

Script generated by TTT

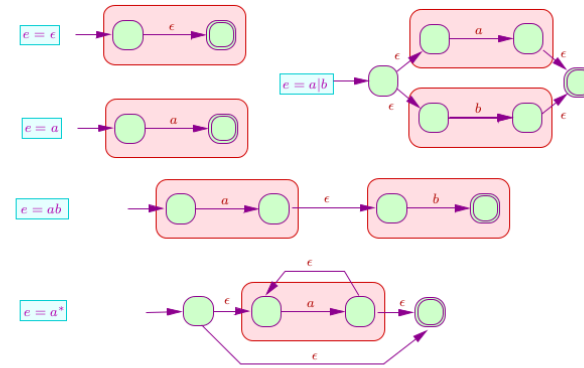
Title: Petter: Compilerbau (18.04.2016)

Date: Mon Apr 18 14:26:01 CEST 2016

Duration: 90:04 min

Pages: 35

## In Linear Time from Regular Expressions to NFAs



### Thompson's Algorithm

Produces  $\mathcal{O}(n)$  states for regular expressions of length  $n$ .



Ken Thompson

## Berry-Sethi Approach

### Glushkov Algorithm

Produces exactly  $n + 1$  states without  $\epsilon$ -transitions and demonstrates  $\rightarrow$  *Equality Systems* and  $\rightarrow$  *Attribute Grammars*



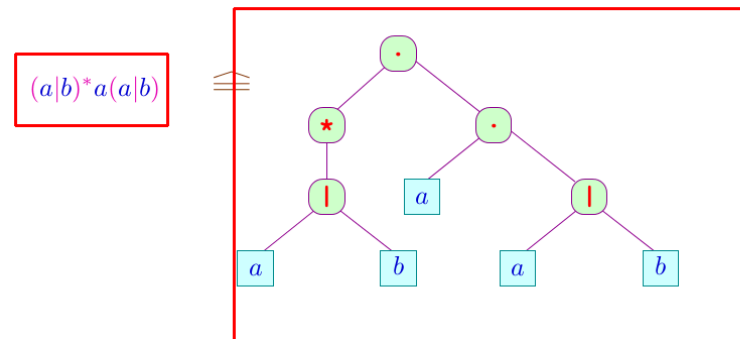
Viktor M. Glushkov

### Idea:

The automaton tracks (conceptionally via a marker " $\bullet$ "), in the syntax tree of a regular expression, which subexpressions in  $e$  are reachable consuming the rest of input  $w$ .

## Berry-Sethi Approach

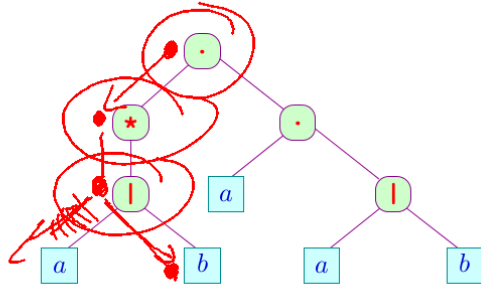
... for example:



## Berry-Sethi Approach

... for example:

$w = bbaa$  :

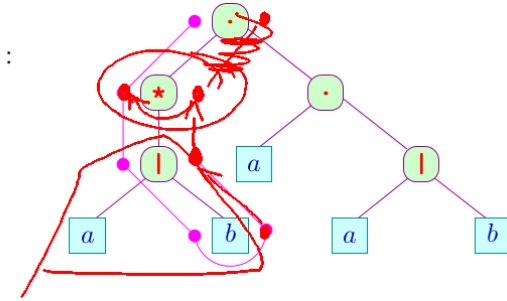


29 / 282

## Berry-Sethi Approach

... for example:

$w = bbaa$  :

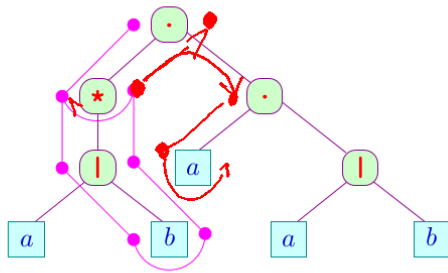


29 / 282

## Berry-Sethi Approach

... for example:

