

Chapitre 2

Analyse lexicale

Partie I

© Compilation / M.MOURCHID, SMI/2021-2022

1

Sommaire

Rôle de l'analyseur lexical

Les unités lexicales

Attributs des unités lexicales

Erreurs lexicales

Langages réguliers

Diagramme de transition

Automates finis

© Compilation / M.MOURCHID, SMI/2021-2022

2

Rôle de l'analyseur lexical

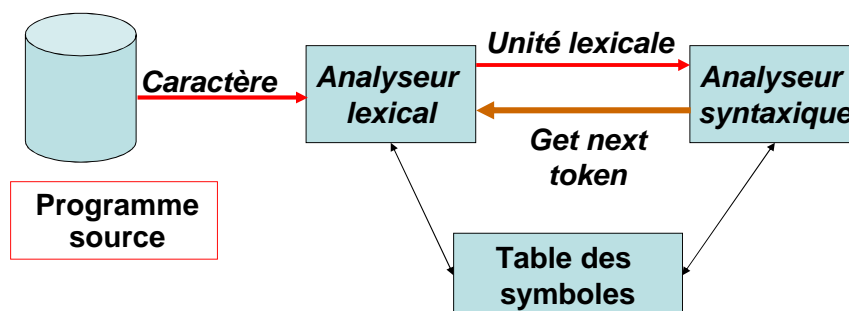
- L'analyse lexicale constitue la première étape d'un compilateur.
- Elle consiste à lire des **caractères d'entrée** et de produire une suite de **mots** que l'on appelle des **unités lexicales**.
- Elle permet également la détection et l'élimination des **caractères superflus** (commentaires, tabulations, fin de ligne, ...).
- Elle gère aussi les numéros de lignes dans le programme source pour pouvoir associer à chaque erreur son numéro de ligne.

© Compilation / M.MOURCHID, SMI/2021-2022

3

Interaction avec les autres modules

- Interaction entre un analyseur lexical et un analyseur syntaxique.



© Compilation / M.MOURCHID, SMI/2021-2022

4

Terminologie

Trois termes majeurs sont utilisés en analyse lexicale :

- **Unité lexicale** : un couple constitué d'un **nom d'unité lexicale** et d'une **valeur d'attribut** optionnelle.

Exemple : Mot clé, identificateur, opérateur, ...

Les **noms d'unité lexicale** sont les **symboles d'entrée** de l'analyseur syntaxique.

- **Motif** : est une description de la forme que les lexèmes d'une unité lexicale peuvent prendre. Il peut être simple ou complexe.
- **Lexème** : est une séquence de caractères dans le programme source qui est reconnue par le **motif** d'une **unité lexicale**.

© Compilation / M.MOURCHID, SMI/2021-2022

5

Les unités lexicales

- Exemples d'unités lexicales :

- **Mots clés** : ce sont des mots réservés au langage de programmation.
 - Exemples (langage C) : if, else, while, do, for, main, char, int, float, double, sizeof, ...etc.
- **Identificateurs** : un identificateur est un nom propre donné par le programmeur pour désigner une entité de son programme (une variable, une fonction ou une étiquette). En C, un identificateur est une suite de caractères qui commence par une lettre ([a...z] ou [A...Z] ou « _ ») et ne contient que des lettres et/ou des chiffres [0...9]).

© Compilation / M.MOURCHID, SMI/2021-2022

6

Les unités lexicales

- Exemples d'unités lexicales (suite) :
 - Symboles spéciaux ou délimiteurs : ce sont des séparateurs d'unités lexicales.
 - Exemples (langage C):
 - point virgule « ; »
 - caractère blanc ou espace « »
 - tabulation
 - saut de ligne
 - accolades « { } »
 - guillemets « " »
 - apostrophe « ' »

© Compilation / M.MOURCHID, SMI/2021-2022

7

Les unités lexicales

- Exemples d'unités lexicales (suite) :
 - Opérateurs : des symboles pour des opérations de base offertes par le langage de programmation.
 - Exemples (opérateurs du langage C) : +, -, *, /, %, >, >=, ==, <=, <, !=, &, |, ^, ~, >>, <<, &&, ||, !, ?:, =, +=, -=, *=, /=, %=, &=, |=, ^=, ~=, >>=, <<=, &&=, ||=, [], (), ., ->, & (adresse), sizeof.
 - Nombres : suite de chiffres avec éventuellement un signe (+ ou -), un point décimal « . ».
 - Commentaires : chaînes de caractères délimitées entre deux symboles spéciaux. Un commentaire est un texte à ignorer.

© Compilation / M.MOURCHID, SMI/2021-2022

8

Attributs des unités lexicales

- L'analyseur lexical regroupe les informations sur les unités lexicales dans des **attributs** qui leur sont associés.
- Les unités lexicales influent sur les décisions de l'analyseur syntaxique.
- Les attributs influent sur la traduction des unités lexicales.
- En général, une unité lexicale a un seul attribut implanté comme un pointeur vers l'entrée de la table des symboles dans laquelle l'information sur l'unité lexicale est stockée.

© Compilation / M.MOURCHID, SMI/2021-2022

9

Attributs des unités lexicales

- Comme exemples d'attributs, considérons l'instruction du langage Fortran : « $E = M * C ** 2$ ». Nous avons les unités lexicales et les attributs suivants :
 - <**ID**, pointeur vers l'entrée de la TS associée à E>
 - <**OP_AFFECT**>
 - <**ID**, pointeur vers l'entrée de la TS associée à M>
 - <**OP_MULT**>
 - <**ID**, pointeur vers l'entrée de la TS associée à C>
 - <**OP_PUISS**>
 - <**NUMBER**, valeur entière 2>

© Compilation / M.MOURCHID, SMI/2021-2022

10

Attributs des unités lexicales

- Remarque :
 - Pour certaines unités lexicales, aucune valeur d'attribut n'est nécessaire : le premier composant suffit pour identifier le lexème en question.
- Exemples :
 - Pour ce qui est des symboles comme $>$, $>=$, $=$, $<=$, $<$ et $<>$ (langage Pascal), l'analyseur syntaxique a juste besoin de savoir que cela correspond à l'unité lexicale OP_REL (opérateur relationnel).
 - C'est seulement lors de la génération de code que l'on aura besoin de distinguer $<$ de $>=$ (par exemple).

