

Script generated by TTT

Title: Simon: Compilerbau (14.04.2014)

Date: Mon Apr 14 14:15:06 CEST 2014

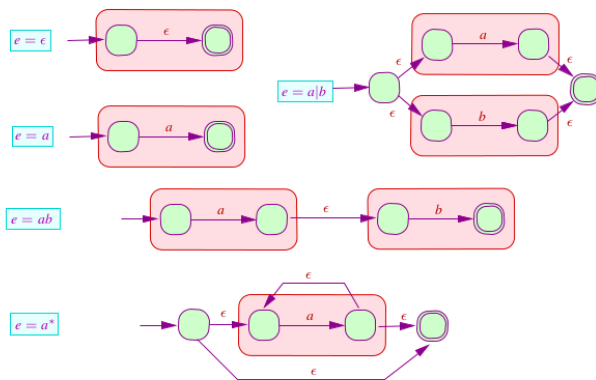
Duration: 94:06 min

Pages: 59

Chapter 3:

Converting Regular Expressions to NFAs

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



Berry-Sethi Algorithm

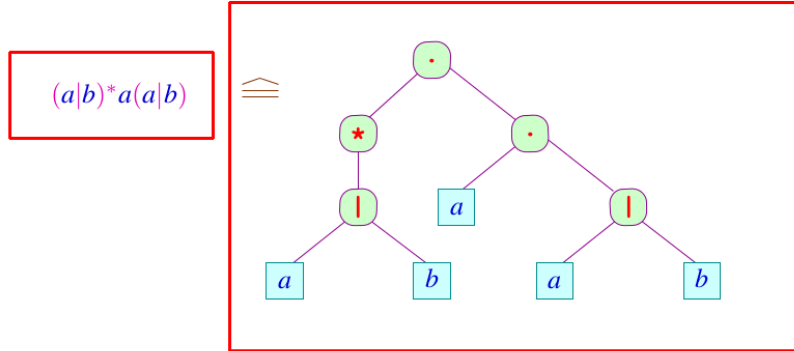
Produces exactly $n + 1$ states without ϵ -transitions and demonstrates \rightarrow *Equality Systems* and \rightarrow *Attribute Grammars*

Idea:

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

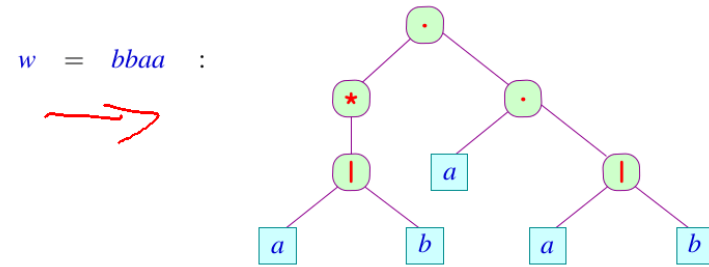
Berry-Sethi Approach

... for example:



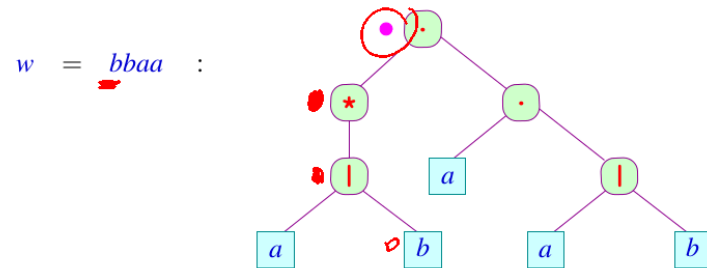
Berry-Sethi Approach

... for example:



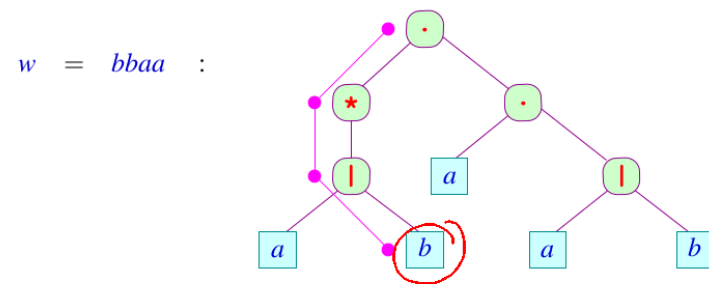
Berry-Sethi Approach

... for example:



Berry-Sethi Approach

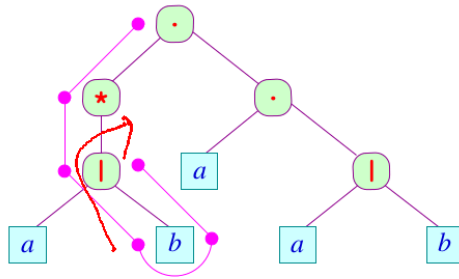
... for example:



Berry-Sethi Approach

... for example:

$$w = \underline{baa} :$$

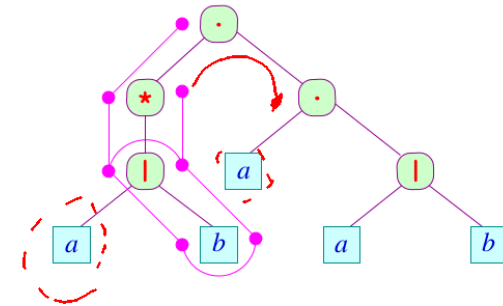


30 / 59

Berry-Sethi Approach

... for example:

$$w = \underline{aa} :$$

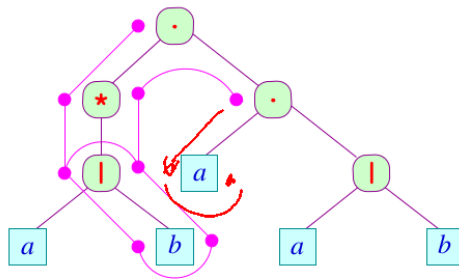


30 / 59

Berry-Sethi Approach

... for example:

$$w = aa :$$

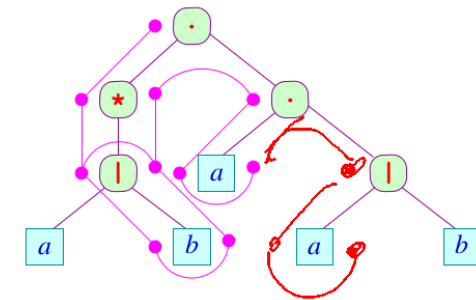


30 / 59

Berry-Sethi Approach

... for example:

$$w = \underline{a} :$$

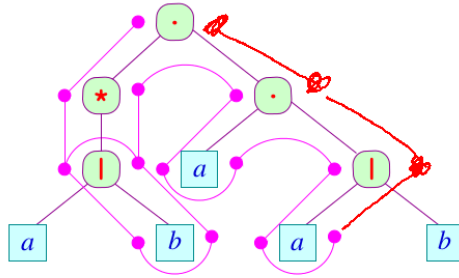


30 / 59

Berry-Sethi Approach

... for example:

$w =$:

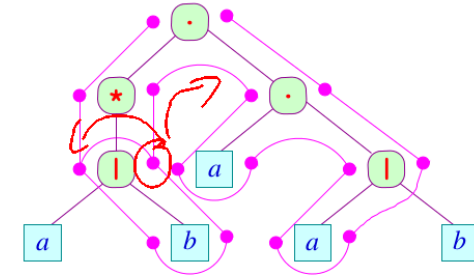


30 / 59

Berry-Sethi Approach

... for example:

$w =$:



30 / 59

Berry-Sethi Approach

In general:

- Input is only consumed by the leaves.
- Navigation in the tree is done without consuming input, i.e. via ϵ -transition.
- For a formal construction we need **identifiers** for states.
- Therefore we use the **subexpression**, corresponding to the subtree, dominated by the particular node.
- There are possibly identical subexpressions in one regular expression.

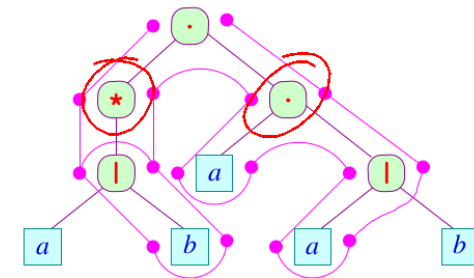
\implies we enumerate the leaves ...

31 / 59

Berry-Sethi Approach

... for example:

$w =$:



30 / 59

Berry-Sethi Approach

In general:

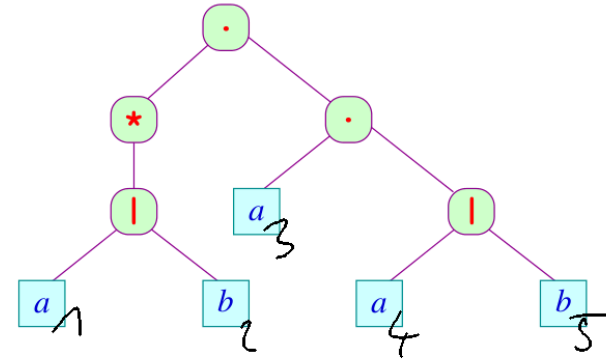
- Input is only consumed by the leaves.
- Navigation in the tree is done without consuming input, i.e. via ϵ -transition.
- For a formal construction we need **identifiers** for states.
- Therefore we use the **subexpression**, corresponding to the subtree, dominated by the particular node.
- There are possibly identical subexpressions in one regular expression.

