### Reconnaissance de motifs

UMLV ©

### **Problème**

localiser les segments d'un texte décrits par une expression rationnelle (régulière) r

texte t

segment *x* décrit par *r* 

### **Applications**

édition : vi, emacs , sed, ed, ... recherche : grep, egrep , ... traduction : lex, sed , ... langages : awk, perl, ...

compression: compress, gzip, ...

701

# GREP

UMLV ©

unix > grep toto t.txt

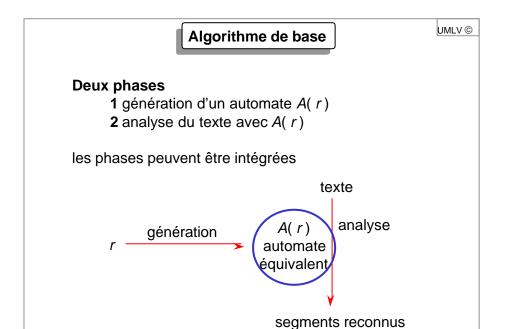
produit les lignes de t.txt qui contiennent le mot toto

### **Applications**

recherche de fichiers par le contenu contrôle du contenu d'un fichier

### **Versions**

egrep, fgrep, ...



703

### Expression / langage

UMLV ©

```
Langage représenté L( r )
Expression rationnelle r
                                                   {a} , {b} , ...
                                                                                   a, bÎ A
          a, b, ...
          Æ, e

        \, E \, , {mot vide}
                                                   L(u) , L(v) ]
          [u, v]
          (u)
                                                   L(u)
                                                   L(u) \stackrel{.}{\to} L(v)
          U + V
                                                    \left\{ \begin{array}{ll} xy \mid x \ \widehat{\mathbf{1}} & L(u), \ y \ \widehat{\mathbf{1}} & L(v) \end{array} \right\} \\ \left\{ \begin{array}{ll} x_1 x_2 ... x_k \mid k \ ^3 \ 0, \ x_i \ \widehat{\mathbf{1}} & L(u) \end{array} \right\} 
          uv
          u*
                                                   {écritures binaires des entiers pairs}
          1(0+1)*0
          (a*b)*a*
                                                   {suites finies de a et b}
           (c*1)*c*f
                                                   {fichiers textes}
```

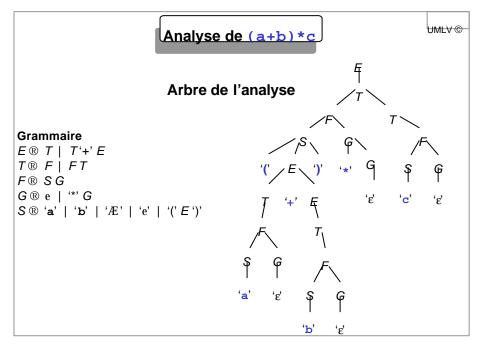
### Analyse des expressions

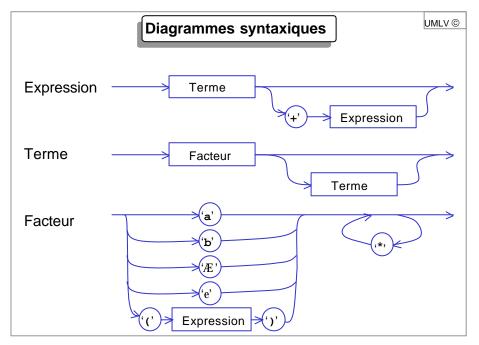
UMLV ©

Une grammaire des expressions rationnelles

*E* expression, *T* terme *F* facteur, *S* facteur simple

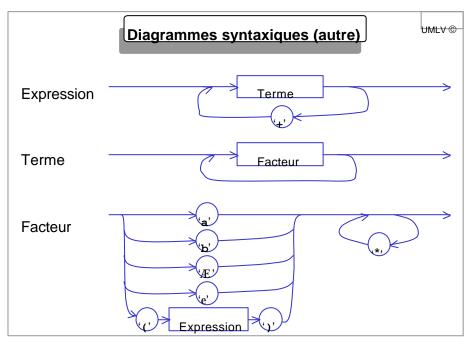
(a+b)\*c {suites finies de a et b terminées par c}





```
UMLV ©
            Algorithme d'analyse
caractère car; /* caractère suivant, global */
Analyse(expression rationnelle r){
     car - premier caractère de r ;
      Expression();
      si(car ≠ fin d'expression)
            erreur();
Expression(){
      Terme();
      si(car = '+'){
            car - caractère suivant;
            Expression();
}
Terme(){
      Facteur();
      {\tt si(car~\hat{I}~\{`a',`b',`E',`e',`(')})} \ {\tt Terme();}
```

