3.5.3	Partie règles de la grammaire . . . . .	86
3.5.4	Partie code additionnel . . . . .	86
3.6	Gestion des erreurs . . . . .	87
3.6.1	Récupération en mode panique . . . . .	89
3.6.2	Récupération au niveau du syntagme . . . . .	89
3.6.3	Production d'erreurs . . . . .	90
Exercices	. . . . .	91
	<b>Examens corrigés</b>	<b>101</b>

# Liste des figures

1	Code hexadécimal du programme d'Eclid calculant le PGCD de deux entiers [27] . . . . .	4
2	Portion de code en C++ et l'équivalent en assembleur . . . . .	5
3	Aire d'un triangle abc en Assembleur et Java [8] . . . . .	5
4	Architecture générale d'un compilateur . . . . .	7
5	Programme source du PGCD en C [27] . . . . .	8
6	Quelques unités lexicales du code source PGCD . . . . .	8
7	Arbre syntaxique d'une portion du code PGCD . . . . .	9
1.1	Illustration de la hiérarchie de Chomsky . . . . .	15
1.2	Exemple d'un AFD . . . . .	17
1.3	Exemple d'un AFN . . . . .	18
1.4	AFD obtenu par la construction des sous-ensemble appliquée à l'AFN précédent . . . . .	20
1.5	Exemple d'un $\varepsilon$ -AFN . . . . .	20
1.6	Un AFD qui reconnaît le langage dénoté par l'expression $(a + b)^*a$ .	23
2.1	Interaction entre l'analyses lexicale et syntaxique . . . . .	28
2.2	Portion d'un code source fourni à l'analyseur lexical . . . . .	29
2.3	Le code précédent vu par l'analyseur lexical comme un flot de caractères . . . . .	29
2.4	Suite de tokens émis par l'analyseur lexical . . . . .	30
2.5	Construction de Thompson pour les cas de base : $\phi$ , $\varepsilon$ et $a$ . . . . .	39
2.6	Construction de Thompson pour les cas d'induction : $E F$ , $EF$ et $E^*$ .	39
2.7	Construction de Thompson pour l'expression $(a b)^*c$ . . . . .	40
2.8	Automate de $(a   b)^*bab$ par la méthode des dérivés. . . . .	41
2.9	Fonctionnement de Flex . . . . .	42
2.10	Un exemple complet de code Flex . . . . .	44
3.1	Interaction entre les différents modules de la partie analyse d'un compilateur . . . . .	52
3.2	Deux arbres syntaxiques pour dériver la chaîne : $w = a + b * c$ . . . .	54

---

3.3	Deux arbres syntaxiques de la même chaîne dans la grammaire des instructions conditionnelles . . . . .	56
3.4	Code de l'analyseur par descente récursive de $G_1$ . . . . .	61
3.5	Code de l'analyseur par descente récursive de $G_2$ . . . . .	62
3.6	Fonctionnement de l'analyse LL . . . . .	64
3.7	Automate LR(0) de $G_{13}$ . . . . .	75
3.8	Automate LR(0) de $G_{14}$ . . . . .	76
3.9	Automate LR(0) de $G_{15}$ . . . . .	78
3.10	Automate LR(0) de $G_{16}$ . . . . .	79
3.11	Automate LR(1) de $G_{17}$ . . . . .	82
3.12	Structure d'un fichier de spécification Bison . . . . .	85
3.13	Code Bison complet . . . . .	88

## Liste des tables

2.1	Expressions régulières étendues communes . . . . .	33
3.1	Exemple d'analyse par la descente récursive . . . . .	61
3.2	Exemple non concluant d'une analyse par descente récursive d'une chaîne appartenant au langage de la grammaire . . . . .	63
3.3	Table LL de $G_7$ . . . . .	69
3.4	Analyse LL de $i*(i+i)$ . . . . .	70
3.5	Table LL de $G_{12}$ . . . . .	71
3.6	Table LR(0) de $G_{14}$ . . . . .	77
3.7	Table SLR de $G_{15}$ . . . . .	79
3.8	Table LR(1) de $G_{17}$ . . . . .	83
3.9	Déroulement de l'analyse LR(1) de $G_{17}$ sur la chaîne <i>abb</i> . . . . .	83
3.10	Table LALR(1) de $G_{17}$ . . . . .	84

## Liste des algorithmes

1	Pseudo code d'un mini-analyseur lexical codé en dur . . . . .	37
2	Implémentation d'un automate par des fonctions . . . . .	37
3	Implémentation d'un automate fini via sa table de transitions . . . . .	38
4	Structure générale d'un code Flex . . . . .	43
5	Procédure Non-terminal N() . . . . .	60
6	Calcul de l'ensemble DEB(X) . . . . .	65
7	Calcul de l'ensemble SUIV . . . . .	66
8	Construction de la table LL . . . . .	69
9	Fermeture d'un ensemble d'items LR(0) . . . . .	74
10	Successeur d'un ensemble d'items LR(0) . . . . .	74
11	Construction de l'automate LR(0) . . . . .	75
12	Construction de la table d'analyse SLR . . . . .	78
13	Fermeture d'un ensemble d'items LR(1) . . . . .	80
14	Successeur d'un ensemble d'items LR(1) . . . . .	80

15	Calcul de la collection des items LR(1) . . . . .	81
16	Construction de la Table LR(1) . . . . .	81

# Introduction

## À l'attention des lecteurs

Ce document se veut une synthèse du cours sur les compilateurs que je dispense à l'université de Ghardaïa. L'écriture des compilateurs est à la fois un sujet fascinant, fastidieux et complexe. En effet, la conception/réalisation de compilateurs est classée dans le top topics les plus inextricables en informatique suivant immédiatement celui de l'écriture des systèmes d'exploitation. De ce fait, l'étudiant doit être pourvu d'un ensemble complet et varié d'outils qui va du pur fondamental (théorie des langages, théories des graphes, des flots de données, optimisation, etc.), au plus palpable (structures de données, processeurs, jeux d'instructions, caches, etc.). Ainsi pour maîtriser cette complexité élevée, ce cours est de coutume scindé en deux grandes parties. La première concerne le travail antérieur accompli par un compilateur dit aussi la partie analyse ou le front-end, qui comprend les phases d'analyse lexicale, d'analyse syntaxique et l'analyse sémantique. La deuxième partie, quant à elle s'occupe de la seconde moitié des tâches d'un compilateur, baptisée la partie synthèse ou le back-end, qui prend en charge : l'optimisation indépendante à la machine, la génération de code et en fin son optimisation au égard de la machine cible.

Dans cette optique, ce rapport aborde la front-end d'un compilateur, ce qui justifie son titre : « Analyses frontales dans les compilateurs ». Il est destiné aux étudiants de la troisième années LMD de la mention Mathématiques et Informatique. La deuxième partie (le back-end) est présentement enseignée au niveau première année Master où un fascicule à part est en cours d'élaboration.

La documentation sur la construction des compilateurs est abondante. Néanmoins, ce manuscrit puise de quelques références de base tel que le fameux dragon book [1] de l'équipe du professeur Jeff Ullman, des classiques livres de J. P. Tremblay [30], de M.L Scott [27] et de Cooper & Torcson [7]. J'ai aussi beaucoup inspiré et exploité le cours en ligne du professeur Alex Aiken de l'université de Stanford [2].

Le lecteur trouvera également dans ce document une liste intéressante d'exercices et de projets de programmation par chapitre dont l'accomplissement et l'élaboration confirme le niveau de maîtrise des points abordés. Dans la liste bibliographiques, j'ai veillé à citer les "seminal" papiers sur les travaux fondateurs de l'ana-

lyse des langages de programmation, et les liens vers les outils les plus connus dans ce domaine tel que Flex, Bison, CUP, JavaCC, etc.

L'imperfection est humaine et ce travail ne fait pas l'exception cela d'une part. D'autre part, je dois avouer que la compilation de ce rapport fait intervenir plusieurs compétences que parfois m'échappent d'une manière ou d'une autre. Je ne puis nier, vu la charge qui m'incombe, la précipitation et le stress qui a accompagné la préparation de ce travail. Par conséquent, je souhaiterais aimablement inviter les lecteurs de ce document de me signaler toute erreur quelque soit sa nature (orthographe, typographie, méthodologie-style, fond, etc.)

## Contexte d'apparition

La révolution industrielle de la moitié du 20<sup>e</sup> siècle a fait naître les premiers instruments et calculateurs triviaux. Ces machines sont caractérisées par leurs coût et taille souvent prohibitifs mais aussi par des capacités de mémorisation et de calcul simples. Un autre aspect qui compliqua leur exploitation est leur mode d'exploitation qui reste soit manuel ou repose sur des langages machines (suite de bits) très difficile à mener et à maintenir. La figure suivante montre un code machine allégé et codé en hexadécimal pour un Pentium x86.

```
1 55 89 e5 53 83 ec 04 83 e4 f0 e8 31 00 00 00 89 c3 e8 2a 00
2 00 00 39 c3 74 10 8d b6 00 00 00 00 39 c3 7e 13 29 c3 39 c3
3 75 f6 89 1c 24 e8 6e 00 00 00 8b 5d fc c9 c3 29 d8 eb eb 90
```

FIGURE 1 – Code hexadécimal du programme d'Eclid calculant le PGCD de deux entiers [27]

Les gens ont rapidement constaté que les programmes en codes machines sont difficiles à lire, écrire et maintenir car ils sont sujet à l'erreur. Le besoin de l'introduction d'une notation plus gérable est devenu une nécessité. Ainsi les langages d'assemblage sont créés, qui permettent de surpasser le code machine vers des codes mnémoniques plus simples.

Travailler en assembleur garde toujours le programmeur dépendant des détails de bas niveau de la machine cible (jeux d'instructions, de registres, etc.), où le programme utilise un nombre important d'instructions primitives. Les progrès constants ont permis l'introduction d'instructions familières à l'utilisateur par l'invention des

```

1  /* Code en assembleur de X = (Y+4)*3;*/
2  MOV EAX, Y
3  ADD EAX, 4
4  MOV EBX, 3
5  IMUL EBX
6  MOV X, EAX

```

FIGURE 2 – Portion de code en C++ et l'équivalent en assembleur

premiers langages de programmation de haut niveau comme : Fortran (destiné aux calculs scientifiques), Cobol (dédié à la gestion) et Lisp (pour le calcul symbolique). Plus tard, des langages indépendants de la structure de la machine offrant des abstractions de haut niveau fut introduits tel que : C++, Java et Python.

```

1  /* Aire d'un triangle abc en Assembleur*/
2  LOAD R1 a ; ADD R1 b ; ADD R1 c ; DIV R1 #2 ; LOAD R2 R1 ;
3  LOAD R3 R1 ; SUB R3 a ; MULT R2 R3 ;
4  LOAD R3 R1 ; SUB R3 b ; MULT R2 R3 ;
5  LOAD R3 R1 ; SUB R3 c ; MULT R2 R3 ;
6  LOAD R0 R2 ; CALL sqrt
7
8  /* Aire d'un triangle abc en Java */
9
10 public double aire_triangle (double a, double b, double c)
11 {
12 double s = (a+b+c)/2 ;
13 return Math.sqrt (s*(s-a)*(s-b)*(s-c)) ;
14 }

```

FIGURE 3 – Aire d'un triangle abc en Assembleur et Java [8]

Si l'introduction de langages de haut niveau a simplifié la tâche au programmeur, les machines cibles demeurent uniquement communicables via leur propres langages (machines). Cette situation a donc parallèlement augmenté le fossé entre les deux niveaux de langages (évolués et machine), ce qui a compliqué la fonction des programmes chargés d'assurer le passage entre les deux. Ces derniers sont appelés compilateurs.

## Structure générale d'un compilateur

Avant de décrire la structure globale d'un compilateur, différencions d'abord des programmes proches de ceux-ci : Interpréteur, Traducteur, Éditeur et Environnement

de développement intégré [1, 8, 27].

**Définition 1** (Éditeur). *Une application qui permet la saisie, la modification et la sauvegarde de textes de programmes. Ex. : Edit, Bloc Notes, vi, emacs...*

**Définition 2** (Interpréteur). *Traduit et exécute instruction par instruction un programme écrit en un langage évolué : Python, PHP, shell d'Unix, Perl...*

**Définition 3** (Traducteur). *Traduit un programme écrit en un langage évolué L1 dit langage source en un programme équivalent en un autre langage évolué L2 dit langage cible.*

**Définition 4** (Compilateur). *Traduit un programme écrit en langage évolué en un code exécutable équivalent en code machine. Ex : C, Pascal...*

**Définition 5** (Environnement de développement intégré (EDI)). *Un programme complet qui intègre éditeur, compilateur avec d'autres outils de mise au point. Ex. : JBuilder, Delphi, Eclipse, IDLE, IntelliJ, etc.*

La complexité du travail d'un compilateur a suscité sa décomposition en plusieurs phases avec des interfaces bien définies comme le montre la Figure 4. Chaque phase est à son tour décortiquée en tâches plus spécifiques. Nous présentons brièvement dans cette section un aperçu général des différentes phases. Considérons le code de calcul du PGCD de deux entiers montré dans la Figure 5, et expliquons les différents traitements qu'il (ou ses transformés) va subir durant ce long processus.

## Analyse lexicale

La mission de l'analyse lexicale est de regrouper les caractères du texte source du programme en unités ayant une signification du point de vue spécification lexicale du langage. Ainsi la succession de i,n,t, et espace par exemple est perçue comme une unité lexicale mots-clé "int". De cette manière, l'analyse lexicale transforme le code source d'une suite de caractères en une suite d'unités lexicales qui constitue l'entrée de la phase suivante. Il est du ressort de cette phase aussi l'insertion et le codage des unités lexicales reconnues dans une structure centrale appelée la table de symboles, la suppression des espaces, blancs et commentaires et en fin de signaler les erreurs d'ordre lexical détectées (la mauvaise formation des unités lexicales).

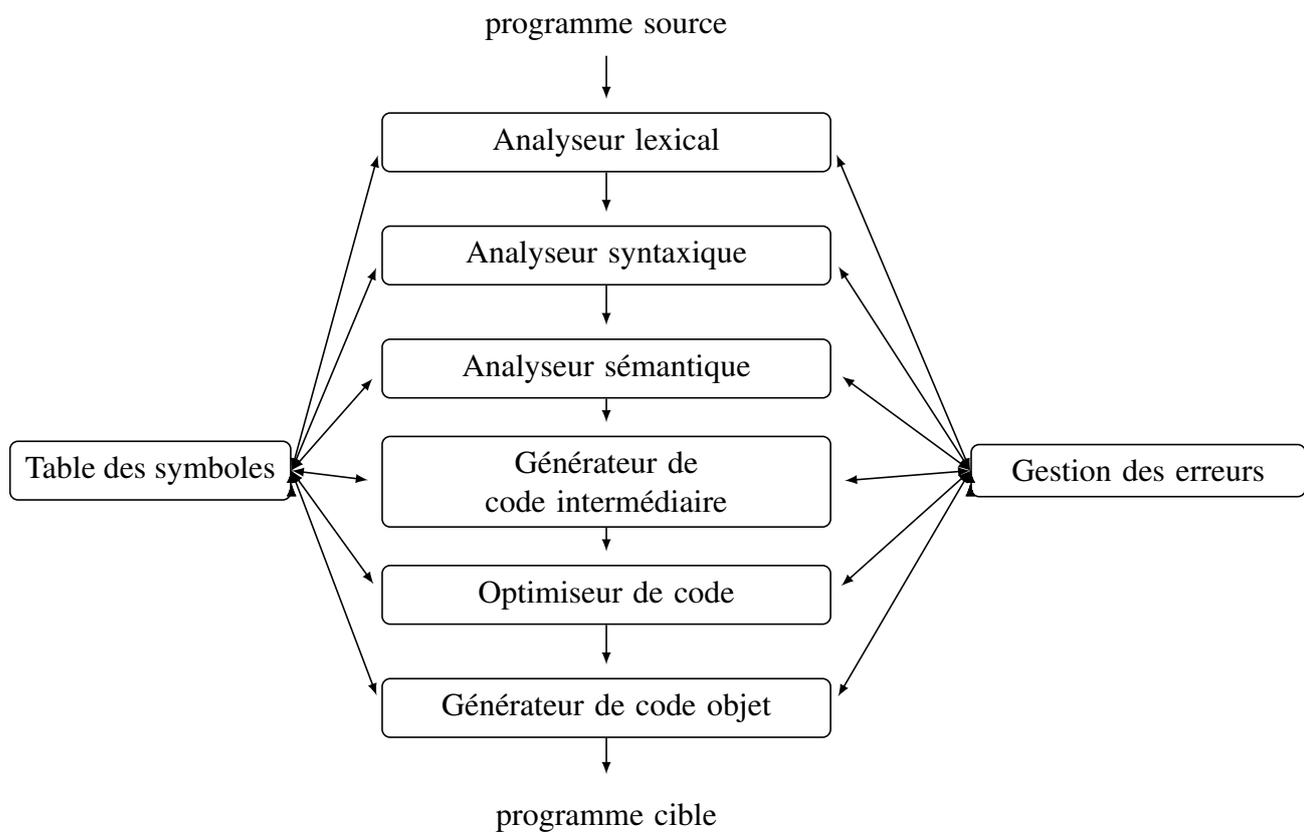


FIGURE 4 – Architecture générale d'un compilateur

```
1 int main() {
2     int i = getint(), j = getint();
3     while (i != j) {
4         if (i > j) i = i - j;
5         else j = j - i;
6     }
7     putint(i);
8 }
```

FIGURE 5 – Programme source du PGCD en C [27]

```
1 int main int ( ) ...while...putint...}
```

FIGURE 6 – Quelques unités lexicales du code source PGCD

## Analyse syntaxique

À la réception des unités lexicales de la phase précédente, le rôle de l'analyse syntaxique est de vérifier la bonne disposition des unités reçues pour former des instructions valides du langage c-à-d conformes à sa définition syntaxique donnée par une grammaire (collection de règles du genre  $A \rightarrow \alpha$ ). Pour justifier la validité syntaxique du code, cette phase crée un arbre dit syntaxique dont les nœuds internes sont les opérations et les fils associés les arguments. En cas d'incohérence des erreurs d'ordre syntaxiques seront émises. La figure ci-dessous montre l'arbre de l'instruction **while** du code précédent.

## Analyse sémantique

Ici l'objectif est de vérifier des aspects ayant trait à la signification du code. Cette phase utilise l'arbre syntaxique fournit par la précédente et la table de symboles pour l'annoter avec les attributs consignées pour effectuer des vérifications liées au contexte comme : les déclarations, la compatibilité des types, le respect des portés, du nombre/types des arguments dans les fonctions, etc. Souvent ces aspects sont complexes à traiter avec des grammaires hors contextes comme dans la phase précédente. On fait recours à des techniques (grammaires attribuées) permettant d'opérer des actions sémantiques par des routines à invoquer dans des endroits/moments appropriés. L'analyse sémantique transforme l'arbre syntaxique en arbre abrégé et annoté ou en une forme intermédiaire pour une machine abstraite et la transmet à la suite du processus de compilation. Les erreurs d'ordre sémantiques sont signalées

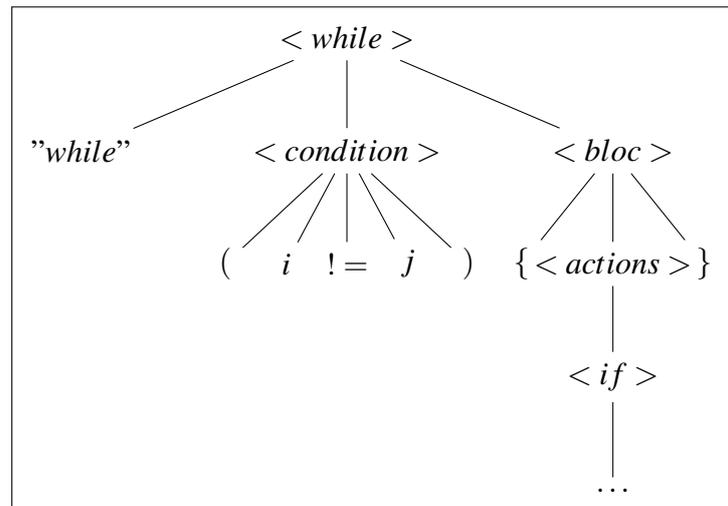


FIGURE 7 – Arbre syntaxique d'une portion du code PGCD

dans cette phase moyennant le module de gestion des erreurs.

## Optimisation

Dans la première opération d'optimisation (indépendante de la machine), le code intermédiaire subi plusieurs transformations, qui préservent son logique, permettant d'améliorer son efficacité. Le gain obtenu concerne divers aspects tel que : le temps d'exécution du programme, sa taille, voire même les communications et l'énergie qu'il dispense. Après la génération du code machine, ce dernier fera l'objet d'une deuxième phase d'optimisation mais cette fois dépendant de l'architecture cible.

## Génération de code

Cette phase produit le code objet en prenant en compte l'architecture et le jeu d'instruction du processeur cible. C'est une traduction du code intermédiaire optimisé vers une séquence d'instructions machine qui réalisent le même travail. Des questions à résoudre sont par exemple : attribuer des adresses (ou un mécanisme d'adressage) pour les objets manipulés, allouer judicieusement les registres de la machine souvent en nombre limité, bien sélectionner/ordonner les instructions, etc.

# Rappel sur les langages formels

---

## 1.1 Fondements

La théorie des langages formels se fonde sur plusieurs autres théories mathématiques. On fait usage ici notamment de la théorie des ensembles, de la logique mathématique et de quelques aspects de l’algèbre et ses structures élémentaires. Des notions de la théorie des graphes sont souvent exploitées également dans ce cours. Le lecteur avide est invité à consulter les diapositifs du cours dispensé en présentiel en deuxième année licence informatique ou d’exploiter les références bibliographiques sur chaque volet sus-mentionné. Le livre [17] est un bon survol de ces concepts.

Nous rappelons succinctement dans ce chapitre les bases afférentes directement aux langages formels. Plus particulièrement, nous nous focalisons sur les langages réguliers. Pour plus d’exemples, de démonstrations et d’exercices veuillez consulter le livre du groupe du professeur Jeff Ullman [12]

## 1.2 Alphabet, mots et langages

On définit un **alphabet** comme étant un ensemble fini non vide de symboles (lettres ou caractères) qu’on note en utilisant souvent une lettre grecque majuscule comme  $\Sigma$ .

Sur un alphabet  $\Sigma$  on appelle un **mot** (chaîne ou string)  $w$  toute suite finie et ordonnée (on dit aussi séquence) de symboles de  $\Sigma$ . Le mot **miroir** ou transposé d’un mot  $w$  est noté  $w^R$  ou  $\tilde{w}$  est obtenu en inversant tous les symboles de  $w$ .

$|w|$  dénote la longueur d’un mot  $w$ . il représente le nombre de symboles qui constituent le mot  $w$ . Le mot ayant une longueur nulle est dit le mot vide qu’on note  $\varepsilon$ . La longueur indiquée par un symbole est le nombre d’occurrences de ce dernier dans le mot considéré. Par exemple  $|2020|_0 = 2$  car 0 apparaît deux fois dans la chaîne 2020.

$w, x, y$  étant trois mots sur un alphabet  $\Sigma$ . On dit que :

- $u$  est **préfixe** de  $w$  s'il existe  $x \mid w = ux$
- $u$  est **suffixe** de  $w$  s'il existe  $x \mid w = xu$
- $u$  est **facteur** de  $w$  s'il existe  $x$  et  $y \mid w = xuy$

Un préfixe, (resp. suffixe ou facteur) est dit **propre** lorsque il est différent du mot vide  $\varepsilon$  et du mot  $w$  lui même

Pour un alphabet  $\Sigma$  l'ensemble  $\Sigma^k$  inclut tous les mots sur  $\Sigma$  de longueur  $k$ . À titre d'exemple sur l'alphabet  $\Sigma = \{a, b\}$ ,  $\Sigma^2 = \{aa, ab, ba, bb\}$ . Par convention, pour tout alphabet  $\Sigma$ , nous admettons que  $\Sigma^0 = \{\varepsilon\}$ .

L'ensemble de tous les mots sur un alphabet  $\Sigma$  sera noté  $\Sigma^*$ . Ainsi, nous avons :

$$\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots = \bigcup_{k \geq 0} \Sigma^k$$

Sur  $\Sigma^*$  est défini une opération dite de **concaténation** entre des mots. si  $u$  et  $v$  sont deux mots de  $\Sigma$  la concaténation de  $u$  et de  $v$  donne le mot  $uv$  formé des symboles de  $u$  suivis de ceux de  $v$ . Cette opération est associative et admet  $\varepsilon$  comme élément neutre. Notez bien que cette opération n'est pas en général commutative.

Un **langage** sur un alphabet  $\Sigma$  est tout ensemble de mots sur cet alphabet. Ceci dit, un langage  $L$  sur  $\Sigma$  est un sous ensemble de  $\Sigma^*$ . Voici quelques remarques à bien avaler sur les langages :

1. Un langage étant un ensemble, il peut être fini ou infini
2. Il n'est pas exigé que les mots d'un langage sur un alphabet inclut tous les symboles de cet alphabet
3. Tout langage d'un alphabet est aussi un langage sur son super-alphabet
4.  $\emptyset$  le langage vide, et le langage formé uniquement du mot vide  $\varepsilon$  sont deux langages sur tout alphabet
5.  $\Sigma^*$  est un langage sur tout alphabet  $\Sigma$

Les langages étant des ensembles, les opérations ensemblistes usuelles (union, intersection, complément, etc.) y sont définies de la même manière. Nous insistons ici en particulier sur les opérations de concaténation et d'itération.

La **concaténation** de deux langages  $M$  et  $N$  est le langage formé des mots obtenus par concaténation des mots des deux langages en respectant l'ordre de concaté-

nation. Cette opération est associative mais, comme pour les mots, n'est pas commutative. Elle admet le langage vide comme absorbant et le langage formé seulement du mot vide comme élément neutre.

**L'itération** ou l'étoile d'un langage offre un moyen de formé des mots du même langage par plusieurs succession de concaténation. En général on note :

$$\begin{cases} L^0 = \{\varepsilon\} \\ L^n = L.L^{n-1} \end{cases}$$

Il existe plusieurs propriétés et résultats sur les opérations sur les langage, consultez la bibliographie pour plus de détail. Nous donnons ci-après des propriétés remarquables de l'opération d'itération.

$$\left\{ \begin{array}{l} L^* = \bigcup_{i \geq 0} L^i \\ L^* = L^{**} \\ L^* = L^*L^* \\ L(ML)^* = (LM)^*L \\ (L \cup M)^* = (M^*L)^*M^* \\ \emptyset^* = \{\varepsilon\} \\ \{\varepsilon\}^* = \{\varepsilon\} \\ L^+ = LL^* = L^*L \end{array} \right.$$

### 1.3 Grammaires

Les grammaires sont un outil permettant de générer des mots en offrant un système de réécriture sous la forme de règles dites de production. Elles permettent ainsi de définir des langages. Une grammaire est spécifiée comme un quadruplet  $G = \langle T, N, S, P \rangle$  constituée d'un ensemble  $T$  de symboles **terminaux** entrant dans la composition des mots appelé aussi alphabet, d'un ensemble  $N$  de symboles **non terminaux** également désignés de variables servant comme intermédiaires dans le processus de génération, d'un symbole spécial  $S$  élément de  $N$  appelé l'**axiome** qui joue le rôle d'initiateur et en fin d'un ensemble  $P$  de **règles de réécriture**. Il est à noter que les ensembles  $T$ ,  $N$  et  $P$  sont non vides et que les deux premiers sont évidemment disjoints.

Les éléments de  $P$  sont des couples  $(\alpha, \beta)$  tel que  $\alpha \in (T \cup N)^* N (T \cup N)^*$  et  $\beta \in (T \cup N)^*$  qu'on note en insérant une flèche entre les deux :  $\alpha \rightarrow \beta$  pour signifier  $\alpha$  peut être remplacé par  $\beta$ . La première partie est aussi appelé membre gauche de production et alors que la deuxième est qualifiée de membre droit de production.

On appelle une **dérivation** l'opération de remplacement d'un membre gauche par un de ses membres droits selon les règles de production de la grammaire.

Si  $A \rightarrow \gamma$  est une règle de  $P$  : alors  $\alpha A \beta \Rightarrow \alpha \gamma \beta$  est une dérivation.

Si un mot  $g$  dérive de  $f$  en une seule étape ( $f \Rightarrow g$ ) on dit que  $g$  dérive directement de  $f$ . Dans le cas contraire (dérive en plusieurs étapes) on dit qu'il dérive indirectement de  $f$  et on écrit ( $f \stackrel{\pm}{\Rightarrow} g$ ). Par ailleurs, un mot  $f$  de  $T^*$  est généré par une grammaire  $G$  s'il peut être dérivé à partir son axiome  $S$ . autrement dit :

$$G \text{ génère } f \text{ ssi } S \stackrel{\pm}{\Rightarrow} f$$

Ceci nous conduit au concept du langage généré par une grammaire qui est donc l'ensemble de mots générés par cette grammaire. On écrit :

$$L(G) = \{w \in T^* \mid S \stackrel{\pm}{\Rightarrow} w\}$$

Pour illustrer le processus de production de mots, les étapes sont schématisées par un arbre dit **arbre de dérivation**. Cet arbre possède l'axiome comme une racine et les éléments de  $N$  comme nœuds internes. Les feuilles quand à elles sont représentées par les terminaux de  $T$ . Les branches de l'arbre de dérivation schématisent les membres droits de production. à la fin d'un processus de génération la concaténation des symboles des feuilles de gauche à droite représente le mot produit associé à cette dérivation.

Une grammaire est **ambiguë** s'il existe un mot dans son langage ayant deux arbres de dérivation différents. Dans la cas contraire, elle qualifiée de **non ambiguë**.

Durant son étude sur les langage et leur structures syntaxiques **Noam Chomsky** aboutira à une classification des langages et des grammaires associées selon le degré de complexité connue sous le terme **hiérarchie de Chomsky** [6]. Cette hiérarchie inclut quatre classes de langages, les voici :

1. **Langages réguliers**, rationnels ou linéaire à droite ou encore de type 3. Ce

sont les grammaires où toutes les productions sont de la forme :

$A \rightarrow \alpha B$  ou  $A \rightarrow \alpha$  avec :  $A$  et  $B$  des éléments de  $N$  et  $\alpha$  une chaîne terminale de  $T^*$

2. **Langages à contexte libre**, non contextuels ou algébriques ou encore de type 2. Si toutes les règles sont de la forme :  $A \rightarrow \alpha$  avec  $A \in N$  et  $\alpha \in (T \cup N)^*$
3. **Langages à contexte lié** ou contextuels si toutes les règles de  $P$  la forme :  $\alpha A \beta \rightarrow \alpha w \beta$  où :  $A \in N$ ,  $\alpha, \beta \in (T \cup N)^*$  et  $w \in (T \cup N)^+$
4. **Langages récursivement énumérables** ou sans restriction dans tous les autres cas.

Il est facile de constater, comme le montre la figure 1.1, la relation d'inclusion entre les classes de grammaires (et les langages correspondants). Les plus larges sont celles de type 0 et les plus restrictives sont les grammaires (langages) de type 3.

## 1.4 Automates finis

Un automate fini (Finite Automaton, or Finite State Machine) est une machine abstraite qui permet, en un temps fini, d'accepter ou de rejeter des mots fournis à son entrée. Sa réponse est alors binaire, *i.e.*, il retourne pour chaque chaîne soumise à lui oui (accepte) ou non (rejette). Cette décision est formulée suivant une fonction de transitions qui joue le rôle de son moteur. Cette fonction est définie souvent par une table qui associe les transitions possibles des différents états selon le symbole du mot actuellement présent devant une sorte de tête de lecture. Le processus de reconnaissance démarre à partir d'un état spécial dit initial, et progresse en respectant la dynamique définie par la table de transitions pour signaler une acceptation du mot si la machine s'arrête sur un état dit final ou signaler un rejet dans le cas contraire. Notez qu'un rejet est aussi émis si l'automate se bloque durant le processus de reconnaissance (absence de décision pour une configuration donnée matérialisée par l'inexistence d'une transition par l'état et le symbole actuels). Il existe différents sorts d'automates selon les modalités et les commodités qu'ils offrent.

Les automates sont un outil formel très puissant qui trouve des applications variées dans plusieurs domaines :

- La recherche de motifs : qui exploite l'essence même des automates à savoir reconnaître un motif au sens large.