\implies we enumerate the leaves ...

31 / 59

Berry-Sethi Approach

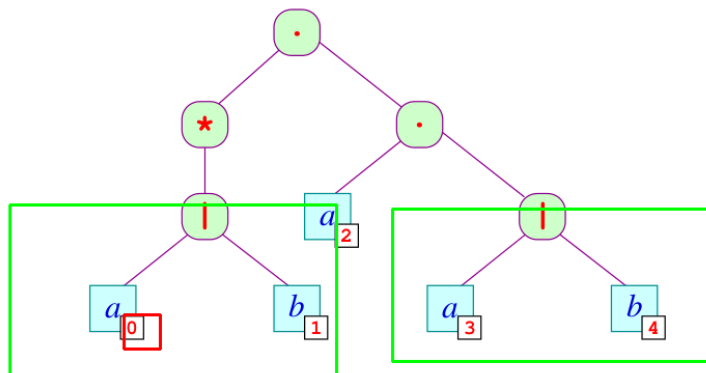
... for example:



32 / 59

Berry-Sethi Approach

... for example:



32 / 59

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 $r \equiv \boxed{i} \boxed{x} \boxed{j}$ we require: $\boxed{\bullet r} \boxed{x} \boxed{r\bullet}$.

The leftover transitions are:

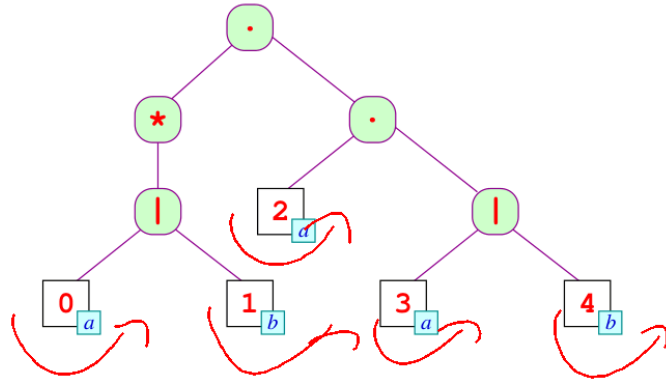
r	Transitions
$r_1 \mid r_2$	$(\bullet r, \epsilon, \bullet r_1)$ $(\bullet r, \epsilon, \bullet r_2)$ $(r_1\bullet, \epsilon, r\bullet)$ $(r_2\bullet, \epsilon, r\bullet)$
$r_1 \cdot r_2$	$(\bullet r, \epsilon, \bullet r_1)$ $(r_1\bullet, \epsilon, \bullet r_2)$ $(r_2\bullet, \epsilon, r\bullet)$

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

33 / 59

Berry-Sethi Approach

... for example:



32 / 59

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

33 / 59

Berry-Sethi Approach

Discussion:

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

34 / 59

Berry-Sethi Approach

Discussion:

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

⇒ Strategy for the **sophisticated version:**

Avoid generating ϵ -transitions

34 / 59

Berry-Sethi Approach

Discussion:

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

⇒ Strategy for the sophisticated version:

Avoid generating ϵ -transitions

Necessary node-attributes:

- empty** can the subexpression r consume ϵ ?
- first** the set of read states below r , which **may** be reached **first**, when descending into r .
- next** the set of read states on 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 .

34 / 59

Berry-Sethi Approach

Discussion:

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

⇒ Strategy for the sophisticated version:

Avoid generating ϵ -transitions

Necessary node-attributes:

- empty** can the subexpression r consume ϵ ?
- first** the set of read states below r , which **may** be reached **first**, when descending into r .
- next** the set of read states on 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 .

Idea:

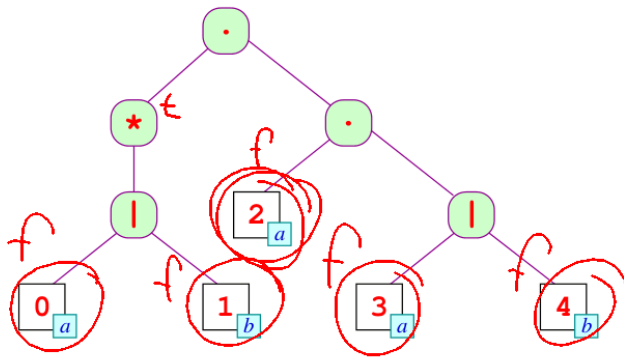
Pre-compute the attributes during **D(ept)H(irst)S(earch)!**

34 / 59

Berry-Sethi Approach: 1st step

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

... for example:

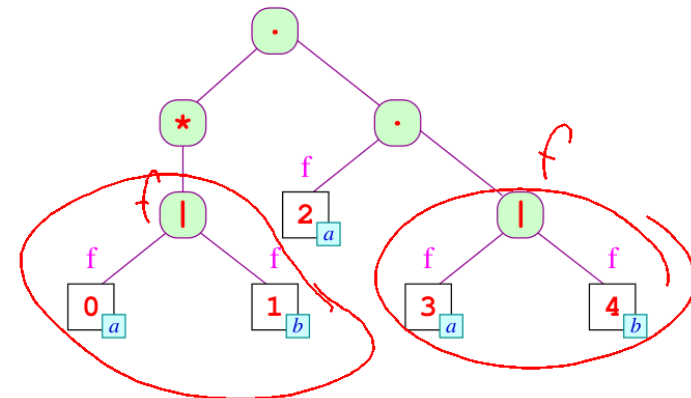


35 / 59

Berry-Sethi Approach: 1st step

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

... for example:

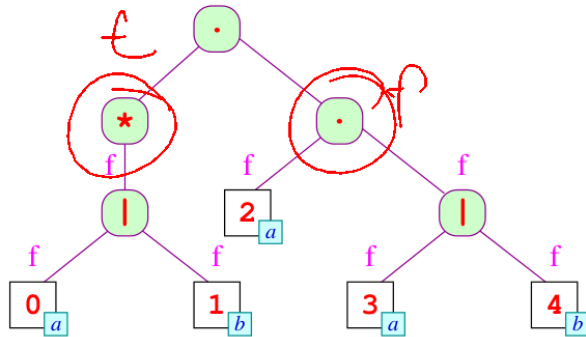


35 / 59

Berry-Sethi Approach: 1st step

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

... for example:

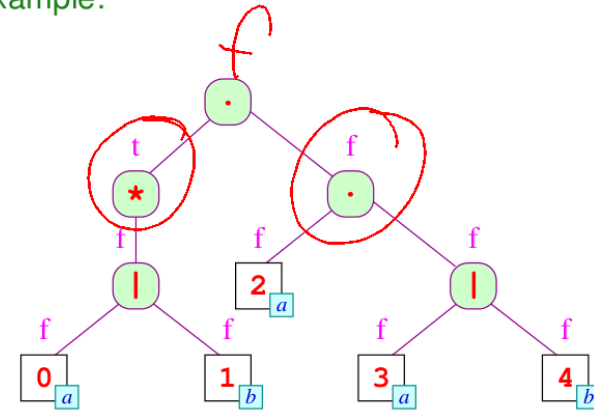


35 / 59

Berry-Sethi Approach: 1st step

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

... for example:

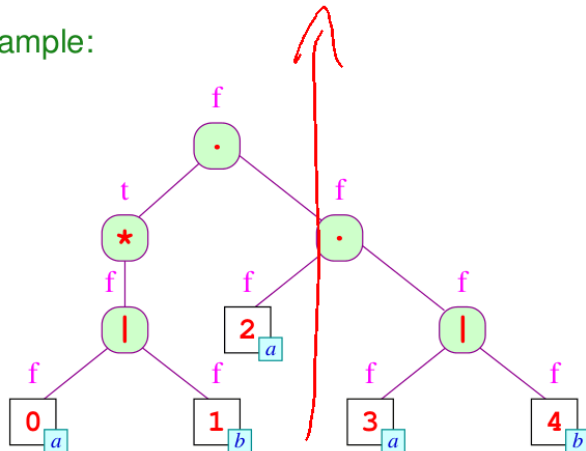


35 / 59

Berry-Sethi Approach: 1st step

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

... for example:



35 / 59

Berry-Sethi Approach: 2nd step

Implementation:

DFS **post-order** traversal

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

Otherwise:

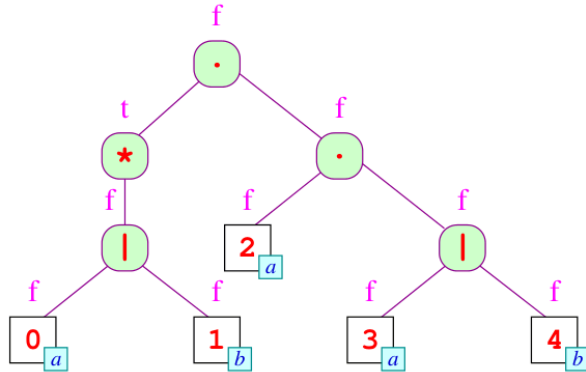
$$\begin{aligned} \text{empty}[r_1 \mid r_2] &= \text{empty}[r_1] \vee \text{empty}[r_2] \\ \text{empty}[r_1 \cdot r_2] &= \text{empty}[r_1] \wedge \text{empty}[r_2] \\ \text{empty}[r_1^*] &= t \\ \text{empty}[r_1?] &= t \end{aligned}$$

