

# Derivatives of Regular Expressions and an Application

Haiming Chen<sup>1</sup> and Yu Shen<sup>2</sup>

<sup>1</sup>State Key Laboratory of Computer Science, ISCAS

<sup>2</sup>Department of Computer Science, University of  
Western Ontario

WTCS 2012, Auckland 2012.2.22

Derivatives  $w^{-1}(E)$  Brzozowski, 1964

→ (left) quotient of language  $w^{-1}L(E)$

Berry and Sethi 's Result 1986

Derivatives of  $E$  → classes of *similar* derivatives

Why?  $E$  linear

Our work

A characterization of the structure of derivatives of linear  $E$  implies Berry and Sethi 's Result

## Derivatives

Berry and Sethi's result

Structure of derivatives

Properties of repeating terms

An application

Conclusion

## Regular expressions

$$E ::= \phi \mid \varepsilon \mid a \in \Sigma \mid E + E \mid EE \mid E^*$$

### ACI-similar

$$E_1 \sim_{\text{aci}} E_2$$

Associativity  $(E_1 + E_2) + E_3 = E_1 + (E_2 + E_3)$

Commutativity  $E_1 + E_2 = E_2 + E_1$

Idempotence  $E + E = E$

## Marked expressions

$$(a+b)^* ab(a+b) E \quad (a_1+b_1)^* a_2 b_2 (a_3+b_3) \bar{E}$$

The same notation used for dropping of subscripts:

$$\bar{\bar{E}} = E$$

### Note

Marking is not unique

For example  $(a_1+b_2)^* a_3 b_4 (a_5+b_6)$

(left) quotient set of a language  $L$

$$w^{-1}(L) = \{u \mid wu \in L\}$$

$$L = wL(w^{-1}(L))$$

Derivatives (Brzozowski)

$$a^{-1}(\emptyset) = a^{-1}(\varepsilon) = \emptyset$$

$$a^{-1}(b) = \begin{cases} \varepsilon, & \text{if } b = a \\ \emptyset, & \text{otherwise} \end{cases}$$

$$a^{-1}(F + G) = a^{-1}(F) + a^{-1}(G)$$

$$a^{-1}(FG) = \begin{cases} a^{-1}(F)G + a^{-1}(G), & \text{if } \varepsilon \in L(F) \\ a^{-1}(F)G, & \text{otherwise} \end{cases}$$

$$a^{-1}(F^*) = a^{-1}(F)F^*$$

$$\varepsilon^{-1}(E) = E, (wa)^{-1}(E) = a^{-1}(w^{-1}(E))$$

$$L(w^{-1}(E)) = w^{-1}(L(E))$$

Derivatives

Berry and Sethi's result

Structure of derivatives

Properties of repeating terms

An application

Conclusion

Regular expressions with distinct symbols (linear):  
One symbol occurs only once

Next we consider this kind of expressions



Derivatives

Berry and Sethi's result

Structure of derivatives

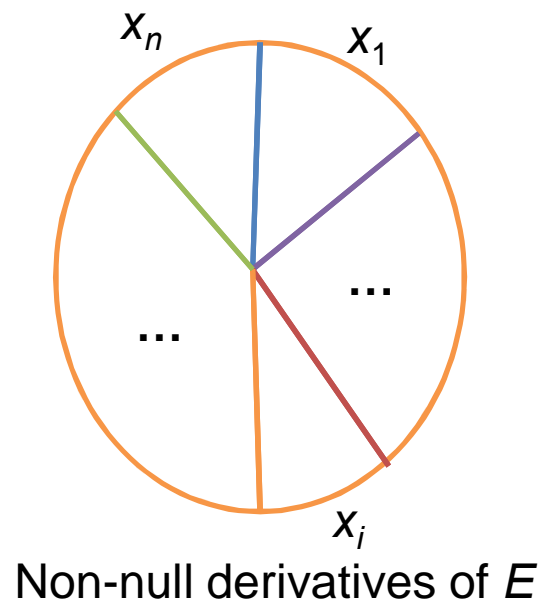
Properties of repeating terms

An application

Conclusion

## Berry and Sethi proved that

*Let all symbols in  $E$  be distinct. Given a fixed  $x \in \Sigma_E$ ,  $(wx)^{-1}(E)$  is either  $\emptyset$  or unique modulo  $\sim_{aci}$  for all words  $w$ .*



$$|\Sigma_E| = n$$

**Theorem 1** Let all symbols in  $E$  be distinct. Given a fixed  $x \in \Sigma_E$ , for all words  $w$ , each non-null  $(wx)^{-1}(E)$  must be of one of the following forms:  $F$  or  $F + \dots + F$ , where  $F$  is a non-null regular expression called the **repeating term** of  $(wx)^{-1}(E)$  which does not contain  $+$  at the top level.

