

# From KAT expressions to automata

Marcello Bonsangue<sup>1,2</sup>   Jan Rutten<sup>1,3</sup>   Alexandra Silva<sup>1</sup>

<sup>1</sup>Centrum voor Wiskunde en Informatica

<sup>2</sup>LIACS - Leiden University

<sup>3</sup>Vrije Universiteit Amsterdam

Arco Meeting, May 2009

# Motivation

- Efficient algorithms exist to compile regular expressions into automata (Berry-Sethi, Thompson, ...)
- Alternatives to the elegant Brzozowski algorithm
- Kozen recently extended Brzozowski algorithm to KAT

Can we extend the efficient algorithms to KAT as well?

# Motivation

- Efficient algorithms exist to compile regular expressions into automata (Berry-Sethi, Thompson, ...)
- Alternatives to the elegant Brzozowski algorithm
- Kozen recently extended Brzozowski algorithm to KAT

Can we extend the efficient algorithms to KAT as well?

# Recall : Berry-Sethi for RE

Basic idea: Assume that all occurrences of letters are different.

$$(ab + b)^*ba \rightarrow (a_1b_2 + b_3)^*b_4a_5$$

## Definition

Let  $E$  be a regular expression and  $\bar{E}$  the corresponding marked expression.

$$\begin{aligned} \text{first}(E) &= \{i \mid a_i w \in L(\bar{E})\} \\ \text{follow}(E, i) &= \{j \mid ua_i a_j v \in L(\bar{E})\} \\ \text{last}(E) &= \{i \mid wa_i \in L(\bar{E})\} \end{aligned}$$

The position automaton for  $E$  is defined as :

$$\mathcal{A}_{\text{pos}}(E) = (\text{pos}(E), \Sigma, t_{\text{pos}}, 0, \text{last}_0(E))$$

where

$$\text{last}_0(E) = \begin{cases} \text{last}(E) \cup \{0\} & \delta(E) = 1 \\ \text{last}(E) & \text{oth.} \end{cases} \quad t_{\text{pos}} = \{(i, a, j) \mid j \in \text{follow}(E, i), a = \bar{a}_j\} \cup \{(0, a, j) \mid j \in \text{first}(E), a = \bar{a}_j\}$$

# Recall : Berry-Sethi for RE

Basic idea: Assume that all occurrences of letters are different.

$$(ab + b)^*ba \rightarrow (a_1b_2 + b_3)^*b_4a_5$$

## Definition

Let  $E$  be a regular expression and  $\bar{E}$  the corresponding marked expression.

$$\begin{aligned} \text{first}(E) &= \{i \mid a_i w \in L(\bar{E})\} \\ \text{follow}(E, i) &= \{j \mid ua_i a_j v \in L(\bar{E})\} \\ \text{last}(E) &= \{i \mid wa_i \in L(\bar{E})\} \end{aligned}$$

The position automaton for  $E$  is defined as :

$$\mathcal{A}_{\text{pos}}(E) = (\text{pos}(E), \Sigma, t_{\text{pos}}, 0, \text{last}_0(E))$$

where

$$\text{last}_0(E) = \begin{cases} \text{last}(E) \cup \{0\} & \delta(E) = 1 \\ \text{last}(E) & \text{oth.} \end{cases} \quad t_{\text{pos}} = \{(i, a, j) \mid j \in \text{follow}(E, i), a = \bar{a}_j\} \cup \{(0, a, j) \mid j \in \text{first}(E), a = \bar{a}_j\}$$

# Algorithm

Berry and Sethi provided a algorithmic description to compute the follow sets.

## Proposition

*For a regular expression  $E$ ,  $F(E, \{!\})$ , as defined below, yields pairs of the form  $\langle i, \text{follow}(E!, i) \rangle$ .*

$$\begin{aligned} F(E_1 + E_2, S) &= F(E_1, S) \cup F(E_2, S) \\ F(E_1 E_2, S) &= F(E_1, \text{first}(E_2) \cup \delta(E_2) \cdot S) \cup F(E_2, S) \\ F(E_1^*, S) &= F(E_1, \text{first}(E_1) \cup S) \\ F(a_i, S) &= \{\langle i, S \rangle\} \quad F(1, S) = \emptyset \quad F(0, S) = \emptyset \end{aligned}$$

*! is just a trick to avoid computing last:*

*$i \in \text{last}(E) \iff ! \in \text{follow}(E!, i)$ .*

Example in the board –  $(a_1 b_2 + b_3)^* b_4 a_5$ .

