

ANALYSE LEXICALE

THÉORIE DES LANGAGES ET DE COMPILATION

R. EZZAHIR¹

¹Département d'Informatique
Ecole Nationale des Sciences Appliquées
ENSA-Agadir
Université Ibn Zohr

4^{ème} année Génie Informatique, 2013/2014

SOMMAIRE

- 1 INTRODUCTION
- 2 NOTION DE LANGAGE
- 3 AUTOMATES FINIS DÉTERMINISTES
- 4 AUTOMATES FINIS NON-DÉTERMINISTES
 - Détermination d'AFN
 - Des expressions régulières aux AFNs
- 5 AMBIGUÏTÉ LEXICALE
- 6 MINIMISATION D'UN AFD

POURQUOI FAIRE UNE ANALYSE LEXICALE?

- Simplifier considérablement l'analyse.
 - Éliminer les espaces.
 - Éliminer les commentaires.
 - Convertir les données au début.
- Convertir la description physique d'un programme en une séquence de jetons.
 - Chaque jeton est associé à un lexème.
 - Chaque jeton peut avoir des attributs optionnels.
 - Le flux de jetons sera utilisé dans l'analyseur pour récupérer la structure du programme.

ANALYSE LEXICALE

Convertir la description physique d'un programme en une séquence de jetons.

```
While(ip<z)  
++ip;
```

ANALYSE LEXICALE

Convertir la description physique d'un programme en une séquence de jetons.

```
While(ip<z)  
++ip;
```

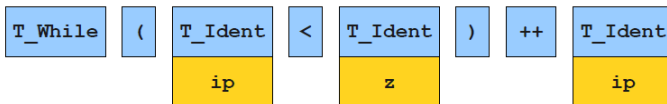
w	h	i	l	e		(i		<		z	\	n	\	t	+	+	i	p	;
---	---	---	---	---	--	---	---	--	---	--	---	---	---	---	---	---	---	---	---	---

ANALYSE LEXICALE

Convertir la description physique d'un programme en une séquence de jetons.

```
While(ip<z)
++ip;
```

w	h	i	l	e		(i		<		z	\	n	\	t	+	+	i	p	;
---	---	---	---	---	--	---	---	--	---	--	---	---	---	---	---	---	---	---	---	---

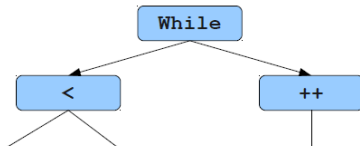
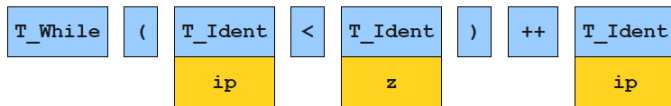


ANALYSE LEXICALE

Convertir la description physique d'un programme en une séquence de jetons.

```
While(ip<z)
++ip;
```

w h i l e (i < z \ n \ t + + i p ;



DÉFIS DE L'ANALYSE LEXICALE

DÉFIS DE L'ANALYSE LEXICALE

- Comment faire pour découper le programme en lexèmes ?
- Comment étiqueter chaque lexème correctement ?

FORTRAN : LES ESPACES NE SONT PAS PERTINENTS

DO 5 | = 1,25

DO5| = 1,25

Il peut être difficile de dire où découper l'entrée.

LES MOTS CLÉS PEUVENT ÊTRE UTILISÉS COMME IDENTIFICATEURS

IF THEN THEN THEN = ELSE ; ELSE ELSE = IF

Il peut être difficile de déterminer la façon d'étiqueter des lexèmes.

DÉFINITION D'UNE ANALYSE LEXICALE

- Définir un ensemble de jetons.
- Définir l'ensemble des lexèmes associé à chaque jeton.
- Définir un algorithme pour résoudre les conflits qui se posent entre lexèmes.

PROBLÉMATIQUE

- But de la théorie des langages
- Un modèle d'aide
- Problématique

Le français est un langage, Java également. Le but de la théorie des langages est de donner un modèle de ce qu'est un langage.

PROBLÉMATIQUE

- But de la théorie des langages
- **Un modèle d'aide**
- Problématique

- pour pouvoir décrire un langage ;
- pour fabriquer une machine capable de reconnaître les textes qui appartiennent à un langage donné.

PROBLÉMATIQUE

- But de la théorie des langages
- Un modèle d'aide
- **Problématique**

il faut donner **une description finie** d'un objet en général **infini** : il y a en effet une infinité de textes français, une infinité de programmes Java

ALPHABETS, MOTS, ET LANGAGES

ALPHABET

Un alphabet Σ (ou vocabulaire, ou lexique) est un ensemble fini de symboles. On note Σ^* l'ensemble des mots sur Σ .

MOT

Un mot (ou phrase) sur un vocabulaire Σ est une séquence finie d'éléments de Σ .

ALPHABETS

1. $\Sigma = \{a, b\} \Rightarrow \Sigma^* = \{a, b, ab, \dots, aaaa, aabb, \dots\}$
2. $\{\text{if, then, else, begin, end, :=, ,, (,), A, B, C, 1, 2}\}$

MOTS

1. aaaa, aabb
2. if A then B := 1 else C := 2
if if if A begin

ALPHABETS, MOTS, ET LANGAGES

ALPHABET

Un alphabet Σ (ou vocabulaire, ou lexique) est un ensemble fini de symboles. On note Σ^* l'ensemble des mots sur Σ .

MOT

Un mot (ou phrase) sur un vocabulaire Σ est une séquence finie d'éléments de Σ .

ALPHABETS

1. $\Sigma = \{a, b\} \Rightarrow \Sigma^* = \{a, b, ab, \dots, aaaa, aabb, \dots\}$
2. $\{\text{if, then, else, begin, end, :=, ,, (,), A, B, C, 1, 2}\}$

MOTS

1. aaaa, aabb
2. if A then B := 1 else C := 2
 if if if A begin

ALPHABETS, MOTS, ET LANGAGES

ALPHABET

Un alphabet Σ (ou vocabulaire, ou lexique) est un ensemble fini de symboles. On note Σ^* l'ensemble des mots sur Σ .

MOT

Un mot (ou phrase) sur un vocabulaire Σ est une séquence finie d'éléments de Σ .

ALPHABETS

1. $\Sigma = \{a, b\} \Rightarrow \Sigma^* = \{a, b, ab, \dots, aaaa, aabb, \dots\}$
2. $\{\text{if, then, else, begin, end, :=, ,, (,), A, B, C, 1, 2}\}$

MOTS

1. aaaa, aabb
2. if A then B := 1 else C := 2
 if if if A begin

ALPHABETS, MOTS, ET LANGAGES

ALPHABET

Un alphabet Σ (ou vocabulaire, ou lexique) est un ensemble fini de symboles. On note Σ^* l'ensemble des mots sur Σ .

MOT

Un mot (ou phrase) sur un vocabulaire Σ est une séquence finie d'éléments de Σ .

ALPHABETS

1. $\Sigma = \{a, b\} \Rightarrow \Sigma^* = \{a, b, ab, \dots, aaaa, aabb, \dots\}$
2. $\{\text{if, then, else, begin, end, :=, ,, (,), A, B, C, 1, 2}\}$

MOTS

1. aaaa, aabb
2. if A then B := 1 else C := 2
if if if A begin

ALPHABETS, MOTS, ET LANGAGES (2)

CONCATÉNATION

• : $\Sigma^* \times \Sigma^* \mapsto \Sigma^*$

soit $x = x_1x_2...x_k \in \Sigma^*$ et

$y = y_1y_2...y_l \in \Sigma^*$ alors

$x.y = x_1x_2...x_ky_1y_2...y_l$ est la concaténation de x et y .

ε : mot vide, neutre pour la concaténation ($\varepsilon.w = w.\varepsilon = w$).

$|u|$ longueur du mot u .

$|\cdot| : \Sigma^* \rightarrow \mathbb{N}$ est le morphisme défini par $|a| = 1$ pour $a \in \Sigma$, $|\varepsilon| = 0$.

$|u|_a$: nombre de a dans le mot u .

$W = w_1w_2...w_k / w_i \in \Sigma \rightarrow |W| = k$

DÉFINITION : LANGAGE

Soit $L \subseteq \Sigma^*$ on dira que L est un langage sur l'alphabet Σ .

On note $\{\}$ ou \emptyset le langage vide.

LANGAGES SUR $\Sigma = \{a, b\}$

$L_1 = \{a, abba, abbb, b\}$

$L_2 = \{w / |w| = 2k \ k \in \mathbb{N}\}$

$= \{\varepsilon, aa, ab, ba, bb, aaaa, aaab, aaba, aabb, \dots\}$

ALPHABETS, MOTS, ET LANGAGES (2)

CONCATÉNATION

• : $\Sigma^* \times \Sigma^* \mapsto \Sigma^*$

soit $x = x_1x_2...x_k \in \Sigma^*$ et

$y = y_1y_2...y_l \in \Sigma^*$ alors

$x.y = x_1x_2...x_ky_1y_2...y_l$ est la concaténation de x et y .

ε : **mot vide**, neutre pour la concaténation ($\varepsilon.w = w.\varepsilon = w$).

$|u|$ longueur du mot u .

$|\cdot| : \Sigma^* \rightarrow \mathbb{N}$ est le morphisme défini par $|a| = 1$ pour $a \in \Sigma$, $|\varepsilon| = 0$.

$|u|_a$: nombre de a dans le mot u .

$W = w_1w_2...w_k / w_i \in \Sigma \rightarrow |W| = k$

DÉFINITION : LANGAGE

Soit $L \subseteq \Sigma^*$ on dira que L est un **langage** sur l'alphabet Σ .

On note $\{\}$ ou \emptyset le langage vide.

LANGAGES SUR $\Sigma = \{a, b\}$

$L_1 = \{a, abba, abbb, b\}$

$L_2 = \{w / |w| = 2k \ k \in \mathbb{N}\}$

$= \{\varepsilon, aa, ab, ba, bb, aaaa, aaab, aaba, aabb, \dots\}$

ALPHABETS, MOTS, ET LANGAGES (2)

CONCATÉNATION

• : $\Sigma^* \times \Sigma^* \mapsto \Sigma^*$

soit $x = x_1x_2...x_k \in \Sigma^*$ et

$y = y_1y_2...y_l \in \Sigma^*$ alors

$x.y = x_1x_2...x_ky_1y_2...y_l$ est la concaténation de x et y .

ε : **mot vide**, neutre pour la concaténation ($\varepsilon.w = w.\varepsilon = w$).

$|u|$ longueur du mot u .

$|\cdot| : \Sigma^* \rightarrow \mathbb{N}$ est le morphisme défini par $|a| = 1$ pour $a \in \Sigma$, $|\varepsilon| = 0$.

$|u|_a$: nombre de a dans le mot u .

$W = w_1w_2...w_k / w_i \in \Sigma \rightarrow |W| = k$

DÉFINITION : LANGAGE

Soit $L \subseteq \Sigma^*$ on dira que L est un **langage** sur l'alphabet Σ .

On note $\{\}$ ou \emptyset le langage vide.

LANGAGES SUR $\Sigma = \{a, b\}$

$L_1 = \{a, abba, abbb, b\}$

$L_2 = \{w / |w| = 2k \text{ } k \in \mathbb{N}\}$

$= \{\varepsilon, aa, ab, ba, bb, aaaa, aaab, aaba, aabb, \dots\}$

ALPHABETS, MOTS, ET LANGAGES (2)

CONCATÉNATION

• : $\Sigma^* \times \Sigma^* \mapsto \Sigma^*$

soit $x = x_1x_2...x_k \in \Sigma^*$ et

$y = y_1y_2...y_l \in \Sigma^*$ alors

$x.y = x_1x_2...x_ky_1y_2...y_l$ est la concaténation de x et y .

ε : **mot vide**, neutre pour la concaténation ($\varepsilon.w = w.\varepsilon = w$).

$|u|$ longueur du mot u .

$|\cdot| : \Sigma^* \rightarrow \mathbb{N}$ est le morphisme défini par $|a| = 1$ pour $a \in \Sigma$, $|\varepsilon| = 0$.

$|u|_a$: nombre de a dans le mot u .

$W = w_1w_2...w_k / w_i \in \Sigma \rightarrow |W| = k$

DÉFINITION : LANGAGE

Soit $L \subseteq \Sigma^*$ on dira que L est un **langage** sur l'alphabet Σ .

On note $\{\}$ ou \emptyset le langage vide.

LANGAGES SUR $\Sigma = \{a, b\}$

$L_1 = \{a, abba, abbb, b\}$

$L_2 = \{w / |w| = 2k \text{ } k \in \mathbb{N}\}$

$= \{\varepsilon, aa, ab, ba, bb, aaaa, aaab, aaba, aabb, \dots\}$

OPÉRATIONS SUR LES ENSEMBLES

- Union

- Intersection

- Complément

- Produit

Cartésien

- Ensemble des parties

- Fermeture de Kleene

- $R1 \mid R2$ est une expression régulière correspondant à la **disjonction** de langues.

- $R1^*$ est une expression régulière correspondant à la fermeture de **Kleene** du langage.

- (R) est une expression régulière correspondant à R .

$$A \cup B = \{w \mid w \in A \text{ ou } w \in B\}$$

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

$L1 = \{\varepsilon, a, aa\}$ et $L2 = \{a, bb\}$.