$w = a$  :

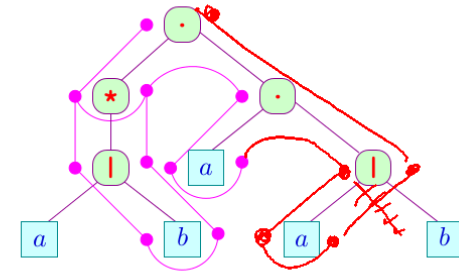


29 / 282

## Berry-Sethi Approach

... for example:

$w = a$  :

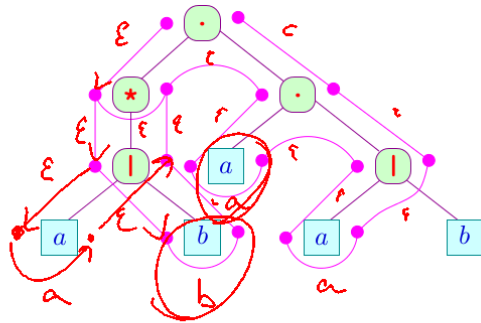


29 / 282

## Berry-Sethi Approach

... for example:

$w =$  :



29 / 282

## Berry-Sethi Approach

In general:

- Input is only consumed at the leaves.
- Navigating the tree does not consume input  $\rightarrow \epsilon$ -transitions
- For a formal construction we need **identifiers** for states.
- For a node  $n$ 's **identifier** we take the **subexpression**, corresponding to the subtree dominated by  $n$ .
- There are possibly **identical subexpressions** in one regular expression.

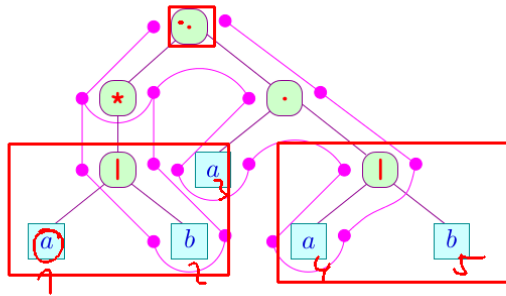
$\implies$  we enumerate the leaves ...

30 / 282

## Berry-Sethi Approach

... for example:

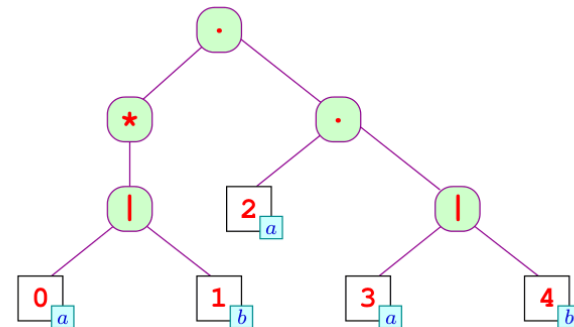
$w =$  :



29 / 282

## Berry-Sethi Approach

... for example:



31 / 282

## Berry-Sethi Approach (naive version)

### Construction (naive version):

States:  $\bullet r, r \bullet$  with  $r$  nodes of  $e$ ;

Start state:  $\bullet e$ ;

Final state:  $e \bullet$ ;

Transitions: for leaves  $r \equiv \boxed{i \mid x}$  we require:  $(\bullet r, x, r \bullet)$ .

The leftover transitions are:

$r$	Transitions
$r_1 \mid r_2$	$(\bullet r, \epsilon, \bullet r_1)$ $(\bullet r, \epsilon, \bullet r_2)$ $(r_1 \bullet, \epsilon, r \bullet)$ $(r_2 \bullet, \epsilon, r \bullet)$
$r_1 \cdot r_2$	$(\bullet r, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, \bullet r_2)$ $(r_2 \bullet, \epsilon, r \bullet)$



$r$	Transitions
$r_1^*$	$(\bullet r, \epsilon, r \bullet)$ $(\bullet r, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, r \bullet)$
$r_1?$	$(\bullet r, \epsilon, r \bullet)$ $(\bullet r, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, r \bullet)$

