

Regular Pattern Matching

Maxime Crochemore
King's College London
<http://www.dcs.kcl.ac.uk/staff/mac/>

1

Regular Pattern Matching

Problem

locate segments of a text t
described by a regular expression r

text t



segment x
described by r

Applications

text editing : vi, emacs , sed, ed, ...
search : grep, egrep , ...
translation : lex, sed , ...
languages : awk, perl, ...
compression : compress, gzip, ...
virus detection, etc.

2

GREP

```
unix > grep toto t.txt
```

produces lines of `t.txt`
containing the word `toto`

Applications

searching files by contents
checking the content of a file

Versions

`egrep`, `fgrep`, ...

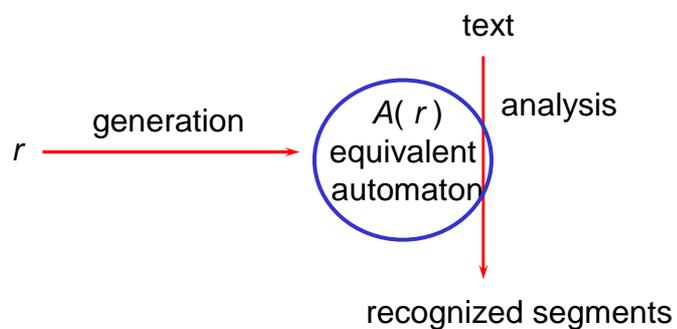
3

Schema of algorithm

Two phases

- 1 generation of automaton $A(r)$
- 2 text analysis with $A(r)$

phases can be integrated (pipelining)



4

Expression / language

Regular expression r	Corresponding language $L(r)$
a, b, \dots	$\{a\}, \{b\}, \dots \quad a, b \in A$
\emptyset, ϵ	$\emptyset, \{\text{empty word}\}$
$[u, v]$	$L(u), L(v)$
(u)	$L(u)$
$u + v$	$L(u) \cup L(v)$
uv	$\{xy \mid x \in L(u), y \in L(v)\}$
u^*	$\{x_1 x_2 \dots x_k \mid k \geq 0, x_i \in L(u)\}$
$1(0+1)^*0$	{binary representations of even integers}
$(a*b)^*a^*$	{finite strings of a and b }
$(c*1)^*c*f$	{textual files}

5

Description of expressions

A grammar of regular expressions

$$\begin{aligned}
 E &\rightarrow T \mid T '+' E \\
 T &\rightarrow F \mid FT \\
 F &\rightarrow SG \\
 G &\rightarrow \epsilon \mid '*' G \\
 S &\rightarrow 'a' \mid 'b' \mid '\emptyset' \mid '\epsilon' \mid '(' E ')'
 \end{aligned}$$

Variables

E expression T term
 F factor S simple factor G for stars

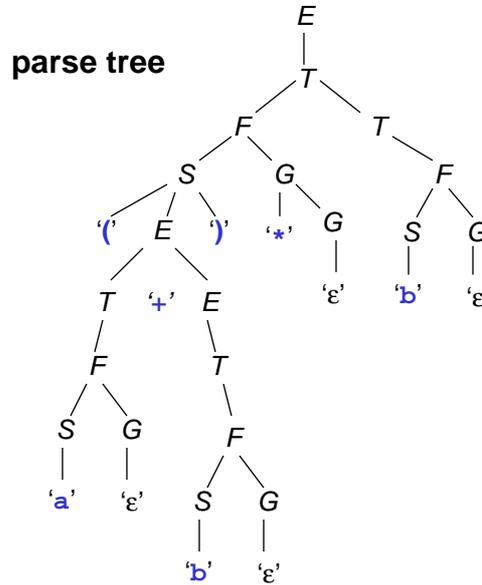
$(a+b)^*b$ {finite strings of **a** and **b** ending with **b**}