$L1 \cup L2 = \{\varepsilon, a, aa, bb\}$

OPÉRATIONS SUR LES ENSEMBLES

- Union
- Intersection
- Complément
- Produit
Cartésien
- Ensemble des
parties
- Fermeture de
Kleene

$$A \cap B = \{w \mid w \in A \text{ et } w \in B\}$$

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

$L1 = \{\varepsilon, a, aa\}$ et $L2 = \{a, bb\}$.

$$L1 \cap L2 = \{a\}$$

- $R1 \mid R2$ est une expression régulière correspondant à la **disjonction** de langues.
- $R1^*$ est une expression régulière correspondant à la fermeture de **Kleene** du langage.
- (R) est une expression régulière correspondant à R .

OPÉRATIONS SUR LES ENSEMBLES

- Union
- Intersection
- Complément
- Produit
Cartésien
- Ensemble des
parties
- Fermeture de
Kleene

$$\overline{A} = \{w \mid w \notin A\}$$

plus précisément, soit U l'ensemble de référence,
alors $\overline{A} = \{w \mid w \in U \text{ et } w \notin A\}$

Ex. Soit $\Sigma = \{a, b\}$ et $L_1 = \{\varepsilon, a, aa\}$
 $\overline{L_1} = \{b, ab, ba, bb, aaa, aab, aba, baa, bba, \dots\}$

- $R_1 \mid R_2$ est une expression régulière correspondant à la **disjonction** de langues.
- R_1^* est une expression régulière correspondant à la fermeture de **Kleene** du langage.
- (R) est une expression régulière correspondant à R .

OPÉRATIONS SUR LES ENSEMBLES

- Union
- Intersection
- Complément
- **Produit Cartésien**
- Ensemble des parties
- Fermeture de Kleene

$$A \times B = \{x \in A \text{ et } y \in B\}$$

$$\text{Note : } A \times \emptyset = \emptyset; \quad |A \times B| = |A| |B|$$

Ex. Soit $\Sigma = \{a, b\}$, $L_1 = \{\varepsilon, a, aa\}$ et $L_2 = \{a, bb\}$.

$$L_1 \times L_2 = \{(\varepsilon, a), (a, a), (a, bb), (aa, a), (aa, bb)\}$$

- $R_1 \mid R_2$ est une expression régulière correspondant à la **disjonction** de langues.
- R_1^* est une expression régulière correspondant à la fermeture de **Kleene** du langage.
- (R) est une expression régulière correspondant à R .

OPÉRATIONS SUR LES ENSEMBLES

- Union
- Intersection
- Complément
- Produit
Cartésien
- Ensemble des
parties
- Fermeture de
Kleene

$$\mathcal{P}(A) = \{B \mid B \subseteq A\}$$

Note : $|\mathcal{P}(A)| = 2^{|A|}$

Ex. $L_1 = \{\varepsilon, a, aa\}$

$$\mathcal{P}(L_1) =$$

$$\{\emptyset, \{\varepsilon\}, \{a\}, \{aa\}, \{\varepsilon, a\}, \{\varepsilon, aa\}, \{a, aa\}, \{\varepsilon, a, aa\}\}$$

- $R_1 \mid R_2$ est une expression régulière correspondant à la **disjonction** de langues.
- R_1^* est une expression régulière correspondant à la fermeture de **Kleene** du langage.
- (R) est une expression régulière correspondant à R .

OPÉRATIONS SUR LES ENSEMBLES

- Union
- Intersection
- Complément
- Produit
Cartésien
- Ensemble des
parties
- Fermeture de
Kleene

Soit L un langage sur l'alphabet Σ alors

$L^* = \{w = x_1.x_2.....x_k | k \geq 0 \text{ et } x_i \in L\}$ est la
fermeture de Kleene de L , dénotée par l'ER $(a)^* |$
 $a \in L$. Ex. $L_2 = \{a, bb\}$.

$L_2^* = \{\varepsilon, a, bb, aa, abb, bba, bbbb, aaa, aabb, \dots\}$

Note : pour tout langage $L \neq \emptyset$ on a $|L^*| = \infty$

- $R1 | R2$ est une expression régulière correspondant à la **disjonction**
de langues.
- $R1^*$ est une expression régulière correspondant à la fermeture de
Kleene du langage.
- (R) est une expression régulière correspondant à R .

EXPRESSIONS RÉGULIÈRES EN ACTION

EXAMPLE :

Nombres paire $(+|-| \epsilon) (0|1|2|3|4|5|6|7|8|9)^* (0|2|4|6|8)$

SIMPLIFICATION DE NOTATION : DÉFINITIONS RÉGULIÈRES

Signe = $+|-|$

OptSigne = $\text{Signe} | \epsilon$

chiffre = $0|1|2|3|4|5|6|7|8|9$

chiffrePaire = $0|2|4|6|8$

nombrePaire = $\text{OptSigne} \text{chiffre}^* \text{chiffrePaire}$

EXPRESSIONS RÉGULIÈRES EN ACTION

EXAMPLE :

Nombres paire $(+|-| \epsilon) (0|1|2|3|4|5|6|7|8|9)^* (0|2|4|6|8)$

SIMPLIFICATION DE NOTATION : DISJONCTION MULTI-MANIÈRES

Signe = $+|-$

OptSigne = $\text{Signe} | \epsilon$

chiffre = $[0123456789]$

chiffrePaire = $[02468]$

nombrePaire = $\text{OptSigne} \text{chiffre}^* \text{chiffrePaire}$

EXPRESSIONS RÉGULIÈRES EN ACTION

EXAMPLE :

Nombres paire $(+|-| \epsilon) (0|1|2|3|4|5|6|7|8|9)^* (0|2|4|6|8)$

SIMPLIFICATION DE NOTATION : **RANGES**

Signe = $+|-$

OptSigne = $\text{Signe} | \epsilon$

chiffre = **[0-9]**

chiffrePaire = **[02468]**

nombrePaire = $\text{OptSigne} \text{chiffre}^* \text{chiffrePaire}$

EXPRESSIONS RÉGULIÈRES EN ACTION

EXAMPLE :

Nombres paire $(+|-| \epsilon) (0|1|2|3|4|5|6|7|8|9)^* (0|2|4|6|8)$

SIMPLIFICATION DE NOTATION : **RANGES**

Signe = $+|-$

OptSigne = $\text{Signe} | \epsilon$

chiffre = **[0-9]**

chiffrePaire = **[02468]**

nombrePaire = $\text{OptSigne} \text{chiffre}^* \text{chiffrePaire}$

EXPRESSIONS RÉGULIÈRES EN ACTION

EXAMPLE :

Nombres paire (+|-| ϵ) (0|1|2|3|4|5|6|7|8|9)* (0|2|4|6|8)

SIMPLIFICATION DE NOTATION : ZERO-OU-UN

Signe = + | -

OptSigne = Signe?

chiffre = [0-9]

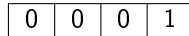
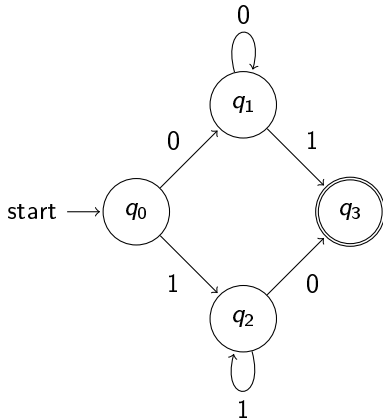
chiffrePaire = [02468]

nombrePaire = OptSigne chiffre* chiffrePaire

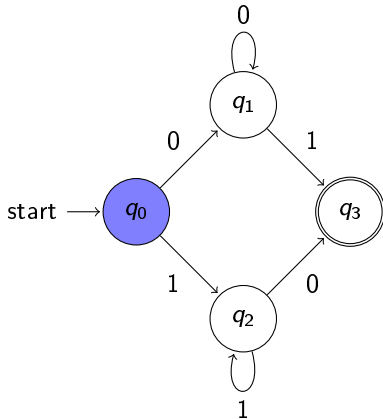
LA MISE EN ŒUVRE DES EXPRESSIONS RÉGULIÈRES

- Les expressions régulières peuvent être mis en œuvre à l'aide **automates finis**.
- Il existe deux types d'automates finis :
 - AFD (**automates finis déterministes**), que nous le verrons dans un peu, et
 - AFN (**automates finis non déterministe**), que nous allons voir plus loin ;
- Les automates sont mieux expliquées par l'exemple ...

EXEMPLE D'UTILISATION D'AF

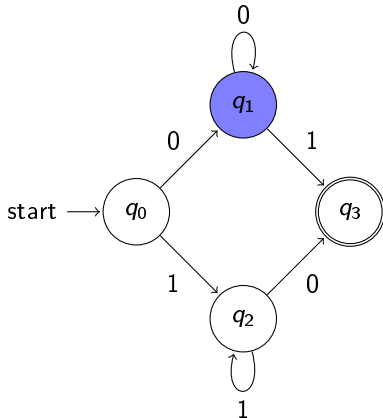


EXEMPLE D'UTILISATION D'AF



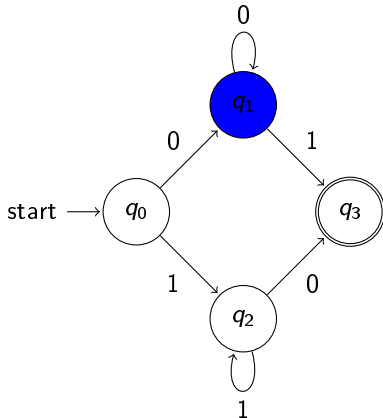
0	0	0	1
↑↑			

EXEMPLE D'UTILISATION D'AF



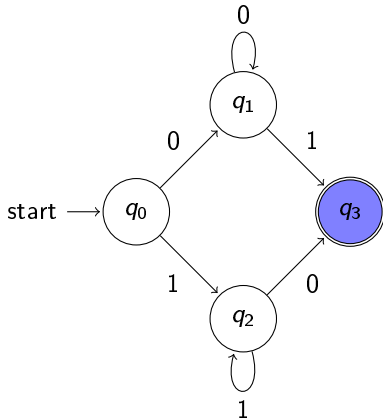
0	0	0	1
	↑		

EXEMPLE D'UTILISATION D'AF



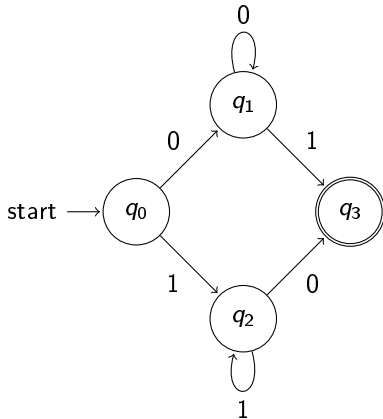
0	0	0	1
		↑↑	

EXEMPLE D'UTILISATION D'AF



0	0	0	1
			↑↑

EXEMPLE D'UTILISATION D'AF



0	0	0	1
---	---	---	---

 OK

AUTOMATES FINIS DÉTERMINISTES

AUTOMATE FINI (AFD)

Un automate fini (AFD) est un quintuplet $(Q, \Sigma, \delta, q_0, F)$

- Q , un ensemble fini d'états ;
- Σ un alphabet (fini) ;
- $\delta : Q \times \Sigma \cup \{\varepsilon\} \rightarrow Q$ une fonction de transition ;
- $q_0 \in Q$ un état initial ;
- $F \subseteq Q$, des états finaux (accepteurs).

EXEMPLE

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

$L(A) = \{w : |w|a \text{ est pair}\}$.

$|L(A)| = \infty$

$A = (Q, \Sigma, \delta, q_0, F)$ où

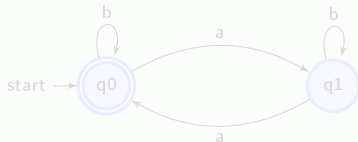
$Q = \{q_0, q_1\}$

$F = \{q_0\}$

et δ est donnée formellement par :

	a	b
$\Rightarrow q_0$	q1	q0
q1	q0	q1

L'automate prend en entrée un mot et l'accepte ou le rejette. On dit aussi qu'il le reconnaît ou ne le reconnaît pas.



q0 : état d'acceptation

AUTOMATES FINIS DÉTERMINISTES

AUTOMATE FINI (AFD)

Un automate fini (AFD) est un quintuplet $(Q, \Sigma, \delta, q_0, F)$

- Q , un ensemble fini d'états ;
- Σ un alphabet (fini) ;
- $\delta : Q \times \Sigma \cup \{\varepsilon\} \rightarrow Q$ une fonction de transition ;
- $q_0 \in Q$ un état initial ;
- $F \subseteq Q$, des états finaux (accepteurs).

EXEMPLE

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

$L(A) = \{w : |w|a \text{ est pair}\}.$

$|L(A)| = \infty$

$A = (Q, \Sigma, \delta, q_0, F)$ où

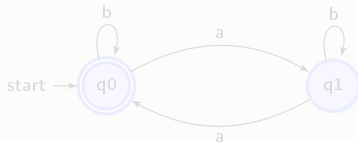
$Q = \{q_0, q_1\}$

$F = \{q_0\}$

et δ est donnée formellement par :

	a	b
$\Rightarrow q_0$	q1	q0
q1	q0	q1

L'automate prend en entrée un mot et l'accepte ou le rejette. On dit aussi qu'il le reconnaît ou ne le reconnaît pas.



q0 : état d'acceptation

AUTOMATES FINIS DÉTERMINISTES

AUTOMATE FINI (AFD)

Un automate fini (AFD) est un quintuplet $(Q, \Sigma, \delta, q_0, F)$

- Q , un ensemble fini d'états ;
- Σ un alphabet (fini) ;
- $\delta : Q \times \Sigma \cup \{\varepsilon\} \rightarrow Q$ une fonction de transition ;
- $q_0 \in Q$ un état initial ;
- $F \subseteq Q$, des états finaux (accepteurs).