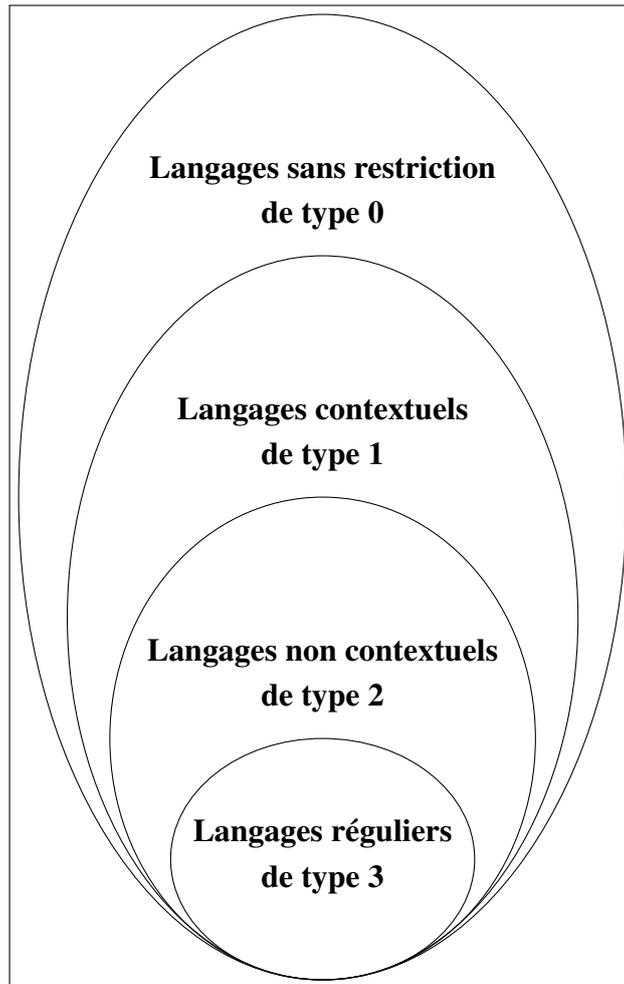


FIGURE 1.1 – Illustration de la hiérarchie de Chomsky

- Forme une brique de base qui aide dans la synthèse des circuits numériques
- Un outil fondamental de vérification pour plusieurs domaines en particuliers les protocoles de communication
- Différents types d'automates sont utilisés dans la réalisation des compilateurs (analyseurs lexicaux)

Ci-après nous présentons de façon formelle les plus communs des automates finis.

### 1.4.1 Automates finis déterministes

Un automate fini déterministe **AFD** est un quintuplet  $A = \langle Q, \Sigma, \delta, q_0, F \rangle$  où :

- $Q$  est un ensemble fini d'états.
- $\Sigma$  est un ensemble fini de symboles dit l'alphabet d'entrée.
- $\delta : Q \times \Sigma \rightarrow Q$  est la fonction de transitions.  $\delta(p, a) = q$  signifie que de l'état  $p$  on passe à l'état  $q$  à la lecture du symbole  $a$
- $q_0 \in Q$  est l'état initial
- $F \subseteq Q$  est l'ensemble des états finaux

Un automate est qualifié par **déterministe** (AFD) s'il existe au plus une transition par état et par lettre de l'alphabet. Dans le cas contraire, il est dit non déterministe.

Notons qu'un automate est représenté soit en donnant sa table de transitions qui élucide le lien entre états et symboles de l'alphabet ou sous la forme d'un graphe qui montre explicitement cette fonction de transition.

**Exemple 1.** Soit l'automate défini par la donnée des composants suivants :

- $Q = \{q_0, q_1, q_2\}$
- $\Sigma = \{a, b\}$
- $\delta = \{(q_0, a, q_0), (q_0, b, q_1), (q_1, a, q_0), (q_1, b, q_2), (q_2, a, q_2), (q_2, b, q_2)\}$
- $q_0$  est l'état initial
- $F = \{q_1\}$

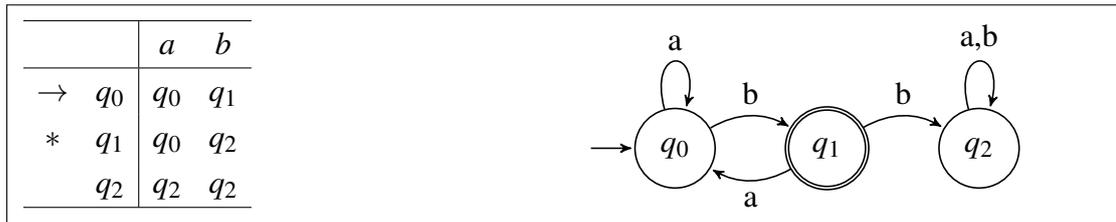


FIGURE 1.2 – Exemple d'un AFD

### Configuration d'un automate

Le couple  $(q, w)$  représente une configuration de l'automate. où  $q$  est un état de l'automate et  $w$  un mot de  $\Sigma^*$ . Ainsi  $(q_0, w)$  est dit configuration initiale et  $(q_f, \varepsilon)$  une configuration finale avec  $q_f \in F$

Un mot  $w$  est accepté par un automate si on peut passer de la configuration initiale de mot  $w$  et on termine vers une configuration finale avec le mot vide. C'est à dire s'il existe un chemin qui part de l'état initial et termine dans un des états finaux ou les étiquettes de ses transitions sont les symboles qui composent le mot dans l'ordre. Plus précisément, un automate accepte  $w = a_1 a_2 \dots a_n$  s'il existe une séquence d'états  $s_0, s_1, \dots, s_n$  tel que :

- $s_0 = q_0$
- $\delta(s_i, a_{i+1}) = s_{i+1}$
- $s_n \in F$

Dans l'exemple précédent, l'automate accepte bien le mot *abab* mais rejette *abb*. La fonction de transitions  $\delta$  est aisément étendue au mots de la manière suivante :

$$\hat{\delta} : Q \times \Sigma^* \rightarrow Q$$

$$\begin{cases} \hat{\delta}(q, \varepsilon) = q \\ \hat{\delta}(q, aw) = \hat{\delta}(\delta(q, a), w) \end{cases}$$

L'ensemble des mots acceptés par un automate forme le **langage** accepté par cet automate. Formellement, on écrit :

$$L(A) = \{w \in \Sigma^* \mid \hat{\delta}(q_0, w) \in F\}$$

Un automate est **complet** si pour tout état il existe exactement une transition pour

chaque symbole de l'alphabet  $\Sigma$ . Deux automates sont **équivalents** s'ils acceptent le même langage.

### 1.4.2 Automates finis non déterministes

Un automate fini est qualifié de non déterministe (AFN) si sa fonction de transitions autorise le passage d'un état par un symbole vers un ensemble d'états au lieu d'un état unique comme dans le cas déterministe. Autrement dit, un automate non déterministe peut se retrouver sur plusieurs états au même moment. Cette propriété offre une souplesse dans la conception d'automates modélisant un espace très large de problèmes ; mais n'améliore pas pourtant son pouvoir de reconnaissance.

Un AFN est donc toujours un quintuplet avec les mêmes composants avec une fonction de transition  $\Delta$  définie comme suit :

$$\Delta : Q \times \Sigma \rightarrow 2^Q$$

**Exemple 2.** Soit l'AFN défini ainsi  $A = \langle \{q_0, q_1, q_2\}, \{0, 1\}, \Delta, q_0, F \rangle$  avec la fonction  $\Delta$  et l'ensemble  $F$  schématisés dans la table et la figure ci-dessous :

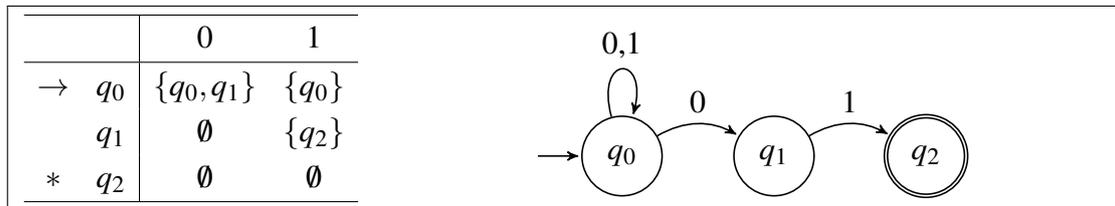


FIGURE 1.3 – Exemple d'un AFN

Le langage d'un AFN est l'ensemble des mots dont l'exécution par l'automate conduit à un ensemble contenant au moins un état final. *i.e* :

$$L(A) = \{w \in \Sigma^* \mid \hat{\Delta}(q_0, w) \cap F \neq \emptyset\}$$

Cette dernière définition n'est correcte qu'en étendant la fonction de transitions aux mots de la manière suivante :

$$\left\{ \begin{array}{l} \hat{\Delta}(q, \varepsilon) = \{q\} \forall q \in Q \text{ cas de base} \\ \hat{\Delta}(q, u) = \{p_1, \dots, p_k\} \text{ Pour } p_i \in Q, a \in \Sigma \text{ et} \\ \bigcup_{i=1}^k \hat{\Delta}(p_i, a) = \{r_1, \dots, r_m\} \text{ alors :} \\ \hat{\Delta}(q, w) = \{r_1, \dots, r_m\} \text{ et } w = ua \end{array} \right.$$

### Équivalence entre AFDs et AFNs

M.O. Rabin & D.S. Scott introduisirent en 1959 [24] un théorème intéressant qui stipule que tout langage accepté par un AFN est aussi accepté par un AFD. Ce résultat donne également un algorithme de passage d'un AFN vers un AFD équivalent dit la construction des sous-ensembles (subset construction).

**Théorème 1.4.1.**  $\mathcal{A}_N$  un AFN : Il existe un AFD  $\mathcal{A}_D$  |  $L(\mathcal{A}_N) = L(\mathcal{A}_D)$

Donc à tout AFN  $A_N = \langle Q_N, \Sigma, \Delta, q_0, F_N \rangle$  correspond un AFD  $A_D = \langle Q_D, \Sigma, \delta, q_1, F_D \rangle$  avec :

- $Q_D = 2^{Q_N}$  c-à-d les parties de  $Q_N$  sans les états inaccessibles (ne possédant pas un chemin de l'état initial)
- $q_1 = \{q_0\}$
- $\delta(S, a) = \bigcup_{p \in S} \Delta(p, a)$  avec  $S \in Q_N$
- $F_D = \{S \in Q_N \mid S \cap F_N \neq \emptyset\}$

**Exemple 3.** L'application de cette construction à l'AFN de l'exemple 2 précédent à la Figure 1.3 donne l'AFD ci-dessous dans la Figure 1.4 :

### 1.4.3 Automates finis non déterministes avec $\varepsilon$ -transitions

Une  $\varepsilon$ -transition dite aussi spontanée ou instantanée est une transition qui ne consomme pas de symbole de l'entrée, elle est effectuée sur le mot vide  $\varepsilon$ . Ce type de transitions est une sorte d'indéterminisme et constitue une commodité qui offre une souplesse dans la conception des automates mais n'améliore pas la puissance de ces derniers. Formellement, la définition ajoute seulement à la fonction de transitions des entrées pour  $\varepsilon$ .

Un  $\varepsilon$ -AFN est un quintuplet incluant les mêmes composants :

$$A_\varepsilon = \langle Q, \Sigma, \Delta_\varepsilon, q_0, F \rangle \text{ avec : } \Delta_\varepsilon : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$$

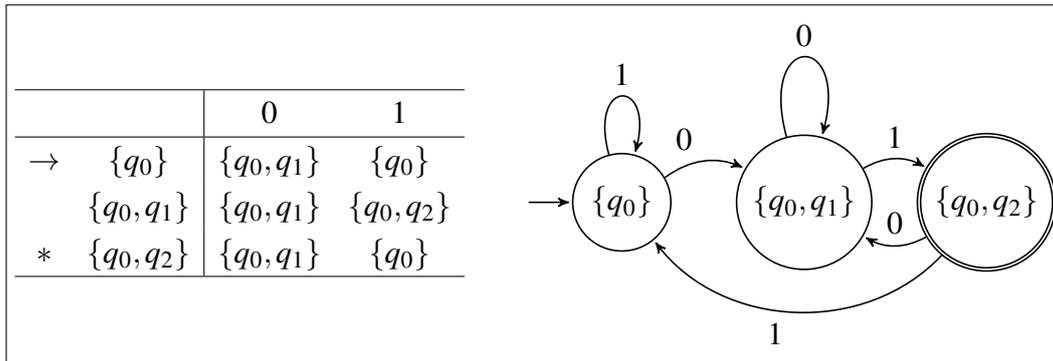


FIGURE 1.4 – AFD obtenu par la construction des sous-ensemble appliquée à l’AFN précédent

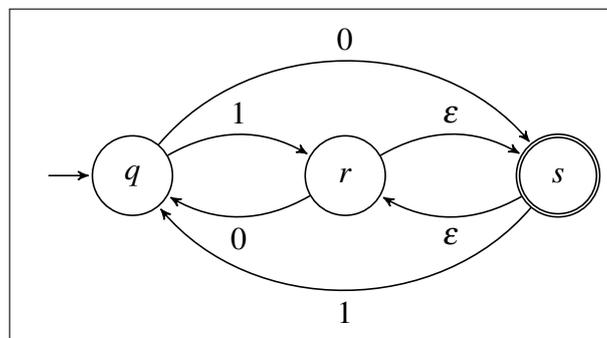


FIGURE 1.5 – Exemple d’un  $\epsilon$ -AFN

L’acceptation d’un mot est toujours l’existence d’un chemin de l’état initial vers un état final. Par exemple, l’ $\epsilon$ -AFN de la Figure 1.5 accepte 001 mais pas 01.

Comme pour les AFNs, à tout  $\epsilon$ -AFN, il existe un AFN sans  $\epsilon$ -transitions équivalent. Cette transformation est basée sur la notion de  $\epsilon$ -fermeture (ou clôture).

**Théorème 1.4.2.**  $\mathcal{A}_\epsilon$  un  $\epsilon$ -AFN : Il existe un AFN sans  $\epsilon$ -transitions  $\mathcal{A}_N \mid L(\mathcal{A}_\epsilon) = L(\mathcal{A}_N)$

### $\epsilon$ -clôture

l’ $\epsilon$ -clôture d’un état  $p$  noté  $C_\epsilon(p)$  est l’ensemble de tous les états atteignables à partir de  $p$  par zéro, un ou plusieurs  $\epsilon$ -transitions. Formellement de façon récursive

la fermeture en  $\varepsilon$  est spécifiée comme suit :

$$\left\{ \begin{array}{l} p \in C_\varepsilon(p) \\ \text{si } q \in C_\varepsilon(p) \ \& \ q' \in \Delta_\varepsilon(q, \varepsilon) \ \text{Alors } q' \in C_\varepsilon(p) \\ C_\varepsilon(E) = \bigcup_{s \in E} C_\varepsilon(s) \end{array} \right.$$

Dans l' $\varepsilon$ -AFN précédent, nous avons  $C_\varepsilon(q) = \{q\}$  ; et  $C_\varepsilon(r) = C_\varepsilon(s) = \{r, s\}$

### Élimination des $\varepsilon$ -transitions

Plusieurs applications requièrent des automates sans transitions vides. Il existent plusieurs algorithmes pour rendre un automates avec  $\varepsilon$ -transitions libre de celles-ci. La procédure la plus directe est similaire à celle de la déterminisation. Il faut seulement prendre en compte les  $\varepsilon$ -transitions par le mécanisme des  $\varepsilon$ -clôtures. Le résultat de cette dernière procédure est un automate fini déterministe.

En général, on peut concevoir un algorithme qui élimine les transitions spontanées et produit un automate non déterministe. Voici la procédure de suppressions des transitions nulles d'un état  $p$  vers un autre état  $q$  :

- Trouver toutes les transitions partantes de  $q$
- Refaire ces transitions (en gardant les mêmes étiquettes) cette fois émanant de  $p$
- Si  $p$  est un état initial ; rendre  $q$  aussi initial
- Si  $q$  est un état final ; rendre  $p$  final aussi.

## 1.5 Expressions régulières

Les expressions régulières ou rationnelles sont un autre moyen pour dénoter des langages. Elles représentent l'équivalent algébrique des automates finis. Voir la section 2.5 pour les techniques de passage entre ces expressions et les automates.

Si  $r$  est une expressions régulière  $L(r)$  est le langage dénoté ou spécifié par l'expression  $r$ . Trois opérateurs dit réguliers sont utilisés pour former ces expressions : l'union (+), la concaténation (.) et l'itération ou l'étoile de Kleene (\*).

Une définition récursive des expressions régulières sur une alphabet  $\Sigma$  est la suivante :

- $\emptyset$  est une expression et  $L(\emptyset) = \emptyset$
- $\varepsilon$  est une expression régulière et  $L(\varepsilon) = \{\varepsilon\}$
- Pour tout  $a \in \Sigma$  :  $a$  est une expression régulière et  $L(a) = \{a\}$
- Si  $r$  et  $s$  sont deux expressions régulières, alors :
  1.  $r + s$  est une expressions régulièrè et  $L(r + s) = L(r) \cup L(s)$
  2.  $r.s$  ou simplement  $rs$  est une expressions régulièrè et  $L(rs) = L(r)L(s)$
  3.  $r^*$  est une expressions régulièrè et  $L(r^*) = (L(r))^*$
  4.  $(r)$  est une expressions régulièrè et  $L((r)) = L(r)$ . cette dernière clause n'est pas un opérateur, mais permet de définir des expressions régulières parenthésées.

Dans les expressions complexes, il faut tenir compte de l'ordre de priorité des opérateurs réguliers. Voici cet ordre du plus fort vers le plus faible :

1.  $()$  : les parenthèses
2.  $*$  : l'étoile de Kleene
3.  $.$  : la concaténation
4.  $+$  : l'union

**Exemple 4.** *Ci-après des exemples d'expressions régulières avec les langages associés :*

$$(0 + 1)1 \quad : \quad \{01, 11\}$$

$$0 + 10^* \quad : \quad \{0, 1, 10, 100, 1000, \dots\}$$

$$(a + b)^* \quad : \quad \{\varepsilon, a, b, aa, ab, ba, bb, aaa, \dots\} \text{ toutes les chaînes sur } \{a, b\}$$

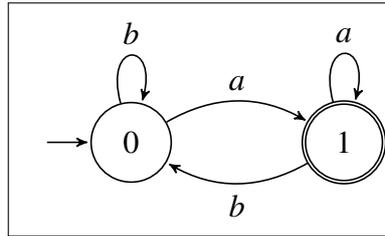
$$(0(0 + 1))^* \quad : \quad \{\varepsilon, 00, 01, 0000, 0001, 0100, 0101, \dots\} \text{ les mots sur } \{0, 1\} \\ \text{de longueur pair dont } 1 \text{ est toujours en position pair}$$

Attention à la priorité sus-indiquée. Par exemple les expressions suivantes ne sont pas équivalentes :

$$ab \quad \neq \quad ba$$

$$ab + c \quad \neq \quad a(b + c)$$

$$ab^* \quad \neq \quad (ab)^*$$

FIGURE 1.6 – Un AFD qui reconnaît le langage dénoté par l’expression  $(a + b)^*a$ 

Pour finir, voici quelques propriétés utiles sur les expressions régulières  $E$ ,  $F$  et  $G$  étant des expressions régulières :

$$\begin{aligned}
 E + F &= F + E \\
 E + E &= E \\
 (E + F)G &= EG + EF \\
 E(F + G) &= EF + EG \\
 E + \emptyset &= \emptyset + E = E \\
 E\emptyset &= \emptyset E = \emptyset \\
 E\varepsilon &= \varepsilon E = E \\
 \emptyset^* &= \varepsilon^* = \varepsilon \\
 E(FE)^* &= (EF)^*E \\
 (E + F)^* &= (E^*F)^*E^* = (F^*E)^*F
 \end{aligned}$$

## 1.6 Langages réguliers et propriétés de fermeture

Un langage  $L$  sur un alphabet  $\Sigma$  est reconnaissable s’il est accepté par un automate fini. Il est qualifié de régulier ou rationnel s’il est dénoté par une expression régulière. Par exemple, le langage sur  $\{a, b\}$  des mots se terminant par un  $a$  est reconnaissable (reconnu par l’automate de la Figure 1.6) et aussi régulier (dénoté par l’expression régulière suivante :  $(a + b)^*a$ ).

S.C. Kleene [14] introduit le théorème d’équivalence entre automates finis et expressions régulières qui porte son nom. Il prouva que la classe des langages reconnaissables est égale à la classe des langages réguliers. La preuve donne une méthode pour passer d’un automate à une expressions régulières équivalente et inversement. Cette preuve introduit pour le premier sens une construction dite de Thompson. Le

deuxième sens est possible via le lemme d'Arden. Afin d'alléger ce support ne nous détaillons pas ces constructions ici. Pour de plus ample information, veuillez vous référer aux slides du cours et les références bibliographiques.

En résumé les modèles présentés dans ce chapitre sont équivalent : tout ce que peut être reconnu par un AFD l'est aussi par : un AFN,  $\epsilon$ -AFN, une expression régulière et évidemment une grammaire régulière.

**Théorème 1.6.1.** *Les automates finis, les expressions régulières et les grammaires linéaires spécifient la même classe de langages (les langages réguliers)*

La classe des langages réguliers jouit de quelques propriétés dites de fermeture. Autrement dit, des opérations sur des éléments de cette classe qui préserve la régularité et produit un langage de la même classe. Nous les citons ici sans démonstration :

**Théorème 1.6.2.** *La famille des langages réguliers est fermée par :*

- *Union*
- *Intersection*
- *Concaténation*
- *Itération*
- *Complémentation*
- *Différence*
- *Transposé*
- *Homomorphisme*

Enfin, notons qu'il existe des langages qui ne sont pas réguliers. Le lemme de pompage est un outil qui valide la régularité d'un langage, veuillez vous référer aux références bibliographiques pour plus de détail [12].

Dans la partie analyse lexicale nous utiliserons les langages réguliers, par contre nous nous recourons aux langages hors contextes dans la partie analyse syntaxique. Les langages contextuels trouvent une application dans l'analyse sémantique.

# **Exercices du chapitre 1**

Ministère de l'Enseignement Supérieur et de la Recherche Scientifique  
 Université de Ghardaia  
 Faculté des Sciences et de la Technologie  
 Département des Mathématiques et d'Informatique

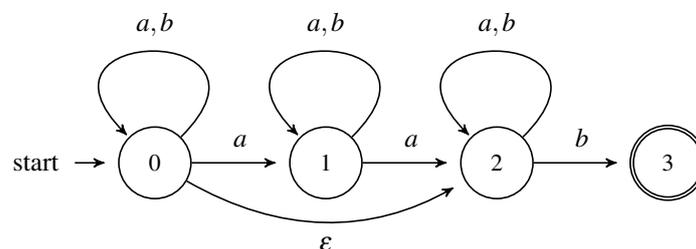
## RAPPEL

### Exercice 1

$L = \{ab, aa, baa\}$  un langage. Lesquels des mots suivants sont dans  $L^*$  :  
 $abaabaaabaa$ ,  $aaaabaaaa$ ,  $baaaaabaaaab$ ,  $baaaaabaa$

### Exercice 2

Soit l'AFN  $\mathcal{A}$  suivant :



- Indiquer tous les chemins étiquetés par  $aabb$
- $\mathcal{A}$  accepte-t-il  $aabb$
- Donner sa table de transition
- Déterminiser  $\mathcal{A}$

### Exercice 3

Construire un automate fini qui accepte le langage  $L$  des mots sur  $\{a, b\}$  ayant exactement deux  $a$  et plus de deux  $b$ . Modifier cet automate de manière à accepter le langage  $\bar{L}$ . Trouver des constructions pour accepter :  $L^2$ ,  $L^2 \setminus L$ .

### Exercice 4

Concevoir un AFN ayant :

- Trois états pour  $\{ab, abc\}^*$
- Au plus cinq états pour la langage  $\{abab^n \mid n \geq 0\} \cup \{aba^n \mid n \geq 0\}$
- Un seul état final pour  $\{a\} \cup \{b^n \mid n \geq 1\}$

### Exercice 5

Donner un AFD pour chacun des langages sur  $\{0, 1\}$  suivants :

- Tout 00 est suivi immédiatement par un 1 (01, 0010, 0010011001 sont dans ce langage et 00100 et 0001 ne le sont pas)
- Les mots où le symbole le plus à gauche est différent de celui le plus à droite.
- Le langage dénoté par l'ER :  $aa^* \mid aba^*b^*$

### Exercice 6

Donner les grammaires sur  $\{a, b\}$  qui génèrent les langages suivants :

- Les mots avec exactement un seul a
- Les mots ayant au moins un a
- Les mots avec au plus trois a
- $L_1 = \{a^n b^m \mid n \geq 0, m > n\}$
- $L_2 = \{a^n b^{2n} \mid n \geq 0\}$
- $L_3 = L_1 L_2$
- $L_4 = L_1 \cup L_2$
- $L_5 = \{w \mid |w| \% 3 = 0\}$
- $L_6 = \{ww^R \mid w \in \{a, b\}^+\}$

**Exercice 7**

Soit l'expression régulière suivante :  $(a^*bc^+ | aacb^+)^+$

- (a) Donner une grammaire régulière à droite qui engendre le même langage
- (b) Déduire l'AFD équivalent
- (c) Analyser les chaînes  $aabcaacb$  et  $bcacb$

**Exercice 8**

Utiliser la construction de Thompson pour concevoir un AFN correspondant à l'expression :  $what|who|why$ , puis donner l'équivalent déterministe et minimal.

**Exercice 9**

Donner une grammaire régulière à gauche qui engendre le langage formé de 0 et de 1 contenant au moins une séquence 010 ou une séquence 000, puis déduire l'AFD équivalent et analyser les chaînes 0110101 et 01101100.

**Exercice 10**

Considérons le langage sur  $\{a, b\}$  formé de mots n'ayant pas deux caractères consécutifs identiques.

- (a) Ce langage est-t-il régulier ? justifier
- (b) Trouver une expression rationnelle qui le dénote

**Exercice 11**

Décrire les langages dénotés par les expressions rationnelles suivantes :

1.  $a(a | b)^*a$
2.  $((\varepsilon | a)b^*)^*$
3.  $(a | b)^*a(a | b)(a | b)$
4.  $a^*ba^*ba^*ba^*$

**Exercice 12**

Minimiser l'automate ci-dessous défini par sa table de transitions (0 est l'état initial et 2 l'unique état final). Donner l'ER du langage reconnu par cet automate.

	a	b
→0	1	5
1	6	2
*2	0	2
3	2	6
4	7	5
5	2	6
6	6	4
7	6	2

# Analyse lexicale

## 2.1 Rôle de l'analyseur lexical

L'analyseur lexical (scanner, lexer ou encore tokenizer en anglais) lit le texte du programme source de gauche à droite<sup>1</sup> caractère par caractère afin d'identifier les différentes entités lexicales qu'il comporte [1, 7]. Pour cela, il découpe le code source pour reconnaître des sous chaînes (dites lexèmes) ayant un rôle dans le langage de programmation du code à compiler et détermine pour chacune sa catégorie ou classe. Le couple (lexème, classe) est appelé token.

La suite des tokens trouvée par l'analyseur lexical est alors communiquée à l'analyseur syntaxique qui représente la phase suivante du processus de compilation. Cette communication est initiée par l'analyseur syntaxique qui demande à chaque fois le token suivant de l'analyseur lexical. Cette interaction est illustrée dans la Figure 2.1.

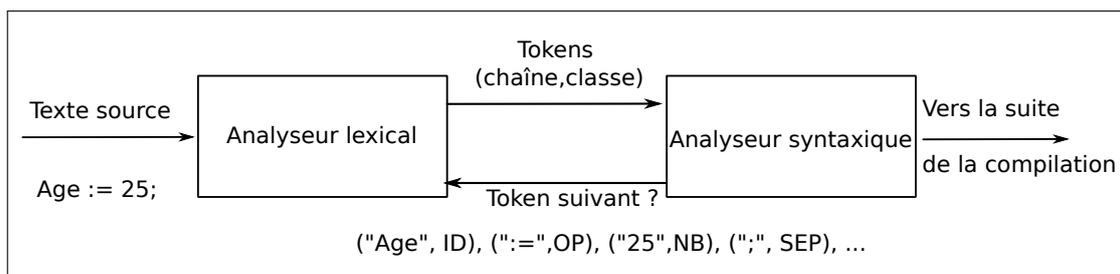


FIGURE 2.1 – Interaction entre l'analyse lexicale et syntaxique

Pour chaque unité lexicale on associe [1] :

- Une **lexème** : la chaîne de caractères correspondante trouvée dans le code source

1. Les langages sont issus du latin. Un scan de droite à gauche pourrait être envisagé dans les cas de l'arabe par exemple

```
1 if (moyenne >= 10) {  
2   Credits = 5;  
3 } else {  
4   Credits = 0;  
5 }
```

FIGURE 2.2 – Portion d'un code source fourni à l'analyseur lexical

```
1 if (moyenne >= 10) {\n\tCredits = 5;\n} else {\n\t Credits = 0;\n}
```

FIGURE 2.3 – Le code précédent vu par l'analyseur lexical comme un flot de caractères

- Une **classe** : la catégorie à laquelle appartient cette unité parmi plusieurs définies par le langage source
- Un **motif** : un modèle permettant de décrire la forme (composition) de l'unité lexicale
- Des **attributs** : pour certaines types d'unités lexicales, toute information complémentaire pour la caractériser (type, valeur, position dans le code source, etc.)

Une implémentation d'un analyseur lexical devrait donc prendre en charge les tâches ci-dessous :

1. Reconnaître les sous-chaînes et leurs classes (les différents tokens).
2. Interagir avec la table de symboles pour la gestion des unités lexicales (codage, insertion, recherche et mise à jour éventuelle dans la suite du processus de compilation)
3. Éliminer l'ensemble de parties du code inutiles à sa traduction, telles que les commentaires et les blancs.
4. Signaler les erreurs lexicales, ou dans les implémentations actuelles entamer une phase de récupération en supprimant une partie de l'entrée jusqu'à ce que la reprise de l'analyse soit possible, ou en proposant des corrections à ces erreurs (insérer, supprimer, permuter des caractères par exemple), voir la Section 3.6.

```

<"if", Mot-clé> <"(", Séparateur> <"moyenne", Identificateur>
<">=", Opérateur> <"10", Nombre> <"", Séparateur> <"{" , Séparateur>
<"Credits", Identificateur> <"=", Opérateur> <"5", Nombre>
<";", Séparateur> <"}", Séparateur> <"else", Mot-clé> <"{" , Séparateur>
<"Credits", Identificateur> <"=", Opérateur> <"0", Nombre>
<";", Séparateur> <"}", Séparateur>

```

FIGURE 2.4 – Suite de tokens émis par l’analyseur lexical

## 2.2 Table des symboles

La table des symboles est une structure de données centrale pour les compilateurs. Elle sert à consigner de manière incrémentale selon l’avancement du processus d’analyse les tokens reconnus et les éventuels attributs associés. Cette structure offre au moins les directives d’ajout, de mise à jour et de recherche des unités lexicales.

Une table de symbole peut être implémentée par une simple structure linéaire comme un tableau ou une liste linéaire chaînée. Toutefois, la complexité des constructions des langages de programmation et les tailles des codes manipulés souvent exigent d’autres implémentations plus efficaces. Pour répondre à ces défis une table de symboles peut être implémentée via des structures composées pour la gestion des portées par exemple. Aussi, les arbres binaires de recherche et les tables de hachages sont alors exploitées dans ces structures.

## 2.3 Spécification lexicale

Tout langage de programmation est construit d’un ensemble de mots (unités lexicales) qu’on peut regrouper en plusieurs catégories. Voici les plus communes.

### 2.3.1 Classes des unités lexicales

Un langage de programmation typique inclut par exemple les catégories ci-dessous :

1. **Mots clés** : begin, if, then, for, ...
2. **Identificateurs** : i, somme1, age, Note\_compil, ...
3. **Nombres** (entiers ou flottants) : 14, -65, 12.33e-4, ...

4. **Opérateurs** : +,-,\*,/,<,<=,>,>=,==,!=,&&,...
5. **Séparateurs** : ;,{, }, ,...
6. **Blancs** : " ", tabulation, nouvelle ligne, ...