6

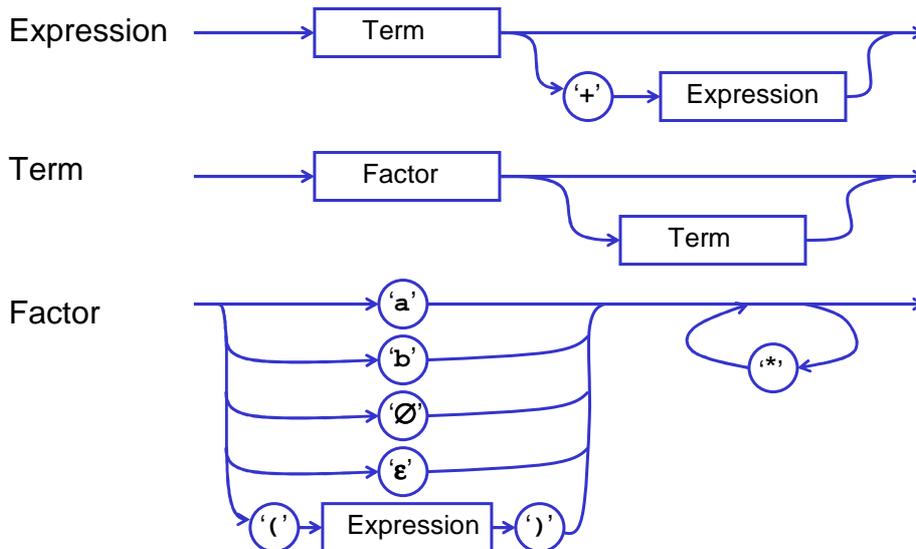
Analysis of (a+b)*b

Grammar

$E \rightarrow T \mid T '+' E$
 $T \rightarrow F \mid FT$
 $F \rightarrow SG$
 $G \rightarrow \epsilon \mid '*' G$
 $S \rightarrow 'a' \mid 'b' \mid \emptyset \mid '\epsilon' \mid '(E)'$



Syntactic diagrams



Parsing algorithm

```
character car; /* next character, global */

Parser(regular expression r){
    car ← first character of r ;
    Expression();
    if(car ≠ end of expression)
        erreur();
}
Expression(){
    Term();
    if(car = '+'){
        car ← next character;
        Expression();
    }
}
Term(){
    Factor();
    if(car ∈ {'a', 'b', '∅', 'ε', '('})
        Term();
}
}
```

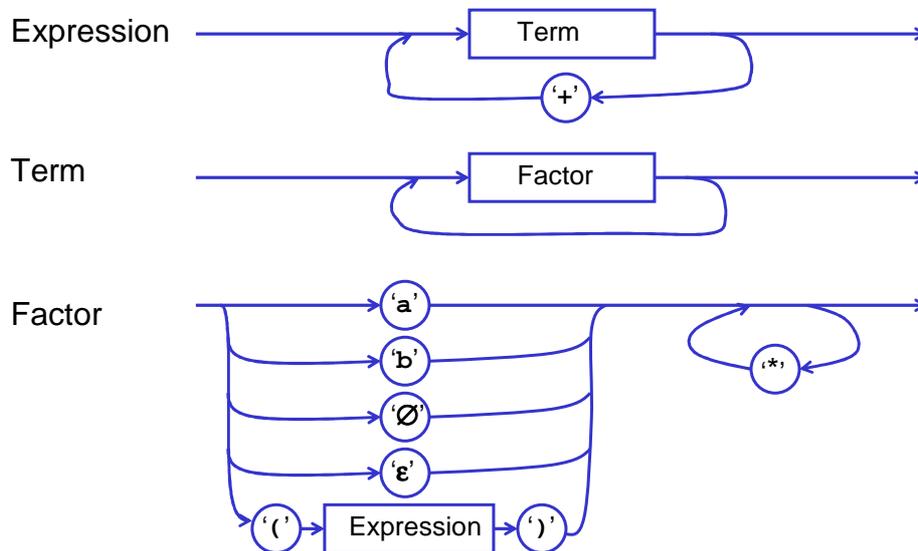
9

Parsing algorithm (followed)

```
Factor(){
    if(car ∈ {'a', 'b', '∅', 'ε'}){
        car ← next character;
    }else if(car = '('){
        car ← next character;
        Expression();
        if(car = ')')
            car ← next character;
        else
            error();
    }else
        error();
    while(car = '*')
        car ← next character;
}
}
```

10

Syntactic diagrams (2)