EXEMPLE

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

$L(A) = \{w : |w| \text{ est pair}\}$.

$|L(A)| = \infty$

$A = (Q, \Sigma, \delta, q_0, F)$ où

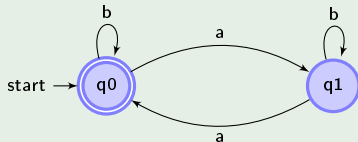
$Q = \{q_0, q_1\}$

$F = \{q_0\}$

et δ est donnée formellement par :

	a	b
$\Rightarrow q_0$	q1	q0
q1	q0	q1

L'automate prend en entrée un mot et l'accepte ou le rejette. On dit aussi qu'il le reconnaît ou ne le reconnaît pas.



q0 : état d'acceptation

LANGAGES RÉGULIERS

Soit $A = (Q, \Sigma, \delta, q_0, F)$ un AFD

DÉFINITION (ÉTAT ACCESSIBLE)

Un état $q \in Q$ est **accessible** s'il existe $w \in \Sigma^*$ tel que $(w, q_0) \xrightarrow{*} (\varepsilon, q)$.

DÉFINITION (MOT ACCEPTÉ)

Soit $w \in \Sigma^*$. On dira que A **accepte** w ssi $\exists q \in F$ tel que $(w, q_0) \xrightarrow{*} (\varepsilon, q)$

DÉFINITION (LANGAGE ACCEPTÉ)

Soit A un AF. $L(A) = \{w : A \text{ accepte } w\}$ est le **langage accepté** par A .

DÉFINITION (LANGAGE RÉGULIER)

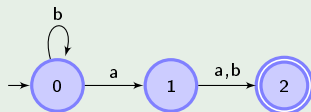
Un langage Y est **régulier (R)** ssi il existe un AF tel que $Y = L(A)$.

AF NON DÉTERMINISME

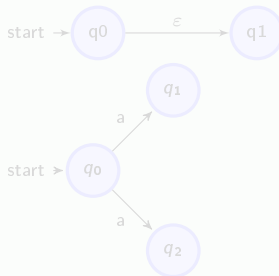
- Un automate A est **déterministe** si pour toute configuration de A, il existe au plus un mouvement possible.

- Un automate est **non déterministe** s'il existe des configurations, pour lesquelles plus d'un mouvement est possible.
- Deux origines :
 - transitions- ϵ
 - + d'1 Mouvement possible

DÉTERMINISTE



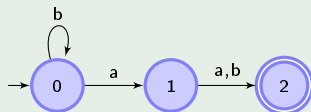
NON DÉTERMINISTE



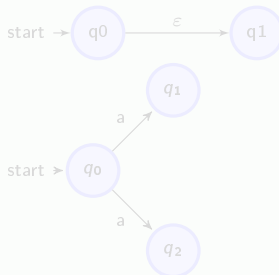
AF NON DÉTERMINISME

- Un automate A est **déterministe** si pour toute configuration de A, il existe au plus un mouvement possible.
- Un automate est **non déterministe** s'il existe des configurations, pour lesquelles plus d'un mouvement est possible.
- Deux origines :
 - transitions- ϵ
 - + d'1 Mouvement possible

DÉTERMINISTE



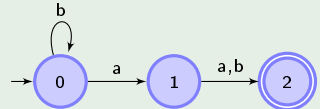
NON DÉTERMINISTE



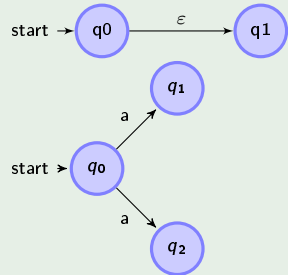
AF NON DÉTERMINISME

- Un automate A est **déterministe** si pour toute configuration de A, il existe au plus un mouvement possible.
- Un automate est **non déterministe** s'il existe des configurations, pour lesquelles plus d'un mouvement est possible.
- Deux origines :
 - transitions- ϵ
 - + d'1 Mouvement possible

DÉTERMINISTE

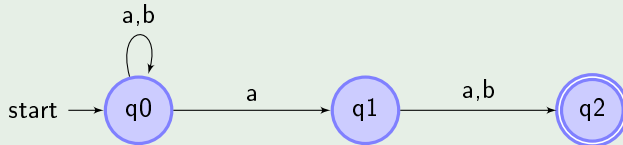


NON DÉTERMINISTE



EXEMPLE D'UN AUTOMATE NON DÉTERMINISTE

AUTOMATE NON DÉTERMINISTE (AFN)



(Attention !) δ est donné formellement par :

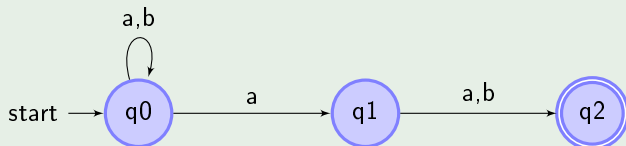
	a	b
→ q0	{q0,q1}	{q0}
q1	{q2}	{q2}
← q2	\emptyset	\emptyset

REMARQUE

Un automate fini déterministe (AFD) est un cas particulier d'un AFN :

EXEMPLE D'UN AUTOMATE NON DÉTERMINISTE

AUTOMATE NON DÉTERMINISTE (AFN)



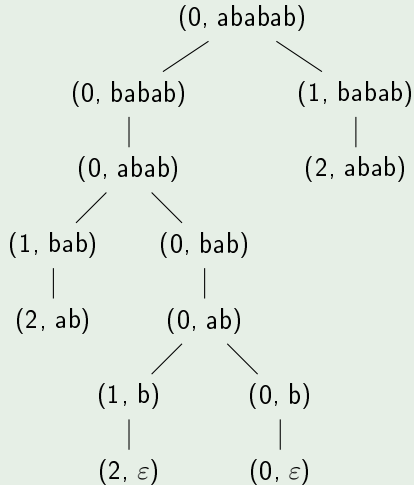
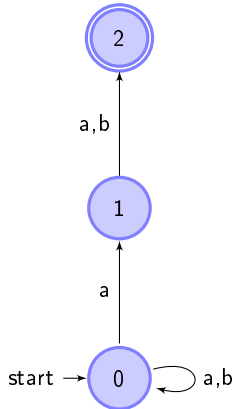
(Attention !) δ est donné formellement par :

	a	b
→ q0	{q0,q1}	{q0}
q1	{q2}	{q2}
← q2	\emptyset	\emptyset

REMARQUE

Un automate fini déterministe (AFD) est un cas particulier d'un AFN :

RECONNAISSANCE D'UN MOT PAR UN AFN



LIGNES DIRECTRICES

- 1 INTRODUCTION
- 2 NOTION DE LANGAGE
- 3 AUTOMATES FINIS DÉTERMINISTES
- 4 AUTOMATES FINIS NON-DÉTERMINISTES**
 - Déterminisation d'AFN
 - Des expressions régulières aux AFNs
- 5 AMBIGUÏTÉ LEXICALE
- 6 MINIMISATION D'UN AFD

DÉTERMINISATION

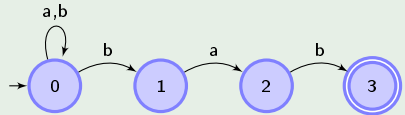
THÉORÈME DE RABIN-SCOTT

Considérons un automate fini **non-déterministe** $A_n = (Q_n, \Sigma, \delta_n, q_0, F_n)$ et construisons un automate fini **déterministe** $A_d = (Q_d, \Sigma, \delta_d, \{q_0\}, F_d)$ qui reconnaît exactement le même langage.

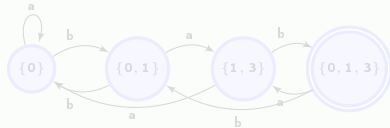
- Les alphabets de A_n et de A_d sont identiques.
- Les états de départ sont respectivement q_0 et le singleton $\{q_0\}$.
- Q_d est constitué de tous les sous-ensembles de Q_n .
- F_d est l'ensemble des sous-ensembles de Q_n qui contiennent au moins un élément de F_n .
- Étant donné un sous ensemble S de Q_n et un symbole $a \in \Sigma$, on définit la fonction de transition $\delta_d(S, a)$ de la manière suivante :
$$\delta_d(S, a) = \bigcup_{q \in S} \delta_n(q, a)$$

DÉTERMINISATION (EXEMPLE)

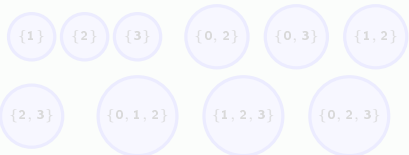
Soit A l'automate fini non-déterministe reconnaissant les mots de l'alphabet a, b qui terminent par bab .



Pour déterminer A en construisant un nouvel état à partir de chaque sous ensemble d'état possible.



Remarquons que les états $\{1\}$, $\{2\}$, $\{3\}$, $\{0, 2\}$, $\{0, 3\}$, $\{1, 2\}$, $\{2, 3\}$, $\{0, 1, 2\}$, $\{1, 2, 3\}$, $\{0, 2, 3\}$ sont **inatteignables** et peuvent être **"retirés"** de l'automate.

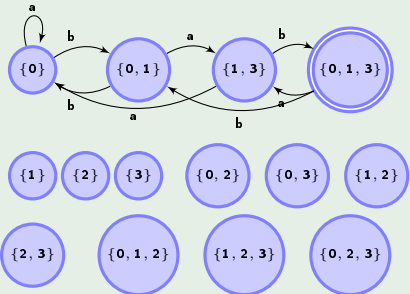
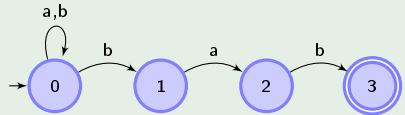


DÉTERMINISATION (EXEMPLE)

Soit A l'automate fini non-déterministe reconnaissant les mots de l'alphabet a, b qui terminent par bab .

Pour déterminer A en construisant un nouvel état à partir de chaque sous ensemble d'état possible.

Remarquons que les états $\{1\}$, $\{2\}$, $\{3\}$, $\{0, 2\}$, $\{0, 3\}$, $\{1, 2\}$, $\{2, 3\}$, $\{0, 1, 2\}$, $\{1, 2, 3\}$, $\{0, 2, 3\}$ sont **inatteignables** et peuvent être **"retirés"** de l'automate.

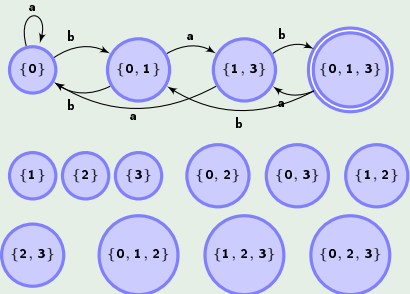
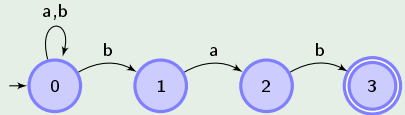


DÉTERMINISATION (EXEMPLE)

Soit A l'automate fini non-déterministe reconnaissant les mots de l'alphabet a, b qui terminent par bab .

Pour déterminer A en construisant un nouvel état à partir de chaque sous ensemble d'état possible.

Remarquons que les états $\{1\}$, $\{2\}$, $\{3\}$, $\{0,2\}$, $\{0,3\}$, $\{1,2\}$, $\{2,3\}$, $\{0,1,2\}$, $\{1,2,3\}$, $\{0,2,3\}$ sont **inatteignables** et peuvent être **"retirés"** de l'automate.



DÉTERMINISATION (EN PRATIQUE)

Lors de la détermination d'un AFN, on ne crée pas immédiatement tous les états de l'AFD.

Les états “utiles” sont créés quand on en a besoin en suivant la méthode de construction ci-dessous :

- Q_d est initialisé à \emptyset et soit E un ensemble d'états initialisé à $E = \{\{q_0\}\}$
- Tant que E est non vide,
 - choisir un élément S de E (S est donc un sous ensemble de Q_n),
 - ajouter S à Q_d ,
 - pour tout symbole $a \in \Sigma$,
 - calculer l'état $S' = \bigcup_{q \in S} \delta_n(q, a)$
 - si S' n'est pas déjà dans Q_d , l'ajouter à E
 - ajouter un arc sur l'automate entre S et S' et la valuer par a.

DÉTERMINISATION (MISE EN ŒUVRE)

- Dans la pratique, lorsque l'on veut construire un automate déterministe AFD à partir d'un automate non déterministe sans transitions- ϵ AFN, on ne commence pas par créer tous les états de AFD, car ils peuvent être nombreux : $|p(Q)| = 2^{|Q|}$!
- On construit plutôt les états de AFD au fur et mesure de leur création en partant de l'état initial.