32 / 282

## Berry-Sethi Approach

### Discussion:

- Most transitions navigate through the expression
- The resulting automaton is in general **nondeterministic**

33 / 282

## Berry-Sethi Approach

### Discussion:

- Most transitions navigate through the expression
- The resulting automaton is in general **nondeterministic**

⇒ Strategy for the sophisticated version:  
Avoid generating  $\epsilon$ -transitions



33 / 282

## Berry-Sethi Approach

### Discussion:

- Most transitions navigate through the expression
- The resulting automaton is in general **nondeterministic**

⇒ Strategy for the sophisticated version:  
Avoid generating  $\epsilon$ -transitions

### Idea:

Pre-compute helper attributes during  $D(\text{epth})F(\text{irst})S(\text{earch})!$



33 / 282

## Berry-Sethi Approach

### Discussion:

- Most transitions navigate through the expression
- The resulting automaton is in general **nondeterministic**

⇒ Strategy for the sophisticated version:  
Avoid generating  $\epsilon$ -transitions

### Idea:

Pre-compute helper attributes during **D**(epth)**F**(irst)**S**(earch)!

### Necessary node-attributes:

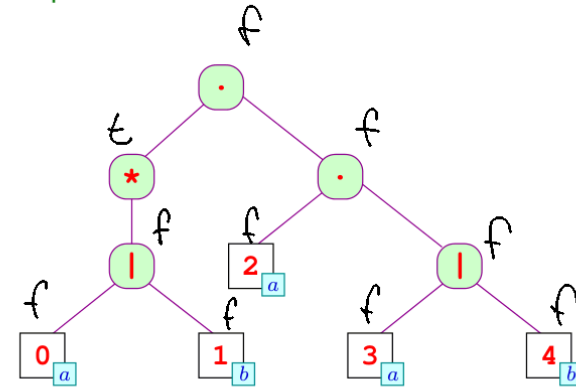
- first** the set of read states below  $r$ , which may be reached first, when descending into  $r$ .
- next** the set of read states to the right of  $r$ , which may be reached first in the traversal after  $r$ .
- last** the set of read states below  $r$ , which may be reached last when descending into  $r$ .
- empty** can the subexpression  $r$  consume  $\epsilon$ ?

33/282

## Berry-Sethi Approach: 1st step

$\text{empty}[r] = t$  if and only if  $\epsilon \in [r]$

... for example:



34/282

## Berry-Sethi Approach: 1st step

### Implementation:

DFS **post-order** traversal

for leaves  $r \equiv [i \ x]$  we find  $\text{empty}[r] = [x \equiv \epsilon]$ .

Otherwise:

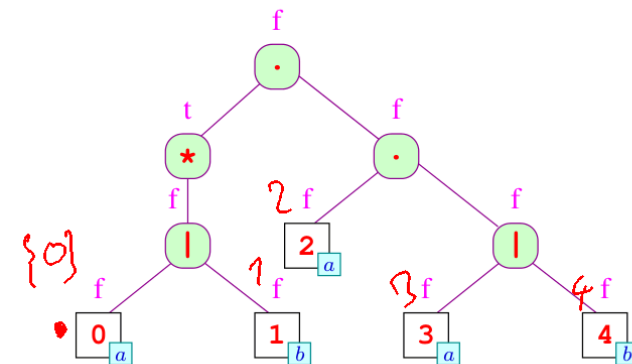
$$\begin{aligned} \text{empty}[r_1 \mid r_2] &= \text{empty}[r_1] \vee \text{empty}[r_2] \\ \text{empty}[r_1 \cdot r_2] &= \text{empty}[r_1] \wedge \text{empty}[r_2] \\ \text{empty}[r_1^*] &= t \\ \text{empty}[r_1^?] &= t \end{aligned}$$

35/282

## Berry-Sethi Approach: 2nd step

The **may-set** of first reached read states: The set of read states, that may be reached from  $\bullet r$  (i.e. while descending into  $r$ ) via sequences of  $\epsilon$ -transitions:  $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, [i \ x]) \in \delta^*, x \neq \epsilon\}$