# Algorithm

Berry and Sethi provided a algorithmic description to compute the follow sets.

## Proposition

*For a regular expression  $E$ ,  $F(E, \{!\})$ , as defined below, yields pairs of the form  $\langle i, \text{follow}(E!, i) \rangle$ .*

$$\begin{aligned}F(E_1 + E_2, S) &= F(E_1, S) \cup F(E_2, S) \\F(E_1 E_2, S) &= F(E_1, \text{first}(E_2) \cup \delta(E_2) \cdot S) \cup F(E_2, S) \\F(E_1^*, S) &= F(E_1, \text{first}(E_1) \cup S) \\F(a_i, S) &= \{\langle i, S \rangle\} \quad F(1, S) = \emptyset \quad F(0, S) = \emptyset\end{aligned}$$

*! is just a trick to avoid computing last:*

$$i \in \text{last}(E) \iff ! \in \text{follow}(E!, i).$$

Example in the board –  $(a_1 b_2 + b_3)^* b_4 a_5$ .

# Algorithm

Berry and Sethi provided a algorithmic description to compute the follow sets.

## Proposition

*For a regular expression  $E$ ,  $F(E, \{!\})$ , as defined below, yields pairs of the form  $\langle i, \text{follow}(E!, i) \rangle$ .*

$$\begin{aligned} F(E_1 + E_2, S) &= F(E_1, S) \cup F(E_2, S) \\ F(E_1 E_2, S) &= F(E_1, \text{first}(E_2) \cup \delta(E_2) \cdot S) \cup F(E_2, S) \\ F(E_1^*, S) &= F(E_1, \text{first}(E_1) \cup S) \\ F(a_i, S) &= \{\langle i, S \rangle\} \quad F(1, S) = \emptyset \quad F(0, S) = \emptyset \end{aligned}$$

*! is just a trick to avoid computing last:*

$$i \in \text{last}(E) \iff ! \in \text{follow}(E!, i).$$

Example in the board –  $(a_1 b_2 + b_3)^* b_4 a_5$ .



$$\begin{aligned} \text{Exp} \ni e, f &::= p \in \Sigma \mid b \in \text{BExp} \mid e.f \mid e + f \mid e^* \\ \text{BExp} \ni b &::= t \in T \mid b_1.b_2 \mid b_1 + b_2 \mid \overline{b} \mid 0 \mid 1 \end{aligned}$$

$\mathcal{B} = \text{Bexp} / \equiv, \equiv$  logical equivalence, is a boolean algebra.

$$\begin{aligned} \text{Exp} \ni e, f &::= p \in \Sigma \mid b \in \text{BExp} \mid e.f \mid e + f \mid e^* \\ \text{BExp} \ni b &::= t \in T \mid b_1.b_2 \mid b_1 + b_2 \mid \overline{b} \mid 0 \mid 1 \end{aligned}$$

$\mathcal{B} = \text{Bexp} / \equiv, \equiv$  logical equivalence, is a boolean algebra.

KAT are regular expressions over  $\Sigma \cup T$

**Observation 1:** We can apply the Berry-Sethi algorithm directly and get:

$$\begin{aligned} \text{Exp} \ni e, f &::= p \in \Sigma \mid b \in \text{BExp} \mid e.f \mid e + f \mid e^* \\ \text{BExp} \ni b &::= t \in T \mid b_1.b_2 \mid b_1 + b_2 \mid \overline{b} \mid 0 \mid 1 \end{aligned}$$

$\mathcal{B} = \text{BExp} / \equiv, \equiv$  logical equivalence, is a boolean algebra.

KAT are regular expressions over  $\Sigma \cup T$

**Observation 1:** We can apply the Berry-Sethi algorithm directly and get:

- A non-deterministic automaton

$$\mathcal{A}_e = (S, \delta : S \rightarrow 2 \times (\mathcal{P}S)^{\Sigma \cup T}), S \subseteq \text{Exp}$$

- $L(\mathcal{A}_e) = L(e)$

$$\begin{aligned} \text{Exp} \ni e, f &::= p \in \Sigma \mid b \in \text{BExp} \mid e.f \mid e + f \mid e^* \\ \text{BExp} \ni b &::= t \in T \mid b_1.b_2 \mid b_1 + b_2 \mid \bar{b} \mid 0 \mid 1 \end{aligned}$$

$\mathcal{B} = \text{Bexp}/\equiv, \equiv$  logical equivalence, is a boolean algebra.