EXEMPLE

	a	b
$\rightarrow 0$	$\{0, 1\}$	0
1	2	2
$\leftarrow 2$	\emptyset	\emptyset

On part de l'état initial 0 dont on calcule les transitions :

	a	b
$\rightarrow 0$	01	0

un nouvel état 01 a été créé, dont on calcule les transitions.

Ce calcul mène à la création des deux nouveaux états 01 et 012 qui ne donneront pas naissance à de nouveaux états, ce qui nous donne l'automate déterministe suivant :

	a	b
$\rightarrow 0$	01	0
01	012	02
$\leftarrow 02$	01	0
$\leftarrow 012$	012	02

DÉTERMINISATION (MISE EN ŒUVRE)

- Dans la pratique, lorsque l'on veut construire un automate déterministe AFD à partir d'un automate non déterministe sans transitions- ϵ AFN, on ne commence pas par créer tous les états de AFD, car ils peuvent être nombreux : $|p(Q)| = 2^{|Q|}$!
- On construit plutôt les états de AFD au fur et mesure de leur création en partant de l'état initial.

EXEMPLE

	a	b
$\rightarrow 0$	$\{0, 1\}$	0
1	2	2
$\leftarrow 2$	\emptyset	\emptyset

On part de l'état initial 0 dont on calcule les transitions :

	a	b
$\rightarrow 0$	01	0

un nouvel état **01** a été créé, dont on calcule les transitions.

Ce calcul mène à la création des deux nouveaux états 01 et 012 qui ne donneront pas naissance à de nouveaux états, ce qui nous donne l'automate déterministe suivant :

	a	b
$\rightarrow 0$	01	0
01	012	02
$\leftarrow 02$	01	0
$\leftarrow 012$	012	02

DÉTERMINISATION (MISE EN ŒUVRE)

- Dans la pratique, lorsque l'on veut construire un automate déterministe AFD à partir d'un automate non déterministe sans transitions- ϵ AFN, on ne commence pas par créer tous les états de AFD, car ils peuvent être nombreux : $|p(Q)| = 2^{|Q|}$!
- On construit plutôt les états de AFD au fur et mesure de leur création en partant de l'état initial.

EXEMPLE

	a	b
$\rightarrow 0$	$\{0, 1\}$	0
1	2	2
$\leftarrow 2$	\emptyset	\emptyset

On part de l'état initial 0 dont on calcule les transitions :

	a	b
$\rightarrow 0$	01	0

un nouvel état **01** a été créé, dont on calcule les transitions.

Ce calcul mène à la création des deux nouveaux états 01 et 012 qui ne donneront pas naissance à de nouveaux états, ce qui nous donne l'automate déterministe suivant :

	a	b
$\rightarrow 0$	01	0
01	012	02
$\leftarrow 02$	01	0
$\leftarrow 012$	012	02

LIGNES DIRECTRICES

- 1 INTRODUCTION
- 2 NOTION DE LANGAGE
- 3 AUTOMATES FINIS DÉTERMINISTES
- 4 AUTOMATES FINIS NON-DÉTERMINISTES**
 - Détermination d'AFN
 - Des expressions régulières aux AFNs
- 5 AMBIGUÏTÉ LEXICALE
- 6 MINIMISATION D'UN AFD

CONVERSION D'UNE ER EN AFN

- Il existe une procédure (magnifique) qui permet de convertir une expression régulière en AFN.
- Associez chaque expression régulière à un AFN avec les propriétés suivantes :
 - Il y a exactement un état acceptant.
 - Il n'y a pas de transitions de l'état acceptant.
 - Il n'y a pas de transitions dans l'état de départ.



LES CAS DE BASE

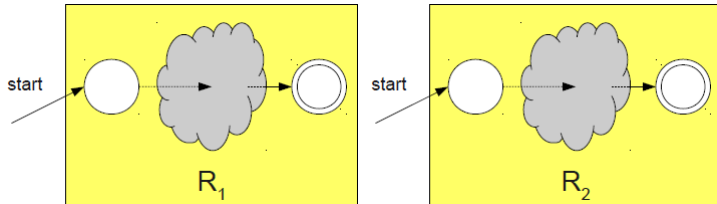
pour ϵ



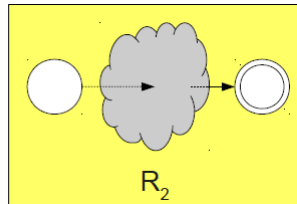
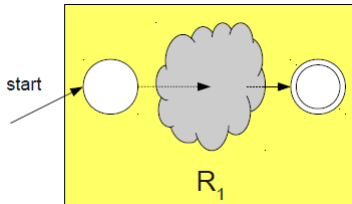
pour a



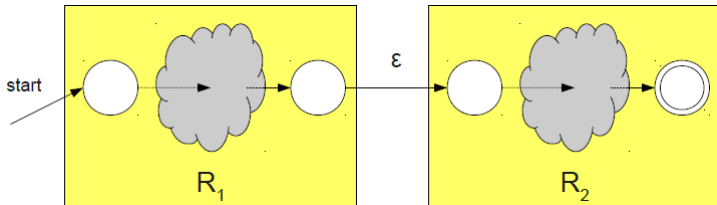
CONSTRUCTION POUR R1R2



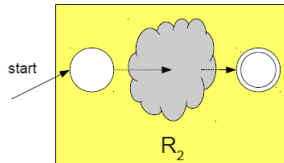
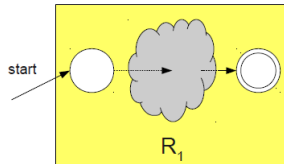
CONSTRUCTION POUR R1R2



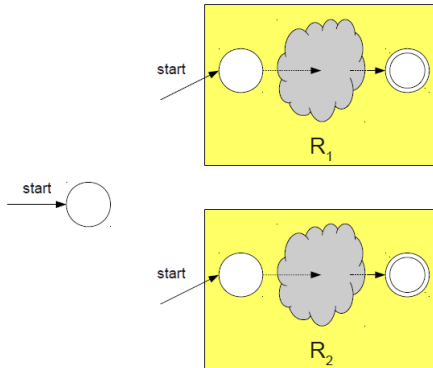
CONSTRUCTION POUR R1R2



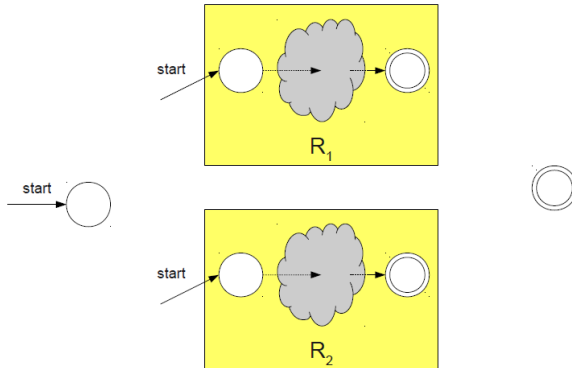
CONSTRUCTION POUR $R_1|R_2$



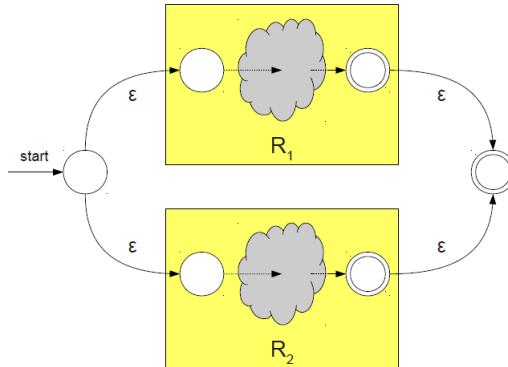
CONSTRUCTION POUR $R_1|R_2$



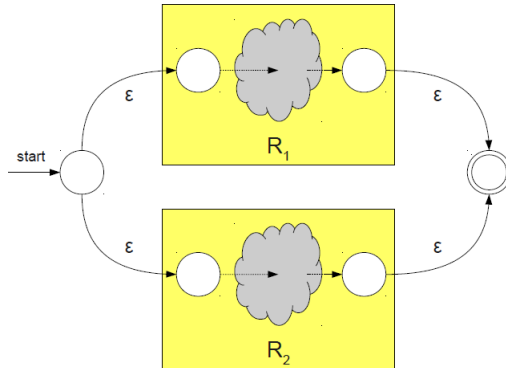
CONSTRUCTION POUR R1|R2



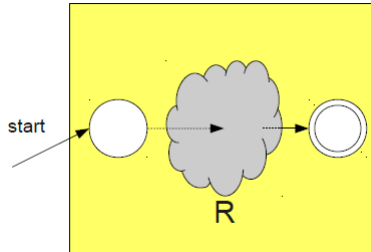
CONSTRUCTION POUR $R_1|R_2$



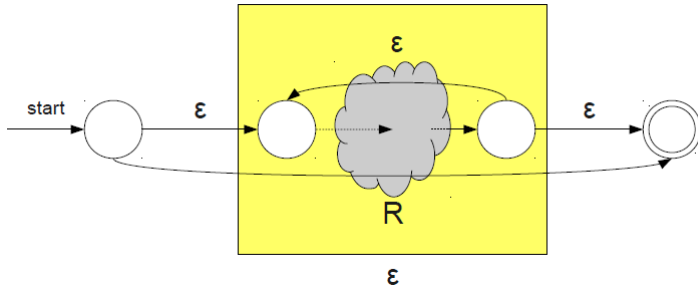
CONSTRUCTION POUR $R_1|R_2$



CONSTRUCTION POUR R^*



CONSTRUCTION POUR R^*



AMBIGUÏTÉ LEXICALE

AMBIGUÏTÉ LEXICALE

T_For for
 T_Identifier [A-Za-z_][A-Za-z0-9_]*

Mot : fort

f	o	r	t		
f	o	r		t	
f	o	r		t	
f	o		r	t	
f	o		r		t

f		o	r	t		
f	o		r		t	
f		o		r	t	
f		o	r		t	
f		o		r		t

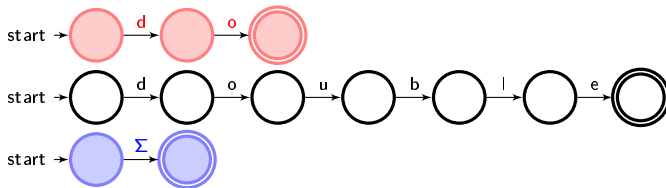
AMBIGUÏTÉ LEXICALE (LONGEST MATCH)

Solution : Maximal Munch

T_For for
T_Identifier [A-Za-z_][A-Za-z0-9_]*
for t

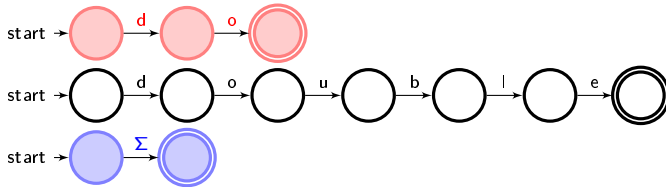
IMPLANTATION DU MAXIMAL MUNCH (1)

T_Do	do
T_Double	double
T_Mystere (inconnu)	[A-Za-z]



IMPLANTATION DU MAXIMAL MUNCH (1)

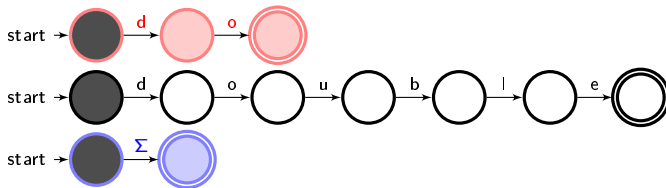
T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



D	O	U	B	D	O	U	B	L	E
---	---	---	---	---	---	---	---	---	---

IMPLANTATION DU MAXIMAL MUNCH (2)

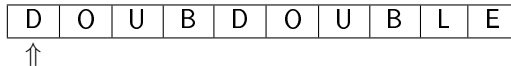
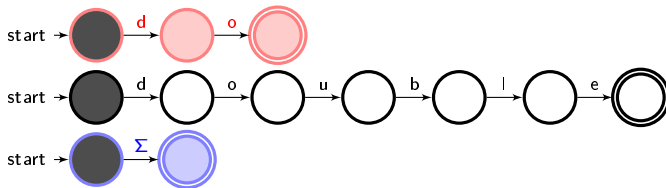
T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



D	O	U	B	D	O	U	B	L	E
---	---	---	---	---	---	---	---	---	---

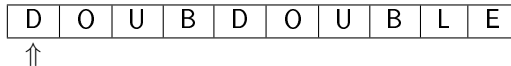
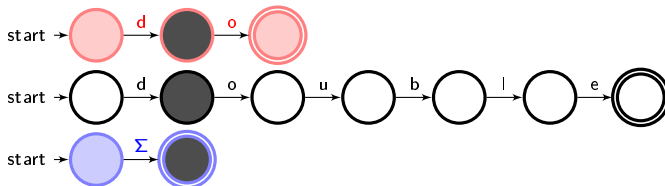
IMPLANTATION DU MAXIMAL MUNCH (2)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