© Compilation / M.MOURCHID, SMI/2021-2022

11

Attributs des unités lexicales

- Exemples (suite) :
 - Pour ce qui est des identificateurs, l'analyseur syntaxique a juste besoin de savoir que c'est l'unité lexicale ID.
 - Mais le générateur de code, lui, aura besoin de l'adresse mémoire de la variable correspondant à cet identificateur.
 - L'analyseur sémantique aura aussi besoin du type de la variable pour vérifier que les expressions sont sémantiquement correctes.

© Compilation / M.MOURCHID, SMI/2021-2022

12

Erreurs lexicales

- Peu d'erreurs sont détectables au seul niveau lexical, car un analyseur lexical a une vision très locale du programme source (il ne voit que les unités lexicales).
- Les erreurs se produisent lorsque l'analyseur est confronté à une suite de caractères qui ne correspond à aucun des motifs d'unités lexicales qu'il a à sa disposition : on dit qu'il y a une erreur lexicale.
- Par exemple, dans un programme C, si l'analyseur rencontre la chaîne de caractères « **esle** », s'agit-il du mot clé « **else** » mal orthographié ou d'un simple identificateur ?
 - Dans ce cas précis, comme « esle » correspond au motif de l'unité lexicale ID, l'analyseur lexical ne détectera pas d'erreur, et transmettra à l'analyseur syntaxique qu'il a reconnu un ID.
 - S'il s'agit du mot clé mal orthographié, c'est l'analyseur syntaxique qui détectera alors une erreur.

© Compilation / M.MOURCHID, SMI/2021-2022

13

Alphabet, mots et langage

- Les lexèmes sont avant tout des chaînes de caractères, c'est-à-dire des mots construits à partir de l'alphabet du langage de programmation.
- Les unités lexicales sont les noms des classes des lexèmes.
- Chaque unité lexicale est définie par une règle (un motif).
- Cette règle est utilisée par l'analyseur lexical pour faire correspondre un lexème à son unité lexicale : c'est le processus de reconnaissance d'unités lexicales.

© Compilation / M.MOURCHID, SMI/2021-2022

14

Alphabet

- On appelle **alphabet** Σ un ensemble fini non vide.
Les éléments d'un alphabet sont appelés **lettres (ou symboles)**
- Exemple d'alphabets
 - $\Sigma = \{a, b, \dots, z\}$
 - $\Sigma = \{\alpha, \beta, \delta, \dots, \varphi, +, *, /, =\}$
 - $\Sigma = \{\text{if, then, else, while, begin, end}\}$

Mots

- Une suite de symboles, appartenant à un alphabet Σ , mis bout à bout est appelé un **mot** (ou une chaîne) sur Σ .
Par exemple : 01001100 , est un mot construit sur l'alphabet $B = \{0, 1\}$.
- Le nombre de symboles entrant dans la composition d'un mot α est appelé la **longueur** de α , notée $|\alpha|$.
Par exemple : La longueur du mot 01001100 est $|01001100| = 8$.
- On distingue un mot particulier, dont la longueur vaut zéro, que l'on appellera le **mot vide**. Ce mot est conventionnellement représenté par le symbole ϵ .

Langages formels

| Alphabet d | Mots de d |
|--|---|
| $d = \{a, b, \dots, z\}$ | « abab », « surf », « cascade », ... |
| $d = \{r, s, u, \dots, \{, +, *, /, =\}$ | « $r = s * 2 + \{$ », « $+r u = *$ », ... |
| $d = \{\text{if, then, else, while, begin, end}\}$ | (if then else), (begin begin end) |

© Compilation / M.MOURCHID, SMI/2021-2022

17

Langages formels

- La **concaténation** de deux mots r et s notée rs est le mot obtenu en juxtaposant les symboles de s à la suite de ceux de r .
- Exemple : Si $r = abra$ et $s = cadabra$ alors $rs = abracadabra$.
- Le mot m^n représente la concaténation de m avec lui même n fois :

$$\overbrace{mm \dots m}^n = m^n$$

© Compilation / M.MOURCHID, SMI/2021-2022

18

Langages formels

- Étant donné trois mots r , s et x définis sur un alphabet Σ , on dit que r est un **préfixe** du mot rs , s un **suffixe** du mot rs et s une sous-chaîne de rsx . v est un préfixe, un suffixe et une sous-chaîne de tout mot. Si $r \neq s$ et r est un préfixe (ou suffixe) de s , alors on dit que r est un **préfixe** (ou **suffixe**) **propre** de s .
- L'ensemble de tous les mots que l'on peut construire sur un alphabet Σ , y compris ϵ , est noté Σ^* .

© Compilation / M.MOURCHID, SMI/2021-2022

19

Langages formels

- Un **langage** sur un alphabet Σ est un ensemble de mots construits sur Σ . Tout langage défini sur Σ est donc une partie de Σ^* .
- On distingue deux langages particuliers qui sont définis indépendamment de tout alphabet. Il s'agit du **langage vide** (noté \emptyset) et du **langage composé** uniquement du **mot vide** ($\{\epsilon\}$). Σ^* peut être vu comme le langage composé des symboles de Σ .

© Compilation / M.MOURCHID, SMI/2021-2022

20

Exemples de Langages

- Exemple de langages sur $d = \{a, b, \dots, z\}$
 - $L_0 = \{a\}, \quad L_1 = \{aa, ab\},$
 - $L_2 = \{waw \mid w \in d^*\},$
 - $L_3 = \{w \in d^* \mid |w|_a \leq 10 \}$

Où $|w|_x \in \mathbb{N} : d^* \rightarrow \mathbb{N}$ est la fonction qui calcule le nombre d'occurrence de la lettre x de d dans un mot de d^* .