Une classe d'unités lexicales englobe donc un ensemble de chaînes obéissant à un motif précis. Pour les catégories en haut on peut lister les motifs suivants :

1. **Mots clés** : la liste des mots réservés du langage.
2. **Identificateurs** : les chaînes non vides formées de lettre et de chiffres ou de l'underscore commençant par une lettre.
3. **Nombres** : les chaînes non vides de chiffres pour les entiers, et un motif plus élaborés pour les réels (voir plus loin)
4. **Opérateurs** : l'ensemble des symboles représentant les différents opérateurs : arithmétiques, logiques, relationnels, etc.
5. **Séparateurs** : l'ensemble des symboles représentant les caractères utilisés comme séparateurs : le point-virgule, les accolades, les parenthèses, etc.
6. **Blancs** : les chaînes non vide de espace, tabulation et nouvelle ligne.

### 2.3.2 Expressions régulières notation étendue

Au niveau lexical, la question est de reconnaître des unités lexicales représentées par des chaînes de caractères (finies). Il est claire que ces ensembles forment des langages réguliers, qu'on peut spécifier par des expressions régulières.

L'analyse lexicale se base sur ces expressions qui forment un outil élégant pour dénoter des ensembles de mots. Elles utilisent les trois opérations régulières par ordre croissant de priorité (+ : l'union, . : la concaténation et \* pour l'itération). Veuillez se référer à la Section 1.5 pour un rappel. Nous donnons ci-après quelque exemples.

**Exemple 5.**

1.  $(0 + 1)1 = \{01, 11\}$
2.  $0^* = \bigcup_{i \geq 0} 1^i = \{\epsilon, 1, 11, 111, \dots\}$
3.  $0^* + 1^* = \{\epsilon, 0, 00, 000, \dots, 1, 11, 111, \dots\}$
4.  $(0 + 1)^* = \bigcup_{i \geq 0} (0 + 1)^i = \{0, 1, 00, 01, 10, 11, 000, \dots\}$  Tous les mots sur  $\{0, 1\}$
5.  $(0(0 + 1))^* = \{\epsilon, 00, 01, 0000, 0001, 0100, 0101, \dots\} = \{\omega \in \{0, 1\}^* \mid |\omega| = 2k, k \in \mathbb{N} \text{ et } 1 \text{ en positions paires}\}$
6.  $a + a(a + b)^* a$  : Les mots sur  $\{a, b\}$  commençants et finissants par un  $a$

**Exercice**

L'expression régulière  $(0 + 1) * 1(0 + 1)^*$  est équivalente à quelle(s) expression(s) parmi les suivantes :

1.  $(01 + 11) * (0 + 1)^*$
2.  $(0 + 1) * (10 + 11 + 1)(0 + 1)^*$
3.  $(1 + 0) * 1(1 + 0)^*$
4.  $(0 + 1) * (0 + 1)(0 + 1)^*$

Depuis les trois opérations régulières de Kleene [14], plusieurs autres extensions et implémentations ont été proposées afin de simplifier et améliorer l'usage des expressions régulières. Voici, ci-dessous, les notations les plus fréquentes [1, 18].

Pour une exploration des expressions régulières dans des langages variés tel que : Java, PHP, Python, etc., le lecteur est invité à consulter le livre [28].

TABLE 2.1 – Expressions régulières étendues communes

$a$	le caractère $a$
$\backslash a$	$a$ littéralement. utilisée pour les métacaractères (les caractères ayant un rôle dans les expressions régulières tel que : $\backslash$ , $ $ , $($ , $)$ , $[$ , $\{$ , $\$$ , $\wedge$ , $*$ , $+$ , $?$ )
$.$	tout caractère sauf fin de ligne
$[abc]$	un caractère parmi $a$ , $b$ ou $c$ . Mettre les caractères $-$ et $]$ en premier pour les inclure dans cette expression. Les méta-caractères n'ont aucune signification ici sauf les séquences reconnues par le langage C (caractères précédés d'un $\backslash$ )
$[a-z]$	le caractère tiret ( $-$ ) indique ici une plage contigue de caractères. Ici toutes les lettres minuscules par exemple, $[0-9]$ signifie tous les chiffres décimaux de 0 à 9
$[\^abc]$	tout caractère différent de ceux de la séquence ici : $a$ , $b$ ou $c$ sauf fin de ligne
$A B$	exprime l'union ou l'alternative. Équivalente à $A + B$
$A^+$	au moins une occurrence de $A$ . Équivalente à $AA^*$
$A?$	0 ou une occurrence de $A$ . Équivalente à $A \epsilon$ , exprime l'option
" $s$ "	la chaîne $s$ littéralement
$\^r$	l'expression $r$ mais uniquement en début de ligne
$r\$$	l'expression $r$ mais uniquement en fin de ligne
$r\{m\}$	$m$ occurrences de $r$ . Si un nom est met à l'intérieur des accolades, elle référence une expression ayant le même nom.
$r\{m,n\}$	de $m$ à $n$ occurrences de $r$
$r\{n,\}$	au moins $n$ occurrences de $r$
$r/s$	l'expression $r$ est reconnue uniquement lorsqu'elle est suivie de $s$ . Cette dernière est ensuite renvoyée à l'entrée
$()$	utilisés pour regrouper des expressions régulières
$\{nom\}$	remplacer $nom$ par l'expression qu'il représente préalablement définie. Utilisée pour simplifier l'écriture d'expressions complexes

### Exercice

- Donner l'expression régulière des nombres réels en notation scientifique (signe et partie décimale optionnelle) [18] :

—  $[-+]?[0-9.]+$  reconnaît plus 1.2.3.4

—  $[-+]?[0-9]+\backslash.[0-9]+$  reconnaît peu : rejette .12 et 12.

—  $[-+]?[0-9]*\backslash.[0-9]+$  n'accepte pas 12.

—  $[-+]?[0-9]+\backslash.[0-9]*$  idem pour .12

—  $[-+]?[0-9]^*\backslash\.[0-9]^*$  accepte le vide et le point seul.

**Solution :**  $[-+]?([0-9]^*\backslash\.[0-9]^+|[0-9]^+\backslash\.)((eE)[-+]?[0-9]^+)?$

Une expression régulière pour les identificateurs pourrait être :

$$[a-zA-Z][a-zA-Z0-9]^*$$

2. Donner une expression régulière pour dénoter des adresses électroniques du genre : `quelquun@serveur.université.xy`

## 2.4 Reconnaissance des unités lexicales

Après avoir spécifié les différentes unités lexicales d'un langage en utilisant des expressions régulières, la question maintenant est de vérifier si une chaîne  $s$  appartient au langage ainsi défini dénoté par ces expressions.

### 2.4.1 Quelques problèmes dans la reconnaissance des tokens

Avant de voir l'approche de reconnaissance, nous présentons ci-dessus quelques difficultés inhérentes à l'analyse lexicale dans des exemples de langages de programmation. Le but est de cerner les problèmes qui peuvent compliquer la tâche des scanners [2].

#### 1. Fortran

Les blancs sont non significatifs. Il en résulte que par exemple l'identificateur `val1` et `val 1` sont identiques. Par ailleurs, l'analyseur lexical ne peut distinguer les deux instructions suivantes qu'après avoir lu le point (.) ou la virgule (,).

— `DO 5 I = 1,25`

— `DO 5 I = 1.25`

Ce genre de problème nécessite l'anticipation (lookahead en anglais) pour décider du partitionnement opportun des unités lexicales.

#### 2. PL/I

Les mots clés dans cet ancien langage ne sont pas réservés. Ceci dit, un programmeur peut coder :

`if else then then = else ; else else = then`

3. **C++**

Des conflits peuvent surgir avec l'usage des templates imbriqués et l'opérateurs de flux ».

4. **Python**

Les blocs sont définis par indentations, ce qui exige de mémoriser les augmentations et retraits de celles-ci.

### 2.4.2 Approche de reconnaissance

L'analyseur lexical scanne l'entrée de gauche à droite caractère par caractère afin de la partitionner à ses tokens. L'algorithme suivant est adopté [2].

1. Écrire les expressions régulières pour les lexèmes de chaque classe (mots clés ( $R_{kw}$ ), identificateurs ( $R_{id}$ ), nombres ( $R_{nb}$ ),...)
2. L'expressions régulière reconnaissant toute les lexèmes du langage est donc l'union de celles-ci.  
$$R = R_{kw} \mid R_{id} \mid R_{nb} \mid \dots$$
3. Pour le préfixe  $a_1 \dots a_i$  ( $1 \leq i \leq n$ ), où  $n$  est la taille de l'entrée, vérifier si  $a_1 \dots a_i \in L(R)$
4. Si succès (oui), alors la chaîne  $a_1 \dots a_i \in L(R_k)$  pour un certain  $k$ ; sinon déclarer le préfixe comme une unité lexicale erronée.
5. Supprimer  $a_1 \dots a_i$  et recommencer à 3.

### 2.4.3 Conventions

Il est aisé de constater des ambiguïtés dans l'approche précédente. Voici quelques-unes et les solutions adoptées.

— **Quel préfixe prendre ?**

Autrement dit à quelle position il faut s'arrêter pour déclarer la détection d'une unité lexicale. Durant le partitionnement, il arrive que  $a_1 \dots a_i$  et  $a_1 \dots a_j$  avec  $j > i$  sont deux préfixes de notre langage. Dans de pareils cas, la convention est de prendre toujours **le préfixe le plus long**. Cette convention est appelée la règle du *maximal-munch* [25]. Exemples : 325.66e-4 est un seul token au lieu de : 325 suivi de .66, puis e et enfin -4. Les opérateurs doubles sont aussi des

exemples :  $\leq$  (inférieur ou égal au lieu de inférieur puis égal),  $\geq$ ,  $==$ ,  $:=$ ,  $\langle \rangle$ , etc.

— **Quelle classe de tokens choisir ?**

Il est souvent possible de tomber sur un préfixe  $a_1 \dots a_i$  qui peut appartenir à plus d'une classe  $a_1 \dots a_i \in L(R_m)$  et  $a_1 \dots a_i \in L(R_n)$  avec  $m \neq n$ . L'exemple courant est celui des mots clés qui sont aussi des identificateurs. Dans ces situations, **on définit un ordre de priorité entre les classes du langage, et on opte pour celle ayant la plus haute priorité**. Dans les implémentations (voir la section 4) l'ordre d'apparition dans la spécification des classes définit cette priorité.

— **Aucune classe ne s'applique !**

Si le préfixe en cours n'appartient pas au langage, il faut signaler une erreur. Pour éviter le blocage de l'analyseur lexical, il est commode d'intercepter ces cas dans une classe **erreur** ayant la plus **faible priorité**.

## 2.5 Implémentation d'analyseurs lexicaux

Après avoir acquis le fondement de l'analyse lexicale d'un code source, son implémentation peut être réalisée par plusieurs méthodes. Nous décrivons dans cette section les plus courantes.

### 2.5.1 Analyseur lexical en dur

Cette méthode repose sur l'écriture manuelle (ad-hoc) de l'analyseur. On prépare l'ensemble des unités lexicales à reconnaître et on écrit le code de reconnaissance pour chacune d'elles. Le squelette du scanner sera composé essentiellement d'une boucle de lecture guidée par une succession de tests sous la forme d'une pile de `if then else` ou par un `switch case` par catégorie de tokens. La simplicité et l'efficacité de cette méthode est évidente. Cependant, elle reste exposée à l'erreur et difficile à maintenir ou étendre. L'Algorithme 1, donne un exemple d'un analyseur simple permettant de reconnaître les trois tokens : **end**, **else**, et des **identificateurs**.

---

**Algorithme 1** Pseudo code d'un mini-analyseur lexical codé en dur

---

```
1 c = car_suiv();
2 if (c == 'e'){
3   c = car_suiv();
4   if (c == 'n'){
5     c = car_suiv();
6     if (c == 'd'){
7       c = car_suiv();
8       if (!lettre(c) && !chiffre(c)) {
9         return KW_END;
10      }
11     else return ID;}
12   else return ID;}
13 else if(c == 'l'){
14 //... continuer pour reconnaître else ou un ID
15 }
```

---

## 2.5.2 Implémentation par des automates finis

Ici, on construit l'automate (déterministe et minimal) [12] global de reconnaissance de l'ensemble des unités lexicales à partir de leur expressions régulières les spécifiant. L'automate (union des automates des classes de ces unités) est ensuite implémenté soit itérativement ou par des procédures (ou fonctions) mutuellement récursives associées aux différents états de l'automate. Un pseudo-code est montré dans l'Algorithme 2.

---

**Algorithme 2** Implémentation d'un automate par des fonctions

---

```
1 int etat0(void){
2   switch(fgetc(input)){
3     case car1 : return etat1(); break;
4     case car2 : return etat2(); break;
5     ...
6     case EOF : return final(); break;
7     default : return 0;
8   }
9 }
```

---

Une autre approche d'implémentation d'un automate consiste à déclarer et initialiser sa table de transitions. La reconnaissance est essentiellement une boucle de lecture et de passages entre états. Cette méthode est relativement efficace si la taille de la table des transitions reste raisonnable. Un pseudo-code général de reconnaissance d'une chaîne par un automate via sa table de transitions est présenté dans

l'Algorithme 3.

Il convient de préciser que l'implémentation d'un analyseur lexical par l'automate associé exige souvent des techniques d'anticipation pour assurer la reconnaissance du préfixe le plus long d'une part, et d'autre un moyen pour restaurer un ou plusieurs caractères à l'entrée afin de reprendre l'analyse et reconnaître les tokens suivants jusqu'à atteindre la fin de l'entrée (EOF) qui doit être prise en considération comme un caractère spécial.

---

**Algorithme 3** Implémentation d'un automate fini via sa table de transitions

---

```
1  etat = e0;
2  lexeme = '';
3  c = fgetc(input);
4  while (c != EOF && etat != ERREUR)
5  {
6      etat = TTrans[etat][c];
7      c = fgetc(input);
8      lexeme += c;
9  }
10 return (etat != ERREUR && final(etat));
```

---

Rappelons que la passage d'une expression régulière à l'automate équivalent peut être effectué par plusieurs méthodes. Nous donnons ici la construction de Thompson qui génère un automate non déterministe AFN (avec  $\varepsilon$ -transitions) et la méthode des dérivées permettant d'avoir directement un automate déterministe AFD.

### 2.5.2.1 D'une expression régulière à un AFN : la construction de Thompson

Cette construction [29] suit la définition récursive des expressions régulières et donne pour chacun des six cas (3 de base et 3 de d'induction) l'automate de Thompson associé. Les automates des cas de base sont schématisés dans la Figure 2.5, et ceux des cas d'induction dans la figure 2.6. La Figure 2.7, quant à elle, montre l'automate obtenu après application de la construction de Thompson correspondante à l'expression  $(a|b)*c$ .

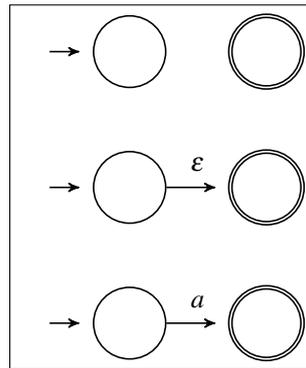


FIGURE 2.5 – Construction de Thompson pour les cas de base :  $\phi$ ,  $\epsilon$  et  $a$

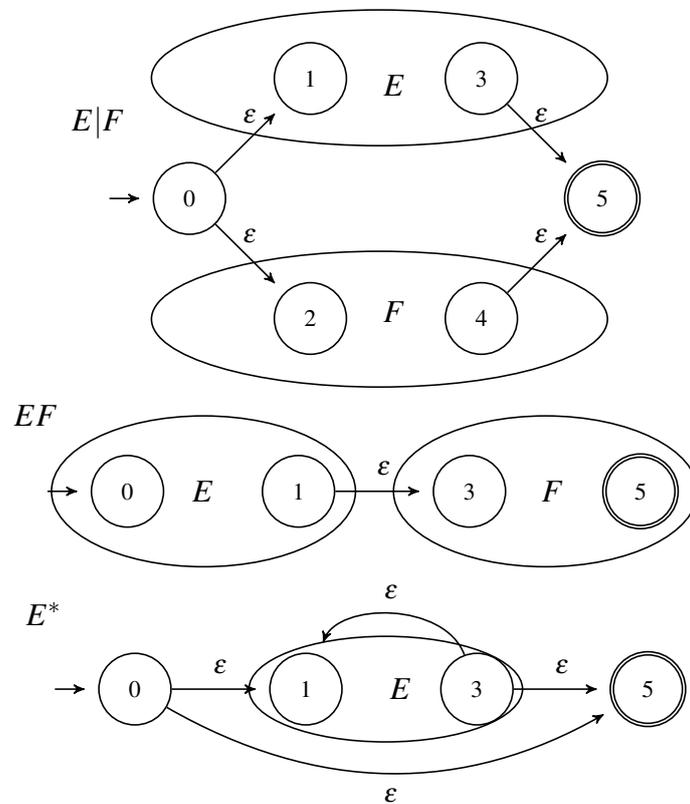
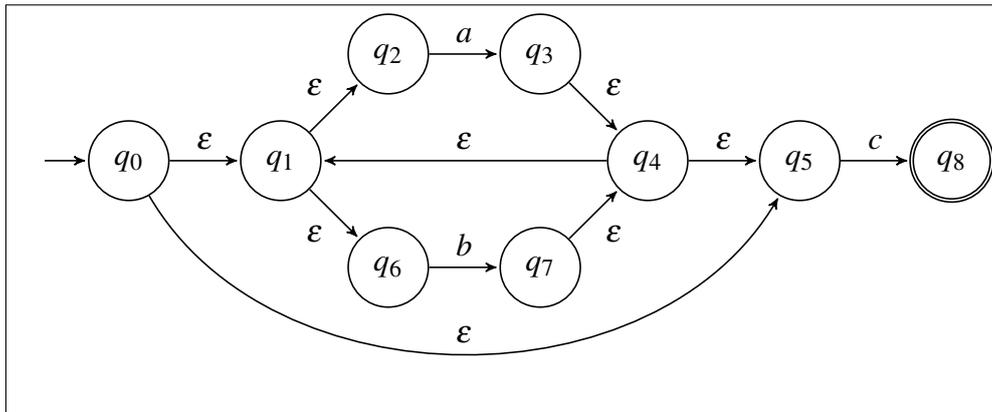


FIGURE 2.6 – Construction de Thompson pour les cas d'induction :  $E|F$ ,  $EF$  et  $E^*$

**2.5.2.2 D'une expression régulière à un AFD : les dérivées de Brzozowski**

La dérivée (dûe à Janusz A. Brzozowski) [5, 22] d'une expression régulière  $E$  sur un alphabet  $\Sigma$  par rapport à un symbole  $a \in \Sigma$ , notée  $a^{-1}(E)$ , est définie comme

FIGURE 2.7 – Construction de Thompson pour l'expression  $(a|b)^*c$ 

suit :

$$\left\{ \begin{array}{l} a^{-1}(\phi) = a^{-1}(\varepsilon) = a^{-1}(b) = \phi \text{ pour } b \neq a \\ a^{-1}(a) = \varepsilon \\ a^{-1}(E | F) = a^{-1}(E) | a^{-1}(F) \\ a^{-1}(EF) = a^{-1}(E)F | \varepsilon(E)a^{-1}(F) \\ a^{-1}(E^*) = a^{-1}(E)E^* \end{array} \right.$$

On peut aussi définir la dérivée par

$$a^{-1}E = \{w \mid aw \in L(E)\}$$

Avec  $\varepsilon(E)$  est une fonction qui aide dans le calcul des dérivées. Elle exprime l'appartenance ou non du mot vide à  $E$ .

$$\varepsilon(E) = \begin{cases} \varepsilon & \text{si } \varepsilon \in L(E) \\ \phi & \text{sinon} \end{cases}$$

Intuitivement, la dérivée d'un langage dénotée par l'expression régulière  $E$  par rapport à un symbole  $a$  est l'ensemble des mots générés en supprimant le premier  $a$  des mots de  $E$  qui commencent par  $a$  (ou encore tout ce qui peut suivre  $a$  dans l'expression  $E$ ). Par exemple, la dérivée de  $ab^*$  par rapport à  $a$  est  $b^*$ , par contre sa dérivée par rapport à  $b$  est le langage vide.

Les états de l'automate résultat sont les différentes dérivées. L'état initial est celui

représentant l'expression globale. Une transition de l'état correspondant à l'expression  $s$  par le symbole  $a$  est créée entre un état correspondant à  $s$  et celui correspondant à l'expression dérivée  $a^{-1}(s)$ . Un état est final si l'expression qu'il représente inclut le mot vide  $\varepsilon$ . Le lecteur intéressé pourra trouver les détails et les preuves (correction, terminaison) de cette construction dans [5]. Un exemple de cette construction pour l'expression  $(a | b)^* bab$  est montré dans la Figure 2.8.

### Exercice

Donner l'AFD qui reconnaît le langage de l'expression :  $(ab)^* | (a|b)^*$

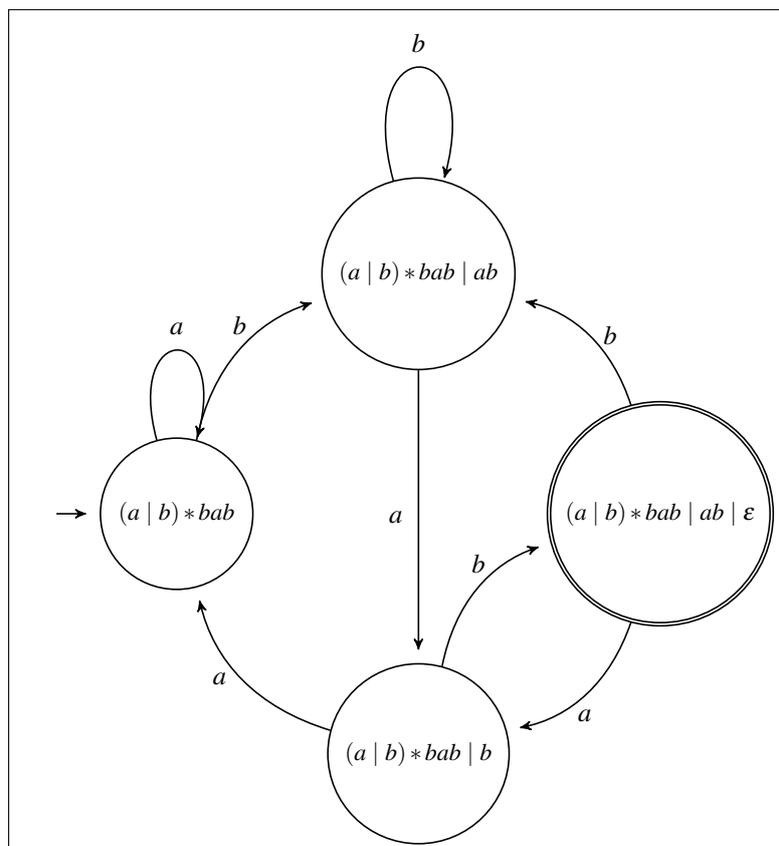


FIGURE 2.8 – Automate de  $(a | b)^* bab$  par la méthode des dérivés.

### 2.5.3 Générateurs d'analyseurs lexicaux cas de Flex

La maîtrise des outils et des modèles théoriques inhérents à la compilation a permis la mise en œuvre de générateurs de différents types d'analyseurs pour la plupart des environnements et langages de programmation. Un générateur d'analyseur lexical est un outil qui accepte en entrée la spécification des unités lexicales sous formes d'expressions régulières et produit en sortie le code permettant la reconnaissance de ces unités. Le cœur de ce code est l'automate (déterministe et minimal) global de reconnaissance. Parmi ces outils, nous citons : lex d'unix [21] et sa version libre de gnu sous linux (Flex) [23], Jlex pour Java [15] et PLY pour Python [4], etc. Dans cette section, nous présentons l'outil Flex.

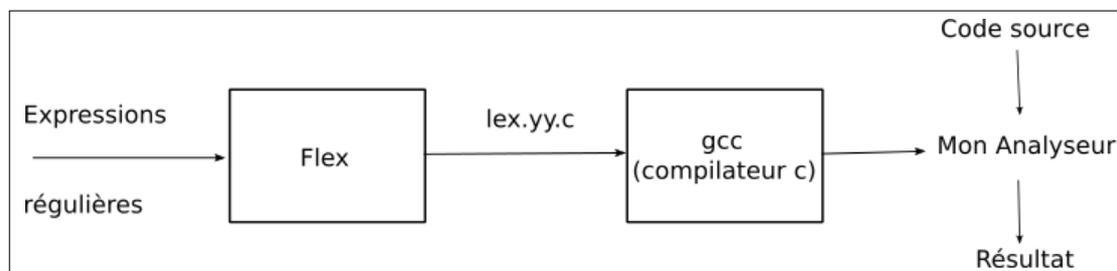


FIGURE 2.9 – Fonctionnement de Flex

#### 2.5.3.1 Structure du fichier de spécification flex

Flex [18], dont le principe est exhibé par la Figure 2.9, génère un code en C (dans le fichier *lex.yy.c* par défaut où est définie la routine *yyllex*) à partir d'un fichier texte de spécification lexicale (souvent ayant *l* ou *lex* comme extension). Le code produit est à compiler pour obtenir l'exécutable de l'analyseur. Le fichier d'entrée de flex est composé des parties ci-dessous comme l'illustre l'Algorithme 4.

#### 2.5.3.2 Définitions régulières

Elles associent des noms à des expressions qui permettent ensuite de simplifier l'écriture du code. Elles ont la forme suivante : **nom** *exp*

Voici des exemples :

*lettre* [a – zA – Z]

*chiffre* [0 – 9]

**Algorithme 4** Structure générale d'un code Flex

```

1 %{
2  Declaration de variables, constantes, includes, etc.
3  /* partie optionnelle qui sera copiee sans modification au debut
   du code source produit */
4  %{
5  Declaration des definitions regulieres    (optionnelle aussi)
6  %%
7  Regles de traduction
8  %%
9  Bloc principal et fonctions auxiliaires
10 /* cette partie est optionnelle et sera copiee a la fin du code
   source produit */

```

Une expression définie dans cette partie peut être ensuite utilisée dans la section suivante relative aux règles de traduction en l'entourant entre une paire d'accolades `{}`. Ainsi, l'expression d'un identificateur est formulée comme suit en utilisant les deux précédentes : `{lettre}({lettre}|{chiffre})*`

**2.5.3.3 Règles de traduction**

C'est la partie principale d'un code flex. Elle définit des expressions régulières (partie gauche) et les actions associées (partie droite comme un code en langage C) à exécuter dans le cas de leur reconnaissance selon la forme ci-dessous.

```

1  exp_1  {action_1}
2  exp_2  {action_2}
3  ...

```

**Exemple 6.** `(0|1)+ { printf("%s est un nombre binaire",yytext); }`

Un exemple complet de l'utilisation de flex pour générer un analyseur d'un mini-langage est présenté dans la Figure 2.10. Mentionnons pour finir que Flex offre plusieurs fonctions et variables prédéfinies :

- **yylex()** la fonction de l'analyseur,
- **yytext** la chaîne reconnue dans le texte source lors de l'action en cours,
- **yylen** représente la longueur de la chaîne reconnue.
- **yyval** les attributs du token
- **yyloc** position (ligne et colonne) du token dans le code source

- **ECHO** l'action par défaut de flex,
- **yyin et yyout** le fichier d'entrée/de sortie de flex (stdin et stdout par défaut).

```

1 %{
2 /* pour la fonction atof() */
3 #include <math.h>
4 %}
5 NB      [0-9]
6 ID      [a-z][a-z0-9]*
7 %%
8 {NB}+ { printf( " Un Entier : %s (%d)\n", yytext, atoi( yytext ) );}
9 {NB}+"."{NB}* {printf(" Un flottant : %s (%g)\n", yytext, atof(
   yytext));}
10 if|then|begin|end|procedure|function { printf( " Mot clef : %s\n",
   yytext );}
11 {ID} printf( " Un identificateur : %s\n", yytext );
12 "+"|"-"|"*"|"/" printf( " Un Operateur : %s\n", yytext );
13 "/"|".*" /* ignorer les commentaires sur une ligne */
14 [ \t\n]+ /* ignorer les espaces et les blanc */
15 . printf( " Caractere invalide : %s\n", yytext );
16 %%
17 main(argc, argv)
18 int argc;
19 char **argv;
20 {
21   if ( argc > 0 )
22     yyin = fopen( argv[0], "r" );
23   else
24     yyin = stdin;
25   yylex();
26 }

```

FIGURE 2.10 – Un exemple complet de code Flex

### 2.5.3.4 Autres usages de Flex

On peut faire plusieurs tâche une fois une chaîne suivant un motif donné a été reconnue. Cela permet d'exploiter Flex pour effectuer non seulement de l'analyse lexicale, mais aussi des actions variées telles que : recopier, supprimer, convertir, etc.

## **Exercices du chapitre 2**

Ministère de l'Enseignement Supérieur et de la Recherche Scientifique  
Université de Ghardaia  
Faculté des Sciences et de la Technologie  
Département des Mathématiques et d'Informatique

---

## ANALYSE LEXICALE

---

### Exercice 1

Soit la portion du code source suivante :

```
x = 0 ;\n\twhile (x < 10) { \n\ttx++;\n }
```

Déterminer le nombre d'unités lexicales pour chaque classe de tokens (Mots-clés, Identificateurs, Nombres, Opérateurs, Blancs, Sep, etc.)

### Exercice 2

Soit le code C++ ci-dessous

```
float limitedSquare(float x) {  
/* return x-squared, but never more than 100 */  
return (x <= -10.0 || x >= 10.0)? 100 : x*x; }
```

- Diviser ce code à ses tokens appropriés
- À quels lexèmes doit-t-on associer des valeurs ? Quelles sont ces valeurs ?

### Exercice 3

Proposer des spécifications de l'heure donner au format 12H exemple : "05:14 PM". Les minutes sont toujours des nombres à deux chiffres, mais l'heure peut être un chiffre unique.

### Exercice 4

En SQL, les mots clés et les identificateurs ne sont pas sensibles à la casse. Écrire un programme Flex qui reconnaît les mots clés SELECT, FROM, WHERE (avec n'importe quelle combinaison de minuscules et de majuscules), ainsi que l'unité lexicale ID.

### Exercice 5

Donner des expressions régulières pour dénoter :

- Une chaîne quotée en C (exemple : "compilation"), sachant qu'une chaîne quotée ne peut s'étendre sur plusieurs lignes
- Un commentaire en C++ introduit par //
- Un commentaire multi-lignes en C++ autorisant ou non le caractère \* dans le texte du commentaire

### Exercice 6

Écrire un analyseur lexical qui reconnaît les unités X, Y et Z suivantes. On supposera que les unités du langage sont séparées par un ou plusieurs blancs.

- L'unité X est une suite de chiffres décimaux qui peut commencer par le caractère '-'. Sa valeur est inférieur à 32768.
- L'unité Y commence par une lettre suivie par une suite de lettres, de chiffres et de tirets. Elle ne se termine jamais par un tiret et ne comporte pas de tirets qui se suivent. Sa longueur maximale est de 31 caractères.

- (c) L'unité Z est une suite de chiffres qui peut commencer par le caractère '-' et qui comporte le point décimal. Le nombre de chiffres est inférieur ou égal à 9.

### Exercice 7

Un commentaire en PASCAL est une suite de 80 caractères au plus délimitée par (\* et \*) ne contenant pas la chaîne \*) à moins que celle-ci n'apparaisse dans une chaîne bien quotée.

- (a) Construire un automate déterministe qui accepte un commentaire en PASCAL  
 (b) Écrire un algorithme de reconnaissance d'un commentaire

### Exercice 8

Considérons la chaîne *abbbaacc*. Quelle(s) spécification(s) lexicale(s) produi(sen)t la tokenisation : *ab/bb/a/acc* :

<i>b+</i>	<i>c*</i>	<i>a(b c*)</i>	<i>ab</i>
<i>ab*</i>	<i>b+</i>	<i>b+</i>	<i>b+</i>
<i>ac*</i>	<i>ab</i>		<i>ac*</i>
	<i>ac*</i>		

### Exercice 9

Pour les entrées ci-après énumérées, donner la liste des tokens reconnus par un analyseur implémentant la spécification suivante :

---

*a(ba)\**  
*b\*(ab)\**  
*abd*  
*d+*

---

1. *dddabbabab*
2. *ababddababa*
3. *babad*
4. *ababdddd*

### Exercice 10

- (a) Étant donné la spécification lexicale ci-dessous.

```
%%
aa      printf("1");
b?a+b?  printf("2");
b?a*b?  printf("3");
.       printf("4");
```

Indiquer pour les chaînes suivantes la sortie de l'analyseur et les tokens produits :

- *bbbabaa*
  - *aaabbbbbbaabanalyse*
- (b) Donner un exemple d'une entrée pour laquelle cet analyseur produit la sortie **123**, ou justifier, le cas échéant, l'inexistence de tel exemple.
- (c) Quelle sera pour les mêmes chaînes données en (1) la sortie de l'analyseur, si on remplace seulement la partie expression régulière de la dernière règle par *.+*

## Mini projet 1 de Compilation

---

### 1 Objectif

Le but de ce TP est de réaliser la première phase de la compilation (analyse lexicale) pour un mini-langage de programmation de deux manières différentes.