### Algorithme d'analyse (suite) UMLV © Facteur(){ $\mathbf{si}(\texttt{car}~\widehat{I}~\{\text{`a'},\text{`b'},\text{`$\rlap{$\cal E$}'},\text{`e'}\})\,\{$ car ¬ caractère suivant; }sinon si(car = '('){ car ¬ caractère suivant; Expression(); **si**(car = ')') car - caractère suivant; sinon erreur(); }sinon erreur(); tant que(car = '\*') car - caractère suivant; }



# Algorithme d'analyse (autre)

UMLV ©

```
caractère car; /* caractère suivant, global */
Analyse(expression rationnelle r){
     car - premier caractère de r ;
     Expression();
     si(car ≠ fin d'expression)
           erreur();
Expression(){
     Terme();
     tant que(car = '+'){
           car - caractère suivant;
           Terme();
Terme(){
     répéter
           Facteur();
     tant que(car \hat{I} {'a', 'b', 'Æ', 'e', '('})
}
```

711

### Expression / automate

UMLV ©

A(r) automate associé à l'expression rationnelle r sur l'alphabet  ${\bf A}$ 

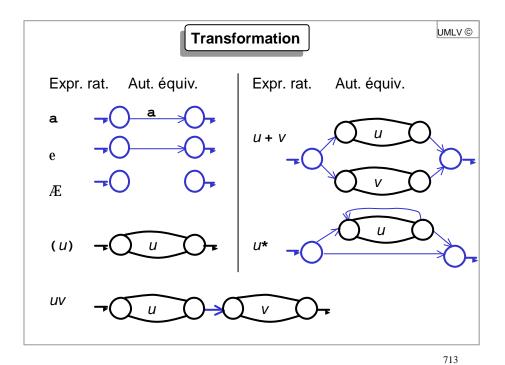
### Invariants de la construction

A(r) possède un seul état initial i un seul terminal t aucune flèche n'entre dans i aucune flèche ne sort de t

Graphique

A(r):





### \_\_\_\_

UMLV ©

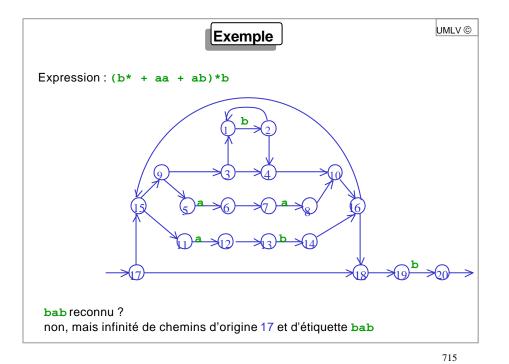
### **Proposition**

A(r) reconnaît le langage L(r)

Automate obtenu

### Propriétés supplémentaires

A(r) possède 2./r/ états, sans compter '('ni ')' de chaque état sortent au plus 2 flèches si exactement 2, elles sont sans étiquette sur chaque état arrivent au plus 2 flèches



UMLV © **Implantation adaptée** état = indice table des flèches 2 struct{ 3 caractère lettre; 5 état suiv1, suiv2; 6 } Trans[MAX]; 7 automate (identifié à 2 états) 8 struct{ 9 état initial; 10 état terminal; 11 } automate; 12 13 14 15 16 17 18 р initial = 1719 terminal = 20

Algorithme de traduction

UMLV ©

# Production de l'automate A( r) associé à l'expression rationnelle r

717

### Algorithme de traduction (suite)

UMLV ©

```
Expression(){
    A ¬ Terme();
    tant que(car = '+'){
        car ¬ caractère suivant;
        B ¬ Terme();
        A ¬ Union(A,B);
    }
    retour A;
}

Terme(){
    A ¬ Facteur();
    tant que(car Î {'a', 'b', 'Æ', 'e', '('})}{
        B ¬ Facteur();
        A ¬ Produit(A,B);
    }
    retour A;
}
```