... for example:

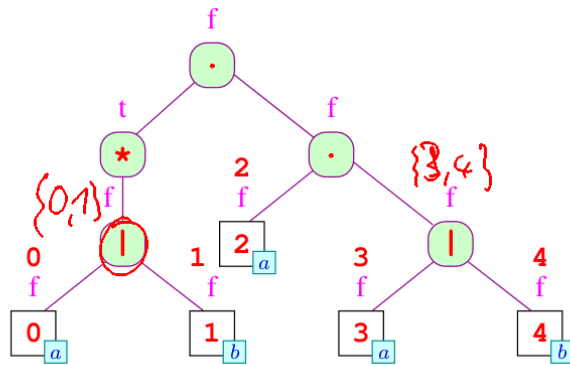


36/282

## Berry-Sethi Approach: 2nd step

The **may-set** of **first reached read states**: The set of read states, that may be reached from  $\bullet r$  (i.e. while descending into  $r$ ) via sequences of  $\epsilon$ -transitions:  $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i} \boxed{x}) \in \delta^*, x \neq \epsilon\}$

... for example:

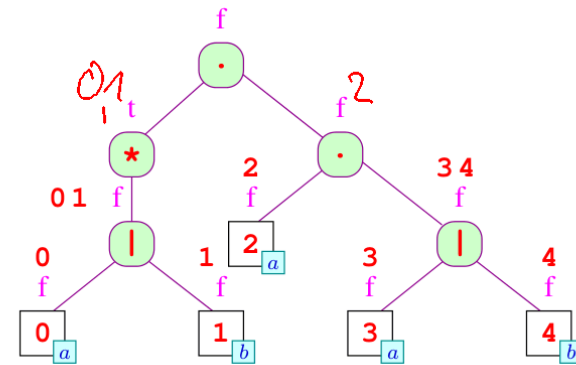


36 / 282

## Berry-Sethi Approach: 2nd step

The **may-set** of **first reached read states**: The set of read states, that may be reached from  $\bullet r$  (i.e. while descending into  $r$ ) via sequences of  $\epsilon$ -transitions:  $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i} \boxed{x}) \in \delta^*, x \neq \epsilon\}$

... for example:

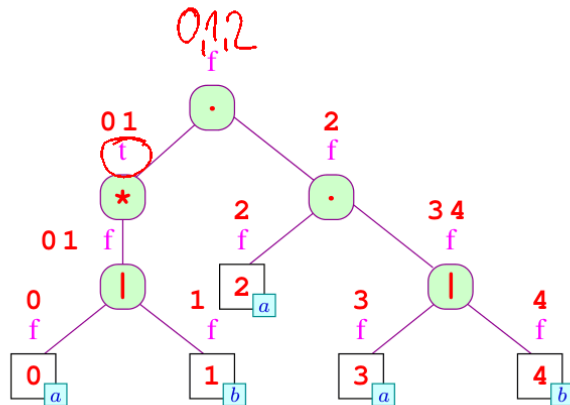


36 / 282

## Berry-Sethi Approach: 2nd step

The **may-set** of **first reached read states**: The set of read states, that may be reached from  $\bullet r$  (i.e. while descending into  $r$ ) via sequences of  $\epsilon$ -transitions:  $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i} \boxed{x}) \in \delta^*, x \neq \epsilon\}$

... for example:

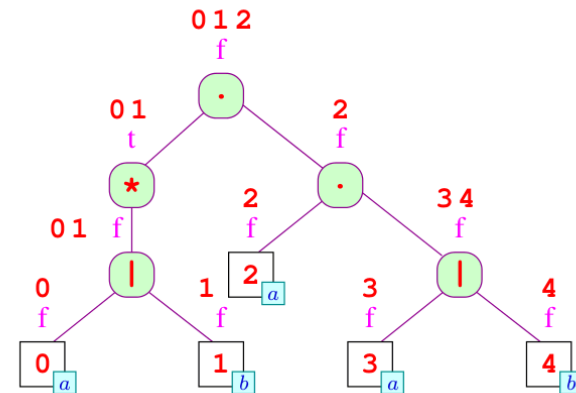


36 / 282

## Berry-Sethi Approach: 2nd step

The **may-set** of **first reached read states**: The set of read states, that may be reached from  $\bullet r$  (i.e. while descending into  $r$ ) via sequences of  $\epsilon$ -transitions:  $\text{first}[r] = \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet \boxed{i} \boxed{x}) \in \delta^*, x \neq \epsilon\}$