### 2 Spécification lexicale

Nous avons défini le mini-langage de programmation **Slim** composé des éléments suivants :

- Les mots clés (**KWD**) insensibles à la casse : **si, sino, is, tq, qt, rpt, jsq, lre, ecr, vrai, faux**
- Les entiers naturels (**INT**) non signés composés par des suites non vides de chiffres
- Les identificateurs (**ID**) : les suites de lettres ou de chiffres commençant par une lettre
- Les opérateurs (**OPR**) : + - \* / = < <= > >= == != && || !
- Séparateurs (**SP**) : ; , ( ) :
- Les blancs (**BLC**) : les suites non vides de : espace, tabulation et fin de ligne

### 3 Travail à faire

1. Écrire un analyseur lexical ad-hoc (codé en dur) pour **Slim**. Vous êtes libre dans le choix du langage de réalisation (C, Java, Python, etc.). L'analyseur lexical doit recevoir le code source dans un fichier texte et fournir son résultat, après avoir retourné le token rencontré, dans un autre fichier de sortie en respectant le format suivant :

Classe	Ligne	Colonne	Lexème
KWD	1	4	si
INT	2	5	123
ID	3	1	somme1
...	...	...	...

Nous enrichissons maintenant notre mini-langage par les nouveaux tokens ci-après :

- (a) Les flottants (**FLT**) : les nombre réels signés en notation scientifique

(b) Les chaînes de caractères (**CHN**) : elles sont incluses entre double quotes "...". Dans une chaîne une séquence '\c' dénote le caractère c avec les exceptions suivantes :

- \b backspace
- \t tab
- \n newline
- \f formfeed

Vous devez respecter les restrictions suivantes sur les chaînes :

- Une chaîne ne peut contenir le nul (le caractère \0) ni le EOF
- Une chaîne ne peut dépasser les limites des fichiers
- Une chaîne ne peut inclure de nouvelles lignes non escapées. Par exemple :  
"Cette chaîne \  
est correcte "  
par contre la suivante est invalide :  
"Cette chaîne  
est incorrecte "

(c) les commentaires (**CMT**) au style du C : sur une seule ligne introduits par "//" ou ceux multi-lignes encadrés par les couples "/\*" et "\*/". Cette dernière forme peut être imbriquée.

2. Réécrire le même analyseur lexical en utilisant cette-fois l'outil **flex** (même entrée/sortie en fichiers textes).

## 4 Gestion des erreurs

Toutes les erreurs doivent être retournées au analyseur syntaxique. Elles seront communiquées en retournant un token spécial **ERROR** et suivant les exigences ci-après :

- Si un caractère invalide est rencontré, la chaîne contenant ce caractère doit être retournée comme chaîne d'erreur,
- Si une chaîne inclut un newline non escapé reporter l'erreur : "Constante chaîne de caractères non terminée", et reprendre l'analyse au début de la ligne suivante,
- Dans le cas d'une chaîne trop longue reporter "Constante chaîne de caractères trop longue" dans la chaîne erreur du token **ERROR**. Si la chaîne contient un caractère invalide (le nul par exemple) reporter l'erreur "Chaîne contenant le caractère nul", dans les deux cas l'analyse reprendra après la fin de la chaîne. La fin d'une chaîne est définie soit comme :
  1. le début de la prochaine ligne si un newline non escapé est rencontré après ces erreurs ou
  2. après les deux quotes de fermeture "
- Si un commentaire reste ouvert alors qu'un EOF est rencontré reporter l'erreur : "Fin de fichier dans un commentaire" ne pas considérer alors ce token étant donné que la terminaison est manquante. De même, si EOF est rencontrés avant la quote de fermeture d'une chaîne reporter "Fin de fichier dans une constante chaîne"
- Si vous rencontrez "\*/" en dehors d'un commentaire reporter "\*/ sans correspondance" au lieu de la tokeniser en \* et /.

## 5 Construction de l'analyseur

1. Pour les codages en dur utiliser votre IDE favori
2. Compiler la spécification Flex par : **Flex -o scanner.c scanner.l** : produira le fichier **scanner.c** à partir de scanner.l
3. Compiler le tout par : **gcc -o analyseur scanner.c -lfl** : produira l'exécutable **analyseur** à partir de l'ensemble des fichiers
4. Pour le codage des tokens vous pouvez soit les définir comme de simples entiers, ou utiliser l'outil **Bison**. Dans ce dernier, on définit un token en le précédant de l'option **%token**. Compiler le fichier bison en utilisant l'option **-d** pour produire la définition de vos tokens en fichier entête **.h** à inclure dans votre fichier flex.
5. Invoquer votre analyseur par : **./analyseur** : où vous pouvez taper votre code source en **Slim** ou en le passant à travers un fichier texte et invoquer l'analyseur par la commande **./analyseur < votre-fichier-source**

# Analyse syntaxique

---

## 3.1 Rôle de l'analyseur syntaxique

La mission de l'analyseur syntaxique (parser en anglais) est de vérifier la conformité de la suite de tokens délivrés par l'analyseur lexical à la définition syntaxique du langage source souvent donnée par une grammaire à contexte libre (de type 2 dans la hiérarchie de Chomsky). Autrement dit, valider que l'ordre des tokens peut être produit par la grammaire du langage.

Pour les programmes bien formés, *i.e* syntaxiquement corrects, l'analyseur syntaxique construit un arbre syntaxique et le transmet aux phases suivantes du processus de compilation (voir figure 3.1). Dans le cas contraire, il doit signaler une erreur syntaxique et entamer une procédure de rattrapage [1, 32].

Deux approches sont utilisées pour vérifier la validité et construire l'arbre syntaxique d'un programme source. La première adopte une méthode descendante qui consiste à démarrer à partir de l'axiome de la grammaire pour aboutir au code source en effectuant une suite de dérivations ; c-à-d des remplacements successives de membres gauches de productions par des membres droits correspondants. La deuxième procède inversement et est qualifiée d'ascendante, car elle part du texte source et remonte à l'axiome par des suites valides de réductions ou de remplacements cette fois de membres droits de production par des membres gauches associés. Ces deux approches seront discutées dans ce chapitre. Commençons par rappeler quelques notions utiles pour la suite du cours.

## 3.2 Spécification syntaxique

Au niveau de l'analyse lexicale, les lexèmes des tokens, étant des chaînes simples, sont décrits alors par des expressions régulières car ils forment un langage régulier. Le niveau syntaxique des langages de programmation inclut de nombreuses

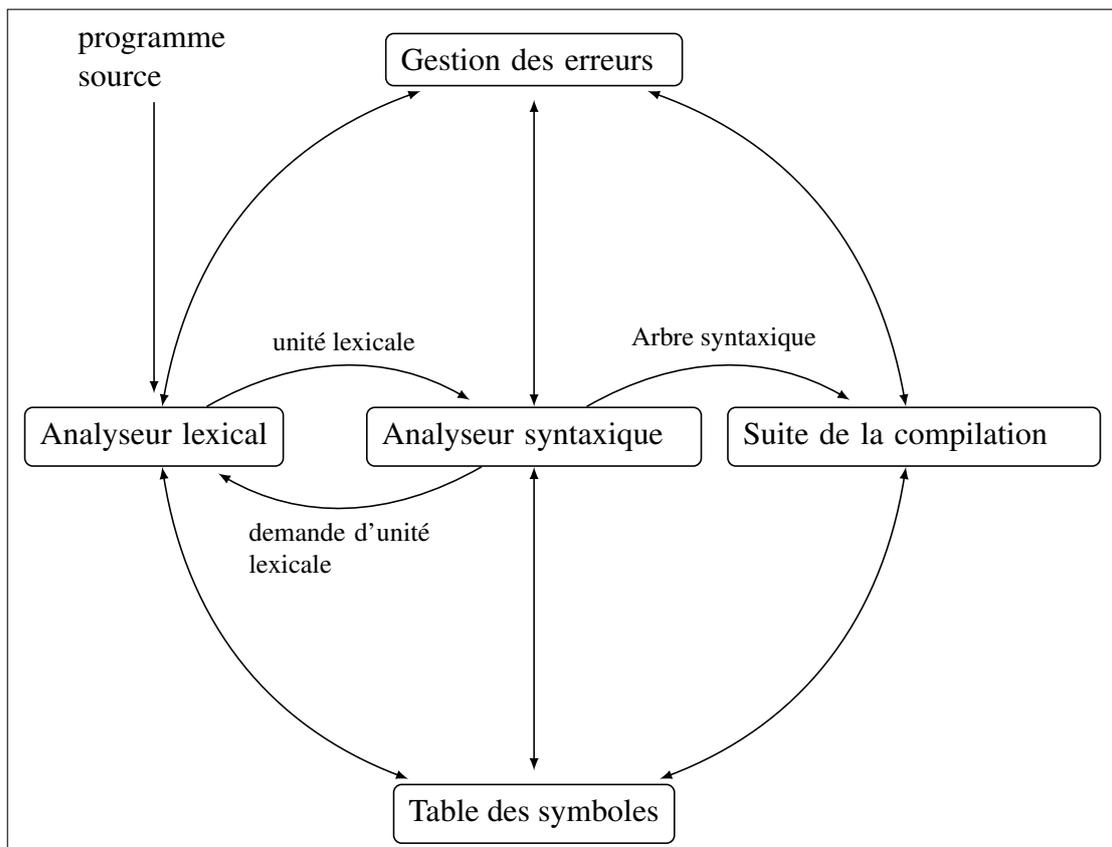


FIGURE 3.1 – Interaction entre les différents modules de la partie analyse d'un compilateur

constructions complexes qui ne sont pas régulières. Ces constructions ne peuvent être décrites par des expressions régulières. Pour l'analyse syntaxique des langages de programmation, nous devons opter pour un autre modèle ou outil plus puissant qui est les grammaires.

### Exemple

- $L = \{(n)^n \mid n \geq 0\}$  langage des expressions bien parenthésées.  $()$ ,  $(())$ , tel que dans  $((1 + 2) * 3) \dots$
- Délimiteurs de blocs  $\{\dots\}$  de Java et C ou les **begin end** de Pascal.
- Les suites de **if then else**

Pour que l'analyseur syntaxique puisse détecter les suites valides d'unités lexicales de celles invalides, il a besoin de :

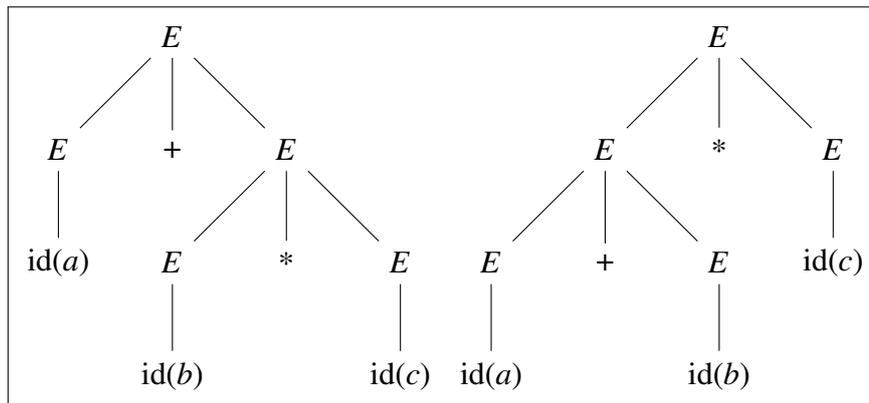
- Un outil décrivant les suites valides d'unités lexicales : les **grammaires**
- Une méthode ou un **algorithme** pour assurer cette séparation.

### 3.2.1 Grammaires hors contexte et dérivation vs réduction

Les grammaires à contexte libre (non contextuelles, hors contexte ou encore algébriques) sont adaptées à la description de la plupart des constructions des langages de programmation.

Rappelons que les règles de production ou de réécriture d'une grammaire à contexte libre sont de la forme  $A \rightarrow \alpha$ , où  $A$  est un non terminal  $\in \mathbb{N}$  et  $\alpha \in (\mathbb{T} \cup \mathbb{N})^*$  une chaîne quelconque de terminaux et de non terminaux. La suite de tokens fournie par l'analyseur lexical forme une chaîne représentant la forme du texte source vu au niveau syntaxique. La question ici est de montrer qu'une telle chaîne disons  $w$  appartient bien au langage généré par la grammaire  $G$  en question. Autrement dit, vérifier si la proposition  $w \in \mathbb{L}(G)$  ? est vraie.

La réponse à cette question peut être donnée par deux principales méthodes : la première est descendante : en partant de l'axiome de la grammaire et opérant une succession de substitutions dite **dérivations** (remplacement du membre gauche d'une production ici toujours un non terminal par le membre droit) jusqu'à aboutir à la chaîne. Cette méthode est à l'origine des méthodes d'analyse syntaxique descendante. Une alternative consiste à adopter le chemin inverse. En commençant à

FIGURE 3.2 – Deux arbres syntaxiques pour dériver la chaîne :  $w = a + b * c$ 

partir de la chaîne du texte source, nous effectuons cette fois-ci une séquence de **réductions** (remplacement du membre droit d'une règle de production par le membre gauche) de proche en proche jusqu'à arriver à l'axiome. Cette méthode donne naissance aux méthodes d'analyse syntaxique ascendantes que nous allons traiter ultérieurement.

Dans un processus de dérivation, si toujours le non terminal le plus à gauche (resp. à droite) est remplacé en premier en parle de **dérivation la plus à gauche (DPG)** (resp. à droite (DPD)). **Leftmost** vs **rightmost** derivation en anglais.

**Exemple 7.** Soit la grammaire ci-dessous :

$$\begin{cases} S \rightarrow aTb \mid c & \text{et la chaîne } w = accacbb \\ T \rightarrow cSS \mid S \end{cases}$$

DPG :  $S \rightarrow aTb \rightarrow acSSb \rightarrow accSb \rightarrow accaTbb \rightarrow accaSbb \rightarrow accacbb$

DPD :  $S \rightarrow aTb \rightarrow acSSb \rightarrow acSaTbb \rightarrow acSaSbb \rightarrow acSacbb \rightarrow accacbb$

**Exemple 8.**

$$E \rightarrow E + E \mid E * E \mid id \quad \text{et la chaîne : } w = a + b * c$$

Avant d'entamer les méthodes d'analyse syntaxique, nous discutons dans la section suivante quelques propriétés des grammaires qui peuvent influencer la qualité ou voire même l'existence d'une telle méthode d'analyse. Pour un analyseur syntaxique efficace, la grammaire doit posséder certaines propriétés.

## 3.2.2 Qualités et formes particulières des grammaires

### 3.2.2.1 Ambiguïté

Pour une grammaire  $G$ , s'il existe un mot  $w$  de  $L(G)$  ayant plusieurs arbres syntaxiques, on dit que  $G$  est **ambiguë**. Nous pouvons rencontrer l'ambiguïté dans plusieurs cas. Voici deux exemples typiques de grammaires ambiguës indispensables dans tous les langages de programmation.

— Les expressions arithmétiques :

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

— Les expressions conditionnelles : ( $I$  et  $C$  sont deux non terminaux permettant de générer les instructions et les conditions respectivement )

$$\begin{cases} I \rightarrow \text{if } C \text{ then } I \\ I \rightarrow \text{if } C \text{ then } I \text{ else } I \end{cases}$$

**Exemple 9.** *if  $x > 10$  then if  $y < 0$  then  $a := 1$  else  $a := 0$*

L'ambiguïté est une propriété indésirable dans la construction d'analyseurs syntaxiques. Malheureusement, il n'existe aucune méthode ni pour détecter, ni pour transformer l'ambiguïté d'une grammaire [1, 2, 12]. Comment faire ?

Devant une grammaire ambiguë, nous devons :

- Soit la transformer, *i.e* construire une grammaire équivalente non ambiguë
- Soit la retenir mais introduire en plus des mécanismes de désambiguïsation. On peut fixer par exemple des règles de priorité, ou définir des conventions de précedence ou d'associativité (à gauche ou à droite), etc.

La grammaire ambiguë des expressions arithmétiques précédente peut être transformée comme suit :

$$E \rightarrow T + E \mid T \qquad T \rightarrow id * T \mid id \mid (E) * T \mid (E)$$

En utilisant la grammaire transformée il n'y a qu'une seule façon de dériver la chaîne de l'exemple celle de l'arbre syntaxique à gauche dans la figure 3.2.

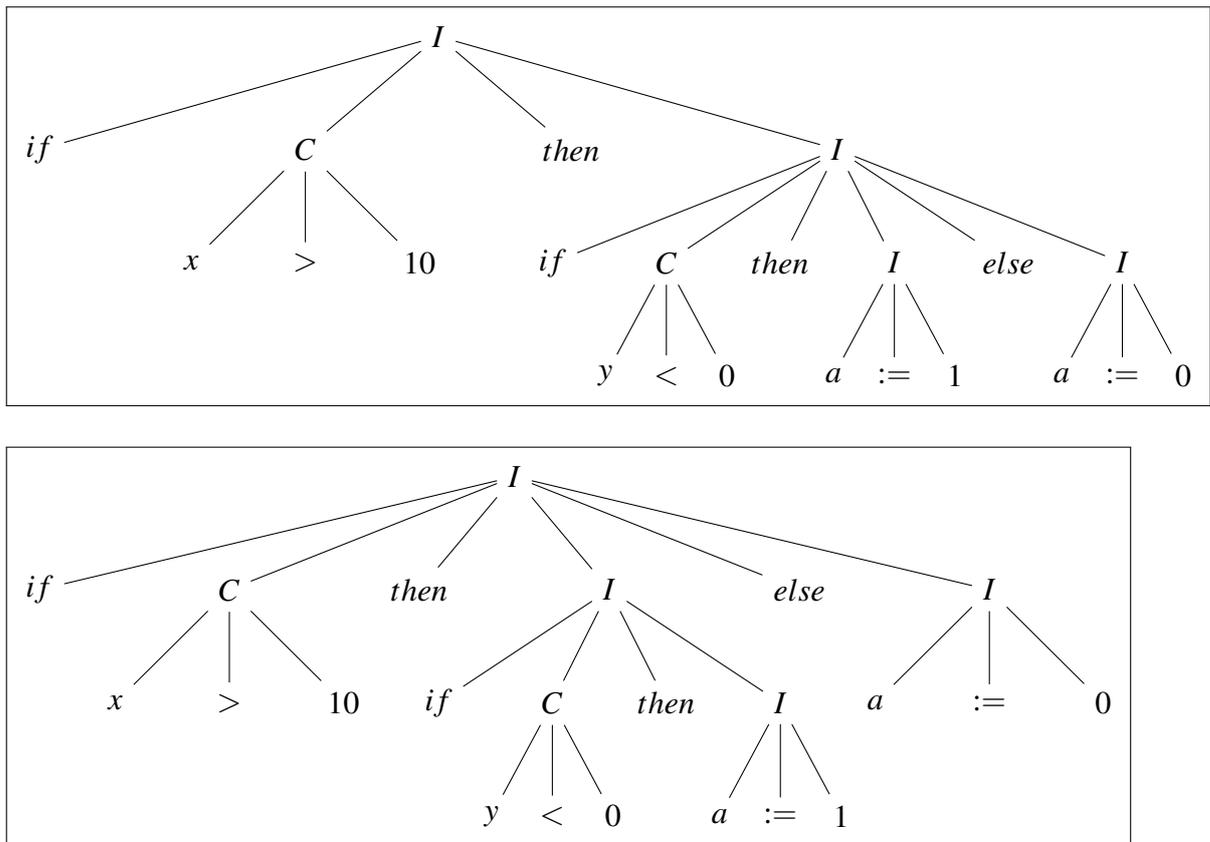


FIGURE 3.3 – Deux arbres syntaxiques de la même chaîne dans la grammaire des instructions conditionnelles

### 3.2.2.2 Récursivité à gauche

Dans les langages de programmation les productions sont souvent définies en terme d'elles mêmes. par exemple la production utilisée dans les déclarations d'une liste de variables :

$$L \rightarrow v \mid L, v$$

Une grammaire est dite réursive à gauche (RG) de façon directe si elle contient une règle de la forme :

$$A \rightarrow A\alpha \mid \beta \quad A \in \mathbb{N}; \alpha, \beta \in (\mathbb{T} \cup \mathbb{N})^*$$

Parfois la RG n'est pas explicite ou apparente. Une grammaire est dite réursive gauche de façon indirecte s'il existe une règle de la forme :

$$A \rightarrow B\alpha \mid \beta \quad A, B \in \mathbb{N} \quad B \rightarrow^+ A\gamma \mid \Gamma \quad \alpha, \beta, \gamma, \Gamma \in (\mathbb{T} \cup \mathbb{N})^*$$

On peut définir de la même façon la récursivité à droite.

Certaines méthodes d'analyse syntaxique ne peuvent fonctionner avec une grammaire réursive à gauche (l'analyse descendante), ce qui nécessite une transformation. Le schéma de l'élimination de la RG repose sur l'introduction d'un non terminal intermédiaire qui prend en charge la récursivité mais à droite en commençant par générer la chaîne la plus à gauche.

$$A \rightarrow A\alpha \mid \beta \Leftrightarrow \begin{cases} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \varepsilon \end{cases}$$

Dans le cas général, on procédera de la même manière :

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_m \Leftrightarrow \begin{cases} A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A' \\ A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \varepsilon \end{cases}$$

**Exemple 10.** *Récursivité gauche directe :*

$$\begin{cases} S \rightarrow aS \mid b \mid Ac \\ A \rightarrow Aa \mid c \end{cases} \Rightarrow \begin{cases} S \rightarrow AS \mid b \mid Ac \\ A \rightarrow cA' \\ A' \rightarrow aA \mid \varepsilon \end{cases}$$

*RGI : Faire des substitutions de manière à faire apparaître une RGD puis appliquer les règles :*

$$\begin{cases} S \rightarrow Ac \mid b \\ A \rightarrow Bd \mid a \\ B \rightarrow Ac \mid d \end{cases} \Rightarrow \begin{cases} S \rightarrow Ac \mid b \\ A \rightarrow Acd \mid dd \mid a \\ B \rightarrow Ac \mid d \end{cases} \Rightarrow \begin{cases} S \rightarrow Ac \mid b \\ A \rightarrow Acd \mid dd \mid a \end{cases} \begin{cases} S \rightarrow Ac \mid b \\ A \rightarrow aA' \mid ddA' \\ A' \rightarrow cdA' \mid \varepsilon \end{cases}$$

### 3.2.2.3 Factorisation

L'analyseur syntaxique lit de gauche à droite la liste des unités lexicales fournies par l'analyseur lexical. Il est souhaitable qu'il puisse décider à la rencontre d'une

unité lexicale quelle production doit utiliser. Ceci est impossible si les membres droits des productions ont les mêmes préfixes.

L'idée de la factorisation est de réécrire la grammaire de sorte à différer l'indéterminisme jusqu'à avoir lu suffisamment d'unités lexicales de l'entrée pour pouvoir faire le bon choix.

**Exemple 11.**

$$I \rightarrow \text{if } C \text{ then } I \text{ else } I \\ | \text{if } C \text{ then } I$$

Ici. On préfère réécrire la grammaire :

$$I \rightarrow \text{if } C \text{ then } I T \\ T \rightarrow \text{else } I | \varepsilon$$

En général,

$$A \rightarrow aX | aY | aZ | \alpha | \beta \quad \Rightarrow \quad \begin{cases} A \rightarrow aB | \alpha | \beta \\ B \rightarrow X | Y | Z \end{cases}$$

**Exemple 12 (Factorisation).** Soit la grammaire

$$\begin{cases} A \rightarrow abC | aBd | aAD \\ B \rightarrow bB | \varepsilon \\ C \rightarrow d | \varepsilon \\ D \rightarrow a | b | \varepsilon \end{cases} \quad \Rightarrow \quad \begin{cases} A \rightarrow aX \\ X \rightarrow bC | Bd | AD \\ B \rightarrow bB | \varepsilon \\ C \rightarrow d | \varepsilon \\ D \rightarrow a | b | \varepsilon \end{cases} \quad \Rightarrow$$

$$\begin{cases} A \rightarrow aX \\ X \rightarrow bC | bBd | d | AD \\ B \rightarrow bB | \varepsilon \\ C \rightarrow d | \varepsilon \\ D \rightarrow a | b | \varepsilon \end{cases} \quad \Rightarrow \quad \begin{cases} A \rightarrow aX \\ X \rightarrow bY | d | AD \\ Y \rightarrow C | Bd \\ B \rightarrow bB | \varepsilon \\ C \rightarrow d | \varepsilon \\ D \rightarrow a | b | \varepsilon \end{cases}$$

### 3.2.2.4 Grammaire $\varepsilon$ -libre

Une grammaire est  $\varepsilon$ -libre si aucune production ne contient  $\varepsilon$  sauf l'axiome, qui ne doit pas apparaître en partie droite d'aucune production.

$$G = \begin{cases} S \rightarrow a \mid A \\ A \rightarrow AB \mid \varepsilon \\ B \rightarrow aAbB \mid b \end{cases} \Leftrightarrow \begin{cases} S \rightarrow a \mid \varepsilon \\ A \rightarrow AB \mid B \\ B \rightarrow aAbB \mid abB \mid b \end{cases}$$

Les productions vides posent souvent des difficultés dans l'analyse syntaxique ; les isoler est en général un remède aux difficultés qu'elles génèrent.

### 3.2.2.5 Grammaire sous la forme normale de Chomsky

Dans une grammaire sous cette forme normale, toutes les règles sont de la forme

$$\begin{cases} A \rightarrow BC & A, B, C \in \mathbb{N} \\ A \rightarrow a & a \in \mathbb{T} \\ S \rightarrow \varepsilon \end{cases}$$

Les arbres de dérivation des grammaires sous la forme normale de Chomsky sont binaires ce qui peut offrir une souplesse dans l'analyse.

### 3.2.2.6 Grammaire sous la forme normale de Greibach

Une grammaire est sous forme normale de Greibach si elle est  $\varepsilon$ -libre et toutes les règles sont de la forme :

$$A \rightarrow a\alpha \quad A \in \mathbb{N}, \alpha \in \mathbb{N}^*, a \in \mathbb{T} \text{ ou } S \rightarrow \varepsilon$$

Les productions dans une grammaire sous la forme normale de Greibach commencent toujours par un terminal, ce qui simplifie l'analyse.

## 3.3 Analyse syntaxique descendante

Dans cette section, nous décrivons les deux principales méthodes d'analyse syntaxique descendante à savoir : la descente récursive et LL. La première est un ensemble de modules mutuellement récursifs alors que la deuxième est itérative et exploite une table dite la table LL.

### 3.3.1 Analyse par la descente récursive

Dans cette méthode [1, 20], nous associons à chaque non terminal de la grammaire une procédure ou une fonction. L'analyse se fait à l'aide d'appels de procédures. L'analyseur syntaxique est donc un ensemble de procédures mutuellement récursives dont la structure rassemble à celle de la grammaire. Son exécution commence par le lancement de la procédure associée à l'axiome.

Dans la procédure d'un non terminal, nous vérifions l'entrée à la rencontre d'un terminal et appelons les procédures associées aux non terminaux selon leur apparitions dans les règles de productions. Si une discordance est trouvée nous levons une erreur.

---

#### Algorithme 5 Procédure Non-terminal $N()$

---

Choisir une production  $N \rightarrow X_1 X_2 \dots X_n$

**Pour**  $i$  de 1 à  $n$  **Faire**

**Si**  $X_i \in \mathbb{N}$  **Alors**

    appeler la procédure  $X_i()$

**Sinon Si**  $X_i =$  terme courant de la chaîne **Alors**

    avancer()

**Sinon** erreur

---

**Exemple 13.** Soit la grammaire augmentée ci-dessous :

$$G_1 = \begin{cases} Z \rightarrow S\# \\ S \rightarrow cAd \\ A \rightarrow ab|c \end{cases}$$

Analysons par cette méthode la chaîne :  $cabd\#$ . Le déroulement du code de la Figure 3.4 est montré dans la Table 3.1.

**Exemple 14.** Donner un analyseur en utilisant la descente récursive de la grammaire suivante et analyser la chaîne :  $bcaa\#$

$$G_2 : \begin{cases} Z \rightarrow S\# \\ S \rightarrow bSa \mid A \\ A \rightarrow cAa \mid aS \mid \varepsilon \end{cases}$$

```

1 Z(){          /* Programme principal de l'analyseur */
2   S()
3   si tc = '#' alors
4     Accepter()
5   sinon erreur()
6 }
7 S(){          /* procedure associee a S */
8   si tc='c' alors {
9     Avancer()
10    A()
11    si tc = 'd' alors
12      Avancer()
13    sinon erreur()
14   }
15   sinon erreur()
16 }
17 A(){         /* procedure associee a A */
18   si tc = 'a' alors {
19     Avancer()
20     si tc = 'b' alors
21       Avancer()
22     sinon erreur()
23   }
24   sinon si tc = 'c' alors
25     Avancer()
26     sinon erreur()
27 }

```

FIGURE 3.4 – Code de l'analyseur par descente récursive de  $G_1$ 

TABLE 3.1 – Exemple d'analyse par la descente récursive

Pile des appels	chaîne
Z	cabd#
ZS	abd#
ZSA	bd#
ZS	d#
Z	# ⇒ chaîne acceptée

### Remarques

1.  $Z$  joue le rôle du main de l'analyseur
2. *Avancer()* est la fonction permettant de lire (consommer) le token courant et passer au suivant
3.  $tc$  est le terminal ou le token courant, et  $\#$  la marque de fin de la chaîne

```
1 Z(){          /* Programme principal de l'analyseur */
2   S()
3   si tc = '#' alors
4     Accepter()
5   sinon erreur()
6 }
7 S(){          /* procedure associee a S */
8   si tc='b' alors {
9     Avancer()
10    S()
11    si tc = 'a' alors
12      Avancer()
13    sinon erreur() }
14  sinon A()
15 }
16 A(){          /* procedure associee a A */
17  si tc = 'c' alors {
18    Avancer()
19    A()
20    si tc = 'a' alors
21      Avancer()
22    sinon erreur() }
23  sinon si tc = 'a' alors {
24    Avancer()
25    S()
26  }
27  sinon ;      /* rien */
28 }
```

FIGURE 3.5 – Code de l'analyseur par descente récursive de  $G_2$ 

4. *erreur()* est la routine de gestion des erreurs
5. Méthode simple à implémenter et relativement efficace
6. Elle n'est pas générale, car liée à la grammaire
7. Pour faire une analyse performante déterministe il faut :
  - Éliminer la récursivité gauche afin d'éviter les boucles infinies
  - Factoriser
  - Tenir compte d'autres indéterminismes qui sont liés aux productions notamment celles générant la chaîne vide

TABLE 3.2 – Exemple non concluant d’une analyse par descente récursive d’une chaîne appartenant au langage de la grammaire

Pile des appels	chaîne
Z	bcaa#
ZS	bcaa#
ZSS	caa#
ZSSAA	aa#
ZSSAAS	a#
ZSSAASASA	#
ZSSAASAS	#
ZSSAASA	#
ZSSAASA	#
ZSSAAS	#
ZSSAA	#
ZSSA	# ⇒ erreur ? pourtant la chaîne appartient au langage

### 3.3.2 Analyse LL( $k$ )

Le problème de l’analyse descendante est comment déterminer les productions appropriées afin d’arriver à la chaîne en partant de l’axiome. Dans la méthode LL l’analyseur, en regardant  $k$  caractères de la chaîne, prédit par contre à chaque étape la bonne production à appliquer afin d’éviter les retours arrières.

Cette méthode prédictive et non récursive utilise une pile explicite et effectue la recherche dans une table dite **table d’analyse LL( $k$ )** [19, 26]. Souvent  $k = 1$  et l’analyse est dite **LL(1)**. Une grammaire est LL(1) si on arrive à accepter une chaîne en regardant un seul symbole de prévision.

LL( $k$ ) stands for **left to right scanning** performing a **leftmost derivation** using  $k$  tokens of lookahead.