Opérations sur les langages formels

Les langages étant des ensembles, on peut leur appliquer les opérations définies sur ces derniers:

- L'**union** de deux langages L_1 et L_2 , est le langage, noté $L_1 \cup L_2$ constitué des mots appartenant à L_1 ou à L_2 .

$$L_1 \cup L_2 = \{x \mid x \in L_1 \text{ ou } x \in L_2\}$$

- L'**intersection** de L_1 et L_2 , est le langage, noté $L_1 \cap L_2$ constitué des mots appartenant à L_1 et à L_2 .

$$L_1 \cap L_2 = \{x \mid x \in L_1 \text{ et } x \in L_2\}$$

- La **différence** de L_1 et L_2 est le langage, noté $L_1 - L_2$, constitué des mots appartenant à L_1 et n'appartenant pas à L_2 .

$$L_1 - L_2 = \{x \mid x \in L_1 \text{ et } x \notin L_2\}$$

Opérations sur les langages formels

- La différence de Σ^* et de L , noté \overline{L} est appelé le **complément** du langage L .

$$\overline{L} = \{x \in \Sigma^* \text{ et } x \notin L\}$$

- On définit de plus l'opération de **concaténation de deux langages** : la concaténation de deux langages L_1 et L_2 est le langage noté L_1L_2 composé des mots xy tels que $x \in L_1$ et $y \in L_2$.

$$L_1L_2 = \{xy \mid x \in L_1 \text{ et } y \in L_2\}$$

- on note L^n la concaténation de L avec lui même n fois :

$$L^n = \{x_1x_2 \dots x_n \mid x_i \in L \text{ } 1 \leq i \leq n\}$$

Opérations sur les langages formels

- On définit enfin la **fermeture de Kleene** du langage L , notée L^* de la façon suivante :

$$L^* = \bigcup_{k \geq 0} L^k$$

$$L^+ = L^* - \{\varepsilon\}$$

Opérations sur les langages formels

Proposition

Les identités suivantes sont vraies pour tous les langages X, Y, Z :

- $X(YZ) = (XY)Z$
- $Xv = vX = X$
- $X(Y \hat{\cup} Z) = XY \hat{\cup} XZ$ et $(X \hat{\cup} Y)Z = XZ \hat{\cup} YZ$
- $XX^* \hat{\cup} v = X^*$
- $(X^*)^* = X^*$

Langages réguliers

- Le problème préliminaire qui se pose à propos des langages est:

comment décrire un langage ?

- La question qui se pose ensuite est:

étant donné une phrase (mot) est-ce qu'elle appartient au langage ?

Dans cette partie nous définissons les expressions régulières qui permettent de décrire un certain type de langages : les langages réguliers.

Langages réguliers

- Un **langage régulier** L sur un alphabet Σ est défini récursivement comme suit:
 - ϵ est un langage régulier sur Σ
 - $a \in \Sigma$, a est un langage régulier sur Σ
 - si R est un langage régulier sur Σ alors R^k et R^* sont des langages réguliers sur Σ
 - si R_1 et R_2 sont des langages réguliers alors $R_1 \cup R_2$ et $R_1 R_2$ sont des langages réguliers sur Σ .

© Compilation / M.MOURCHID, SMI/2021-2022

27

Expressions régulières (ER)

- Une expression régulière est une notation permettant de spécifier un modèle.
- Le modèle reconnaît un ensemble de mots (un langage) sur un alphabet.
- L'analyseur lexical est chargé de reconnaître les unités lexicales du texte source.
- Chacun de ces lexèmes est un langage dans la plupart des cas, un langage régulier.

© Compilation / M.MOURCHID, SMI/2021-2022

28

Expressions régulières (ER)

Exemple:

Soit $\Sigma = \{a, b\}$

- $a \mid b$ dénote l'ensemble $\{a, b\}$
- $(a \mid b)(a \mid b)$ dénote $\{aa, ab, bb, ba\}$
- a^* dénote l'ensemble des mots formés d'un nombre quelconque (éventuellement nul) de a : $\{, a, aa, aaa, \dots\}$
- $(a \mid b)^* = (b \mid a)^* = (a^* \mid b^*)^* = \Sigma^*$ dénote l'ensemble des mots formés d'un nombre quelconque (éventuellement nul) de a ou de b
- $a \mid a^*b$ dénote l'ensemble contenant le mot a et tous les mots constitués d'un nombre quelconque (éventuellement nul) de a suivi de b .

© Compilation / M.MOURCHID, SMI/2021-2022

29

Expressions régulières (ER)

- Soit Σ un alphabet, les expressions régulières et les langages qu'elles décrivent sont définis comme suit :
 - \emptyset , a est une expression régulière qui décrit le langage ' a ',
 - ϵ (mot vide) est une expression régulière qui décrit le langage ' ϵ '
 - \emptyset (ensemble vide) est une expressions régulières qui décrit le langage vide.
 - Si r et s sont des expressions régulières décrivant respectivement les langages A et B alors
 - $(r \mid s)$ (ou $r + s$ est l'expression régulière décrivant le langage $A \cup B$
 - (rs) est l'expression régulière décrivant le langage AB ,
 - r^* est l'expression régulière décrivant le langage A^*
- Nulle autre expression n'est une expression régulière

© Compilation / M.MOURCHID, SMI/2021-2022

30

Exemples d'expressions régulières

- **Alphabétique** = ('a' | 'b' | ... | 'z' | 'A' | 'B' | ... | 'Z')
- **Numérique** = (0 | ... | 9)
- **Opérateurs** = (+ | - | / | * | = | <= | >= | < | >)
- **Naturel** = Numérique⁺
- **Entier** = (+ | - | v) Naturel
- **Identificateur** =
(Alphabétique(Alphabétique|Numérique)*)

Spécification des unités lexicales

- Les **motifs** de lexèmes sont représentés par des **expressions régulières**.
- Les détails sur les **expressions régulières** se trouvent dans le cours **Théorie des langages**.
- **Exemple** : On décrit l'ensemble des identificateurs C par l'expression régulière :
 - **lettre_** = [A-Za-z_]
 - **chiffre** = [0-9]
 - **id** = lettre_(lettre_|chiffre)*

Reconnaissance des unités lexicales

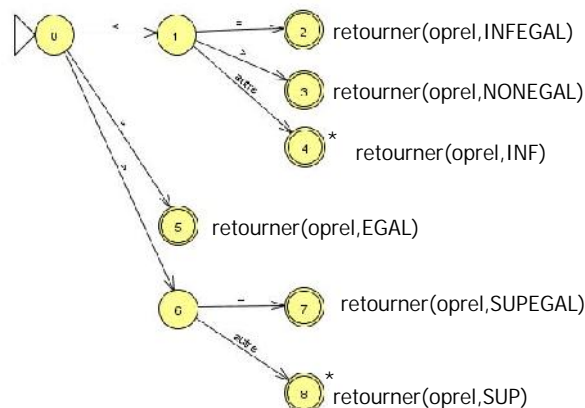
- Grammaire pour les instructions de branchement :
 - instr \rightarrow if expr then instr
 | if expr then instr else instr
 - expr \rightarrow terme oprel terme
 | terme
 - terme \rightarrow id
 | nbre
- Motif de l'unité lexicale oprel
 - oprel = (< | > | <= | >= | = | <>)
- On commence par convertir les **motifs** en **diagrammes de transition**.
- Il peut y avoir plusieurs diagrammes de transitions, chacun spécifiant une unité lexicale.