36 / 59

Berry-Sethi Approach: 1st step

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

... for example:

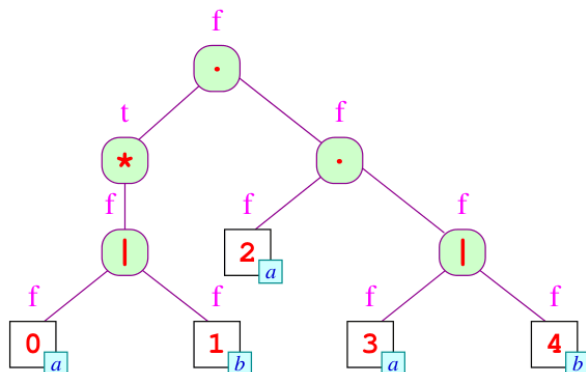


35 / 59

Berry-Sethi Approach: 2nd step

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

... for example:



37 / 59

Berry-Sethi Approach: 2nd step

Implementation:

DFS **post-order** traversal

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

Otherwise:

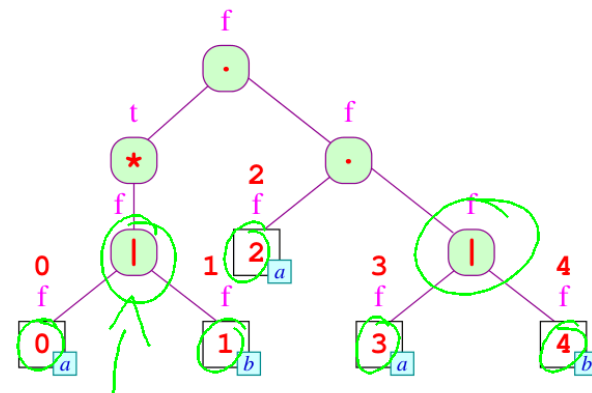
$$\begin{aligned} \text{empty}[r_1 \mid r_2] &= \text{empty}[r_1] \vee \text{empty}[r_2] \\ \text{empty}[r_1 \cdot r_2] &= \text{empty}[r_1] \wedge \text{empty}[r_2] \\ \text{empty}[r_i^*] &= t \\ \text{empty}[r_i?] &= t \end{aligned}$$

36 / 59

Berry-Sethi Approach: 2nd step

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

... for example:

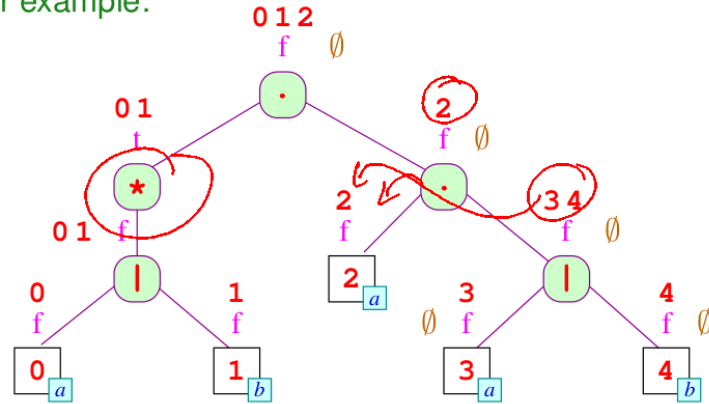


37 / 59

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 ϵ -transitions. $\text{next}[r] = \{i \mid (r\bullet, \epsilon, \bullet \boxed{i} x) \in \delta^*, x \neq \epsilon\}$

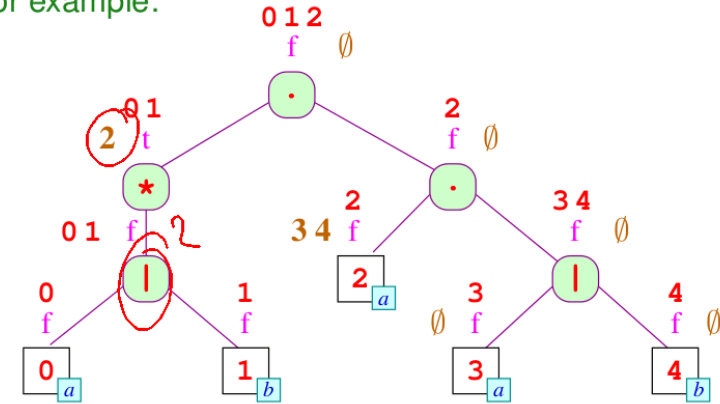
... for example:



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 ϵ -transitions. $\text{next}[r] = \{i \mid (r\bullet, \epsilon, \bullet \boxed{i} x) \in \delta^*, x \neq \epsilon\}$