... for example:



36 / 282

## Berry-Sethi Approach: 2nd step

### Implementation:

DFS post-order traversal

for leaves  $r \equiv \boxed{i \ x}$  we find  $\text{first}[r] = \{i \mid x \neq \epsilon\}$ .

Otherwise:

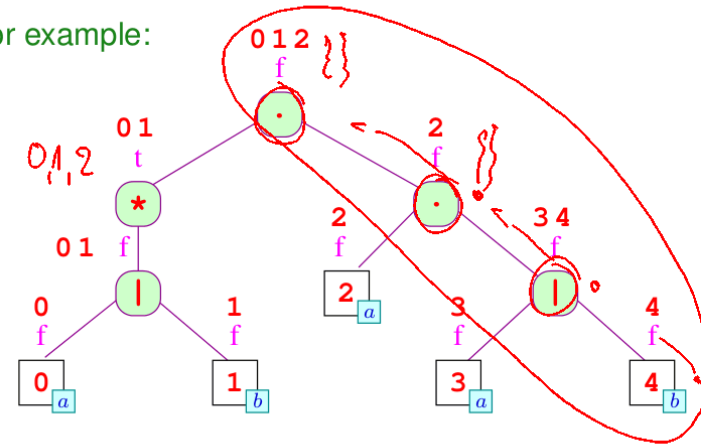
$$\begin{aligned} \text{first}[r_1 \mid r_2] &= \text{first}[r_1] \cup \text{first}[r_2] \\ \text{first}[r_1 \cdot r_2] &= \begin{cases} \text{first}[r_1] \cup \text{first}[r_2] & \text{if } \text{empty}[r_1] = t \\ \text{first}[r_1] & \text{if } \text{empty}[r_1] = f \end{cases} \\ \text{first}[r_1^*] &= \text{first}[r_1] \\ \text{first}[r_1^?] &= \text{first}[r_1] \end{aligned}$$

37/282

## Berry-Sethi Approach: 3rd step

The **may-set** of **next read states**: The set of read states within the subtrees right of  $r \bullet$ , that may be reached next via sequences of  $\epsilon$ -transitions.  $\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \ x}) \in \delta^*, x \neq \epsilon\}$

... for example:

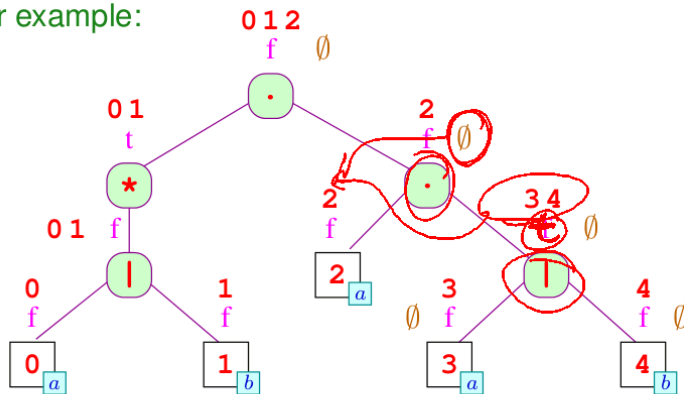


38/282

## Berry-Sethi Approach: 3rd step

The **may-set** of **next read states**: The set of read states within the subtrees right of  $r \bullet$ , that may be reached next via sequences of  $\epsilon$ -transitions.  $\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \ x}) \in \delta^*, x \neq \epsilon\}$

... for example:

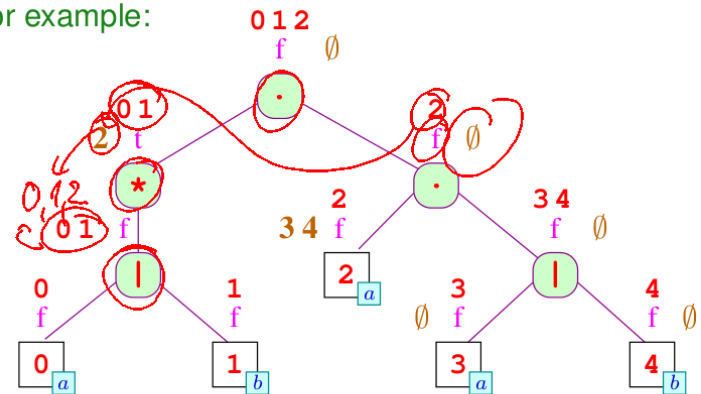


38/282

## Berry-Sethi Approach: 3rd step

The **may-set** of **next read states**: The set of read states within the subtrees right of  $r \bullet$ , that may be reached next via sequences of  $\epsilon$ -transitions.  $\text{next}[r] = \{i \mid (r \bullet, \epsilon, \bullet \boxed{i \ x}) \in \delta^*, x \neq \epsilon\}$

... for example:



38/282

## Berry-Sethi Approach: 3rd step

### Implementation:

DFS pre-order traversal

For the root, we find:  $next[e] = \emptyset$

Apart from that we distinguish, based on the context:

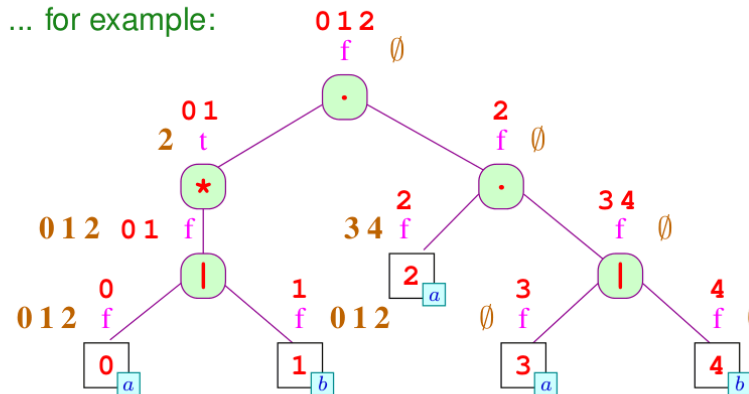
$r$	Equalities
$r_1 \mid r_2$	$next[r_1] = next[r]$ $next[r_2] = next[r]$
$r_1 \cdot r_2$	$next[r_1] = \begin{cases} first[r_2] \cup next[r] & \text{if } empty[r_2] = t \\ first[r_2] & \text{if } empty[r_2] = f \end{cases}$ $next[r_2] = next[r]$
$r_1^*$	$next[r_1] = first[r_1] \cup next[r]$
$r_1^?$	$next[r_1] = next[r]$

39/282

## Berry-Sethi Approach: 4th step

The may-set of last reached read states: The set of read states, which may be reached last during the traversal of  $r$  connected to the root via  $\epsilon$ -transitions only:  $last[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$

... for example:

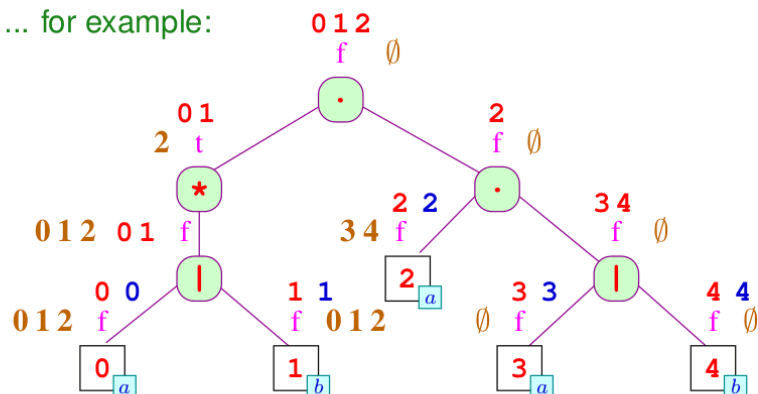


40/282

## Berry-Sethi Approach: 4th step

The may-set of last reached read states: The set of read states, which may be reached last during the traversal of  $r$  connected to the root via  $\epsilon$ -transitions only:  $last[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$