### Fonctionnement

Nous disposons de la chaîne d’entrée terminée par # et une table d’analyse à deux dimensions. Initialement, la pile contient les symboles de la grammaire avec # au fond au dessous de l’axiome. Les lignes de la table pour les non terminaux, et les colonnes représentent les terminaux avec #. La case à la ligne de  $N$  et la colonne de

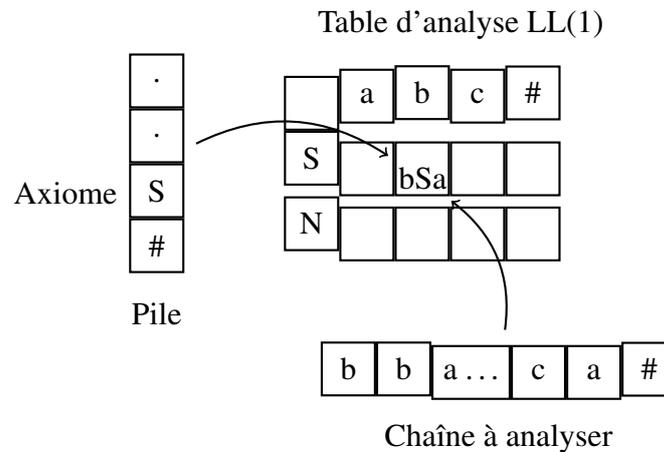


FIGURE 3.6 – Fonctionnement de l'analyse LL

$a$  donne la production à appliquer lorsque  $N$  est au sommet de la pile et  $a$  le terminal courant.

En général dans cette méthode, nous avons trois cas :

1. Si le sommet de la pile est un terminal qui coïncide avec le terminal courant : on dépile et on avance dans la chaîne
2. Si le sommet correspond à un non terminal, on consulte la table : si la case à la ligne de ce non terminal et la colonne du terminal courant est non vide : on dépile et on empile le miroir du membre droit de la production trouvée. Si l'entrée est vide, une erreur est signalée.
3. Si la lecture est terminée (terme courant est #), et le sommet est aussi # l'analyse est réussie et terminée.

### 3.3.2.1 Ensembles DÉBUT et SUIVANT

Afin de construire un analyseur descendant prédictif, on utilise un caractère d'anticipation (lookahead) pour décider de l'action adéquate. Les débuts ou les préfixes de tous les membres droits des production d'un non terminal doivent être disjoints. Par ailleurs, la grammaire doit être exempte de toute forme de récursivité à gauche afin d'éviter les boucles interminables. Nous définissons ici deux concepts qui aident dans la construction de la table d'analyse LL.

## DÉBUT

L'ensemble **DEB** constitue la collection de tous les terminaux qui peuvent commencer une chaîne ou ses dérivées. Formellement :

**Définition 6.**

$$DEB(X) = \{a \in \mathbb{T} \mid X \rightarrow^* a\alpha\} \cup \{\varepsilon \mid X \rightarrow^* \varepsilon\} \text{ Avec } \alpha \in (\mathbb{T} \cup \mathbb{N})^*$$

Pour calculer les ensembles **DEB** d'une chaîne  $X$ , appliquer les règles ci-après récursivement jusqu'à stabilité des ensembles.

---

**Algorithme 6** Calcul de l'ensemble  $DEB(X)$

---

1. Si  $X$  est un terminal ou commence par un terminal Ajouter ce terminal à  $DEB(X)$
  2. Si  $X$  est  $\varepsilon$  ou dérive directement ou indirectement la chaîne vide ajouter  $\varepsilon$  à  $DEB(X)$
  3. Si  $X$  est un non terminal et  $X \rightarrow Y_1 \dots Y_n$  une production de  $G$  avec  $Y_i \in (\mathbb{T} \cup \mathbb{N})$ ;  $i \in [1, n]$  :
    - Ajouter  $DEB(Y_1)$  à  $DEB(X)$
    - Ajouter  $DEB(Y_i)$  si  $\varepsilon \in DEB(Y_1), \dots, \varepsilon \in DEB(Y_{i-1})$
    - Ajouter  $\varepsilon$  à  $DEB(X)$  si  $\varepsilon \in DEB(Y_1), \dots, \varepsilon \in DEB(Y_n)$
- 

Il est aisé d'étendre cet ensemble à un préfixe de longueur au plus  $k$  (on dit un  $k$ -préfixe). Ainsi :

$$DEB_k(X) = \{u|_k \in \mathbb{T}^* \mid X \rightarrow^* u\alpha\} \cup \{\varepsilon \mid X \rightarrow^* \varepsilon\} \text{ Avec } \alpha \in (\mathbb{T} \cup \mathbb{N})^*$$

avec  $u|_k$  est le  $k$ -préfixe de  $u$ . Ce dernier est spécifié, pour un mot  $u = a_1 a_2 \dots a_n$  ( $a_i \in \mathbb{T}$ ) par la définition ci-dessous :

$$w|_k = \begin{cases} a_1 a_2 \dots a_n & \text{Si } n \leq k \\ a_1 a_2 \dots a_k & \text{Sinon} \end{cases}$$

## SUIVANT

De même, l'ensemble **SUIV** est la collection des terminaux qui peuvent suivre (apparaître immédiatement à droite) un symbole non terminal.

**Définition 7.**

$$SUIV(X) = \{a \in \mathbb{T} \mid S \rightarrow^* \beta X a \gamma\} \text{ Avec } \beta, \gamma \in (\mathbb{T} \cup \mathbb{N})^*, X \in \mathbb{N}$$

Suivre les règles suivantes pour calculer les ensembles SUIVANT de manière récursive jusqu'à stabilité des ensembles SUIVANT :

**Algorithme 7** Calcul de l'ensemble SUIV

1. Ajouter # à  $SUIV(S)$  ou  $S$  est l'axiome de la grammaire (découle de la production ajoutée  $Z \rightarrow S\#$ )
2. Si  $A \rightarrow \alpha B \beta$  : Ajouter  $DEB(\beta) \setminus \{\varepsilon\}$  à  $SUIV(B)$
3. si  $\begin{cases} A \rightarrow \alpha B \\ \text{ou} \\ A \rightarrow \alpha B \beta \wedge \beta \rightarrow^* \varepsilon \end{cases}$  : Ajouter les  $SUIV(A)$  dans  $SUIV(B)$

De manière analogue l'ensemble  $SUIV$  peut être étendu au suivants d'ordre  $k$  comme suit :

$$SUIV_k(X) = \{w \in \mathbb{T}^* \mid S \xrightarrow{\pm} \beta X \gamma \text{ et } w \in DEB_k(\gamma\#)\}$$

Dans les exemples suivants, nous calculons les ensembles DEB et SUIV.

**Exemple 15.**

$$G_3 : \begin{cases} S \rightarrow aAS \mid b \\ A \rightarrow a \mid bSA \end{cases}$$

	DEB	SUIV
$S$	$a, b$	$\#, a, b$
$A$	$a, b$	$a, b$

**Exemple 16.**

$$G_4 : \begin{cases} S \rightarrow d \mid XYS \\ Y \rightarrow c \mid \varepsilon \\ X \rightarrow a \mid Y \end{cases}$$

	DEB	SUIV
$S$	$a, c, d$	$\#$
$X$	$a, c, \varepsilon$	$a, c, d$
$Y$	$c, \varepsilon$	$a, c, d$

**Exemple 17.**

$$G_5 : \begin{cases} S \rightarrow ABCe \\ A \rightarrow aA \mid \varepsilon \\ B \rightarrow bB \mid cB \mid \varepsilon \\ C \rightarrow de \mid da \mid dA \end{cases}$$

	<i>DEB</i>	<i>SUIV</i>
<i>S</i>	<i>a,b,c,d</i>	<i>#</i>
<i>A</i>	<i>a,ε</i>	<i>b,c,d,e</i>
<i>B</i>	<i>b,c,ε</i>	<i>d</i>
<i>C</i>	<i>d</i>	<i>e</i>

**Exemple 18.**

$$G_6 : \begin{cases} S \rightarrow aSb \mid cd \mid SAe \\ A \rightarrow aAdB \mid \varepsilon \\ B \rightarrow bb \end{cases}$$

	<i>DEB</i>	<i>SUIV</i>
<i>S</i>	<i>a,c</i>	<i>a,b,e,#</i>
<i>A</i>	<i>a,ε</i>	<i>d,e</i>
<i>B</i>	<i>b</i>	<i>d,e</i>

**Exemple 19.**

$$G_7 : \begin{cases} E \rightarrow TX \\ X \rightarrow +E \mid \varepsilon \\ T \rightarrow iY \mid (E) \\ Y \rightarrow *T \mid \varepsilon \end{cases}$$

	<i>DEB</i>	<i>SUIV</i>
<i>E</i>	<i>(,i</i>	<i>#,)</i>
<i>X</i>	<i>+,ε</i>	<i>#,)</i>
<i>T</i>	<i>(,i</i>	<i>+,#,)</i>
<i>Y</i>	<i>*,ε</i>	<i>+,#,)</i>

**Exemple 20.**

$$G_8 : \begin{cases} S \rightarrow ABCD \\ A \rightarrow a \mid \varepsilon \\ B \rightarrow CD \mid b \\ C \rightarrow c \mid \varepsilon \\ D \rightarrow Aa \mid d \mid \varepsilon \end{cases}$$

	<i>DEB</i>	<i>SUIV</i>
<i>S</i>	<i>a,b,c,d,ε</i>	<i>#</i>
<i>A</i>	<i>a,ε</i>	<i>a,b,c,d,#</i>
<i>B</i>	<i>a,b,c,d,ε</i>	<i>a,c,d,#</i>
<i>C</i>	<i>c,ε</i>	<i>a,c,d,#</i>
<i>D</i>	<i>a,d,ε</i>	<i>a,c,d,#</i>

**Exemple 21.**

$$G_9 : \begin{cases} S \rightarrow ASSA \mid \varepsilon \\ A \rightarrow aSB \mid b \\ B \rightarrow bB \mid \varepsilon \end{cases}$$

	<i>DEB</i>	<i>SUIV</i>
<i>S</i>	<i>a,b,ε</i>	<i>a,b,#</i>
<i>A</i>	<i>a,b</i>	<i>a,b,#</i>
<i>B</i>	<i>b,ε</i>	<i>a,b,#</i>

**Exemple 22.**



	DEB	SUIV
S	a,b, $\varepsilon$	a,b,#
A	a,b	a,b,c
B	b,c, $\varepsilon$	a,b

Les résultats suivants donnent les conditions nécessaires et suffisantes pour qu'une grammaire soit  $LL(k)$ .

**Théorème 3.3.1.** Une grammaire  $G = \langle T, N, S, P \rangle$  est  $LL(k)$  si et seulement si pour toute paire de productions différentes de  $P$  :  $A \rightarrow \alpha$  et  $A \rightarrow \beta$  nous avons :

$$DEB_k(\alpha\gamma) \cap DEB_k(\beta\gamma) = \emptyset \text{ pour tout } \gamma \text{ tel que } S \xrightarrow{\pm} \omega A \gamma$$

**Théorème 3.3.2.** Une grammaire  $G = \langle T, N, S, P \rangle$  est  $LL(k)$  si et seulement si pour toute pair de productions différentes de  $P$  :  $A \rightarrow \alpha$  et  $A \rightarrow \beta$  nous avons :

$$DEB_1(\alpha) \odot_1 SUIV_1(A) \cap DEB_1(\beta) \odot_1 SUIV_1(A) = \emptyset$$

Où  $\odot_k$  est l'opération de concaténation à  $k$ -préfixe

**Corollaire 3.3.3.** Une grammaire est  $LL(1)$  si elle est non réursive à gauche et si pour toute production  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$ , nous avons pour tout  $i \neq j$  :

- $DEB(\alpha_i) \cap DEB(\alpha_j) \neq \emptyset$
- Si  $\alpha_i \rightarrow^* \varepsilon$  alors  $\alpha_j \not\rightarrow^* \varepsilon$
- Si  $\alpha_i \rightarrow^* \varepsilon$  alors  $DEB(\alpha_j) \cap SUIV(A) = \emptyset$

**Exemple 23.** La grammaire ci-dessous est-elle  $LL$

$$G_{11} : \begin{cases} S \rightarrow aAb \mid bS \\ A \rightarrow bB \mid \varepsilon \\ B \rightarrow b \end{cases}$$

$G_{11}$  n'est pas réursive gauche et elle est factorisée. Cependant les règles de  $A$  posent un problème. En effet l'alternative  $bB$  commence par un  $b$  qui est aussi un suivant de  $A$  ( $DEB(bB) \cap SUIV(A) \neq \emptyset$ ). Donc  $G$  n'est pas  $LL$ .

### 3.3.2.2 Construction de la table d'analyse LL

Pour l'élaboration de la table d'analyse  $LL$ , suivre le pseudo algorithme suivant. Les entrées vides de la table correspondent à des erreurs. La Table 3.3 montre celle

associée à la grammaire  $G_7$  de l'Exemple 18.

---

**Algorithme 8** Construction de la table LL
 

---

**Pour** chaque non terminal  $A$  de  $G$  **Faire**

**Pour** chaque production  $A \rightarrow \alpha$  **Faire**

**Pour** chaque  $a \in DEB(\alpha) \setminus \{\varepsilon\}$  **Faire**

$T[A,a] \cup = A \rightarrow \alpha$

**Si**  $\varepsilon \in DEB(\alpha)$  **Alors**

**Pour** chaque  $a \in SUIV(A)$  **Faire**

$T[A,a] \cup = A \rightarrow \alpha$

---

**Exemple 24.** Construire pour la grammaire  $G_7$  de l'exemple 19 la table d'analyse LL puis analyser la chaîne :  $i^*(i+i)\#$ .

TABLE 3.3 – Table LL de  $G_7$

	(	+	*	$i$	)	#
$E$	$TX$			$TX$		
$X$		$+E$			$\varepsilon$	$\varepsilon$
$T$	$(E)$			$iY$		
$Y$		$\varepsilon$	$*T$		$\varepsilon$	$\varepsilon$

TABLE 3.4 – Analyse LL de  $i * (i + i)$ 

Pile	Chaîne	Action
# $E$	$i * (i + i) \#$	Remplacer $E$ par $TX$
# $XT$	$i * (i + i) \#$	Remplacer $T$ par $iY$
# $XYi$	$i * (i + i) \#$	Dépiler et avancer
# $XY$	$*(i + i) \#$	Remplacer $Y$ par $*T$
# $XT*$	$*(i + i) \#$	Dépiler et avancer
# $XT$	$(i + i) \#$	Remplacer $T$ par $(E)$
# $X)E($	$(i + i) \#$	Dépiler et avancer
# $X)E$	$i + i) \#$	Remplacer $E$ par $TX$
# $X)XT$	$i + i) \#$	Remplacer $T$ par $iY$
# $X)XYi$	$i + i) \#$	Dépiler et avancer
# $X)XY$	$+i) \#$	Remplacer $Y$ par $\varepsilon$ (dépiler)
# $X)X$	$+i) \#$	Remplacer $X$ par $+E$
# $X)E+$	$+i) \#$	Dépiler et avancer
# $X)E$	$i) \#$	Remplacer $E$ par $TX$
# $X)XT$	$i) \#$	Remplacer $T$ par $iY$
# $X)XYi$	$i) \#$	Dépiler et avancer
# $X)XY$	)#	Remplacer $Y$ par $\varepsilon$ (dépiler)
# $X)X$	)#	Remplacer $X$ par $\varepsilon$ (dépiler)
# $X)$	)#	Dépiler et avancer
# $X$	#	Remplacer $X$ par $\varepsilon$ (dépiler)
#	#	Accepter et arrêter

**Exemple 25.** En utilisant la table d'analyse LL vérifier que la grammaire  $G_{12}$  sui-

$$\text{vante n'est pas LL. } G_{12} : \begin{cases} S \rightarrow iEtSR \mid a \\ R \rightarrow eS \mid \varepsilon \\ E \rightarrow b \end{cases}$$

TABLE 3.5 – Table LL de  $G_{12}$ 

	a	i	b	t	e	#
S	a	iEtSR	b	t	e	
R					eS,ε	ε
E		b				

La table d'analyse étant multi-définie (case  $T[R, e]$ ), cette grammaire n'est pas LL.

### 3.4 Analyse syntaxique ascendante

Bien que les méthodes d'analyse syntaxique descendantes sont simples et efficaces, il existe hélas des grammaires qui ne peuvent être analysées en utilisant ces méthodes (les grammaires récursives à gauche par exemple). Nous présentons ici des méthodes plus générales et puissantes qui procèdent de façons ascendante (des feuilles en remontant à la racine) dans la construction de l'arbre syntaxique d'un code source afin de vérifier sa syntaxe.

Les méthodes d'analyse syntaxiques ascendantes lisent toujours l'entrée de gauche à droite, elle cherchent, par contre, à trouver une suite de réductions valides qui permettent de transformer le texte source formés d'une séquences de tokens vers l'axiome de la grammaires décrivant le langage de programmation étudié. Ces méthodes dites  $LR(k)$  [16] (pour Left to right scanning performing the Reverse of the right most derivations using  $k$  symbols of lookahead) réalisent l'inverse d'une suite de dérivations la plus à droite et effectuent deux actions primitives : **Décaler** (empiler un symbole dans la pile) ou **Réduire** (remplacer un membre droit d'une production présent au sommet de la pile par le membre gauche correspondant) (Shift/reduce in english). Elles exploitent donc une pile et un automate déterministe (l'ensemble est en fait un automate à pile) dans la recherche de membres droits candidats (sous-chaînes susceptibles à conduire à une suite valide de réductions). Ces chaînes candi-

dates sont appelées manches ou handles. l'entier  $k$  désigne le nombre de caractères de prédiction à utiliser pour effectuer la bonne action. Ce nombre est rarement supérieur à 1, en raison de la complexité impliquée avec l'augmentation de la taille de cette information sur le fonctionnement de l'analyseur. Commençons par décrire la base de ces méthodes en ne regardant aucun symbole de prédiction, c-à-d avec  $k$  nul, après avoir examiné un exemple.

### Exemple introductif

Avant d'aller en profondeur, donnons un exemple illustratif de cette famille d'approches et le genre de décisions qu'elles sont confrontées à prendre.

#### Exemple 26.

$$G_{13} : \begin{cases} Z \rightarrow E \\ E \rightarrow E + T \mid T \\ T \rightarrow T * i \mid i \end{cases}$$

Prenons la chaîne de tokens  $i + i * i$  générée par la grammaire précédente  $G_{13}$  et essayons de montrer les étapes de déroulement de sa réduction en la lisant de gauche à droite. Testons deux scénarios différents afin de montrer la problématique de l'analyse ascendante. La barre verticale indique la position atteinte durant la lecture de la chaîne d'entrée.

1.  $|i + i * i \xRightarrow{\text{décaler}} i | + i * i \xRightarrow{\text{réduire}} T | + i * i \xRightarrow{\text{réduire}} E | + i * i \xRightarrow{\text{décaler}} E + | i * i \xRightarrow{\text{décaler}} E + i | * i \xRightarrow{\text{réduire}} E + T | * i \xRightarrow{\text{réduire}} E | * i \xRightarrow{\text{décaler}} E * | i \xRightarrow{\text{décaler}} E * i | \xRightarrow{\text{réduire}} E * T | \xrightarrow{?} \text{Blocage!}$

Cette suite d'actions n'a pas réussi à réduire une chaîne qui appartient bel et bien au langage de la grammaire. Réessayons d'une autre manière.

2.  $|i + i * i \xRightarrow{\text{décaler}} i | + i * i \xRightarrow{\text{réduire}} T | + i * i \xRightarrow{\text{réduire}} E | + i * i \xRightarrow{\text{décaler}} E + | i * i \xRightarrow{\text{décaler}} E + i | * i \xRightarrow{\text{réduire}} E + T | * i \xRightarrow{\text{décaler}} E + T * | i \xRightarrow{\text{décaler}} E + T * i | \xRightarrow{\text{réduire}} E + T | \xRightarrow{\text{réduire}} E |$  Une réduction réussie!

### 3.4.1 Analyse LR(0)

Une méthode *LR* doit implémenter un mécanisme pour reconnaître l'apparition de membres droits de production et de choisir la bonne action à effectuer (décaler /réduire) au bon moment. D. Knuth [16] remarqua que les préfixes possibles (appelés préfixes viables : ceux qui ne dépassent pas les handles) des membres droits

des productions d'une grammaire forment un langage régulier, ce qui signifie qu'il est possible de concevoir un AFD pour les reconnaître. Nous introduisons la notion d'item, qui nous aide à la construction de cet automate.

**Définition 8** (Manche (poignée ou handle)). *Pour une grammaire  $G$ , une dérivation la plus à droite  $Z \xRightarrow{+} \gamma'A\omega \Rightarrow \gamma\alpha\omega$  et une production  $A \rightarrow \alpha$ ,  $\alpha$  est la handle de la partie à droite  $\gamma\alpha\omega$*

**Définition 9** (Préfixe viable). *Un préfixe viable est un préfixe d'une chaîne résultat d'une dérivation la plus à droite qui n'excède pas (ne va pas au delà) le handle de cette chaîne.*

**Définition 10** (Item). *Un item pour une grammaire  $G$  est une production dont le membre droit est marqué par un point. Si  $A \rightarrow \alpha\beta \in P$  alors  $[A \rightarrow \alpha.\beta]$  est un item.*

Dans cette notation les crochets servent à distinguer les items et la règle de production associée. Généralement à une règle de production correspondent plusieurs items. Par exemple, quatre items sont associés à la production :  $E \rightarrow E - T$  qui sont :  $[E \rightarrow .E - T]$ ,  $[E \rightarrow E. - T]$ ,  $[E \rightarrow E - .T]$ ,  $[E \rightarrow E - T.]$ . La production  $A \rightarrow \epsilon$  quant à elle donne un seul item :  $[A \rightarrow .]$ . Un ensemble d'items est qualifié de final (terminal) s'il contient un item dans la marque termine le membre droit de l'item (un point à la fin).

L'intuition derrière un item est simple, à titre d'exemple pour l'item  $[E \rightarrow E. - T]$  : l'analyseur vient juste de reconnaître une chaîne dérivable de  $E$  et espère rencontrer le symbole  $-$  suivi d'une chaîne dérivable de  $T$ .

Afin de construire notre AFD, l'objectif est d'associer items et préfixes viables. Ainsi  $[A \rightarrow \alpha.\beta]$  est valide pour un préfixe viable  $\phi\alpha$  si et seulement si il existe une dérivation la plus à droite  $S \xRightarrow{+} \phi Aw \Rightarrow \phi\alpha\beta w$  où  $w$  est un mot terminal [32]. Nous donnons une construction d'un AFD à partir des règles de production de la grammaire qui reconnaît les préfixes viables. Les items seront regroupés en ensembles qui forment les états de cet AFD. Définissons les opérations suivantes pour achever cette construction : **fermeture** qui donne les états de l'automate et **successeur** qui en spécifie les transitions entre ces états.

### Fermeture d'un ensemble d'items

La fermeture d'un ensemble d'items  $I$  est l'ensemble d'items  $C(I)$  construit par application de l'algorithme suivant :

**Algorithme 9** Fermeture d'un ensemble d'items LR(0)**Entrée :**  $I$  : Un ensemble d'items LR(0)**Sortie :** La fermeture  $C(I)$  $C(I) \leftarrow I$ **Répéter****Pour tout** item  $[A \rightarrow \alpha.B\gamma] \in C(I)$  **Faire****Pour tout**  $B \rightarrow \beta \in P$  **Faire****Si**  $[B \rightarrow \cdot\beta] \notin C(I)$  **Alors** $C(I) \leftarrow C(I) \cup \{[B \rightarrow \cdot\beta]\}$ **Fin Si****Fin Pour****Fin Pour****Jusqu'à** Stabilité de  $C(I)$ **Retourner**  $C(I)$ **Successeur  $\delta$** 

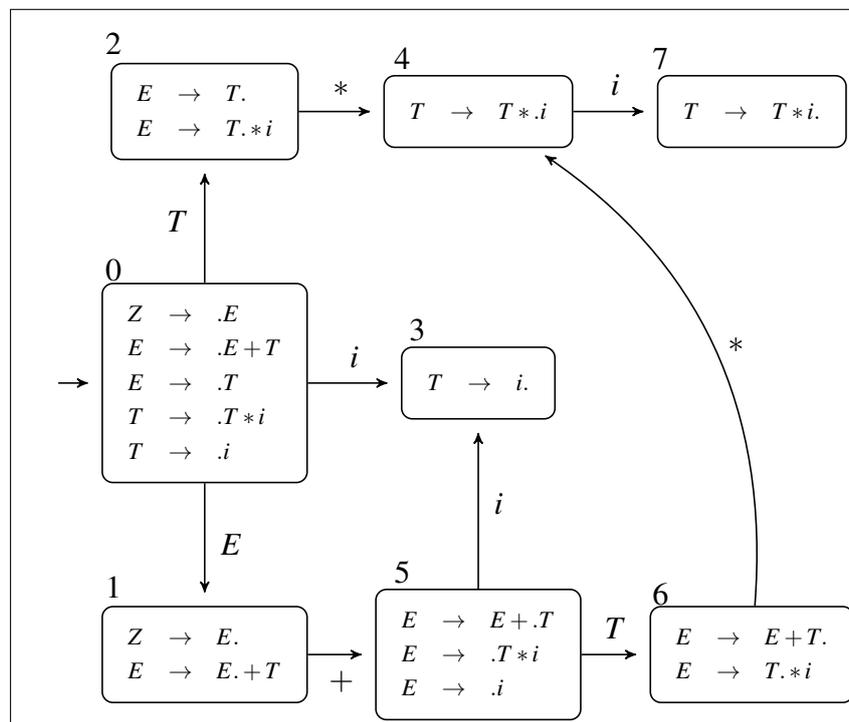
Pour tout  $X \in (T \cup N)$ , si  $[A \rightarrow \alpha.X\beta]$  est un item de l'ensemble  $I$  alors  $\Delta(I, X) = C(J)$  où  $J$  comprend tous les items  $[A \rightarrow \alpha X \cdot \beta]$

**Algorithme 10** Successeur d'un ensemble d'items LR(0)**Entrée :**  $I$  : Un ensemble d'items LR(0) et  $X \in (T \cup N \cup \{\#\})$ **Sortie :** La fonction  $\delta(I, X)$  $J \leftarrow \phi$ **Pour tout** item  $[A \rightarrow \alpha.X\beta] \in I$  **Faire** $J \leftarrow J \cup \{[A \rightarrow \alpha X \cdot \beta]\}$ **Fin Pour****Retourner**  $C(J)$ **Exemple 27.** Pour la grammaire  $G_{13}$ , nous avons : $I_0 = C(\{[Z \rightarrow \cdot E]\}) = \{[Z \rightarrow \cdot E], [E \rightarrow \cdot E + T], [E \rightarrow \cdot T], [T \rightarrow \cdot T * i], [T \rightarrow \cdot i]\}$  $I_1 = \delta(I_0, E) = \{[Z \rightarrow E \cdot], [E \rightarrow E \cdot + T]\}$  $\delta(I_1, +) = \{[E \rightarrow E + \cdot T], [T \rightarrow \cdot T * i], [T \rightarrow \cdot i]\}$ **Construction de l'automate LR(0)**

La décision (décaler/réduire) dans un analyseur  $LR$  peut parfois être non déterministe. On peut rencontrer des conflits, où on parle d'état ambiguë ou invalide. deux types de conflit peuvent surgir :

**Algorithme 11** Construction de l'automate LR(0)

1. Augmenter les règles de la grammaire par la production  $Z \rightarrow S$  où  $S$  est l'axiome
2. Générer l'ensemble initial d'items correspondant à  $C(\{[Z \rightarrow .S]\})$
3. Calculer les autres ensembles en utilisant la fonction  $\delta(I, X)$  avec  $I$  un nouveau ensemble obtenu et  $X \in (T \cup N)$  un symbole qui précède la marque dans les ensembles  $I$  trouvés.
4. Continuer tant qu'il y a de nouveaux ensembles d'items (états)
5. Tous les états sont finaux

FIGURE 3.7 – Automate LR(0) de  $G_{13}$

- décaler/réduire : si l'état contient un item final et un autre non final à transition par un symbole terminal
- réduire/réduire : si l'état inclut deux items finals différents.

**Lemme 3.4.1.** Une grammaire  $G$  est  $LR(0)$  si l'automate  $LR(0)$  associé ne contient aucun état ambiguë

### Une grammaire $LR(0)$

La grammaire  $G_{14}$  ci-dessous est  $LR(0)$ , son automate est montré dans la Figure 3.8 :

$$G_{14} : S \rightarrow aSb \mid c$$

La méthode  $LR(0)$  repose dans sa décision uniquement sur le contenu de la pile et n'utilise aucun token de prédiction. C'est une méthode faible car peu de grammaires sont  $LR(0)$ . Par exemple la grammaire  $G_{13}$  en haut n'est pas  $LR(0)$ , il inclut des états ambiguës (2 et 6). Voyons son automate.

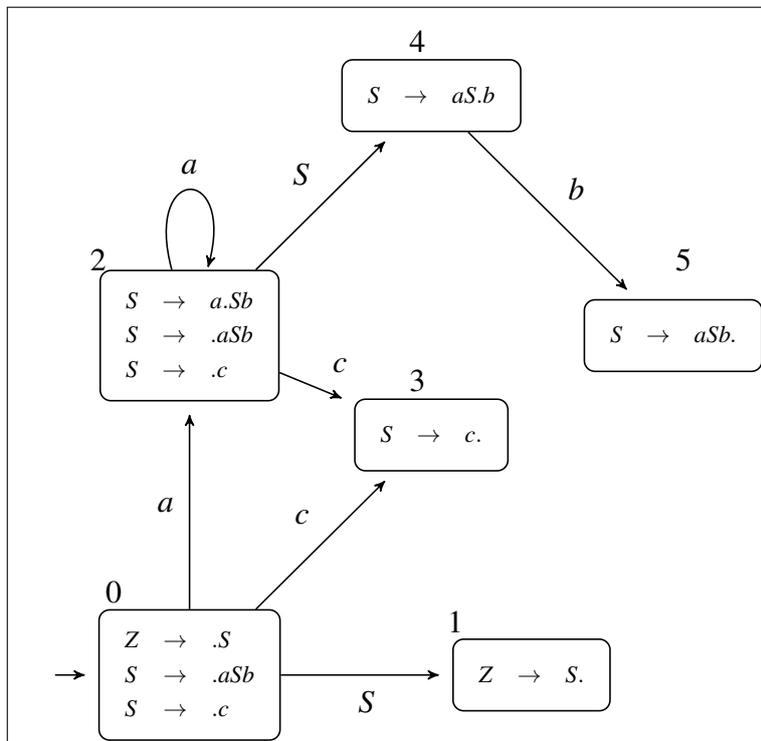


FIGURE 3.8 – Automate  $LR(0)$  de  $G_{14}$

TABLE 3.6 – Table LR(0) de  $G_{14}$ 

	a	b	c	#	S
0	$D_2$		$D_3$		1
1				Acc	
2	$D_2$		$D_3$		4
3	$R_{S \rightarrow c}$	$R_{S \rightarrow c}$	$R_{S \rightarrow c}$	$R_{S \rightarrow c}$	
4		$D_5$			
5	$R_{S \rightarrow aSb}$	$R_{S \rightarrow aSb}$	$R_{S \rightarrow aSb}$	$R_{S \rightarrow aSb}$	

### 3.4.2 Analyse SLR

SLR( for Simple LR) [9] est une amélioration légère sur la méthode précédente dans les cas de conflits décaler/réduire. L'heuristique stipule que devant un état ambiguë à cause d'un conflit décaler/réduire il faut choisir de réduire si le caractère de prédiction fait partie des suivants du non terminal de l'item final, sinon il faut opter pour le décalage. Si malgré cette astuce le conflit persiste, la grammaire est donc non SLR et le recours à d'autres méthodes plus puissantes s'impose alors.