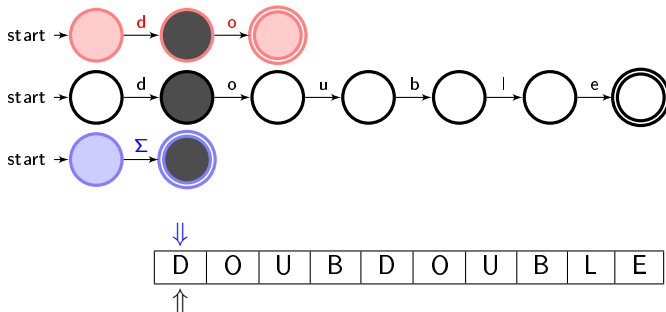
IMPLANTATION DU MAXIMAL MUNCH (3)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



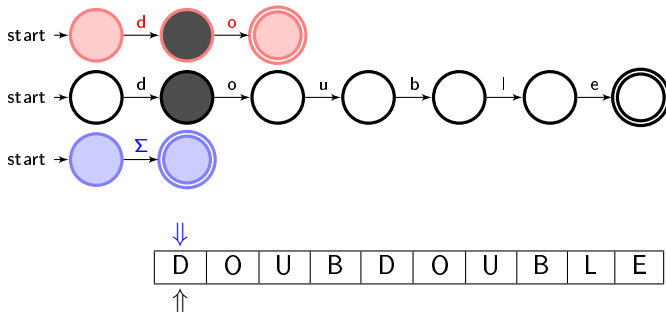
IMPLANTATION DU MAXIMAL MUNCH (3)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



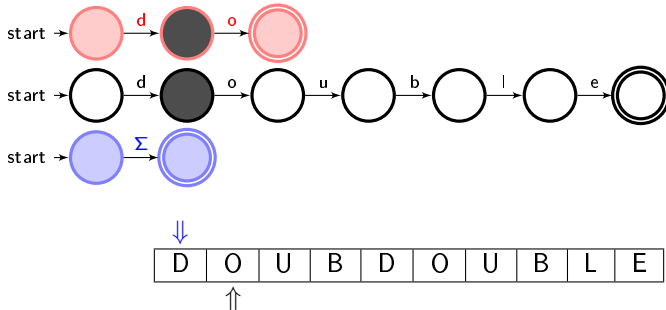
IMPLANTATION DU MAXIMAL MUNCH (4)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



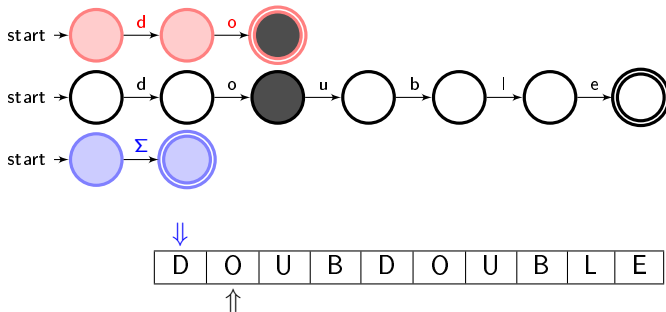
IMPLANTATION DU MAXIMAL MUNCH (4)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



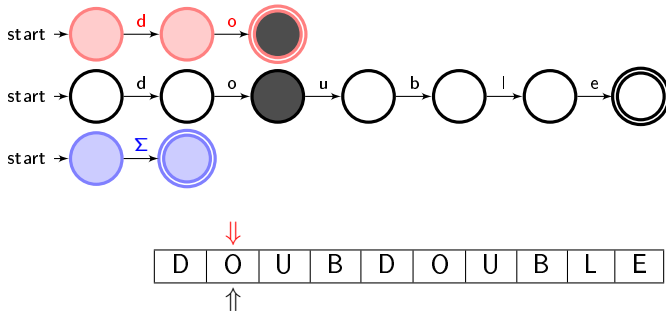
IMPLANTATION DU MAXIMAL MUNCH (5)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



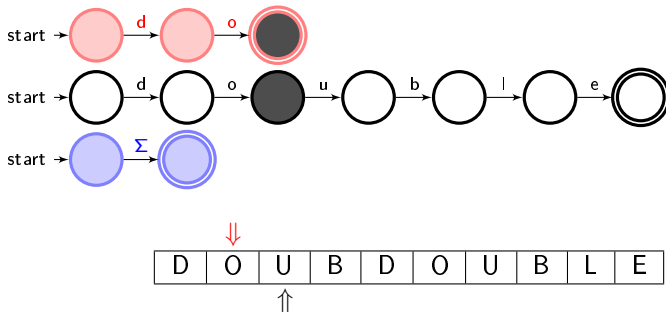
IMPLANTATION DU MAXIMAL MUNCH (5)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



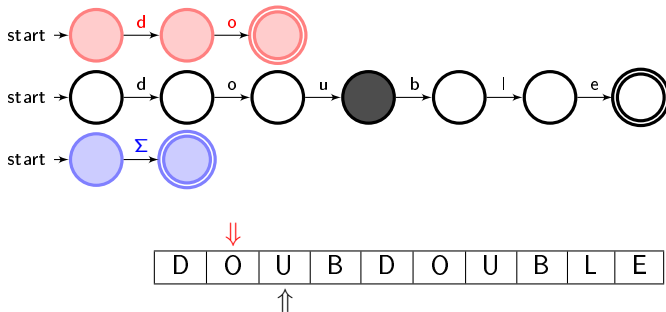
IMPLANTATION DU MAXIMAL MUNCH (5)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



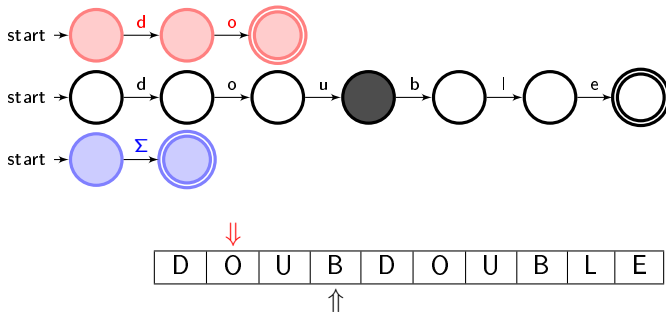
IMPLANTATION DU MAXIMAL MUNCH (6)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



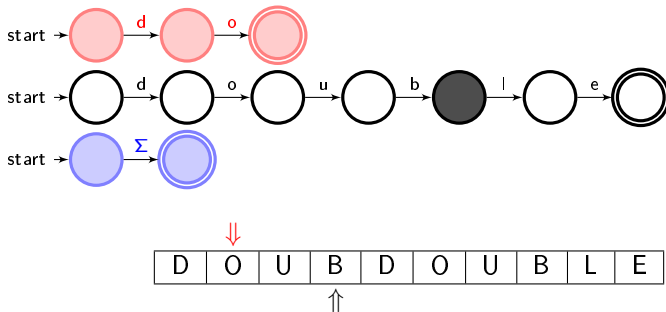
IMPLANTATION DU MAXIMAL MUNCH (6)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



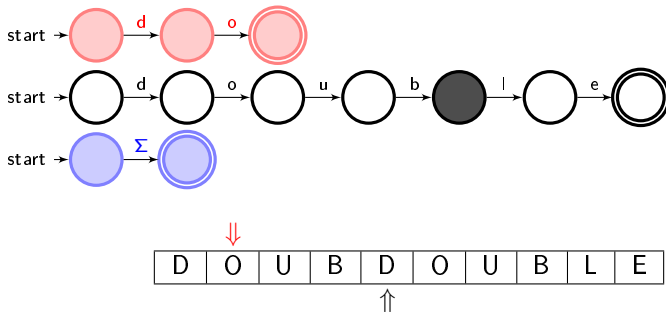
IMPLANTATION DU MAXIMAL MUNCH (7)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



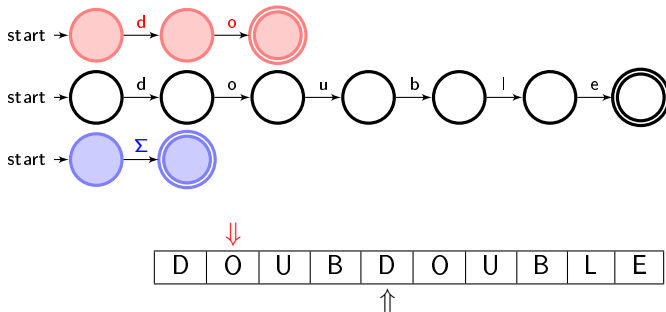
IMPLANTATION DU MAXIMAL MUNCH (7)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



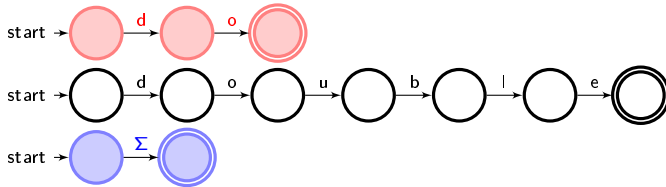
IMPLANTATION DU MAXIMAL MUNCH (8)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



IMPLANTATION DU MAXIMAL MUNCH (9)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]

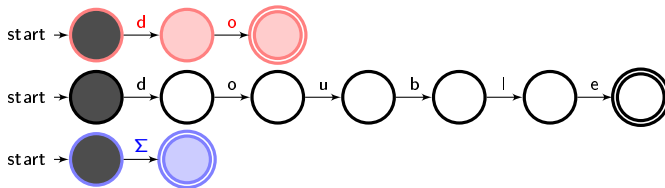


D O

U B D O U B L E
 ↑

IMPLANTATION DU MAXIMAL MUNCH (10)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]

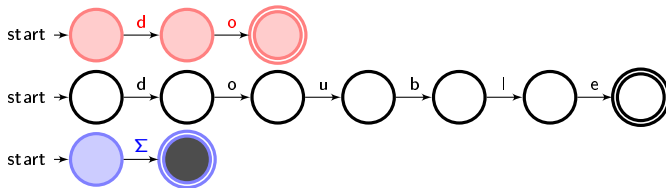


D O

U B D O U B L E
 ↑↑

IMPLANTATION DU MAXIMAL MUNCH (11)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]

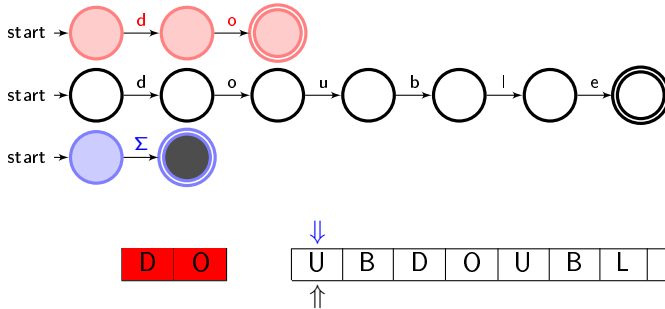


D O

U B D O U B L E
 ↑↑

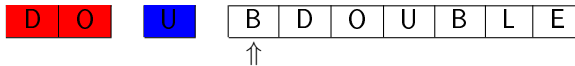
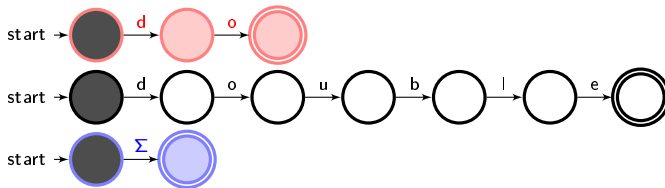
IMPLANTATION DU MAXIMAL MUNCH (11)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



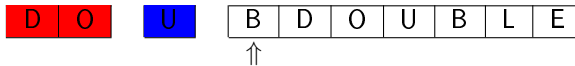
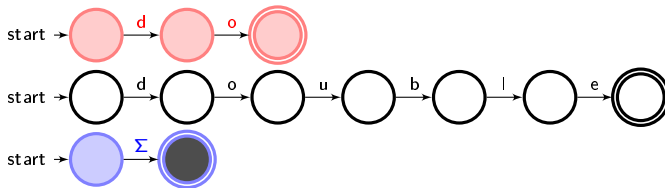
IMPLANTATION DU MAXIMAL MUNCH (12)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



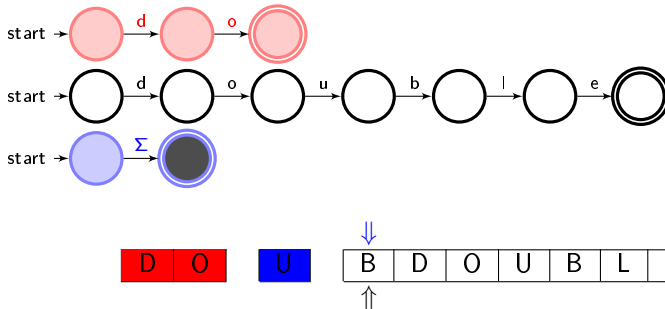
IMPLANTATION DU MAXIMAL MUNCH (12)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



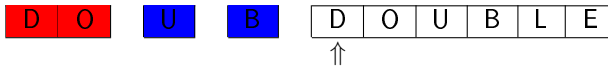
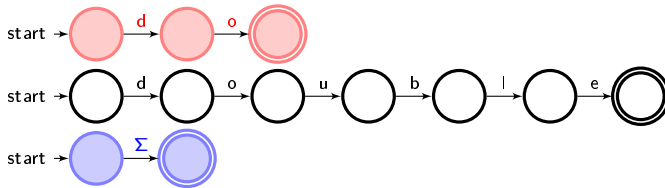
IMPLANTATION DU MAXIMAL MUNCH (12)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