**Example 1** Let  $E = (a + b)(a^* + ba^* + b^*)^*$ , then

$$\begin{aligned} \bar{E} &= (a_1 + b_2)(a_3^* + b_4a_5^* + b_6^*)^*, \\ a_1^{-1}(\bar{E}) &= (a_3^* + b_4a_5^* + b_6^*)^* = \tau_1, \\ (a_1a_3)^{-1}(\bar{E}) &= a_3^{-1}(\tau_1) = a_3^*\tau_1 = \tau_2, \\ (a_1a_3a_3)^{-1}(\bar{E}) &= a_3^{-1}(\tau_2) = \tau_2 + \tau_2, \\ &\dots \end{aligned}$$

repeating term for  
 $(wa_1)^{-1}(E): \tau_1$   
 $(wa_3)^{-1}(E): \tau_2$

Denote by  $rt_x(E)$  the repeating term of  $(wx)^{-1}(E)$

**Corollary 1** Let all symbols in  $E$  be distinct. If  $(wx)^{-1}(E)$  is non-null, then  $(wx)^{-1}(E) \sim_{\text{aci}} rt_x(E)$ .

a more precise version of Berry and Sethi's result

**Q:** For each  $x \in \Sigma_E$ , whether there is a non-null  $(wx)^{-1}(E)$  containing **one**  $rt_x(E)$ , that is,  $rt_x(E)$  is a derivative of  $E$ .

**A:** positive see below

**The first appearance**  $F_x(E)$

$ind : \Sigma_E \rightarrow \{1, \dots, \|E\|\}$ :  $ind(x) = d$  if  $x$  is the  $d$ th occurrence of symbols from left to right in  $E$

$x < y$  iff  $ind(x) < ind(y)$   $x, y \in \Sigma_E$

$w_1 \prec w_2$  if either  $|w_1| < |w_2|$ , or  $|w_1| = |w_2|$  and

let  $w_1 = x_1 \dots x_n, w_2 = x'_1 \dots x'_n, 1 \leq k \leq n,$

$x_t = x'_t$  for  $t = 1, \dots, k - 1$ , and  $x_k < x'_k$

A non-null  $(wx)^{-1}(E)$  is called the *first appearance* of derivative of  $E$  w.r.t.  $x$  if for any other non-null  $(w_1x)^{-1}(E)$  it has  $w \prec w_1$

**Example 2** For  $E = (a+b)(a^*+ba^*+b^*)^*$ ,  $\overline{E} = (a_1+b_2)(a_3^*+b_4a_5^*+b_6^*)^*$

$$\begin{array}{ll}
\underline{a_1}^{-1}(\overline{E}) = (a_3^*+b_4a_5^*+b_6^*)^* = \tau_1, & \underline{b_2}^{-1}(\overline{E}) = (a_3^*+b_4a_5^*+b_6^*)^* = \tau_1, \\
(\underline{a_1a_3})^{-1}(\overline{E}) = a_3^{-1}(\tau_1) = a_3^*\tau_1 = \tau_2, & (\underline{a_1b_4})^{-1}(\overline{E}) = b_4^{-1}(\tau_1) = a_5^*\tau_1 = \tau_3, \\
(\underline{a_1b_6})^{-1}(\overline{E}) = b_6^{-1}(\tau_1) = b_6^*\tau_1 = \tau_4, & (\underline{b_2a_3})^{-1}(\overline{E}) = a_3^{-1}(\tau_1) = \tau_2, \\
(\underline{b_2b_4})^{-1}(\overline{E}) = b_4^{-1}(\tau_1) = \tau_3, & (\underline{b_2b_6})^{-1}(\overline{E}) = b_6^{-1}(\tau_1) = \tau_4, \\
(\underline{a_1a_3a_3})^{-1}(\overline{E}) = a_3^{-1}(\tau_2) = \tau_2 + \tau_2, & (\underline{a_1a_3b_4})^{-1}(\overline{E}) = b_4^{-1}(\tau_2) = \tau_3, \\
(\underline{a_1a_3b_6})^{-1}(\overline{E}) = b_6^{-1}(\tau_2) = \tau_4, & (\underline{a_1b_4a_3})^{-1}(\overline{E}) = a_3^{-1}(\tau_3) = \tau_2, \\
(\underline{a_1b_4b_4})^{-1}(\overline{E}) = b_4^{-1}(\tau_3) = \tau_3, & (\underline{a_1b_4a_5})^{-1}(\overline{E}) = a_5^{-1}(\tau_3) = \tau_3.
\end{array}$$

**Proposition 1** Let all symbols in  $E$  be distinct. Given a fixed  $x \in \Sigma_E$ , the first appearance  $F_x(E)$  consists of only one repeating term.