#### Une grammaire non LR(0) mais SLR

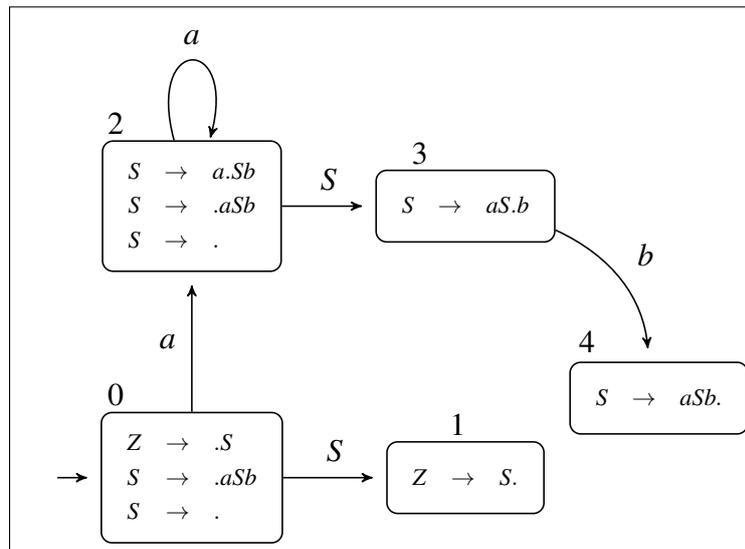
**Exemple 28.** La grammaire  $G_{15}$  présente des conflits décaler/réduire dans deux états.

$$G_{15} : \rightarrow aSb \mid \varepsilon$$

Comme  $\text{Suiv}(S) = \{b, \#\}$ , l'ambiguïté est levée : On décale à la rencontre d'un  $a$  et on réduit à la lecture de  $b$  ou  $\#$  et donc  $G_{15}$  est SLR.

#### Construction de la table d'analyse SLR

Cette table à deux dimensions consigne les actions (décaler, réduire, accepter ou erreur). Elle possède des lignes au nombre des états de l'automate ; les symboles de  $T \cup N \cup \{\#\}$  représentent quant à eux ses colonnes. Cette table est scindée en deux parties : une partie pour les actions de l'analyseur et la deuxième aux différents transitions de l'automates entre ses états. Ci-après les étapes de sa construction [12].

FIGURE 3.9 – Automate LR(0) de  $G_{15}$ **Algorithme 12** Construction de la table d'analyse SLR

- Calculer l'ensemble des items LR(0)
- L'état  $i$  est construit à partir de l'ensemble d'items  $I_i$ . Ses actions sont :
  1. Si  $[A \rightarrow \alpha.a\beta] \in I_i$  et  $\delta(I_i, a) = I_j$  pour  $a \in T$  alors  $Tab[i, a] = D_j$
  2. Si  $[A \rightarrow \alpha.] \in I_i$  pour tout  $c \in SUIV(A)$  et  $A \neq Z$  alors  $Tab[i, c] = R_{A \rightarrow \alpha}$
  3. Si  $[Z \rightarrow S.] \in I_i$  alors  $Tab[i, \#] = Accepter$
- Si  $\delta(I_i, A) = I_j$  alors :  $Tab[i, A] = j$
- L'état initial est celui qui contient l'item  $[Z \rightarrow .S]$
- Les entrées vides de la table correspondent à des erreurs.

Pour une grammaire  $G$ , si cette construction donne une table mono-définie (chaque case contient au plus une action),  $G$  est donc SLR et inversement. La table 3.7 montre la table SLR de la grammaire  $G_{15}$  qui remarquons-le est mono-définie.

**Une grammaire non SLR**

**Exemple 29.** La grammaire ci-dessous n'est pas SLR. Son automate comprend un état ambiguë ayant un conflit décaler/réduire. L'analyseur ne sait pas comment agir à la lecture d'un  $b$ , car l'état a un item de réduction  $[A \rightarrow c.]$  et un autre de décalage  $[S \rightarrow c.b]$ . L'heuristique SLR ne résout pas ce conflit étant donné que les suivants

TABLE 3.7 – Table SLR de  $G_{15}$

	a	b	#	S
0	$D_2$	$R_{S \rightarrow \epsilon}$	$R_{S \rightarrow \epsilon}$	1
1			Accepter	
2	$D_2$	$R_{S \rightarrow \epsilon}$	$R_{S \rightarrow \epsilon}$	3
3		$D_4$		
4		$R_{S \rightarrow aSb}$	$R_{S \rightarrow aSb}$	

de  $A$  incluent  $b$ .

$$G_{16} : \begin{cases} Z \rightarrow S \\ S \rightarrow A \mid cb \\ A \rightarrow aAb \mid c \end{cases}$$

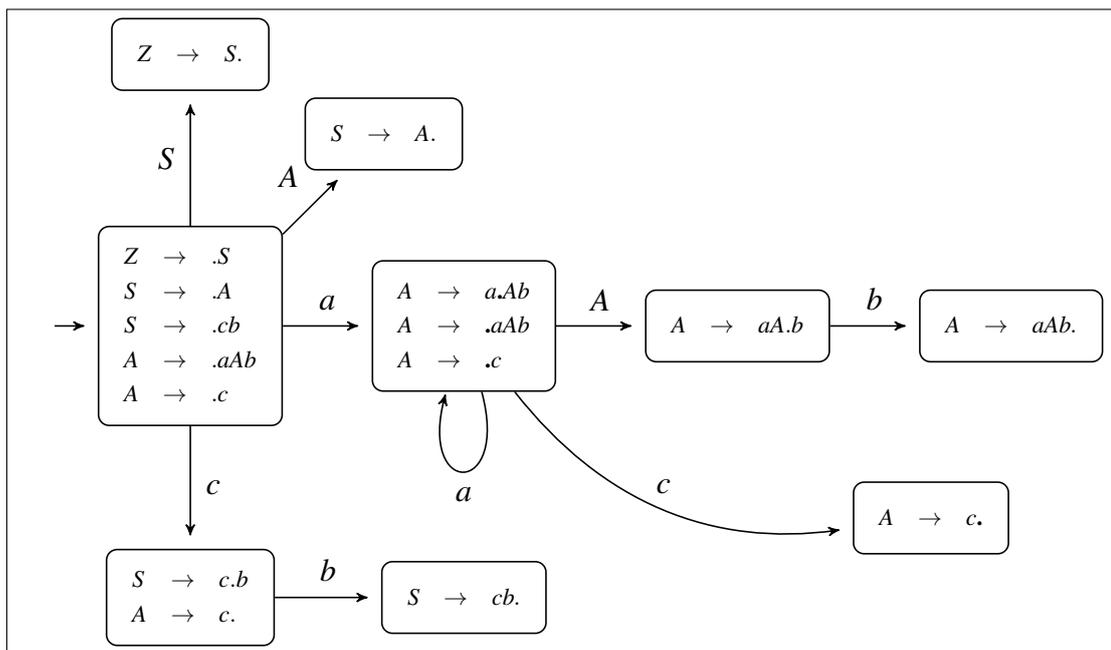


FIGURE 3.10 – Automate LR(0) de  $G_{16}$

### 3.4.3 Analyse LR(1) canonique

Pour doter les méthodes LR d'avantage de puissance, on étend l'ensemble des items LR(0) par l'information de prévision (le caractère lookahead). Un item LR(1)

inclut désormais deux composants : l'item de la règle comme dans LR(0) et le symbole de prévision.

**Définition 11.** *Un item LR(1) est un objet à deux composants : une production marquée et un caractère de prévision  $[A \rightarrow \alpha.\beta, a]$  où  $a \in T \cup \{\#\}$*

Un analyseur LR(1) (automate et table) peut être élaboré de manière similaire à celle de l'analyseur SLR. Il est cependant réputé plus efficace même s'il implique une complexité supérieure. Revenons sur les deux opérations de base **Fermeture** et **Successeur**.

---

**Algorithme 13** Fermeture d'un ensemble d'items LR(1)

---

**Entrée :**  $I$  : Un ensemble d'items LR(1)

**Sortie :** La fermeture  $C(I)$

$C(I) \leftarrow I$

**Répéter**

**Pour tout** item  $[A \rightarrow \alpha.B\beta, a] \in I$  **Faire**

**Pour tout**  $B \rightarrow \gamma \in P$  **Faire**

**Pour tout**  $b \in DEB(\beta a)$  **Faire**

$C(I) \leftarrow C(I) \cup \{[B \rightarrow \gamma, b]\}$

**Fin Pour**

**Fin Pour**

**Fin Pour**

**Jusqu'à** Stabilité de  $C(I)$

**Retourner**  $C(I)$

---



---

**Algorithme 14** Successeur d'un ensemble d'items LR(1)

---

**Entrée :**  $I$  : Un ensemble d'items LR(1) et un symbole  $X \in (T \cup N \cup \{\#\})$

**Sortie :** La fonction  $\delta(I, X)$

$J \leftarrow \phi$

**Pour tout** item  $[A \rightarrow \alpha.X\beta, a] \in I$  **Faire**

$J \leftarrow J \cup \{[A \rightarrow \alpha.X.\beta, a]\}$

**Fin Pour**

**Retourner**  $C(J)$

---

Le principe reste identique, il faut seulement prendre en considération le symbole de prévision. Ci-dessous la procédure pour le calcul des items LR(1).

Voici comment élaborer la table d'analyse LR(1).

Les entrées non définies après application de l'algorithme précédent correspondent à des erreurs. L'état initial  $I_0$  est celui contenant l'item  $[Z \rightarrow \cdot S, \#]$ . Dans

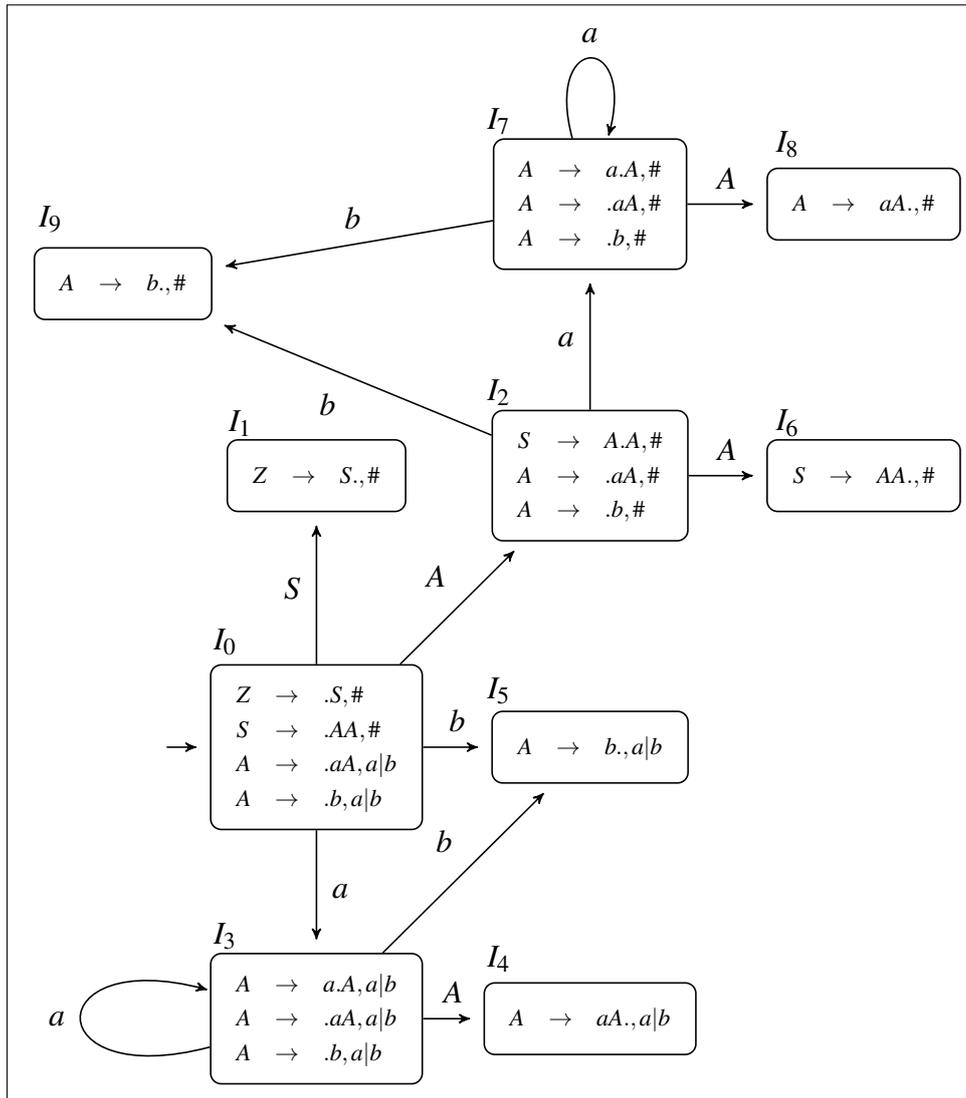
**Algorithme 15** Calcul de la collection des items LR(1)**Entrée :**  $\hat{G}$  : Une grammaire hors contexte augmentée**Sortie :**  $\mathcal{L}$  : La collection des items LR(1) $\mathcal{L} \leftarrow \{[Z \rightarrow .S, \#]\}$ **Répéter****Pour tout** item  $I \in \mathcal{L}$  **Faire****Pour tout**  $X \in T \cup N \cup \{\#\}$  **Faire****Si**  $\delta(I, X) \neq \emptyset \wedge \notin \mathcal{L}$  **Alors** $\mathcal{L} \leftarrow \mathcal{L} \cup \delta(I, X)$ **Fin Si****Fin Pour****Fin Pour****Jusqu'à** Stabilité de  $\mathcal{L}$ **Algorithme 16** Construction de la Table LR(1)**Entrée :** La collection des items LR(1)**Sortie :** La table LR(1) renseignée**Si**  $[A \rightarrow \alpha.a\beta, b] \in I_i \wedge \delta(I_i, a) = I_j$  **Alors** $Tab[i, a] \leftarrow D_j$ **Fin Si****Si**  $[A \rightarrow \alpha., a] \in I_i$  **Alors** $Tab[i, a] \leftarrow R_{A \rightarrow \alpha}$  Avec  $A \neq Z$ **Fin Si****Si**  $[Z \rightarrow S., \#] \in I_i$  **Alors** $Tab[i, \#] \leftarrow \text{Accept}$ **Fin Si****Si**  $\delta(I_i, A) = I_j$  **Alors** $Tab[i, A] \leftarrow j$ **Fin Si**

ce qui suit nous prenons un exemple d'une grammaire LR(1), nous élaborons l'automate et la table LR(1) et appliquons l'analyse sur une chaîne.

**Exemple 30.** Soit la grammaire

$$G_{17} : \begin{cases} Z \rightarrow S \\ S \rightarrow AA \\ A \rightarrow aA \mid b \end{cases}$$

Ci-après la collection des items LR(1) et la table associée.

FIGURE 3.11 – Automate LR(1) de  $G_{17}$ 

Pour illustrer le fonctionnement de l'analyseur examinons la chaîne :  $abb$ .

TABLE 3.8 – Table LR(1) de  $G_{17}$ 

	a	b	#	S	A
0	$D_3$	$D_5$		1	2
1			Accepter		
2	$D_7$	$D_9$			6
3	$D_3$	$D_5$			4
4	$R_{A \rightarrow aA}$	$R_{A \rightarrow aA}$			
5	$R_{A \rightarrow b}$	$R_{A \rightarrow b}$			
6			$R_{S \rightarrow AA}$		
7	$D_7$	$D_9$			8
8			$R_{A \rightarrow aA}$		
9			$R_{A \rightarrow b}$		

TABLE 3.9 – Déroulement de l'analyse LR(1) de  $G_{17}$  sur la chaîne  $abb$ 

Pile	Chaîne	Action
0	$abb\#$	Décaler et passer à 3
0a3	$bb\#$	Décaler et passer à 5
0a3b5	$b\#$	Réduire $b$ en $A$
0a3A	$b\#$	Passer à 4
0a3A4	$b\#$	Réduire $aA$ en $A$
0A	$b\#$	Passer à 2
0A2	$b\#$	Décaler et passer à 9
0A2b9	$\#$	Réduire $b$ en $A$
0A2A	$\#$	Passer à 6
0A2A6	$\#$	Réduire $AA$ en $S$
0S	$\#$	Passer à 1
0S1	$\#$	Accepter

### 3.4.4 Analyse LALR(1)

LALR [10] (pour Look-Ahead LR) est la méthode la plus utilisée en pratique. En effet, la méthode LR(1) bien que puissante, elle implique un calcul prohibitif et requière un espace mémoire très important. LALR a été développée comme un compromis entre simplicité de SLR et puissance de LR(1). En effet, pour toute grammaire le nombre des états de table LALR associée est identique à celui de la table SLR correspondantes.

La table d'analyse LALR(1) est obtenue directement à partir des items LR(1) par fusion des états ayant la même partie gauche (cœur). Si cette astuce fournit une table mono-définie, la grammaire est donc LALR. Dans le cas contraire elle est LR(1), mais non LALR(1).

L'application de cette astuce à notre dernier exemple de la grammaire  $G_{17}$  permet de réduire les états de 10 à 7. Il est aisé de voir que les états 3 et 7, 5 et 9, 4 et 8 peuvent être fusionnés en 3 états seulement. La table d'analyse 3.10 demeure mono-définie ce qui signifie qu'on peut l'analyser par la méthode LALR.

TABLE 3.10 – Table LALR(1) de  $G_{17}$

	a	b	#	S	A
0	$D_{37}$	$D_{59}$		1	2
1			Accepter		
2	$D_{37}$	$D_{59}$			6
37	$D_{37}$	$D_{59}$			48
48	$R_{A \rightarrow aA}$	$R_{A \rightarrow aA}$	$R_{A \rightarrow aA}$		
59	$R_{A \rightarrow b}$	$R_{A \rightarrow b}$	$R_{A \rightarrow b}$		
6			$R_{S \rightarrow AA}$		

## 3.5 Générateurs d'analyseurs syntaxiques cas de Bison

À l'instar des générateurs de scanners, une panoplie de générateurs d'analyseurs syntaxiques sont développés en utilisant différents langages de programmation et plateformes. Ces outils reçoivent la spécification syntaxique du langage à implémenter et produisent le code du parser correspondant. Certains de ces outils adoptent

l'approche descendante (JavaCC [31], Coco/R [11], etc. ), d'autre reposent sur la méthode LALR et sont conformes de ce fait à l'approche ascendante (Bison [18], PLY [4], CUP [3], etc.). Dans ce qui suit nous présentons l'essentiel pour exploiter Bison afin de générer un analyseur syntaxique.

Bison [18] est une version améliorée et libre du fameux générateur Yacc d'Unix [13]. C'est un générateur d'analyseurs syntaxiques qui reçoit la spécification syntaxique et produit le code source de l'analyseur en langage C. L'entrée de Bison est la spécification syntaxique du langage c-à-d une grammaire à contexte libre LALR en format BNF facile à lire par un programme. Bison compile les règles et les déclarations dans le fichier d'entrée afin de produire la table d'analyse LALR codée en langage C sous la forme d'instructions switch dans la routine `yyparse()`.

Similairement à Flex, un fichier (d'extension `.y` par convention) d'entrée à Bison inclut, comme illustre la Figure 3.12, quatre parties : une pour les déclarations globales, une deuxième pour les déclarations propres à Bison, la partie règles de production de la grammaire et en fin un code supplémentaire.

```
1  %{
2      Declaration C
3  %}
4      Declaration Bison
5  %%
6      Regles de la grammaire
7  %%
8      Code C additionnel
```

FIGURE 3.12 – Structure d'un fichier de spécification Bison

### 3.5.1 Partie déclaration C

Cette partie facultative et délimitée par la couple `%{` et `%}` englobe les macros, les includes, les définitions et les déclarations des variables et fonctions utilisées dans la parties actions des règles de la grammaire.

### 3.5.2 Partie déclaration Bison

Ici on définit les symboles de la grammaire (terminaux et non terminaux), les priorités et précédences des opérateurs et les types de données pour les actions. Ci-dessous sont résumés les directives de Bison.

### 3.5.3 Partie règles de la grammaire

Cette section contient exclusivement les règles de la grammaire sous la forme :  $A : MDP$  où  $A$  est un non terminal. Si la règle possède de multiples alternatives, elles seront séparées par une barre verticale. L'ensemble des règles est terminé par un point virgule. Voici un exemple :

```
exp  :  exp PLUS exp
      |  exp MOINS exp
      |  exp FOIS exp
      ;
```

Pour assurer une sémantique au langage analysé, des actions sémantiques écrites en langage C peuvent être insérées après chaque règle de production. Un exemple serait d'ajouter les valeurs associées aux symboles pour la première règle.

```
exp  :  exp PLUS exp    {$$ = $1 + $3 ;}
```

### 3.5.4 Partie code additionnel

Cette dernière est optionnelle comme la première et seront copiées (une à la fin et l'autre au début) telle quelles dans le code source de l'analyseur. On peut y ajouter tout ce que nous voulons accomplir en plus du code du parser proprement dit.

Voici ci-après un résumé des principales directives de Bison.

- Les unités lexicales seront déclarées par la directive **%token nom\_de\_l'unité**
- Les directives **%left (resp. %right)** permettent de définir une règle d'associativité gauche (resp. droite), par exemple **%left +** : signifie que l'expression  $a+b+c$  sera évaluée :  $(a+b)+c$ . Noter qu'un token déclaré par **%left (resp. %right)** n'a pas besoin d'être introduit par une directive **%token**. **%nonassoc** : pour les opérateurs non associatifs, et **%start** pour préciser l'axiome, sinon le premier dans la partie règles sera considéré comme tel.
- L'ordre de précedence est déterminé par l'ordre d'apparition des tokens (les premiers listés ont la priorité la plus faible). Si plusieurs tokens possèdent le même ordre de précedence, ils seront déclarés en même temps sur la même ligne.

- On introduit la grammaire dans la partie délimitée par `%%`. Les productions d'un non terminal seront séparées sur différentes lignes par `|` comme dans la notation BNF. Une ligne vide remplace l'alternative vide ( $\epsilon$ )
- Si une règle de production a été reconnue l'action associée sera exécutée. Comme en Flex, une action est un bloc de code en C. les  $\$n$  désignent les numéros des symboles comme ils apparaissent dans le MDP et `$$` fait référence au MGP. Exemple :  $E \rightarrow E + E$  `{$$ = $1 + $3;}`.
- La fonction `main()` de l'outil appelle `yyparse()` pour traiter l'entrée, qui appelle à son tour la fonction `yylex()` afin de récupérer les tokens de l'analyseur lexical. Si une erreur est rencontrée, l'analyseur appelle `yyerror(m)`, où `m` est le message d'erreur devant être renvoyé.
- La déclaration `%union` définit la collection des types de données utilisés dans les actions de l'analyseur. Par exemple la déclaration suivante indique que deux types (un réel et une chaîne seront exploités).

```

1  \%union {
2      double val;
3      char* nom;
4  }
```

Ces types peuvent être utilisés avec les directives `%token` et `%type` afin d'indiquer les types pour les symboles de la grammaire.

Voici comme résumé un exemple complet d'un code Bison pour la grammaire des expressions arithmétiques :

## 3.6 Gestion des erreurs

Différentes sortes d'erreurs peuvent être commises dans les codes sources. Celles-ci peuvent s'intercaler dans les différentes phases d'analyse. Quelque soit l'erreur, le compilateur doit la prendre en charge efficacement. Ainsi la gestion des erreurs inclut :

- La détection rapide et précise de l'erreur (spécifier sa nature et indiquer la position exacte dans le texte source)
- Le traitement efficace de l'erreur de façon à ne pas compromettre le processus d'analyse

```

1  %{ #include <stdio.h>
2     #include<string.h>
3     #include "calc.h"
4     int yylex(void);
5     void yyerror(char *);
6  }
7  %union { double fl;
8         char car; }
9  %type <fl> expr exp condition
10 %token <fl> NB
11 %token <car> var egal diff et ou si alors
12 %token '<'
13 %token '>'
14 %token sortie
15 %token '(' ')'
16 %left '+'
17 %left '*'
18 %right '='
19 %nonassoc na
20 %%
21 prg: /* rien (vide) */
22 expr';' '\n' { printf("%f\n> ", $1); }
23 | si '(' condition ')' alors expr {if ($3) { $1;}}
24 | var '=' expr { inserer(symbole, $1 , $3);}
25 | sortie { printf("> Au revoir !");exit(0);}
26 ;
27 expr : NB { $$ = $1 ; }
28 | var { $$ = recuperer(symbole, $1);}
29 | expr '+' expr { $$ = $1 + $3; }
30 | expr '*' expr { $$ = $1 * $3; }
31 | '-' expr %prec na { $$ = - $2; }
32 | '(' expr ')' { $$ = $2; }
33 ;
34 condition : exp ou exp { $$ = $1 && $3;}
35 | exp et exp { $$ = $1 && $3;}
36 ;
37 exp : expr '<' expr { $$ = $1 < $3; }
38 | expr '>' expr { $$ = $1 > $3; }
39 | expr "==" expr { $$ = $1 == $3; }
40 | expr "!=" expr { $$ = $1 != $3; }
41 ;
42 %%
43 void yyerror(char *s) {
44     fprintf(stderr, "%s\n> ", s); }
45 void inserer (struct table **t, char c, double val){
46     /* inserer une variable dans une llc */
47     while (t != NULL){
48         ...
49     }
50 double recuperer (struct table **t, char c){
51     /* rechercher une variable et recuperer sa valeur*/
52     ....
53     return val;
54 }
55 int main(void) {
56     printf("> "); *symbole = NULL; yyparse(); return 0;
57 }

```

FIGURE 3.13 – Code Bison complet

Les erreurs communes se répartissent en quatre types :

- **Lexicales** comme la malformation des unités lexicales (mots-clés, identificateurs, etc.) qui résultent souvent des coquilles de frappes
- **Syntaxiques** : le non respect de la syntaxe du langage, comme les séparateurs déplacés ou oubliés, les parenthèses non équilibrées ; ...
- **Sémantiques** qui n'observent pas à titre d'exemple le systèmes de types, les déclarations manquantes, etc.
- **Logiques** représentant les erreurs de raisonnement de la part du programmeur comme les boucles infinies

Généralement, les deux premières classes d'erreurs et une partie de la troisième peuvent être détectées sans failles. Cela est dû aux outils des analyseurs (automates, tables, préfixes viables, etc.) qui offrent une détection quasi-instantanée. Cependant, beaucoup d'erreurs sémantiques et la totalité des erreurs logiques sont loin d'être maîtrisées par les compilateurs en raison de leur nature extrêmement complexe.

Plusieurs approches ont été adoptées pour traiter les erreurs. Le principe est surtout la correction et l'efficacité. Les plus connues des techniques de récupération suite à une erreur sont : le mode panique, le mode au niveau du syntagme et la production des erreurs [12].

### 3.6.1 Récupération en mode panique

Cette technique est la plus simple des stratégiques. À la rencontre d'une erreur, le compilateur supprime généralement plusieurs tokens jusqu'à la détection d'un qui appartient à un ensemble dit de synchronisation ce qui permet de reprendre l'analyse promptement. Les symboles de synchronisation englobent par exemple les séparateurs comme les point-virgules ou les parenthèses, les délimiteurs de blocs, etc.

### 3.6.2 Récupération au niveau du syntagme

À l'opposé de la précédente stratégie, ici le compilateur n'élimine pas des entités mais plutôt propose des corrections locales du genre remplacer un , par un ; inverser fi à if, etc. Bien entendu, les rectifications introduites ne doivent pas générer d'autres situations ingérables ou faire coïncider l'analyse dans une boucle infinie par exemple.

### 3.6.3 Production d'erreurs

Ici l'approche est différente. La grammaire est augmentée dès la conception par des constructions qui gèrent les erreurs prévisibles. Cette manière d'anticiper les erreurs susceptibles de surgir offre un traitement plus lisible et compréhensible à l'utilisateur.

## **Exercices du chapitre 3**

Ministère de l'Enseignement Supérieur et de la Recherche Scientifique  
 Université de Ghardaia  
 Faculté des Sciences et de la Technologie  
 Département des Mathématiques et d'Informatique

## ANALYSE SYNTAXIQUE

### Exercice 1

Quel est le nombre de mots et d'arbres syntaxiques distincts générés par la grammaire ci-

$$\text{contre : } G_1 : \begin{cases} S \rightarrow A1 \mid 1B \\ A \rightarrow 10 \mid C \mid \varepsilon \\ B \rightarrow C1 \mid \varepsilon \\ C \rightarrow 0 \mid 1 \end{cases}$$

### Exercice 2

Ci-après une grammaire des expressions lisp simplifiées :

$$G_2 : \begin{cases} L \rightarrow A \mid M \\ A \rightarrow nb \mid id \\ M \rightarrow (N) \\ N \rightarrow N L \mid L \end{cases}$$

- (a) Écrire la dérivation la plus à gauche/à droite de la chaîne :  $(a \ 23 \ (m \ x \ y))$   
 (b) Donner son arbre syntaxique

### Exercice 3

Éliminer la récursivité gauche dans les grammaires suivantes :

$$G_3 : \begin{cases} S \rightarrow S \setminus S \mid SS \mid A \\ A \rightarrow S^* \mid (S) \mid a \mid b \end{cases} \quad G_4 : \begin{cases} S \rightarrow AA \mid 0 \\ A \rightarrow SS \mid 1 \end{cases}$$

$$G_5 : \begin{cases} S \rightarrow AS \mid b \\ A \rightarrow SA \mid a \end{cases} \quad G_6 : \begin{cases} A \rightarrow Cd \\ B \rightarrow Ce \\ C \rightarrow A \mid B \mid f \end{cases}$$

### Exercice 4

Rendre la grammaire suivante  $\varepsilon$ -libre :

$$G_7 = \begin{cases} S \rightarrow AB \mid AaB \mid aS \\ A \rightarrow aA \mid \varepsilon \\ B \rightarrow bB \mid \varepsilon \end{cases} \quad G_8 = \begin{cases} S \rightarrow A \mid B \\ A \rightarrow aB \mid bS \mid \varepsilon \\ B \rightarrow AB \mid BCa \\ C \rightarrow AS \mid \varepsilon \end{cases}$$

### Exercice 5

Transformer les grammaires ci-dessous sous FNG puis sous FNC :

$$G_9 = \begin{cases} S \rightarrow AaA \mid Sa \\ A \rightarrow Sb \mid \varepsilon \end{cases} \quad G_{10} = \begin{cases} S \rightarrow AS \mid a \\ A \rightarrow Sb \mid Aa \end{cases}$$

### Exercice 6

Soit la grammaire :  $G_{11} : S \rightarrow 0S1S \mid 2S \mid \varepsilon$

- (a) Écrire un analyseur syntaxique pour  $G_{11}$  en utilisant la descente récursive.  
 (b) Analyser la chaîne : 012

**Exercice 7**

Soit la grammaire  $G_{12}$  ci-contre.

$$G_1 : \begin{cases} S \rightarrow aB \mid bA \mid c \\ A \rightarrow aS \mid bA \\ B \rightarrow b \end{cases}$$

- (a) Montrer que  $G_{12}$  est LL(1)  
 (b) Écrire un analyseur syntaxique pour  $G_{12}$  en utilisant la descente récursive.  
 (c) Analyser la chaîne : *babbac*

**Exercice 8**

Soit la grammaire :  $G_{13} : \begin{cases} S \rightarrow a \mid b \mid (T) \\ T \rightarrow T, S \mid S \end{cases}$

- (a) Éliminer la récursivité à gauche. Montrer que la grammaire transformée est LL  
 (b) Donner la table d'analyse LL, et analyser la chaîne :  $(a, (b, a), a)$

**Exercice 9**

Soit la grammaire :  $G_{14} : \begin{cases} S \rightarrow (L) \mid a \\ L \rightarrow SM \mid \varepsilon \\ M \rightarrow ;SM \mid \varepsilon \end{cases}$

- (a) Cette grammaire est-elle LL; dessiner sa table d'analyse  
 (b) Analyser la chaîne :  $(a; a)$

**Exercice 10**

Les grammaires suivantes sont-elles LL(k) ?

$$G_{15} = \begin{cases} S \rightarrow AB \mid aSb \mid CSB \\ A \rightarrow bA \mid \varepsilon \\ B \rightarrow dB \mid \varepsilon \\ C \rightarrow cC \mid e \end{cases} \quad G_{16} = \begin{cases} S \rightarrow ASB \mid \varepsilon \\ A \rightarrow aAc \mid c \\ B \rightarrow bBA \mid \varepsilon \end{cases} \quad G_{17} = \begin{cases} S \rightarrow A \mid B \mid a \\ A \rightarrow bBA \\ B \rightarrow dB \mid \varepsilon \end{cases}$$

**Exercice 11**

On s'intéresse à la construction d'un analyseur syntaxique pour la grammaire  $G_{18}$  ci-après,

où  $\mathbb{T} = \{., id, @\}$ ,  $\mathbb{N} = \{S, N\}$ :  $G_{18} : \begin{cases} S \rightarrow N@N.id \\ N \rightarrow id \mid id.N \end{cases}$

- (a) Donner l'arbre de dérivation pour la chaîne **id@id.id.id**  
 (b) Peut-on analyser cette grammaire par une méthode descendante ? Pourquoi ? Si votre réponse est négative, la transformer donc.  
 (c) La grammaire transformée obtenue dans la question précédente est-elle LL(1)? Justifier.  
 (d) En étudiant le langage généré par  $G_{18}$ , donner une grammaire équivalente qui soit LL(1), puis tracer sa table d'analyse LL.  
 (e) Analyser la chaîne introduite en (a).