11

Parsing algorithm (2)

```
character car; /* next character, global */

Parser(regular expression r){
    car ← first character of r ;
    Expression();
    if(car ≠ end of expression)
        error();
}

Expression(){
    Term();
    while(car = '+'){
        car ← next character;
        Term();
    }
}

Term(){
    repeat
        Factor();
    while(car ∈ {'a', 'b', '∅', 'ε', '('})
}


```

12

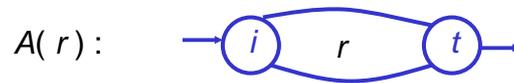
Expression / automaton

$A(r)$ automaton associated with regular expression r on the alphabet A

Invariants of the construction

- $A(r)$ has only one initial state i
- only one terminal state t
- no arc enters i
- no arc leaves t

Graphic representation



Transformation

Reg expr. Equiv. aut.



Reg expr.

Equiv. aut.



Generated automaton

Proposition

$A(r)$ accepts the language $L(r)$

Additional properties

$A(r)$ has $2 \cdot |r|$ states [not counting '(' nor ')']

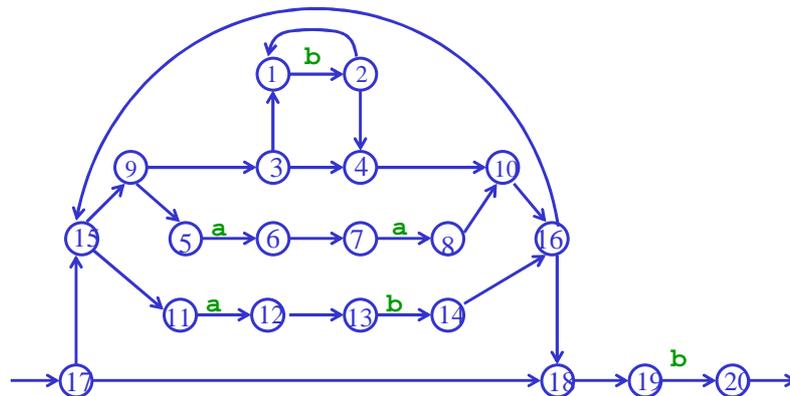
output degrees are at most 2

if exactly 2: arcs have empty labels

on each arc arrive at most 2 arcs

Example

Expression : $(b^* + aa + ab)^*b$



is bab accepted?

no, but infinitely many paths starting at 17 and labelled by bab

Special implementation

```

state = index on table
table of arcs
struct{
    char lettre;
    state succ1, succ2;
} Trans[MAX];
automaton (identified to 2 states)
struct{
    state initial;
    state terminal;
} automaton;
    
```

Trans

p	'a'	p1
q		q1

initial = 17
terminal = 20

1	'b'	2	
2		1	4
3		1	4
4		10	
5	'a'	6	
6		7	
7	'a'	8	
8		10	
9		3	5
10		16	
11	'a'	12	
12		13	
13	'b'	14	
14		16	
15		9	11
16		15	18
17		15	18
18		19	
19	'b'	20	
20			

Translation algorithm

Production of automaton $A(r)$
associated with the regular expression r

```

Translation(regular expression r){
    car ← first character of r ;
    A ← Expression();
    if(car ≠ end of expression)
        erreur();
    else
        return A;
}
    
```

Translation algorithm (followed)

```
Expression(){
  A ← Term();
  while(car = '+'){
    car ← next character;
    B ← Term();
    A ← Union(A,B);
  }
  retour A;
}
Term(){
  A ← Factor();
  while(car ∈ {'a', 'b', 'Ø', 'ε', '('}){
    B ← Factor();
    A ← Product(A,B);
  }
  return A;
}
```

19

Translation algorithm (followed)

```
Factor(){
  if(car ∈ {'a', 'b', 'Ø', 'ε'}){
    A ← Elementary-Automaton(car);
    car ← next character;
  }else if(car = '('){
    car ← next character;
    A ← Expression();
    if(car = ')')
      car ← next character;
    else
      error();
  }else
    error();
  while(car = '*'){
    A ← Star(A);
    car ← next character;
  }
  return A ;
}
```

20

Recognition (with non-deterministic automaton)

Function

`boolean accept(automaton A, string x)`

Value is true iff there exists a path labelled by x
from $A.initial$ to $A.terminal$

Programming:

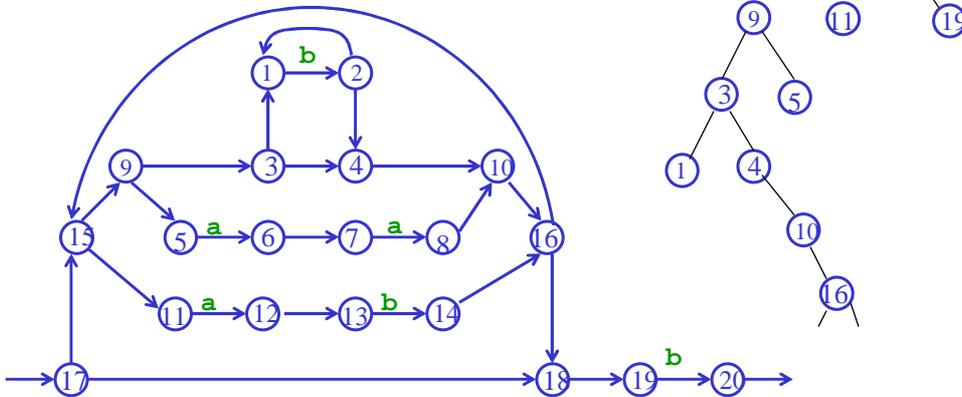
backtracking (mind cycles in the automaton!)
dynamic programming (storing possible states)

Closure

$closure(p) = \{ \text{states accessible from } p \text{ using paths having empty labels} \}$

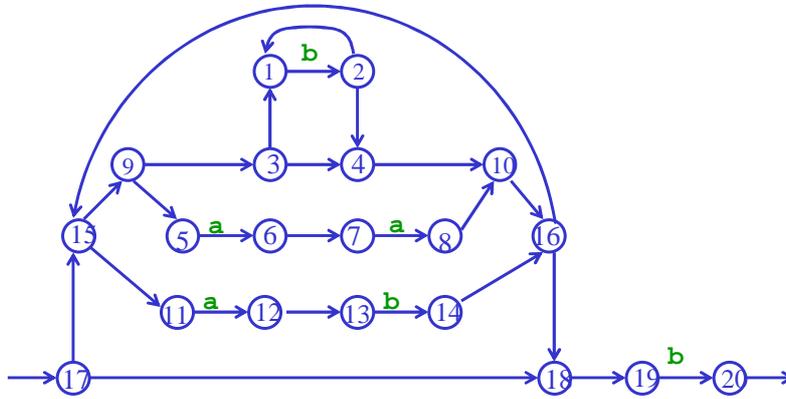
Closure

$closure(17) = \{ 1, 3, 4, 5, 9, 10, 11, 15, 16, 17, 18, 19 \}$



Recognition

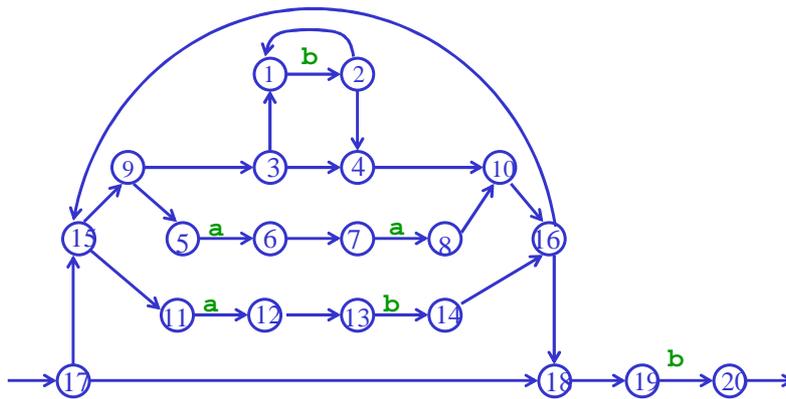
Expression : $(b^* + aa + ab)^*b$ Recognition of bab



Initial state 17
 Closure {1, 3, 4, 5, 9, 10, 11, 15, 16, 17, 18, 19}

Recognition (followed)