© Compilation / M.MOURCHID, SMI/2021-2022

33

Diagramme de transition

Diagramme de transition pour **oprel**



© Compilation / M.MOURCHID, SMI/2021-2022

34

Architecture d'un analyseur lexical

- But : construire un analyseur lexical sur la base d'un diagramme de transition.
- Une variable `etat` sert à conserver le numéro de l'état courant.
- Une fonction `car_suiv()` récupère le caractère suivant dans l'entrée.
- Si le caractère suivant n'est pas reconnu, on appelle la fonction `echec()`.
- Si l'état contient * alors il faut reculer le pointeur d'une position (`reculer()`).

© Compilation / M.MOURCHID, SMI/2021-2022

35

Implémentation de oprel

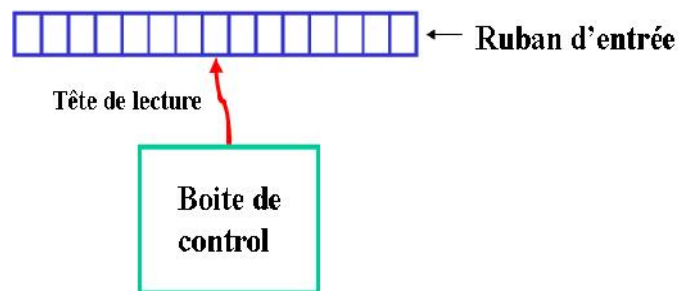
```
int lire_op_rel() {
    while(1) {
        switch(etat) {
            case 0 : c=car_suiv();
                    if(c=='<') etat=1;
                    else if(c=='=') etat=5;
                    else if(c=='>') etat=6;
                    else echec();
                    break;
            . . . . .
            case 8 : reculer();
                    return SUP; }
        }
    }
```

© Compilation / M.MOURCHID, SMI/2021-2022

36

Automates finis

Un automate «lit» un mot écrit sur son ruban d'entrée. Il part d'un état initial et à chaque lettre lue, il change d'état. Si, à la fin du mot, il est dans un état final, on dit qu'il reconnaît le mot lu, ou qu'il l'accepte.



© Compilation / M.MOURCHID, SMI/2021-2022

37

Automates finis déterministes

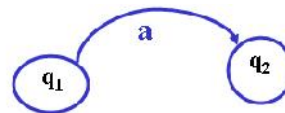
- Un automate d'états finis A est défini comme un 5-uplet $A = \langle Q, d, u, q_0, F \rangle$ où :
 - Q est l'ensemble des états,
 - d est un alphabet (un ensemble de symboles),
 - $u : Q \times d \rightarrow Q$ est la fonction de transition de l'automate,
 - q_0 est l'état initial de l'automate,
 - $F \subseteq Q$ est l'ensemble des états finaux de l'automate.
- Les automates peuvent être représentés par des graphes où les états sont les noeuds du graphe et les arcs représentent la fonction de transition.



État initial



État final



Si $\delta(q_1, a) = q_2$

© Compilation / M.MOURCHID, SMI/2021-2022

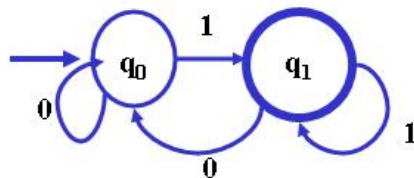
38

Exemple d'automate fini

Soit $A = \langle Q, d, u, q_0, F \rangle$ avec

- $d = \{0, 1\}$
- $S = \{q_0, q_1\}$
- $u = \{(q_0, 0, q_0), (q_0, 1, q_1), (q_1, 1, q_1), (q_1, 0, q_0)\}$
- q_0 est l'état initial
- $F = \{q_1\}$

A peut être représenté par le graphe suivant:



| | 0 | 1 |
|-------|-------|-------|
| q_0 | q_0 | q_1 |
| q_1 | q_0 | q_1 |

© Compilation / M.MOURCHID, SMI/2021-2022

39

Langage accepté par un AFD

- On dit qu'un automate d'états finis $A = \langle Q, d, u, q_0, F \rangle$ **accepte** le mot x , composé des symboles $x_1x_2x_3 \dots x_n$ ssi $u(u(\dots u(u(q_0, x_1), x_2), x_3), \dots), x_{n-1}), x_n) \in F$
- Intuitivement $x_1x_2x_3 \dots x_n$ correspond aux étiquettes des arcs formant un chemin à travers les transitions de l'automate, commençant par q_0 et se terminant par $q_k \in F$.
- A chaque étape, x_i correspond à l'étiquette du $i^{\text{ème}}$ arc du chemin.
- L'ensemble des mots pouvant être acceptés par un automate A est appelé le **langage accepté par A** , noté $L(A)$.

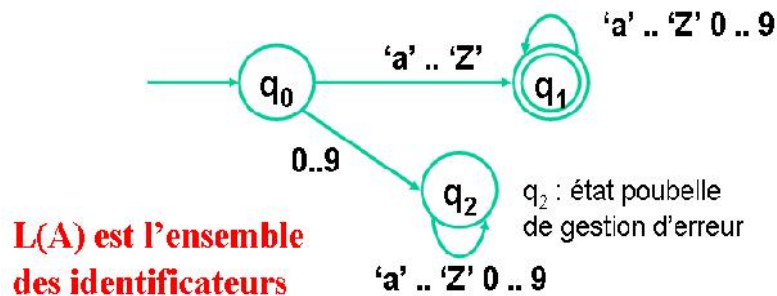
© Compilation / M.MOURCHID, SMI/2021-2022

40

Exemple: Langage accepté par un AFD

Soit A l'automate donné par:

- $d = \{ 'a', \dots, 'z', 'A', \dots, 'Z', 0, \dots, 9 \}$
- $Q = \{q_0, q_1, q_2\}$ avec comme état initial q_0
- $F = \{q_2\}$



© Compilation / M.MOURCHID, SMI/2021-2022

41

Automates non-déterministes (AFN)

- Il est difficile d'écrire simplement les automates correspondant à une ER compliquée, l'automate non-déterministe simplifie cette tâche
- Un **automate non-déterministe** $A = \langle Q, d, u, q_0, F \rangle$ est un automate qui présente la particularité suivante :

Il possède une fonction de transition étendue

$$u : Q \times d \rightarrow \{v\} \subseteq Q$$

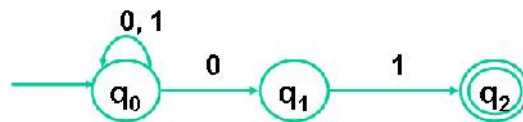
1. plusieurs arcs reconnaissant le même symbole peuvent « sortir » du même état
2. Il peut y avoir des transitions, appelées (ϵ -transitions), qui ne correspondent à aucune reconnaissance de caractères.