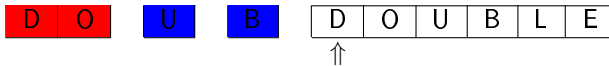
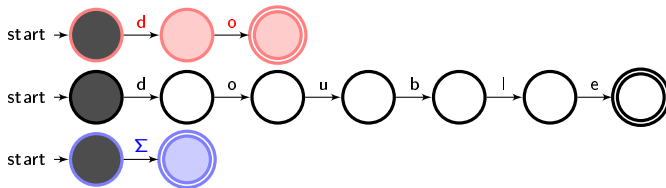
IMPLANTATION DU MAXIMAL MUNCH (13)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



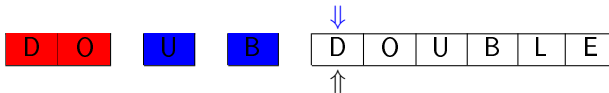
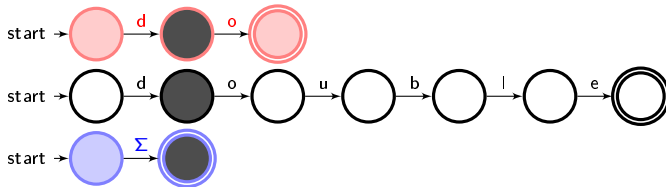
IMPLANTATION DU MAXIMAL MUNCH (14)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



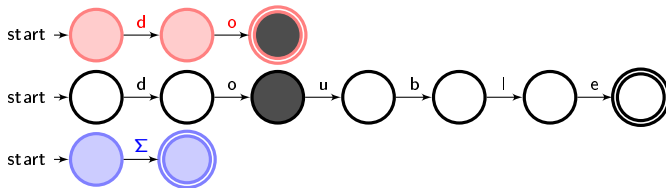
IMPLANTATION DU MAXIMAL MUNCH (15)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



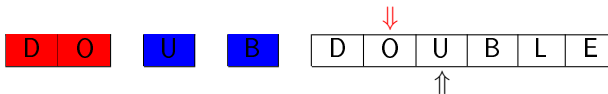
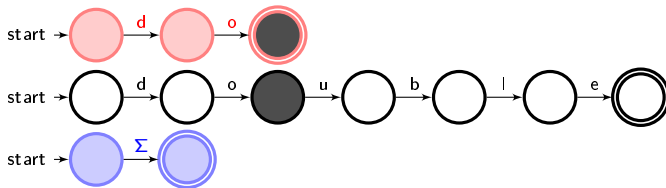
IMPLANTATION DU MAXIMAL MUNCH (16)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



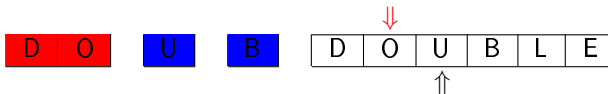
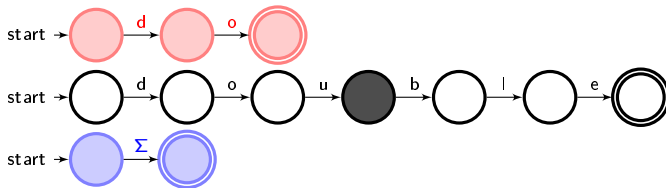
IMPLANTATION DU MAXIMAL MUNCH (16)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



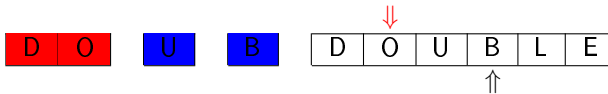
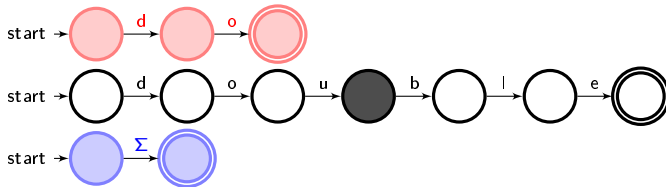
IMPLANTATION DU MAXIMAL MUNCH (17)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



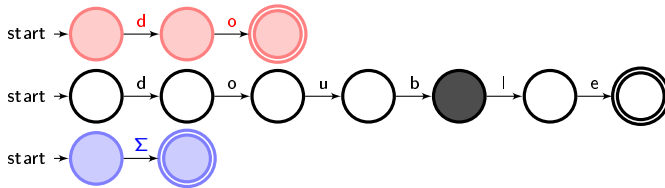
IMPLANTATION DU MAXIMAL MUNCH (17)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



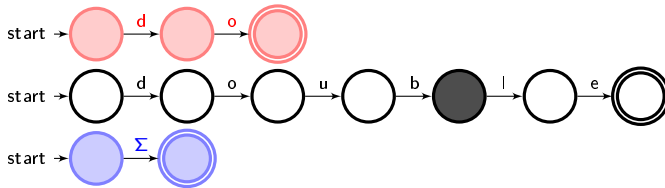
IMPLANTATION DU MAXIMAL MUNCH (18)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



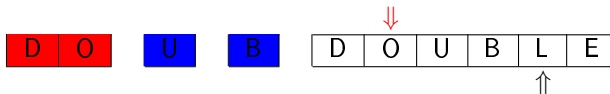
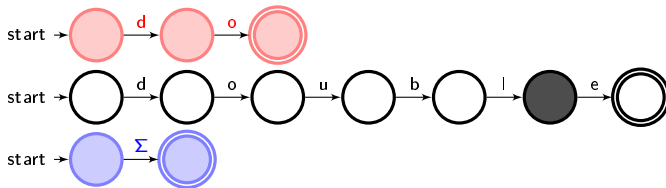
IMPLANTATION DU MAXIMAL MUNCH (18)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



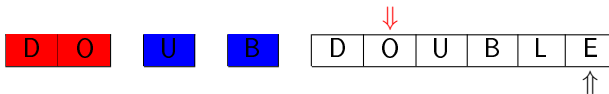
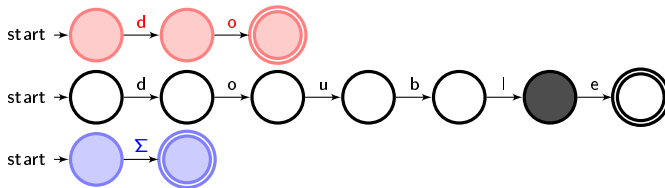
IMPLANTATION DU MAXIMAL MUNCH (19)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



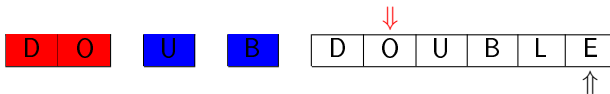
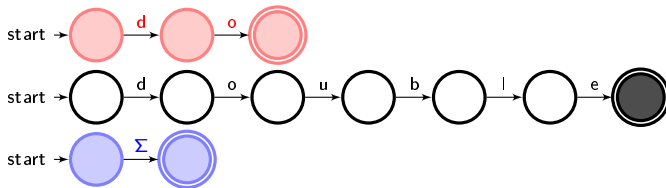
IMPLANTATION DU MAXIMAL MUNCH (19)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



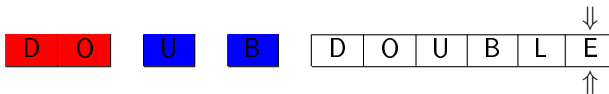
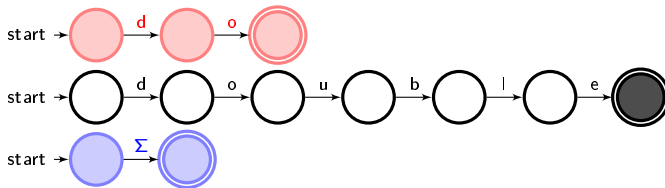
IMPLANTATION DU MAXIMAL MUNCH(20)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



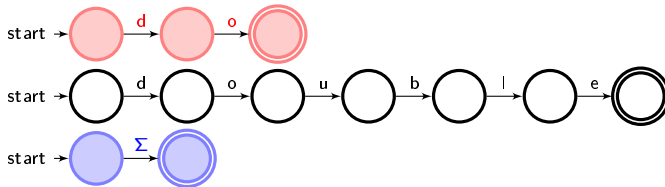
IMPLANTATION DU MAXIMAL MUNCH(20)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



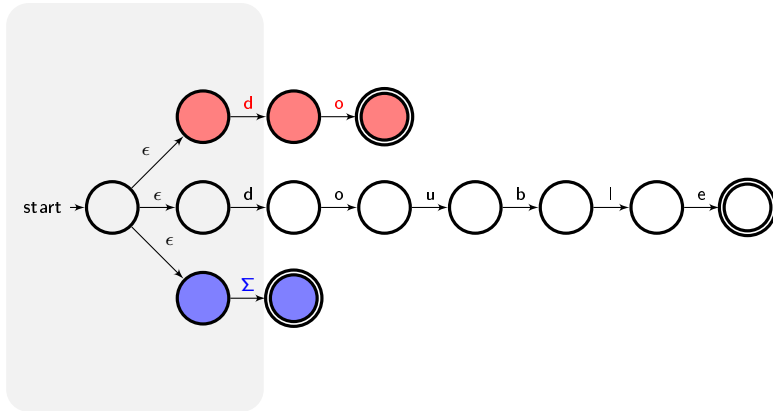
IMPLANTATION DU MAXIMAL MUNCH(20)

T_Do do
 T_Double double
 T_Mystere (inconnu) [A-Za-z]



D	O	U	B	D	O	U	B	L	E
---	---	---	---	---	---	---	---	---	---

PETITE SIMPLIFICATION



UN EXEMPLE DE CLASSE DE LEXÈMES DE JAVA EST :

IDENT	foo, main,
NUMBER	0, 123, 1000
FLOAT	0.5, 1.0e+2

ANALYSEUR LEXICAL EN DUR (INTRODUCTION)

POUR LE PROGRAMME

```
class Foo {  
    void bar {  
        println ("hello world\n");  
    }  
}
```

Sortie du Scanner \Rightarrow entrée du Parser.

L'interface entre l'analyseur lexical et l'analyseur syntaxique doit être une fonction (par exemple `int nextSym()`), qui renvoie à chaque appel l'unité lexicale (symbole) suivante trouvée dans le texte source.

LA SORTIE DU SCANNER

```
CLASS IDENT(Foo) LBRACE VOID IDENT(bar)  
LPAREN RPAREN LBRACE IDENT(println)  
LPAREN STRING ("hello world\n") RPAREN  
SEMICOLON RBRACE RBRACE EOF
```

NEXTSYM()

```
Token sym;  
void nextSym ()  
    "ignore les espaces blancs et assigne le  
    prochain lexème à sym"  
    ...  
}
```

LES DEUX ANALYSEURS PARTAGENT LA DÉFINITIONS DE SYMBOLES.

```
interface Symbols {  
    static final int ERROR = 0, EOF = ERROR + 1, IDENT = EOF + 1,  
        LITERAL = IDENT + 1, LPAREN = LITERAL + 1, RPAREN = LPAREN + 1,
```

ANALYSEUR LEXICAL EN DUR (INTRODUCTION)

POUR LE PROGRAMME

```
class Foo {
    void bar {
        println ("hello world\n");
    }
}
```

Sortie du Scanner \Rightarrow entrée du Parser.

L'interface entre l'analyseur lexical et l'analyseur syntaxique doit être une fonction (par exemple *int nextSym()*), qui renvoie à chaque appel l'unité lexicale (symbole) suivante trouvée dans le texte source.

LA SORTIE DU SCANNER

```
CLASS IDENT(Foo) LBRACE VOID IDENT(bar)
LPAREN RPAREN LBRACE IDENT(println)
LPAREN STRING ("hello world\n") RPAREN
SEMICOLON RBRACE RBRACE EOF
```

NEXTSYM()

```
Token sym;
void nextSym ()
    "ignore les espaces blancs et assigne le
    prochain lexème à sym"
    ...
}
```

LES DEUX ANALYSEURS PARTAGENT LA DÉFINITIONS DE SYMBOLES.

```
interface Symbols {
    static final int ERROR = 0, EOF = ERROR + 1, IDENT = EOF + 1,
        LITERAL = IDENT + 1, LPAREN = LITERAL + 1, RPAREN = LPAREN + 1,
```

ANALYSEUR LEXICAL EN DUR (INTRODUCTION)

POUR LE PROGRAMME

```
class Foo {  
    void bar {  
        println ("hello world\n");  
    }  
}
```

Sortie du Scanner \Rightarrow entrée du Parser.

L'interface entre l'analyseur lexical et l'analyseur syntaxique doit être une fonction (par exemple *int nextSym()*), qui renvoie à chaque appel l'unité lexicale (symbole) suivante trouvée dans le texte source.

LA SORTIE DU SCANNER

```
CLASS IDENT(Foo) LBRACE VOID IDENT(bar)  
LPAREN RPAREN LBRACE IDENT(println)  
LPAREN STRING ("hello world\n") RPAREN  
SEMICOLON RBRACE RBRACE EOF
```

nextSym()

```
Token sym;  
void nextSym ()  
    "ignore les espaces blancs et assigne le  
    prochain lexème à sym"  
    ...  
}
```

LES DEUX ANALYSEURS PARTAGENT LA DÉFINITIONS DE SYMBOLES.

```
interface Symbols {  
    static final int ERROR = 0, EOF = ERROR + 1, IDENT = EOF + 1,  
        LITERAL = IDENT + 1, LPAREN = LITERAL + 1, RPAREN = LPAREN + 1,
```

ANALYSEUR LEXICAL EN DUR (INTRODUCTION)

POUR LE PROGRAMME

```
class Foo {  
    void bar {  
        println ("hello world\n");  
    }  
}
```

Sortie du Scanner \Rightarrow entrée du Parser.