The choice of the order is not significant.

**Proposition 2** Let all symbols in  $E$  be distinct. Given any words  $w_1, w_2 \in \Sigma_E^*$  and  $x \in \Sigma_E$ , if  $|w_1| = |w_2|$  and  $(w_1x)^{-1}(E), (w_2x)^{-1}(E) \neq \emptyset$ , and there is no  $w$ , such that  $|w| < |w_1|$  and  $(wx)^{-1}(E) \neq \emptyset$ , then  $(w_1x)^{-1}(E) = (w_2x)^{-1}(E)$ .

**Proposition 3** Let all symbols in  $E$  be distinct. There exists a word  $w \in \Sigma_E^*$  for each  $x \in \Sigma_E$ , such that  $(wx)^{-1}(E) = rt_x(E)$ .

Thus repeating terms are derivatives of  $E$ , and any non-null derivative of  $E$  is built from one of them.

Derivatives

Berry and Sethi's result

Structure of derivatives

Properties of repeating terms

An application

Conclusion

**Proposition 7** *Let all symbols in  $E$  be distinct. For each  $x \in \Sigma_E$ ,*

*(1)  $rt_x(E)$  exists, and  $rt_x(E) \neq \emptyset$ .*

*(2)  $rt_x(E)$  is unique.*

There cannot be two  $rt_x(E)$  for  
 $(wx)^{-1}(E)$



**Example 3** For  $E = (a + b)(a^* + ba^* + b^*)^*$ ,  $\bar{E} = (a_1 + b_2)(a_3^* + b_4a_5^* + b_6^*)^*$ .

$$\begin{aligned} rt_{a_1}(\bar{E}) &= rt_{a_1}(a_1 + b_2)(a_3^* + b_4a_5^* + b_6^*)^* = rt_{a_1}(a_1)(a_3^* + b_4a_5^* + b_6^*)^* \\ &= \varepsilon(a_3^* + b_4a_5^* + b_6^*)^* = (a_3^* + b_4a_5^* + b_6^*)^* = \tau_1, \end{aligned}$$

$$rt_{b_2}(\bar{E}) = \varepsilon(a_3^* + b_4a_5^* + b_6^*)^* = \tau_1,$$

$$\begin{aligned} rt_{a_3}(\bar{E}) &= rt_{a_3}(a_3^* + b_4a_5^* + b_6^*)^* = rt_{a_3}(a_3^* + b_4a_5^* + b_6^*)\tau_1 = rt_{a_3}(a_3^*)\tau_1 \\ &= rt_{a_3}(a_3)a_3^*\tau_1 = a_3^*\tau_1 = \tau_2, \end{aligned}$$

$$rt_{b_4}(\bar{E}) = rt_{b_4}(a_3^* + b_4a_5^* + b_6^*)^* = rt_{b_4}(b_4a_5^*)\tau_1 = a_5^*\tau_1 = \tau_3,$$

$$rt_{a_5}(\bar{E}) = rt_{a_5}(a_3^* + b_4a_5^* + b_6^*)^* = rt_{a_5}(b_4a_5^*)\tau_1 = a_5^*\tau_1 = \tau_3, \text{ and}$$

$$rt_{b_6}(\bar{E}) = rt_{b_6}(a_3^* + b_4a_5^* + b_6^*)^* = rt_{b_6}(b_6^*)\tau_1 = b_6^*\tau_1 = \tau_4.$$

**Proposition 8** Let all symbols in  $E$  be distinct. If there are non-null  $(w_1 x_1)^{-1}(E)$  and  $(w_2 x_2)^{-1}(E)$ , such that  $(w_1 x_1)^{-1}(E) \sim_{\text{aci}} (w_2 x_2)^{-1}(E)$ , then  $rt_{x_1}(E) = rt_{x_2}(E)$ , and vice versa.

**Corollary 2** Let all symbols in  $E$  be distinct. If  $rt_{x_1}(E) \sim_{\text{aci}} rt_{x_2}(E)$ , then  $rt_{x_1}(E) = rt_{x_2}(E)$ .