KAT are regular expressions over  $\Sigma \cup T$

**Observation 1:** We can apply the Berry-Sethi algorithm directly and get:

- A non-deterministic automaton

$$\mathcal{A}_e = (S, \delta : S \rightarrow 2 \times (\mathcal{P}S)^{\Sigma \cup T}), S \subseteq \text{Exp}$$

- $L(\mathcal{A}_e) = L(e)$
- $GS(\mathcal{A}) = GS(e)$  (due to Kozen). **Good!**

# Examples

- $(b + p)qr^*$
- $b_1 + pb_2(b_3 + q)^*rb_4$

# Can we do “better”?

- Less states

$S \rightarrow 2 \times (\mathcal{P}S)^{\Sigma \cup T}$     Original automata on GS, too many states

$S \rightarrow \mathcal{B} \times S^{\Sigma \times At}$     Deterministic version, less intuitive

# Can we do “better”?

- Less states

$S \rightarrow 2 \times (\mathcal{P}S)^{\Sigma \cup T}$     Original automata on GS, too many states

$S \rightarrow \mathcal{B} \times (\mathcal{P}S)^{\Sigma \times \mathcal{B}}$     Compromise

$S \rightarrow \mathcal{B} \times S^{\Sigma \times At}$     Deterministic version, less intuitive

# Can we do “better”?

- Less states

$S \rightarrow 2 \times (\mathcal{P}S)^{\Sigma \cup T}$     Original automata on GS, too many states

$S \rightarrow \mathcal{B} \times (\mathcal{P}S)^{\Sigma \times \mathcal{B}}$     Compromise

$S \rightarrow \mathcal{B} \times S^{\Sigma \times At}$     Deterministic version, less intuitive

$At$ : minimal elements of  $\mathcal{B}$  (truth assignments or valuations,  
 $At \cong 2^T$ ).



# Can we do “better”?

- Less states

$S \rightarrow 2 \times (\mathcal{P}S)^{\Sigma \cup T}$  Original automata on GS, too many states

$S \rightarrow \mathcal{B} \times (\mathcal{P}S)^{\Sigma \times \mathcal{B}}$  Compromise

$S \rightarrow \mathcal{B} \times S^{\Sigma \times At}$  Deterministic version, less intuitive

$At$ : minimal elements of  $\mathcal{B}$  (truth assignments or valuations,  $At \cong 2^T$ ).

- Example in the board.

# Adapted Berry-Sethi for KAT

Basic idea: Assume that all occurrences of letters in  $\Sigma$  are different.

$$(b + p)qr^* \rightarrow (b + p_1)q_2r_3^*$$

## Definition

Let  $e$  be a KAT expression and  $\bar{e}$  the corresponding marked expression.

$$\begin{aligned} \text{first}(e) &= \{ \langle b, i \rangle \mid b_1 b_2 \dots b_n p x \in L(\bar{e}), b = \bigvee (b_1 \wedge b_2 \wedge \dots b_n) \} \\ \text{follow}(e, i) &= \{ \langle b, j \rangle \mid x p_i b_1 b_2 \dots b_n q_j y \in L(\bar{e}), b = \bigvee (b_1 \wedge b_2 \wedge \dots b_n) \} \\ \text{last}(e) &= \{ \langle b, i \rangle \mid x p_i b_1 b_2 \dots b_n \in L(\bar{e}), b = \bigvee (b_1 \wedge b_2 \wedge \dots b_n) \} \end{aligned}$$

The automaton corresponding to  $e$  is defined as :

$$\mathcal{A}_e = (\text{pos}(e), \Sigma, \langle o, \delta \rangle, 0)$$

where

$$o(i) = \begin{cases} E(e) & i = 0 \\ b & \langle b, i \rangle \in \text{last}(e) \\ \perp & \text{otherwise} \end{cases} \quad \delta = \{ \langle i, \langle b, p \rangle, j \rangle \mid \langle b, j \rangle \in \text{follow}(e, i), p = \bar{p}_i \} \\ \cup \{ \langle 0, \langle b, a \rangle, j \rangle \mid \langle b, j \rangle \in \text{first}(e), a = \bar{a}_j \}$$