**Exercice 12**

Soit la grammaire  $G_{19}$  ci-après, où  $\mathbb{T} = \{a, b\}$ ,  $\mathbb{N} = \{S, N\}$ :

$$G_{19} : \begin{cases} S \rightarrow abN \mid \varepsilon \\ N \rightarrow Saa \mid b \end{cases}$$

(a) Calculer les ensembles  $DEB_1$  et  $SUIV_1$  de  $S$  et  $N$

(b) Montrer que  $G_{19}$  n'est pas LL(1).

On étend les ensembles DEB et SUIV à  $k$  caractères  $k > 1$  comme suit :

$$DEB_k(X) = \{w \in \Sigma^* \mid X \rightarrow^* w\alpha \text{ et } |w| = k, \text{ ou } X \rightarrow^* w \text{ et } |w| < k\}$$

$$SUIV_k(X) = \{w \in \Sigma^* \mid S \rightarrow^* \alpha X \beta \text{ et } |w| \in DEB_k(\beta)\}$$

(c) Calculer les ensembles  $DEB_2$  et  $SUIV_2$  de  $S$  et  $N$ .

Sachant que la construction d'une table LL( $k$ ) et la même que pour une analyse LL(1), en remplaçant le terminal de prévision par  $DEB_k$  (ou  $SUIV_k$ ) terminaux.

(d) Tracer la table LL(2) de  $G_{19}$ , puis analyser la chaîne :  $ababbaa$ .

### Exercice 13

On veut construire un analyseur syntaxique pour la grammaire  $G_{19}$  suivante :

$$G_{20} : \begin{cases} S \rightarrow aAS \mid bA \\ A \rightarrow cA \mid d \end{cases}$$

(a) Trouver les ensembles DEB et SUIV des non terminaux de  $G_{20}$

(b) Calculer la collection des items LR[0].  $G$  est-elle SLR(1) ? Justifier.

(c) Construire sa table d'analyse, puis la simuler sur la chaîne : **acdbd**

### Exercice 14

$$\text{Soit la grammaire } G_{21} \text{ ci-contre : } G_{21} : \begin{cases} S \rightarrow A \mid cb \\ A \rightarrow aAb \mid B \\ B \rightarrow c \end{cases}$$

(a) Construire l'automate des items LR(1) et la table d'analyse associée

(b) Analyser la chaîne  $acbb$

### Exercice 15

Montrer que la grammaire suivante est LR(1) mais non LALR(1):

$$G_{22} : \begin{cases} S \rightarrow aAd \mid bBd \mid aBe \mid bAe \\ A \rightarrow c \\ B \rightarrow c \end{cases}$$

### Exercice 16

$$\text{Soit la grammaire } G \text{ suivante : } G_{23} : \begin{cases} S \rightarrow AS \mid b \\ A \rightarrow SA \mid a \end{cases}$$

(a)  $G_{23}$  est-elle SLR. Si oui, construire sa table d'analyse.

(b) Cette grammaire est elle LR ? LALR ?

### Exercice 17

Étant donné la grammaire  $G_{19}$  ci-contre:  $S \rightarrow S \star S \mid \blacktriangle$

Discuter en justifiant si  $G_{19}$  est :

(a) Ambiguë

(b) LL

(c) SLR

(d) LALR

(e) Simuler l'exécution de l'analyseur approprié de  $G$  sur la chaîne :  $\blacktriangle \star \blacktriangle \star \blacktriangle$

Ministère de l'Enseignement Supérieur et de la Recherche Scientifique  
Université de Ghardaïa  
Faculté des Sciences et de la Technologie  
Département des Mathématiques et d'Informatique  
Troisième Année Licence Informatique

---

## Mini projet 2 : Expressions arithmétiques étendues

---

### 1 Objectif

Après avoir exploré les différentes techniques d'analyse syntaxique. On vous propose de réaliser un petit analyseur syntaxique pour les expressions arithmétiques en utilisant les outils **Bison** et **Flex**. le point de démarrage sera deux fichiers de spécification lexicale(Flex) et syntaxique(Bison) à enrichir graduellement par des constructions simples.

### 2 Bison

Bison est la version gnu (libre) du fameux générateur automatique d'analyseurs syntaxique yacc. Ce dernier crée un analyseur **LALR(1)** extensible par des actions sémantiques et des règles de désambiguïsation pour les grammaires non LALR(1). Le format du fichier de spécification est semblable à celui de Flex à savoir : une partie déclaration, une deuxième de définition de règles et une dernière pour le code supplémentaire.

#### 2.1 Spécification Bison

- Les unités lexicales seront déclarées par la directive **%token**, qui donne lieu à l'insertion d'une déclaration par **#define** dans le fichier entête calc.tab.c produit par Bison.
- Les directives **%left** (resp. **%right**) permettent de définir une règle d'associativité gauche (resp. droite), par exemple **%left +** : signifie que l'expression  $a+b+c$  sera évaluée :  $(a+b)+c$ . Noter qu'un token déclaré par **%left** (resp. **%right**) n'a pas besoin d'être introduit par une directive **%token**. **%nonassoc** : pour les opérateurs non associatifs.
- L'ordre de précedence est déterminé par l'ordre d'apparition des tokens(les premiers listés ont la priorité la plus faible). Si plusieurs tokens possèdent le même ordre de précedence, il seront déclarés en même temps sur la même ligne.
- On introduit la grammaire avec les actions sémantiques associées aux différentes productions. Ces dernières seront séparées sur différentes lignes par **'|'** comme dans la notation BNF. Une ligne vide remplace l'alternative vide ( $\epsilon$ )
- Si une règle de production a été reconnue l'action associée sera exécutée. Comme en Flex, une action est un bloc de code en C. les  $\$n$  désignent les numéros des symboles comme ils apparaissent dans le MDP et **\$\$** fait référence au MGP. Exemple :  $E \rightarrow E + E$  { $$$ = \$1 + \$3;$ }

- (f) La fonction **main()** de l'outil appelle **yyparse()** pour traiter l'entrée, qui appelle à son tour la fonction **yylex()** afin de récupérer les tokens. Si une erreur est rencontrée, l'analyseur appelle **yyerror(m)**, où m est le message d'erreur devant être renvoyé.
- (g) Bison maintient **2 piles séparées** : une première pour l'analyse syntaxique qui contient les symboles de la grammaire et une deuxième de valeurs contenant les attributs de ces symboles.
- (h) La communication entre Bison et Flex est assurée via la variable **yyval**, son type est le même que celui des éléments de la pile des valeurs

## 2.2 Construction de l'analyseur

1. Compiler la spécification Flex par : **Flex -o calc.c calc.l** : produira le fichier calc.c à partir de calc.l
2. Compiler la spécification Bison par : **Bison -d calc.y** : produira calc.tab.h et calc.tab.c à partir de calc.y
3. Lier le tout par : **gcc calc.c calc.tab.y -o calculateur** : produira calculateur à partir de l'ensemble des fichiers .c et.h
4. Invoquer votre analyseur par : **./calculateur** : où vous pouvez taper vos expressions.

## 2.3 Travail demandé

- (a)
  1. A partir du groupe de travail télécharger les deux fichiers de spécification calc.l et calc.y (dossier : 3I/compilation/TP)
  2. Compiler et tester votre mini-calculatrice sur des exemples d'expressions simples avec les opérateurs + et -
  3. Compléter la grammaire par l'introduction des opérateurs binaires : \*, /
  4. Remarquer les messages des conflits détectés par Bison, puis introduire des règles d'associativité/précédence et réexaminer le résultat.
  5. Ajouter l'opérateur unaire -
  6. Ajouter à votre grammaire la possibilité d'introduction de variable dont le nom est composé d'une seule lettre (autoriser toutes les lettres de l'alphabet sauf le Q : en majuscule !), ainsi que l'instruction d'affectation simple (associative à droite). Exemple :  
 $x = 14$
  7. Ajouter les expressions parenthésées, du genre  $((1+2)*x/(1-3))/12$
  8. Ajouter une production pour terminer les calcul et quitter la calculatrice si on tape la lettre Q ( en majuscule)
  9. Modifier la grammaire pour que toute expression soit terminée par un point-virgule suivie d'une nouvelle ligne
  10. Le type des attribut étant ENTIER par défaut modifier-le pour permettre des expressions avec des flottants (le type double en C)
- (b) S'il vous reste du temps

1. Permettre cette fois-ci des noms de variables quelconques dont la longueur n'excède pas 5 caractères.
2. Ajouter la structure conditionnelle simple : **Si (condition) Alors Action.**

### 3 Exemples

```
> 1+1
> 2
> (2+5)
> 7
> -6-3
> -9
> x=y=12
> 12
> x+y
> 24
> x+3*(y/2-1)
> 27
> 1+2-
> syntax error
> 1+2:5
> caractere invalide syntaxe error
> Q
> Au revoir !
```

Ministère de l'Enseignement Supérieur et de la Recherche Scientifique  
Université de Ghardaïa  
Faculté des Sciences et de la Technologie  
Département des Mathématiques et d'Informatique  
Troisième Année Licence Informatique

---

## Mini projet 1 : Analyseur lexical et syntaxique de Slim

---

### 1 Objectif

Dans ce mini-projet nous envisageons d'écrire un analyseur pour le mini-langage **Slim**. Votre mission est d'élaborer un scanner pour ce langage en fournissant les ERs et les règles de reconnaissance des différents tokens valides du langage et les actions appropriées associées tel que retourner le token du type correct, enregistrer la valeur d'un lexème si nécessaire, ou reporter une erreur.

### 2 Structure lexicale de Slim

- Les identificateurs (**ID**) sont des chaînes (autres que les mots clés) composées de lettres, chiffres et le blanc souligné. Un Id ne se termine pas par un blanc souligné et ne contient pas deux blancs soulignés consécutifs,
- Les entiers (**INT**) sont des chaînes non vides de chiffres de 0 à 9,
- Les nombres flottants (**REAL**) contenant obligatoirement le point décimal avec le signe et la notation scientifique en option
- Les mots clés (**KWD**) insensibles à la casse : **if, els, fi, wle, elw, rpt, unt, rd, wr, tr, fl, int, ft, bol, str**
- Les chaînes de caractères (**STR**) : suite de caractères encadrés par les couples "bonjour" n'excédant pas une ligne. Dans une chaîne une séquence '\c' dénote le caractère c avec les exceptions suivantes :
  - \b backspace
  - \t tab
  - \n newline
  - \f formfeed

Certaines restrictions sont aussi appliquées aux chaînes de caractères. Celles-ci sont :

- Une chaîne ne peut contenir le nul (le caractère \0) ni le EOF
- Une chaîne ne peut dépasser les limites des fichiers

- Une chaîne ne peut inclure de nouvelles lignes non escapées. Par exemple :  
 "Cette chaîne \  
 est correcte "  
 par contre la suivante est invalide :  
 "Cette chaîne  
 est incorrecte "

- Les opérateurs (**OPR**) : + - \* / < <= > >= := == != && || !
- Séparateurs (**SP**) : ; , ( ) :
- Les commentaires (**CMT**) de deux formes : la première inclut tous les caractères entre "#"  
 et la prochaine retour à la ligne ( ou EOF s'il n'y a pas de newline ) ou ceux multi-lignes  
 encadrés par les couples "(\*" et "\*)". Cette dernière forme peut être imbriquée.
- Les blancs (**BLC**): les suites non vides de : blanc (ascii 32), \n (newline, ascii 10), \f  
 (form feed, ascii 12), \r (carriage return, ascii 13), \t (tab, ascii 9), \v (vertical tab, ascii  
 11).

### 3 Syntaxe de Slim

La grammaire du langage **Slim** est donnée ci-après : les symboles en gras sont des terminaux, ceux entre les couples < et > représentent les non terminaux.

```

< prgm > → < decl > < bloc >
< decl > → < type > < list > ; < decl > | ε
< type > → int | flt | bol | str
< list > → < list > , < id > | < id >
< bloc > → < bloc > ; < instr > | < instr >
< instr > → < ass > | < if > | < wle > | < rpt > | < rd > | < wr >
< ass > → < id > := < exp >
< if > → if < exp > : < bloc > < sif >
< sif > → fi | els : < bloc > fi
< wle > → wle < exp > : < bloc > elw
< rpt > → rpt < bloc > unt < exp >
< rd > → rd < id >
< wr > → wr < exp >
< exp > → < exp > < op > < term > | < term > | < str >
< term > → < term > < opb > < fact > | < fact >
< fact > → < nb > | < id > | ( < exp > ) | < opu > < exp > | tr | fl
< op > → + | - | || | < | <= | > | >= | == | !=
< opb > → * | / | &&
< opu > → ! | -
< nb > → INT | REAL
< id > → ID
< str > → STR

```

## 4 Gestion des erreurs

Toutes les erreurs doivent être retournées au analyseur syntaxique, vous ne devez rien imprimer. Les erreurs seront communiquées en retournant un token spécial **ERROR** et suivant les exigences ci-après :

- Si un caractère invalide est rencontré, la chaîne contenant ce caractère doit être retournée comme chaîne d'erreur,
- Si une chaîne inclut un newline non escapé reporter l'erreur : "Constante chaîne non terminée", et reprendre l'analyse au début de la ligne suivante,
- Dans le cas d'une chaîne trop longue reporter "Constante chaîne trop longues" dans la chaîne erreur du token **ERROR**. Si la chaîne contient un caractère invalide ( le nul par exemple) reporter l'erreur comme "Chaîne contenant le caractère null", dans les deux cas l'analyse reprendra après la fin de la chaîne. Cette fin de chaîne est définie soit comme :
  1. le début de la prochaine ligne si un newline non escapé est rencontré après ces erreurs ou
  2. après les deux cotes de fermeture "
- Si un commentaire reste ouvert alors qu'un EOF est rencontré reporter l'erreur : "EOF dans un commentaire" ne pas considérer alors ce token étant donné que la terminaison est manquante. De même, si EOF est rencontrés avant la quote de fermeture d'une chaîne reporter "EOF dans une constante chaîne"
- Si vous rencontrez "\*" en dehors d'un commentaire reporter "\*" non reconnu" au lieu de la tokeniser en \* et )
- Pour les erreurs syntaxiques indiquer la ligne contenant l'erreur ainsi que l'UL concernée.

## 5 Travail à faire

1. Après avoir codifié les différentes ULs. Écrire un analyseur lexical pour le langage **Slim** en utilisant le générateur flex et la machine virtuelle **Compil** vue en TP. L'analyseur lexical recevra un code source dans un fichier texte et fournira son résultat (suite des codes des unités lexicales reconnues, ou un code d'erreur) dans un autre fichier texte de sortie
2. Écrire un analyseur syntaxique pour ce mini-langage par la **descente récursive** (Vous serez amené à transformer certaines productions de la grammaire du langage **Slim** !).
3. Fournir vos fichiers codes ainsi qu'un rapport détaillé augmenté d'exemples illustratifs **avant le 07/12/2022 à minuit** à l'adresse électronique : **s.ouladnaoui@gmail.com**. Vous pouvez travailler en binômes. Aucun retard ne sera toléré.

Bon courage

# **Examens corrigés**

**Exercice I (QC & Analyse Lexicale, 4 points)**

- (1) Répondre par **VRAI** ou **FAUX** (1)
  - a. La compilation est un processus de traduction d'un langage vers un autre
  - b. **gcc** est un interpréteur
  - c. On peut dénoter les mots **palindromes** sur un alphabet  $\Sigma$  par une expression rationnelle
  - d. Une grammaire ambiguë **ne constitue pas** un problème pour un analyseur **LR(k)**
- (2) Donner une expression rationnelle spécifiant un identificateur. Un identificateur est une chaîne composée de lettres, chiffres décimaux et le caractère `_`. Il doit commencer par une lettre et ne doit pas se terminer par `_` ni contenir deux `_` consécutifs. (1)
- (3) Quelle(s) chaîne(s) parmi les suivantes **ne sont pas traitée(s) correctement** (de façon réussie) par la spécification lexicale ci-dessous: (2)
  - (xx)\*
  - xy+
  - yx+
  - a. xxxxxx
  - b. xyxxxyx
  - c. xxxxyxy
  - d. xyxxxyx

**Exercice II (Analyse Syntaxique I, 8 points)**

Soit la grammaire  $G = \langle T, N, S, P \rangle$  suivante :

$$\begin{cases} S \rightarrow [SA] \mid a \\ A \rightarrow +SB \mid Bb \mid \varepsilon \\ B \rightarrow -SAC \mid \varepsilon \end{cases}$$

- (1) En calculant les ensembles DEBUT et SUIVANT des éléments de  $N$ , montrer que  $G$  est LL(1) (4)
- (2) Construire la table d'analyse LL pour  $G$  (2)
- (3) Analyser la chaîne **[a+a-ac]** (2)

**Exercice III (Analyse Syntaxique II, 8 points)**

On s'intéresse à la construction d'un analyseur syntaxique pour la grammaire  $G$  ci-contre

$$S \rightarrow aSc \mid SS \mid b$$

- (1) Quelle approche doit-on adopter ? (1)
- (2) Calculer la collection des items LR[0].  $G$  est-elle SLR(1) ? Justifier. (4)
- (3) Discuter les autres alternatives de construction d'un analyseur syntaxique pour cette grammaire. (1)
- (4) En choisissant la piste la plus simple, dire combien d'opérations de décalage et de réduction sont nécessaires pour accepter la chaîne : **abbbc** (2)

Bon courage

**Exercice I (Analyse Lexicale, 5 points)**

- (1) Étant donné la spécification lexicale ci-dessous. (2)

```
%%  
aa      printf("1");  
b?a+b?  printf("2");  
b?a*b?  printf("3");  
.       printf("4");
```

Indiquer pour les chaînes suivantes la sortie de l'analyseur et les tokens produits :

- bbbabaa
  - aaabbbbbaabanalyse
- (2) Donner un exemple d'une entrée pour laquelle cet analyseur produit la sortie **123**, où justifier, le cas échéant, l'inexistence de tel exemple. (2)
- (3) Quelle sera pour les mêmes chaînes données en (1) la sortie de l'analyseur, si on remplace seulement la partie expression régulière de la dernière règle par **.+** (1)

**Exercice II (Analyse Syntaxique I, 6 points)**

Soit la grammaire transformée des expressions arithmétiques vue en cours :

$$G_1 : \begin{cases} E \rightarrow TP \\ P \rightarrow +TP \mid -TP \mid \varepsilon \\ T \rightarrow FQ \\ Q \rightarrow \times FQ \mid \div FQ \mid \varepsilon \\ F \rightarrow (E) \mid int \mid id \end{cases}$$

- (1) Écrire un analyseur syntaxique pour  $G_1$  en utilisant la descente récursive. (5)
- (2) Analyser la chaîne:  $int \times (id - id)$  (1)

**Exercice III (Analyse Syntaxique II, 9 points)**

On s'intéresse à la construction d'un analyseur syntaxique pour la grammaire  $G$  ci-après, où  $\mathbb{T} = \{., id, @\}$ ,  $\mathbb{N} = \{S, N\}$ :

$$G_2 : \begin{cases} S \rightarrow N@N.id \\ N \rightarrow id \mid id.N \end{cases}$$

- (1) Donner un arbre de dérivation pour la chaîne **id@id.id.id** (1)
- (2) Peut-on analyser cette grammaire par une méthode descendante ? Pourquoi ? Si votre réponse est négative, la transformer donc. (2)
- (3) La grammaire transformée obtenue dans la question précédente est-elle LL(1)? Justifier. (2)
- (4) En étudiant le langage généré par  $G_2$ , donner une grammaire équivalente qui soit LL(1), puis tracer sa table d'analyse LL. (3)
- (5) Analyser la chaîne introduite en (1). (1)

Bon courage

**Exercice I (Analyse Lexicale, 10 points)**

Considérons un sous ensemble du langage de programmation Pascal composé des unités lexicales suivantes :

- **Mots clés insensibles à la casse**  $\Rightarrow$  program, var, integer, begin, end, if, then
- **Identificateurs**  $\Rightarrow$  les suites de lettres ou de chiffres commençant par une lettre
- **Nombres**  $\Rightarrow$  les suites non vides de chiffres
- **Opérateurs**  $\Rightarrow$  +, -, \*, >, >=, <, <=, =, <>, :=
- **Séparateurs**  $\Rightarrow$  , ; ( ) : .

- (1) Donner le contenu du code Flex pour reconnaître les différentes unités lexicales de ce mini-langage. À la reconnaissance d'une unité lexicale, veuillez afficher sa classe ainsi que le lexème reconnu. **Signaler tout autre caractère différent des blancs ( ' ', Tab, Nouvelle ligne) comme illégal.** (5)
- (2) Écrire l'algorithme d'un analyseur lexical **codé en dur** permettant la reconnaissance des unités lexicales de ce même langage. L'algorithme doit lire, **caractère par caractère**, le code source consigné dans un **fichier d'entrée** et produire les tokens (classe, lexème) dans un deuxième **fichier de sortie**. Vous pouvez supposer l'existence de fonctions de traitements de caractères telles que : Lettre(), Chiffre(), Minuscule(), etc. (5)

**Exercice II (Analyse Syntaxique, 10 points)**

Soit la grammaire  $G$  ci-contre.

$$G : \begin{cases} S \rightarrow 0AB5 \\ A \rightarrow 2 \mid A1 \\ B \rightarrow CD \\ C \rightarrow 4 \mid \varepsilon \\ D \rightarrow 3 \mid \varepsilon \end{cases}$$

- (1) Donner les arbres de dérivation des chaînes : 0215 et 02435 (1)
- (2) Peut-on analyser syntaxiquement cette grammaire par la descente récursive ? Justifier (1)
- (3) Transformer  $G$  si nécessaire (1)
- (4) Calculer les ensembles DEB et SUIV des non terminaux de la grammaire obtenue (3)
- (5) Tracer la table d'analyse LL de cette grammaire et montrer qu'elle est LL(1) (3)
- (6) Analyser la chaîne : 025# (½)
- (7) Trouver le langage généré par  $G$ . (½)

Bon courage

Corrigé type de l'examen final de compilation

EXO I

code Flex

)

```

program [pP][rR][oO][gG][rR][aA][mM]
var [rv][aA][rR]
integer [iI][nN][tT][eE][gG][eE][rR]
begin [bB][eE][gG][iI][nN]
end [eE][nN][dD]
if [iI][fF]
then [tT][hH][eE][nN]
lettre [a-zA-z]
chiffre [0-9]
Blancs [ \t\n]+
%%
[- + * > < =] | "<=" | ">=" | "<>" | "!=" printf ("lop, %s)\n", yytext);
[ , ; ( ) = . ] printf ("(sep, %s)\n", yytext);
{ program } | { var } | { integer } | { begin } | { end } | { if } | { then }
printf ("(mot-clé, %s)\n", yytext);
{ lettre } ( { lettre } | { chiffre } ) * printf ("(Id, %s)\n", yytext);
{ chiffre } + printf ("(nb, %s)\n", yytext);
{ Blancs } ;

```

1/4 .

```
printf ("caract. illegal, %c\n", yytext[0]);
```

## 2) Analyseur en dur

Algorithme lexer

Etiquette test  
Constante KW = Tableau [7] de chaîne = ('program', 'var', 'integer', 'begin', 'end', 'if', 'then')

variables blancs = [ " , \t, \n ]  
in, out = Text  
c = caractère  
lexeme = chaîne

debut  
Assigner (in, '---'); ouvrir (in)  
Assigner (out, '---'); créer (out)

lexeme ← "

1/4 Tant que Non (FDF (in)) faire

FDF = fin de fichier

debut lire (in, c)

test : selon Minuscule (c) faire

1/2 '+' , '-' , '\*' , '=' : Entité Lexicale ('OP', c)

1/2 ',' , ':' , ';' , '(' , ')' : Entité Lexicale ('SEP', c)

1/2 ':' : debut lire (in, c)  
si c = '=' alors Entité Lexicale ('OP', '=')  
sinon debut Entité Lexicale ('SEP', ':')  
Aller à test  
fin

'<' : début  
 lire(in, c)  
 si c='=' alors EntitéLexicale('OP', '<=')  
 |  
 | sinon si c='>' alors EntitéLexicale('OP', '<>')  
 |  
 | sinon début EntitéLexicale('OP', '<')  
 |  
 | Aller à test  
 | fin  
 |  
 | fin

'>' : début  
 lire(in, c)  
 si c='=' alors EntitéLexicale('OP', '>=')  
 |  
 | sinon début EntitéLexicale('OP', '>')  
 |  
 | Aller à test  
 | fin  
 |  
 | fin

'a', 'b', ..., 'z' : début  
 Tantque Non FDF(in) et (lettre(c) ou chiffre(c)) faire  
 |  
 | début lexème ← lexème + c  
 |  
 | lire(in, c)  
 |  
 | fin  
 |  
 | si lexème dans KW alors EntitéLexicale('mot clé',  
 |  
 | | lexème)  
 |  
 | | sinon EntitéLexicale('Id', lexème)  
 |  
 | | Aller à test  
 |  
 | fin

'0', ..., '9' : début  
 Tantque chiffre(c) et Non FDF(in) faire



G' :  $\begin{cases} S \rightarrow 0AB5 \\ A \rightarrow 2X \\ X \rightarrow 1X | \epsilon \\ B \rightarrow CD \\ C \rightarrow 4 | \epsilon \\ D \rightarrow 3 | \epsilon \end{cases}$

①

4)

	DEB	SUIV
S	0	#
A	2	3, 4, 5
X	1, $\epsilon$	3, 4, 5
B	3, 4, $\epsilon$	5
C	4, $\epsilon$	3, 5
D	3, $\epsilon$	5

5)

Table LL

	0	1	2	3	4	5	#
S	0AB5						
A			2X				
X		1X		$\epsilon$	$\epsilon$	$\epsilon$	
B				CD	CD	$\epsilon$	
C				$\epsilon$	4	$\epsilon$	
D				3		$\epsilon$	

G' est factorisée, et non RG et  
 La table d'analyse LL est monodéfinie  
 $\Rightarrow G'$  est LL(1)

6)

Pile	chaîne	Action
#S	025 #	remplacer S par 0AB5
#5BA0	025 #	Avancer et dépiler
#5BA	25 #	remplacer A par 2X
#5BX2	25 #	avancer et dépiler
#5BX	5 #	remplacer X par $\epsilon$
#5B	5 #	" B " "
#5	5 #	Avancer et dépiler
#	#	Accepter

7)  $L(G) = 021*4?3?5$

Mard. 8/01/19

**EXAMEN FINAL DE COMPILATION**

Durée : 1<sup>h</sup>30

**DOCUMENTS, CALCULETTES ET TÉLÉPHONES PORTABLES INTERDITS**

**Exercice I (Analyse lexicale, 6 points)**

Les identificateurs d'un langage imaginé sont composés de lettres (minuscules ou majuscules), de chiffres décimaux et du symbole blanc souligné \_ (underscore). Un identificateur de ce langage est soumis aux contraintes suivantes :

- Commence par une lettre minuscule ou un underscore et ne se termine jamais par un underscore
- La succession de deux chiffres, deux lettres, ou deux underscores est interdite
- Sa longueur est divisible par 5

- (1) Concevoir un AFD ayant moins de 6 états qui aide à reconnaître ce token (dans un premier temps ignorer dans votre AFD la contrainte de longueur) (3)
- (2) En utilisant l'automate précédent, écrire la fonction de reconnaissance correspondante (retourne vrai si la chaîne donnée en paramètre est un identificateur correct, *i.e* respecte toutes les conditions) (2)
- (3) On élimine maintenant la dernière condition et on autorise la répétition de lettres et des chiffres seulement (les conditions sur le underscore restent inchangées). Donner un programme **flex** minimal qui permet de reconnaître cette nouvelle classe d'identificateurs. (1)

**Exercice II (Analyse Syntaxique I, 7 points)**

Considérons la grammaire ci-contre  $G$  :

$$G : \begin{cases} X \rightarrow T \mid G \\ T \rightarrow Xa \\ G \rightarrow YZb \mid YZc \\ Y \rightarrow Yd \mid e \\ Z \rightarrow Zf \mid \epsilon \end{cases}$$

- (1) On veut construire un analyseur syntaxique descendant de  $G$ . Quels problèmes présente cette grammaire. (2)
- (2) Corriger les problèmes trouvés dans (1) (transformer et éventuellement simplifier  $G$ ) (2)
- (3) Écrire un analyseur pour  $G$  en utilisant la descente récursive. (3)

**Exercice III (Analyse syntaxique II, 7 points)**

On se propose de réaliser une analyse syntaxique LL pour la grammaire suivante  $G$  :

$$G : \begin{cases} S \rightarrow ABb \\ A \rightarrow CD \\ B \rightarrow dB \mid \epsilon \\ C \rightarrow aCb \mid \epsilon \\ D \rightarrow cDd \mid \epsilon \end{cases}$$

- (1) Calculer les ensembles **DEB** et **SUIV** des non terminaux de  $G$  (2½)
- (2) Tracer sa table d'analyse LL (2½)
- (3) Analyser la chaîne  $abdb$  (2)

Bon courage

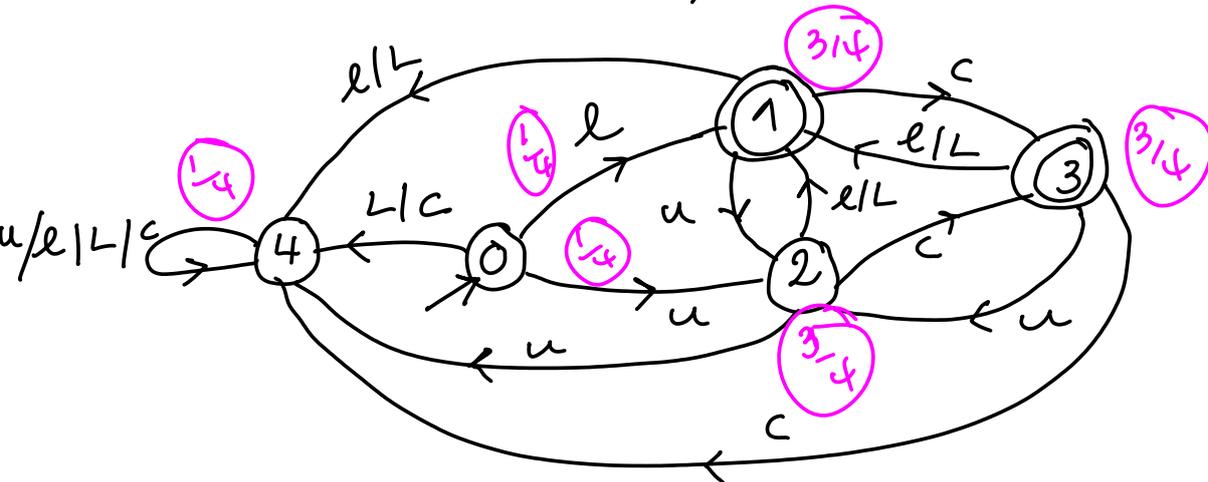
7/01/19

# Corrigé type de l'examen final compil

## Exo I

1) Soit pour simplifier la notation suivante

$l \rightarrow [a-z]$  /\* lettres minuscules \*/  
 $L \rightarrow [A-Z]$  /\* " majuscules \*/  
 $c \rightarrow [0-9]$  /\* chiffres \*/  
 $u \rightarrow [-]$  /\* underscore \*/



2) On implémente l'automate par sa table de transition

soit T cette Table pré-déclarée et initialisée

```
bool reconnaître_id ( ch: chaîne )  
{  
    (1/2) int etat = 0;  
    char c  
    Tantque ( etat != 4 && ! FDC (ch) )  
    {  
        c = lire_car ();  
        (1) etat ← T [ etat, c ]  
    }  
}
```

$\frac{1}{2}$  return (état  $\in \{1,3\}$  & longueur(ch) % 5 == 0)  
}

3) % {

define ID 257;  $\frac{1}{4}$

% }

% %

([a-z] | [-][a-zA-Z0-9]) ([a-zA-Z0-9] | [-][a-zA-Z0-9])\* return ID;  $\frac{1}{4}$

### Exo II

1)  $\frac{1}{2}$  les non terminaux  $X, Y, Z$  sont récursifs gauche  
( $X$  de façon indirecte)

$\frac{1}{2}$  -  $X$  n'est pas factorisé (de " " " " )

1) -  $DEB(Zf) \cap suiv(Z) \neq \emptyset$  ( $= \{f\}$ )

2) - remplacer  $T$  et  $G$  dans  $X$ :

$X \rightarrow Xa \mid YZb \mid YZc$

- Élimination de la RG de  $X$ :

$X \rightarrow YZbA \mid YZcA$   $\frac{1}{2}$

$A \rightarrow aA \mid \epsilon$

- Factorisation de  $X$ :

$X \rightarrow YZM$

$M \rightarrow bA \mid cA$   $\frac{1}{2}$

$A \rightarrow aA \mid \epsilon$

- Élimination de la RG de Y :

$$\begin{aligned} Y &\rightarrow eN \\ N &\rightarrow dN | \epsilon \end{aligned} \quad \left(\frac{1}{2}\right)$$

- Élimination de la RG de Z :

Il suffit de remarquer que Z génère  $f^*$   
 $\Rightarrow$  réécrire la RG de Z en RD :

$$Z \rightarrow fZ | \epsilon \quad \left(\frac{1}{2}\right)$$

G devient :

$$G: \begin{cases} \overline{X} \rightarrow YZM \\ M \rightarrow bA | cA \\ A \rightarrow aA | \epsilon \\ Y \rightarrow eN \\ N \rightarrow dN | \epsilon \\ Z \rightarrow fZ | \epsilon \end{cases}$$

3)

S() {  
 \_\_\_\_\_  
 X()  
 si tc = '#' : accepter()  $\left(\frac{1}{4}\right)$   
 | sinon erreur()  
 }

X() {  
 \_\_\_\_\_  
 Y()  $\left(\frac{1}{4}\right)$   
 Z()  
 M()  
 }

Y()  
 {  
   si tc = 'e' : { Avancer() } N() (1/2)  
   |  
   sinon erreur()  
 }

Z()  
 {  
   si tc = 'f' : { Avancer() } Z() (1/2)  
   |  
   sinon ;  
 }

M()  
 {  
   si tc = 'b' ou tc = 'c' : { Avancer() } A() (1/2)  
   |  
   sinon erreur()  
 }

A()  
 {  
   si tc = 'a' : { Avancer() } A() (1/2)  
   |  
   sinon ;  
 }

N()  
 {  
   si tc = 'd' : { Avancer() } N() (1/2)  
   |  
   sinon ;  
 }

Exo III 1)

	DEB	Suiv
S	a, b, c, d	#
A	a, c, ε	b, d
B	d, ε	b
C	a, ε	b, c, d
D	c, ε	b, d

$\frac{1}{2}$   
 $\frac{1}{2}$   
 $\frac{1}{2}$   
 $\frac{1}{2}$   
 $\frac{1}{2}$

2)

	a	b	c	d	#
S	1	1	1	1	
A	2	2	2	2	
B		4		3	
C	5	6	6	6	
D		8	7	8	

$\frac{1}{2}$   
 $\frac{1}{2}$   
 $\frac{1}{2}$   
 $\frac{1}{2}$   
 $\frac{1}{2}$

3)

Pile	chaîne	Action
S #	abdb #	S → ABb
ABb #	abdb #	A → CD
CD B b #	abdb #	C → aCb
aCbDBb #	abdb #	Avancer
CbDBb #	bdb #	C → ε
bDBb #	bdb #	avancer
DBb #	db #	D → ε
Bb #	db #	B → dB
dBb #	db #	avancer
Bb #	b #	B → ε
b #	b #	Avancer
#	#	Accepter

1

$\frac{1}{2}$

$\frac{1}{2}$

**DOCUMENTS INTERDITS**

**Exercice I (AF, ER et UL, 5 points)**

On définit ici l'unité lexicale ID, comme une suite de symboles commençant par \$ suivi d'une suite quelconque d'un flot alterné de chaîne non vide de lettres latines minuscules suivie d'une chaîne quelconque de chiffres décimaux. Sa longueur est supérieure ou égale à 2.