Expression : $(b^* + aa + ab)^*b$ Recognition of bab

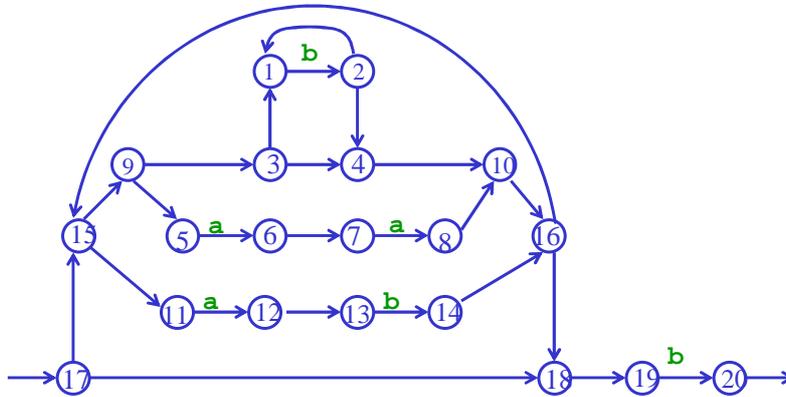


Transition by b {1, 3, 4, 5, 9, 10, 11, 15, 16, 17, 18, 19}
 {2, 20}
 Closure {1, 2, 3, 4, 5, 9, 10, 11, 15, 16, 18, 19, 20}

Recognition (followed)

Expression : $(b^* + aa + ab)^*b$

Recognition of bab

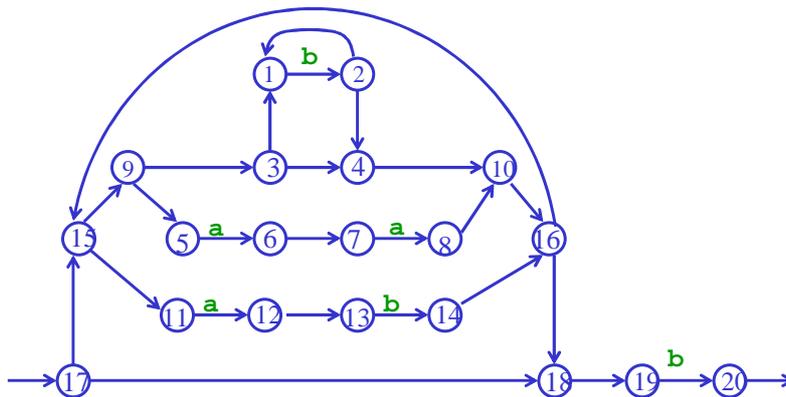


Transition by **a** {1, 2, 3, 4, 5, 9, 10, 11, 15, 16, 18, 19, 20}
 {6, 12}
 Closure {6, 7, 12, 13}

Recognition (followed)

Expression : $(b^* + aa + ab)^*b$

Recognition of bab



Transition by **b** {6, 7, 12, 13}
 {14}
 Closure {1, 3, 4, 5, 9, 10, 11, 14, 15, 16, 18, 19}
 bab not recognized since 20 not reached

Recognition algorithm

```
boolean accept(automaton A, string x){
    set E;
    E ← closure(A.initial);
    while(not end of x){
        car ← next letter of x;
        E ← closure(transition(E,car));
    }
    if (A.terminal ∈ E) return(true);
    else return(false);
}
```

On-line treatment of x : buffer of only one letter on x .

Running time: $O(\#\{\text{states of } A\} \cdot |x|)$

(with adequate implementation of E)

27

Regular pattern matching

Pattern specified by regular expression r

$A(r)$ automaton associated with r

Pattern in text

text t

segment x
accepted by $A(r)$, $x \in L(r)$

Searching principles

like recognition except:

- initial state is added at each step
- search terminates as soon as terminal state is reached

28

Matching algorithm

```
boolean regular-matching(automaton A, text t){
  set E;
  E ← closure(A.initial);
  while(not end of t){
    car ← next letter of t;
    E ← closure({A.initial} ∪ transition(E, car));
    if(A.terminal ∈ E) return(true);
  }
  return(false);
}
```