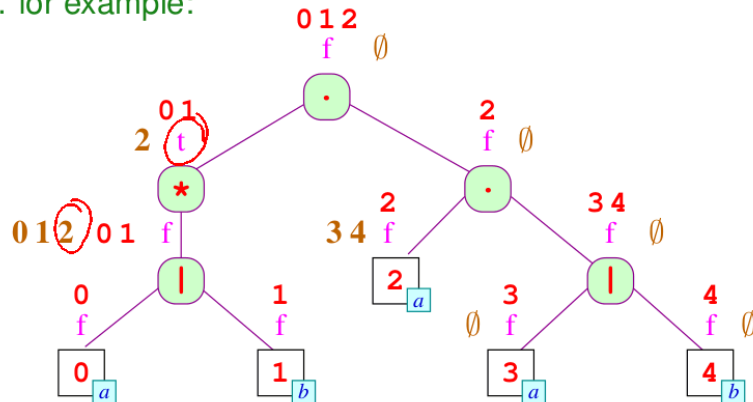
... for example:



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 ϵ -transitions. $\text{next}[r] = \{i \mid (r\bullet, \epsilon, \bullet \boxed{i} x) \in \delta^*, x \neq \epsilon\}$

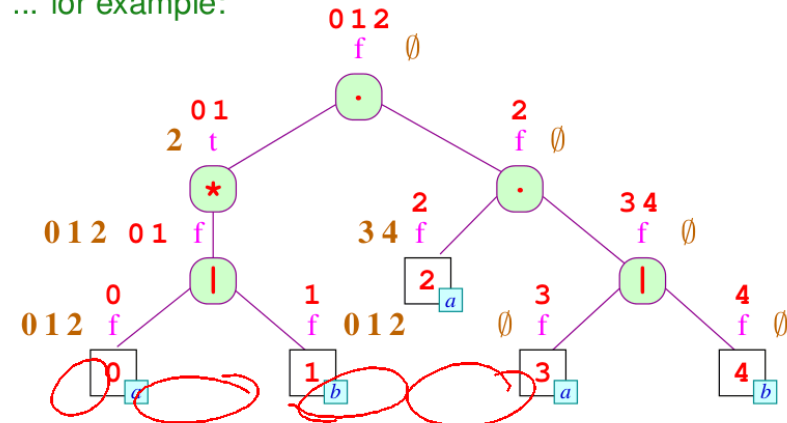
... for example:



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 ϵ -transitions. $\text{next}[r] = \{i \mid (r\bullet, \epsilon, \bullet \boxed{i} x) \in \delta^*, x \neq \epsilon\}$

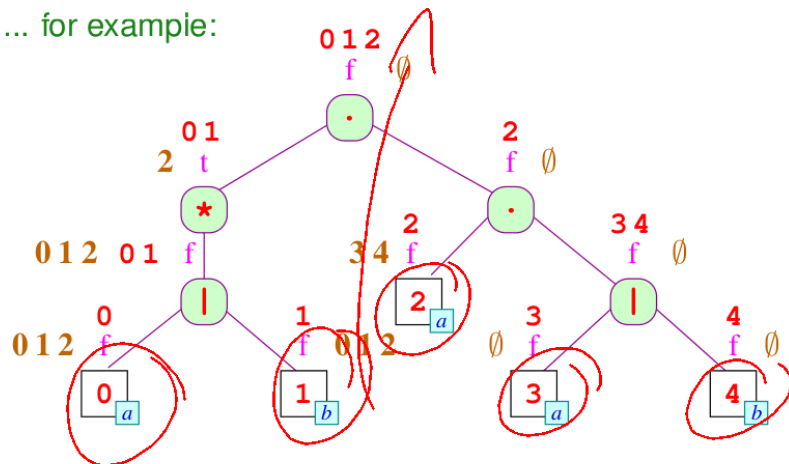
... for example:



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 ϵ -transitions only: $\text{last}[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$

... for example:

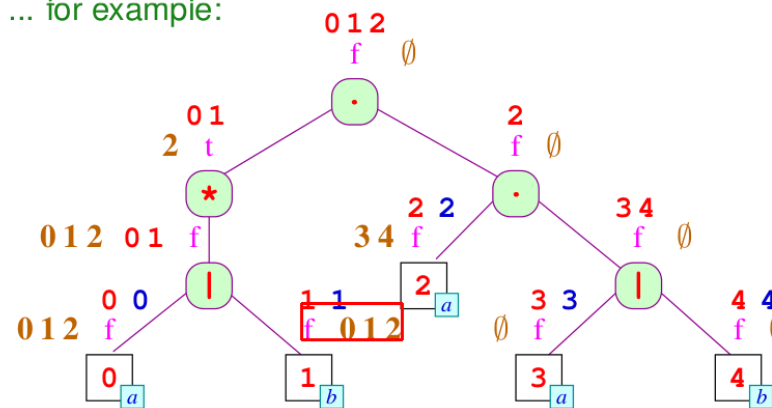


41 / 59

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 ϵ -transitions only: $\text{last}[r] = \{i \text{ in } r \mid (\boxed{i \ x} \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$