UMLV ©

### Algorithme de traduction (suite)

```
Facteur(){
     si(car \hat{I} \{ a', b', A', e' \}) 
           A ¬ Automate-élémentaire(car);
           car - caractère suivant;
      }sinon si(car = '('){
           car - caractère suivant;
           A ¬ Expression();
           si(car = ')')
                car - caractère suivant;
           sinon
                 erreur();
     }sinon
          erreur();
      tant que(car = '*'){
          A ¬ Etoile(A);
           car ¬ caractère suivant;
     retour A ;
```

719

# Reconnaissance (avec automate non-déterministe)

UMLV ©

### Écrire une fonction

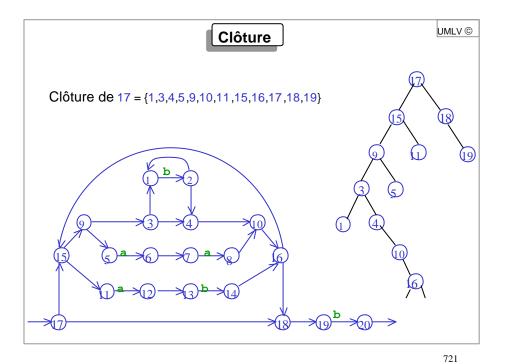
booléen reconnaît ( automate A, mot x ) qui prend la valeur vraie ssi il existe un chemin d'étiquette x depuis A.initial jusqu'à A.terminal

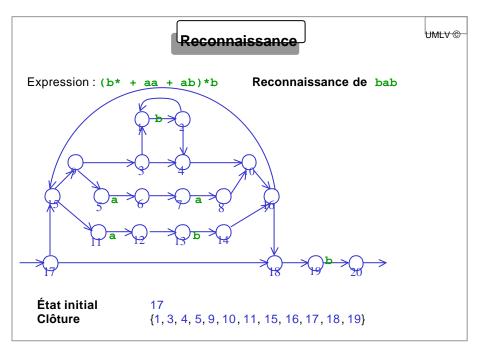
### **Programmation par**

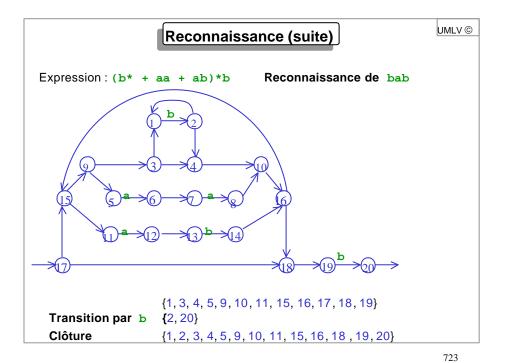
essais et erreurs (attention aux cycles) programmation dynamique (mémorisation des états possibles)

### Clôture

clôture(p) = { états accessibles à partir p par chemins sans étiquette }







Reconnaissance (suite)

Expression: (b\* + aa + ab)\*b Reconnaissance de bab

(1, 2, 3, 4, 5, 9, 10, 11, 15, 16, 18, 19, 20)

Transition par a {6, 12}
Clôture {6, 7, 12, 13}

# Expression: (b\* + aa + ab)\*b Reconnaissance de bab

725

### Algorithme de reconnaissance

{1, 3, 4, 5, 9, 10, 11, 14, 15, 16, 18, 19}

UMLV ©

```
booléen reconnaît(automate A, mot x){
    ensemble E;
    E ¬ clôture(A.initial);
    tant que(non fin de x){
        car ¬ lettre suivante de x;
        E ¬ clôture(transition(E,car));
    }
    si(A.terminal Î E) retour(vrai);
    sinon retour(faux);
}
```

Lecture séquentielle du mot.

Transition par b {14}

bab non reconnu car 20 non atteint

Mémorisation (tampon) d'une seule lettre de x.