L'interface entre l'analyseur lexical et l'analyseur syntaxique doit être une fonction (par exemple *int nexSym()*), qui renvoie à chaque appel l'unité lexicale (symbole) suivante trouvée dans le texte source.

LA SORTIE DU SCANNER

```
CLASS IDENT(Foo) LBRACE VOID IDENT(bar)  
LPAREN RPAREN LBRACE IDENT(println)  
LPAREN STRING ("hello world\n") RPAREN  
SEMICOLON RBRACE RBRACE EOF
```

NEXTSYM()

```
Token sym;  
void nextSym ()  
    "ignore les espaces blancs et assigne le  
    prochain lexème à sym"  
    ...  
}
```

LES DEUX ANALYSEURS PARTAGENT LA DÉFINITIONS DE SYMBOLES.

```
interface Symbols {  
    static final int ERROR = 0, EOF = ERROR + 1, IDENT = EOF + 1,  
        LITERAL = IDENT + 1, LPAREN = LITERAL + 1, RPAREN = LPAREN + 1,
```

ANALYSEUR LEXICAL EN DUR (EXEMPLE)

Soit la syntaxe des lexèmes en EBNF :

- $\text{symbol} = \{\text{blank}\} (\text{identifier} \mid \text{literal} \mid "(" \mid ")" \mid "[" \mid "]" \mid "{" \mid "}" \mid "|" \mid "=" \mid "." \mid " ")$.
- $\text{Identifier} = \text{letter} \{ \text{letter} \mid \text{digit} \}$.
- $\text{literal} = "\" \{ \text{stringchar} \} "\"$.
- $\text{stringchar} = \text{escapechar} \mid \text{plainchar}$.
- $\text{escapechar} = "\" \backslash \" \text{char}$.
- $\text{plainchar} = \text{charNoQuote}$.

NB : Forme de Backus Naur étendue (EBNF)

ANALYSEUR LEXICAL EN DUR (SUITE)

Interface API de notre Analyseur lexical est de la forme suivante :

```
1  class Scanner implements Symbols {
2      /** Constructor */
3      Scanner (InputStream in)
4      /** The symbol read = tokenclass last */
5      int sym;
6      /** The symbol's character representation */
7      String chars;
8      /** Read next token into sym and chars */
9      void nextSym()
10     /** Close input stream */
11     void close()
12 }
```

ANALYSEUR LEXICAL EN DUR (SUITE)

```
1  import java.io.*;
2  class Scanner implements Symbols {
3  public int sym;
4  public String chars;
5  /** the character stream being tokenized */
6  private InputStream in;
7  /** the next unconsumed character */
8  private char ch;
9  /** a buffer for assembling strings */
10 private StringBuffer buf = new StringBuffer();
11 /** the end of file character */
12 private final char eofCh = (char) -1;
13 public Scanner(InputStream in) { this.in = in; }
14 public static void error(String msg) {
15     System.out.println("**** error: " + msg);
16     System.exit(-1);
17 }
18 /** print current character and read next character */
```

ANALYSEUR LEXICAL EN DUR (SUITE)

```
1  ...
2  /** print current character and read next character */
3  private void nextCh() {
4      System.out.print(ch);
5      try {
6          ch = (char) in.read();
7      } catch (IOException ex) {
8          error("read failure: " + ex.toString());
9      }
10 }
11 /** read next symbol */
12 ...
13 }
```


ANALYSEUR LEXICAL EN DUR (SUITE)

```
1  /** read next symbol */
2  public void nextSym() {
3      while (ch <= ' ') nextCh();
4      switch (ch) {
5          case 'a':    case 'b': ... case 'z':
6          case 'A':    case 'B': ... case 'Z':
7              buf.setLength(0);  buf.append(ch);
8              nextCh();
9              while ('a' <= ch && ch <= 'z' ||
10                 'A' <= ch && ch <= 'Z' ||
11                 '0' <= ch && ch <= '9'){
12                  buf.append(ch); nextCh();
13              }
14              sym = IDENT;
15              chars = buf.toString(); break;
16
17
18 }
```

ANALYSEUR LEXICAL EN DUR (SUITE)

```
1  case '\': nextCh(); buf.setLength(0);
2      while ( ' ' <= ch && ch != eofCh && ch != '\') {
3          if (ch == '\\') nextCh();
4          buf.append(ch); nextCh();
5      }
6      if (ch == '\n') nextCh();
7      else error("unclosed string literal");
8      sym = LITERAL; chars = buf.toString(); break;
9  case '(': sym = LPAREN; nextCh(); break;
10 case '[': sym = LBRACK; nextCh(); break;
11 case '{': sym = LBRACE; nextCh(); break;
12 ...
13 case '|': sym = BAR; nextCh(); break;
14 case '=': sym = EQL; nextCh(); break;
15 case '.': sym = PERIOD; nextCh(); break;
16 case eofCh: sym = EOF; break;
17 default: error("illegal character: " + ch + "(" + (int)
18 }
```

ANALYSEUR LEXICAL EN DUR (SUITE)

```
1
2  /** the string representation of a symbol */
3  public static String representation(int sym) {
4      switch (sym) {
5          case ERROR:      return "<error>";
6          case EOF:        return "<eof>";
7          case IDENT:      return "identifier";
8          case LITERAL:    return "literal";
9          case LPAREN:     return "'('";
10         case RPAREN:     return "')'";
11         case EQL:        return "'='";
12         case PERIOD :    return " '.'";
13         default:         return "<unknown>";
14     }
15 }
16
17 public void close() throws IOException {
18     in.close(); }
```

PROGRAMME DE TESTE

```
1  class ScannerTest implements Symbols {
2      static public void main(String[] args) {
3          try {
4              Scanner s = new Scanner(
5                  new FileInputStream(args[0]));
6              s.nextSym();
7              while (s.sym != EOF) {
8                  System.out.println("[ " +
9                      Scanner.representation(s.sym) + " ]");
10                 s.nextSym();
11             }
12             s.close();
13         } catch (IOException ex) {
14             ex.printStackTrace();
15             System.exit(-1);
16         } } }
```

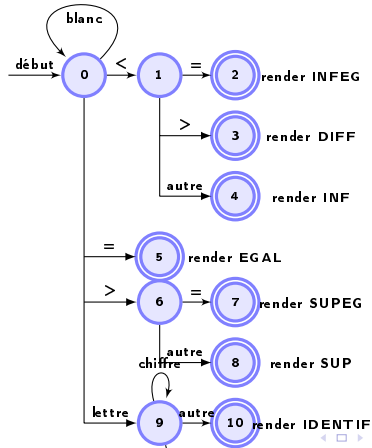
ANALYSEUR LEXICAL À L'AIDE D'AUTOMATES (1)

Il suffira d'implémenter la fonction de transition de l'automate reconnaissant le langage.

	"	'\t'	'\n'	'<'	'='	'>'	lettre	chiffre	autre
0	0	0	0	1	5	6	9	erreur	erreur
1	4*	4*	4*	4*	2	3	4*	4*	4*
6	8*	8*	8*	8*	7	8*	8*	8*	8*
9	10*	10*	10*	10*	10*	10*	9	9	10

ANALYSEUR LEXICAL À L'AIDE D'AUTOMATES

Diagramme de transition pour les opérateurs de comparaison et les identificateurs



ANALYSEUR LEXICAL À L'AIDE D'AUTOMATES

Dans cet exemple, on peut obtenir un analyseur qui peut être plus rapide que celui de la première manière puisque l'essentiel du travail de l'analyseur se réduira à répéter bêtement l'action état = transit[état][nextCh()] jusqu'à tomber sur un état final.

Un tableau identifié par *terminal* et indexé par les états, est défini par :

- $\text{terminal}[e] = 0$ si e n'est pas un état final (vu comme un boolean, $\text{terminal}[e]$ est faux)
- $\text{terminal}[e] = U + 1$ si e est final, sans étoile et associé à l'unité lexicale U (en tant que booléen, $\text{terminal}[e]$ est vrai, car les unités lexicales sont numérotées au moins à partir de zéro),
- $\text{terminal}[e] = -(U + 1)$ si e est final, étoilé et associé à l'unité lexicale U (en tant que booléen, $\text{terminal}[e]$ est encore vrai).
- Nous avons ajouté un état supplémentaire, ayant le numéro `NBR_ETATS`, qui correspond à la mise en erreur de l'analyseur lexical, et une unité lexicale `ERREUR` pour signaler cela.

PROGRAMME SCANNER À L'AIDE D'AUTOMATES

```
1  class ScannerAuomata{
2      static final int NBR_ETATS = 10
3      static final int NBR_CARS = 256
4      int transit[NBR_ETATS][NBR_CARS];
5      int terminal[NBR_ETATS + 1];
6      public ScannerAuomata(){
7          transit = new int[NBR_ETATS][NBR_CARS];
8          terminal = new int[NBR_ETATS+1];
9          int i,j;
10         for (i = 0; i < NBR_ETATS; i++)
11             terminal[i] = 0;
12
13         terminal[2] = INFEG + 1; terminal[3] = DIFF + 1;
14         terminal[4] = - (INF + 1); terminal[5] = EGAL + 1;
15         terminal[7] = SUPEG + 1; terminal[8] = - (SUP + 1);
16         terminal[10] = - (IDENTIF + 1);
17         terminal[NBR_ETATS] = ERREUR + 1;
```


SCANNER À L'AIDE D'AUTOMATES (SUITE)

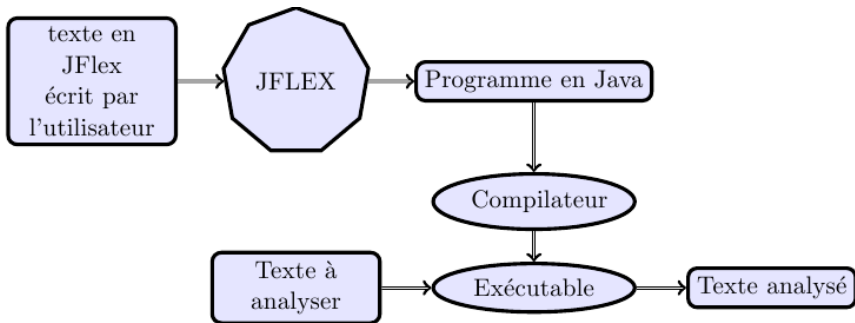
```
1  for (i = 0; i < NBR_ETATS; i++)
2      for (j = 0; j < NBR_CARS; j++)
3          transit[i][j] = NBR_ETATS;
4
5  transit[0][' '] = 0; transit[0]['\t'] = 0;
6  transit[0]['\n'] = 0; transit[0]['<'] = 1;
7  transit[0]['='] = 5; transit[0]['>'] = 6;
8
9  for (j = 'A'; j <= 'Z'; j++)
10     transit[0][j] = 9;
11  for (j = 'a'; j <= 'z'; j++)
12     transit[0][j] = 9;
13  for (j = 0; j < NBR_CARS; j++)
14     transit[1][j] = 4;
15  transit[1]['='] = 2; transit[1]['>'] = 3;
16  for (j = 0; j < NBR_CARS; j++)
17     transit[6][j] = 8;
18  transit[6]['='] = 7;
```

SCANNER À L'AIDE D'AUTOMATES (SUITE)

```
1  for (j = 0; j < NBR_CARS; j++) transit[9][j] =10;
2      for (j = 'A'; j <= 'Z'; j++) transit[9][j] =9;
3      for (j = 'a'; j <= 'z'; j++) transit[9][j] =9;
4      for (j = 0; j < 9; j++) transit[9][j] =9;
5  }
6
7  int nextSym() {
8      char ch;
9      int etat = etatInitial;
10     while(!terminal[etat]){
11         ch = nextCh();
12         etat = transit[etat][ch];
13     }
14     if (terminal[etat] < 0)
15         ch = ' ' ; //blanc
16     return abs(terminal[etat]) - 1;
17 }
18 } //end
```

JFLEX UN GÉNÉRATEUR D'ANALYSEUR LEXICAUX

Principe général de génération d'analyseur lexical par JFlex



JFLEX UN GÉNÉRATEUR D'ANALYSEUR LEXICAUX

MINIMISATION D'UN AFD

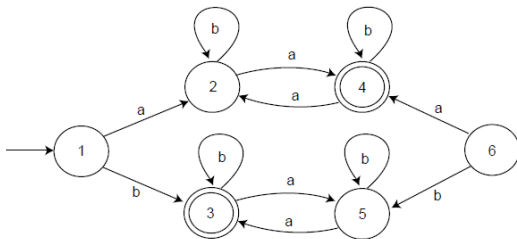
ÉTATS INACCESSIBLES (1)

- L'état 6 est inaccessible à partir de l'état initial.
- Si on enlève cet état, on obtient l'automate M_2 de la figure suivante

MINIMISATION D'UN AFD

ÉTATS INACCESSIBLES (1)