... for example:



41 / 59

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: $\bullet e$

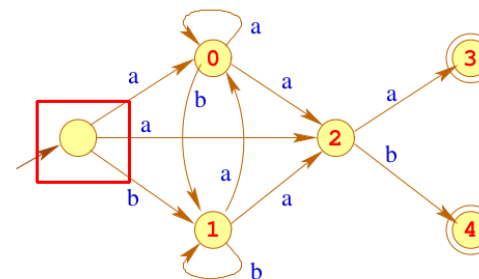
Final states: $\text{last}[e]$ if $\text{empty}[e] = f$
 $\{\bullet e\} \cup \text{last}[e]$ otherwise

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 .

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 / 59

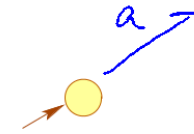
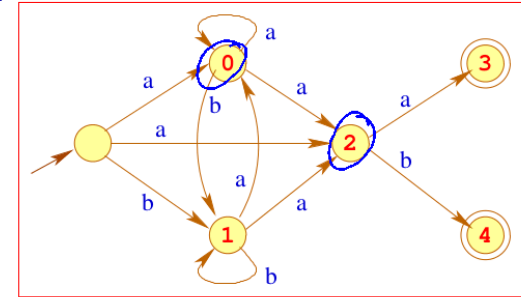
44 / 59

Chapter 4: Turning NFAs deterministic



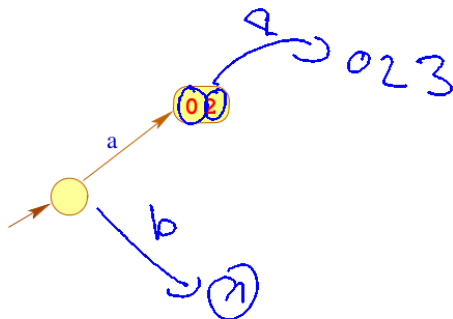
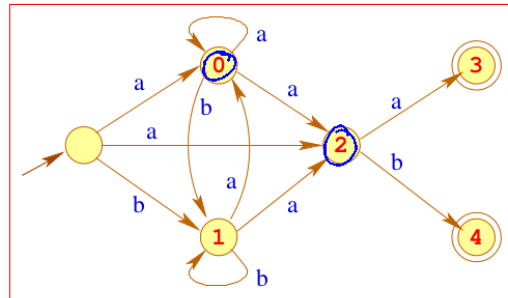
Powerset Construction

... for example:



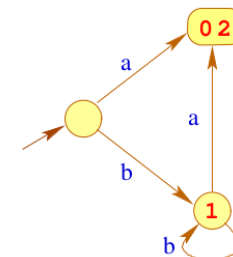
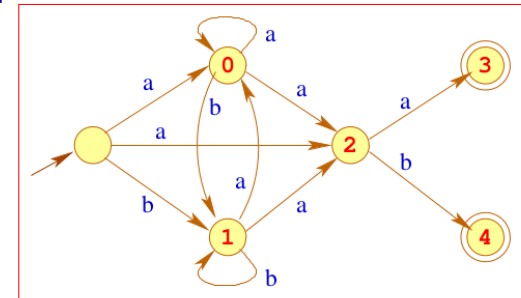
Powerset Construction

... for example:



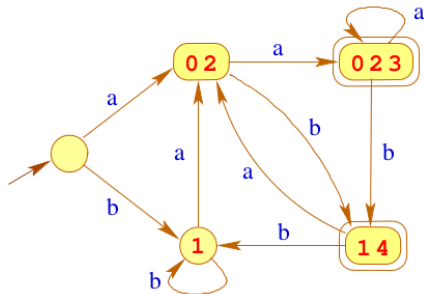
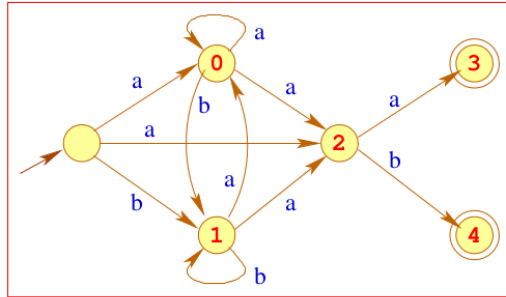
Powerset Construction

... for example:



Powerset Construction

... for example:



47 / 59

Powerset Construction

Bummer!