Temps : O(nb états de A . | x |)
(si bonne gestion de l 'ensemble E)

Recherche de motif

UMLV ©

Motif défini par l'expression rationnelle r A(r) automate associé à r

Motif dans le texte

texte t

segment x reconnu par A(r),  $x \hat{1} L(r)$ 

### Principe de la recherche

comme reconnaissance avec : ajout l'état initial à chaque position arrêt de la recherche dès que l'état terminal est atteint

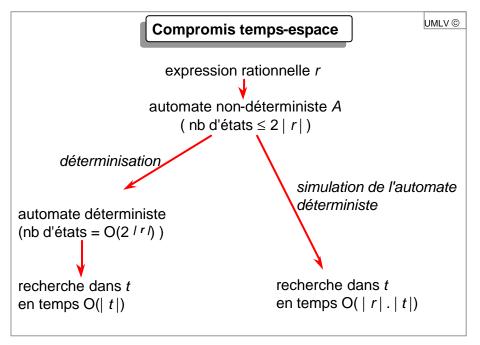
727

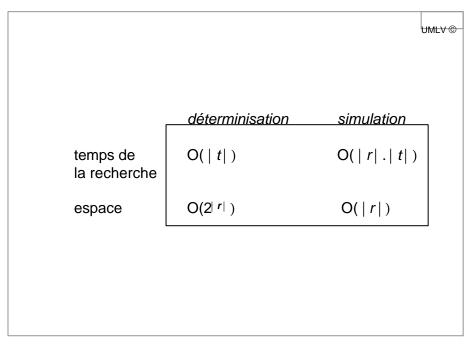
### Algorithme de recherche

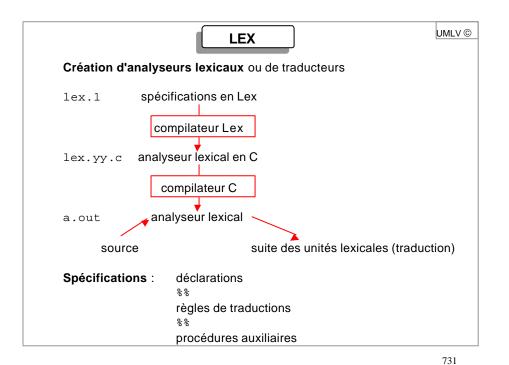
UMLV ©

```
booléen recherche(automate A, mot t){
    ensemble E;
    E ¬ clôture(A.initial);
    tant que(non fin de t){
        car ¬ lettre suivante de t;
        E ¬ clôture({A.initial}Ètransition(E,car));
        si(A.terminal Î E) retour(vrai);
    }
    retour(faux);
}

Lecture séquentielle du texte.
Mémorisation (tampon) d'une seule lettre de t.
Temps: O(nb états de A. | t|)
    (si bonne gestion de l'ensemble E)
```







```
UMLV ©
% {/* définitions des constantes littérales */
         PPQ, PPE, EGA, DIF, PGQ, PGE,
         SI, ALORS, SINON, ID, NB, OPREL
/* définitions régulières */
délim
           [ \t \n]
           {delim}+
b1
           [A-Za-z]
lettre
           [0-9]
chiffre
id
           {lettre}+({lettre}|{chiffre})*
           {chiffre}+(\.{chiffre}+)?(E[+\-]?{chiffre}+)?
nombre
응응
{bl}
            { \ \ \ } /* pas d'action et pas de retour */ { \ \ }
si
            {return (SI); }
            {return (ALORS); }
alors
sinon
            {return (SINON); }
{id}
            {yylval = RangerId ( ) ; return (ID) ; }
{nombre}
            {yylval = RangerNb ( ) ; return (NB) ; }
            {yylval = PPQ ; return (OPREL) ; }
" < "
            {yyval = PPE ; return (OPREL) ; }
" <= "
" = "
            {yyval = EGA ; return (OPREL) ; }
" <> "
            {yyval = DIF ; return (OPREL) ;
            {yyval = PGQ ; return (OPREL) ;
            {yyval = PGE ; return (OPREL) : }
```

RangerId() {

/\* procédure qui range dans la table des symboles
le lexème dont le premier caractère est pointé
par yytext et dont la longueur est yyleng,
et retourne un pointeur sur son entrée \*/
...
}

RangerNb() {

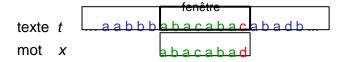
/\* procédure similaire pour ranger un lexème
qui est un nombre \*/
...
}

733

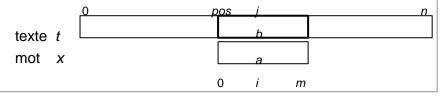
### Recherche d'un mot

UMLV ©

Recherche des positions du mot *x* dans le texte *t* Par balayages (gauche-droite) et décalages