© Compilation / M.MOURCHID, SMI/2021-2022

42

Exemple de AFN (sans ε -transitions)

Soit l'automate donné par le graphe:



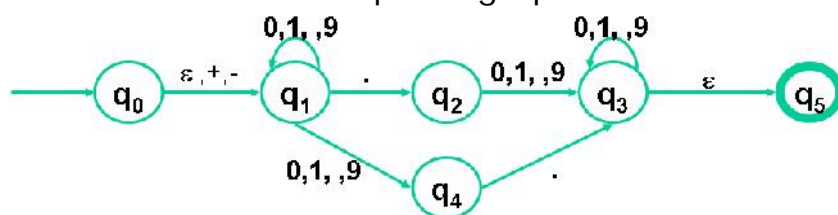
| Q | 0 | 1 |
|----------------|------------------------------------|-------------------|
| q ₀ | {q ₀ , q ₁ } | {q ₀ } |
| q ₁ | à | {q ₂ } |
| q ₂ | à | à |

© Compilation / M.MOURCHID, SMI/2021-2022

43

Exemple de AFN (avec ε -transitions)

Soit l'automate donné par le graphe:



| Q | v | +, - | . | 0,1,,9 |
|----------------|-------------------|-------------------|-------------------|------------------------------------|
| q ₀ | {q ₁ } | {q ₁ } | à | à |
| q ₁ | à | à | {q ₂ } | {q ₁ , q ₄ } |
| q ₂ | à | à | à | {q ₃ } |
| q ₃ | {q ₅ } | à | à | {q ₃ } |
| q ₄ | à | à | {q ₃ } | à |
| q ₅ | à | à | à | à |

© Compilation / M.MOURCHID, SMI/2021-2022

44

Transformer un AFN en un AFD

Deux cas sont possibles:

- 1^{er} cas : l'automate non-déterministe est sans v-transition
- 2^{ème} cas : l'automate non-déterministe contient des v-transitions

Transformer un AFN (sans ε -transitions) en un AFD

Soit $A_N = \langle Q_N, d_N, u_N, q_{N0}, F_N \rangle$ un automate non-déterministe sans v-transitions, il existe un automate déterministe

$A_D = \langle Q_D, d_D, u_D, q_{D0}, F_D \rangle$ tel que:

$$L(A_N) = L(A_D)$$

$$Q_D \subseteq c(Q_N)$$

tous les éléments de $c(Q_N)$ qui sont accessibles depuis l'état initial q_{D0}

$$d_D = d_N$$

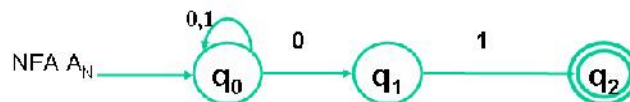
$$u_D(q_D \in Q_D, a \in d_D) = \hat{a}_{q_N \in q_D} u_N(q_N, a)$$

$$q_{D0} = \{q_{N0}\}$$

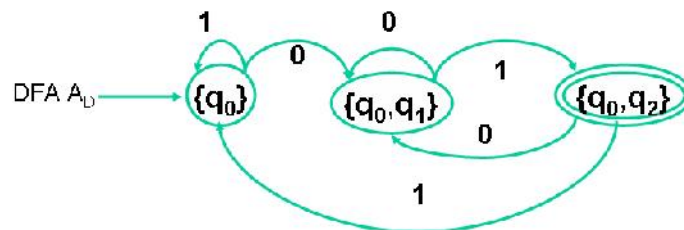
$$F_D = \{q \in Q_D \mid (q \in F_N) \vee \dots\}$$

Exemple : Transformer un AFN (sans ε -transitions) en un AFD

L'automate donné par le graphe suivant:



est équivalent à celui donné par le graphe suivant

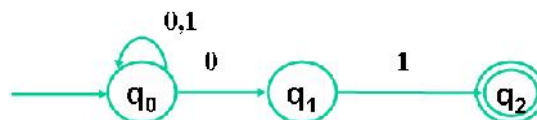


$$L(A_N) = L(A_D) = \{w = (0|1)^*01\}$$

© Compilation /
M.MOURCHID, SMI

47

Exemple : Transformer un AFN (sans ε -transitions) en un AFD



D

■ état accessible depuis q_{D0}

| $\tilde{Y}(Q_N)$ | 0 | 1 |
|----------------------|----------------|----------------|
| \emptyset | \emptyset | \emptyset |
| $q_{D0} = \{q_0\}$ | $\{q_0, q_1\}$ | $\{q_0\}$ |
| $\{q_1\}$ | \emptyset | $\{q_2\}$ |
| $\{q_2\}$ | \emptyset | \emptyset |
| $\{q_0, q_1\}$ | $\{q_0, q_1\}$ | $\{q_0, q_2\}$ |
| $F_D = \{q_0, q_2\}$ | $\{q_0, q_1\}$ | $\{q_0\}$ |
| $\{q_1, q_2\}$ | \emptyset | $\{q_2\}$ |
| $\{q_0, q_1, q_2\}$ | $\{q_0, q_1\}$ | $\{q_0, q_2\}$ |

© Compilation /
M.MOURCHID, SMI

48

Transformer un AFN (avec ε -transitions) en un AFD

On définit la fonction **v-fermeture** : $Q \rightarrow \mathcal{P}(Q)$ comme suit:

1. $q \in v\text{-fermeture}(q)$
2. Si $q_1 \in u(q, v)$ alors $q_1 \in v\text{-fermeture}(q)$
3. Si $q_1 \in v\text{-fermeture}(q)$ ó $q_2 \in v\text{-fermeture}(q_1)$ alors $q_2 \in v\text{-fermeture}(q)$