- (1) Donner un AFD qui accepte l'unité lexicale ID (3)
- (2) Écrire le code flex minimal qui le reconnaît (2)

**Exercice II (AL & flex, 4 points)**

Étant donné l'analyseur implémentant la spécification lexicale suivante :

```
0(10)*      printf("a");  
1*(01)*    printf("b");  
(00|11)*   printf("c");  
[^2-9\n]{2,5} printf("d");
```

- (1) Donner la sortie de l'analyseur et la tokenisation pour la chaîne 11111001111100010 (2)
- (2) Étudier l'existence d'une chaîne pour laquelle cet analyseur produit la sortie *bac*? (2)

**Exercice III (Analyse Syntaxique, 11 points)**

Considérons la grammaire ci-contre où :  $N = \{S, A, B, C, D\}$ ,  $T = \{a, b, c, d, e\}$ , et  $p$  un paramètre élément de  $T$ .

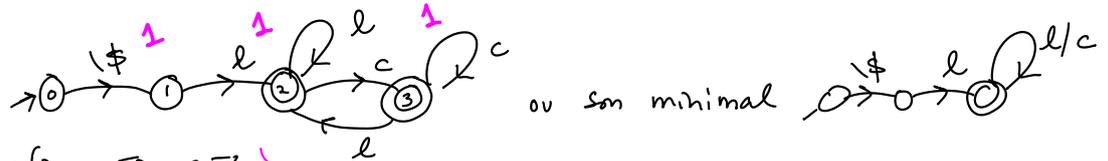
$$G_p : \begin{cases} S \rightarrow DCeBA \\ A \rightarrow a \mid \varepsilon \\ B \rightarrow DC \mid p \mid \varepsilon \\ C \rightarrow c \mid \varepsilon \\ D \rightarrow CbD \mid d \end{cases}$$

- (1) Calculer les ensembles DEB et SUIV des éléments de  $N$  (5)
- (2) Montrer qu'il existe un  $p$  tel que  $G_p$  est LL (2)
- (3) Tracer la table LL de  $G_p$  pour la valeur de  $p$  trouvée en (2). (3)
- (4) Analyser la chaîne *cbdeea* (1)

Bon courage

Exo I

1) Soit  $l \in [a-z]$  et  $c \in [0-9]$



2) 

```
%{
%}
define ID 256
%{
%}
%{
%}
\$\{ [a-z] [a-z0-9]*
return(ID);
```

Exo II

1)  $11111 | 001111 | 100010$   
 $\swarrow \quad \swarrow \quad \swarrow \quad \swarrow$   
 $b \quad c \quad d \quad a$

$0(10)^*$  — 'a'  
 $1^*(01)^*$  — 'b'  
 $(00|11)^*$  — 'c'  
 $[2-9|n]\{2,5\}$  — 'd'

2) Pour produire bac la chaîne doit correspondre au motif  $1^*(01)^* 0(10)^* (00|11)^*$

Le 1<sup>er</sup> problème est que le motif  $0(10)^*$  sera toujours consommé par la 2<sup>e</sup> règle avec un

1  $0$  restant ce qui donne  $1^*(01)^* 0(10)^* (00|11)^* = 1^*(01)^* 0(00|11)^*$   
 pour produire "a" les chaînes de "c" commencent soit par 00 ou 11

et dans les 2 cas cela ne correspond pas au motif de "a" qui est  $0(10)^* \perp$

1 donc cette chaîne n'existe pas.

Exo III

	DEB	SUIV
S	b, c, d	#
A	a, e	#
B	b, c, d, p, e	a, #
C	c, e	a, b, e, #
D	b, c, d	a, c, e, #

2) Conditions pour que GP soit LL

- Prefixes 2 à 2 disjoints:  $DEB(CBD) \cap \{d\} = \{c, b\} \cap \{d\} = \emptyset \quad \checkmark$

$\checkmark$   $DEB(DC) \cap \{p\} = \{b, c, d\} \cap \{p\} = \emptyset \Rightarrow p \notin \{b, c, d\}$

- En présence de productions vides: les préfixes et les suivants du MGP doivent être disjoints

-  $a \notin \text{suiV}(A) \quad \checkmark$

$\checkmark$  -  $c \notin \text{suiV}(C) \quad \checkmark$

$\checkmark$  -  $p \notin \text{suiV}(B) = \{a, \#\} \Rightarrow p \in T$  et  $p \notin \{a, b, c, d\}$

$\checkmark \Rightarrow p \in \{e\} \quad \perp$

Donc  $G \in LR(1)$ .

3)

	a	b	c	d	e	#	
S		DCeBA	DCeBA	DCeBA			6/10
A	a					E	6/10
B	e	DC	DC	DC	e	E	6/10
C	e	E	c		E	E	6/10
D		CbD	CbD	d			6/10

4)

Pile	chaîne	Action
#S	cbdeea #	S → DCeBA 1/3
#ABeCD	cbdeea #	D → CbD
#ABeCDbC	cbdeea #	C → c
#ABeCDbc	cbdeea #	dépiler et avancer
#ABeCDb	bdeea #	" "
#ABeCD	d eea #	D → d
#ABeCd	d eea #	dépiler et avancer
#ABeC	eea #	C → e
#ABe	eea #	dépiler et avancer
#AB	ea #	B → e
#Ae	ea #	dépiler et avancer
#A	a #	A → a
#a	a #	dépiler et avancer
#	#	Accepter.

**Exercice I (Analyse lexicale, 5 points)**

On veut construire un analyseur lexical pour un langage  $L$  composé des tokens suivants :

- les mots réservés sensibles à la casse ci-après : **si, sin, do, done, in, int**
- les **identificateurs** usuels (suites non vides de lettres **minuscules** et de chiffres décimaux commençant par une lettre)
- les **nombres entiers signés**
- les **opérateurs** suivants : +, -, =, <, <=, >, >=, <>

- (1) Donner un AFD qui accepte l'ensemble des tokens ci-dessus cités. (3)
- (2) En supposant que les codes des token sont déjà prédéfinis, écrire le code **flex** (la partie règles seulement) pour produire cet analyseur. (2)

**Exercice II (TP, 4 points)**

Pratiquement tout langage de programmation inclut la classe commentaires. Dans votre Mini-projet :

- (1) Rappeler l'ER d'un commentaire multilignes au style de C (délimités par /\* et \*/), n'incluant pas la chaîne '\*/' mais autorisant les étoiles. (1)
- (2) Comment un analyseur usuel tokenise le motif '\*/' lorsque il le rencontre en dehors d'un commentaire. (1)
- (3) Les ER étant des descripteurs pour les langages réguliers; expliquer comment vous avez pu reconnaître les commentaires multilignes **imbriqués**. (2)

**Exercice III (Analyse syntaxique, 11 points)**

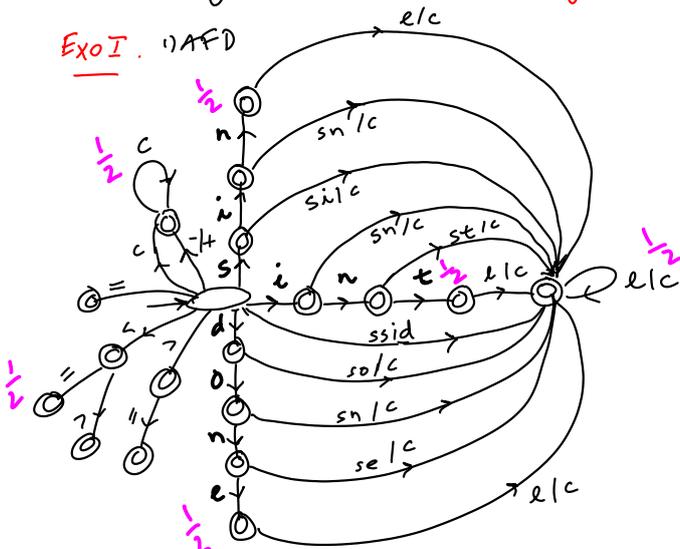
Considérons la grammaire ci-contre :

$$G : \begin{cases} S \rightarrow S[S] \mid Tb \mid c \\ T \rightarrow Ta \mid \varepsilon \end{cases}$$

- (1)  $G$  est-elle LL(1). Justifier. (1)
- (2) Transformer  $G$  pour qu'elle soit LL(1) (2)
- (3) Calculer les ensembles DEB et SUIV des non terminaux de la transformée de  $G$  (3)
- (4) Écrire, en utilisant la descente récursive, un analyseur syntaxique pour cette grammaire (3)
- (5) Analyser la chaîne :  $ab[c]$  (2)

Bon courage

Exo I 1) AFD



e	[a-z]	toutes les lettres
c	[0-9]	tous les chiffres
ss	[a-rt-z]	les lettres sauf s
si	[a-hj-z]	_____ i
sd	[a-ce-z]	_____ d
sn	[a-mo-z]	_____ n
st	[a-su-z]	_____ t
se	[a-df-z]	_____ e
so	[a-np-z]	_____ o
ssdi	[a-ce-hj-rt-z]	_____ s, d, i

```

2) Flex
%%
"si" return csi;
"sin" — csin;
"do" — cdo;
"done" — cdone;
"in" — cin;
"int" — cint;
"<" — cint;
">" — csup;
"=" — cegal;
"+" — cplus;
"-" — cmoins;
"<=" — cmfegal;
">=" — csupegal;
"<>" — cdiff;
[a-z][a-z0-9]* — cid;
[-+]?[0-9]+ — cnb;
    
```

cxx → code de xx

Exo II

1) `"/*" ([^*] | [*]+[^\s/]) * [*] * "*/"`

- un analyseur usuel tokenise `'*/'` → `(*, op)`, `(/, op)`
- Flex inclut le mécanisme des conditions; on l'exploite avec un compteur incrémenté à la rencontre de `/*` et décrémente à la lecture de `*/`

```

%x comment déclaration (INITIAL est la condition par défaut)
%%
<INITIAL> "/*" { BEGIN(comment); nb++; }
<comment> "/*" { nb++; }
<comment> "*/" { nb--; }
if (nb == 0) BEGIN (INITIAL); }
    
```

Exo III 1) Non, G n'est pas LL(1), car elle RG en S ( $S \rightarrow S[S]$ ) et T ( $T \rightarrow Ta$ )

2) Élimination de la RG:  $S \rightarrow TbX \mid cX$   
 $X \rightarrow [S]X \mid \epsilon$   
 $T \rightarrow aT \mid \epsilon$

(on peut faire  $T \rightarrow Y$   
 $Y \rightarrow aY \mid \epsilon$   
 pour constater que  $T=Y$ )

	DÉB	Soix	
S	a, b, c	#, ]	5/4
X	[, ε	#, ]	4/4
T	a, ε	b	3/4

4)  $\underline{z(c)}$  {  $s(c)$ ; si  $tc = \#$  : accepter(); sinon erreur() }  $\frac{1}{2}$

$\underline{s(c)}$  {  
 si  $tc = 'c'$  : { avancer() }  
 }  
 |  
 sinon { T() }  
 |  
 si  $tc = 'b'$  : { avancer() }  
 |  
 sinon erreur() }  
 }

$\underline{x(c)}$  {  
 si  $tc = '['$  : { avancer() }  
 |  
 S(c)  
 |  
 si  $tc = ']'$  : { avancer() }  
 |  
 sinon erreur() }  
 }

$\underline{T(c)}$  {  
 si  $tc = 'a'$  : { avancer() }  
 |  
 T(c)  
 |  
 sinon ;  
 }

z	chaîne	Action
z	ab[c] #	appel à S
zS	ab[c] #	" a T
zST	ab[c] #	avancer() + appel à T
zSTT	b[c] #	retour à T
zST	b[c] #	" a S
zS	b[c] #	avancer + appel à X
zSX	[c] #	" " S
zSXS	c] #	" " X
zSXSX	] #	retour à S
zSXS	] #	" à X
zSX	] #	avancer et appel à X
zSXX	#	retour à X
zSX	#	" à S
zS	#	" à z
z	#	Accepter

2/3

2/3

2/3

**Exercice I (AL I, 5 points)**

Un identificateur d'un langage **L** est composé de lettres latines minuscules, de chiffres décimaux et du symbole \$ avec les conditions suivantes :

- Commence par une lettre minuscule et ne se termine jamais par un \$,
- La succession de deux symboles de la même catégorie (deux lettres, deux chiffres ou deux \$ ) est interdite
- Sa longueur est divisible par 3

(1) Concevoir un AFD qui accepte ce token (identificateur) (3)

(2) Implémenter cet AFD comme reconnaisseur des identificateurs de **L** (2)

**Exercice II (AL II, 4 points)**

Étant donné l'analyseur lexical implémentant la spécification suivante :

```
%%  
[0-2]+0?1          printf("D");  
[02] [^12]+       printf("L");  
[1] [2]* | ([01] [10])* [2] [2] printf("M");
```

(1) Donner la sortie de l'analyseur et la tokenisation pour la chaîne : **21033311112122** (2)

(2) Donner une chaîne, si elle existe, pour laquelle cet analyseur produit la sortie : **LMD**. Justifier (2)

**Exercice III (Analyse Syntaxique I, 14 points)**

Considérons la grammaire ci-contre  $G : \left\{ \begin{array}{l} E \rightarrow R \mid E + R \\ R \rightarrow S \mid SR \\ S \rightarrow T \mid S* \\ T \rightarrow 0 \mid 1 \mid (E) \end{array} \right.$

(1) Relever les problèmes qui rendent  $G$  inadaptée à une analyse syntaxique descendante. Justifier (1)

(2) Transformer  $G$  pour qu'elle devienne **LL**. (3)

(3) Calculer les ensembles **DEB** et **SUIV** des non terminaux de la nouvelle grammaire. (3½)

(4) Construire sa table **LL**. (4)

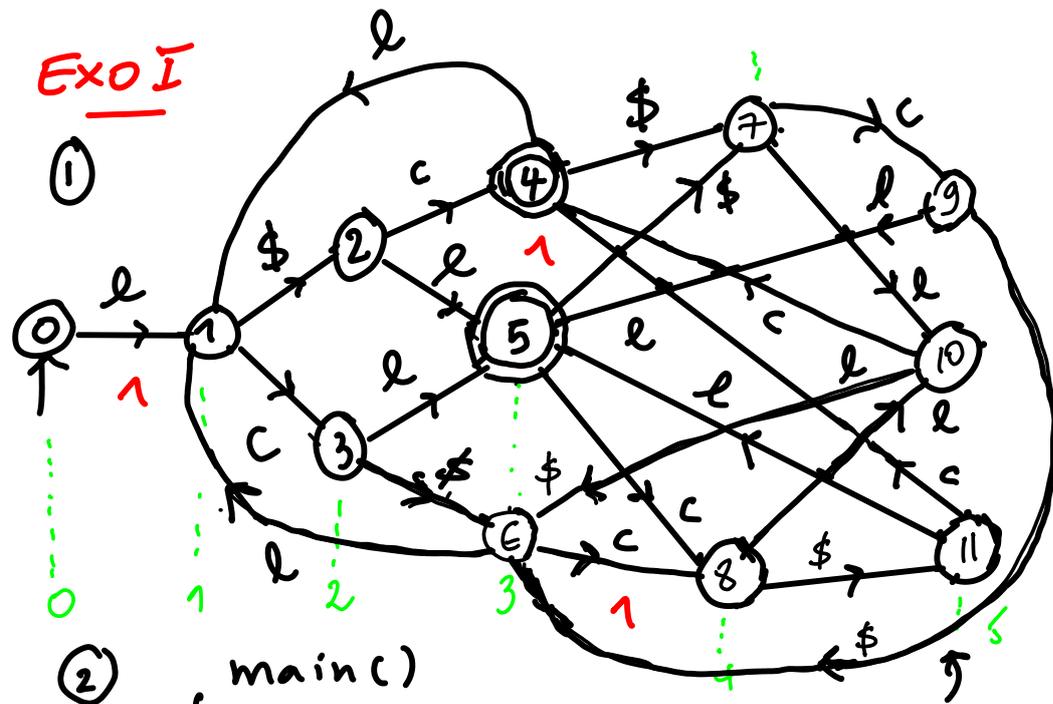
(5) Analyser la chaîne  $0*+1$  (2½)

Bon courage

Corrigé type de l'examen final de  
compilation I (L3)

Exo I

l [a-z]  
c [0-9]



12 états + 37 symboles

```
main()
{ int T[12][37] = { ... } /* initialisation */
```

```
int c, etat = 0, erreur = 0;
```

1/4

```
TQ ( ! erreur && ! FDC )
```

1/2

```
line(c)
```

```
selon c {
  'a' .. 'z' : c = c - 'a' + 10
```

```
  '0' .. '9' : c = c - '0'
```

```
  $ : c = 36
```

1/4

```
  }
  etat = T[etat][c]
```

1/2

}

si état  $\in \{4,5\}$  alors Ecrire ('succès')  $\frac{1}{2}$

sinon Ecrire ('erreur')

}

## Exo II

①

21 | 0333 | 111121 | 22

D      L                      D                      M

$\frac{1}{2}$        $\frac{1}{2}$                        $\frac{1}{2}$                        $\frac{1}{2}$

②

- Pour générer le L, il suffit de choisir une chaîne qui commence par 0 ou 2 suivi d'un symbole autre que 1 ou 2 soit (23)

- Pour produire M choisir une chaîne qui commence par (0|1|2) et se termine par c + 1

$\frac{1}{2}$  c  $\neq$  1 (éviter 0) et  $\in \{1,2\}$  (éviter L)

soit (12)

- Pour générer le dernier D : impossible  $\frac{1}{2}$

car la règle commence par  $[0-2]^+$  qui

est équivalente à  $(0112)(0112)^*$   
 et va par conséquent consommer ce qui a  
 reconnu la troisième (celle de  $M$ )  
 la sortie devient donc LD (au lieu de LM)

Exo III.

①  $G$  st RG dans  $E$  ( $E \rightarrow E+R$ ) et dans  $\frac{1}{3}$   
 $S$  ( $S \rightarrow S^*$ )  $\frac{1}{3}$

et elle n'est pas factorisée dans  $R$  ( $R \rightarrow S \mid SR$ )  $\frac{1}{3}$

② transformation

$$\frac{1}{2} E \rightarrow RX$$

$$\frac{1}{2} X \rightarrow +RX \mid E$$

$$\frac{1}{2} R \rightarrow SY$$

$$\frac{1}{2} Y \rightarrow R \mid E$$

$$\frac{1}{2} S \rightarrow TZ$$

$$\frac{1}{2} Z \rightarrow *Z \mid E$$

$$T \rightarrow 0 \mid 1 \mid (E)$$

$\frac{1}{4} \times 7 \quad \frac{1}{4} \times 7 \quad ) \Rightarrow 3 \frac{1}{2}$

③

	DEB			Suiv		
E	0	1	(	#	)	
X	+	ε		#	)	
R	0	1	(	#	)	+
Y	0	1	(	#	)	+
S	0	1	(	#	)	+ 0 1 (
Z	*	ε		#	)	+ 0 1 (
T	0	1	(	#	)	+ 0 1 ( *

④

	0	1	+	*	(	)	#
E	RX	RX			RX		
X			+RX			ε	ε
R	SY	SY			SY		
Y	R	R	ε		R	ε	ε
S	TZ	TZ			TZ		
Z	ε	ε	ε	*Z	ε	ε	ε
T	0	1			(ε)		

$\frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2}$

⑤

Pile	chaîne	Action
# E	0* + 1 #	E → RX
# X R	0* + 1 #	R → SY
# X Y S	0* + 1 #	S → TZ
# X Y Z T	0* + 1 #	T → 0
# X Y Z 0	0* + 1 #	av + dépiler
# X Y Z	* + 1 #	Z → *Z
# X Y Z *	* + 1 #	av + dépiler
# X Y Z	+ 1 #	Z → ε

$\frac{1}{8} \times 20$

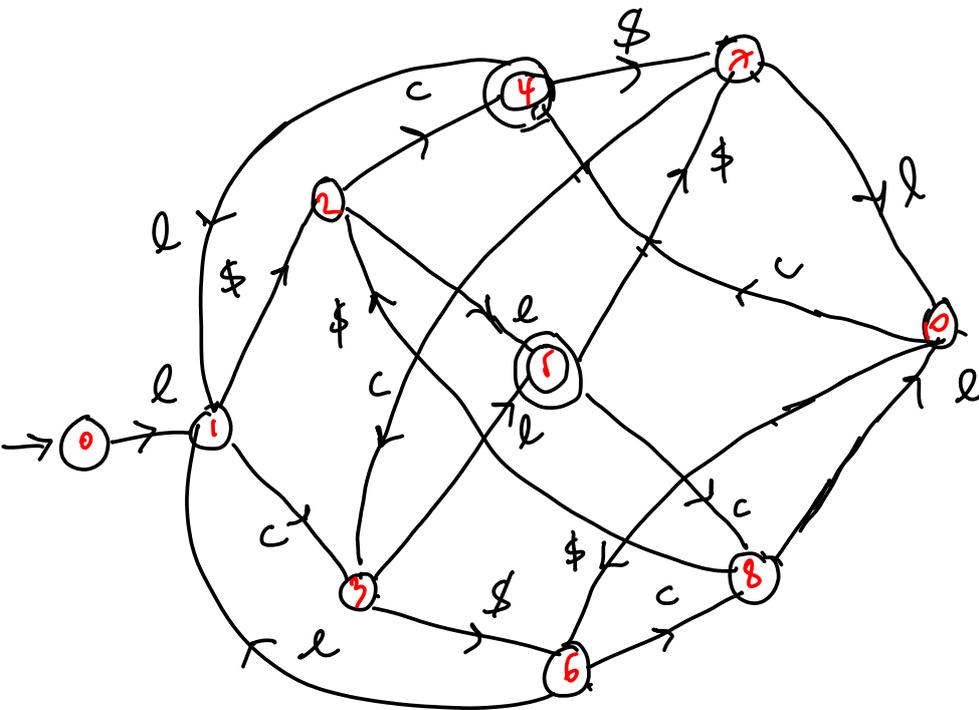
$2 \frac{1}{2}$

$\# x y$   
 $\# x$   
 $\# x R +$   
 $\# x R$   
 $\# x y S$   
 $\# x y z T$   
 $\# x y z 1$   
 $\# x y z$   
 $\# x y$   
 $\# x$   
 $\#$

$+ 1 \#$   
 $+ 1 \#$   
 $+ 1 \#$   
 $1 \#$   
 $1 \#$   
 $1 \#$   
 $1 \#$   
 $\#$   
 $\#$   
 $\#$   
 $\#$

$y \rightarrow \epsilon$   
 $x \rightarrow + R x$   
 $A v + \text{dépiler}$   
 $R \rightarrow S y$   
 $S \rightarrow T z$   
 $T \rightarrow 1$   
 $A v + \text{dépiler}$   
 $z \rightarrow \epsilon$   
 $y \rightarrow \epsilon$   
 $x \rightarrow \epsilon$   
 Acceptor

Bonus l'AFO de l'exo I minimal



# Bibliographie

- [1] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers : Principles, Techniques, and Tools*. Addison-Wesley series in computer science / World student series edition. Addison-Wesley, 1986. (Cité en pages 3, 6, 28, 32, 51, 55 et 60.)
- [2] Alex Aiken. Compilers : Online course. <http://openclassroom.stanford.edu/MainFolder/CoursePage.php?course=Compilers>, 2020. (Cité en pages 3, 34, 35 et 55.)
- [3] C Scott Ananian, Frank Flannery, Dan Wang, Andrew W Appel, and Michael Petter. Cup : Construction of useful parsers. <http://www2.cs.tum.edu/projects/cup/>, 2020. (Cité en page 85.)
- [4] David Beazley. Ply (python lex-yacc). <https://www.dabeaz.com/ply/>, 2020. (Cité en pages 42 et 85.)
- [5] Janusz A. Brzozowski. Derivatives of regular expressions. *J. ACM*, 11(4) :481–494, 1964. (Cité en pages 39 et 41.)
- [6] Noam Chomsky. Three models for the description of language. *IRE Trans. Inf. Theory*, 2(3) :113–124, 1956. (Cité en page 13.)
- [7] Keith D. Cooper and Linda Torczon. *Engineering a Compiler*. Morgan Kaufmann, 2004. (Cité en pages 3 et 28.)
- [8] Robert W. Sebesta David Watt, Deryck Brown. *Programming Language Processors in Java : Compilers and Interpreters AND Concepts of Programming Languages*. Prentice Hall Press, USA, 2007. (Cité en pages iii, 5 et 6.)
- [9] Frank DeRemer. Simple lr(k) grammars. *Commun. ACM*, 14(7) :453–460, 1971. (Cité en page 77.)
- [10] Frank DeRemer and Thomas J. Pennello. Efficient computation of LALR(1) look-ahead sets. In Stephen C. Johnson, editor, *Proceedings of the 1979 SIG-PLAN Symposium on Compiler Construction, Denver, Colorado, USA, August 6-10, 1979*, pages 176–187. ACM, 1979. (Cité en page 84.)
- [11] Albrecht Wöß Hanspeter Mössenböck, Markus Löberbauer. The compiler generator coco/r. <http://www.ssw.uni-linz.ac.at/coco/>, 2020. (Cité en page 85.)
- [12] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to automata theory, languages, and computation, 3rd Edition*. Pearson international edition. Addison-Wesley, 2007. (Cité en pages 10, 24, 37, 55, 77 et 89.)
- [13] Stephen C. Johnson. Yacc : Yet another compiler compiler. Technical report, Bell Laboratories, 1979. (Cité en page 85.)

- [14] S. C. Kleene. Representation of events in nerve nets and finite automata. In Claude Shannon and John McCarthy, editors, *Automata Studies*, pages 3–41. Princeton University Press, Princeton, NJ, 1956. (Cité en pages 23 et 32.)
- [15] Gerwin Klein, Steve Rowe, and Régis Décamps. Jflex. <https://jflex.de/>, 2020. (Cité en page 42.)
- [16] Donald E. Knuth. On the translation of languages from left to right. *Inf. Control.*, 8(6) :607–639, 1965. (Cité en pages 71 et 72.)
- [17] Eric Lehman. *Mathematics for Computer Science*. Samurai Media Limited, London, GBR, 2017. (Cité en page 10.)
- [18] John R. Levine. *flex and bison - Unix text processing tools*. O'Reilly, 2009. (Cité en pages 32, 33, 42 et 85.)
- [19] P. M. Lewis and R. E. Stearns. Syntax-directed transduction. *J. ACM*, 15(3) :465–488, 1968. (Cité en page 63.)
- [20] P. Lucas. The structure of formula-translators. *ALGOL Bull.*, (Suppl. 16) :1–27, 1961. (Cité en page 60.)
- [21] E. Schmidt M. E. Lesk. Lex - a lexical analyzer generator. Technical report, AT&T Bell Laboratories, Murray Hill, New Jersey 07974, 1975. (Cité en page 42.)
- [22] Scott Owens, John H. Reppy, and Aaron Turon. Regular-expression derivatives re-examined. *J. Funct. Program.*, 19(2) :173–190, 2009. (Cité en page 39.)
- [23] Vern Paxson, Will Estes, and John Millaway. Flex : fast lexer; main page. <https://www.gnu.org/software/flex/>, 2020. (Cité en page 42.)
- [24] Michael O. Rabin and Dana S. Scott. Finite automata and their decision problems. *IBM J. Res. Dev.*, 3(2) :114–125, 1959. (Cité en page 19.)
- [25] Thomas W. Reps. "maximal-munch" tokenization in linear time. *ACM Trans. Program. Lang. Syst.*, 20(2) :259–273, 1998. (Cité en page 35.)
- [26] Daniel J. Rosenkrantz and Richard Edwin Stearns. Properties of deterministic top-down grammars. *Inf. Control.*, 17(3) :226–256, 1970. (Cité en page 63.)
- [27] Michael L. Scott. *Programming Language Pragmatics (3. ed.)*. Academic Press, 2009. (Cité en pages iii, 3, 4, 6 et 8.)
- [28] Tony Stubblebine. *Regular expression - pocket reference : regular expressions for Perl, C, PHP, Python, Java, and .NET*. O'Reilly, 2003. (Cité en page 32.)
- [29] Ken Thompson. Regular expression search algorithm. *Commun. ACM*, 11(6) :419–422, 1968. (Cité en page 38.)
- [30] Jean-Paul Tremblay and Paul G. Sorenson. *Theory and Practice of Compiler Writing*. McGraw-Hill, Inc., USA, 1st edition, 1985. (Cité en page 3.)
- [31] Sreeni Viswanadha and Sriram Sankar. Javacc, the most popular parser generator for use with java applications. <https://javacc.github.io/javacc/faq.html>, 2020. (Cité en page 85.)