### Variables de travail



### Algorithme naïf

UMLV ©

### Temps maximal O(m.n)

» m.n comparaisons au pire

### Temps moyen O(n)

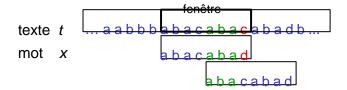
» 2.*n* comparaisons moyennes sur alphabet binaire (équiprobabilité et indépendance)

735

### Méthode de Knuth, Morris & Pratt

UMLV ©

KMP : balayages gauche-droite et décalages compatibles

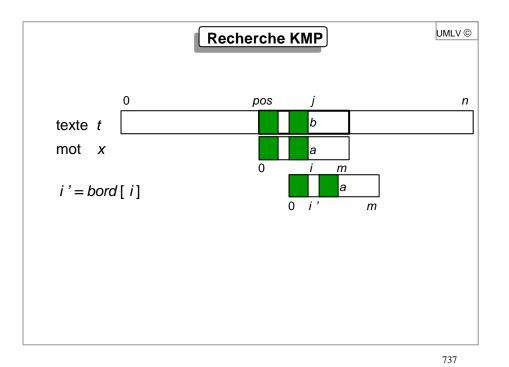


0

Bords: Bord(abacaba) = aba

0

Longueur des bords des préfixes : bord[i] = |Bord(x[0..i-1])|i 0 1 2 3 4 5 6 7 8 x[i]b d a а С а b а



# Algorithme

UMLV ©

# Temps O(n)

moins de 2.*n* comparaisons recherche séquentielle (*j* ne décroît jamais) recherche linéaire mais pas en temps réel

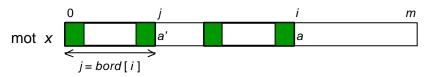
### Calcul des bords

UMLV ©

**Bords** : Bord(u) = plus long préfixe et suffixe propre de u

Note: le bord d'un bord est un bord

```
Bord(a b a c a b a c) = a b a c
Bord(a b a c a b a b) = a b
Bord(a b a c a b a d) = \epsilon
```



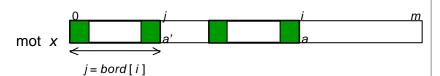
si a=a' bord [ i+1 ] = bord [ i ] +1 = j+1sinon faire le calcul avec a sur Bord(x[0 .. j-1])

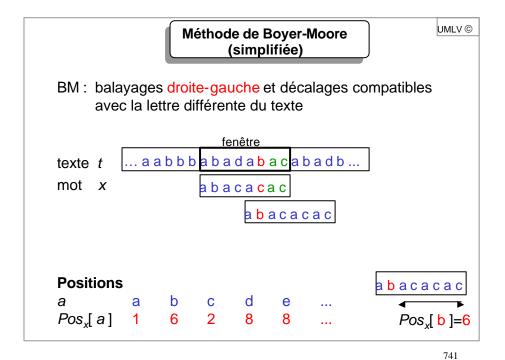
739

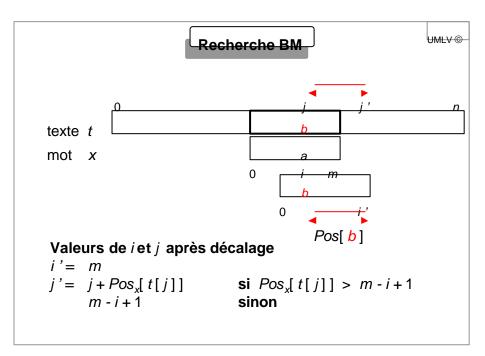
# BORDS(mot x, longueur m) {

UMLV ©

**Temps** O(m), moins de 2.m comparaisons} Comme la recherche!







# Algorithme BM

UMLV ©

### Temps maximal O(m.n)

**Temps**, sur des textes en langue naturelle, proche du minimal :  $\Theta(n/m)$ 

743

### Calcul de Pos

UMLV ©

```
Positions(mot x, longueur m){
    pour chaque lettre a faire
        Pos[a] ¬ m ;
    pour i ¬ 0 à m-2 faire
        Pos[x[i]] ¬ m-i-1 ;
    retour Pos ;
}
```



Temps O(m + card A)