On peut définir par extension

v-fermeture : $\mathcal{P}(Q) \rightarrow \mathcal{P}(Q)$ telle que:

$$v\text{-fermeture}(Q' \subseteq Q) = \bigcup_{q' \in Q'} v\text{-fermeture}(q')$$

$v\text{-fermeture}(q)$ regroupe tous les états accessibles depuis q sur des chemins constitués seulement de v -transitions

© Compilation /
M.MOURCHID, SMI

49

Transformer un AFN (avec ε -transitions) en un AFD

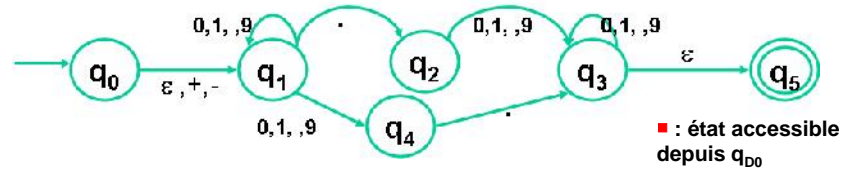
Soit $A_E = \langle Q_E, d_E, u_E, q_{E0}, F_E \rangle$ un automate non-déterministe avec v -transitions (v-NFA), il existe un automate déterministe (DFA) $A_D = \langle Q_D, d_D, u_D, q_{D0}, F_D \rangle$ tel que $L(A_E) = L(A_D)$

- $Q_D \subseteq v\text{-fermeture}(c(Q_E))$
tous les éléments de $v\text{-fermeture}(c(Q_E))$ qui sont accessibles depuis l'état initial q_{D0}
- $d_D = d_E$
- $u_D(q_D \in Q_D, a \in d_D) = v\text{-fermeture}(\bigcup_{q_E \in q_D} u_E(q_E, a))$
- $q_{D0} = \{v\text{-fermeture}(q_{E0})\}$
- $F_D = \{q \in Q_D \mid (q \supseteq F_E) \text{ ó } \emptyset\}$

© Compilation /
M.MOURCHID, SMI

50

Exemple : Transformer un ε -AFN en AFD



D

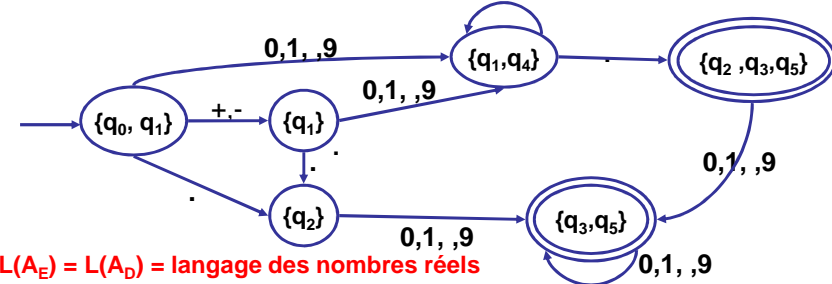
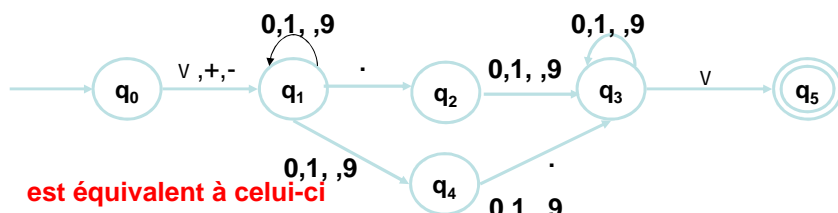
| Q_E | v-fermeture | $+, -$ | v-fermeture | . | v-fermeture | $0, 1, 9$ | v-fermeture |
|-----------|------------------------------|---------------|---------------|----------------|---------------------|----------------|----------------|
| $\{q_0\}$ | $q_{D0} = \{q_0, q_1\}$ | $\{q_1\}$ | $\{q_1\}$ | $\{q_2\}$ | $\{q_2\}$ | $\{q_1, q_4\}$ | $\{q_1, q_4\}$ |
| $\{q_1\}$ | $\{q_1\}$ | \rightarrow | \rightarrow | $\{q_2\}$ | $\{q_2\}$ | $\{q_1, q_4\}$ | $\{q_1, q_4\}$ |
| $\{q_2\}$ | $\{q_2\}$ | \rightarrow | \rightarrow | \rightarrow | \rightarrow | $\{q_3\}$ | $\{q_3, q_5\}$ |
| $\{q_3\}$ | $\{q_3, q_5\}$ | \rightarrow | \rightarrow | \rightarrow | \rightarrow | $\{q_3\}$ | $\{q_3, q_5\}$ |
| $\{q_4\}$ | $\{q_4\}$ | \rightarrow | \rightarrow | $\{q_3\}$ | $\{q_3, q_5\}$ | \rightarrow | \rightarrow |
| $\{q_5\}$ | $\{q_5\}$ | \rightarrow | \rightarrow | \rightarrow | \rightarrow | \rightarrow | \rightarrow |
| | $F_{D1} = \{q_1, q_4\}$ | \rightarrow | \rightarrow | $\{q_2, q_3\}$ | $\{q_2, q_3, q_5\}$ | $\{q_1, q_4\}$ | $\{q_1, q_4\}$ |
| | $F_{D2} = \{q_2, q_3, q_5\}$ | \rightarrow | \rightarrow | \rightarrow | \rightarrow | $\{q_3\}$ | $\{q_3, q_5\}$ |

© Compilation /
M.MOURCHID, SMI

51

Exemple : Transformer un ε -AFN en AFD

L'automate suivant



© Compilation /
M.MOURCHID, SMI

52

Minimiser un AFD

- Minimiser un AFD A revient à regrouper ses états équivalents de telle sorte à réduire son nombre d'états
- Deux états p, q de A sont dits équivalents ssi:
 $\exists w \in \Sigma^* \text{ tq } u(p, w) \text{ est un état final } \vee u(q, w) \text{ est un état final}$
 1. $w = \epsilon$, si $p \in F$ et $q \notin F$ alors p et q sont distinguables (non équivalents)
 2. $w \neq \epsilon$, ($u(p, w) = r$ et $u(q, w) = s$) \wedge (r et s sont distinguables) alors p et q sont distinguables

© Compilation /
M.MOURCHID, SMI

53

Minimiser un DFA

L'algorithme de minimisation construit une partition P de l'ensemble Q des états de l'automate déterministe, telle que : chaque état se trouvant dans le même sous ensemble a le même comportement sur les mêmes entrées.