... for example:

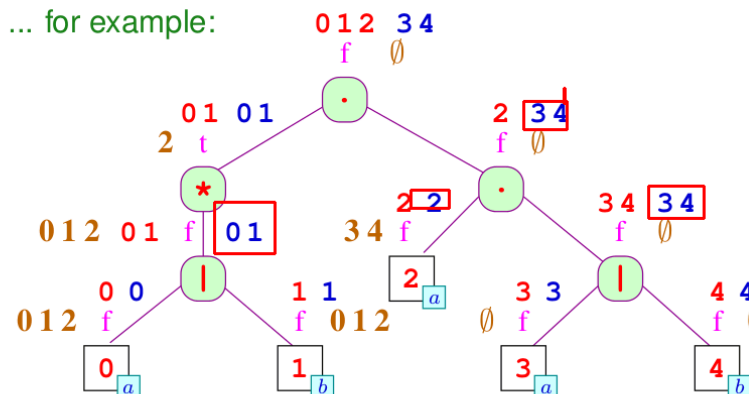


40/282

## Berry-Sethi Approach: 4th step

The may-set of last reached read states: The set of read states, which may be reached last during the traversal of  $r$  connected to the root via  $\epsilon$ -transitions only:  $last[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$

... for example:



40/282



## Berry-Sethi Approach: 4th step

### Implementation:

DFS **post-order** traversal

for leaves  $r \equiv \boxed{i \ x}$  we find  $\text{last}[r] = \{i \mid x \neq \epsilon\}$ .

Otherwise:

$$\begin{aligned} \text{last}[r_1 \mid r_2] &= \text{last}[r_1] \cup \text{last}[r_2] \\ \text{last}[r_1 \cdot r_2] &= \begin{cases} \text{last}[r_1] \cup \text{last}[r_2] & \text{if } \text{empty}[r_2] = t \\ \text{last}[r_2] & \text{if } \text{empty}[r_2] = f \end{cases} \\ \text{last}[r_1^*] &= \text{last}[r_1] \\ \text{last}[r_1^?] &= \text{last}[r_1] \end{aligned}$$

41/282

## Berry-Sethi Approach: (sophisticated version)

### Construction (sophisticated version):

Create an automaton based on the syntax tree's new attributes:

States:  $\{\bullet e\} \cup \{i \bullet \mid i \text{ a leaf}\}$

Start state:  $\bullet e$

Final states:  $\text{last}[e]$  if  $\text{empty}[e] = f$   
 $\{\bullet e\} \cup \text{last}[e]$  otherwise

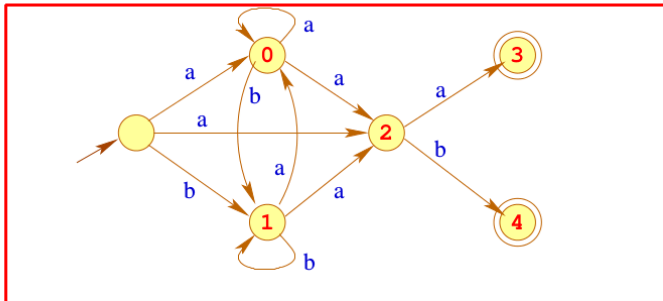
Transitions:  $(\bullet e, a, i \bullet)$  if  $i \in \text{first}[e]$  and  $i$  labeled with  $a$ .  
 $(i \bullet, a, i' \bullet)$  if  $i' \in \text{next}[i]$  and  $i'$  labeled with  $a$ .

We call the resulting automaton  $A_e$ .

42/282

## Berry-Sethi Approach

... for example:



### Remarks:

- This construction is known as **Berry-Sethi-** or **Glushkov-**construction.
- It is used for **XML** to define **Content Models**
- The result may not be, what we had in mind...

43/282

## Lexical Analysis

## Chapter 4: Turning NFAs deterministic

44/282