**Remark**  $rt_x(E)$ 's are 'atomic' **building blocks**

- (1) Each non-null  $(wx)^{-1}(E)$  is uniquely decomposed into a sum of  $rt_x(E)$ , that is,  $(wx)^{-1}(E) = \sum rt_x(E)$ .
- (2)  $rt_x(E)$  and  $rt_y(E)$  are either identical, or not equivalent modulo  $\sim_{\text{aci}}$ , if  $x \neq y$ .

Derivatives

Berry and Sethi's result

Structure of derivatives

Properties of repeating terms

An application

Conclusion

Solves an issue in using Berry and Sethi's result: find a **unique representative** for  $(wx)^{-1}(E)$

## Glushkov automaton

$$M_{\text{pos}}(E) = (Q_{\text{pos}}, \Sigma, \delta_{\text{pos}}, q_E, F_{\text{pos}}),$$

where

1.  $Q_{\text{pos}} = \Sigma_{\overline{E}} \cup \{q_E\}$ ,  $q_E$  is a new state not in  $\Sigma_{\overline{E}}$
2.  $\delta_{\text{pos}}(q_E, a) = \{x \mid x \in \text{first}_{\Sigma_{\overline{E}}}(\overline{E}), \bar{x} = a\}$  for  $a \in \Sigma$
3.  $\delta_{\text{pos}}(x, a) = \{y \mid y \in \text{follow}(\overline{E}, x), \bar{y} = a\}$  for  $x \in \Sigma_{\overline{E}}$  and  $a \in \Sigma$
4.  $F_{\text{pos}} = \begin{cases} \text{last}(\overline{E}) \cup \{q_E\}, & \text{if } \varepsilon \in L(E), \\ \text{last}(\overline{E}), & \text{otherwise} \end{cases}$

Berry and Sethi showed the class of derivatives  $\{(wx)^{-1}(E)\}$  corresponds to a state  $x$  of  $M_{\text{pos}}(E)$ ,  $x \in \Sigma_{\overline{E}}$

In many cases, however, one needs a unique representative for the class of  $\{(w x)^{-1}(E)\}$  to correspond to a state  $x$

By the work, the representatives are obtained immediately

# An improvement of Ilie and Yu's proof presented in (Ilie & Yu 2003)

A proof about the quotient relation between Glushkov and partial derivative automata

## Partial derivatives

$$\begin{aligned}\partial_a(\emptyset) &= \partial_a(\varepsilon) = \emptyset \\ \partial_a(b) &= \begin{cases} \{\varepsilon\}, & \text{if } b = a \\ \emptyset, & \text{otherwise} \end{cases} \\ \partial_a(F + G) &= \partial_a(F) \cup \partial_a(G) \\ \partial_a(FG) &= \begin{cases} \partial_a(F)G \cup \partial_a(G), & \text{if } \varepsilon \in L(F) \\ \partial_a(F)G, & \text{otherwise} \end{cases} \\ \partial_a(F^*) &= \partial_a(F)F^* \\ \partial_\varepsilon(E) &= \{E\}, \quad \partial_{wa}(E) = \bigcup_{p \in \partial_w(E)} \partial_a(p) \\ PD(E) &= \bigcup_{w \in \Sigma^*} \partial_w(E)\end{aligned}$$



## Partial derivative automaton

$$M_{\text{pd}}(E) = (PD(E), \Sigma, \delta_{\text{pd}}, E, \{q \in PD(E) \mid \varepsilon \in L(q)\}),$$

where  $\delta_{\text{pd}}(q, a) = \partial_a(q)$ , for any  $q \in PD(E)$ ,  $a \in \Sigma$ .

$M_{\text{pd}}(E)$  is a quotient of  $M_g(E)$

## Ilie and Yu's proof

- . The central issue is to find a unique representative for a class of derivatives
- . The proof fails to find the correct representatives

It is claimed in the proof that, by using the rules ( $\phi\varepsilon$ -rules), for a fixed  $x \in \Sigma_{\overline{E}}$  and for all words  $w$ ,  $(wx)^{-1}(\overline{E})$  is either  $\phi$  or unique.  
 incorrect

$$E + \emptyset = \emptyset + E = E$$

$$E\emptyset = \emptyset E = \emptyset$$

$$E\varepsilon = \varepsilon E = E$$

Rules ( $\phi$ -rules)

Example. In [Example 1](#),  $(a_1 a_3)^{-1}(\overline{E})$  and  $(a_1 a_3 a_3)^{-1}(\overline{E})$  are distinct

## An improved proof

Use  $rt_x(\overline{E})$  as the unique representative.

See our paper

Derivatives

Berry and Sethi's result

Structure of derivatives

Properties of repeating terms

An application

Conclusion

A characterization of the structure of derivatives

Several properties

An application

A useful technique

Thanks!