# Adapted Berry-Sethi for KAT

Basic idea: Assume that all occurrences of letters in  $\Sigma$  are different.

$$(b + p)qr^* \rightarrow (b + p_1)q_2r_3^*$$

## Definition

Let  $e$  be a KAT expression and  $\bar{e}$  the corresponding marked expression.

$$\begin{aligned} \text{first}(e) &= \{ \langle b, i \rangle \mid b_1 b_2 \dots b_n p x \in L(\bar{e}), b = \bigvee (b_1 \wedge b_2 \wedge \dots b_n) \} \\ \text{follow}(e, i) &= \{ \langle b, j \rangle \mid x p_i b_1 b_2 \dots b_n q_j y \in L(\bar{e}), b = \bigvee (b_1 \wedge b_2 \wedge \dots b_n) \} \\ \text{last}(e) &= \{ \langle b, i \rangle \mid x p_i b_1 b_2 \dots b_n \in L(\bar{e}), b = \bigvee (b_1 \wedge b_2 \wedge \dots b_n) \} \end{aligned}$$

The automaton corresponding to  $e$  is defined as :

$$\mathcal{A}_e = (\text{pos}(e), \Sigma, \langle o, \delta \rangle, 0)$$

where

$$o(i) = \begin{cases} E(e) & i = 0 \\ b & \langle b, i \rangle \in \text{last}(e) \\ \perp & \text{otherwise} \end{cases} \quad \delta = \{ \langle i, \langle b, p \rangle, j \rangle \mid \langle b, j \rangle \in \text{follow}(e, i), p = \bar{p}_i \} \cup \{ \langle 0, \langle b, a \rangle, j \rangle \mid \langle b, j \rangle \in \text{first}(e), a = \bar{a}_j \}$$

# Adapted Berry-Sethi for KAT

## Proposition

$$GS(\mathcal{A}_e) = GS(e)$$

## Proposition

For a KAT expression  $e$ ,  $F(e, \{\langle \top, ! \rangle\})$ , as defined below, yields pairs of the form  $\langle i, \text{follow}(E!, i) \rangle$ .

$$\begin{aligned} F(e_1 + e_2, S) &= F(e_1, S) \cup F(e_2, S) \\ F(e_1 e_2, S) &= F(e_1, \text{first}(e_2) \cup E(e_2) \triangleleft S) \cup F(E_2, S) \\ F(e_1^*, S) &= F(e_1, \text{first}(E_1) \cup S) \\ F(a_i, S) &= \{\langle i, S \rangle\} \quad F(b, S) = \emptyset \end{aligned}$$

$$b \triangleleft S = \{\langle b \wedge b', p \rangle \mid \langle b', p \rangle \in S\} \text{ (with } \perp \triangleleft S = \emptyset)$$

# Adapted Berry-Sethi for KAT

## Proposition

$$GS(\mathcal{A}_e) = GS(e)$$

## Proposition

For a KAT expression  $e$ ,  $F(e, \{\langle \top, ! \rangle\})$ , as defined below, yields pairs of the form  $\langle i, \text{follow}(E!, i) \rangle$ .

$$\begin{aligned} F(e_1 + e_2, S) &= F(e_1, S) \cup F(e_2, S) \\ F(e_1 e_2, S) &= F(e_1, \text{first}(e_2) \cup \textcolor{red}{E}(e_2) \triangleleft S) \cup F(e_2, S) \\ F(e_1^*, S) &= F(e_1, \text{first}(E_1) \cup S) \\ F(a_i, S) &= \{\langle i, S \rangle\} \quad F(\textcolor{red}{b}, S) = \emptyset \end{aligned}$$

$$b \triangleleft S = \{\langle b \wedge b', p \rangle \mid \langle b', p \rangle \in S\} \text{ (with } \perp \triangleleft S = \emptyset)$$

# Examples revisited

- $(b + p)qr^*$
- $b_1 + pb_2(b_3 + q)^*rb_4$

# Conclusions

- Extended Berry-Sethi's algorithm to KAT.
- New type of automaton acceptor of guarded strings.
- Still need complexity analysis.

Thanks!

# Conclusions

- Extended Berry-Sethi's algorithm to KAT.
- New type of automaton acceptor of guarded strings.
- Still need complexity analysis.

Thanks!