L'algorithme est le suivant :

© Compilation /
M.MOURCHID, SMI

54

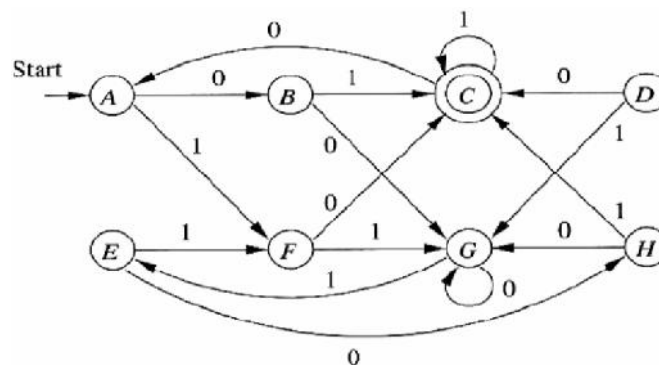
Minimiser un DFA

→ $P \in \{F, Q-F\}$ (où F est l'ensemble des états finaux)
 Tant que (P change)
 $T \in \emptyset$
 Pour chaque ensemble $p \in P$
 $T \leftarrow T \hat{\cup} \text{partition}(p)$
 $P \leftarrow T$
 Fin tant
 → Partition(p)
 Pour chaque $c \in p$
 Si c sépare p en ' p_1, \dots, p_k '
 Alors return (' p_1, \dots, p_k ')
 Fin pour
 Return(p)

© Compilation /
M.MOURCHID, SMI

55

Exemple: Minimiser un AFD



- $A \hat{=} E$, car $u(A,1)=u(E,1)=F$ et $u(A,0)=B$; $u(E,0)=H$ avec $B \hat{=} H$
- $B \hat{=} H$, car $u(B,1)=u(H,1)=C$ et $u(B,0)=u(H,0)=G$
- $\neg(A \hat{=} B)$, car $u(B,1)=C \in F$ et $u(A,1) \notin F$
- etc.

© Compilation /
M.MOURCHID, SMI

56

Exemple: Minimiser un DFA

| | | | | | | | |
|----------|----------|----------|----------|----------|----------|----------|----------|
| <i>B</i> | X | | | | | | |
| <i>C</i> | X | X | | | | | |
| <i>D</i> | X | X | X | | | | |
| <i>E</i> | | X | X | X | | | |
| <i>F</i> | X | X | X | | X | | |
| <i>G</i> | X | X | X | X | X | X | |
| <i>H</i> | X | | X | X | X | X | X |
| | <i>A</i> | <i>B</i> | <i>C</i> | <i>D</i> | <i>E</i> | <i>F</i> | <i>G</i> |

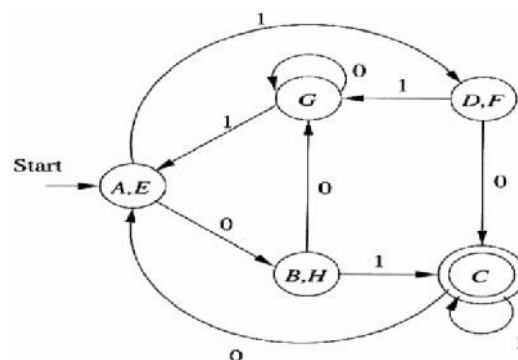
- $[s_i, s_j] = \times$ signifie que $\neg (s_i \hat{=} s_j)$ (i.e. s_i et s_j sont distinguables)

© Compilation /
M.MOURCHID, SMI

57

Exemple: Minimiser un AFD

L'automate minimal regroupant les états équivalents est



© Compilation /
M.MOURCHID, SMI

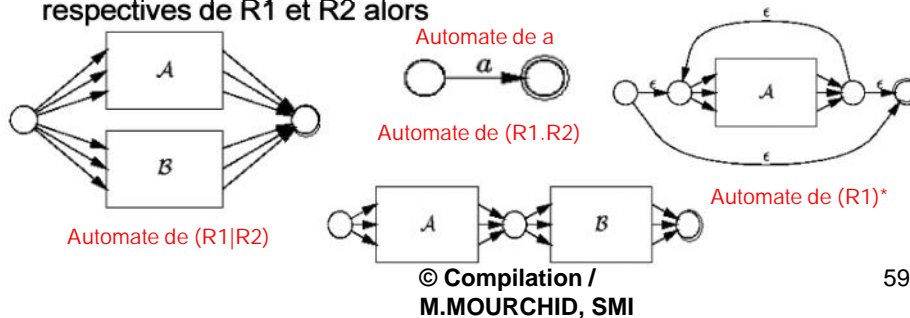
58

Equivalence des automates finis aux expressions régulières

Théorème

Un langage est régulier si et seulement s'il est généré par un automate fini.

- On suppose que les automates finis n'ont qu'un seul état final q_f avec $u(q_f, a) = \emptyset$ pour tout $a \in \Sigma$. De même, on peut supposer que $q_0 \notin u(q, a)$ pour tout $q \in Q$ et $a \in \Sigma$.
- Alors les automates finis reconnaissent les expressions régulières par les constructions suivantes : Si A et B sont des automates respectives de R1 et R2 alors



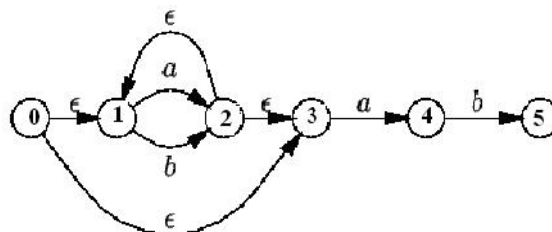
59

Equivalence des automates finis aux expressions régulières

D'après ce qui précède, toute expression régulière est reconnue par un automate.