Soit l'AFDM₁ = {1, 2, 3, 4, 5}; {a, b}; 1; {1, 34}; δ}

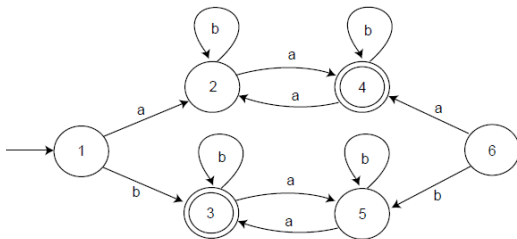


- L'état 6 est inaccessible à partir de l'état initial.
- Si on enlève cet état, on obtient l'automate M_2 de la figure suivante

MINIMISATION D'UN AFD

ÉTATS INACCESSIBLES (1)

Soit l'AFDM₁ = {1, 2, 3, 4, 5}; {a, b}; 1; {1, 34}; δ}

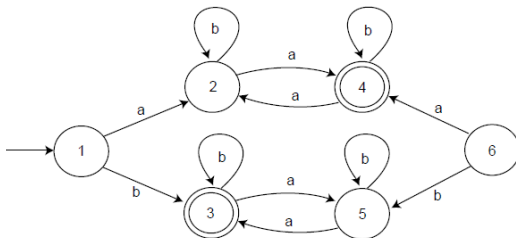


- L'état 6 est inaccessible à partir de l'état initial.
- Si on enlève cet état, on obtient l'automate M_2 de la figure suivante

MINIMISATION D'UN AFD

ÉTATS INACCESSIBLES (1)

Soit l'AFD $M_1 = \{1, 2, 3, 4, 5\}; \{a, b\}; 1; \{1, 34\}; \delta\}$

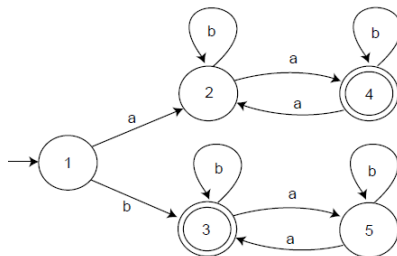


- L'état 6 est inaccessible à partir de l'état initial.
- Si on enlève cet état, on obtient l'automate M_2 de la figure suivante

MINIMISATION D'UN AFD

ÉTATS INACCESSIBLES (2)

L'AFD M_2 obtenu en enlevant les états accessibles de M_1



REGLE 1

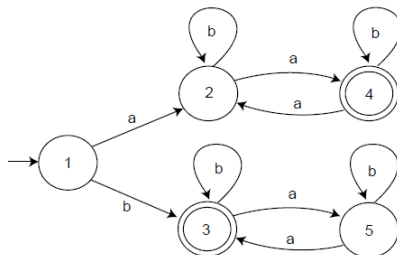
Soit M un AFD et soit M_0 l'AFD obtenu en enlevant tous les états accessibles de M . Alors $L(M) = L(M_0)$.

Peut-on encore réduire le nombre d'états de M_2 ?

MINIMISATION D'UN AFD

ÉTATS INACCESSIBLES (2)

L'AFD M_2 obtenu en enlevant les états accessibles de M_1



REGLE 1

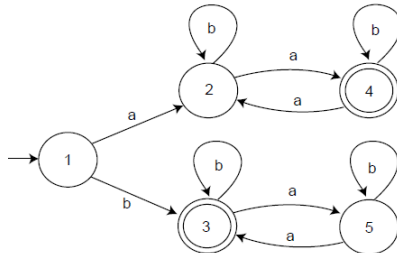
Soit M un AFD et soit M_0 l'AFD obtenu en enlevant tous les états accessibles de M . Alors $L(M) = L(M_0)$.

Peut-on encore réduire le nombre d'états de M_2 ?

MINIMISATION D'UN AFD

ÉTATS INACCESSIBLES (2)

L'AFD M_2 obtenu en enlevant les états accessibles de M_1



REGLE 1

Soit M un AFD et soit M_0 l'AFD obtenu en enlevant tous les états accessibles de M . Alors $L(M) = L(M_0)$.

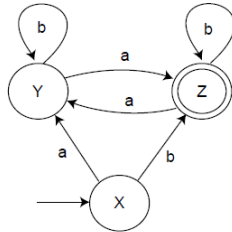
Peut-on encore réduire le nombre d'états de M_2 ?

MINIMISATION D'UN AFD

L'AFD MINIMAL

La réponse est encore oui !

L'AFD minimal M_3 équivaut à M_1



Cet automate est le plus petit AFD reconnaissant $L(M_1)$.

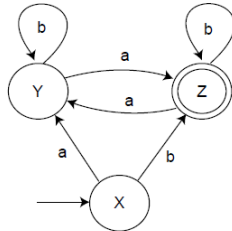
Comment a-t-on obtenu cet automate ?

MINIMISATION D'UN AFD

L'AFD MINIMAL

La réponse est encore oui !

L'AFD minimal M_3 équivaut à M_1



Cet automate est le plus petit AFD reconnaissant $L(M_1)$.

Comment a-t-on obtenu cet automate ?

MINIMISATION D'UN AFD

Dénotons par L_i ($i = 1, 2, 3, 4, 5$) l'ensemble des mots menant à un état final à partir de l'état i . Plus formellement, définissons :

$$L_i = \{w \in \{a, b\}^* \mid \delta(i, w) \in F\}$$

- $L_2 = L_5 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre impair de } a\}$
- $L_3 = L_4 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre pair de } a\}$

À chaque fois que l'on se retrouve dans l'état 2, on pourrait poursuivre l'exécution à partir de l'état 5 sans rien changer au résultat.

DÉFINITION (ÉTATS ÉQUIVALENTS)

Deux états i et j sont équivalents si $L_i = L_j$. En d'autres termes, i et j sont équivalents si pour tout $w \in A^*$ on a :
 $\delta(i, w) \in F$ si et seulement si $\delta(j, w) \in F$

MINIMISATION D'UN AFD

Dénotons par L_i ($i = 1, 2, 3, 4, 5$) l'ensemble des mots menant à un état final à partir de l'état i . Plus formellement, définissons :

$$L_i = \{w \in \{a, b\}^* \mid \delta(i, w) \in F\}$$

- $L_2 = L_5 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre impair de } a\}$
- $L_3 = L_4 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre pair de } a\}$

À chaque fois que l'on se retrouve dans l'état 2, on pourrait poursuivre l'exécution à partir de l'état 5 sans rien changer au résultat.

DÉFINITION (ÉTATS ÉQUIVALENTS)

Deux états i et j sont équivalents si $L_i = L_j$. En d'autres termes, i et j sont équivalents si pour tout $w \in A^*$ on a :
 $\delta(i, w) \in F$ si et seulement si $\delta(j, w) \in F$

MINIMISATION D'UN AFD

Dénotons par L_i ($i = 1, 2, 3, 4, 5$) l'ensemble des mots menant à un état final à partir de l'état i . Plus formellement, définissons :

$$L_i = \{w \in \{a, b\}^* \mid \delta(i, w) \in F\}$$

- $L_2 = L_5 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre impair de } a\}$
- $L_3 = L_4 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre pair de } a\}$

À chaque fois que l'on se retrouve dans l'état 2, on pourrait poursuivre l'exécution à partir de l'état 5 sans rien changer au résultat.

DÉFINITION (ÉTATS ÉQUIVALENTS)

Deux états i et j sont équivalents si $L_i = L_j$. En d'autres termes, i et j sont équivalents si pour tout $w \in A^*$ on a :
 $\delta(i, w) \in F$ si et seulement si $\delta(j, w) \in F$

MINIMISATION D'UN AFD

Dénotons par L_i ($i = 1, 2, 3, 4, 5$) l'ensemble des mots menant à un état final à partir de l'état i . Plus formellement, définissons :

$$L_i = \{w \in \{a, b\}^* \mid \delta(i, w) \in F\}$$

- $L_2 = L_5 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre impair de } a\}$
- $L_3 = L_4 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre pair de } a\}$

À chaque fois que l'on se retrouve dans l'état 2, on pourrait poursuivre l'exécution à partir de l'état 5 sans rien changer au résultat.

DÉFINITION (ÉTATS ÉQUIVALENTS)

Deux états i et j sont équivalents si $L_i = L_j$. En d'autres termes, i et j sont équivalents si pour tout $w \in A^*$ on a :
 $\delta(i, w) \in F$ si et seulement si $\delta(j, w) \in F$

MINIMISATION D'UN AFD

Dénotons par L_i ($i = 1, 2, 3, 4, 5$) l'ensemble des mots menant à un état final à partir de l'état i . Plus formellement, définissons :

$$L_i = \{w \in \{a, b\}^* \mid \delta(i, w) \in F\}$$

- $L_2 = L_5 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre impair de } a\}$
- $L_3 = L_4 = \{w \in \{a, b\}^* \mid w \text{ contient un nombre pair de } a\}$

À chaque fois que l'on se retrouve dans l'état 2, on pourrait poursuivre l'exécution à partir de l'état 5 sans rien changer au résultat.

DÉFINITION (ÉTATS ÉQUIVALENTS)

Deux états i et j sont équivalents si $L_i = L_j$. En d'autres termes, i et j sont équivalents si pour tout $w \in A^*$ on a :
 $\delta(i, w) \in F$ si et seulement si $\delta(j, w) \in F$

MINIMISATION D'UN AFD

ÉTATS ÉQUIVALENTS

Ainsi, dans l'AFD M_2 les états 2 et 5 sont équivalents et les états 3 et 4 sont équivalents. L'état 1 n'est équivalent à aucun autre état.

REMARQUE

si i et j sont deux états équivalents d'un automate donne, alors à chaque fois que l'on se retrouve dans l'état i , on pourrait poursuivre l'exécution à partir de l'état j sans rien changer au résultat. Cela conduit à l'observation suivante :

OBSERVATION

Il est toujours possible d'ajouter des ε -transitions entre deux états équivalents sans rien changer au langage reconnu.

MINIMISATION D'UN AFD

ETATS ÉQUIVALENTS

Ainsi, dans l'AFD M_2 les états 2 et 5 sont équivalents et les états 3 et 4 sont équivalents. L'état 1 n'est équivalent à aucun autre état.

REMARQUE

si i et j sont deux états équivalents d'un automate donne, alors à chaque fois que l'on se retrouve dans l'état i , on pourrait poursuivre l'exécution à partir de l'état j sans rien changer au résultat. Cela conduit à l'observation suivante :

OBSERVATION

Il est toujours possible d'ajouter des ε -transitions entre deux états équivalents sans rien changer au langage reconnu.

MINIMISATION D'UN AFD

ETATS ÉQUIVALENTS

Ainsi, dans l'AFD M_2 les états 2 et 5 sont équivalents et les états 3 et 4 sont équivalents. L'état 1 n'est équivalent à aucun autre état.

REMARQUE

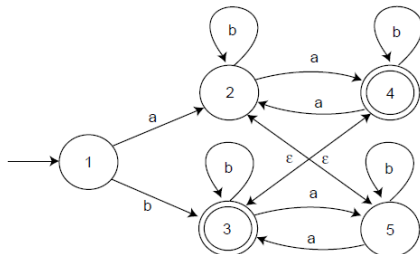
si i et j sont deux états équivalents d'un automate donne, alors à chaque fois que l'on se retrouve dans l'état i , on pourrait poursuivre l'exécution à partir de l'état j sans rien changer au résultat. Cela conduit à l'observation suivante :

OBSERVATION

Il est toujours possible d'ajouter des ε -transitions entre deux états équivalents sans rien changer au langage reconnu.

MINIMISATION D'UN AFD

L'AFN M_4 obtenu a partir de M_2 en ajoutant des ε -transitions entre les paires d'états équivalents.

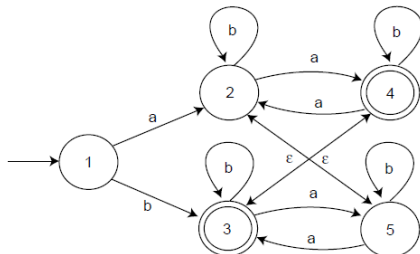


REMARQUE

L'automate obtenu n'est plus déterministe. Cependant, on peut obtenir un AFD M_5 à partir de L'AFN M_4 en utilisant la méthode de détermination.

MINIMISATION D'UN AFD

L'AFN M_4 obtenu a partir de M_2 en ajoutant des ε -transitions entre les paires d'états équivalents.

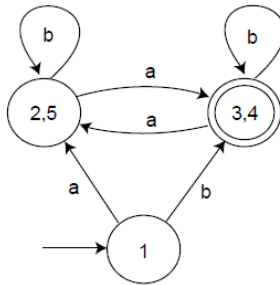


REMARQUE

L'automate obtenu n'est plus déterministe. Cependant, on peut obtenir un AFD M_5 à partir de L'AFN M_4 en utilisant la méthode de détermination.

MINIMISATION D'UN AFD

L'AFD M_5 obtenu a partir de L'AFN M_4



LA SUITE

Analyse Santaxique

LECTURES COMPLÉMENTAIRES I

-  Modern compilers in C/Java/ML (3 livres jumeaux). de A. Appel. Cambridge University Press, 1998.
-  Compilateurs : Principes, techniques et outils. de A. Aho, R. Sethi, and J. Ullman. Dunod, 2000.
-  Programming Language Processors - Compilers and Interpreters de David Watt, Prentice-Hall International Series in Computer Science.