There are exponentially many powersets of Q

- **Idea:** Consider only **contributing** powersets. Starting with the set $Q_{\mathcal{P}} = \{I\}$ we only add further states **by need** ...
- i.e., whenever we can reach them from a state in $Q_{\mathcal{P}}$
- Even though, the resulting automaton can become enormously **huge** ... which is (sort of) not happening in **practice**

49 / 59

Powerset Construction

Theorem:

For every non-deterministic automaton $A = (Q, \Sigma, \delta, I, F)$ we can compute a deterministic automaton $\mathcal{P}(A)$ with

$$\mathcal{L}(A) = \mathcal{L}(\mathcal{P}(A))$$

Construction:

States: Powersets of Q ;

Start state: I ;

Final states: $\{Q' \subseteq Q \mid Q' \cap F \neq \emptyset\}$;

Transitions: $\delta_{\mathcal{P}}(Q', a) = \{q \in Q \mid \exists p \in Q' : (p, a, q) \in \delta\}$

48 / 59

Powerset Construction

Bummer!

There are exponentially many powersets of Q

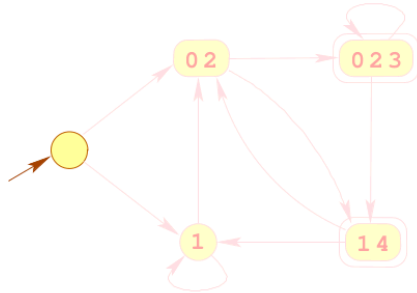
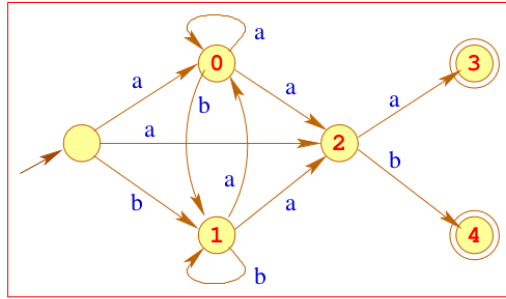
- **Idea:** Consider only **contributing** powersets. Starting with the set $Q_{\mathcal{P}} = \{I\}$ we only add further states **by need** ...
- i.e., whenever we can reach them from a state in $Q_{\mathcal{P}}$
- Even though, the resulting automaton can become enormously **huge** ... which is (sort of) not happening in **practice**
- Therefore, in tools like **grep** a regular expression's **DFA** is never created!
- Instead, only the sets, directly necessary for interpreting the input are generated **while processing the input**

49 / 59

Powerset Construction

... for example:

a b a b

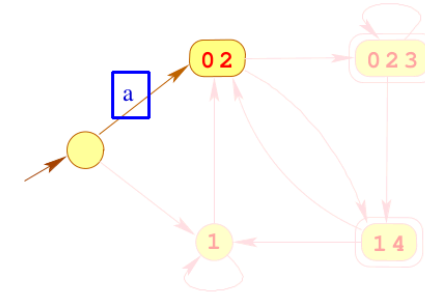
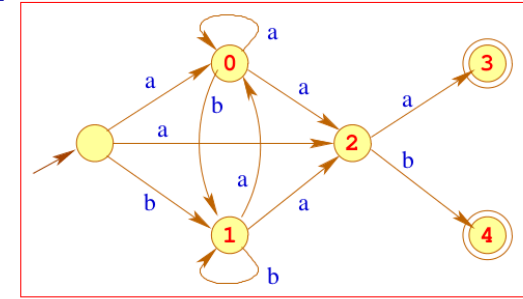


50 / 59

Powerset Construction

... for example:

a b a b

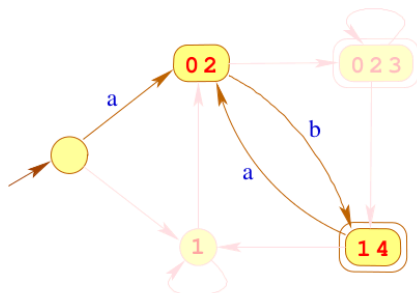
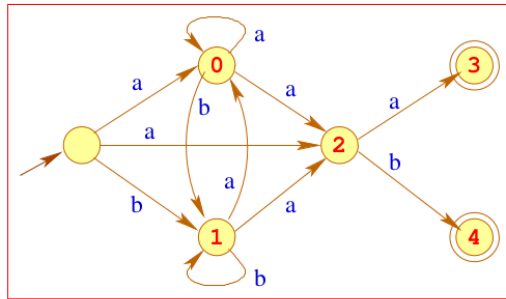


50 / 59

Powerset Construction

... for example:

a b a b



50 / 59

Remarks:

- For an input sequence of length n , maximally $\mathcal{O}(n)$ sets are generated
- Once a set/edge of the DFA is generated, they are stored within a hash-table.
- Before generating a new transition, we check this table for already existing edges with the desired label.

51 / 59