Exemple

Soit l'expression régulière $(a+b)^*ab$, l'automate associé est:

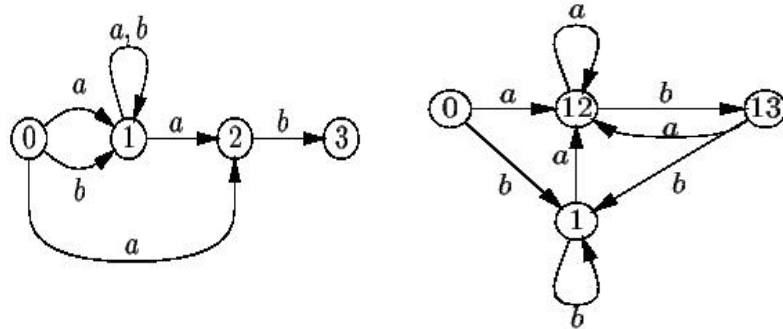


© Compilation / M.MOURCHID, SMI

60

Equivalence des automates finis aux expressions régulières

L'automate ainsi obtenu est non déterministe et il possède beaucoup de ϵ -transitions. En supprimant les ϵ -transitions et en déterminisant, on obtient:



© Compilation /
M.MOURCHID, SMI

61

Détermination du langage accepté par un automate fini

Réciproquement, pour obtenir le langage reconnu par un automate on utilise l'algorithme suivant :

Algorithme de Mac Naughton et Yamada

Principe : L'ensemble des transitions menant d'un état s à un état t est, pour un état q distinct, l'union de deux ensembles de chemins

- d'une part, les chemins menant de s à t sans passer par q ;
- d'autre part, les chemins menant de s à t en passant par q .

© Compilation /
M.MOURCHID, SMI

62

Détermination du langage accepté par un automate fini

ALGORITHME :

Soient $Aut = \langle \Sigma, Q, q_0, F, u \rangle$, $s \in Q$, $t \in Q$ et $P \subseteq Q$ (un sous-ensemble des états).

On note $W_{s,P,t}$ l'ensemble de tous les mots (chemins) $u \in \Sigma^*$ tels que u est l'étiquette d'un chemin menant de s à t et n'utilisant pas d'autres états que ceux de P comme intermédiaires.

Donc le langage reconnu par l'automate est :

$L(Aut) = \bigcup_{s \in D, t \in F} W_{s,Q,t}$ où $s \in D$ et $t \in F$ qui est l'union de tous les $W_{s,Q,t}$ possibles. D étant l'ensemble des états départ.

Détails

On démarre avec $P = Q$ et on cherche, pour tous les couples (s,t) tels que $s \in D$ et $t \in F$, le langage $W_{s,P,t}$. Ce langage est défini grâce à la formule générale suivante :

$$W_{s,P,t} = W_{s,P',t} \hat{\cup} W_{s,P',q} \cdot (W_{q,P',q})^* \cdot W_{q,P',t}$$

où P' est défini par $P' = P \setminus \{q\}$ (q étant un état que l'on choisit arbitrairement dans P afin de pouvoir avancer.

© Compilation /
M.MOURCHID, SMI

63

Détermination du langage accepté par un automate fini

Méthode de Kleene

Soit $A = \langle \Sigma, Q, q_0, F, u \rangle$ un automate fini non déterministe ne contenant pas d'v-transitions.

Pour tout état q , notons X_q l'ensemble des mots menant de l'état q à un état final. Chaque langage X_q satisfait l'équation linéaire

$$X_q = \bigcup_{q' \in N_u(q,a)} a X_{q'} \quad \text{si } q \in F$$

Pour obtenir une expression régulière du langage, il suffit de résoudre le système des équations linéaires ci-dessus par élimination de variables : toutes les équations qu'on manipule peuvent être mises sous la forme $X_q = K X_q + L$ où K et L sont des expressions linéaires des variables $X_{q'}$, $q' \neq q$;

© Compilation /
M.MOURCHID, SMI

64

Détermination du langage accepté par un automate fini

K peut être le langage vide (si le membre droit de l'équation ne contient pas la variable X_q), mais il ne contient pas le mot vide (car l'automate ne possède pas d'è-transitions) ; d'après le lemme d'Arden, l'équation est équivalent à $X_q = K^*L$. En remplaçant X_q par K^*L dans les autres équations, on obtient encore un système d'équations linéaires, équivalent au système d'origine.

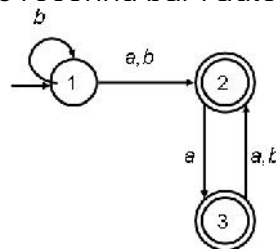
© Compilation /
M.MOURCHID, SMI

65

Détermination du langage accepté par un automate fini

Exemple

Soit L le langage reconnu par l'automate suivant :



Le système d'équations linéaires associé à cet automate est :

$$X_1 \equiv bX_1 + (a + b)X_2 \quad (1)$$

$$X_2 \equiv aX_3 \quad (2)$$

$$X_3 \equiv (a + b)X_2 \quad (3)$$

© Compilation /
M.MOURCHID, SMI

66

Détermination du langage accepté par un automate fini

On peut éliminer X_3 dans l'équation (2), en remplaçant X_3 par $(a+b)X_2+v$: on obtient

$$X_2 = a((a+b)X_2+v) + v = a(a+b)X_2 + (a+v)$$

Ce qui est équivalent (d'après le lemme d'Arden) à :

$$X_2 = (a(a+b))^*(a+v)$$

On peut maintenant éliminer X_2 dans l'équation (3) :

$$X_3 = (a+b)(a(a+b))^*(a+v) + v$$

Ainsi que dans l'équation (1) :

$$X_1 = bX_1 + (a+b)(a(a+b))^*(a+v)$$

D'où d'après le lemme d'Arden $X_1 = b^*(a+b)(a(a+b))^*(a+v)$

Le langage reconnu par l'automate est alors :

$$X_1 = b^*(a+b)(a(a+b))^*(a+v)$$

© Compilation /
M.MOURCHID, SMI

67

Automates finis déterministes

Un AFD est un cas particulier d'AFN :

- Il n'y a aucun arc étiqueté par ϵ .
- $\forall (s, a) \in S \times \Sigma$, il y a un seul arc sortant de s étiqueté par a .

Implémentation d'un AFD.

```
s=e0;
c=car_suiv();
tant que (c!=eof) {
    s=trans(s,c);
}
si (s∈F)
    retourner "oui"
sinon
    retourner "non"
```

© Compilation /
M.MOURCHID, SMI

68

Lemme d'Arden (Résolution d'équation linéaire)

Lemme d'Arden

Soient X , K et L des langages sur un alphabet A . On suppose que $v \in K$.

Alors

$X = KX \hat{=} L$ si et seulement si $X = K^*L$