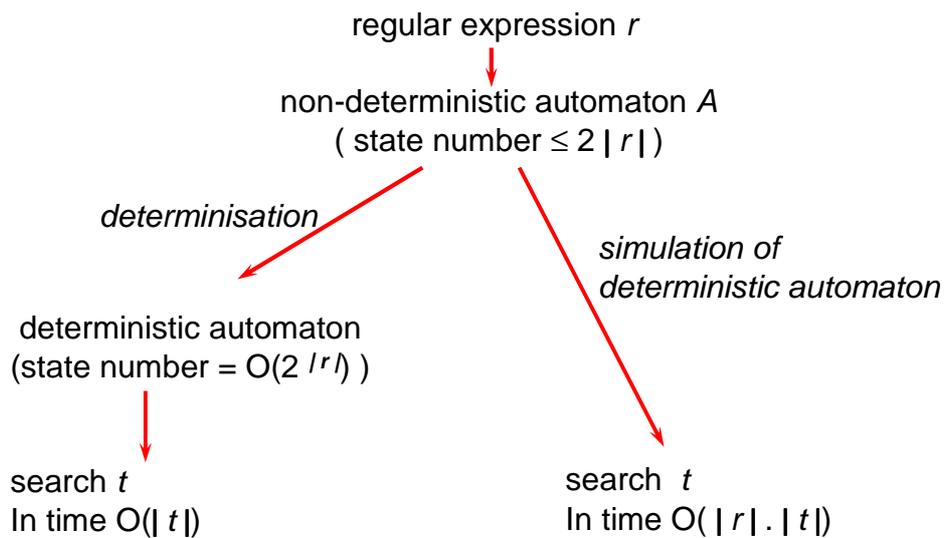
On-line treatment of the text: buffer of only one letter on t

Running time: $O(\#\{\text{states of } A\} \cdot |t|)$

(with adequate implementation of E)

29

Time-space tradeoff

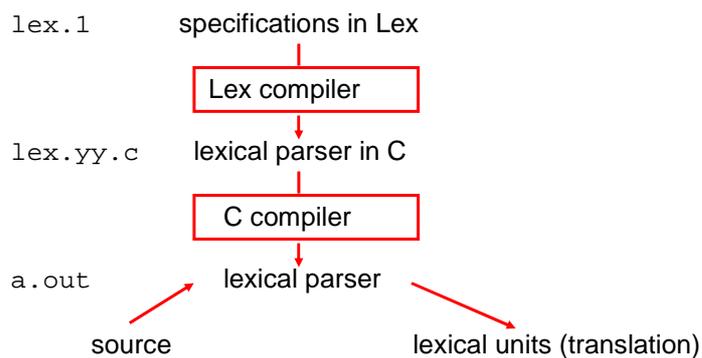


30

	<i>determinisation</i>	<i>simulation</i>
Search time	$O(t)$	$O(r \cdot t)$
space	$O(2^{ r })$	$O(r)$

LEX

production of lexical parsers or translators



Specifications : declarations
 %%
 translation rules
 %%
 auxiliary procedures

```

% { /* definitions of constants */
      LT, LE, EQ, NEQ, GT, GE,
      IF, THEN, ELSE, ID, NB, RELOP
% }
/* definitions of regular expressions */
spacer      [ \t\n]
bl          {delim}+
letter      [A-Za-z]
digit       [0-9]
id          {lettre}+({lettre}||{chiffre})*
number      {chiffre}+(\.{chiffre}+)?(E[+\-]?{chiffre}+)?
%%
{bl}        { /* no action, no return */ }
if          {return (IF) ; }
then        {return (THEN) ; }
else        {return (ELSE) ; }
{id}        {yyval = RangerId ( ) ; return (ID) ; }
{number}    {yyval = RangerNb ( ) ; return (NB) ; }
"<"        {yyval = LT ; return (RELOP) ; }
"<="       {yyval = LE ; return (RELOP) ; }
"="         {yyval = EQ ; return (RELOP) ; }
">"        {yyval = NEQ ; return (RELOP) ; }
">"        {yyval = GT ; return (RELOP) ; }
">="       {yyval = GE ; return (RELOP) ; }

```

33

```

%%
StoreId(){
    /* procedure that store in the table of symbols
       a lexical unit which first letter is pointed by
       yytext and which length is yyleng;
       returns a pointer on table entry */
    ...
}

StoreNb(){
    /* similar procedure to store a number (constant) */
    ...
}

```

34